

Toward Automatic Summarization of Hospital Discharge Notes

by

Paul Landes

B.S. (University of North Texas) 1999

B.S. (University of Illinois Chicago) 2016

Thesis submitted in partial fulfillment of the requirements
for the degree of Doctor of Philosophy in Computer Science
in the Graduate College of the
University of Illinois Chicago, 2024

Chicago, Illinois

Defense Committee:

Barbara Di Eugenio, Chair and Advisor

Cornelia Caragea

Anastasios Sidiropoulos

Andrew Boyd, Department of Biomedical and Health Information Sciences

Martha Palmer, University of Colorado Boulder

Copyright by

Paul Landes

2024

To Steph.

ACKNOWLEDGMENTS

I wish to thank my advisor, Prof. Barbara Di Eugenio, for her invaluable assistance in shaping me into the scientist, researcher and author that I have become. Her dedication, insight and moral compass are impeccable.

I am profoundly thankful to Prof. Andrew Boyd for his insightful recommendation to undertake the Herculean task of discharge summarization. His medical direction and resourcefulness were fundamental to this undertaking.

My gratitude extends to Prof. Anastasios Sidiropoulos for his exceptional ability to teach both theoretical concepts and unconventional algorithms, which greatly contributed to the solution of this work's problem.

The guidance from Prof. Cornelia Caragea insofar as scientific thinking and writing within the framework of state-of-the-art methods and practices has been invaluable to me. I am also grateful for her encouragement and belief in my abilities.

I am very fortunate to have met Prof. Martha Palmer, whose generosity with her time and energy was only matched by her expert guidance in linguistics and graph theory.

I am also thankful to have had the pleasure and opportunity to work such great physicians and colleagues Kunal Patel, Sean Huang, Aaron Chaise, Adam Webb, and Sitara Rao.

PL

CONTRIBUTION OF AUTHORS

All chapters of the thesis were edited and revised by my advisor Prof. Barbara Di Eugenio with additional input by Prof. Andrew Boyd. Chapter 1 introduces the work and Chapter 2 discusses related and previous work in context to the research topic. Prof. Martha Palmer provided significant feedback and edits for Chapters 4, 5, 6 and 8. The referenced manuscripts in Chapters 3, 4, 5 (Landes and Di Eugenio, 2024) and 7 (Landes et al., 2023) were advised and edited by Prof. Barbara Di Eugenio. Prof. Cornelia Caragea co-advised and edited the manuscript referenced in Chapter 3 and Chapter 7. I was the primary author and research lead on the work presented in Chapters 3, 4 and 6. However, medical informatics fellows contributed to the clinical research of the datasets. Kunal Patel, Sean Huang, and Adam Webb annotated and analyzed MedSecId dataset in Chapter 3 (Landes et al., 2022). Chapter 4 (Landes et al., 2023) was annotated and analyzed by Kunal Patel, Sean Huang, and Aaron Chaise. The clinical analysis and annotation in Chapter 6 was contributed by Aaron Chaise and Sitara Rao.

TABLE OF CONTENTS

<u>CHAPTER</u>		<u>PAGE</u>
1	INTRODUCTION	1
1.1	Clinical Summarization	2
1.2	Contributions	5
1.2.1	Summarization Approach	6
1.2.2	Methods	9
1.2.2.1	Medical Note Section Identification	10
1.2.2.2	Graph Alignment	10
1.2.2.3	Learning to Summarize	11
1.2.2.4	Natural Language Generation	12
1.3	Corpora	12
1.4	Outline	13
2	RELATED WORK	15
2.1	Summarization	15
2.1.1	A Brief Chronology	15
2.1.2	Neural Network Summarization Methods	18
2.1.3	Summarization Scoring Methods	19
2.1.4	Summarization using Graphs and Abstract Meaning Representation	21
2.2	Relevant Clinical Applications and Summarization	23
2.2.1	Abstract Meaning Representation in Medical Research	23
2.2.2	Medical Note Section Identification	24
2.2.3	Provenance of Data	25
2.2.4	Medical and Discharge Note Summarization	26
2.3	Graph Alignment	27
2.3.1	Flow Networks	27
2.3.2	Alignment	29
2.4	Frameworks	30
3	MEDICAL SECTION IDENTIFICATION	33
3.1	Motivation	33
3.2	Dataset	36
3.2.1	Annotation Process	38
3.2.2	Final Annotation and IAA Computation	41
3.2.3	Data Analysis	44
3.3	Limitations	45
3.4	Baseline Models	47

TABLE OF CONTENTS (Continued)

<u>CHAPTER</u>		<u>PAGE</u>
	3.4.1 Implementation Details	49
	3.5 Results	50
	3.6 Conclusion	51
4	DISCHARGE SUMMARY PROVENANCE OF DATA	53
	4.1 Motivation	55
	4.2 Dataset	56
	4.2.1 Limitations	59
	4.2.2 Data Analysis	59
	4.2.2.1 Note Category	61
	4.2.2.2 Section by Discharge Summary	62
	4.2.2.3 Section by Note Antecedent	63
	4.3 Methods	64
	4.3.1 Evaluation	65
	4.3.2 Word Mover	65
	4.3.3 Hybrid Semantic Positional Token Clustering	66
	4.4 Results	69
	4.5 Conclusion	72
5	AMR GRAPH ALIGNMENT	73
	5.1 Introduction	73
	5.1.1 Abstract Meaning Representation	74
	5.1.2 Flow Networks	76
	5.2 Motivation	77
	5.3 Alignment Method	79
	5.4 AMR Graph Embedding	82
	5.4.1 PropBank	82
	5.4.2 Node to Text Embeddings	84
	5.4.3 Document Nodes	85
	5.4.4 Sentence Nodes	86
	5.4.5 Concept Nodes	86
	5.4.6 Attribute Nodes	91
	5.4.7 Network Neighborhood	91
	5.5 Capacity Calculation	94
	5.6 Alignment Graph Construction	99
	5.6.1 Graph Component Construction	101
	5.6.2 Graph Component Connection	104
	5.6.3 Complete Flow Network	104
	5.7 Alignment Graph Algorithm	106
	5.7.1 Max Flow	107
	5.7.2 Flow Normalization	108
	5.7.3 Capacity Constriction	109

TABLE OF CONTENTS (Continued)

<u>CHAPTER</u>		<u>PAGE</u>
5.7.4	Summary Max Flow	111
5.7.5	Final Alignment Graph	112
5.8	Scoring Method	112
5.9	Experiment Design and Setup	113
5.9.1	Summarization	114
5.9.2	Similarity Scoring	114
5.9.3	Baselines	115
5.10	Results	116
5.10.1	Summarization	116
5.10.2	Alignment	117
5.10.3	Previous Methods	118
5.11	Reentrancies	119
5.12	Conclusion	122
5.13	Alignment Graph Examples	122
6	SUMMARIZATION	127
6.1	Datasets	129
6.2	Methods	131
6.2.1	Admission Graph	132
6.2.2	Concept Variable Renaming	134
6.2.3	Sentence Matching Algorithm	135
6.2.4	Source Section Dataset	137
6.3	Discharge Summary Generation	139
6.3.1	Source Section Model	142
6.3.2	Experimental Setup	143
6.4	Results	144
6.4.1	Parser Evaluation	144
6.4.2	Source Section Model Performance	146
6.4.3	Discussion	147
6.5	Conclusion	148
7	SOFTWARE FRAMEWORK	153
7.1	Motivation	154
7.2	Library Design	155
7.3	NLP-Focused Abstractions and Features	156
7.4	Vectorization	158
7.5	Batching	158
7.6	Technology Stack	161
7.7	Execution	162
7.8	Runtime Analysis	164
7.9	Conclusion	167

TABLE OF CONTENTS (Continued)

<u>CHAPTER</u>		<u>PAGE</u>
8	CONCLUSIONS AND FUTURE WORK	168
8.1	Contributions	168
8.2	Future Work	170
8.2.1	Abstractive Summarization	171
8.2.1.1	Training	172
8.2.1.2	Testing	175
8.2.1.3	Generation	177
8.2.2	Model Tuning and Optimization	177
8.2.3	Application to Large Language Models	178
	APPENDICES	180
	Appendix A	181
	Appendix B	186
	Appendix C	187
	Appendix D	189
	Appendix E	190
	CITED LITERATURE	191
	VITA	210

LIST OF TABLES

<u>TABLE</u>		<u>PAGE</u>
I	Annotated Medical Note Categories	37
II	Most Frequently Annotated Medical Sections	40
III	MedSecId Inter-annotator Agreement	42
IV	MedSecId Sentence Distribution	48
V	MedSecId Performance	50
VI	MedSecId BiLSTM-CRF _{tok} Performance	52
VII	DSProv Corpus Statistics	59
VIII	Statistics by MIMIC Note Category	61
IX	Statistics Grouped by Discharge Summary Section Type	63
X	Statistics Grouped by Note Antecedent Section Type	64
XI	Note Antecedent Score	71
XII	Document Summarization Scores	117
XIII	Parser Alignment Scoring	118
XIV	Aligned Node Text Comparison	119
XV	Summary Alignment Coverage	120
XVI	Reentrancy Repair Statistics	121
XVII	Graph Alignment Statistics	132
XVIII	MIMIC-III Matched Sentence Notes	138
XIX	MIMIC-III Matched Sentence Sections	139
XX	UI Health Matched Sentence Sections	140
XXI	Informal Evaluation Questions	143
XXII	Source Section Model Results	147
XXIII	Informal Evaluation	148
XXIV	DeepZensols Efficiency Benchmarks	166

LIST OF FIGURES

<u>FIGURE</u>		<u>PAGE</u>
1	Discharge Summary	4
2	Process Overview	8
3	Incremental Tree Compression	17
4	Summarized Text Overlap	20
5	Soviet Rail Network, 1955	28
6	Discharge Summary	36
7	Concept Unique Identifier Plots	45
8	MedSecId Baseline Models	46
9	A Note Match Annotation	54
10	Hybrid Semantic Positional Token Clustering	68
11	Parsed Abstract Meaning Graph	75
12	Penman Example	76
13	Flow Network	77
14	AMR Graph Components	79
15	Graph Construction	80
16	Graph Embeddings	88
17	Network Neighborhood Weights	92
18	Sentence Scaled Capacities	98
19	Graph Construction	103
20	Bipartite Graph Connection	105
21	Flow Network Construction	106
22	Network Flow	109
23	Reentrancies Alignment Failure	120
24	Disconnected Source and Summary AMR Graphs	123
25	Source Component Alignment	124
26	Summary Component Alignment	125
27	Final Flow	126
28	Aligned Admission Graph	131
29	Admission Graph	133
30	Admission Graph Variable Renaming	135
31	Sentence Matching Algorithm	136
32	MIMIC-III Note to Section Alignment Contingency	141
33	UI Health Note to Section Alignment Contingency	150
34	Source Section Model	151
35	Generated Discharge Summary	152
36	Network Architecture	157
37	Batch Process	159
38	Batch Reassembly	160

LIST OF FIGURES (Continued)

<u>FIGURE</u>		<u>PAGE</u>
39	Batch Decoding	161
40	Framework Library Stack	163
41	Validation and Training Loss	164
42	Summarization Iteration	173
43	Prediction with the Summarization Network	176
44	Gold Discharge Summary	188

LIST OF ABBREVIATIONS

API	advanced programming interface p. 30
BiLSTM	bi-directional long-short term memory p. 9
BiLSTM-CRF	BiLSTM with a CRF output layer p. 47
CNN	convolutional neural network p. 25
CRF	conditional random field p. 24
DL	deep learning p. 2
EHR	electronic health record p. 3
DAG	directed acyclic graph p. 74
GloVe	Global Vectors for Word Representation p. 47
GraphRNN	deep auto-regressive model that generates graphs p. 171
PHI	protected health information p. 146
HPI	history of present illness p. 24
IAA	inter-annotator agreement p. 41
ICU	intensive care unit p. 36
ILP	integer linear programming p. 6
IOB	in, out, begin p. 24
IO	in, out p. 47
IR	information retrieval p. 12
JAMR	Discriminative Graph-Based Parser for AMR p. 115
LDA	latent Dirichlet allocation p. 25
LLM	large language model p. xv
LSA	latent semantic analysis p. 157
LSTM	long-short term memory p. 23
MaxEnt	Maximum Entropy p. 24
ML	machine learning p. 10
MT	automatic machine translation p. 21
NLG	natural language generation p. 21
NLP	natural language processing p. 6
NMT	neural machine translation p. 18
NN	neural network p. 2

LIST OF ABBREVIATIONS (Continued)

OSCE	objective structured clinical examination p. 34
PCA	principal component analysis p. 44
POS	part of speech p. 18
QA	question/answer systems p. 21
RNN	recurrent neural network p. 18
RVU	relative value unit p. 33
SI	section identification p. 5
S-LSTM	segment pooling LSTM p. 25
SoTA	state of the art p. 1
SRL	semantic role labeling p. 22
SVM	support vector machine p. 24
TF/IDF	term frequency/inverse document frequency p. 16
THYME	Temporal Histories of Your Medical Events p. 133
UI Health	University of Illinois Chicago p. 13
UMLS	Unified Medical Language System p. 27
word2vec	word-to-vector p. 25

SUMMARY

Large language models have recently have become pervasively used for automatically generated summaries. However, large language models often generate summarized text not found in the source or text that is factually incorrect. The goal of this research is to discover new methods to ground automatic summaries using language-based graphics and deep learning models. The medical domain is of particular importance in creating factual summaries. Specifically, this work focuses on generating discharge summaries, which are narrative reports written by physicians. These medical notes summarize a hospital in-patient visit including historical information about the patient and what happened during the hospital stay. These summaries should not only be faithful to the source electronic health record data, but also allow for traceability of the summarized text back to the source. The method explored in this dissertation is one that uses abstract meaning representation graphs to model the source data and the summarization, which allows for more transparency and grounding as they are human readable and aligned with their source text.

This work reframes discharge summaries as multi-document multi-summarization since the output consists of more than one medical section’s summary, and specifically, all human authored sections of the discharge summary. This is necessary since recent literature frames their automatic generation as a multi-document summarization task that includes only a single section. However, little work has been accomplished to summarize medical notes into sections of the discharge summary as these sections are highly diverse in content and contain formatted

SUMMARY (Continued)

structured data. This will be the first work that generates multiple sections discharge summaries framed as a multi-document summarization task using a novel graph alignment method.

The hypothesis this dissertation explores is that abstract meaning representation graphs are not only performant, but also provide summarization that is both faithful and traceable, and thus, meets the requirements of a discharge summary usable by hospitals and other clinical settings. In this dissertation, I explore this hypothesis by generating non-structured sections of discharge summaries.

CHAPTER 1

INTRODUCTION

Automatic summarization is the task of using machines to summarize natural language text (Luhn, 1958; Kupiec et al., 1995; Li et al., 2014; Ranjitha and Kallimani, 2017; Maynez et al., 2020). Automatic summarization is broadly classified as either extractive or abstractive: the extractive method copies selected sentences from the source text verbatim to the summary while abstractive methods generate unique text not always found in the source text. Generally, abstractive summaries have shown to be higher quality more readable (See et al., 2017).

However, extractive summarization provides text that is more faithful (how accurate the summary is) to the source text and traceable (if the summary can be traced back to its source content) since each sentence is appended as is to the summary (Dang and Owczarzak, 2008; Liu et al., 2015). Even though extractive summaries are at times choppy, lack coherence and are generally hard to read as surrounding context varies from the source (Knight and Marcu, 2000), they are a faithful representation of the source text.

Automatic summarization has developed across many areas in computational linguistics for more than six decades (Luhn, 1958). Recently LLMs have been influential and shown to achieve state of the art performance in summarization. However, these models can not summarize a large number of long documents given memory constraints. Hallucination, which is erroneous and nonfactual text automatically generated by a model, presents additional challenges to automatic summarization when using LLMs.

I address faithfulness and traceability by utilizing AMR (abstract meaning representation), which is a semantic representation language that describes the abstract meaning of a sentence. My method uses AMR in a novel algorithm to align the graph representation of the source and summary text. This strategy combines a deterministic algorithm to create the alignment with deep learning as input to a neural network supervised algorithm. This method provides traceable summaries since each graph represents a specific sentence with nodes aligned to source tokens from the source text (Lyu and Titov, 2018) that are faithful as a “linguistically-grounded semantic formalism” (Liao et al., 2018).

These alignments are then used by an extractive summarization method for the summarization. The choice of an extractive summarization method aids in generating both faithful and traceable summaries, which is of paramount importance, since this work’s domain is clinical medical notes. The summaries are implicitly faithful as they are copied from the source text. They are traceable because selected sentences are aligned via the AMR graphs, which are in turn aligned to the source text.

All methods provided in this dissertation are compatible with summarization in any domain. In fact, the only difference in methodology across any domain is the structure of the data such as placeholders for document and section types.

1.1 Clinical Summarization

Automatic summarization is becoming as important for personal decision-making as it is for academic and professional communities that face issues of “information overload” (Habernal and Gurevych, 2016). The medical field is no exception as the volume of data continues to increase

with the adoption of electronic health records (EHRs), which require constant updating, and can impede a clinician’s ability to carry out their basic duties and often leads to “physician burnout” (Hirsch et al., 2015). Sifting through the data needed to author a discharge summary is time consuming for a physician who could otherwise spend this time with the patient (Sinsky et al., 2016).

A discharge summary is written by physicians when a patient is discharged from the hospital and documents previous medical history, what happened during the hospital stay, and follow up instructions. These documents must be written by a physician 48 hours within the time the patient is discharged. An example of a discharge summary is given in Figure 1.

The motivation for my work is to create a system that obviates, or greatly reduces, the work on the part of the physician to write discharge summaries. Generating patient oriented discharge summaries is another motivation for this work as they are often synthesized from discharge notes (Acharya, 2019). Furthermore, automatically generated summaries could be provided to the patient upon leaving the hospital rather than up to 48 hours afterward when the physician completes chart notation.

The specific needs of clinical summarization differ from that of most other domains. Medical professionals need summarizations that are faithful to previously authored documents. They must also be traceable when cross-referencing additional detail. In terms of performance metrics, recall is key to clinical summarization because the most minor of details in a patient’s record can lead to improper future care. However, precision is equally important in the clinical setting as physicians have little time to search through irrelevant and redundant information.

Admission Date: [**2126-2-7**] Discharge Date: [**2126-2-20**] Date of Birth: [**2069-4-1**] Sex: M
history-of-present-illness HISTORY OF PRESENT ILLNESS: Mr. [**Known lastname **] is a 56-year-old male who experienced chest. . .
past-medical-history PAST MEDICAL HISTORY: Hypertension, former smoker with a 4- pack per day history for which he. . .
social-history SOCIAL HISTORY: He lives alone, and he works at [**Hospital3 2576**] as a cargo transporter.
medication-history MEDICATIONS ON ADMISSION: Aspirin 325 mg p.o. once a day, Toprol-XL 50 mg p.o. once a day.
allergies ALLERGIES: He had no known drug allergies.
labs-imaging PREOPERATIVE LABORATORY DATA: White count 6.0, hematocrit 33.3, platelet count 329,000. . .
He was discharged to home with VNA services in good condition on [**2126-2-20**]...

Figure 1: **Discharge Summary.** A MIMIC-III discharge summary note with the section type in **bold**; omitted text is indicated with ellipses.

To this end, the importance of faithful and traceable summaries takes priority for clinical documentation.

The benefits of clinical automatic summarization are as numerous as are the challenges to generate them. These challenges include condensing the vast amount of data taken from a patient during a hospital stay while eliminating redundancy, errors, and data incoherence (Cohen et al., 2013). A large corpus of copious medical notes paired with a discharge summary for each patient is necessary to sufficiently train a summarizer. One such corpus is Medical Information

Mart for Intensive Care III (MIMIC-III) (Johnson et al., 2016), a large freely accessible hospital database of ICU data from the Beth Israel Deaconess Medical Center in Boston, Massachusetts. However, this corpus is incomplete across some medical note categories adding to the difficulty of the task (Landes et al., 2022). In addition, this corpus is not unlike most discharge summaries in that up to 46% is copied and pasted from previous notes making meaningful summaries less useful (Adams et al., 2021). Creating faithful and traceable summaries, given the need for factual information, is the primary challenge with respect to the needs of medical professionals.

1.2 Contributions

The contributions of my work include:

- Clinical section identification annotations, deep learning methods and a pretrained model (see Chapter 3) for sectioning medical notes. This is an automatic method to identify and demarcate the lexical boundaries of clinical medical notes.
- A summarization feasibility study and annotations for the provenance of data from previously written clinical medical notes to discharge summaries (see Chapter 4). This study uses the section identification annotations for analysis and better section accuracy. The addition of the provenance annotations serves to create a richer dataset for summarization.
- A novel unsupervised method using the infrequently used word mover algorithm (Kusner et al., 2015). This method uses clustering to assist in automatically matching semantically similar, or copied and pasted text, as lexical spans between EHR notes and discharge summaries (see Section 4.3).

- A novel method utilizing the max flow algorithm (a rarely used algorithm in natural language processing) for graph bipartite matching (see Chapter 5) between the source component with the summary component to create more faithful summaries (see Section 6.2)
- Use of AMR graphs for summarization (see Chapter 5). While this summarization algorithm is inspired by the work of Liu et al. (2015) and Liao et al. (2018), it differs in the following ways:
 - Uses the aforementioned max flow based alignment method to find summarized content rather than graph reduction to classify edges for deletion,
 - Partitions the source and summary into sections, which is necessary for the discharge summary application,
 - Utilizes deep learning in place of integer linear programming.
- Summarization using AMR in the medical domain (see Chapter 6). This method adds intermediate nodes that are used as placeholders for document and section types. An additional paragraph node groups AMR sentences.
- A NLP deep learning framework extended to facilitate all of the contributions above (see Chapter 7). This includes models to train and classify section identification, to align graphs for summarized overlap, and to classify sentences for extractive summarization.

1.2.1 Summarization Approach

The end-to-end process used to summarize includes two steps as shown in Figure 2:

1. **Preprocessing:** Human annotated summaries are preprocessed in an intermediate data format (aligned graphs) to be used as training examples for supervised learning.
2. **Summarization:** The training examples created in the preprocessing step are used to train and evaluate a summarization model.

Preprocessing

The preprocessing step takes documents with sentence annotations indicating whether they are members of the source or summary. For this work, the document store is the EHR and the system knows summary sentences only come from discharge summary notes. Preprocessing includes:

1. **Section Documents:** Each document is sectioned. This is an optional step and not applicable to all domains. However, my work on sectioning clinical notes as explained in Chapter 3.
2. **Parse Documents into Graphs:** Each document is parsed into connected AMR graphs. After this step, the source component and summary component are created. The graph construction process is detailed in Chapter 5 and clinical constructed graph is given in Chapter 6.
3. **Align Source and Summary:** The source component and summary component are aligned to create a bipartite graph. The alignment method is explained in Section 5.3.

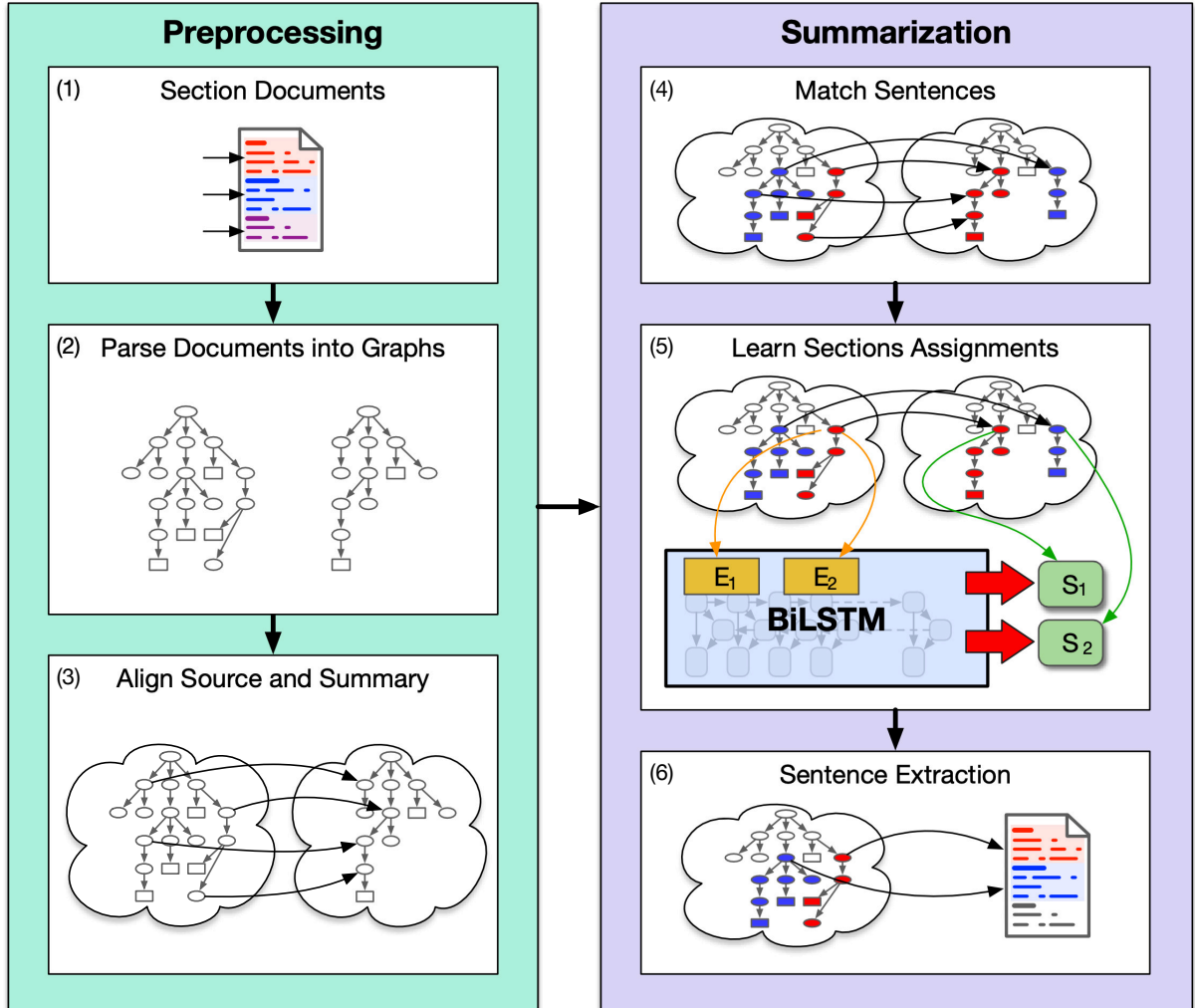


Figure 2: **Process Overview.** The process overview for summarization. The preprocessing step on the left depicts documents that are loaded from a repository, sectioned, parsed into AMR graphs, and aligned. The summarization step on the right shows how alignments are used to match sentences between the source and summary. A BiLSTM is then used to select sentences and classify in which summary section they belong. During inference time, summaries are generated by using the model to select sentences from the source documents.

Summarization

As with preprocessing, the extractive summarization method works with any domain data and not just clinical notes¹. The summarization method follows (see Chapter 6):

4. Match Sentences: The alignments are used to match sentences in the source component with sentences in the summary component based on the information gain between them (see Figure 2).
5. Learn Sections Assignments: Learn which section to slot each sentence, or to not include the sentence at all (see Figure 2). This uses a bi-directional long-short term memory (BiLSTM) neural network model (Graves and Schmidhuber, 2005) to classify sections from sentences.
6. Sentence Extraction: This is done at inference time to generate the summary using the classifications from the section assignment model (see Figure 2). In this work, the generated text is the discharge summary.

1.2.2 Methods

My work presents a novel algorithm to create training examples that provide faithful and traceable discharge summaries using previously human authored medical notes in an EHR. This is done by first creating an alignment between the EHR source text and the discharge summary

¹Documents must have at least one section, in which case, extraction becomes a “yes add the sentence” or “no do not add the sentence” boolean classification.

text. This alignment is then used as training data for a machine learning (ML) algorithm that summarizes free text medical source notes into discharge summaries.

1.2.2.1 Medical Note Section Identification

The process by which sections in a document are demarcated and labeled is known as section identification. Such sections are helpful to the reader when searching for information and contextualizing specific topics. Examples of sections include a patient’s family history or the patient’s recent illness. Because sections are structured and pertinent to some aspect of care, discharge summaries have been generated on a per section basis (Adams et al., 2021). The pipeline for this work uses the section identification annotations and trained model from MedSecId corpus (Landes et al., 2022) for the summarization process. The section identification module identifies section types and demarcates their lexical spans. These sections are added as document nodes to the source component and summary component to allow a per section summarization. Section identification is also important to identify sections to omit since many of them lack quality summarizable material or have formatted structured data rather than human written notes.

1.2.2.2 Graph Alignment

Flow networks have been used for applications such as matching medical students to institutions (Hopcroft and Karp, 1973), baseball team competition elimination (Schwartz, 1966), and airline scheduling (Cormen et al., 2001). The use of flow networks has been shown to solve many complex problems efficiently, which has led to its sustained interest over the decades.

It is surprising that despite the longevity of AMR and other similar graph models, that flow networks have yet to be exploited to solve problems in the NLP domain¹.

The alignment graph algorithm leverages the max flow algorithm using AMR as a linguistic semantic data structure. It takes a source text to be summarized, and the summary text for the given source text as input. It then produces an alignment of their AMR graph forms as an output. Each AMR graph has a root node, which is then made a child to a root node: one for the source text and another for the summary text.

The result is two separate disconnected graph components that represent the described source and summary in AMR having concepts as nodes and the relationship between each concept node pair as role edges. The algorithm then connects the source and summary components into one bipartite graph with alignment edges, which are pruned using max flow. The details of the algorithm are further explained in Chapter 5.

1.2.2.3 Learning to Summarize

Once the alignment graph algorithm has completed, the neural network is trained with the aligned source and summary components of the bipartite graph in a supervised fashion. The summarization algorithm is the algorithm that uses the alignment graph algorithm output to learn to summarize an AMR summary component by selecting sentences from the EHR notes to add to the discharge summary.

¹Guo et al., 2016 used a flow network for semantic matching information retrieval but not in the context of language graphical models. See flow network related work in Section 2.3.1.

Because the summarization algorithm uses AMR component alignment’s non-stochastic deterministic method of creating unsupervised training examples, it produces a summary component with antecedent nodes from the source that provides a traceable summarization. While previous work have used information retrieval (IR) methods with LLMs (Xu et al., 2024), I have found no other work that generates traceable summaries in any state of the art deep learning models. The closed form prediction provides the kind of faithful summaries that eliminate hallucinations present in any deep learning models such as BART. This summarization method is described further in Section 6.2.

1.2.2.4 Natural Language Generation

In the general summarization case, we simply iterate through each sentence level subgraph of the summary component converting each in turn. One strength of the alignment graph algorithm and summarization algorithm is that they allow us to add an arbitrary organization of document nodes that represent the nodes between the component roots and the sentence level. Specifically, this structural flexibility is useful in adding clinical section, which have a large breadth in clinical medical notes. This structural flexibility is the reason why the use of the alignment graph algorithm is well suited for summarizing EHR medical notes into the discharge summary.

1.3 Corpora

My work utilizes two publicly available corpora, a new dataset I created based on MIMIC-III, and one private corpus for training and testing models. They include general summarization and medical domain specific corpora:

- The LDC2020T02 AMR Annotation Release 3.0 (Knight et al., 2021) is a semantic treebank of English natural language sentences from broadcast conversations, newswire, weblogs, web discussion forums, fiction and web text. The corpus contains over 59,255 AMR sentences in Penman format annotated for sentence coreference, named entity annotation, modality, negation, questions, quantities, and represents the semantic structure of a sentence largely independent of its syntax.
- MIMIC-III Version 1.4 (Johnson et al., 2016) is a freely available accessible de-identified critical care database of patients admitted to Beth Israel Deaconess Medical Center in Boston, Massachusetts. The corpus has 46,520 distinct patients, 58,976 hospital admissions and 2,083,180 free form notes (59,652 of which are discharge summaries).
- MedSecId (Landes et al., 2022) is my contribution of a publicly available set of 2,002 fully annotated clinical medical notes from the MIMIC-III corpus annotated for section boundaries and identifiers. Annotations cover 50 section types over five medical note categories including discharge summary, physician notes, radiology, echo, and consult notes.
- The UI Health dataset is an IRB approved (protocol #2024-0109) private dataset of 11,001 admissions and 607,872 notes, which include daily progress, radiology, ECG and a variety of other notes from the University of Illinois Chicago hospital.

1.4 Outline

This dissertation has the following outline:

- Prior work in AMR, flow networks, cross domain summarization, and medical summarization is given in Chapter 2.
- An algorithm and methodology to segment and identify clinical notes is explained in Chapter 3. It also lists the statistics of the annotated dataset and the model results trained on the annotations.
- Chapter 4 details the discharge summary provenance of data study. The span-match method and results are also reported.
- The AMR graph alignment method is provided in Chapter 5. The chapter explains the method and its use as a summarization score.
- The summarization method is described in Chapter 6. This chapter demonstrates the use of the alignments for extractive summarization.
- The natural language processing oriented framework that was used facilitate the pipeline is presented in Chapter 7.
- Chapter 8 summarizes the goals of my dissertation, how they were achieved, and the contributions of each. It ends with a list of future research directions.

CHAPTER 2

RELATED WORK

This chapter is divided into four sections covering general summarization (see Section 2.1), biomedical and clinical specific domain summarization (see Section 2.2), graph alignment (see Section 2.3), and frameworks (see Section 2.4). The previous work in this chapter is bifurcated into clinical domain summarization, and summarization methods such as AMR (abstract meaning representation) and neural network (NN). However, there is cross-over in Section 2.1.2 with NNs for the medical domain.

2.1 Summarization

The task of summarization is defined as reducing the length of natural language source text to a shorter summary text. The general approach to summarization is categorized as either extractive or abstractive. The extractive approach selects text chunks verbatim, usually sentences, to add to the summary. The abstractive approach generates a distinctly different summary. The former creates a summary that is more faithful (how accurate the summary is) to the source, whereas the latter creates summaries that typically flow better and have stronger coherence, which increases readability. My work is best described as an extractive method.

2.1.1 A Brief Chronology

Summarization is a well established area in natural language processing (NLP) stretching back decades and involved statistical linguistic analysis for extractive summarization. The early

work of Luhn (1958) relied on statistical methods with word counts to create a “significance factor” for ranking sentences, which then were used to compose abstracts. Similarly, Salton et al. (1994) constructed simple graphs composed of token nodes and term frequency/inverse document frequency (TF/IDF) statistics to weight edges derived from term/document matrices (Sparck Jones, 1972). Statistical methods persisted for decades with such work as Bayesian classification for sentence matching in academic journals using an interesting set of heuristic features (Kupiec et al., 1995). Such features includes cue word lists and indicator phrases, which are token sequences likely to inform the summary.

Efforts in summarization transitioned to parse tree reduction with probabilistic models for the task of abstractive summarization and includes the seminal work of Knight et al. (2000). Their approach used a noisy channel and decision-based history model for sentence compression of parse trees for the generation step of news articles. The smaller compressed parse trees make up the abstractive generated summary. An example of this tree reduction algorithm is given in Figure 3 having the following operations:

1. Shift operations move words to a stack.
2. Reduce operations pop syntax trees to the top of the stack, then are combined and stacked again.
3. Drop the tokens that form the input that correspond to syntactic constituents.

Other tree editing approaches include syntax and grammar for sentence compression to rewrite the source text by Cohn and Lapata (2008). The authors compare supervised, semi-supervised and supervised learning with integer linear programming (ILP) using the Ziff-Davis

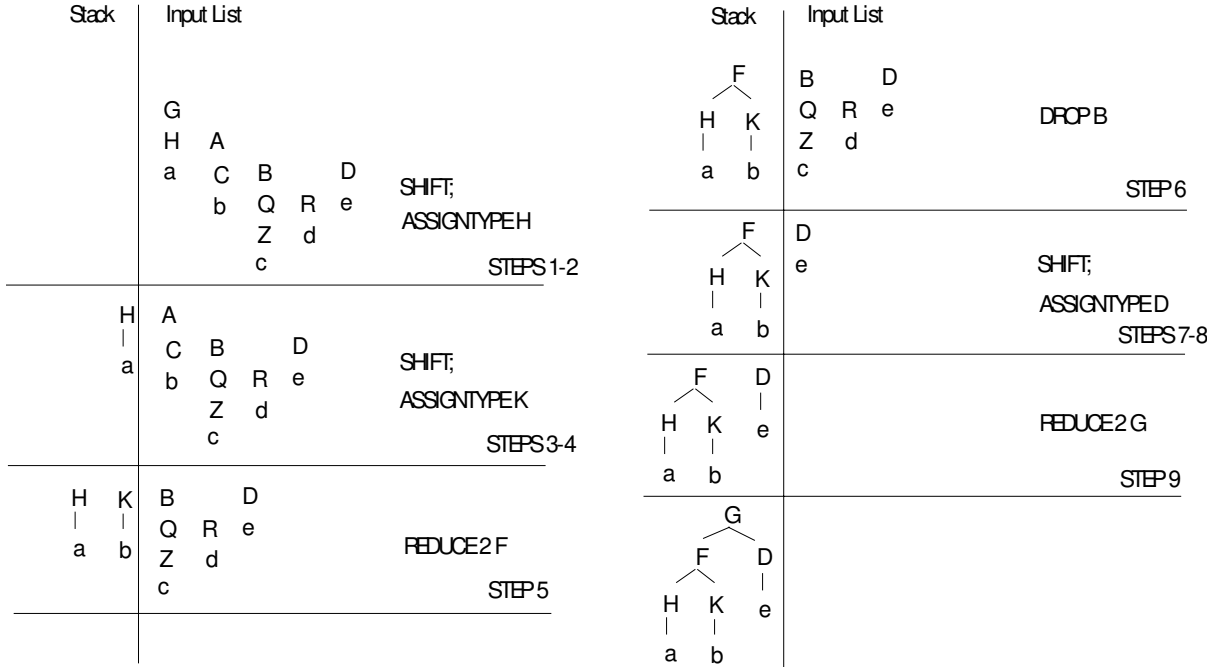


Figure 3: **Incremental Tree Compression.** An example of the tree reduction algorithm using letters as place holders for subtrees (such as verb and noun phrases). The parse tree is reduced into a smaller tree, which is a compressed version of the original.

corpus for analysis (Knight and Marcu, 2002), which is a corpus of newspaper articles about computer products automatically paired with human abstracts. These tree methods were eventually reframed as graph reduction problems as the syntactic representation became richer. One such graph based model incorporated expanding-constituent parse trees (Li et al., 2014) who use ILP with features consisting of tree node data such as distance to other nodes and depth.

Later, hybrid approaches across fields of research enhanced parse tree based methods and became more common. For example, models for multi-document summarization (MDS) also

began to leverage graphs via methods of extractive selection and reduction as a global optimization task (Ranjitha and Kallimani, 2017). Another multi-document summarization (MDS) solution included spectral methods by embedding eigenvector decompositions of their similarity matrices based on TF/IDF, and then clustering for sentence selection (Gupta et al., 2019).

Eventually, a shift from parse trees to graphs became dominant using various editing and reduction strategies. This work involves AMR graph reduction techniques seen in previous work for both single-document summarization (SDS) (Liu et al., 2015) and MDS (Liao et al., 2018). However, the focus shifted to NN methods after breakthrough performance gains were discovered for upstream NLP tasks in part of speech (POS) and named entity recognition (NER) tagging. The work of Rush et al. (2015) made use of a neural machine translation (NMT) language model’s attention layer to generate abstractive summarizations. See Section 2.1.4 for more related work on AMR summarization.

2.1.2 Neural Network Summarization Methods

The trend toward deep learning (DL) models continued given the promising results of Rush et al. (2015). In this work, the authors increase readability by producing text that transitions between topics using sequence to sequence neural models. This is one of the first efforts using NNs for abstractive summarization. Nallapati et al. (2016) built on this work by adding an attention encoder-decoder to a recurrent neural network (RNN) model.

Until this point, only the source document was sampled and used to create the summary. This changed when word embeddings were used to supplement out-of-vocabulary words (See et al., 2017), or “generate” them from the embedding vocabulary. An indicator specifying whether

to use the source text or the vocabulary was jointly learned with a probability distribution over the source text, called the “pointer”, to the next word in the generated summary.

However, these models suffer by producing factually incorrect details or information missing from source (Maynez et al., 2020), and ultimately fail to generate human-like summaries (Wiseman et al., 2017). Cao et al. (2018) brought attention to the lack of faithfulness in NN models and used sequence-to-sequence models experimenting with several network topologies, such as GRUs with dual attention, in an attempt to ameliorate these issues. Nonetheless, like most summarization work, their success was predicated on the reliability and usefulness of ROUGE, which is a measure that assesses the quality of generated summaries¹ (Lin and Hovy, 2003; Lin, 2004).

2.1.3 Summarization Scoring Methods

Maynez et al. (2020) conclude that these scoring metrics, such as ROUGE, do not sufficiently evaluate the quality of the generated summaries since they do not correlate well to summaries that are faithful given their high degree of hallucination (erroneous and nonfactual text automatically generated by a model). They group the more contemporary BERTSCORE² (Zhang et al., 2020a), with ROUGE in the way the metric provides some measure of information, but

¹This quality is assessed by comparing to summaries generated by humans using overlapping n -grams counts or word pairs.

²This quality is assessed by comparing reference and generated tokens’ similarity based on the contextual embeddings of a BERT family model.

still fails to assess quality. The work’s qualitative analysis is detailed and successfully argues the issues of scoring summaries, but provides no recommendation on how to address them.

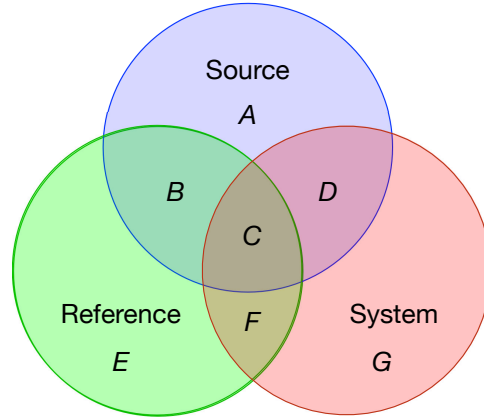


Figure 4: **Summarized Text Overlap.** Shared information between source text, the reference summary, and the system generated summary.

Only the recent work of Shing et al. (2021) has attempted to address these scoring procedures by introducing two new measures, faithfulness-adjusted and hallucination rate, and then used them to evaluate automatically generated discharge summaries. Faithfulness-adjusted precision is the amount of relevant and faithful information in the discharge summary, and is defined as the quotient of the intersection of content in the source text, reference summary and system summary divided by the content of the system (see Figure 4), or more precisely $\frac{C}{G}$. Faithfulness-adjusted recall is the relevant and faithful information found in the discharge summary, and

defined as $\frac{C}{B+C}$. Given that discharge summaries often include information not found in any clinical note (Shing et al., 2021; Adams et al., 2021; Alsentzer and Kim, 2018), these metrics provide a much better summary qualitative assessment.

2.1.4 Summarization using Graphs and Abstract Meaning Representation

AMR (abstract meaning representation), which is a semantic representation language that describes the abstract meaning of a sentence as an acyclic graph, has seen a decade of interest across many tasks (Liu et al., 2015; Bonial et al., 2020a; Naseem et al., 2022). AMR captures “who is doing what to whom” in a sentence (Banarescu et al., 2013) and can also be represented in the context free notation of Penman (Kasper, 1989). AMR (abstract meaning representation) graphs were later enriched with PropBank frames, which greatly enhanced their expressiveness (Palmer et al., 2005). Recent achievements that use AMR models as the primary data representation include work in natural language generation (NLG) (Manning et al., 2020; Gu et al., 2020), automatic machine translation (Blloshmi et al., 2020), question/answer systems (Lim et al., 2020), and building logical forms (Galitsky, 2020).

AMR models have also shown promise for automatic summarization with recent work reminiscent of the early related-map graph reductions of Salton et al. (1994) in the way source component nodes and edges are deleted. Later, Liu et al. (2015) used the same graph reduction methods with AMR graphs. In this work, the authors created a fully connected graph and modeled edge inclusion with ILP, which was then used to heuristically generate abstractive text for SDS. This was later broadened by with a more comprehensive and robust AMR graph based realization algorithm for MDS (Liao et al., 2018).

My thesis is inspired by the work of Liu et al. (2015) and Liao et al. (2018) with regard to AMR graph reduction methods. However, my method is different in how it uses flow networks to induce or predict a new graph, whereas their work builds on the graph reduction methods of Thadani and McKeown (2013) For sentence comprehension, which re-frames the concept of commodity flow (Magnanti and Wolsey, 1995) as indicator flow constraints for edge inclusion.

Furthermore, Liao et al. (2018) used extractive methods for sentence selection where my work has no extraction phase. Their extractive method selects the top ranked sentences using spectral clustering on an eigen vector decomposition (Blei et al., 2003), which are then abstractive summarized using AMR. Unsupervised clustering models still dominate the state of the art (SoTA) to subset the source text (Gupta et al., 2019; Jin and Wan, 2020) for abstractive summarization. Many of these models cluster using contextual sentence embeddings (Miller, 2019; Chaturvedi et al., 2020). While the utility of these models is potentially successful in choosing a salient topic helpful to a summary, they implicitly are not faithful as they do not learn those same topics given in the reference summary with any rigor or presentation order. To my understanding, mine is one of the first works that addresses this issue of faithful summarization by using a flow network connected from the source component to the summary component.

Heterogeneous graphs exploiting the relationships of language part granularities, such as word, sentence and topic, have been used in document summarization (Wei, 2012). During the early 2010s sentence compression was still a popular method for summarization using dependency tree graph and semantic role labeling (SRL) (Thadani and McKeown, 2013; Ranjitha

and Kallimani, 2017). A concept graph constructed using AMR and refined with the PageRank algorithm was used as an abstractive method. The SimpleNLG library (Gatt and Reiter, 2009) was then used for NLG for summarization (Vilca and Cabezudo, 2017). Recent work, once again, uses heterogeneous graphs as graph neural network (GNN) for extractive summarization using AMR (Wang et al., 2020).

2.2 Relevant Clinical Applications and Summarization

Much in the medical domain has been accomplished from visually summarizing clinical patient data (Powsner and Tufte, 1997) to hybrid neural methods (Adams et al., 2021) with multi-level RNNs for extraction (Zhou et al., 2018) and BART (Lewis et al., 2020) to abstractively condense and smooth the summaries. My work is similar to comprehensive clinical summarization systems, such as BabyTalk, as a comprehensive end-to-end pipeline (Reiter, 2007; Portet et al., 2009). Most of this work focuses on patient summarizations or only partial discharge summary generation. However, my work accomplishes a multi-section automatic summarization of the discharge summary for clinical physicians.

2.2.1 Abstract Meaning Representation in Medical Research

AMR subgraphs representing events for modeling molecular events and interactions were learned with long-short term memory (LSTM) neural networks (Rao et al., 2017). Drug interactions were also modeled using AMR while utilizing word and PropBank embeddings in the graph (Wang et al., 2017). My work is similar in that it also uses PropBank, but theirs does not use role embeddings in the same way or edge embeddings at all. In biomedical research literature, AMR was used to model the meaning of a question used in a semantic search for providing

answers related to ultraviolet light in COVID-19 (Bonial et al., 2020b). Biomedical research literature information extraction was also used to narrow large texts to scientific entities by enriching an AMR graph with a biomedical and chemical interaction knowledge graph (Zhang et al., 2021).

2.2.2 Medical Note Section Identification

Most work in section identification (SI) has been applied in academic papers and clinical medical notes. Sectioning MedLINE abstracts was explored by McKnight and Srinivasan (2003) using a support vector machine (SVM). This classifier was used to label sentences as *Introduction*, *Method*, *Result*, or *Conclusion* and showed promising results using a bag-of-words approach. Sequence based approaches (Hirohata et al., 2008) were also used to section scientific abstracts into *Objective*, *Methods*, *Results*, and *Conclusion* labels using a conditional random field (CRF) model producing a sentence level accuracy of 95.5%.

While academic abstract segmentation was a well explored area (Hirohata et al., 2008), Tepper et al. (2012) were the first to apply statistical methods to the medical domain to automatically classify sections of clinical free text into sections. Their method used in, out, begin (IOB) annotation (Ramshaw and Marcus, 1995) with labels to mark named sections. For example, B-HPI indicates a beginning token for the history of present illness (HPI) section. Their dataset consisted of annotating the 2010 i2b2 corpus with a section header and medical ontology label, and obtained an F-measure of 0.92 for the concept extraction task (Uzuner et al., 2011; de Bruijn et al., 2010). A Maximum Entropy (MaxEnt) model (Berger et al., 1996) and beam search were used for classification to produce the IOB sequence for token tagging.

Along with MaxEnt, other non-neural network methods, such as SVM and CRF models continue to be popular with few exceptions as detailed in the comprehensive survey of Pomares-Quimbaya et al. (2019). One such exception (Sadoughi et al., 2018) used a LSTM model with word-to-vector (word2vec) embeddings (Mikolov et al., 2013a; Mikolov et al., 2013b) for a binary classification of section boundaries. Even though the corpus consists of dictated and transcribed notes, they show that neural methods work for the section segmentation task. Other notable NN text segmentation works use convolutional neural networks (CNNs) over sentence embeddings with a softmax over the output of a bi-directional long-short term memory (BiLSTM) layer to demarcate sections as a binary classification across both medical and non-medical datasets (Badjatiya et al., 2018). Barrow et al. (2020) also used a LSTM in a network that aggregates features across fastText word embeddings using a concatenated segment pooling LSTM (S-LSTM) for non-medical Wikipedia articles (Bojanowski et al., 2017).

2.2.3 Provenance of Data

Clinical notes are riddled with redundant information in the form of copied and pasted text and embedded structured data as templates from large electronic health record (EHR) systems. Because this has been a recognized issue for decades, efforts such as redundant text detection in same-category EHR notes using latent Dirichlet allocation (Blei et al., 2003) has been explored (Cohen et al., 2013). Duplication in medical EHR notes was studied in recent work to find root causes of copied text using 10-gram token spans, which found that 58% of a physician note is copied from previous notes (Steinkamp et al., 2022).

Provenance of data as it relates to the task of summarization has included multi-level RNNs used for extractive summarization (Zhou et al., 2018). Soon after, BART (Lewis et al., 2020) was used to abtractively condense and smooth the summaries. A large corpus was analyzed and a hybrid extractive/abstractive neural method was used to summarize the *Brief Hospital Course* section (Adams et al., 2021).

2.2.4 Medical and Discharge Note Summarization

The literature is rich with examples of medical note summarization that include both longitudinal (Hirsch et al., 2015), and non-longitudinal (Pivovarov and Elhadad, 2015) note types, which are just two examples of mutual discipline interest. Furthermore, the shared understanding, agreement, and acknowledgment that faithful summarization is necessary, but lacking, has been thoroughly reviewed (Zhang et al., 2020b). Another example highlighting the interest in summarization across the clinical and NLP communities is exemplified by the effort of corpora publishing such as the Evidence Based Medicine Summarisation corpus (Molla and Santiago-Martinez, 2011) and large text datasets such as MIMIC-III (Johnson et al., 2016).

Doctor-patient SOAP (Subjective, Objective, Assessment and Plan) note generation (a type of physician note) is another example of areas of mutual interest as a NLP work in the medical domain (Krishna et al., 2021). More recently, cross discipline interest has become collaboration (Acharya et al., 2018; Boyd et al., 2018). This thesis includes the multi-discipline research with medical academics, and most closely resembles the discharge summarization work of Adams et al. (2021).

Adams et al. (2021) showed promising results in summarizing the *Brief Hospital Course* section. However, for single section summarization, the summarization of Medical Information Mart for Intensive Care III (MIMIC-III) physician notes is perhaps the most interesting comparison and potentially most impactful (Gao et al., 2022a). Clinical notes were summarized by fine-tuning the T5 (Raffel et al., 2020) and BART (Lewis et al., 2020) state of the art seq2seq models and evaluated using the BERTSCORE (Zhang et al., 2020a) and ROUGE (Lin, 2004) scoring methods. In addition, co-occurring extracted Unified Medical Language System (UMLS) medical concepts (Bodenreider, 2004) were also reported. Bodenreider concludes with their intention of providing a foundation for future summarization work and acknowledgment of the challenges facing the clinical note summarization, which echoes the motivation of this work.

2.3 Graph Alignment

Flow networks and graph alignment prior work is presented in this section. While there is no known work using a flow network with AMR graphs (as noted in Section 1.2.2.2), they have been used to align bipartite graphs for the matching problem for decades.

2.3.1 Flow Networks

A flow network is a graph that models material through a network with two weighted values for each edge: a capacity and a flow. Both are non-negative values and the flow is constrained by the capacity. The study and optimization of a graph as a flow network has been evolving for more than eight decades since the publication of a 1930 article by A. N. Tolstoï on Soviet railroad planning (Schrijver, 2002). Several decades later, the max flow problem was

formalized in a declassified military report (Harris and Ross, 1955), which became the foundation of the seminal Ford-Fulkerson max-cut min-flow algorithm (Ford and Fulkerson, 1962). Figure 5 gives an example of the application to the Soviet railway system. Since, numerous applications of the max-cut min-flow algorithm have been devised including maximum cardinality matching (Hopcroft and Karp, 1973; Irving et al., 1987), baseball elimination (Schwartz, 1966; Kleinberg and Tardos, 2005), airline scheduling (Cormen et al., 2001; Bazaraa et al., 2010); the algorithm has been addressed in the quantum computing field (Krauss et al., 2020).

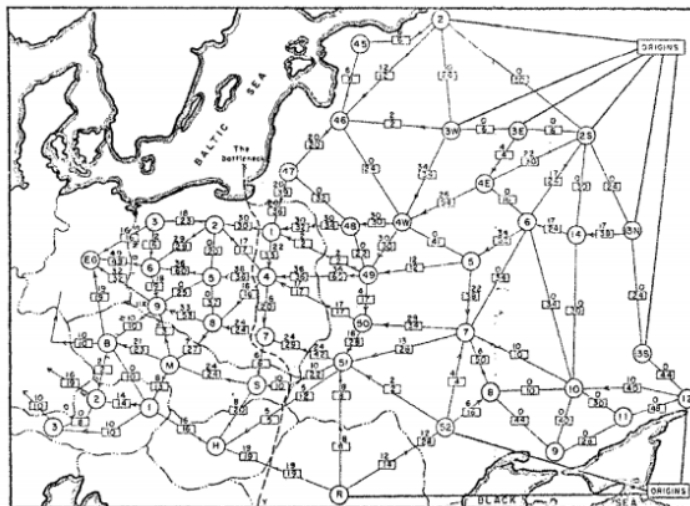


Figure 5: **Soviet Rail Network, 1955.** A diagram of the railway network of the Western Soviet Union and Eastern European countries with a description of the cut “bottleneck” (Harris and Ross, 1955; Ford and Fulkerson, 1962).

The max flow algorithm has seen many optimization improvements over the years, with some specific to graph topologies such as bipartite graphs, that run in $\tilde{O}(m+n^{1.5})$ time on an n -vertex, m -edge graph. Other more general algorithms include the push-relabel algorithm (Goldberg and Tarjan, 1988), which has a running time of $O(nm \log \frac{n^2}{m})$. The fastest algorithm runs in almost linear time $m^{1+o(1)}$ in m edges (Chen et al., 2022).

Flow networks have also been used in NLP for the task of MT framed as a bipartite flow matching task (Gaussier, 1998), and for information retrieval (IR) to semantically match documents to queries (Guo et al., 2016). However, they have not been used with any language based graphical model, AMR or knowledge graphs to my knowledge.

2.3.2 Alignment

This section covers graph alignment in how it relates to the NLP field. A semantic web is an ontology graph with edges that describe relationships between nodes that describe concepts. Semantic web graph alignment has an extensive history of work that includes ontology matching (Fossati et al., 2006) and data aggregation (Faria et al., 2013). There has been work using a flow network to align a semantic web graph for data discovery (Heidary et al., 2010), which is very similar to the general summarization case for this work. However, no node or edge embeddings were used, nor were there any aspects of AMR utilized.

In many cases, a semantic web ontology is domain specific while the alignment algorithm is not. However, domain specific graph alignment is prevalent in the biomedical area. Protein-protein interaction networks have been automatically sequenced and aligned (Kuchaiev et al., 2010; Malod-Dognin and Pržulj, 2015). Genome graphs have been aligned with minimal

edit distance paths using an A-star-like algorithm (Ivanov et al., 2020), and have been a focus of algorithmic optimization (Rautiainen and Marschall, 2020).

As in many domains such as computer vision, graph bipartite alignment transitioned from non-neural methods that employ graph algorithms (Koutra et al., 2013; Zhang and Tong, 2016; Ni et al., 2018) to predominantly neural methods (Goyal and Ferrara, 2018; Trung et al., 2020) that combine graph algorithms with DL. These neural methods typically involve creating node and edge embeddings based on features and graph topology (Chu et al., 2019; Pei et al., 2022).

In addition to specific applications, alignment has been used as a measure of similarity. The most common metric used for such tasks on text-to-graph models use SMATCH (Cai and Knight, 2013), which measures the likeness of two AMR graphs. More recent AMR similarity measures involve local network neighborhood matching using distributed probability distributions (Opitz et al., 2021).

2.4 Frameworks

An additional contribution of my thesis is a NLP deep learning framework, DeepZensols (Landes et al., 2023). A framework is a set of software advanced programming interfaces (APIs) used to facilitate complex processing pipelines common in NLP tasks. These frameworks are typically specialized for a specific purpose or domain and typically implement behavior that is used often and in many places in the pipeline. A well written framework allows the researcher to concentrate on the task at hand rather than rewriting common software modules pervasive in NLP pipelines.

Popular DL frameworks such as TensorFlow¹ have a dashboard that provides metrics, such as training and validation loss. However, these general purpose frameworks offer basic performance metrics and do not provide a means of producing higher abstraction level NLP specific models. More specifically, frameworks such as Keras, supply a very coarse API allowing solely for cookie-cutter models. They lack the ability to easily create and evaluate models past this surface interface.

Frameworks such as PyTorch², which are more common in academia, provide a more straightforward simple API that is similar to the core TensorFlow libraries, and thus have the same shortcomings as a tool to bridge the gap between pure research and reproducibility.

AllenNLP (Gardner et al., 2018) is a flexible configuration driven framework that provides construction of NLP NN architectures and is the closest framework to mine. However, it does not have fast feature swapping (see Section 7.5) and batch creation capability, and lacks most of the components necessary to consistently reproduce results.

Popular packages providing support for transformer architectures such as BERT (Devlin et al., 2019) include HuggingFace³. However, this framework only provides transformer models for contextual word embeddings. Compared to the HuggingFace API DeepZensols provides

¹<https://www.tensorflow.org>

²<https://pytorch.org>

³<https://huggingface.co>

comprehensive access to not only contextual embeddings, but many of the non-contextual embeddings still used such as fastText (Bojanowski et al., 2017).

CHAPTER 3

MEDICAL SECTION IDENTIFICATION

(This chapter expands on the paper “*A New Public Corpus for Clinical Section Identification: MedSecId*” by Landes et al. (2022) in the Proceedings of the 29th International Conference on Computational Linguistics.)

The alignment methods for AMR (abstract meaning representation) graphs I presented in Chapter 5 are designed for sectioned medical documents for the task of summarization. One important and innovative application is discharge summary generation. We turn now to the work on section identification (SI), which grounds generation of discharge summaries. This chapter introduces MedSecId, a contribution of this thesis and a publicly available set of 2,002 fully annotated medical notes from the Medical Information Mart for Intensive Care III (MIMIC-III) corpus.

3.1 Motivation

While discharge summary sectioning helps a physician locate specific information, the primary impetus for the structure and content stems from the ongoing dispute between providers and healthcare insurance companies in the United States. Providers are limited by how much they can bill for relatively simple medical procedures, but increasingly complex procedures garner more revenue with proper documentation. Specifically, medical billing staff and insurance companies use relative value units (RVUs), which is a monetary unit updated annually and

currently set at \$34.30. The number of RVUs billed is based on the composition and number of sections included in the medical notes per guidelines set by the Centers for Medicare and Medicaid Services¹.

For this reason, providers are encouraged to write medical notes to maximize RVUs out of necessity (Barnes et al., 2008) even though physician training lacks such emphasis. In contrast, medical residents and students are evaluated with the objective structured clinical examination (OSCE), which is a student examination that evaluates students based on direct observation (Zayyan, 2011). However, the exam’s evaluation with respect to medical note authoring and structure uses very different criteria and omits RVUs (Gallagher et al., 2020). The necessity of a particular structure in medical notes, for the purpose of patient care and unfortunately more important insurance billing requirements, highlights the need for understanding sectioning.

Most unstructured medical text found in electronic health records (EHRs) written by medical staff have conceptually well defined sections. For example, the discharge summaries consists of named sections, typically in a specific sequence, such as the *History of Present Illness*; this type of section appears both in discharge summaries and in physician notes that describe a chronology of an illness that begins with the admission of the patient.

Since each section contains specific information, SI is often the first step in a medical natural language processing (NLP) pipeline and can lead to downstream propagation errors causing

¹<https://www.cms.gov/Regulations-and-Guidance>

poor task specific results if not properly executed. Examples of downstream tasks that benefit from SI include medical summarization, entity linking and natural language understanding and extraction.

Even though sections usually start with header text, SI is more challenging than simply parsing the first several leading header tokens of the respective section (underlined in Figure 6). While the first several tokens can be helpful in identifying a section, their naming often varies. For example, the 6th section in Figure 6 starts with header tokens *Preoperative Laboratory Data*, but the section type is `labs-imaging`. There are also cases where the header tokens are missing, as shown in the 7th section (`hospital-course`) in the same figure, or where sections have several header text spans placed throughout the section. Adding to this challenge is the non-uniformity of the text, lack of section boundary syntax, and copy-pasted text from other notes or from structured data such as patient vital signs.

Summarization methods for structured text documentation has been accomplished by processing each section separately (Adams et al., 2021). The highly structured nature of medical notes can be exploited by splitting sections in to groups as document nodes in the source component graph to create the corresponding summary component, thereby learning structure of the discharge summaries (see Section 5.6.1). Since the MIMIC-III corpus was used for summarization in this thesis, sectioning was done using the MedSecId annotations (Landes et al., 2022). A pretrained model trained on the MedSecId annotations was used for inferencing sections for notes that have not been annotated.

Admission Date: [**2126-2-7**] Discharge Date: [**2126-2-20**] Date of Birth: [**2069-4-1**] Sex: M
history-of-present-illness HISTORY OF PRESENT ILLNESS: Mr. [**Known lastname **] is a 56-year-old male who experienced chest. . .
past-medical-history PAST MEDICAL HISTORY: Hypertension, former smoker with a 4- pack per day history for which he. . .
social-history SOCIAL HISTORY: He lives alone, and he works at [**Hospital3 2576**] as a cargo transporter.
medication-history MEDICATIONS ON ADMISSION: Aspirin 325 mg p.o. once a day, Toprol-XL 50 mg p.o. once a day.
allergies ALLERGIES: He had no known drug allergies.
labs-imaging PREOPERATIVE LABORATORY DATA: White count 6.0, hematocrit 33.3, platelet count 329,000. . .
He was discharged to home with VNA services in good condition on [**2126-2-20**]...

Figure 6: **Discharge Summary.** (Repeated from Section 1.1.) A MIMIC-III discharge summary note with the section type in **bold**; omitted text is indicated with ellipses.

3.2 Dataset

MedSecId is a subset of the MIMIC-III Version 1.4 (Johnson et al., 2016) corpus (Johnson et al., 2016) that we annotated; MIMIC-III is publicly available and consists of critical care unit EHR records from the Beth Israel Deaconess Medical Center in Boston, Massachusetts. The dataset contains 58,976 hospital admissions across 46,520 patients who were admitted to the intensive care unit (ICU) surgical, medical, and neonatal departments. It includes 2,083,180

unstructured medical text notes handwritten by medical professionals across several disciplines and contains 15 categories, such as discharge summaries and radiology notes.

We created a curated annotation set consisting of text spans taken from a random sample across five categories of MIMIC-III medical notes¹, including discharge summaries, *Radiology*, *Echo*, *Physician*, and *Consult* progress notes (see Table I). Radiology clinical notes provide the diagnosis of images such as X-rays; similarly Echo notes provide diagnoses based on an ultrasound of the heart. Physician notes have the updated progress of a patient and are written on a daily basis, whereas Consult notes are written by an outside consultant physician. Each text span annotation contains the type of the section in MedSecId, such as *History of Present Illness*, with zero-index character offsets of where the span starts and ends in the note.

Category	Count	Proportion
Discharge summary	1,254	62.64%
Physician	288	14.39%
Radiology	205	10.24%
Echo	198	9.89%
Consult	57	2.85%
Total	2,002	100%

TABLE I: **Annotated Medical Note Categories.** Annotated medical notes by category and their distribution in the annotation set.

¹The unstructured medical note data was taken from the NOTEEVENTS table.

While each section contains a single type, sections have zero or more overlapping header text spans (see Figure 6). In most cases, there is a single header span, but vital signs sections can “float” without a physical exam header. These header spans consist of text that identify the section such as *History Of Present Illness*, an alternate spelling or abbreviation such as *HPI*. Even though single header spans usually appear at the beginning of a section, additional section headers are found later in the body indicating subsections in some cases. Since section type inclusion highly varies based on the patient’s age, notes were annotated with an age type (adult, pediatric or neonatal patient), based on the content of the note by the annotator.

3.2.1 Annotation Process

Our annotation process consisted of several preliminary rounds of annotation, that led to our final annotation guidelines and final annotation.

Before annotation began, a custom set of regular expressions were used to pre-annotate, similar to previous work (Shivade et al., 2015); mine were medical note specific and captured header tokens along with the section spans. The application of the regular expressions was only a means to reduce the work of the annotators, who followed the annotation guide regardless of any rule based pre-annotations. The initial rule based automatic annotation process was amended by the work of Alsentzer et al. (2018), who generously shared their *History of Present Illness* annotations to better identify and segment the initial dataset used by my collaborator annotators. These automatic annotations were edited by the annotators after they were imported into INCEpTION (Klie et al., 2018) and saved to later compute an inter-coder agreement between the physicians and rule-based output (see Section 3.2.2).

An attending physician (designated as a primary annotator) co-wrote a preliminary annotation guide with input from a secondary physician annotator. These two annotators engaged in a process of annotation, discussion and revision of the guidelines: they annotated a first set of one hundred notes, revised the guidelines, annotated a second set of one hundred notes, and finalized the guidelines after this second round.

Here we summarize the issues that the annotators faced during these preliminary rounds of annotation. This process was useful for the physicians to reach a consensus on what sections should be annotated and agreed on section types given their experience writing such notes themselves. A set of sections and their relation to notes began to coalesce during this process, which provided the motivation to create an ontology for the purpose of a meta documentation about the annotations and the utilitarian purpose to assist in annotation by importing it as a “knowledge base” in INCEpTION. The ontology consisted of a one-to-many mapping from notes to 50 section types using each section’s header tokens captured by the regular expressions by string massaging. For example, *History of Present Illness* became `history-of-present-illness`. Among the categories, 29 sections were shared across more than one note, such as *History of Present Illness* shared between notes *Discharge summary*, *Consult*, and *Physician* (see Table II for annotated sections).

Each section type was then agreed on by the physicians with many re-typed and regrouped. For example, *Echo* notes contained internal subsections for each chamber of the heart; this was resolved by grouping the entire section as *Findings* to match section types in *Radiology* notes. Other subsections implicitly resulted by physicians copying radiology findings in discharge sum-

Type	Tokens	Spans	Notes
physical-examination	203K (8%)	1,385 (6%)	Consult, Physician
history-of-present-illness	239K (9%)	1,348 (6%)	Consult, Discharge summary, Physician
allergies	9,221 (0%)	1,205 (5%)	Consult, Discharge summary, Physician
hospital-course	692K (26%)	1,165 (5%)	Discharge summary
labs-imaging	416K (16%)	1,155 (5%)	Consult, Discharge summary, Physician
past-medical-history	60K (2%)	1,141 (5%)	Consult, Discharge summary, Physician
discharge-condition	14K (1%)	1,132 (5%)	Discharge summary
discharge-instructions	183K (7%)	1,077 (5%)	Discharge summary
discharge-diagnosis	34K (1%)	1,040 (5%)	Discharge summary
chief-complaint	9,622 (0%)	996 (4%)	Consult, Discharge summary, Physician
discharge-medications	196K (7%)	914 (4%)	Discharge summary
social-history	28K (1%)	912 (4%)	Consult, Discharge summary, Physician
medication-history	49K (2%)	867 (4%)	Consult, Discharge summary, Physician
family-history	11K (0%)	802 (4%)	Consult, Discharge summary, Physician
discharge-disposition	5,602 (0%)	754 (3%)	Discharge summary
major-surg-or-inv-proc	16K (1%)	704 (3%)	Discharge summary
facility	2,668 (0%)	502 (2%)	Discharge summary
reason	5,588 (0%)	458 (2%)	Consult, Radiology
findings	58K (2%)	395 (2%)	Echo, Radiology
assessment-and-plan	131K (5%)	381 (2%)	Consult, Physician
review-of-systems	7,422 (0%)	329 (1%)	Consult, Discharge summary, Physician
image-type	1,820 (0%)	328 (1%)	Radiology
last-dose-of-antibiotics	3,689 (0%)	293 (1%)	Consult, Physician
24-hour-events	16K (1%)	250 (1%)	Physician
code-status	1,879 (0%)	237 (1%)	Physician
impression	8,233 (0%)	224 (1%)	Echo, Radiology
disposition	1,161 (0%)	210 (1%)	Physician
conclusions	28K (1%)	206 (1%)	Echo
communication	1,304 (0%)	199 (1%)	Physician
patient-test-information	13K (1%)	198 (1%)	Echo

TABLE II: **Most Frequently Annotated Medical Sections.** The top 30 most frequently annotated sections.

maries. In an effort to reduce complexity, a flat note-to-section hierarchy without creating a second section level was kept. In some cases this was achieved by combining laboratory results data with radiology findings/diagnosis as a single section by simply re-casting *Labs* to *Labs/Ra-*

diology for sections that included imaging studies. Other sections needed to be combined as not all notes had a clean separation.

To accommodate for a significant variation in how physicians labeled sections in these situations, *Labs* and *Radiology* was combined into a *Labs/Radiology* section. *Labs* and *Imaging* were also combined into *Labs/Imaging*. Since discharge summaries typically incorporate instructions for the patient and follow up information, we categorized these together broadly as *Discharge instructions*. The MIMIC-III pseudo tokens, such as **[**First Name**]** were not annotated unless they were included in the body of the section.

The primary and secondary annotators finished revising the annotation guidelines and then trained the third annotator. A subset of 80 medical notes, chosen from the second batch of 100 that the primary and secondary annotator had annotated and discussed, was used to train the third annotator. Because these first two batches were only used for creating guidelines and training, they were not added in the final annotation set. During this process, the well known Krippendorff’s α coefficient (Krippendorff, 2011), was used to compute inter-annotator agreement (IAA) between this last annotator and the other two, until α became higher than 0.8.

3.2.2 Final Annotation and IAA Computation

Once the guidelines were finalized the final annotation process started. A set of 100 notes (different than the sets discussed in Section 3.2.1) was held out to compute the inter-annotator agreement (IAA) on the final guidelines. The remaining 1,902 notes were divided up among the three annotators, as customary.

Inter-annotator agreement was calculated on the 100 held out notes as exact section character offsets and section types—both the offsets and the section type had to match to be considered correct. This agreement was calculated among the human annotators, and subsequently between each annotator and the regular expressions that were initially used to segment the notes.

Among humans, Krippendorff’s α yielded more than acceptable values of 0.84 to 0.87 on the final set held out for this IAA calculation (see Table III). At this point, these annotations were added to the final dataset by selecting notes with the fewest issues¹ using the primary annotator as the tie-breaker.

	A1	A2	A3	R
A1	1.0	0.81	0.87	0.73
A2		1.0	0.84	0.49
A3			1.0	0.53
R				1.0

TABLE III: **MedSecId Inter-annotator Agreement.** Krippendorff’s α coefficient of inter-annotator agreement between the annotators and the regular expressions. **A1** is the primary annotator, **A2** is the secondary annotator and **A3** is the third annotator, and **R** represents the regular expressions.

¹Annotation issues included placement of header tokens and missing sections. For example, an annotation with a defined section would win over another’s annotation having a missing section header annotation.

While we achieved a high inter-coder agreement among human annotators, we found troubling data in terms of the performance of the regular expression annotation approach. We computed an aggregate Krippendorff’s $\alpha=0.54$ between the human physician annotators and my custom regular expressions (see Section 3.2.1) on the final annotated data, which falls more than 14 points shy of the “lowest conceivable limit” of 0.68 (Krippendorff, 2004). This shows how regular expression’s performance to segment notes falls short of that by human annotators (see Table III), yet regular expressions continue to be the most common methods used for section identification (Pomares-Quimbaya et al., 2019; Shivade et al., 2015).

In part, the regular expressions often failed to demarcate the entire section, especially in text with irregular formatting toward the end. Furthermore, additional analysis shows the α scores between individual annotators and the regular expressions are low as well, albeit with a fairly high variance. Krippendorff suggests that acceptable scores that are “customary to require” have $\alpha > 0.8$ (Krippendorff, 2004). On one hand, an α of 0.73 between physician *A1* and the automatic regular expression annotator *R* clears the minimal limit threshold. However, this metric falls well below the “aimed” score of 0.8. The larger issue is with physician *A2*’s and *A3*’s scores of 0.49 and 0.53, which fall short of the minimum limit by a large margin. From these scores (see Table III) and the low overall α , I conclude regular expressions do not sufficiently segment medical notes, therefore the annotation set I provide should be considered the gold standard for medical note identification and segmentation.

3.2.3 Data Analysis

An interesting discovery concerned projections of medical conditions across sections in embedded space. Concept unique identifiers (CUIs), which are entity links to Unified Medical Language System (a large medical concept ontology) were extracted (Bodenreider, 2004). The CUIs were linked using MedCAT (Kraljevic et al., 2021) and weighted by TF/IDF (Sparck Jones, 1972) across sections. Each CUI was mapped to a vector from *cui2vec*¹, and then reduced to three dimensions using principal component analysis (PCA), shown in Figure 7. The plot was generated without normalizing or standardizing the data so CUI vector magnitudes were retained for analysis. Figure 7 (a) shows the **past-medical-history** section (purple) CUIs on the horizontal axis with **past-surgical-history** (blue) CUIs only on the vertical axis with size proportional to TF/IDF.

The past surgical and medical history sections in discharge summary notes project many medical disease CUIs as orthogonal to surgical CUIs. The medical disease CUIs on the vertical axis are those that do not have surgery as a treatment option, such as hypertension. However, a CUI representing coronary artery disease that plots along the surgical history vertical axis does require surgery. Most of the data points that share the vertical axis along with **past-medical-history** are those that require both medication and surgery, such as cancer.

Not only does this show *cui2vec* being used in practice for the first time, as far as I know, it illustrates an application of how groupings of concepts can be visualized and analyzed to gain

¹A unique concept identifier embedding trained from biomedical text.

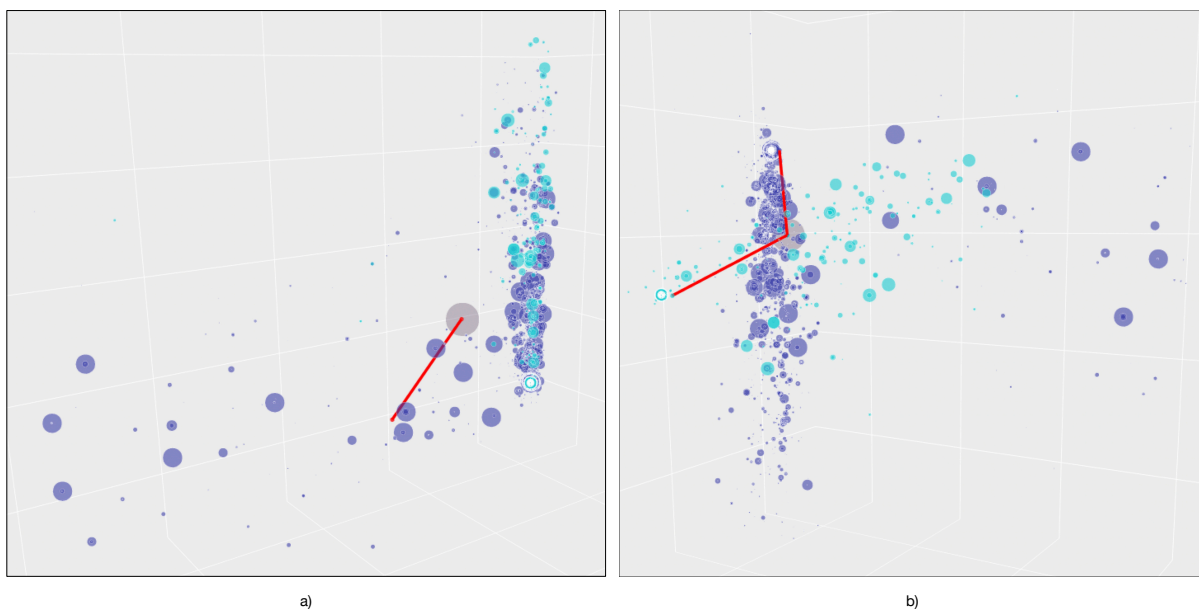


Figure 7: **Concept Unique Identifier Plots.** a) Plot of **past-medical-history** (purple) and **past-surgical-history** (blue) reduced to 3D together as one data set with the first principal component (red line) with data point size as the TF/IDF score, b) Plot of the same sections but reduced to 3D as separate data points with respective first principal components.

intuition and insight in complex medical data. In my data, this includes not only a semantic relationship between concepts, but how those concepts represent the treatments involved based on the section from which they originate. Given this data relationship, I hypothesize that utilization of *cui2vec* embeddings, such as concatenating them to word vectors, will increase performance of task specific models including SI.

3.3 Limitations

MedSecId is limited to notes of patients admitted to an ICU from the MIMIC-III corpus for several note categories with no data included after the patient leaves to a lower severity

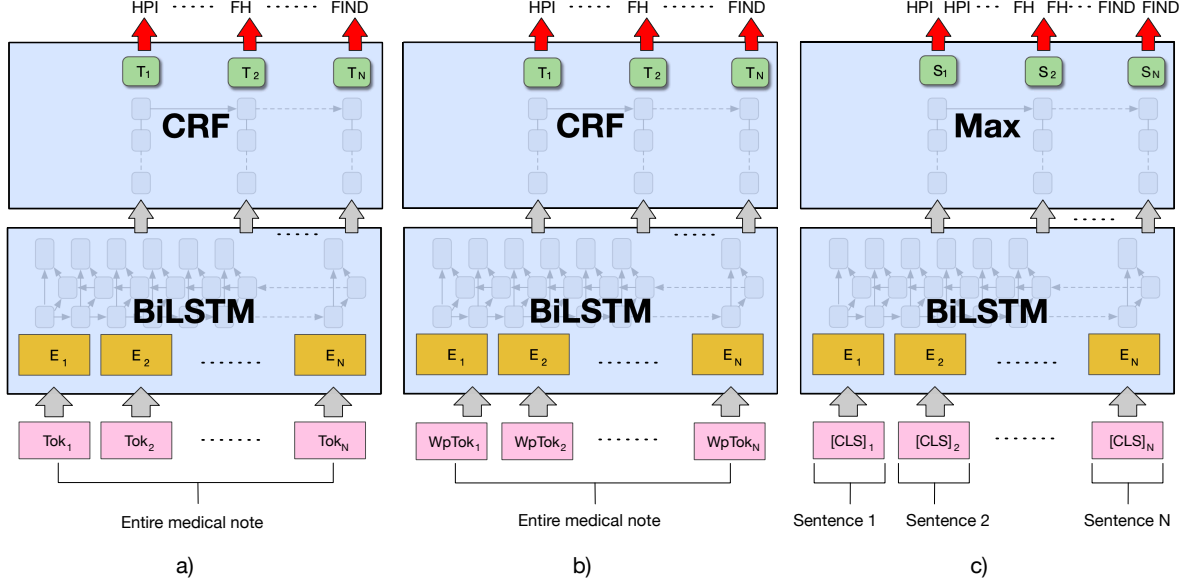


Figure 8: **MedSecId Baseline Models.** a) BiLSTM-CRF_{tok} BiLSTM model with non-contextual token input embeddings, b) BERT-CRF_{tok} BiLSTM model with BERT word piece token fixed input embeddings, c) BERT_{sent} BiLSTM model with [CLS] sentence embeddings using the per sentence majority label.

department¹. Notes written afterward are an essential source of data that provides an aspect of the patients' stay that is otherwise lacking in the corpus, such as daily progress *Physician* notes. However, *Radiology* and *Echo* notes from the MIMIC-III corpus apply to all hospital departments since they are uniform for all patients, regardless of their location, outpatient or inpatient. In addition, discharge summaries entail the entire hospital visit, including the ICU

¹Only five note categories are available (see Table I).

and the remainder of the admission. Since discharge summaries are included in the MIMIC-III dataset, they are also included in the MedSecId annotation set.

3.4 Baseline Models

With this data, we can now build some strong baseline models on MedSecId. Because the section text spans do not break on tokens, we cast our task as a named entity recognition (NER) using in, out (IO) encoding¹ on a 50 way classification including `<none>` for text with no sections (see Table VI). Using this encoding, I created several baselines across two BiLSTM models² for the purpose of future work benchmarking. These baselines include majority label metrics, a token BiLSTM with a CRF output layer (BiLSTM-CRF), and a sentinel BERT embedding (Devlin et al., 2019) long-short term memory (LSTM) model (see Figure 8). Aside from adjusting the LSTM hidden size, gradient clipping, and number of epochs, all parameters were held constant across all experiments.

BiLSTM-CRF_{tok}

The token model consists of a simple non-contextual input word embeddings, a LSTM layer and fully connected linear layer using a conditional random field (CRF) output with labels assigned by the Viterbi algorithm. Several embeddings were used with this model, including word-to-vector (word2vec) (Mikolov et al., 2013a; Mikolov et al., 2013b), Global Vectors for

¹in, out, begin (IOB) encoding was not used as there are no transitions from one section to another and to reduce the label count.

²No models use a BERT transformer, only BERT token and sentinel ([CLS]) embeddings.

Word Representation (GloVe) (Pennington et al., 2014) and fastText (Bojanowski et al., 2017) (Crawl) embeddings.

BERT_{sent}

To address the issue of exploding gradients, I created a sentence-based model using static BERT sentinel embeddings to lower the input length to the LSTM layer. The model assumes sections rarely break mid-sentence since every sentence is assigned one section. Sentences with more than one section annotation will lower end-to-end performance. However, 97.6% of the annotation set contains sentences with a single section for all tokens of the respective sentence as shown in Table IV. The output of the final layer of the first time step was used as the input to a LSTM. The LSTM output forwarded to a dense layer with one output neuron for each label and an output max over the label.

Unique Sections	Count	Proportion
1	253025	97.59%
2	5589	2.16%
3	589	0.23%
4	72	0.03%
5	11	0.00%

TABLE IV: **MedSecId Sentence Distribution.** Distribution of sentences having a single section label across all tokens of the respective sentence.

Both the standard small BERT model and BioBERT embeddings (Lee et al., 2020) are included in the baseline results (see Section 3.5). A ClinicalBioBERT baseline model (Alsentzer et al., 2019) would not provide a fair baseline metric for comparison with future works since it trained on the MIMIC-III corpus so it was excluded.

BERT-CRF_{sent}

Like BERT_{sent}, but adds a CRF layer with Viterbi assigned labels.

3.4.1 Implementation Details

The annotation set was randomly sampled per note and divided as a stratified dataset into training (80%), validation (10%) and test (10%) datasets. The medical note structure ontology (see Section 3.2.1) is distributed as both a RDF Turtle file and a CSV file along with the annotations. The publicly available¹ code to train, validate, and test the model also includes additional APIs to access the annotated data, perform inference with the pretrained model or train a new model. This codebase includes functionality to use the pretrained model or utilize the annotations for experimentation and is ready to easily be installed.² This codebase also references a related project useful for parsing MIMIC-III text, pseudo token replacement, and Postgres database to Python object relational mapping.

¹<https://github.com/uic-nlp-lab/medsecid>

²All that is required in a `pip install`. See the GitHub repo for details.

3.5 Results

The baseline models described in Section 3.4 were each trained until the validation loss converged, then early stopped. The results are summarized in Table V with label specific results in Table VI. I report performance metrics by counting correct predictions when the character span boundaries match exactly and the sections type match. If either do not match, it is counted as an incorrect prediction.

Id	Name	mF1	mP	mR	MF1	MP	MR
1	Majority Label	0.023	0.023	0.023	0.0	0	0.005
2	BERT _{sent}	0.925	0.925	0.925	0.589	0.616	0.6
3	BiLSTM-CRF _{tok} (word2vec)	0.927	0.927	0.927	0.778	0.78	0.801
4	BERT-CRF _{sent}	0.929	0.929	0.929	0.689	0.734	0.7
5	BERT _{sent} BioBERT	0.94	0.94	0.94	0.687	0.73	0.679
6	BERT-CRF _{sent} BioBERT	0.94	0.94	0.94	0.705	0.757	0.704
7	BiLSTM-CRF _{tok} (GloVE 50D)	0.954	0.954	0.954	0.76	0.783	0.765
8	BiLSTM-CRF _{tok} fastText	0.954	0.954	0.954	0.796	0.806	0.806
9	BiLSTM-CRF _{tok} (GloVE 300D)	0.955	0.955	0.955	0.787	0.801	0.788

TABLE V: **MedSecId Performance.** Summarization of performance metrics where mF1 is the micro F1, mP is the micro precision, mR is the micro recall, MF1 is the macro F1, MP is the macro precision, MR is the macro recall.

From the majority label, it is clear the models perform comparatively well as shown in the summary results in Table V. The highest performing GloVe model has a micro F1 of 0.96 and the highest performing fastText model has a macro F1 of 0.8. This 16 point spread is evident from how performance drops off for the bottom 13 section types. Many of these low performers

are those that were re-casted or re-grouped (see Section 3.2.1), and could be regrouped to an umbrella section type like *Labs/Imaging/Radiology* if such a rigorous delineation was not necessary.

The BERT_{sent} does not lag far behind, but its performance using sentinel embeddings does not capture sections as well as the token level models despite long document length. Performance significantly improved and models converged faster with the use of gradient clipping to alleviate issues of LSTM exploding gradients (Bengio et al., 1994).

3.6 Conclusion

MedSecId is a comprehensive dataset of medical annotations from the MIMIC-III corpus across five note types and 50 sections. The dataset contains section types, headers and patient age annotations. My dataset shows promising baseline results from simple models such as BiLSTMs with diverse inputs, but still leaves room for improvement by more sophisticated models.

These models improve pipelines that use rule based methods for SI as mentioned in Section 3.2.2. These pipelines include discharge note summarization Chapter 6, and other downstream tasks that would benefit from having header and non-section text removed such as training word embeddings such as ClinicalBioBERT.

The MedSecId model has been used in other discharge summary generation (Damm et al., 2024), and shown to improve performance for the summarization pipeline (see Chapter 6) for this work. In addition, it was also used to assist and supplement the DSProv annotations as explained in the next chapter.

Id	Label	mF1	mP	mR	MF1	MP	MR	Acc	Count
1	procedure	0	0	0	0	0	0	0	156
2	labs	0	0	0	0	0	0	0	436
3	prenatal-screens	0.276	0.276	0.276	0.216	0.5	0.138	0.276	105
4	imaging	0.357	0.357	0.357	0.263	0.5	0.178	0.357	990
5	comparison	0.414	0.414	0.414	0.195	0.333	0.138	0.414	222
6	code-status	0.513	0.513	0.513	0.226	0.333	0.171	0.513	150
7	wet-read	0.521	0.521	0.521	0.342	0.5	0.26	0.521	121
8	communication	0.556	0.556	0.556	0.179	0.25	0.139	0.556	133
9	impression	0.563	0.563	0.563	0.18	0.25	0.141	0.563	920
10	disposition	0.647	0.647	0.647	0.262	0.333	0.216	0.647	68
11	history	0.688	0.688	0.688	0.272	0.333	0.229	0.688	170
12	past-surgical-history	0.745	0.745	0.745	0.427	0.5	0.372	0.745	145
13	current-medications	0.746	0.746	0.746	0.142	0.167	0.124	0.746	1406
14	contrast	0.8	0.8	0.8	0.444	0.5	0.4	0.8	25
15	none	0.816	0.816	0.816	0.03	0.033	0.027	0.816	6378
16	discharge-disposition	0.83	0.83	0.83	0.151	0.167	0.138	0.83	513
17	addendum	0.833	0.833	0.833	0.151	0.167	0.139	0.833	3106
18	last-dose-of-antibiotics	0.872	0.872	0.872	0.466	0.5	0.436	0.872	397
19	indication	0.88	0.88	0.88	0.468	0.5	0.44	0.88	117
20	physical-examination	0.881	0.881	0.881	0.156	0.167	0.147	0.881	22113
21	image-type	0.884	0.884	0.884	0.313	0.333	0.295	0.884	181
22	discharge-condition	0.904	0.904	0.904	0.317	0.333	0.301	0.904	1490
23	infusions	0.909	0.909	0.909	0.476	0.5	0.455	0.909	99
24	history-of-present-illness	0.924	0.924	0.924	0.137	0.143	0.132	0.924	24950
25	discharge-medications	0.925	0.925	0.925	0.192	0.2	0.185	0.925	25088
26	flowsheet-data-vitals	0.932	0.932	0.932	0.482	0.5	0.466	0.932	2128
27	24-hour-events	0.954	0.954	0.954	0.244	0.25	0.238	0.954	1765
28	past-medical-history	0.959	0.959	0.959	0.163	0.167	0.16	0.959	5990
29	discharge-diagnosis	0.959	0.959	0.959	0.196	0.2	0.192	0.959	3578
30	family-history	0.968	0.968	0.968	0.328	0.333	0.323	0.968	1171
31	chief-complaint	0.968	0.968	0.968	0.492	0.5	0.484	0.968	1142
32	medical-condition	0.971	0.971	0.971	0.328	0.333	0.324	0.971	409
33	review-of-systems	0.977	0.977	0.977	0.494	0.5	0.488	0.977	724
34	labs-imaging	0.981	0.981	0.981	0.142	0.143	0.14	0.981	45855
35	discharge-instructions	0.986	0.986	0.986	0.166	0.167	0.164	0.986	23208
36	social-history	0.988	0.988	0.988	0.249	0.25	0.247	0.988	3114
37	allergies	0.989	0.989	0.989	0.331	0.333	0.33	0.989	891
38	assessment-and-plan	0.99	0.99	0.99	0.199	0.2	0.198	0.99	12728
39	reason	0.992	0.992	0.992	0.332	0.333	0.331	0.992	646
40	conclusions	0.994	0.994	0.994	0.498	0.5	0.497	0.994	2814
41	findings	0.998	0.998	0.998	0.333	0.333	0.333	0.998	6053
42	hospital-course	0.998	0.998	0.998	0.2	0.2	0.2	0.998	78321
43	social-and-family-history	1	1	1	1	1	1	1	52
44	technique	1	1	1	1	1	1	1	22
45	clinical-implications	1	1	1	1	1	1	1	36

TABLE VI: **MedSecId BiLSTM-CRF_{tok} Performance.** By label BiLSTM-CRF_{tok} performance (top 45 sections) where mF1 is the micro F1, mP is the micro precision, mR is the micro recall, MF1 is the macro F1, MP is the macro precision, MR is the macro recall, Acc is the accuracy, count is the the number of tokens encountered in the test set. The **<none>** label is for tokens with no section annotated.

CHAPTER 4

DISCHARGE SUMMARY PROVENANCE OF DATA

(This chapter expands on the paper “*Hospital Discharge Summarization Data Provenance*” by Landes et al. (2023) in the 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks of the Association for Computational Linguistics (ACL) conference.)

The MedSecId annotations were a useful dataset to train a model for the summarization pipeline as described in (see Chapter 6). They were also useful in an automatic discharge summary generation feasibility study, which was motivated by attention this task has seen in recent years. I believe this effort is warranted given the notes’ length and complexity, and that they are often riddled with poorly formatted structured data and redundancy in copy and pasted text. In this chapter, I investigate the feasibility of the summarization task by finding the origin, or data provenance, of the discharge summary’s source text. As a motivation to understanding the data challenges of the summarization task, I collaborated in the creation of a new dataset of 51 hospital admissions annotated by clinical informatics physicians. The dataset is analyzed for semantics and the extent of copied text from human authored electronic health record (EHR) notes.

Physicians often reference or copy EHR notes to the discharge summary. These copied notes include progress notes, consult notes, and test results. I call these previously written medical documents *note antecedents* since they are written prior to the discharge summary and can be used as source text as input to an automatic summarization system. The flow of information

from the note antecedent to the discharge summary is traced by annotating semantically similar content and copied text across notes and sections (such as *Brief Hospital Course*), which can be copied or paraphrased text. These annotations, called *note matches*, “tie” the notes’ discrete lexical spans of text. Each annotation is a text span in the discharge summary that carries a similar or identical meaning to its matching note antecedent as a span, which is a clinical medical note containing the original information that contributed to the discharge summary. In many cases, the text annotated by the note antecedents is later paraphrased or copied into the discharge summary.

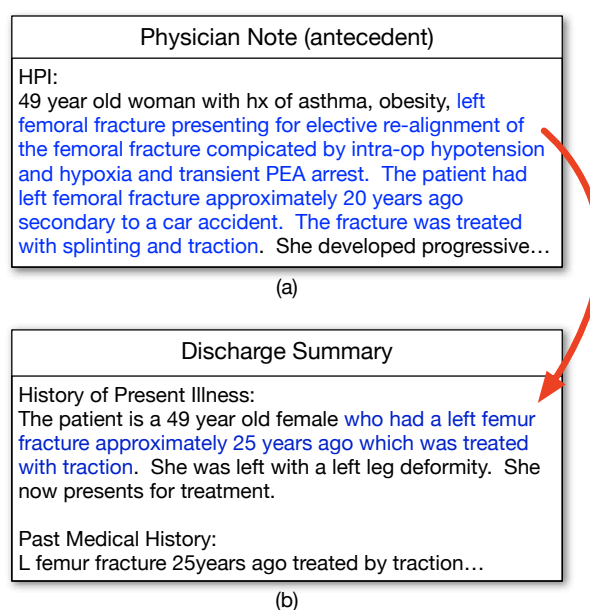


Figure 9: **A Note Match Annotation.** The annotation ties a physician note antecedent to a discharge summary with the matching text spans in blue and their coupling represented with the red arrow.

Figure 9a shows a note match annotation of a physician note and the linked text in the discharge summary given in Figure 9b. In this example, the patient’s encounter is documented in the *History of Present Illness* section of the note antecedent and later paraphrased in the summary at the time of the patient’s discharge. The red arrow represents the connection between these two lexical spans within an admission. This connection always starts from the original text in any clinical note and terminates in the discharge summary for the respective admission. For this work, both ends of the link are text spans written by physicians (see Section 4.2 for more detail on the annotation process).

4.1 Motivation

Understanding the extent of meaningful structured data is key in determining the feasibility and choice of methods needed to generate discharge summaries. In many instances, the copied text comes from high quality sources such as the *History of Present Illness* section of the physician note. However, redundancy, errors and data incoherence is pervasive in the vast amounts of medical data taken from patients during a hospital stay (Cohen et al., 2013). As much as 46% of the discharge summary is copied and pasted, 36% imported from structured data sources, leaving only 18% manually entered (Adams et al., 2021).

The extent of copied text in notes is well known, but the origins of discharge summary text is not. Since it is unclear how concepts arrive into the summary, extrapolating anything more than the measure of copied text and their semantics would necessitate understanding the decisions made by the physician while writing the summary. However, we can infer what is missing from a summarization by subtracting the portion represented by the note antecedents

of existing discharge summaries statistically. Knowing what is missing from the summary, in a probabilistic sense, and unique to the physician’s direct individual experience is fundamental for gauging the summarization performance upper bound.

For this reason, my primary goal for this work is to find and analyze the flow of data from notes written prior to the discharge summary, as shown in Figure 9. It is my position that a better understanding of the provenance of data, by note and section, is a crucial first step before automatic summarization can be successfully applied to discharge summarization future work. A secondary goal is to produce an unsupervised baseline and model for the research community.

To accomplish these goals three clinicians annotated admission records from the freely available MIMIC-III Version 1.4 (Johnson et al., 2016) corpus. These annotations uncovered where notes overlap and offer high quality human examples for supervised learning methods. My contributions are the annotated dataset (see Section 4.2), its analysis (see Section 4.2.2), and a novel unsupervised method using the word mover algorithm (Kusner et al., 2015) with clustering to assist in the annotation process (see Section 4.3).

4.2 Dataset

Three clinical informatics physicians annotated 51 admissions from the MIMIC-III Version 1.4 (Johnson et al., 2016) corpus for overlap with the EHR note antecedents to create a new annotation set called the DSProv dataset. The overlap was annotated by selecting the semantically similar portions of text from the discharge summary as shown in Figure 9. This provides statistics of overlap both at a note and a section level. To find the section overlap, the Med-

SecId corpus was used in most notes. In the few cases where note section annotations were not available, the pretrained MedSecId model was used to automatically section notes. Such sections are helpful to medical clinicians when finding topical information and useful for billing insurance companies. See Chapter 3 for the function and utility of sections in medical notes.

DSProv dataset admission candidates were selected from admissions annotated for sections in MedSecId to curate a richer combined dataset. More specifically, selection $\mathcal{S}_{\mathcal{C}}$ for admission A was chosen as $\mathcal{S}_{\mathcal{C}} = \arg \max_A \sum_{A \in \mathcal{C}} |A|$, where \mathcal{C} is the Medical Information Mart for Intensive Care III (MIMIC-III) corpus and $|A|$ are the number of notes in the admission. My early findings showed very little overlap of data with the exception of Consult (documentation by consulting physicians across departments) and Physician (typically written daily by the physician) notes. For this reason, admissions with these notes took priority in my annotation process to maximize note match overlap. While it could be argued that these notes' statistics may be combined given their likeness in purpose, it was decided to keep them separate for summarization.

Admissions were annotated with note matches using the following process:

1. Extract $\mathcal{S}_{\mathcal{C}}$ admissions from MIMIC-III.
2. For each admission $A \in \mathcal{S}_{\mathcal{C}}$:
 - (a) Read A 's discharge summary.
 - (b) For each note antecedent of the admission A , semantically similar or verbatim copied text was identified and each annotated note match annotated as:
 - i. A single span as character offsets tuples in the note antecedent (call it \mathbf{n}).

- ii. A single span as character offsets tuples in the discharge summary (call it **d**).
- iii. A link between **n** and **d**.

Each note antecedent may have several note matches (as can a discharge summary), but each has a single link across both note documents.

An annotation guide was created by the lead physician annotator, who then trained the other two physician annotators. Agreement on the criteria for note matches was decided on after each annotator completed one admission. The first admission annotations were then updated per the consensus agreed upon by all annotators. The same process of annotation, discovery, and agreement repeated for three additional admissions; one for each annotator. The remaining 39 were then split among the three annotators.

The 51 admission count might give the mistaken impression of a small dataset. However, as shown in Table VII, the extent of the annotation set is comprehensive with a total of 569 note matches from 291 notes that encompassed a little over 3 million tokens and 11.7 million characters (11.65MB). The number of admissions annotated was an aspect of the time consuming nature of the task. Each admission took an average 20 minutes for the physician to review and annotate as some admissions contained up to 494 notes. However, admissions had an average of 11.16 ($\sigma = 7.3$) notes and one admission had a single note antecedent. Statistics across discharge summaries and note antecedents are separate and independent; for example the 240 note antecedent count in Table VII does not include other information from discharge summaries.

Description	Count
Admissions	51
Match pairs	569
Discharge summary notes	51
Antecedent notes	240
Total Notes	291
Discharge summary tokens	1,695,466
Antecedent note tokens	1,323,422
Total tokens	3,018,888
Discharge summary characters	7,872,052
Antecedent characters	3,780,025
Total characters	11,652,077

TABLE VII: **DSProv Corpus Statistics.** The annotation statistics include discharge summary and note antecedent span, token and character counts.

4.2.1 Limitations

The MIMIC-III corpus includes a discharge summary for each admission. However, it is limited to patient’s time in the intensive care unit (ICU), meaning that the patient’s history for any time after transfer from the department is lost. Given most patients progress to lower severity departments as they recover from intensive care, a large cross section of the patient’s notes are missing from my analysis. In Section 4.2.2 I discuss the statistics that justify this conclusion.

4.2.2 Data Analysis

The DSProv annotation dataset was motivated by the summarization task. However, the dataset also provides insight into how EHR note antecedents are used by physicians to write discharge summaries and for the practicality of automatically summarizing them. This analysis

provides a quantitative justification for qualitative hypotheses based on clinician’s experience writing notes with data observed during annotation.

The human annotated note match text spans were tokenized to compute overlap using ROUGEL, which a longest common sequence n -gram metric that assesses the quality of generated summaries (Lin and Hovy, 2003; Lin, 2004). BLEU (a measure that assess the quality of machine generated translations) and BERTSCORE (an evaluation metric for scoring generated text) were also used in my quantitative analysis (Papineni et al., 2002; Zhang et al., 2020a). Each note match annotation includes the unique MIMIC-III note antecedent and discharge summary note identifier, absolute character offset in both notes, and the section they span. An additional set of annotations were automatically generated that break spans that overlap sections.

A normalized Levenshtein edit distance (Levenshtein, 1966) was used to measure the extent of copied and pasted text. Since the distance counts the minimum number of edits, rather than a relative measure to the note match span character length, it was normalized with:

$$\text{levsim}(\mathbf{w}_1, \mathbf{w}_2) = 1 - \frac{\text{lev}(\mathbf{w}_1, \mathbf{w}_2)}{\max |\mathbf{w}_1|, |\mathbf{w}_2|} \quad (4.1)$$

where lev is the Levenshtein edit distance, and $|\mathbf{w}_1|$ and $|\mathbf{w}_2|$ are lengths of the words in characters.

The statistical results by note categories (i.e. Radiology vs. Echo) and by note section. The section statistics are reported by grouping on the discharge summary section type, and the note antecedent section type.

4.2.2.1 Note Category

Table VIII provides statistics and similarity measures on the DSProv annotation set. The portion columns show the token overlap between each category of note antecedent (“Ant”) with the discharge summary (“DS”). The “LS” is the Levenshtein edit similarity as computed with Equation Equation 4.1. All similarity scores (“LS”, “BERTScore”, “ROUGE” and “BLEU”) are computed between the note antecedent and discharge summary span for each match.

Note Category	DS Portion	Ant Portion	Count	LS	BERTScore	RougeL	Bleu
Physician	28%	4%	310	49.50	68.14	44.10	19.34
Radiology	8%	14%	148	70.10	80.75	64.63	41.03
Echo	6%	29%	48	65.72	85.11	66.30	43.65
General	5%	4%	27	50.14	67.71	45.12	19.81
Consult	5%	4%	17	60.24	70.44	54.33	22.36
Nursing	2%	7%	8	23.43	44.56	15.25	0.91
ECG	1%	93%	5	80.94	79.07	74.08	52.33
Nursing/other	19%	5%	3	11.38	54.62	7.02	0.00
Rehab Services	2%	2%	2	26.31	60.04	21.26	0.00
Case Management	0%	4%	1	14.29	53.22	18.18	0.00

TABLE VIII: **Statistics by MIMIC Note Category.** The ratio of (D)ischarge (S)ummary tokens to total discharge summary tokens and the ratio of note (Ant)ecedent tokens to total note antecedent tokens. The Count column gives the number of notes annotated.

The highest discharge summary token overlap is with physician notes (28%), which I consider surprisingly high considering the MIMIC-III corpus only includes ICU notes¹ as mentioned in

¹Daily progress notes are recorded in the ICU as well.

Section 4.2.1. I expect the other discharge summary overlap statistics to be underrepresented for the same reason. This high token overlap supports the conjecture that daily progress notes are highly summarized with little copied text. The relatively low edit distance with a high BERTSCORE further supports this conclusion as surface similarity of copied text is low but the semantic similarity is high.

ECG (electrocardiogram notes), Radiology and Echo (echocardiogram data and analysis) notes show a higher similarity (80.94, 70.1, and 65.72 respectively) with the discharge summary since they are frequently copied and pasted. However, these statistics should be higher than presented since some spans include page breaks that result in counting superfluous header and footer tokens. Also note that ECG has the highest note antecedent portion, implying that most of these short reports make it into the discharge summary.

4.2.2.2 Section by Discharge Summary

Statistics by discharge summary section are given in Table IX. Labs and Imaging (71.55%) is a highly copied and pasted section from Radiology and Echo notes. This section takes into account cultures, blood results, and lab tests, and is copied and pasted as structured data.

Hospital Course has the highest discharge summary (14%) section representation. There is a high likelihood that this section has a high overlap due to transfer notes that describe patients moving between departments. These notes have an impact on summarization as they describe what happened to the patient during the hospital visit, including time of death. Generally speaking, good sections for summarization are those that have a low edit score (minimal copying and pasting) but a high similarity. In this case the comparatively low Levenshtein edit

Section Type	DS	Ant	Count	LS	BERTScore	RougeL	Bleu
Hospital course	14%	4%	125	41.13	63.39	33.35	11.54
Labs imaging	14%	14%	194	71.55	82.59	67.60	44.08
History of present illness	10%	7%	57	58.51	70.94	51.97	30.22
Physical examination	5%	3%	16	66.23	67.50	56.27	36.50
Addendum	5%	10%	2	54.89	74.53	43.05	19.00
Past medical history	4%	4%	64	44.60	68.54	40.22	20.32
Medication history	4%	4%	11	90.64	84.25	80.47	55.35
Imaging	2%	10%	1	97.92	86.99	88.06	54.16
Review of systems	2%	14%	1	24.84	58.46	21.08	2.41
Social history	1%	1%	15	59.55	71.79	49.95	22.66
Major surgical or invasive proc.	1%	0%	7	23.95	55.65	24.13	0.00
Family history	1%	1%	9	56.25	72.34	61.59	35.45
Default	0%	24%	1	25.00	28.58	14.29	0.00
Chief complaint	0%	0%	26	55.65	74.68	51.56	9.25
Discharge diagnosis	0%	1%	1	31.87	57.74	10.26	0.00

TABLE IX: **Statistics Grouped by Discharge Summary Section Type.** The (D)ischarge (S)ummary tokens to total discharge summary tokens and the ratio of note (Ant)ecedent tokens to total note antecedent tokens. The Count column gives the number of notes annotated. Only the top 15 sections with the highest discharge summary overlap are reported.

similarity (41.13) and somewhat high BERTSCORE (63.39) implies this section is a good target for automatic summarization as previously investigated in prior work (Adams et al., 2021).

4.2.2.3 Section by Note Antecedent

Table X shows the overlap from the perspective of the note antecedent. The *Assessment and Plan* section (the overall impression of the patient and how to treat them) has a high (15%) antecedent overlap. I found that the majority of the discharge summary content comes from the *Brief Hospital Course* section. The Echo note’s *Conclusions* section has a high portion of text that is summarized from note antecedents. This section’s echocardiogram information only consists of 8% on average of the notes annotated.

Section Type	DS	Ant	Count	LS	BERTScore	RougeL	Bleu
Indication	3%	28%	1	84.88	83.83	81.00	53.58
Conclusions	8%	26%	33	79.49	87.21	78.73	57.31
Findings	9%	15%	78	81.53	83.11	74.51	48.53
Technique	10%	9%	1	18.38	84.84	29.91	1.15
Impression	5%	8%	54	65.47	81.30	61.70	36.75
Review of systems	5%	8%	2	20.74	57.12	16.31	1.21
Wet read	4%	6%	2	98.28	91.18	76.19	52.69
Addendum	3%	6%	3	25.77	56.82	14.41	0.00
History	7%	5%	2	12.63	82.96	20.55	0.09
Hospital course	6%	5%	12	40.37	61.23	31.60	16.25
History of present illness	11%	5%	58	59.68	71.75	51.03	29.06
Comparison	5%	5%	8	13.79	80.72	21.58	0.55
Labs imaging	10%	4%	10	85.37	79.32	73.21	44.22
Assessment and plan	15%	4%	86	44.90	65.01	38.94	14.06
Discharge instructions	0%	4%	1	14.29	53.22	18.18	0.00

TABLE X: **Statistics Grouped by Note Antecedent Section Type.** (DS) is the ratio of discharge summary tokens to total discharge summary tokens and (Ant) is the ratio of note antecedent tokens to total note antecedent tokens. The Count column gives the number of notes annotated. Only the top 15 sections with the highest discharge summary overlap are reported.

4.3 Methods

Automatic methods to assist in bootstrapping the corpus were considered given the large amount of data the physician needs to sift through to find note matches¹. However, the task is challenging given the pages long document lengths, and precludes transformer pre-training methods. While the task has much in common with paraphrase matching and information retrieval tasks, it is fundamentally different in the way sections of text are matched. For example in question/answer systems, a query matches to an answer in the source text. However, my

¹The highest note count for an admission in MIMIC-III corpus is 1,233 notes.

task requires the correct text span in both note documents. Both the discharge summary and the note antecedent predicted text spans for the respective linked note matches are used when evaluating.

4.3.1 Evaluation

The human annotations were used for comparison since there is no previous baseline for this task. The unsupervised methods were then used to estimate the overlap portion to trace the origins of the discharge summary to the note antecedent, and then compared against the human annotations. The task is framed as a token classification task (without a label) since spans are token boundary demarcated.

My methods were evaluated using the SemEval 2013 Task 9.1 (Segura-Bedmar et al., 2013) entity extraction scoring method since it is flexible in its strictness as a score. Given the novelty and difficulty of the task, I used the *partial boundary matching*, which evaluates spans based on whether they overlap rather than on exact matches. Even though the SemEval measure is flexible for partial token matching, it is not ideal since this task aims to classify single token spans on a match-by-match basis, which is why the ROUGE metric was used as an additional reference point.

4.3.2 Word Mover

The word mover algorithm (Kusner et al., 2015) was used in the first step of my method. My method frames the task as a transportation problem by using the earth mover algorithm (Pele and Werman, 2009). The intuition follows from modeling probability distributions as piles of dirt that are moved from one location to another. The algorithm treats high dimensional word

embeddings of a source document as the probability distribution that “moves” words to the target document embedded distribution. The optimization algorithm minimizes the objective on words $\mathbf{w}_t, \dots, \mathbf{w}_T$:

$$\frac{1}{T} \sum_{t=1}^T \sum_{j \in \mathcal{N}(t)} \log p(\mathbf{w}_j | \mathbf{w}_t) \quad (4.2)$$

where $\mathcal{N}(\mathbf{n})$ are the neighboring words, and $p(\mathbf{w}_j | \mathbf{w}_t)$ is the word vector’s hierarchical softmax values. Word embeddings are made unit vectors and distance measures are euclidean. My experiments include non-contextual word embeddings (Mikolov et al., 2013a; Mikolov et al., 2013b) and static contextual embeddings from transformer architectures (Devlin et al., 2019). Because document word frequencies are used as the histogram weight to the earth mover algorithm, each surface word form was associated with its constituent wordpiece¹ (Wu et al., 2016). The centroid of the constituent wordpiece token(s) for each respective word was used for the embeddings.

4.3.3 Hybrid Semantic Positional Token Clustering

The word mover algorithm maps words from the discharge summary onto the note antecedent efficiently, but it does not help us link the note matches. A naïve approach would be to chunk tokens based on a metric such as cosine similarity. However, words with little similarity would frequently result in too many span breaks. We still need to cluster the word embeddings to group concepts in each document separately. However, this still does not address the “Swiss

¹Wordpieces are token sub-units with associated vectors and provided by the model’s tokenizer.

cheese” problem of note matches with too many “holes” (span breaks). I propose to simply add a component with a normalized scaled value to the word embedding. More specifically, I defined the concatenated position vector with:

$$\text{posemb}(\mathbf{w}) \triangleq \left[\text{emb}(\mathbf{w}); \left(\frac{\alpha_t \cdot i}{|T|} \right) \right] \quad (4.3)$$

where $\text{emb}_i(\mathbf{w})$ s the word embedding as an application of the language model; α_t is the token position component scaler; i is the index the i^{th} word w_i , and $|T|$ is the document token length.

The higher the token position component scaler (α_t) value is set, the more the word position is prioritized. This means that high values of the hyperparameter will create longer contiguous token spans, but at the cost of semantic similarity. This effect can be visually explained as a simple 2D word embedding with an additional token position axis. Figure 10 shows such a coordinate system with an example of an embedded span. On the positional axis, each token is spaced at even intervals. Because the positional components are proportionally scaled up for higher values of α_t , their relative distances shrink. On the other hand, if this value is lowered, their positional components diminish, effectively reverting the word points (word location in the embedded space) to their pre-trained vectors.

Once clustering of each document’s word points is complete, each cluster’s points are assigned to note matches. For each iteration over the Cartesian product, each document’s points are added to matches by associated cluster (see Algorithm 1). The source and target docu-

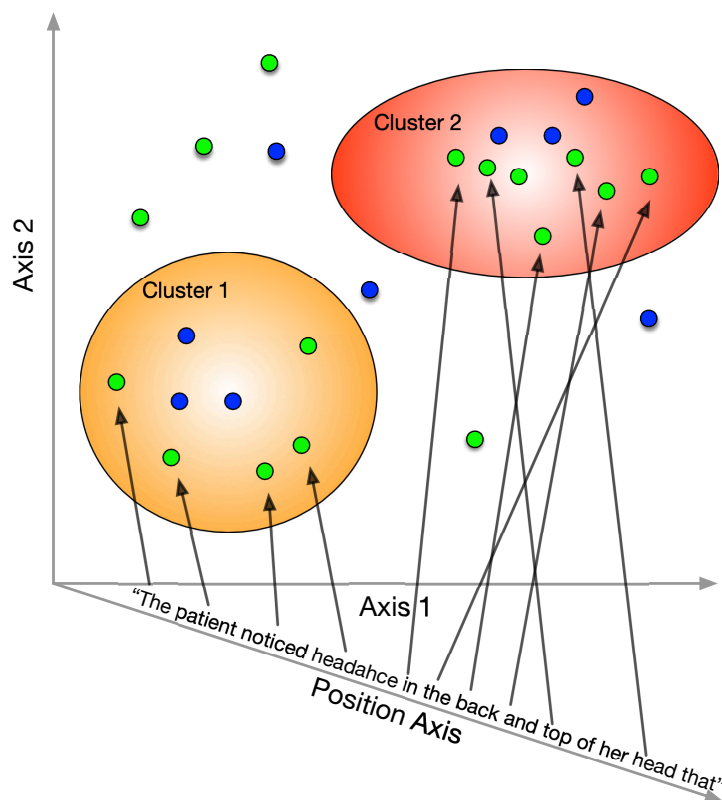


Figure 10: **Hybrid Semantic Positional Token Clustering.** Position embeddings on a third axis shows data blue word embeddings moving from cluster 1 to cluster 2. Cluster spans the discharge summaries (orange), the note antecedent (green) and arrows connecting the tokens to word points. (The example, and misspelling of “headahce”, is taken verbatim from the MIMIC-III corpus.)

ments are swapped and Algorithm 1 is run again to create flows¹ from the target to the source.

¹These are earth mover flows, which is calculated from the transportation plan, and not to be confused with the flow from the CALAMR max flow algorithm in Chapter 5.

Lexical overlapping matches are combined and their flows added together and sorted to create a ranking of matches.

Since each note match is assigned a flow (as a function of work to transport the word embedding) in both directions, they are combined as a single flow, which represents the highest similarities having the most information between the notes. Finally, matches are sorted in descending order by their flow values, so those with the maximum amount of information gain are ranked first.

4.4 Results

The DSProv dataset was provided to the unsupervised algorithm described in Section 4.3 and evaluated against the human annotated note matches. The match spans with the highest flow (see Section 4.3.3) were compared and scored using the measures listed in Section 4.3.1. Before the evaluation, Bayesian hyperparameter optimization (Bergstra et al., 2013) was used on the human annotated dataset on a subset of the data. The model’s hyperparameters were set to the Bayesian optimized values and evaluated. Of the 569 note matches annotated, an additional 359 were optimized on 500 iterations. This process was repeated on each word embedding for each note match.

While the hybrid method explained in Section 4.3.3 had a relatively high SemEval partial recall of 69.06 for matching discharge summaries spans, it suffered a low precision score. This implies finding spans is not an issue, but finding correct span boundaries as more difficult. I report both the good recall but poor precision in Table XI to help explain the kinds of challenges in matching spans between discharge summaries with note antecedents.

```

Input: Documents source A and target B
Output: N note match spans
1 Function MatchNoteSpans(A, B)
2   // assign word vectors and normalize to unit
3    $\mathcal{V}_a \leftarrow \frac{\text{emb}(\mathbf{w})}{\|\text{emb}(\mathbf{w})\|_2}; \forall \mathbf{w} \in A;$ 
4    $\mathcal{V}_b \leftarrow \frac{\text{emb}(\mathbf{w})}{\|\text{emb}(\mathbf{w})\|_2}; \forall \mathbf{w} \in B;$ 
5   // assign position embeddings
6    $\mathcal{P}_a \leftarrow \text{posemb}(\mathcal{V}_a);$ 
7    $\mathcal{P}_b \leftarrow \text{posemb}(\mathcal{V}_b);$ 
8   // assign word flows
9    $(\mathcal{F}_a, \mathcal{F}_b) \leftarrow \text{WordMover}(\mathcal{V}_a, \mathcal{V}_b);$ 
10  // cluster word points
11   $\mathcal{C}_a \leftarrow \text{Cluster}(\mathcal{P}_a);$ 
12   $\mathcal{C}_b \leftarrow \text{Cluster}(\mathcal{P}_b);$ 
13  // add matches
14   $\mathcal{M} \leftarrow \{\emptyset\};$ 
15  for  $f_a \in \mathcal{F}_a$  do
16    for  $f_b \in \mathcal{F}_b$  do
17      // get cluster from flow
18       $c_a \leftarrow \mathcal{C}_a[f_a];$ 
19       $c_b \leftarrow \mathcal{C}_b[f_b];$ 
20      // add the match and
21      // token points
22      if  $\{(c_a, c_b)\} \text{not} \in \mathcal{M}$  then
23         $\mathcal{M} \leftarrow \mathcal{M} \cup \{(c_a, c_b)\};$ 
24      end
25    end
26  end
27  return  $\mathcal{M};$ 
28 end
29 Function BiMatchNoteSpans(A, B)
30    $\mathcal{M}_{a \rightarrow b} \leftarrow \text{MatchNoteSpans}(A, B);$ 
31    $\mathcal{M}_{b \rightarrow a} \leftarrow \text{MatchNoteSpans}(B, A);$ 
32    $\mathcal{M}_{bi} \leftarrow \text{Sort}(\mathcal{M}_{a \rightarrow b}, \mathcal{M}_{b \rightarrow a});$ 
33 end

```

Algorithm 1: Matching Algorithm. The algorithm that matches text spans between documents. Word mover is used to map tokens to an embedded space as flows, positional embeddings are concatenated to word vectors, and documents are clustered. Then spans are created from spacial proximate words, and the top-K match spans are ranked by flow.

Table XI also shows how the performance for the note antecedents matching tells a better story. We see a similar pattern with a low precision, but a high recall with both netting a higher SemEval partial match harmonic mean of 15.85, which indicates that many matches are missed. Surprisingly SapBERT (Liu et al., 2021), which models semantic relationships of biomedical domain entities, performed worse than Sentence-BERT (Reimers and Gurevych, 2019). This suggests models trained for embedding and clustering provide better embeddings for my task. The non-contextual word embeddings do not perform as well with the exception of GloVe (Pennington et al., 2014) having the best ROUGE1 for discharge summaries.

Model	Se P	Se R	Se F1	Se Co	Rouge1	Rouge2	RougeL	EM
BioBERT	7.82	2.05	2.62	168.75	8.43	3.38	7.60	0.000
ClinicalBERT	8.61	2.89	3.91	330.60	10.76	5.07	9.39	0.216
Glove 300D	7.80	5.16	4.88	460.78	12.59	5.99	10.62	0.000
word2vec	9.76	34.15	12.30	3277.37	17.04	13.14	15.43	0.216
SapBERT	10.74	35.90	13.44	3316.16	19.36	14.69	16.95	0.216
SBERT	11.79	50.04	15.85	4635.78	19.62	16.48	18.03	0.216

TABLE XI: **Note Antecedent Score.** Performance of the unsupervised method for each word embedding model. Score methods include (SE)MEVAL-2013 (P)recision, (R)ecall, F1 and (Co)rrect mean. ROUGE1, ROUGE2, and ROUGEL F1 scores also provided with an (E)xact (M)atch score.

As mentioned in Section 4.2.1, only ICU notes and discharge summaries are provided in MIMIC-III. This has the effect of decreasing available information to potentially summarize and also has a negative impact on results as there are fewer examples to match between note

documents. However, it has an even greater impact on summarization since the machine cannot generate what is not available in the EHR.

4.5 Conclusion

DSProv is a new freely available dataset of 569 textual span matches between discharge summaries and note antecedents annotated by clinical informatics physicians. The analysis of my dataset presents new qualitative and quantitative findings of EHR notes. I have also presented a novel unsupervised method for annotating note matches using models tuned on human examples from my dataset. This dataset analysis provided insights necessary for assessing the feasibility of generating traceable and faithful summaries used in Chapter 6. The dataset can also be used to verify and tune the alignment method for medical text as described in the next chapter.

CHAPTER 5

AMR GRAPH ALIGNMENT

(This chapter expands on the paper “*CALAMR: Component Alignment for Abstract Meaning Representation*” by Landes et al. (2024) in the Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024).

The MedSecId annotations and model described in Chapter 3 provide a way to divide notes into smaller chunks of text to summarize and the DSProv annotations from Chapter 4 provide ground truth in how information flows from note antecedents to the discharge summary. These datasets were used to inform and facilitate discharge summarization explained in Chapter 6 using the alignment method explained in this chapter. This chapter gives a high level overview of AMR (abstract meaning representation) and flow networks (see Section 5.1). The motivation to use the alignment method for summarization is described (see Section 5.2) and explains how they were used to find overlap between the source text and the summary (see Section 5.3).

5.1 Introduction

The graph alignment method uses graphs to represent natural language with a graph algorithms not typically used in natural language processing (NLP). The output of the these algorithms are used in subsequent summarization components.

5.1.1 Abstract Meaning Representation

AMR is a semantic representation language that captures “who is doing what to whom” in a sentence (Banarescu et al., 2013). AMR graphs differ from parse trees as represent meaning rather than grammatical structure. For this reason, their representations are conducive to tasks such as summarization since they are balanced in their level of abstraction. AMR represents the abstract meaning of a sentence as a directed acyclic graph (DAG) as shown in Figure 11.

Each node of an AMR graph represents a concept (an idea grounded by natural language) or abstract placeholder (i.e. **person** for names). A concept node is often a predicated element (verb or “verb-like” main event or action core to an AMR or AMR subtree). Examples include verbs (**chase-01**), adjectives (**attract-01**) and nominalizations (the root **destroy-01** is used in place of “destruction”). Concept nominalizations (using the root verb form of a noun in place of a noun) may apply to entire event as the children nodes, or may only refer to the role player of an event. Attribute nodes are “constants” such as cardinals, date parts, and other named entity text (Banarescu et al., 2024).

PropBank (Kingsbury and Palmer, 2002) is a proposition lexicon database consisting of word senses and their arguments as “frames” (Loper et al., 2007). In AMR, predicated elements are modeled as rolesets in PropBank. The children of predicated element concept nodes are concept or attribute role nodes because they “play a role” in a specific way relative to the “sense” of the roleset. Each role is a member of a roleset, semantically labeled with a function tag and an enumerated index used as the edge label that connects it with its parent predicated

element concept node. Edges are also labeled with relationships between nodes that are not roleset concepts such as possessive, conjunction, location, etc.

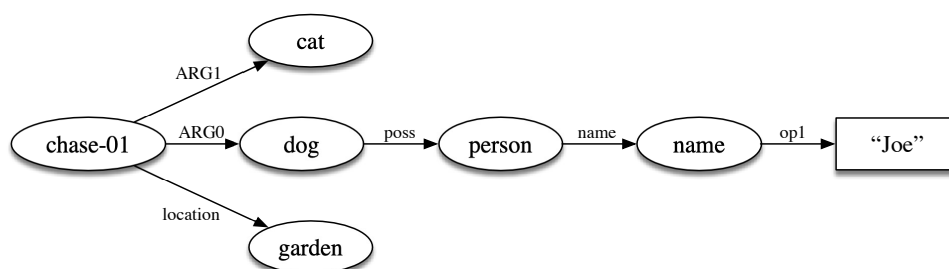


Figure 11: **Parsed Abstract Meaning Graph.** The AMR graph of sentence, “*Joe’s dog was chasing a cat in the garden.*” Example from Liu et al. (2015).

Figure 11 presents the AMR graph of the sentence “*Joe’s dog was chasing a cat in the garden.*”. In this example, **chase-01** is the roleset concept node having the form **<verb infinitive>-DD** where DD is a two digit that indexes the roleset. The **chase-01** node’s predicate is **chase** and **01** is predicating element index. AMR is also represented as a context free notation of Penman, which is a flat text representation widely used for abstract meaning representation graphs (Kasper, 1989). Figure 12 gives the example AMR in Penman format.

```

1 # ::snt Joe's dog was chasing a cat in the garden.
2 (c / chase-01~e.4
3   :ARG0 (d /dog~e.2
4         :poss (p /person
5               :name (n /name
6                     :op1 "Joe"))))
7   :ARG1 (c2 /cat~e.6)
8   :location~e.7 (g /garden~e.9))

```

Figure 12: **Penman Example.** The Penman notation of the sentence, “*Joe’s dog was chasing a cat in the garden.*” Example from Liu et al. (2015).

5.1.2 Flow Networks

A flow network (a.k.a. capacitance network) is a graph with two values associated with each edge: a capacity and a flow¹. They can be conceptualized as a network of pipes carrying water from node to node where the capacity is the maximum amount of water allowed to flow, and the flow is the amount of water traversing a pipe, which is an edge in the graph. While this example is water flowing through a pipe, it can be any material or grouping of objects such as traffic on a road network traversing an edge.

Every flow network has both a source s that produces an infinite amount of flow, and a sink t that accepts an infinite amount of flow. The flow that enters the system from t must be

¹The term “flows” is not to be confused with the earth mover flows that were calculated from the transportation plan in Chapter 4.

the same amount that exits through to the sink node t . Figure 13 shows an example of a flow network.

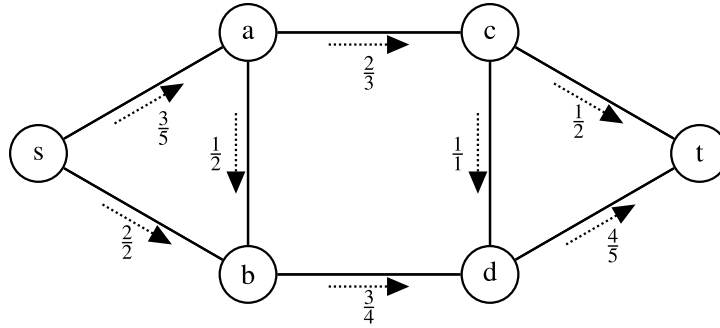


Figure 13: **Flow Network.** An flow network with the flow as the number on top and the capacity on the bottom labeled on each edge. In this example flow network the flow from the source node “s” has values 3 and 2 with capacities 5 and 2. The sink node “t” has incoming flow 1 and 4.

5.2 Motivation

The AMR graph alignments are used as input to the summarization model (see Chapter 6). These alignments are generated using an unsupervised algorithm using AMR graphs with edge and node embeddings. Node level (Cai and Knight, 2013) or network neighborhood (Opitz et al., 2021) alignments appear to perform well in small connected subgraph locales. However, when considering larger contexts, there are cases where we find similarity, but for the wrong

reasons. For example, “the dog chases the cat” has lexical similarity to “the cat chases the dog”, but their meanings are quite different.

This is most apparent with same surface text words in close lexical proximity across sentences that have high similarity but very different meaning. Such an observation highlights why n-gram measures such as ROUGE, which is a measure that assesses the quality of generated summaries (Lin, 2004), provide a poor method of scoring¹ (Schluter, 2017). Furthermore, ROUGE inadequately judges the faithfulness or hallucination rate of generated summaries (Maynez et al., 2020). The newer large language model (LLM) based semantic similarity BERTSCORE (Zhang et al., 2020a) has recently been proposed, but also falls short of quality generated summarized text assessments (Maynez et al., 2020) and lacks an element of transparency that inspires confidence in a scoring method.

My method addresses these issues by evaluating summarized or generated text for faithfulness and traceability using transparent methods by providing clear statistics of AMR graph alignments. The method aligns nodes that are semantically similar using PropBank roleset node, edge and token embeddings locally, and globally with the max flow algorithm. It provides not only a general AMR similarity score, but also a scoring method for summarized text. The method also produces a full node alignment essential to many down stream NLP tasks such as summarization. This alignment output is then later used to generate unsupervised

¹See page 60 for definitions of metrics.

training data for a supervised model. Section 5.13 presents renderings of the output graphs of the method.

5.3 Alignment Method

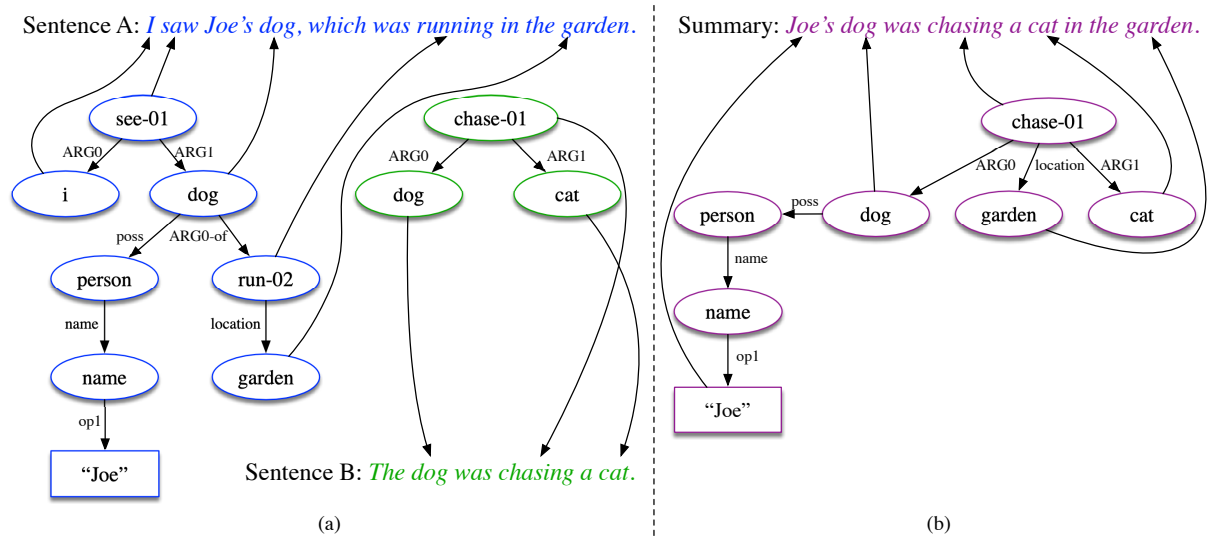


Figure 14: **AMR Graph Components.** Left (a) the source graph. Right (b) the summary graph. The aligned tokens shown as arrows. Example from Liu et al. (2015).

My method uses AMR sentence graphs as building blocks that are iteratively constructed into larger graphs that represent any arbitrary language structure such as paragraphs or documents. Each iteration of this process connects one or more graphs from the previous step; the input are *AMR graphs* that represent a sentence from human annotations (Figure 15a):

1. *Source* and *summary components* (see Figure 14) are each one or more sentence AMR graphs combined with a root node to form the document's source or summary (see Figure 15b).
2. The *alignment graph* is a bipartite graph of the source and summary components (see Figure 15c).
3. The *flow network* models the flow of information through a graph using two weighted edge values: capacities and flow (see Figure 15d).

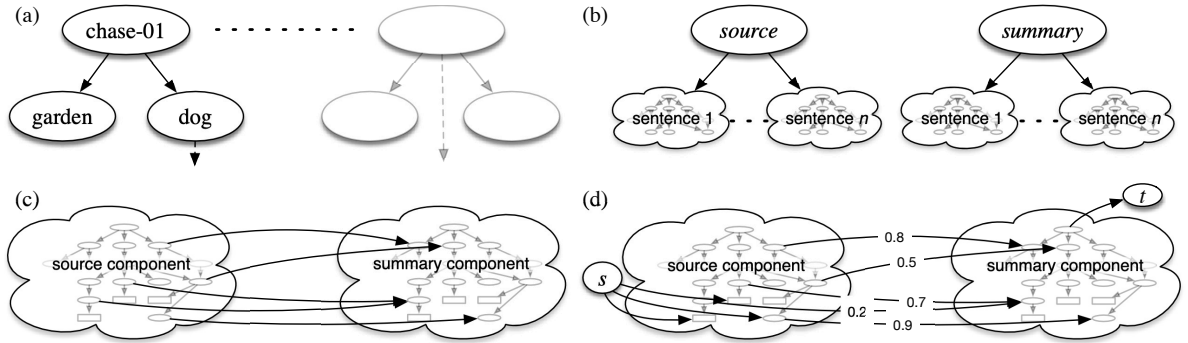


Figure 15: **Graph Construction.** The graph construction process starting with AMR graphs.

In this work, the flow material is information through the connected flow network constructed in Step 3. The direction of the flow goes from the source flow network node¹ to the leaf nodes of the summary component, over the bipartite alignment edges to the source component, and then to the source root. The amount of information flow through each alignment edge provides local scores of the strength of each source to summary node alignment pair and globally for the document.

Two kinds of alignments are described a) text-to-graph aligned tokens (TATs) associate AMR nodes to tokens (see Figure 14), and b) graph alignment edges connect nodes across the flow network. Embeddings generated from the source text of the TATs are attached to AMR graph nodes at inference time. However, static PropBank roleset embeddings (see Section 5.4) are computed before the method begins. The end-to-end pipeline includes:

1. Preprocessing PropBank embeddings (5.4).
2. Constructing the flow network (5.6.1).
3. Attaching graph embeddings from TATs and PropBank to predicating element nodes (5.6.1).
4. Computing alignment edge capacities (5.6.2). Capacity computation is described in Section 5.5.
5. Reducing the alignment (5.7).
 - (a) Run max flow algorithm (5.7.1).

¹The terminology difference between the source node s (flow network) and the source component (AMR graph) is explicitly differentiated since the source node can be connected to the source component.

- (b) Normalize flow-per-node.
- (c) Remove low flow alignment edges (5.7).
- (d) Go to step 5a until convergence.

5.4 AMR Graph Embedding

When the bipartite graph is created (see Section 5.6.1), contextual word embeddings are provided to every node and edge of the bipartite graph and used as input to the similarity measure. Additional embeddings are then computed and combined for local the network neighborhood of the bipartite graph. The process described in this section is used as the input to computing the capacities described in Section 5.5, so we only need compute the embeddings for those nodes that are aligned. Only alignments of the same type of node for each node pair is created, meaning no concept to sentence node alignment edges, for example. This reduces the number of capacities we have to calculate since it reduces the types of node symmetric pairings of: document (Section 5.4.3), sentence (Section 5.4.4), concept (Section 5.4.5), and attribute (Section 5.4.6). But before each of the node’s embeddings are calculated, a preprocessing step is necessary.

5.4.1 PropBank

The frame files are parsed and the hierarchical structured data added to an SQLite normalized database. The embeddings are generated from Sentence-BERT (SBERT), a siamese

network model that captures semantic similarity (Reimers and Gurevych, 2019) of sentences.

The input text to the model is taken from the following description text:¹

- A roleset uses the name (i.e. “follow, pursue” from `chase-01`²).
- A role uses the name (i.e. “follower”, which is the description of the first roleset role of `chase-01`).
- A function tag uses the semantic tag of the role (i.e. “PAG” for “prototypical agent” for the first roleset role of `chase-01`).
- A role edge uses the AMR role descriptions (i.e. “argument frame” for `:ARGO` or “possessive” for `:poss`).

The PropBank frames and embeddings are provided as a freely available³ library in Python with object oriented access via a SQLite database. The large model (768D) SBERT embeddings output are compiled along with the database into a large compressed file that is downloaded and then indexed on first use. These embeddings are then used for computing the capacities of the alignment edges of the bipartite graph as described in Section 5.4.

¹Even though much of this text is short (even a single word long), they are treated as full sentence text with respect to model input.

²PropBank uses a dot (.) index delimiter for roleset identifiers (i.e. `chase.01`), but AMR syntax is used for simplicity.

³<https://github.com/plandes/propbankdb>

5.4.2 Node to Text Embeddings

Each type of node has one or more strings used as input to SBERT. Since document nodes are placeholders for sentences and other aggregated node collections, only sentence, concept and attribute nodes have text to extract. Extracting text is straight forward from sentence nodes since the node has the chunked sentence text parsed from when the graph was constructed. Concept and attribute nodes use the text from the text-to-graph aligned tokens (TATs) if there are any (see Section 5.6.1). If not, the role ID lemma or label text (i.e. “chase” from `chase-01` or “dog”) for concept nodes, or the constant (i.e. “Joe”) for attribute notes (see Figure 11) is extracted.

Once the text is extracted from the node, embeddings are produced by inferencing with a forward pass using the SBERT model. The method of creating embeddings is different for each node. However, the following notation defines an application of the transformer model to generate the embeddings from the text of all nodes:

$$\mathbf{S}_m : \mathcal{S} \rightarrow \mathcal{X}, \mathcal{X} \in \mathbb{R}^{|\mathcal{S}| \times d}, \mathbf{s} \in \mathcal{S} \quad \text{embedding mapping} \quad (5.1)$$

where:

- \mathbf{S}_m is the mapping from input sentences to embeddings using model m
- m is the SBERT large model (`all-mpnet-base-v2`)¹

¹https://www.sbert.net/docs/pretrained_models.html

- \mathcal{S} is the list of sentences provided to the model
- \mathcal{X} is the output embedding with dimension d
- d is the SBERT model output dimension (768)

Equation 5.1 represents an inference as a mapping from sentences to embeddings of dimension d . Equations for each node use this notation to better explain how their content is used to create or combine (average, add or max) embeddings.

5.4.3 Document Nodes

The embeddings for document nodes (see Section 5.6.1) are computed as the mean of their constituent subgraph children node embeddings:

$$\text{emb}_d(\mathbf{n}) \triangleq \frac{1}{|\mathcal{C}(\mathbf{n})|} \sum_{\mathbf{u} \in \mathcal{C}(\mathbf{n})} \text{emb}(\mathbf{u})$$

For paragraph nodes, this is the mean of the sentence nodes (children nodes $\mathcal{C}(\mathbf{n})$ are defined in Section 5.7.2) that compose that paragraph. Similarly, the mean of all electronic health record (EHR) medical notes for a patient's admission at the root level for the discharge summaries case. The $\text{emb}(\mathbf{u})$ notation is used to get the node embeddings of node \mathbf{u} . Each node has their own embedding definition including sentence node embeddings defined in Section 5.4.4

5.4.4 Sentence Nodes

The embeddings of sentence nodes are taken directly from the SBERT [CLS] token model output, and defined as:

$$\text{emb}_s(\mathbf{n}) \triangleq \mathbf{S}_m(\mathbf{s}_n) \in \mathbb{R}^d$$

where:

- $\text{emb}_s(\mathbf{n})$ is a sentence embedding for sentence s
- \mathbf{n} is a sentence node in the graph
- \mathbf{S}_m is the mapping from input sentences to embeddings using model m
- \mathbf{s}_n is the sentence text for node \mathbf{n}
- $\mathbf{S}_m(\mathbf{s}_n)$ is the sentence to embedding mapping given in Equation Equation 5.1

In addition to sentence $\text{emb}_s(\mathbf{n})$, node embeddings include document nodes ($\text{emb}_d(\mathbf{n})$), concept nodes ($\text{emb}_c(\mathbf{n})$), and attribute nodes ($\text{emb}_a(\mathbf{n})$). These embeddings are used in later definitions.

5.4.5 Concept Nodes

Concept nodes have a rich set of information at their disposal, which is leveraged to better contextualize and embed them. As mentioned in Section 5.4.2, their embeddings are created from the graph-to-text token alignments. All predicating element concepts contain PropBank information and embeddings generated from the preprocessing step described in Section 5.4.1.

Aligned Tokens

The TAT embeddings are defined as:

$$\text{emb}_t(\mathbf{n}) \triangleq \begin{cases} \sum_{\mathbf{w}} \sum_i^{|w|} \mathbf{S}_m(\mathbf{s}_n)_i & \text{if node } \mathbf{n} \text{ has alignments} \\ \mathbf{1}_d & \text{otherwise, 1-vector of size } d \end{cases} \quad \text{token} \quad (5.2)$$

where:

- \mathbf{n} is the concept node
- \mathbf{s}_n is the sentence text of the sentence node of which concept node \mathbf{n} is a member
- $\mathbf{S}_m(\mathbf{s}_n)_i$ is the SBERT embedding of i th wordpiece¹ of word \mathbf{w} of sentence text \mathbf{s}_n
- \mathbf{w} is an aligned token, which contains one or more word pieces
- d is the SBERT embedding dimension, which is 768

Equation 5.2 applies to nodes with TATs, which include both concept and attribute nodes. Each aligned token's embeddings are used, but each token can be represented by more than one wordpiece (Wu et al., 2016). Instead of the [CLS] token embedding, the last output layer embedding of the model at position i of word token \mathbf{w} is used. Figure 16 shows an example AMR graph with TAT **garden** at index 9 in the sentence, which gives the 9th index in to the SBERT embedding output, which in this case has dimensionality 10×768 for the model's last output layer assuming each token maps to a single wordpiece.

¹See page 66 for definition of wordpiece.

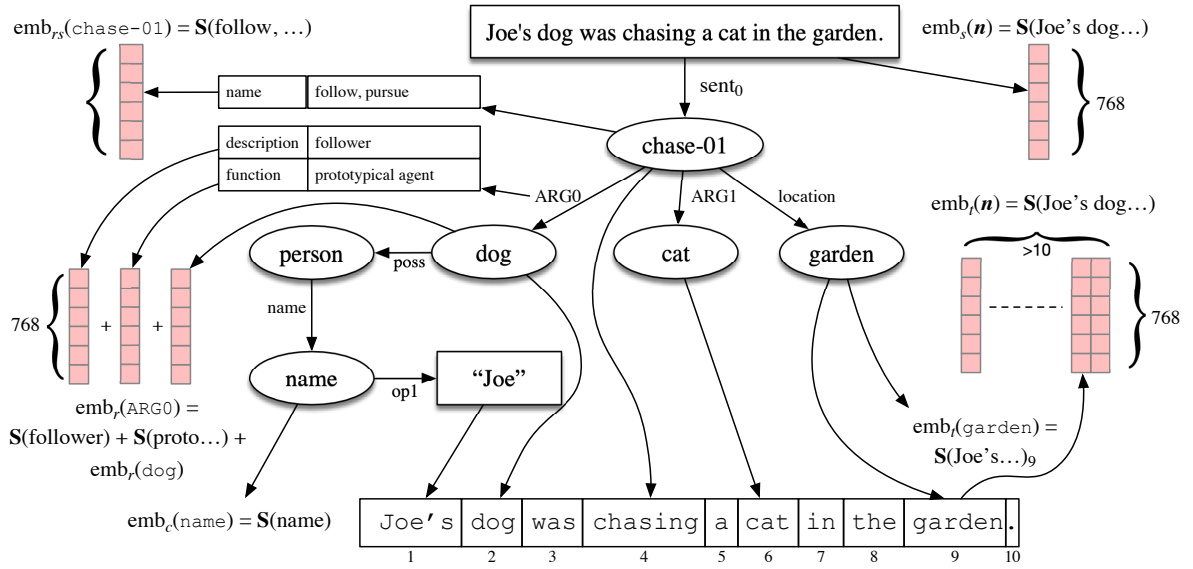


Figure 16: **Graph Embeddings.** An AMR graph of the sentence, “*Joe’s dog was chasing a cat in the garden.*”. The embedding node definitions with the model output are given for TAT **garden** at index 9, the roleset (**chase-01**), the role (**ARG0**) with description “follower” and the function tag “prototypical agent”. Each token is represented as one word piece for simplicity. Example from Liu et al. (2015).

Role

A roleset is a grouping of roles for a verb concept node. It, and all other PropBank entries used for embeddings, are preprocessed (see Section 5.4.1). A concept’s role embedding on edge \mathbf{e} are defined as:

$$emb_r(\mathbf{e}) \triangleq emb(\mathcal{C}(\mathbf{e})) \cdot \rho_c + [S_m(\mathbf{s}_r) + S_m(\mathbf{s}_f)] \cdot \rho_r \quad (5.3)$$

where:

- \mathbf{e} is the role edge
- \mathbf{s}_r is the name of the role (i.e. “follower” from the first role of **chase-01**)
- \mathbf{s}_f semantic tag of the role (i.e. “prototypical agent” from the first role of **chase-01**)
- $\text{emb}(\mathcal{C}(\mathbf{e}))$ is the embedding of role edge node child (see Section 5.7.2)
- role child weight (ρ_c) and role weight (ρ_r) are hyperparameters that allow for assigning importance to the specific terms

The **ARG0** role embedding of **chase-01** is shown in Figure 16 as the component-wise sum of the child node **dog**’s embedding, and the edge’s description and function embedding.

Node

The information that comes from the PropBank roleset can be thought of as “metadata” attached to verb concept nodes and static across all concept verb nodes of the same verb roleset. For example, if **chase-01** is found in two sentences, as in the case of the graph in Figure 20, they will share the same roleset and its embedding. In addition to the roleset, most concept nodes have TATs regardless of verb or non-verb concept node, which are also factored into the node’s embedding. The embedding for a concept verb node \mathbf{n} is defined as the weighted mean of the TATs, the roleset, and the role edges:

$$\text{emb}_{rs}(\mathbf{n}) \triangleq \mathbb{1}_v[\mathbf{n}] \cdot \mathbf{S}_m(\mathbf{s}_{rs}) + (1 - \mathbb{1}_v[\mathbf{n}]) \cdot \mathbf{S}_m(\mathbf{s}_n) \quad (5.4)$$

$$\text{emb}_c(\mathbf{n}) \triangleq \text{emb}_t(\mathbf{n}) \cdot \omega_t + \text{emb}_{rs}(\mathbf{n}) \cdot \omega_{rs} + \mathbb{1}_v[\mathbf{n}] \left[\sum_{\mathbf{e} \in \mathcal{N}(\mathbf{n})} \text{emb}_r(\mathbf{e}) \cdot \omega_r \right] \quad (5.5)$$

where:

- \mathbf{n} is the concept node
- $\text{emb}_t(\mathbf{n})$ is the TAT embedding produced from Equation Equation 5.2
- $\mathbf{S}_m(\mathbf{s}_{rs})$ is the SBERT embedding created from the roleset the name text \mathbf{s}_{rs} (i.e. “follow, pursue” from `chase-01`)
- $\mathbf{S}_m(\mathbf{s}_n)$ is the SBERT embedding created from the text of the non-predicating element node (see Section 5.1.1) text \mathbf{s}_n (i.e. a noun instance might have “dog” or abstract meaning “person”)
- $\text{emb}_{rs}(\mathbf{n})$ is the roleset node embedding definition in Equation Equation 5.4
- $\mathbb{1}_v[\mathbf{n}]$ is the indicator function yielding 1 when \mathbf{n} is a predicating element node and 0 when it is a noun and other abstract meaning node
- $\mathbf{e} \in \mathcal{N}(\mathbf{n})$ are the outgoing role edges of node \mathbf{n} that connect to its children
- $\text{emb}_r(\mathbf{e})$ is the role embedding defined by Equation Equation 5.3
- ω_t , ω_{rs} and ω_r are hyperparameters that allow for assigning importance to the specific terms

$\mathbb{1}_v[\mathbf{n}]$ is the indicator function yielding 1 when \mathbf{n} is a predicating element node and 0 when it is a noun and other abstract meaning node. This “toggles” the use of the roleset embedding for predicating element nodes. Otherwise, it uses the text of the node itself, which could be “dog” for a noun node or “person” for an abstract node. Figure 16 shows the constituent parts of the $\text{emb}_c(\mathbf{n})$ embeddings.

5.4.6 Attribute Nodes

If models that assigned TATs were completely accurate, attribute nodes would always have at least one token alignment. Instead, they leave some attributes with only the surface text given to the node by the text-to-graph model. From what we know about concept nodes, the attribute’s embeddings definition looks much like equations Equation 5.4 and Equation 5.5:

$$\text{emb}_a(\mathbf{n}) \triangleq \min(1, |\mathcal{T}|) \text{emb}_t(\mathbf{n}) + (1 - \min(1, |\mathcal{T}|)) \mathbf{S}_m(\mathbf{s}_n)$$

where:

- \mathbf{n} is the attribute node
- \mathcal{T} is the set of TATs
- $\text{emb}_t(\mathbf{n})$ is the TAT embedding produced from Equation Equation 5.2
- $\mathbf{S}_m(\mathbf{s}_n)$ is the SBERT embedding created from the text of the attribute (i.e. an attribute might have a name constant such as “Joe”)

5.4.7 Network Neighborhood

The nodes adjacent to a target node are traversed to create a local network neighborhood embedding. This local context embedding better contextualizes node alignments. Figure 17 shows the network neighborhood around the target node **run-02**. The dark blue nodes are the k^{th} order neighbor set that are at exactly k hops away from the target node. They are colored according to how much they influence the embeddings scaled by the hyperparameter network neighborhood weights (Λ). This hyperparameter is an array of real values for each concentric

k^{th} order neighbor set. The higher the k the lower the weight, and influence their embeddings have on the target node.

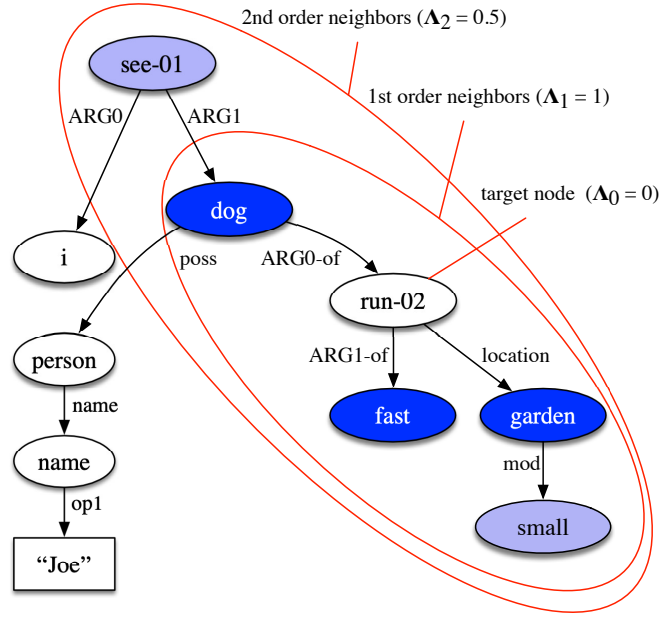


Figure 17: **Network Neighborhood Weights.** The AMR graph for the sentence, “I saw Joe’s dog, which was running fast in the small garden.”, with the network neighborhoods of verb concept node **run-02** and their weights. The Λ weights are shown with default hyperparameters. Example from Liu et al. (2015).

Let the selection of the k^{th} order neighbor set from node \mathbf{n} be $\mathcal{U}(\mathbf{n}, k)$, then the network neighborhood embedding is defined as:

$$\text{emb}_{\mathbf{n}}(\mathbf{n}) = \sum_{\mathbf{i}}^k \sum_{\mathbf{u} \in \mathcal{U}(\mathbf{n}, \mathbf{i})} \text{emb}(\mathbf{u}) \cdot \Lambda_{\mathbf{i}} \quad \text{network neighborhood embedding}$$

$$\forall \mathbf{x} \in \Lambda : \mathbf{x} > 0 \quad \text{hyperparameter weight constraint}$$

where:

- k is the maximum order k^{th} order neighbor set.
- $\mathcal{U}(\mathbf{n}, \mathbf{i})$ is the k^{th} order neighbor set \mathbf{i} hops from \mathbf{n}
- $\text{emb}(\mathbf{u})$ node \mathbf{n} 's embedding
- $\Lambda_{\mathbf{i}}$ is the \mathbf{i}^{th} 's weight hyperparameter that dampens the node \mathbf{n} 's embedding
- $\Lambda_0 = 0$ so that the target node's embedding is not added to $\text{emb}_{\mathbf{n}}(\mathbf{n})$, which is why the target node is not colored in Figure 17

5.5 Capacity Calculation

The flow network capacities of the alignment edges are calculated using the semantic cosine similarity of the node pairings as explained in Section 5.4, which derived from an exponentiated cosine similarity:

$$\begin{aligned} \text{cossim}(\mathbf{n}_1, \mathbf{n}_2) &= \left[\frac{\text{emb}(\mathbf{n}_1) \cdot \text{emb}(\mathbf{n}_2)}{\|\text{emb}(\mathbf{n}_1)\| \|\text{emb}(\mathbf{n}_2)\|} \right] && \text{cosine similarity definition} \\ \text{sim}(\mathbf{n}_1, \mathbf{n}_2, \mu) &= \text{cossim}(\mathbf{n}_1, \mathbf{n}_2)^\mu && \text{similarity measure} \quad (5.6) \\ \mu &> 0 && \text{hyperparameter constraint} \end{aligned}$$

where:

- \mathbf{n}_1 and \mathbf{n}_2 are the bipartite node pair
- $\text{emb}_{s_1}(\mathbf{n})$ and $\text{emb}_{s_2}(\mathbf{n})$ are the embeddings of the node pair
- μ is the hyperparameter that non-linearly adjusts the similarity down ($\mu > 1$) or up ($\mu < 1$)

Equation 5.6 defines the $\text{sim}(\mathbf{n}_1, \mathbf{n}_2, \mu)$ as the alignment edge capacity for a node pair. The μ hyperparameter is used to penalize certain node pairs for over-zealous similarities. As mentioned in Section 5.2, some nodes or network neighborhoods might be similar only locally. This is addressed globally by the max flow algorithm (see Section 5.7) and locally with combined network neighborhood embeddings (see Section 5.4.7). The μ hyperparameter is used

to adjust node neighborhoods (μ_n) and each node type: document (μ_d), sentence (μ_s), concept (μ_c), attribute (μ_a).

The method of obtaining the embedding $\text{emb}(\mathbf{n})$ is done by: a) averaging over child nodes for document nodes (Section 5.4.3), b) pooled sentence embeddings for sentence nodes (Section 5.4.4), c) predicating element text, TATs, role, roleset, and surrounding nodes' embeddings for concept nodes (Section 5.4.5), and d) literal text for attribute nodes (Section 5.4.6).

Capacity Definition

The capacity value assigned to the alignment edges added to the graph described in Section 5.6 are defined as:

$$\begin{aligned} \sigma(x) &= \left(1 + \exp((\tau_x - x) \cdot \tau_c)\right)^{-1} + \tau_y && \text{sigmoid} \\ \text{cap}(\mathbf{n}_1, \mathbf{n}_2, \mu) &\triangleq \min\left(\max(\text{sim}(\mathbf{n}_1, \mathbf{n}_2, \mu) + \sigma(\text{sim}(\mathbf{n}_1, \mathbf{n}_2, \mu)), 0), 1\right) && \text{capacity} \end{aligned} \quad (5.7)$$

where:

- $\sigma(x)$ is the translated and compressed sigmoid function used to adjust the similarity measure
- $\text{sim}(\mathbf{n}_1, \mathbf{n}_2, \mu)$ is the similarity function (see Equation Equation 5.6)
- τ_x, τ_y are the transition hyperparameters and τ_c is the compression hyperparameter.

The sigmoid network neighborhood skew (τ) hyperparameter defaults are $\tau_x = 0.5$, $\tau_y = 0.5$, and $\tau_c = 1$ so that network neighborhood similarities values are “pushed” away in both

directions from 1. The τ_c hyperparameter “squeezes” the curve along the x-axis so capacity values are adjusted at different rates.

Document Node Capacities

The capacities of the document node alignment edges are computed using the document node embeddings (see Section 5.4.3) as $\text{cap}_d(\mathbf{n}_1, \mathbf{n}_2, \mu_d)$ defined in Equation Equation 5.7.

Sentence Node Capacities

Similar to document node capacities, the sentence node alignment edge capacities are computed using the sentence node embeddings (see Section 5.4.4) as $\text{cap}_s(\mathbf{n}_1, \mathbf{n}_2, \mu_s)$ defined in Equation Equation 5.7. However, a sentence skew is also computed for each sentence pair using the hyperparameter sentence dampen (γ), which is used to scale the concept (Section 5.4.5) and attribute (Section 5.4.6) capacities. The sentence skew of two sentences is defined as:

$$\text{sentskew}(\mathbf{s}_1, \mathbf{s}_2) \triangleq \text{cap}_s(\mathbf{n}_1, \mathbf{n}_2, \mu_s) \cdot \gamma + (1 - \gamma) \quad \text{sentence skew} \quad (5.8)$$

$$\gamma \in [0, 1] \quad \text{hyperparameter constraint}$$

The hyperparameter γ is the slope for the linear dampening of nodes under a sentence by sentence cosine similarity. The higher the value the lower the sentence similarity, which leads to lower concept and attribute node similarities.

Concept Node Capacities

Similar to document nodes and sentence nodes, concept node alignment edge capacities are computed using concept node embeddings (see Section 5.4.5) as input to $\text{cap}_c(\mathbf{n}_1, \mathbf{n}_2, \mu_c)$

defined in Equation Equation 5.7. However, the capacities of the sentence node descendants are scaled by their sentence skew produced in the sentence capacity calculation step defined in Sentence Node Capacities. Concept node capacity assignments are defined as:

$$\text{cap}_c(\mathbf{n}_1, \mathbf{n}_2, \mu_c) \triangleq \begin{cases} 1 & \text{if node variable match} \\ \text{cap}(\mathbf{n}_1, \mathbf{n}_2, \mu_c) \cdot \text{sentskew}(\mathbf{s}_1, \mathbf{s}_2) & \text{otherwise} \end{cases} \quad (5.9)$$

where:

- \mathcal{E}_a is the set of alignment edges
- $(\mathbf{s}_1, \mathbf{s}_s)$ is a bipartite aligned sentence node pair
- $(\mathbf{n}_1, \mathbf{n}_s)$ is a bipartite aligned node descendant pair of their respective sentence nodes
- $\text{sentskew}(\mathbf{s}_1, \mathbf{s}_2)$ is the sentence skew (see Equation Equation 5.8)

The capacity in Equation Equation 5.9 is set to the max value 1 when the variable name is the same for both concept nodes (i.e. “c” in “c / chase-01”). This is an example of an AMR reference between the source component and summary component. Otherwise, it is scaled as a function of the sentence embedding to which it belongs using the sentence skew defined in Equation Equation 5.7. Figure 18 shows a strongly semantic similarity with capacity 0.9, and weakly semantically similar sentence with capacity 0.5, between the bipartite components. The strongly similar sentence uses blue to denote the sentence alignment and that sentence alignment’s effect on the concept and attribute alignment edges. For example, the strongly

similar sentence's **person** only loses a flow value of 0.02 after applying the sentence skew, while the weakly similar sentence's **cat** loses half its capacity $(0.4)^1$.

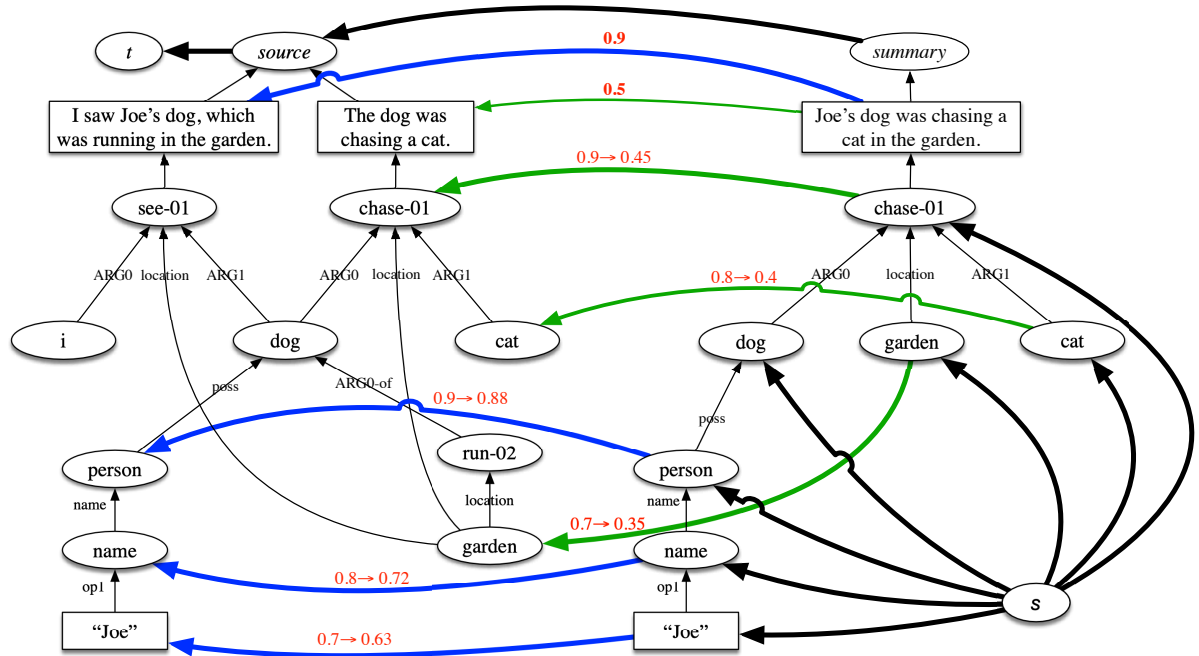


Figure 18: **Sentence Scaled Capacities.** The scaled capacities with alignment edges colored by related sentence with bold red capacities. The concept and attribute node capacity updates are in red with the initial capacity on the left and the sentence skew updated on the right. All capacity values were created to illustrate the concept of sentence dampening and not the real values using these sentences. Example from Liu et al. (2015).

¹All capacity values were created to illustrate the concept of sentence dampening and not the real values using these sentences.

Attribute Node Capacities

Attribute node capacities are calculated in the same way as concept nodes defined in Equation Equation 5.9. However, they have no variable, so the capacity definition given in Equation Equation 5.7 is used directly with the attribute's sentence skew scaled embeddings:

$$\{(\mathbf{n}_1, \mathbf{n}_2) \in \mathcal{E}_a, (\mathbf{s}_1, \mathbf{s}_2) \in \mathcal{E}_a \mid \mathbf{n}_1 \in \mathcal{D}(\mathbf{s}_1) \wedge \mathbf{n}_2 \in \mathcal{D}(\mathbf{s}_2)\},$$

$$\text{cap}_a(\mathbf{n}_1, \mathbf{n}_2, \mu) \triangleq \text{cap}(\mathbf{n}_1, \mathbf{n}_2, \mu_a) \cdot \text{sentskew}(\mathbf{s}_1, \mathbf{s}_2)$$

5.6 Alignment Graph Construction

As described in Section 2.3.1, a flow network models the flow of material through a graph using two values: capacities and flow. More formally, for flow network \mathcal{G} :

$$\mathcal{G} = (\mathcal{V}, \mathcal{E}) \quad \text{a flow network is a graph} \quad (5.10)$$

$$f : \mathcal{E} \rightarrow \mathbb{R} \quad \text{a s-t flow that assigns a flow}$$

$$\forall e \in \mathcal{E}, 0 \leq f(e) \leq c_e \quad \text{capacity constraint} \quad (5.11)$$

$$\forall v \in \mathcal{V} \quad \sum_{e \text{ into } v} f(e) = \sum_{e \text{ leaving } v} f(e) \quad \text{conservation of flow constraint} \quad (5.12)$$

where:

- \mathcal{V} is the set of vertexes
- \mathcal{E} the set of edges
- $f(e)$ is the amount of material flowing through edge e and c_e is a capacity

Conceptually, the capacity is the limit of water pumped through a pipe before it breaks from the pressure. The maximum amount of material (c_e) that can flow through edge e is formalized with the capacity constraint inequality Equation 5.11. Equation 5.12 states that all material going into a node must have the same amount of material leaving that node. This conservation of flow holds at a global level since material flows from source s through the network to t :

$$\begin{aligned} f^{\text{in}}(v) &= \sum_{e \text{ into } v} f(e) && \text{a node's input flow} \\ f^{\text{out}}(v) &= \sum_{e \text{ leaving } v} f(e) && \text{a node's output flow} \end{aligned} \quad (5.13)$$

$$\forall v \notin \{s, t\}, f^{\text{in}}(v) = f^{\text{out}}(v) \quad \text{node flow balance constraint} \quad (5.14)$$

$$\sum_{e \text{ out of } s} f(e_s) = f^{\text{out}}(s) = \sum_{e \text{ into } t} f(e_t) = f^{\text{in}}(t) \quad \text{s-t conservation of flow} \quad (5.15)$$

$$v(f) \triangleq f^{\text{out}}(s) \quad \text{value of flow definition} \quad (5.16)$$

where:

- e_s are the edges connected to the flow network source node s
- e_t are the edges connected to the flow network terminal sink node t
- $v(f)$ the value of flow

The value of flow provided in the graph is given in Equation 5.14, which is the sum of flow leaving the source s constrained by the capacities (c_e) of edges connected to it. Since the amount of material that starts at s is equal to what arrives at t their values of flow must also

be equal per Equation Equation 5.15. The value of flow capable of flowing through the graph \mathcal{G} from source s to sink t is defined in Equation Equation 5.16. Finding the flow value across all edges of the graph to satisfy facilitating the movement of the material is known as the max flow problem (Ford and Fulkerson, 1962), which is a well known problem that has been solved with many well known, tried, and proven methods that have seen great performance increases over the last six decades (Gao et al., 2022b).

In this work, the flow material is information through AMR graphs that have been parsed and constructed from sentences. The AMR sentence graphs are connected with a root node into a graph, one for the source and another for the summary. These two initially disconnected graphs are called the source component and summary component, but then are connected as a bipartite graph. This connected graph is called the alignment graph, and the process of creating it (as a flow network) is called the alignment graph construction phase. Note that the alignment edges of the alignment graph are not to be confused with the text-to-graph aligned tokens (TATs), which align a token’s index in the sentence to a node in the AMR graph. Figure 14 shows two example AMR graphs with their TATs as blue arrows before the alignment graph is constructed.

5.6.1 Graph Component Construction

The first step to creating the graph is to perform Coreference Resolution, which is the task of finding all mentions referring to the same named entity. This is done by identifying the antecedent tokens of a mention, then clustering them by entities they reference. The

Coreference Resolution method uses the `amr_coref`¹ library, which uses the general approach followed by `neuralcoref`² as in previous work (Blodgett and Schneider, 2021; Szubert et al., 2020).

In AMR, coreferenced mentions are clustered across the source and summary components so each entity reference is associated with the set of nodes to the entity referred. This has the effect of connecting and renaming the variables to be identical for those referred nodes to create reentrancies. It is important to note the surface name of the context node and aligned textual tokens might be different in cases of pronominal reference such as “*Obama* was the 44th president...*he* served as...”. A sentence node for each AMR graph and a root node that ties all sentences together are also created as shown in Figure 19. Additional nodes between the sentence and root node are added for more complex summarization scenarios, which include paragraphs, document sections³, and documents as nodes. All nodes from the sentence to the root node are referred to as document nodes.

During graph construction, parsed sentences are associated with sentence nodes, and TATs are associated with concept and attribute nodes (see Figure 14). Concept nodes are also populated with their PropBank rolesets (see Section 5.1.1). Outgoing role edges are populated with their connected concepts using the database generated from the preprocessing step described

¹https://github.com/bjascob/amr_coref

²<https://github.com/huggingface/neuralcoref>

³Medical documents are highly sectioned having, for example “History of Present Illness”

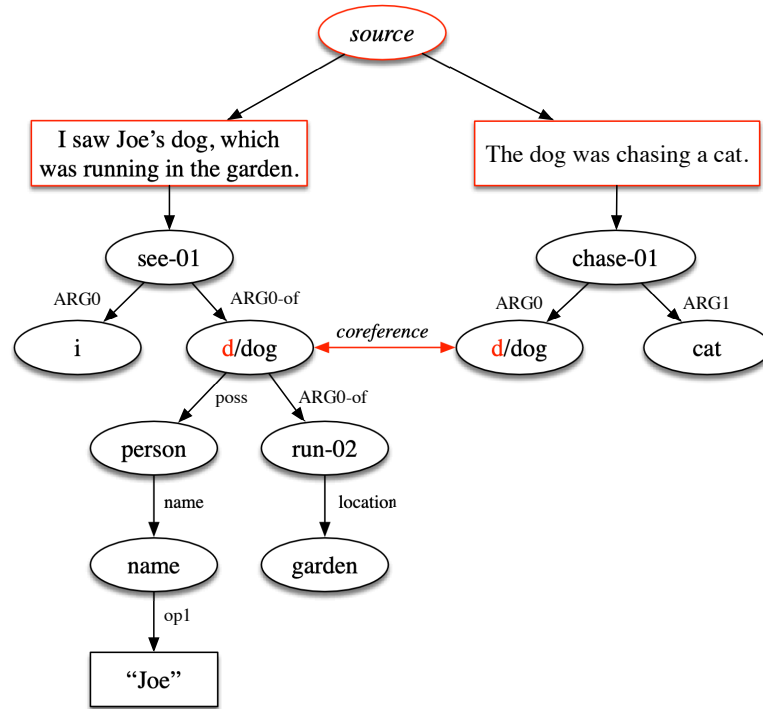


Figure 19: **Graph Construction.** The construction of the summary AMR graph with the changes shown in red. Example from Liu et al. (2015).

in Section 5.4.1. At this point in the construction, all nodes and edges have at least one data structure with embeddings including sentences from their sentence SBERT embedding (Reimers and Gurevych, 2019), concept nodes and edges with their roleset embeddings, and attribute nodes with their text-to-graph aligned token embeddings. An SBERT sentinel embeddings is generated from the label of nodes that have no textual alignments.

Sentence embeddings use the [CLS] token’s “pooled output”, which represents the entire sentence’s embedding. On the other hand, the concept and attribute node embeddings are

computed as the mean of the model’s final output layer of all the wordpieces (Wu et al., 2016) of all TATs¹. This sentence vs. token strategy is conducive to SBERT since the model was trained for clustering and semantic similarity. See Section 5.4 for more detail.

5.6.2 Graph Component Connection

The source and summary components are then connected with the alignment edges to create the bipartite graph. They are created as the Cartesian product of the concept nodes across the components with their capacity (the upper bound constraint on the flow of material through an edge of a flow network graph) values set to the semantic similarity for each bipartite node pair (see Section 5.6). All document nodes across the components are also connected, which are just the root **source** and **summary** nodes in the examples provided in this section. Capacities for all alignment edges are assigned as described in Section 5.5 with edges discarded for those that fall under the similarity threshold (τ_δ) hyperparameter. The edges are then reversed to allow the flow network configuration necessary for alignment assignments as described in Section 5.7. The resulting graph is shown in Figure 20.

5.6.3 Complete Flow Network

The last step of the alignment graph construction phase is to add the source and sink nodes and connect them to the graph. The source node² (**s**) is connected to the summary’s AMR graph leaf nodes and the sink **t** is connected to the source component root node. The role edges

¹An AMR node may represent zero, one, or many tokens.

²The terminology difference between the source node **s** (flow network) and the source component (AMR graph) is explicitly differentiated since the source node can be connected to the source component.

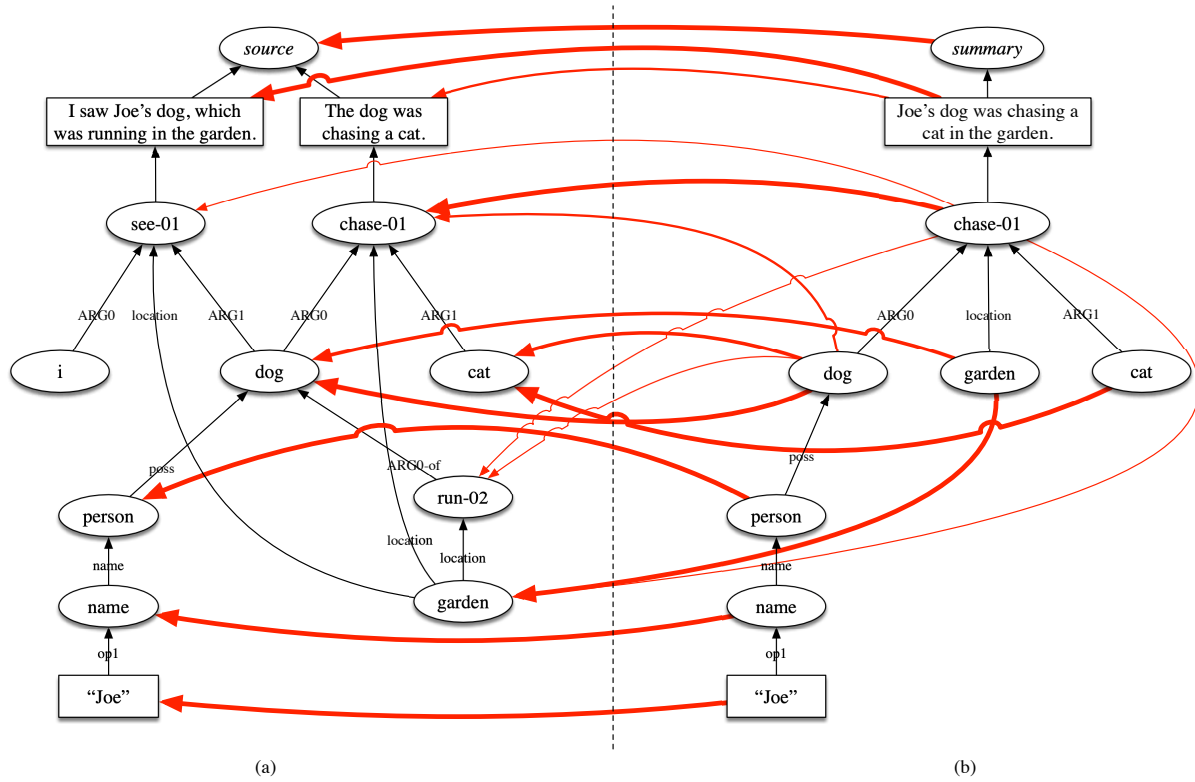


Figure 20: **Bipartite Graph Connection.** The alignment edges added to the bipartite source and summary components with the changes shown in red. The width of the edge represents the value of its capacity as a function of its similarity. Example from Liu et al. (2015).

of the AMR graph that relate the roles of concept nodes also become capacitated edges with values set to infinity. The resulting flow network graph, shown in Figure 21, can be described and constrained by equations Equation 5.10 through Equation 5.12 and used with the max flow algorithm to compute the flow through all edges.

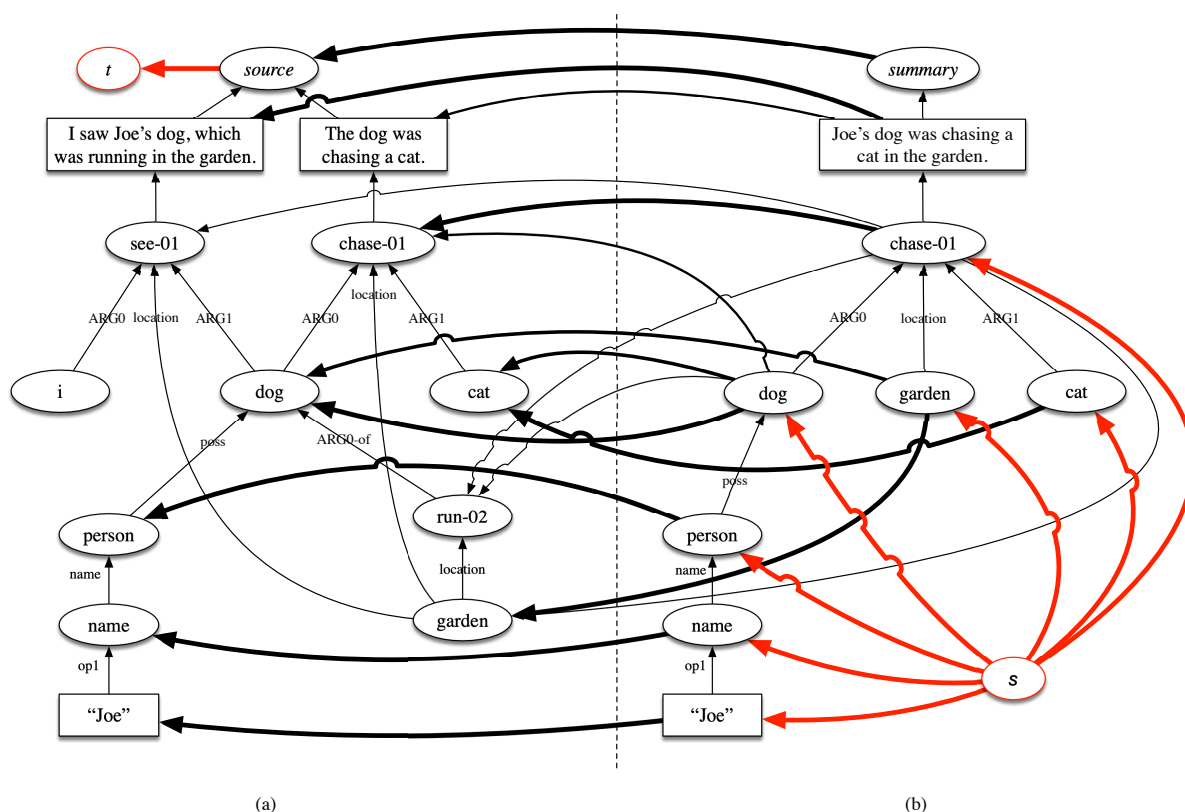


Figure 21: **Flow Network Construction.** The source (s) and sink (t) nodes and connected edges added to the bipartite graph with the changes shown in red. The addition of the terminal nodes completes the flow network. Example from Liu et al. (2015).

5.7 Alignment Graph Algorithm

After the graph is created, the alignment graph algorithm is used to modify the alignment capacities between the concept nodes of source and summary components. A byproduct of the algorithm is a score assigned to each role edge of how well the respective branch of the component is aligned across the bipartite graph. The initial capacities that are set on the

constructed alignment edges (see Section 5.6.2) provide an initial estimate of concept node pair alignments. However, these capacity values only take in to account the local node neighborhood at the sentence, or even document level. Instead, the flow values assigned by the max flow algorithm provide a global representation with a better estimate of the alignment strength. The algorithm goes a step further and iteratively updates capacities informed by flow changes by re-running the max flow algorithm until convergence.

5.7.1 Max Flow

The algorithm starts by computing the max flow, which pushes flow from the flow network's source s node to AMR summary component. It then flows to the alignment edges from the summary to the source component, up toward the root of the source component, and out to the flow network's sink node t . Since the role edges of both graphs are set to infinity, flow goes through to every node and edges of the graph¹ as shown in Figure 20. I used the igraph² implementation of the push-relabel max flow algorithm (Goldberg and Tarjan, 1988) for my experiments. While not mathematically proven for the non-integral case, the algorithm works on experiments with 100+ test graphs.

¹Complete flow saturation occurs in some instances with nodes that have more than one parent. See Section 5.11.

²<https://igraph.org>

5.7.2 Flow Normalization

The algorithm's graph at this point has flow values on each edge, which are then normalized for each node based on the descendants of that node. First, terminology regarding graph semantics must be defined.

Let:

- The edge node parent $\mathcal{P}(e)$ be the node at the source end of a directed edge e in the alignment graph. For example, in Figure 22 (a), the edge node parent of edge **poss** is **dog**.
- The edge node child $\mathcal{C}(e)$ be the node at the target end of a directed edge e .
- The node descendants $\mathcal{D}(v)$ of v be all paths to the terminal leaf nodes (grandchildren). For example, in Figure 22 (a), the $\mathcal{D}(\text{dog})$ are the nodes **person**, **name**, and attribute **Joe**.
- The flow per node $\tilde{f}(e)$ be the flow at an edge divided by all its source node's descendants:

$$\tilde{f}(e) = \frac{f(e)}{|\mathcal{D}(\mathcal{P}(e))|}$$

where $f(e)$ is the amount of flow of material through edge e .

Figure 22 (b) shows the values from a run of the algorithm on the example from Figure 20 with flows (left of arrow) and flow per node values (right of arrow). For example, a query of the flow that leaves **person**, which is one of the role edges labeled **:poss**, starts with traversing node $\mathcal{P}(\text{poss}) = \text{dog}$ with $\mathcal{D}(\mathcal{P}(\text{dog})) = \{\text{dog}, \text{name}, \text{Joe}\}$. Note we use the forward flow graph

for gathering descendants and fetching the edge node parent. The flow on `:poss` is 2.41 so our flow per node = $\frac{2.41}{3} = 0.8$.

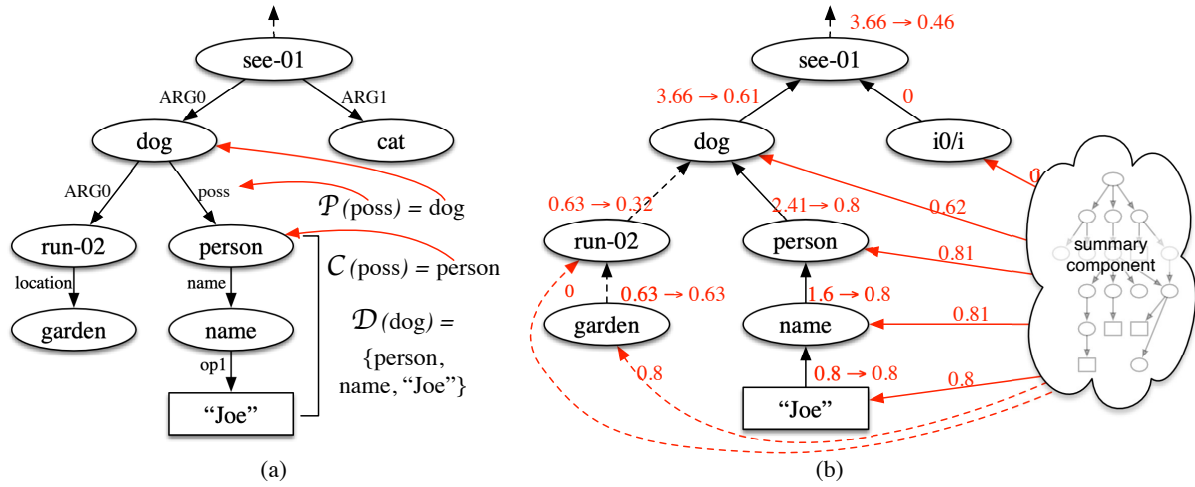


Figure 22: **Network Flow.** Left (a): the edge node parent of edge `poss` is `dog` and its edge node child is `person`. Right (b): example flow from the summary to the source component with flows (left of arrow) and normalized flow per node values (right of arrow) with constricted edges as dotted lines. Example from Liu et al. (2015).

5.7.3 Capacity Constriction

After normalizing the max flow, the next step is to “squeeze” capacities of the alignment edges in the network that will affect the next iteration, which is detailed in Section 5.7.4. The edge capacities are set to zero if they fall under an alignment edge minimum capacity cutoff

(τ_α) , which is tuned as a hyperparameter of the model.¹ In addition, all of the edge node parent's descendant's role edge capacities that fall under the role edge minimum capacity cutoff (τ_ρ) are set to zero. More specifically,

$$\{\forall e \in \mathcal{E} : f(e) < \tau_\rho\} \mathcal{D}(\mathcal{P}(e)) \leftarrow 0$$

where:

- e is the component role edge
- $f(e)$ edge e 's flow value
- $\mathcal{P}(e)$ is the edge node parent
- $\mathcal{D}(\mathcal{P}(e))$ are the descendant nodes of the edge node parent (see Section 5.7.2 for definitions)

For example, if the τ_ρ is set to 0.4, the role edge that connects **run-02** and **dog** in Figure 22 (b) would meet the criteria for its edge node parent and descendants to be set to zero. The path of role edges affected in the figure are represented with dotted lines. Since $\mathcal{P}(e) = \text{dog}$, and $\mathcal{D}(\text{dog}) = \{\text{run-02}, \text{garden}\}$ then the capacities of the alignment edges between the summary component and **run-02** (already set to 0) and **garden** (from 0.8) will be set to 0 (represented with dotted lines). Note that even though there are capacity edges well over the value of the

¹This is not to be confused with edge creation cutoff (τ_δ) for creating component alignment edges during graph creation used in Section 5.6.2.

τ_p , we still set them to zero since each edge is evaluated by the flow per node value that takes into account all flows of descendants at that level.

5.7.4 Summary Max Flow

The non-zero alignment edges remaining after finishing the steps described in sections 5.7.1 – 5.7.3 explain how much (and what part of) the source is included in the summary. Likewise, these steps are repeated on a second graph with its output explaining how much of the source is summarized. This second graph has shared capacities with the first, and alignment edges reversed so it flows from the summary component to the source component. The source flow node s is connected to all of the leaf nodes of the source component and the sink t node is moved to the summary root node.

Before the max flow algorithm is run on the second graph instance, the flow values on all role edges are tracked, and if any change, the steps described in Section 5.7.1 through Section 5.7.3 are run again. The algorithm’s execution is alternated across both graph instances until it converges when no flow value changes. Given the symmetry of the algorithm across both components, there is no need to distinguish them as the source and summary. A harmonic mean of their flow values is useful for measuring semantic similarity for text-to-graph evaluation and individually for summarization scoring and matching. This harmonic mean can not only be used as a scoring method for two AMR graphs, but also provides flow values both globally and locally that inform where the summarized content is found in the source text.

5.7.5 Final Alignment Graph

In summary, the alignment graph is made up of the source and summary components composed of AMR graphs, which are then connected to create the flow network. The alignment edges are then deleted by clamping shared capacities across two flow network instances (each with a reversed flow). Low flow edges result in minimized or deleted alignment edges for iteration of the algorithm.

The final alignment graph looks similar to Figure 21 before the alignment edges are deleted. For example, `chase-01` to `run-02` nodes for $\tau_\delta = 0.3$.

5.8 Scoring Method

My scoring method, Component ALignment for Abstract Meaning Representation (CALAMR), measures AMR graph semantic similarity. This scoring method can be used as a similarity metric similar to SMATCH, which is an evaluation metric for abstract meaning representation similarity that uses a greedy feature overlap method (Cai and Knight, 2013). However, CALAMR includes metrics for summarization overlap. We define the value of flow exiting the source node to the sink as the *source root flow* using Equation Equation 5.16 with $\mathbf{C}_{fc} \triangleq f^{\text{out}}(s_{\text{source}})$.

This metric applies for every subgraph or globally as the source component's root node connected edge. For example, this value is 0.46 in Figure 22b for the subgraph from the `see-01` node to the leaf nodes. Likewise, the value of flow exiting the summary node to the sink is defined as the *summary root flow* with $\mathbf{C}_{fy} \triangleq f^{\text{out}}(s_{\text{summary}})$. Additional CALAMR scoring methods include the portion of nodes in the source component that have at least one alignment with the summary, defined as the *source aligned portion* ($\tilde{\mathbf{C}}_c$) and the portion of nodes in the

summary component that have at least one alignment with the source defined as the *summary aligned portion* ($\tilde{\mathbf{C}}_{\mathbf{y}}$).

Given the symmetry of the algorithm as detailed in Section 5.7.4, we can treat the source as any sentence (and likewise the summary) if we wish to score them in like fashion to previous baseline scoring methods (Cai and Knight, 2013; Opitz et al., 2021); meaning we can score the similarity of two AMRs.

Unlike previous methods, we can also score AMR graphs as multi-sentence graphs by the summarization overlap as a non-negative real value called the *aggregate flow*:

$$\mathbf{C}_f = 2 \frac{\mathbf{C}_{fc} \cdot \mathbf{C}_{fy}}{\mathbf{C}_{fc} + \mathbf{C}_{fy}}, \quad (5.17)$$

which results in higher values for graphs with multiple node alignments. This score is useful for subgraphs but not globally, so we define the *aggregate alignment portion* score as the harmonic mean of the aligned node portion across components:

$$\tilde{\mathbf{C}} = 2 \frac{\tilde{\mathbf{C}}_{\mathbf{c}} \cdot \tilde{\mathbf{C}}_{\mathbf{y}}}{\tilde{\mathbf{C}}_{\mathbf{c}} + \tilde{\mathbf{C}}_{\mathbf{y}}}, \quad (5.18)$$

which is also a real value but has range $[0, 1]$, and is advantageous as it has the same range as established AMR scores such as SMATCH.

5.9 Experiment Design and Setup

We report two kinds of experiments with the first concerning summarization (see Section 5.9.1) and the second similarity scoring (see Section 5.9.2) between human annotated AMRs

and the output of three popular parsers. We run these evaluations to assess how CALAMR (my score) compares to other methods of finding summarized content and judging alignment as a similarity measure. Given the results reported in Section 5.10, I believe CALAMR is a more perspicuous score as it judges semantic similarity both locally and globally. The second set of experiments use the overlap of two multi-sentence AMR graphs as a measure of summarization.

5.9.1 Summarization

The human annotated “Proxy Report” LDC2020T02 AMR Annotation Release 3.0 (Knight et al., 2021) corpus contains news articles with sentences tagged as a date, country, topic, summary and body. Only sentences tagged with summary or body are used in our experiments. The corpus development set has 35 articles and 826 sentences, and its test set has 33 articles and 823 sentences. The alignment graph algorithm was first used to find summarized text through text-to-graph alignments. It was then used to score summarizations in development and test sets.

After scoring, sentences of the source and summary for 33 articles of the two data sets were swapped and scored a second time (we call this the “Mismatch set”). This was accomplished by switching the source sentences from the test set with the summary sentences in the development set. Source sentences in the Proxy Report corpus resulted in 14,108 role edges, 660 concept nodes and the summary totaled 1,153 role edges and 802 concept nodes.

5.9.2 Similarity Scoring

In addition to summarization experiment using the Proxy Report corpus, an AMR graph similarity experiment used parser output to compare CALAMR with previous scoring meth-

ods. JAMR (Flanigan et al., 2014) and two other popular and peer-reviewed parsers were used to generate AMR graphs from natural language sentences. They were then scored with SMATCH (Cai and Knight, 2013), WLK (Opitz et al., 2021)¹ and the my scoring method.

The parsers’ output was scored against the first 20 sentences of the AMR 3.0 Proxy Report corpus and two Information Science Institute² corpora, which include *Little Prince* (1,562 sentences) and *Biomedical* (6,952 sentences). These corpora were selected based on their availability and variety of domain. The amrlib³ library was used with the SPRING (Bevilacqua et al., 2021) and Gsii (Cai and Lam, 2020) parsers as it has been shown to produce good results on text-to-graph tasks (Heinecke and Shimorina, 2022; Opitz and Frank, 2022; Opitz et al., 2021).

5.9.3 Baselines

My primary baseline uses the Liu et al. (2015) graph source-to-summary reduction method of comparing the text-to-graph aligned words from the source to the summary as a bag of words. ROUGE, which is a measure that assesses the quality of generated summaries, is used on the unigrams to evaluate the aligned overlap to gauge how well CALAMR “finds” summarized content.

I also compare the resulting CALAMR edge coverage of alignment outputs with the document-level AMR graph heuristic method that links based on exact match of Liu et al. (2015) on the

¹All experiments used WLK settings of $K = 2$ iterations with “all directional communication”.

²<https://amr.isi.edu/download.html>

³<https://github.com/bjascob/amrlib>

LDC2014T12 AMR Annotation Release 1.0 corpus as a baseline¹. We cannot directly compare coverage as their method is a graph reduction into the summary and mine is an alignment method. Instead, I compare CALAMR positive value flow role edges to their sentence level graph expansion since both methods’ goal is to identify AMR subgraphs used for alignment. Both of my metrics are taken as a quotient of the total number of summary component role edges.

5.10 Results

Section 5.10.1 pertains to CALAMR as a summarization method with respect to scoring, Section 5.10.2 explains the measure of summarization between two AMR graphs, and my results are given in Section 5.10.3.

5.10.1 Summarization

The differences in alignment across the corpora is significant and strongly indicates the method’s effectiveness. The scoring method captures the rate of summarization as evidenced with 86.6% average aligned AMR nodes in the unaltered document set (Proxy Report) compared to the 35.1% in the development and test switched (Mismatch) set shown in Table XII. The source component aligned portion also reveals the difference between the unaltered and Mismatch document set as the portion of the summary is represented in the source (43.2% vs. 14.6%). The high summary root flow (C_{fy}) of 7.21 (see Table XII) in the unaltered set indicates the source graph is strongly aligned compared to the low summary root flow of 2.61

¹The AMR 3.0 corpus was used for other experiments.

in the Mismatch set. These flow metrics represent the degree to which each graph is aligned with the other.

Corpus	$\tilde{\mathbf{C}}_y$	$\tilde{\mathbf{C}}_c$	\mathbf{C}_{fy}	\mathbf{C}_{fc}
Proxy report	86.6%	43.2%	7.21	0.67
Mismatch	35.1%	14.6%	2.61	0.20

TABLE XII: **Document Summarization Scores.** Scoring matched vs. mismatch corpus. See Section 5.8 for CALAMR scoring notation.

5.10.2 Alignment

As explained in Section 5.8, the particularity of the source and summary becomes moot given the symmetry of the algorithm across the components, which motivates the semantic similarity scoring method. We hypothesize that the aggregate alignment portion CALAMR $\tilde{\mathbf{C}}$ score reflects the similarity between any two AMR graphs, and thus, is an indicator of the effectiveness of the score for summarization.

To estimate this effectiveness we compared the text-to-graph parser output of AMR sentence metrics with previous methods as a reference point. This was done by comparing the scoring methods with the aggregate alignment portion ($\tilde{\mathbf{C}}$ from Equation Equation 5.18) using Pearson (ρ) correlations. Table XIII lists the scores between the human annotated AMR sentences and

parser output. We find highly correlated scores between WLK (69.2) and SMATCH (67.7) using JAMR on the Little Prince corpus.

To that end, JAMR produces the highest correlations for the other corpora as well. In fact, this trend extends to SPRING (second most correlated) and Gsii (least correlated). Because the parser is a consistent indicator of highly positively correlated results, regardless of corpus, I conclude that CALAMR is an effective scoring method.

Corpus	Parser	Sent	$\rho \tilde{C}, S$	$\rho \tilde{C}, W$
Biomedical	Gsii	6,644	0.412	0.318
Biomedical	Jamr	5,612	0.662	0.652
Biomedical	Spring	6,617	0.501	0.413
Little prince	Gsii	1,464	0.388	0.357
Little prince	Jamr	1,501	0.677	0.692
Little prince	Spring	1,497	0.413	0.471
Proxy report	Gsii	8,042	0.229	0.308
Proxy report	Jamr	7,781	0.532	0.562
Proxy report	Spring	8,120	0.373	0.482

TABLE XIII: **Parser Alignment Scoring.** AMR Sentence Pearson correlations (ρ) between aggregate alignment portion (\tilde{C})_{ALAMR} (see Equation Equation 5.18) and previous scoring methods (S)MATCH and (W)LK. Metrics are reported only for successfully parsed (Sent)ences.

5.10.3 Previous Methods

Even though CALAMR is not a summarizer, previous AMR summarization methods' edge coverage were compared with alignment edge counts. Table XIV shows the alignment results on

the LDC2014T12 corpus using the same methodology as Liu et al. (2015). The high ROUGE1 score (F1 of 17.5 over the baseline), shows that CALAMR is effective at finding summarized content.

Method	Precision	Recall	F1
Liu et al., 2015	51.9%	39.0%	44.3%
Dohare et al., 2017	52.4%	55.7%	51.3%
Fu et al., 2021	-	-	49.1%
CALAMR	69.0%	68.6%	68.8%

TABLE XIV: **Aligned Node Text Comparison.** Unigram (bag of words) aligned source to summary overlap of text-to-graph tokens.

The Liu et al. (2015) method is compared with role and alignment edge coverage in Table XV as explained in Section 5.9.3. The results show a significantly higher coverage of edges from the expanded sentence-level edges to CALAMR component role edges. CALAMR also has a much higher coverage of summary component alignment with the source compared to the previous method’s document-level expanded edges. This is a strong indicator that my method is a better comparatively.

5.11 Reentrancies

AMR nodes with multiple parents (reentrancies) lead to catastrophic alignment failure. Because the role edge’s capacities are set to infinity, the max flow algorithm redirects all flow

Split	Sent	Doc	$\tilde{\mathbf{C}}$ Role	$\tilde{\mathbf{C}}$ Alignment
Train	75.5%	84.6%	84.3%	94.8%
Dev.	85.4%	91.8%	83.1%	92.5%
Test	75.0%	83.3%	84.3%	96.3%

TABLE XV: **Summary Alignment Coverage.** Edge coverage as a portion within (Sent)ences and (Doc)uments (Liu et al., 2015), and the portion of non-zero flow Role and Alignment edges using the $(\tilde{\mathbf{C}})_{\text{ALAMR}}$ method.

through only one of the parent’s incoming edges starving the other paths from the reentrancy to the root as shown in Figure 23.

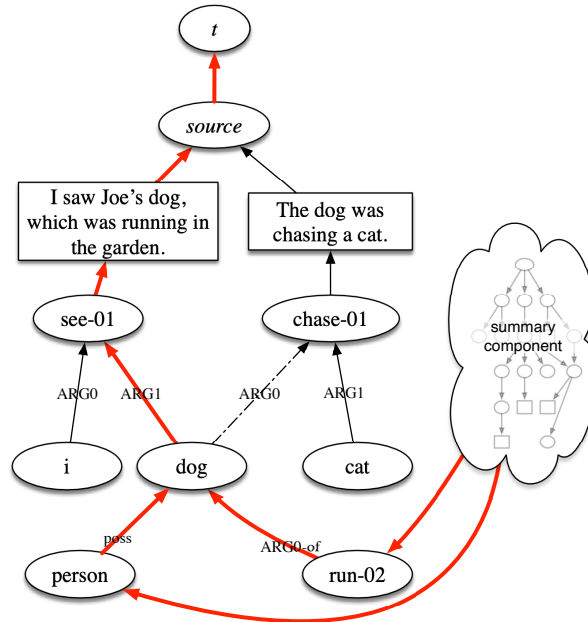


Figure 23: **Reentrancies Alignment Failure.** An example of a catastrophic alignment failure due to a reentrancy (dog) that has two parents (see-01 and chase-01).

I observed 2.6% of the source graphs nodes and 2.18% of the summary graphs nodes to be reentrancies. To address this issue I clamped the capacities of all incoming edges of reentrancy nodes to the normalized sum of the incoming flow so that:

$$\forall e \in \mathcal{P}(\mathbf{n}), c_e \leftarrow \frac{1}{|\mathcal{P}(\mathbf{n})|} \sum_{\mathbf{u} \in \mathcal{P}(\mathbf{n})} f(\mathbf{u}) \quad (5.19)$$

where:

- c_e is the capacity value
- $\mathcal{P}(\mathbf{n})$ are incoming edges to \mathbf{n}
- $f(\mathbf{u})$ is the flow through e into \mathbf{n}

This resulted in all of the summary graphs being repaired (all reentrancy positive edge flows) and 80.3% repaired in the source graphs. This method fixed three catastrophic alignment failures of the 366 Proxy Report documents, but one (2.7%) remained unchanged (see Table XVI).

Component	Total	Repaired	Portion Fixed
source	151,625	4,205	2.8%
summary	14,213	365	2.6%

TABLE XVI: **Reentrancy Repair Statistics.** The statistics on reentrancies that were repaired across components.

5.12 Conclusion

CALAMR is an alignment method that supports scoring metrics that provide a semantic similarity metric for multi-sentence AMR graphs and subgraph level summarization metrics. The method is suitable as both a scoring method capable of determining the portion of overlapping content through alignment and flow metrics. The utilization of determining the overlap provides a method of creating faithful and traceable summaries as demonstrated in the next chapter.

5.13 Alignment Graph Examples

Graphs were exported using the reference implementation of the alignment method (see Section 5.3). A rendering of the graph after each phase of the alignment method is given in the subsequent figures.

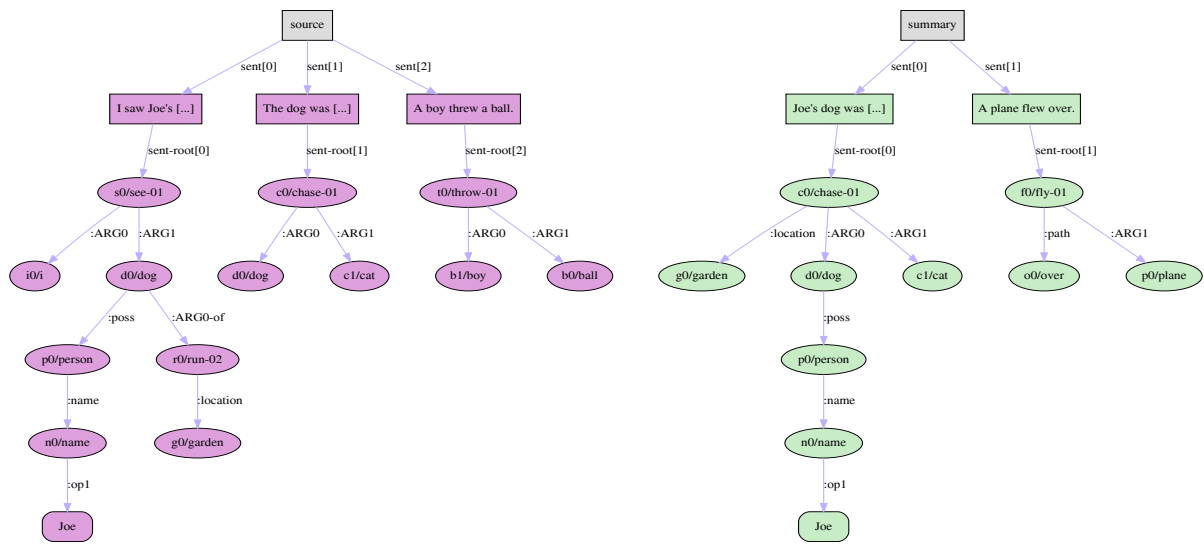


Figure 24: **Disconnected Source and Summary AMR Graphs.** Left: the source component (purple nodes) including the sentences: a) “*I saw Joe’s dog, which was running in the garden.*”, b) “*The dog was chasing a cat.*”, c) “*A boy threw a ball.*”, and d) “*A plane flew over.*”. Examples b) and c) from Liu et al., 2015. Right: the summary component (green nodes) including the sentence “*Joe’s dog was chasing a cat in the garden.*” Each AMR is tied with sentence, then root notes (see Section 5.6.1). Note that the sentence “A plane flew over.” was added to show how the method finds content in the summary not found in the source. This illustrates how the method finds hallucination in LLMs.

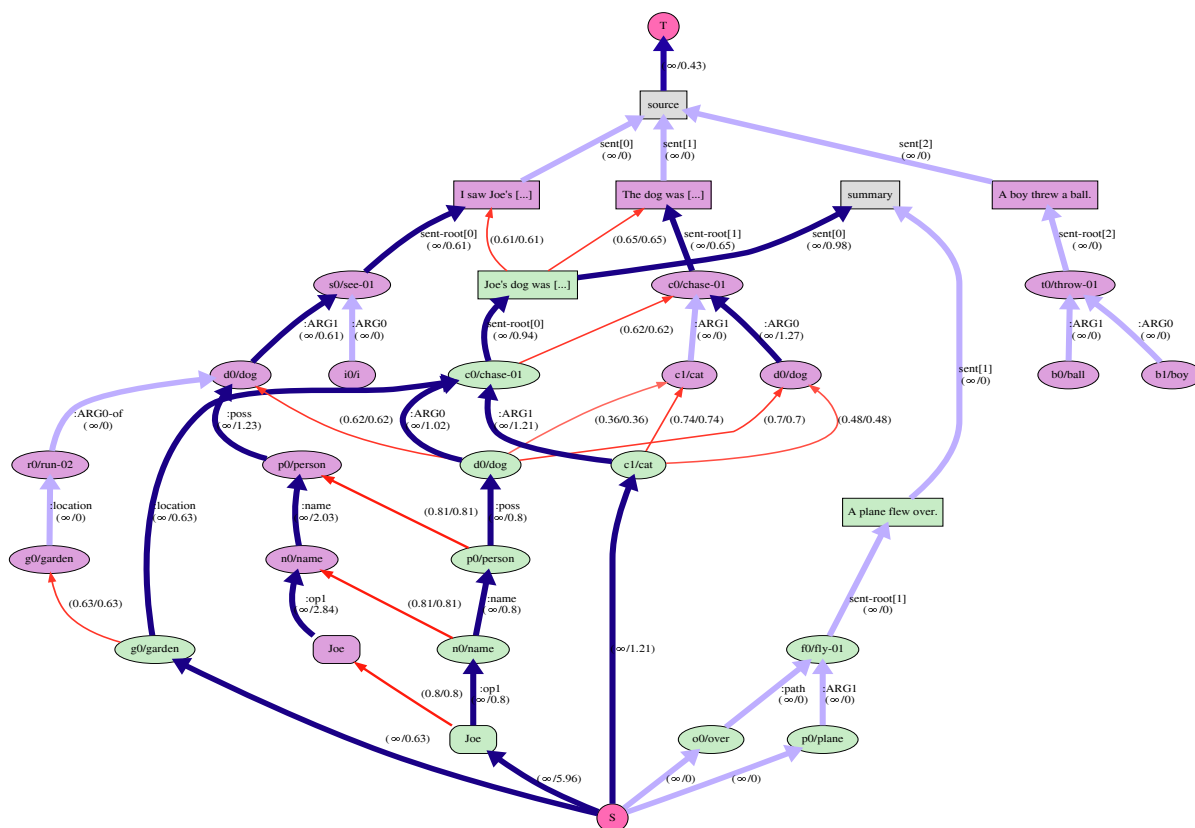


Figure 25: **Source Component Alignment.** The flow from the summary component (green nodes) to the source component (purple nodes) after capacity constriction (see Section 5.7.3). The role edges are in blue and the alignment edges are in red with all edges' width representing the capacity. The flow is represented by the darkness of the edges' color. The capacity and flow is shown in parenthesis.

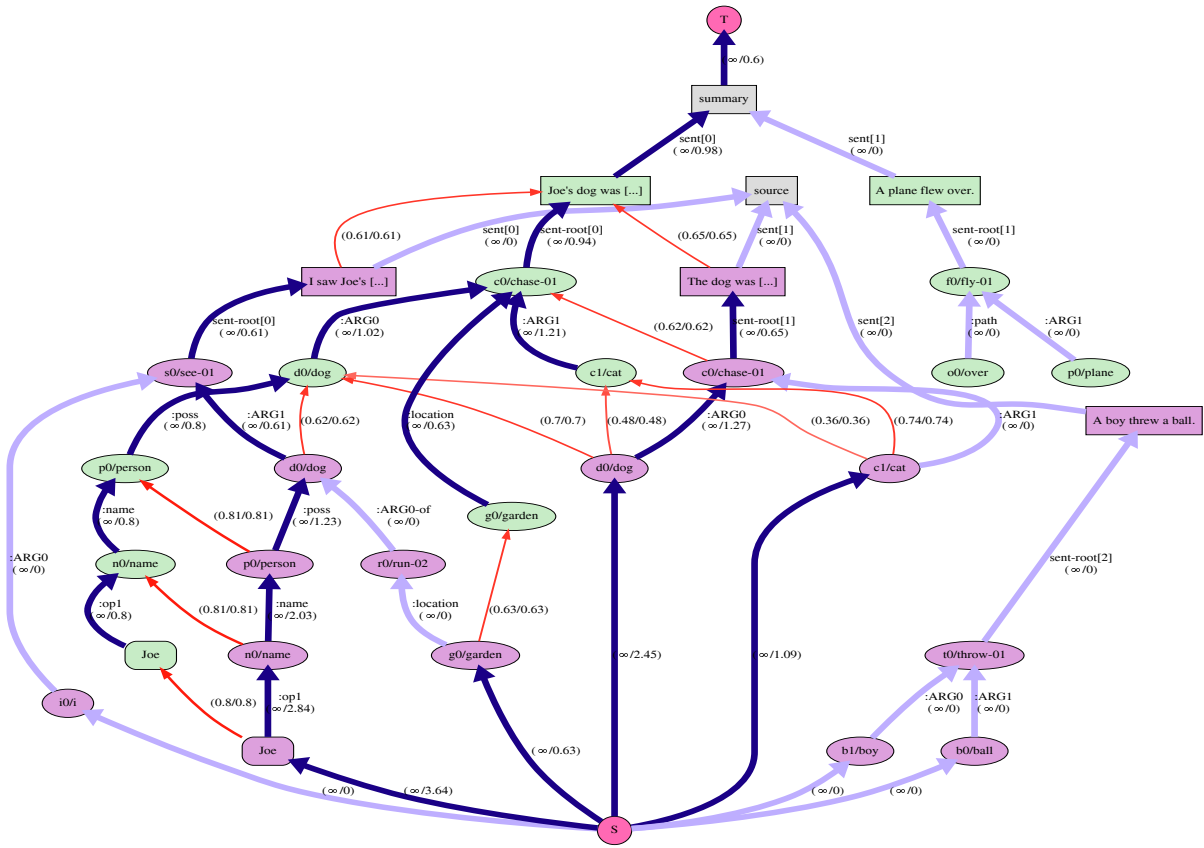


Figure 26: **Summary Component Alignment.** The alternated graph with flow from the summary to the source. Like Figure 27, this is taken after the capacity constriction step in Section 5.7.3.

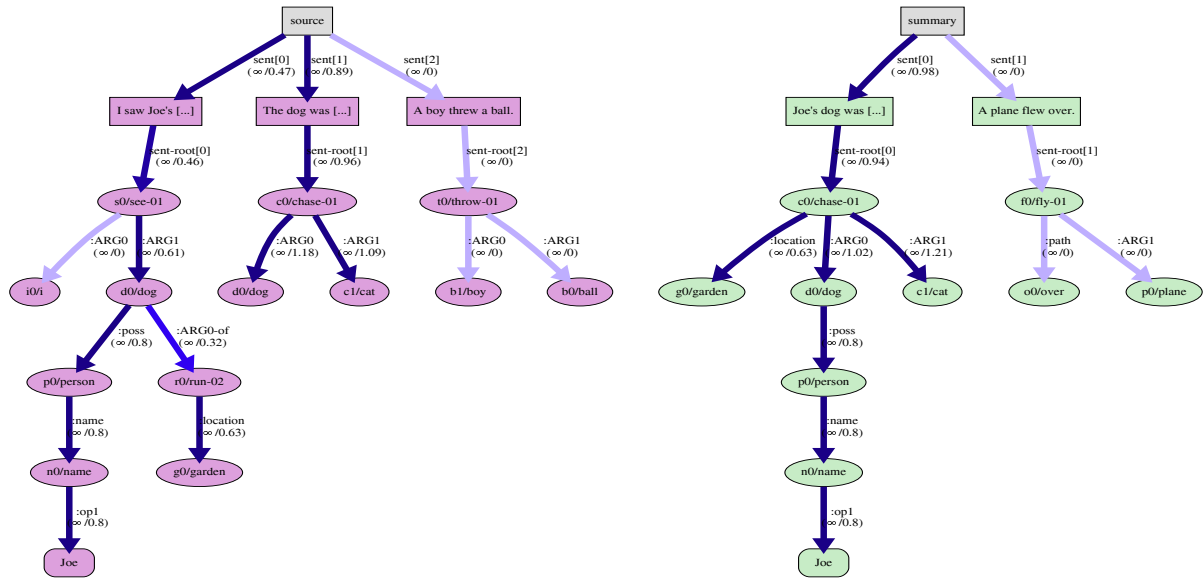


Figure 27: **Final Flow**. The final flow of both components rendered without the capacity edges in the same format as the disconnected components in Figure 24. The sentence “*A boy threw a ball.*” has zero flow, and thus indicates this entire sentence is not summarized. Similarly, the *i0/i* node has no flow so it will also not be summarized. The light blue role edge incoming to *r0/run-02* node has less flow, so its partial summarization follows from the “garden” mention in the summary, but “running” is missing. The solid blue in all edges of the summary component indicates the entire summarization is represented in the source. The fact that the *c1/cat* has a higher flow value than one means it has higher presentation in the source and has more than one alignment.

CHAPTER 6

SUMMARIZATION

The MedSecId model given in Chapter 3 gives a way of separating topical sections of medical notes and the DSProv feasibility study from Chapter 4 gives the upper bound of discharge summarization performance. With the alignment graph algorithm from Chapter 5, we now have an informed strategy, section identification model, and alignment method to match complex clinical sectioned note antecedents with discharge summaries to formulate a summarization method.

Clinical summarizations must be both faithful (how accurate the summary is) and traceable (if the summary can be traced back to its source content). Generating discharge summaries with these requirements highlights the difficulty of the task for large admissions with electronic health record (EHR) notes that number in the hundreds¹. Contemporary state of the art (SoTA) methods, such as fine-tuning large language models (LLMs), render the task nearly impossible with the volume of information for patients with prolonged hospitalizations. In some cases LLMs might feasibly summarize clinical documentation on a per section basis. However, many of these sections easily extend past the limit of large input text open source LLMs, such as the

¹The largest admission has 1,233 notes in the Medical Information Mart for Intensive Care III (MIMIC-III) corpus.

4K limit of Longformer (Beltagy et al., 2020) tokens and the 8K limit of LLaMa 3 (Touvron et al., 2023)¹.

Because of these memory constraints, other solutions are needed for the large multi-document automatic summarization task of generating the discharge summary. Most abstractive methods inherently have a tendency to create less faithful summaries because they typically formulate text as a probability distribution over the vocabulary. They also provide no traceable means of cross-referencing the text of summarized documents. The chosen method for this work is extractive since it is both faithful and traceable, and thus, acceptable for the clinical domain.

While previous methods have shown success at summarizing a single section (Adams et al., 2021), to my knowledge, there is no peer-reviewed work that attempts to generate a complete discharge summary using EHR notes. This motivates the extractive methods formulated in this work and provides a baseline for future abstractive summarization.

This extractive method leverages Component ALignment for Abstract Meaning Representation (CALAMR), described in Chapter 5, to create supervised training examples for a new model: the source section model. The alignment graph is used to match note antecedent source sentences to discharge summary sentences. Each note antecedent source sentence is then assigned to the section of the matched discharge summary sentence. This is then used as the label in a classification neural network (NN) model.

¹One *Labs and Imaging* MedSecId annotated section contains 2,178 tokens.

In addition to utilizing the CALAMR alignments, this summarization method uses the Med-SecId (see Chapter 3) model to section notes and the DeepZensols framework (see Chapter 7) for the pipeline automation. However, the utilization of the DSProv dataset and span matching methods for summarization is left as a future work as discussed in Section 8.2.

An overview of the discharge summary automatic summarization end-to-end pipeline follows:

1. Construct an admission graph from a subset of MIMIC-III Version 1.4 (Johnson et al., 2016) admissions (see Section 6.2.1).
2. Relabel concept variables to be unique (see Section 6.2.2).
3. Use CALAMR to create an alignment graph for each admission (see Section 5.3).
4. Use the admission alignment graphs to create a corpus of matches of note antecedent with discharge summary sentences (see Section 6.2.4).
5. Train the source section model using the corpus created in step 4 (see Section 6.3.1).

6.1 Datasets

The MIMIC-III Version 1.4 (Johnson et al., 2016) corpus and the UI Health dataset (see Section 1.3) were used for all experimentation using the summarization methods described. Admissions from the MIMIC-III dataset were selected starting with those having the lowest note count. The UI Health dataset was also processed in order of the lowest count of notes on a per admission basis from admissions that had at least one of each note category.

Even though the MIMIC-III admission selection totaled 11,957 admissions, a smaller subset of 113 admissions was processed through the entire pipeline due to a variety of factors such as memory limitations, processing latency and runtime errors¹. Admissions were aligned starting with those having the lowest note antecedent count to ameliorate processing failures caused by computational and memory factors. Two of these admissions contained at least one note annotated from the MedSecId (see Section 3.2) dataset. Figure 28 shows a rendering of an aligned admission (see Section 6.2.1) with a single radiology clinical note.

Despite these factors a total of 3,520 of the 11,957 MIMIC-III admissions were aligned using CALAMR, and then used to create the source section dataset (see Section 6.2.4). Table XVII (left) shows the average number of alignments across note antecedent and discharge summary components and the average number of reentrancies per admission. The “alignable” statistics are nodes that are alignment candidates, such as concept and attribute notes. The “aligned” statistics are those nodes with alignment edges.

Aligning the UI Health dataset resulted in additional challenges as shown in Table XVII (right). The dataset has more notes across category types compared to MIMIC-III because the latter has only intensive care unit (ICU) notes available (see Section 4.2.1). The consequence of this more robust note variety is that admission note counts are much higher, and therefore, take much longer to align. There is also a higher risk of missed alignments since there is more risk of multiple reentrancies (more than one in a path from the reentrancy to the root), which

¹Some admissions failed to align due to data issues, such as missing discharge summaries.

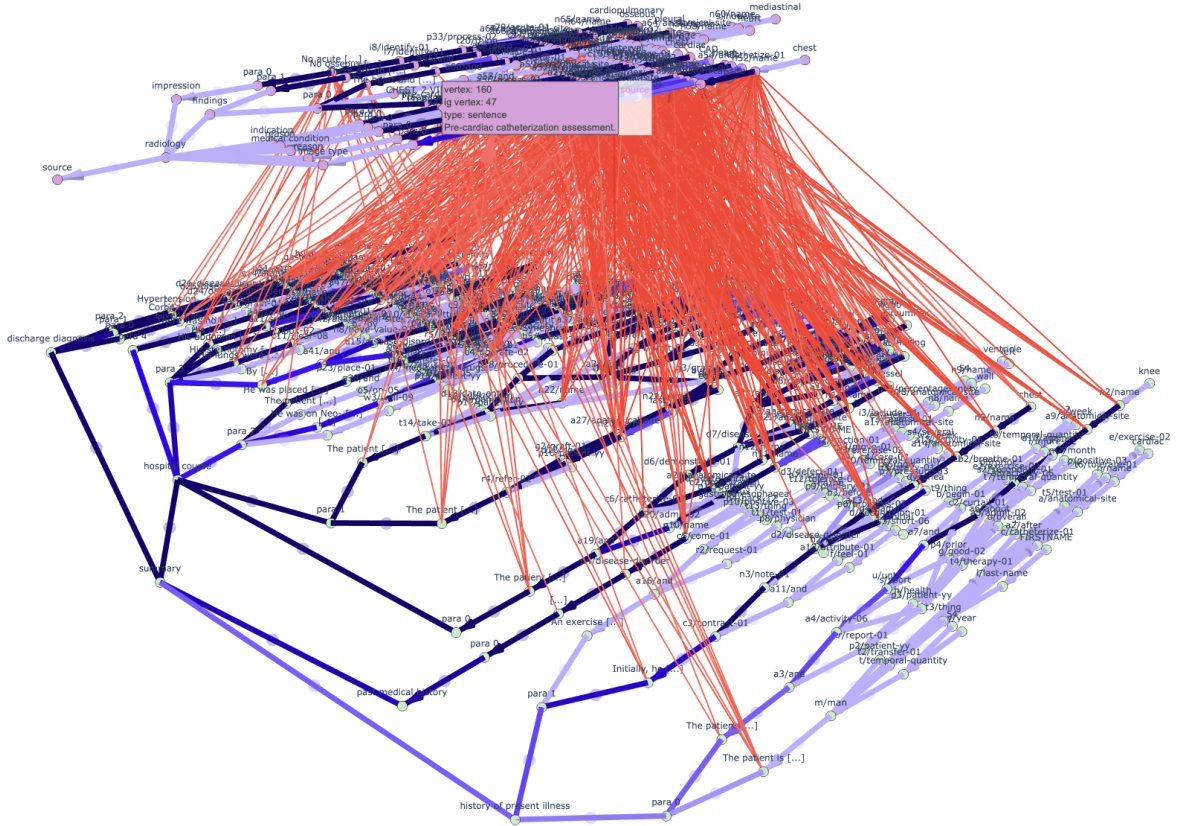


Figure 28: **Aligned Admission Graph.** A rendering of a CALAMR aligned admission graph with a single radiology node.

lead to flow issues (see Section 5.11). Even though the MIMIC-III alignments far outnumber the UI Health dataset, the UI Health dataset have many more reentrancies.

6.2 Methods

The unsupervised method described in Chapter 5 was leveraged to create supervised training examples using the flow data of the alignment graphs to match note antecedent source sentences

Description	Value	Description	Value
Average alignable in source	827	Average alignable in source	1,228
Average aligned in source	570	Average aligned in source	975
Average reentrancies in source	17	Average reentrancies in source	30
Average alignable in summary	819	Average alignable in summary	768
Average aligned in summary	482	Average aligned in summary	632
Average reentrancies in summary	34	Average reentrancies in summary	20
Total of graphs aligned	3,520	Total of graphs aligned	494

TABLE XVII: **Graph Alignment Statistics.** Alignment and reentrancy averages by admission with the number of nodes aligned by component in the admission graph. The MIMIC-III statistics are given on the left and the UI Health dataset on the right.

to discharge summary sentences. The source section model then used the matches to learn what to add to the summary. Each note antecedent source sentence was then assigned the section of the matched discharge summary sentence, which was used as the label in the source section model (see Section 6.2.3).

6.2.1 Admission Graph

The admission graph is a semantic representation of a patient’s hospital stay. It is composed of two disconnected components: the note antecedent source graph and the sections of the discharge summary. These two disconnected components follow the structure of the source and summary components of the bipartite graph described in Chapter 5. However, document nodes that represent note categories, note sections and clinical text paragraphs are used between the roots and their respective AMR (abstract meaning representation) subgraphs as shown in Figure 29. The note antecedent source has a note category level, whereas the summary component’s root represents the discharge summary.

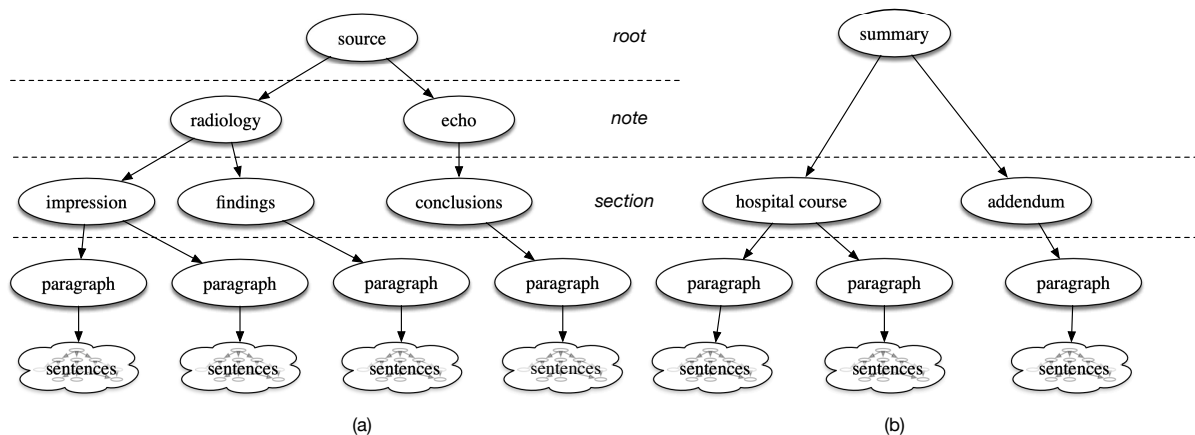


Figure 29: **Admission Graph.** An admission graph of the note antecedents (a), and the discharge summary (b).

The spaCy¹ library was used to tokenize, sentence chunk and tag biomedical and non-biomedical named entities, which are available as graph aligned features in downstream models. MedCAT (Kraljevic et al., 2021) was used to link token spans to concept unique identifiers (CUIs) (see page 44 for a definition), which were used to aid in graph aligning their text-to-graph concept nodes in the CALAMR alignment capacity calculation (see Section 5.5).

A heuristic method was used to chunk sentences into paragraphs. These paragraphs were then used as input to the clinically trained Temporal Histories of Your Medical Events (THYME) parser (Cai et al., 2024) to create AMR graphs². Each AMR sentence root was connected to

¹<https://spacy.io>

²The THYME parser was generously made available by Jon Cai and Martha Palmer at the University of Colorado Boulder.

its corresponding paragraph node in the alignment graph. Edges were added to connect each paragraph with its parent section node, which was classified using the MedSecId model (see Section 3.4). Once both graph components were built, Coreference Resolution alignment edges were added (see Section 5.6.1).

6.2.2 Concept Variable Renaming

Concept variables (i.e. “c” in `c/chase-01`) must be unique because multiple references (represented as multiple incoming edges to a referred concept in a graph) have no equivalent representation in the textual Penman format (see Chapter 5). All AMR parsers enforce this constraint so that concept variables are unique in the context of each sentence they parse. However, this unique assumption across joined graphs is often violated. Because Coreference Resolution across sentences uses the same representation of referring across concepts of a single node to nodes across sentences, concept variables must be unique¹.

A simple algorithm was used to make all concept variables unique. The algorithm indexes all concepts across all sentences of bipartite graph (note antecedent source and the discharge summary components) as shown in Figure 30. The algorithm keeps the convention of using the first letter of the concept followed by a 1-based index. Every concept of every sentence ordered by paragraph and section is enumerated and then replaced. Updated variables are also replaced in referring nodes to maintain referential integrity of the entire graph.

¹An additional reason for the unique concept variable constraint is graph reconstruction from Penman formatted sentences.

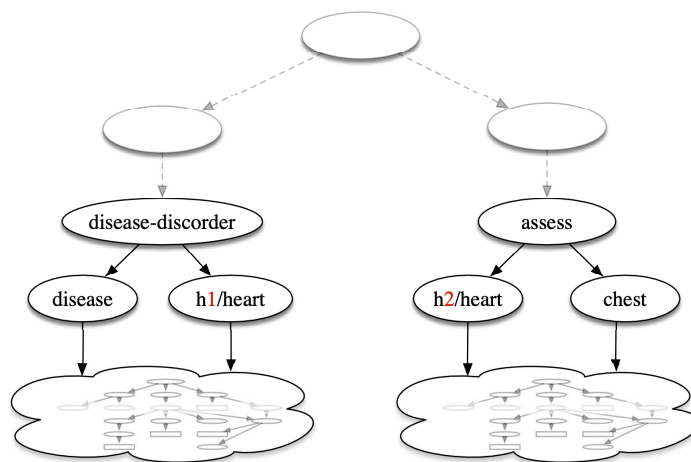


Figure 30: **Admission Graph Variable Renaming.** Once the sentence on the left is joined in the same graph with the sentence on the right, the name collision **h1/heart** is resolved by incrementing the variable index of the concept node in the second sentence.

6.2.3 Sentence Matching Algorithm

The sentence matching algorithm uses the CALAMR alignments to identify the sentences that best represent what is found in the summary. Figure 31 shows how the source sentence, “*Pre-cardiac catheterization assessment.*” matches with the discharge summary sentence “*Coronary artery disease, status post coronary artery bypass grafting,*”¹ by creating paths through the graph from a source sentence to a summary sentence. Each sentence that is connected in this way then becomes a candidate.

The sentence matching algorithm follows:

¹In this example, “catheterization” does not refer to the urinary track—rather it refers to a heart artery and grafting is an additional treatment.

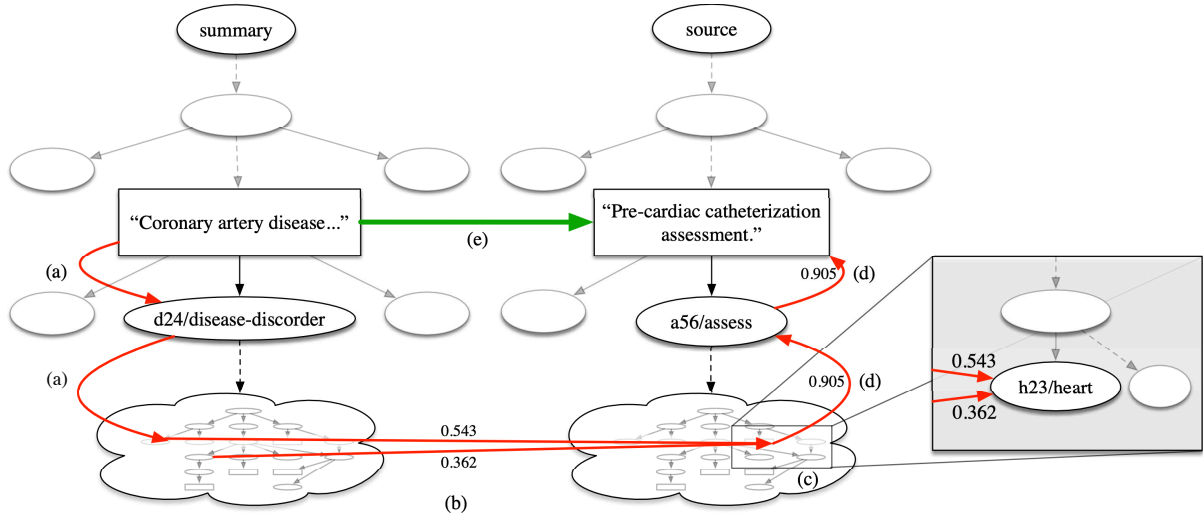


Figure 31: **Sentence Matching Algorithm.** The path (red) of alignment flow from the source to the summary for sentence. The enlarged box shows two incoming alignment flows from the source into the **heart** concept with a combined flow of 0.905. The green arrow represents a match candidate as a result of this alignment flow and the path to their respective sentences.

1. For each discharge summary sentence in the reduced graph (see Section 6.2) with an alignment, use a depth-first search to index all aligned nodes (Figure 31a).
2. For each node indexed in step 1, traverse the alignment edge over to source nodes in the note antecedent component (Figure 31b).
3. Annotate each aligned source node indexed in step 2 with alignment flows from the edges originating from the discharge summary component (Figure 31c).
4. Associate the aligned node summary annotations for each respective sentence in the source component (Figure 31d).

5. Create a sentence match candidate between the source and summary sentences (Figure 31e).
6. Sort the source sentences by the sum of the flow from each summary sentence.
7. Match sentences based on the flow from each summary to source sentence.
8. All remaining unmatched note antecedent sentences are given the **no-section** label.

Once the source sentences are pared with distributions of summary sentences by flow in step 6 each source sentence is matched with zero or more summary sentences. In step 7, a source sentence is matched with the summary sentence that has the maximum flow determined by the minimum sent flow (μ_s) hyperparameter. The matched summary sentence is then eliminated as a candidate for matching with any other source sentence. Finally, the source sentences are tagged with the section of the matched summary sentence.

6.2.4 Source Section Dataset

The set of notes that were successfully aligned, as described in Section 6.1, is referred to as the source section dataset. The sentence matching algorithm described in Section 6.2.3 was used to match the sentences of this dataset. The source section dataset includes only the notes with categories that were selected by clinical informatics physicians for the MedSecId (see Section 3.2) project. However, there were no *Consult* notes with enough alignment to have made it into the dataset. The note counts by categories are given in Table XVIII.

The selected discharge summary sections were based on what were considered most necessary and beneficial for summarization by a physician authoring the note by a clinical informatics

Note Cateogry	Count
Discharge summary	1,434
Echo	483
Physician	70
Radiology	3,652
Total	5,639

TABLE XVIII: **MIMIC-III Matched Sentence Notes.** Counts of notes per admission in the source section dataset.

fellow and a 4th year medical student. The physician-selected discharge summary sections and their counts are given in Table XIX. Most notable is the imbalance between the section labels and `no-section` label. This high disparity leads to a terse generated discharge summary, which is explained further in Section 6.4. The same process was used when selecting sections in the UI Health dataset. The sections and their counts are given in Table XX.

Figure 32 shows the number of matched sentence candidates in a contingency table based source sentences with their associated distribution (in terms of alignment graph flow) of summary sentences (see Section 6.2.3) for the MIMIC-III corpus. We see the largest flow from the discharge summary’s *Hospital Course* section¹ to radiology notes, which is understandable given the importance of the section (see Section 4.2.2.2). The *Hospital Course* section is the first to be read by physicians when reviewing a patient’s history. To my knowledge, it is the only discharge summary section generated using automatic summarization that has been peer

¹The reversed summary to source flow graph was used (see Section 5.7.5).

Section	Train Split	Test Split	Validation Split
no-section	63,558	7,709	7,816
Hospital course	4,585	590	562
History of present illness	2,131	272	297
Discharge instructions	1,267	195	147
Past medical history	1,123	139	137
Discharge diagnosis	459	44	60
Discharge condition	147	21	17
Discharge disposition	128	22	18
Addendum	68	6	9
Facility	58	9	11
Total	73,524	9,007	9,074

TABLE XIX: **MIMIC-III Matched Sentence Sections.** Counts of notes per admission in the source section dataset across splits.

reviewed (Adams et al., 2021). The second highest flow is from the *History Of Present Illness* section to radiology notes, which is not surprising given that this section is often copied and pasted into the discharge summary.

Figure 33 provides the UI Health dataset match sentence candidate contingency tables. This table has many more sections because of the note diversity in the dataset. Once again, we see a lot of *History Of Present Illness* and *Hospital Course* to the *H & P* note (*History and Physical Examination*).

6.3 Discharge Summary Generation

The discharge summaries were generated using the source section model (see Section 6.3.1) trained on the source section dataset (see Section 6.2.4), after which, the MIMIC-III aligned test set of admissions were compared with the machine generated summaries. The note antecedents

Section	Train Split	Test Split	Validation Split
no-section	14,422	2,053	1,203
Physical examination	959	102	116
History of present illness	883	104	92
Hospital course	734	97	71
Addendum	555	66	57
Discharge medications	548	65	84
Discharge instructions	503	48	56
Labs imaging	270	16	17
Assessment and plan	84	5	44
Discharge diagnosis	70	4	12
Discharge disposition	26	6	6
Chief complaint	23	2	1
Discharge condition	8	1	1
Total	19,085	2,569	1,760

TABLE XX: **UI Health Matched Sentence Sections.** Counts of notes per admission in the source section dataset across splits.

of the source section dataset’s test set were used as input to the source section model. Sentences were added to the predicted section in the generated discharge summary, or they were discarded if the special **no-section** label was predicted.

The UI Health dataset was also used to align and match and generate discharge summaries. However, this performance was poor as text output consisted of five or fewer sentences. In an effort to improve performance the process was repeated by tuning the CALAMR k^{th} order neighbor set hyperparameter to include more network neighborhood semantic information. The minimum sent flow (μ_s) hyperparameter was also adjusted down to include more sentence

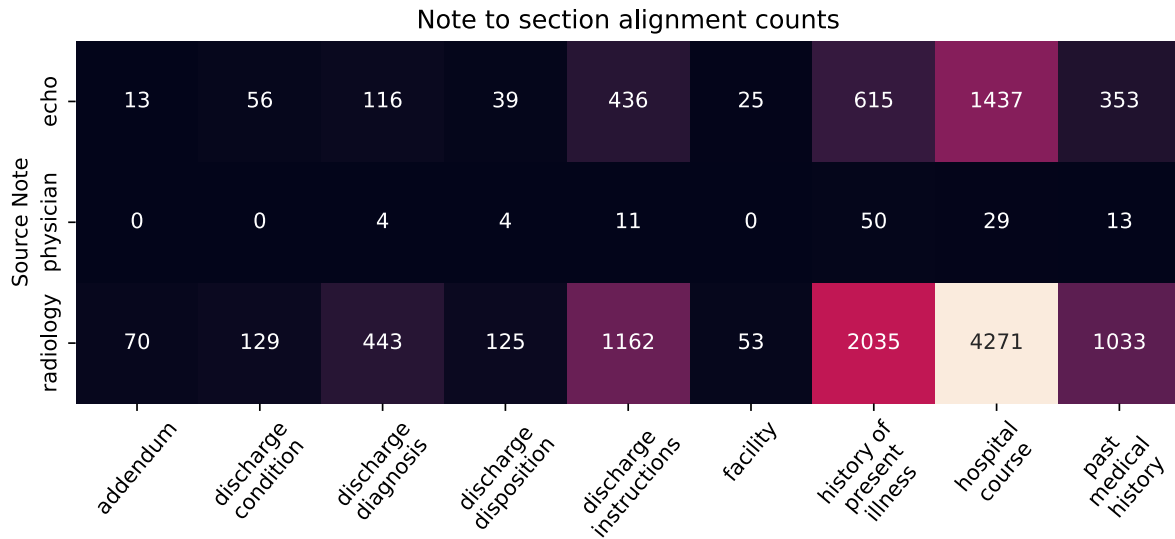


Figure 32: **MIMIC-III Note to Section Alignment Contingency.** The contingency table shows CALAMR alignment flow of data from the note antecedents to the sections of the discharge summary.

matches. These changes yielded 248 aligned admissions and better results and longer generated text on a second run.

The MIMIC-III trained summarization model yielded 133 automatically generated discharge summaries and the UI Health dataset model generated five. The alignment challenges described Section 6.1, such as missing discharge summaries and GPU memory constraints, show the difficulty of hospitalization summarization. Further analysis and discussion of these challenges are described in Section 6.4.

6.3.1 Source Section Model

Once the sentence matching algorithm was used to assign labels to source sentences (see Section 6.2.3) a bi-directional long-short term memory (BiLSTM) was trained to learn the discharge summary section type of each note antecedent source sentence, such as *Hospital Course*, or *no-section*. The full set of MIMIC-III discharge summary section type classes are listed on the x-axis of Figure 32, and UI Health dataset classes are listed in Figure 33. The memory network BiLSTM (Graves and Schmidhuber, 2005) was selected with the goal of remembering what sentences had been assigned to sections, and thus, reduce or eliminate redundancy in the summary.

The features to the model are the GatorTron embeddings¹ (Yang et al., 2022), the note antecedent’s note category, and section. Because of the very large input data size (see Section 6.1) the embedding model cannot be fine-tuned. Instead the static (frozen) model’s embeddings are used as input to the BiLSTM layer. A fully connected feed forward linear was added between the BiLSTM and output layers. The output layer is the set of neurons, each of which represents a the discharge summary section. Figure 34 shows the model with example sentences feeding into the layer from the bottom and forward propagating to the output layer at the top. The BiLSTM layer had a hidden size of 500 parameters, a dropout of $p = 0.15$, a learning rate of 5×10^4 and used gradient clipping. The model was set to train for 30 epochs and converged at 24 epochs.

¹A pretrained embedding model trained on more than 90 billion words de-identified clinical note text.

6.3.2 Experimental Setup

The standard set of quantitative machine learning (ML) performance classification metrics were used to evaluate the source section model. While this provides insights into how well the model performs, it does not inform as to how well CALAMR creates the training data for the source section model.

To bridge this gap, a human quantitative and qualitative evaluation was performed to better understand the effectiveness of the generated summaries. This informal evaluation of generated discharge summaries on the source section dataset’s test set was evaluated by a clinical informatics fellow and a 4th year medical student. Each generated discharge summary was ranked using a Likert scale (Likert, 1932) as an integer value ranking in the range 1 — 5 with five as the highest. Table XXI lists the questions asked for the informal evaluation.

Question Category	Question Text
Preference	Whether you prefer the generated summary?
Readability	Of the data in the generated summary, how readable is it?
Correctness	Of the data that is in the generated summary, how correct is it?
Complete	How complete is the generated summary?
Sections	Of the data in the summary, how well is it sectioned?

TABLE XXI: **Informal Evaluation Questions.** Questions given to medical domain experts for the human feedback informal evaluation of the generated discharge summaries.

6.4 Results

This section reports and provides an analysis of the performance of the source section model (see Section 6.4.2). It also documents the work done in choosing the AMR parser used for the project.

6.4.1 Parser Evaluation

The flow alignments of admissions with low note counts were manually analyzed for content overlap to make hyperparameter adjustments. The models were further tuned based on the generated summary output as described in (see Section 6.3). Prior to this analysis and training and tuning, AMR parsers were evaluated with the goal of finding the best parser for the clinical domain.

The AMR graphs were created from *Physician* notes, and then judged for correctness. An informatics fellow physician was asked to evaluate each AMR on a scale of 1 to 5 for correctness (5 as the most correct) for each AMR. I collaborated with the physician to provide the linguistic understanding for each graph.

Two AMR parsers fine-tuned on the Information Science Institute¹ *Biomedical* corpus (see Section 5.9.2) were evaluated. Eighteen graphs were scored including nine scored by the T5 amrlib² trained parser and nine scored by the Gsii (Cai and Lam, 2020) parser. The T5 parser scored an average 4.1 and the Gsii parser scored 2.9. However, the qualitative analysis of the

¹<https://amr.isi.edu/download.html>

²<https://github.com/bjascob/amrlib>

graphs by the physician showed fundamental flaws, such as products misclassified as molecules, misclassified roles, and added spurious nodes.

This analysis presented sufficient motivation to compare performance with the THYME parser, which was used in an informal evaluation with a clinical informatics physician and a 4th year medical student. Sentences from MIMIC-III discharge summaries were parsed into AMRs by the THYME and then translated back into sentences by a AMR graph-to-text model that was trained on the *Biomedical* corpus.

The physician and medical student were asked the following questions¹:

- How well did the physician write the original text?
- What is the quality of the generated text?
- How correct is medical terminology in the generated text?

Generated sentences for the THYME model were observed with a mean score of 3.42 on 91 sentences. However, for sentences with quality physician authored sentences, this mean increased to 3.7 on 75 sentences. An additional question on the correctness of the medical terms used in the generated output, which received a mean score of 3.5 for all 91 sentences and improved to 3.8 on the 75 originally authored sentences.

In addition to this quantitative analysis, a qualitative review of AMR graphs created from the parser on clinical text with a physician was also used to inform the evaluation on ten graphs.

¹All scores were reported on the Likert scale (Likert, 1932) in the range 1 – 5 where 5 is the best.

Each of these graphs were reviewed for completeness, correctness and CUI attachment (see Section 6.2.1).

The criteria of the parser selection was based on:

- the qualitative assessment that produced high quality AMR graphs by the THYME,
- the graph-to-text model used in the quantitative analysis was trained on out-of-domain Pubmed articles (*Biomedical corpus*) rather than clinical text, and,
- the quality of text and medical terminology of the quantitative study's scores.

Based on this criteria the decision was made to use the THYME for the alignment of all training and test data (see Section 6.1).

6.4.2 Source Section Model Performance

The source section model results are summarized in Table XXII (left). The weighted F1 score of 88.72 on the MIMIC-III trained corpus shows good results for the sentences' discharge summary section classification. However, we see a low macro F1 of 20.41, which shows the classification performing poorly on some labels. Another reason for a high weighted and micro F1 but low macro F1 is that the majority label, the **no-section** label, dominated as shown in Table XVIII.

The model trained on the UI Health dataset show lower results shown in Table XXII (right). This might be partly due to smaller dataset or the higher rate of reentrancies as shown in Table XVII (right) and discussed in Section 6.1. The fact that MIMIC-III is a curated dataset is the most likely reason the results are higher compared to UI Health dataset, which is unmodified and contains protected health information (PHI).

Metric	Value	Metric	Value
Weighted F1	88.72	Weighted F1	83.52
Weighted Precision	88.99	Weighted Precision	84.23
Weighted Recall	89.52	Weighted Recall	84.3
Micro F1	89.52	Micro F1	84.3
Micro Precision	89.52	Micro Precision	84.3
Micro Recall	89.52	Micro Recall	84.3
Macro F1	20.41	Macro F1	19.63
Macro Precision	29.07	Macro Precision	23.11
Macro Recall	21.27	Macro Recall	19.87

TABLE XXII: **Source Section Model Results.** The results as weighted, micro and macro scores of the source section model. The results of the model trained on the MIMIC-III corpus are given on the left and the UI Health dataset on the right.

The informal evaluation of the 133 generated discharge summaries trained on the MIMIC-III corpus is given in Table XXIII (left). The evaluation clearly shows the difficulty of the task and room for growth. The preference for the generated notes is low (1), the summaries are shown to not be complete (1) and the sectioning is incorrect (1). However, the generated summaries provide a readability of just over 3 and a perfect correctness score (5), which is what we expect from a faithful summary.

6.4.3 Discussion

The label imbalance in the source section dataset might be attributed to the sparsity of CALAMR’s alignments. If this were the case, we could adjust the hyperparameters of CALAMR to produce more sentence matches. However, the lack of alignment could be justified by the lack of notes (other than those from the ICU department) present in the MIMIC-III corpus (see

Question Category	Score	Question Category	Score
Preference	1	Preference	1
Readability	3.05	Readability	2.4
Correctness	5	Correctness	5
Complete	1	Complete	1.6
Sections	1	Sections	1.6

TABLE XXIII: **Informal Evaluation.** The average likert scale scores of 133 MIMIC-III generated summaries are given on the left and five generated summaries UI Health dataset on the right.

Section 4.2.1). The misalignment could also be attributed in cases where the physician writes from personal experience with the patient that is otherwise lacking from the EHR notes.

The five discharge summaries produced by the model trained on the UI Health dataset, shown in Table XXIII (right) show better completeness but slightly lower readability. A somewhat strange and interesting statistic is the higher sectioning of the UI Health dataset over the MIMIC-III despite the fact that the MedSecId model was trained on the latter. This implies the MedSecId is able to section the UI Health dataset notes or the source section model is able to predict sections based on other factors such as better alignments. See (see Figure 35) for an example of a de-identified automatically generated discharge summary. The gold discharge summary for this admission is given in Appendix C.

6.5 Conclusion

I have presented a new extractive method of multi-section automatic summarization; a significant step forward in generating faithful and traceable for clinical documentation. Shortcomings of the data used to train the models led to challenges that affected performance and

exposed limitations in the summarization process. In the case of the MIMIC-III data, the issue of missing notes (other than ICU as described in Section 4.2.1) lead to worse summarizations from missing hospitalization notes. The analysis of the impact of this is detailed in (see Section 4.4). A key findings of the DSProv study was that only 28% of *Physician* daily progress notes overlapped with discharge summaries (see Section 4.2.2). Because the *Physician* notes are used for manually summarized content in the discharge summary, the marginal representation reveals the challenge in using this corpus, or any partial dataset for training.

Summarization with the UI Health dataset trained model have shown to be as challenging to the MIMIC-III corpus, but for different reasons. Performance suffered due to memory constraints of the size of data in the UI Health dataset data, and specifically the number of notes compared to the procured MIMIC-III corpus. The poor sectioning of the data on the MIMIC-III trained MedSecId model is another reminder of how “real” EHR data can negatively impact the performance of the summarization task.

Despite these challenges, the source section model show promising results, which show that the generated training data from the CALAMR alignments produces a model that learns how to classify for extractive summarization. This is further qualified by the informal evaluation with high correctness scores and reasonable readability scores. These summaries are faithful in that only content for the source text is added to the summary. They are traceable in how each sentence can be traced back via the CALAMR alignments. Recommendations with respect to recent SoTA breakthroughs also inspire hope to bridging the gap of the difficulty of the summarization task.

	Source Note				Note to section alignment counts
	consults	h p	progress notes	significant event	
24 hour events	1	14	2	0	
addendum	27	555	96	0	
allergies	1	0	0	0	
assessment and plan	15	36	82	0	
chief complaint	0	26	0	0	
communication	0	9	0	0	
current medications	0	0	4	0	
discharge condition	0	8	2	0	
discharge diagnosis	2	77	7	0	
discharge disposition	2	35	1	0	
discharge instructions	22	502	82	1	
discharge medications	25	582	89	1	
family history	0	0	2	0	
findings	0	28	3	0	
history of present illness	53	935	91	0	
hospital course	54	679	168	1	
imaging	0	2	0	0	
labs	30	216	57	0	
imaging major surgical or invasive procedure	0	8	4	0	
medication history	0	5	9	0	
past medical history	0	3	0	0	
past surgical history	0	1	2	0	
physical examination	38	973	163	3	
reason	0	1	0	0	
review of systems	0	0	1	0	
social history	0	0	3	0	

Figure 33: **UI Health Note to Section Alignment Contingency.** The contingency table shows CALAMR alignment flow of data from the note antecedents to the sections of the discharge summary.

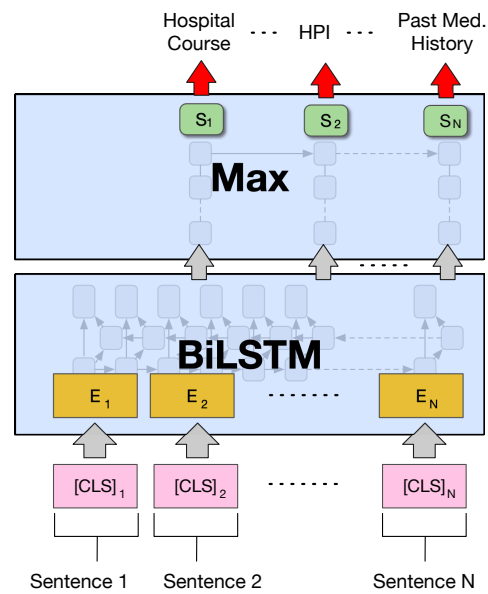


Figure 34: **Source Section Model.** The model classifies sentences' pooled [CLS] token, which is used in a BiLSTM across all sentences of all notes. Each prediction from each sentence is a section in the discharge summary.

History of present illness:

9:36 AM Status: Attested Editor: [**Doctor First Name**] [**Doctor Last Name**], MD (Resident) Related Notes: Original Note by [**Doctor First Name**] [**Doctor Last Name**], MD (Resident) filed at [**Date**] 4:01 PM Cosigner: [**Doctor First Name**] [**Doctor Last Name**], MD at [**Date**] 12:46 PM Consult Orders 1. Inpatient consult to General Neurology [[**Correspondence ID**]] ordered by [**Doctor Last Name**] [**Doctor Last Name**], MD at [**Date**] 0754 Attestation signed by [**Doctor First Name**] [**Doctor Last Name**], MD at [**Date**] 12:46 PM Stroke Attending: Pt eloped prior to being seen. Blurred Vision and Extremity Weakness [**First Name**] [**Last Name**] is a 47 y.o. male with PMH HTN, HLD, previous stroke [**Year**] (R ACA-MCA watershed), presenting to ED with L sided weakness, LLE numbness, and blurry vision b/ l. Patient reports symptoms started acutely on Friday [**Date**] around 4 pm while driving causing him to have to pull over. He decided to try to sleep it off. After waking up the following morning with no improvement, he went to [**Hospital Name**] ED where he was seen by neurology, but left AMA as he felt he was being asked the same questions repeatedly and nothing was getting done. CTH at [**Hospital Name**] was without ICH, reportedly showed wedge shaped hypodensity in frontal lobe likely from chronic infarct. The patient reports that the left sided weakness has improved somewhat today, but he still endorses b/ l blurry vision with constant white floaters, as well as numbness/ tingling in his LLE. CTA Hwith multifocal narrowing of b/ l ACAs as well as R M1 focal narrowing. MRI B w/o showing acute ischemia in medial R temporal lobe involving the posterior limb of R internal capsule as well as possibly the thalamus, and redemonstrating old R frontal ACA-MCA watershed infarct and old L occipital cortical infarct. Patient to be admitted to stepdown under stroke service. Patient with poorly controlled HTN and HLD, not taking any medications for the past 5 months as he says he had trouble getting primary care appointment for prescription renewals.

Physical examination:

NIHSS is 2 (LUQ quadrantopia, LLE numbness). SBP gets up to 200s per patient. Labs today significant for Troponin of 0.77. EKG showing T inversions in V5 and V6. Cardiology consulted. Patient denying CP at this time. Prior stroke/ TIAs (date, description): R ACA-MCA watershed stroke in [**Year**] per ED, worked up at [**Hospital Name**], on DAPT Vascular risk factors: HTN, HLD, prior stroke Past Medical/ Surgical history: Past Medical History: Diagnosis Date Hypertension

Figure 35: **Generated Discharge Summary.** A UI Health de-identified automatically generated discharge summary.

CHAPTER 7

SOFTWARE FRAMEWORK

(This chapter expands on the paper “*DeepZensols: A Deep Learning Natural Language Processing Framework for Experimentation and Reproducibility*” by Landes et al. (2023) in the Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software of the Empirical Methods in Natural Language Processing (EMNLP) conference.)

Summarization was implemented as an end-to-end pipeline with an admission’s electronic health record (EHR) clinical notes as input with a discharge summary as output. The pipeline first prepared the corpus by parsing the text and organizing it into data structures usable by the summarization algorithm.

The task of summarization requires natural language processing (NLP) and deep learning (DL) advanced programming interfaces (APIs) as clinical medical note summarization includes extra tooling and specialized models. To accomplish this task, the pipeline’s framework used the DeepZensols framework, which I wrote specifically for DL NLP applications. This framework reproduces consistent results, allows hot-swapping features and embeddings without further processing and re-vectorizing the dataset, and provides a means of easily creating training and evaluating natural language processing deep learning models with little to no code changes. In addition, it has been extended to include: medical domain specific tasks, such as concept unique identifier (CUI) extraction using the zensols.mednlp API, and PropBank database that include

embeddings for alignment and learning with neural networks (NNs) and alignment using the `zensols.propbankdb` API (see Section 5.4.1).

7.1 Motivation

Given the criticality and difficulty of reproducing machine learning experiments, there have been significant efforts in reducing the variance of these results. The ability to consistently reproduce results effectively strengthens the underlying hypothesis of the work and should be regarded as important as a novel aspect of the research itself. The proof of concept implementation used in this work utilizes DeepZensols, a new framework I developed, for not only reproducibility but also because of its NLP and medical domain extensions.

A key feature that sets this framework apart from others is a novel method to rapidly and easily swap features sets and compare performance across models (see Section 7.5). Other systems must re-parse and re-vectorize each mini-batch over each epoch. While there exist similar frameworks to ours (Gardner et al., 2018; Paszke et al., 2019; Ning et al., 2020; Alberti et al., 2018), none of these provides this batch strategy, vectorization of natural language text features and reproducibility of results across APIs and datasets in one framework. Popular NN architectures are available out of the box and easily configurable with little to no coding necessary (see Section 7.3 for NLP specific framework details) and the framework is freely available¹.

¹<https://github.com/plandes/deepnlp>

7.2 Library Design

DeepZensols is a combination of Python APIs built on top of PyTorch that provide a means of easily and quickly creating NLP task specific pipelines. The framework’s source code and installable libraries are released under the MIT License¹, which is available both on GitHub and as Python pip packages along with extensive and in depth overview and API documentation, tutorials, Jupyter Notebook examples and class diagrams for three NLP reference models and datasets. The framework is validated with 274 unit tests and six integration tests, most of which are automated using continuous integration testing for both functionality and reproducibility.

Reproducibility

All random state, including utility libraries, scientific libraries, PyTorch, and GPU state, is consistent across each run of the interpreter execution of the model’s training, evaluation and testing when using the framework. Results are consistent by saving this random state with the model, then retrieving and resetting it before using the model.

The order of mini-batches, and their constituent data can affect the model performance as an aspect of training or the results of validation and testing (Pham et al., 2020). This performance inconsistency is addressed by recording the order of all data² and tracking the training, validation and test data splits. Not only are mini-batches given in the same order, the ordering in each mini-batch is also preserved. These dataset partitions and their order are

¹<https://opensource.org/licenses/MIT>

²Regardless of any given data pre-processing or shuffling.

saved to the file system so the community has the option of distributing it along with the source code for later experiment duplication.

The framework also saves the configuration used to recreate the same in-memory state along with the model. This duplicates the model structure, parameters, hyper-parameters and all other train-time memory during testing. Reproducibility is built in to unit tests for individual components and integration tests by comparing the validation and training loss across six data sets¹.

7.3 NLP-Focused Abstractions and Features

The framework provides many APIs for natural language tasks, including concatenation of vectorized language features to input embedding (see Figure 36). Vectorization of contextual embeddings such as BERT (Devlin et al., 2019) and non-contextual embeddings such as word2vec (Mikolov et al., 2013b), GLoVe (Pennington et al., 2014) and fastText (Bojanowski et al., 2017) are available.

The framework includes many layer implementations, which are compatible with the PyTorch API as module classes. Examples of layers provided include BiLSTM CRF, BERT transformer models, 1D convolution NN, expanding or contracting DL feed forward linear networks with repeats with input and output feature calculation, word embedding layer for concatenating

¹Data sets include the MNIST, Adult, Iris datasets and those mentioned in Section 7.8.

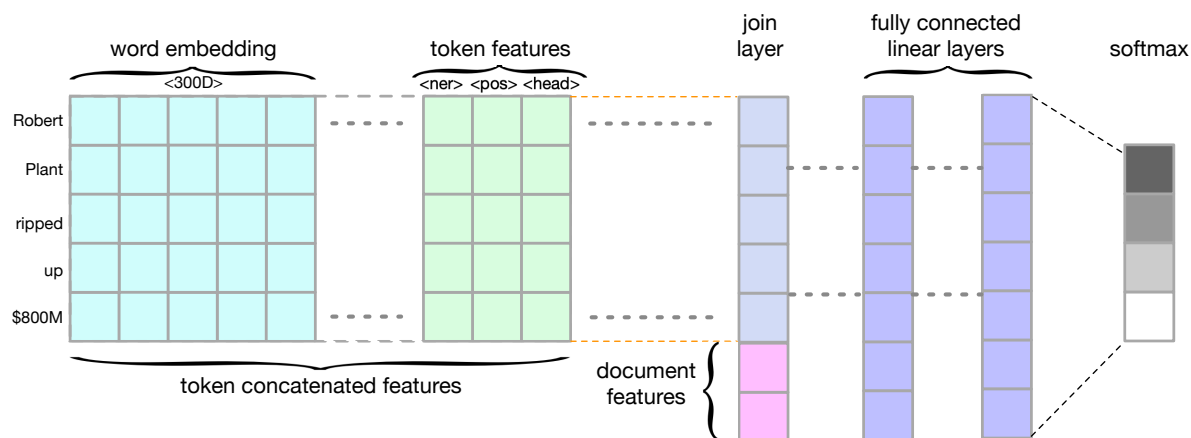


Figure 36: **Network Architecture.** Word embeddings concatenated to vectorized linguistic features, and then joined with vectorized document features constructed using configuration with no coding.

features (see Section 7.4), TF/IDF frequency weighting (Sparck Jones, 1972), and an adaption module for the pytorch-crf¹ CRF implementation.

HuggingFace transformer layers are available as embeddings, document, sentence and token features. The framework also provides direct access to these models' data and utilizes it in a variety of tasks such as text classification, token classification, language generation, latent semantic analysis, etc. A linguistic feature mapper that translates spaCy² to wordpieces, which are token sub-units with associated vectors and provided by the model's tokenizer (Wu et al., 2016), is also accessible as an easy to configure module.

¹<https://github.com/kmkurn/pytorch-crf>

²<https://spacy.io>

7.4 Vectorization

The DeepZensols framework allows for easily configurable components that provide a higher level abstraction that tokenizes, sentence chunks, and vectorizes linguistic features. These *vectorizers* have a class taxonomy based on data they vectorize so their output data can be automatically constructed in various off-the-shelf architectures. For NLP tasks, vectorizers include:

token: spaCy token features such as part of speech (POS) tags, named entity recognition (NER)

tags, dependency head relations and depth.

document: Features taken from the document level, typically added to a join layer such as count sums of spaCy parsed features.

embedding: Vectorizes text into word embeddings, such as sentence or document text.

multi-document: Aggregating and shared features between more than one document, such as overlapping POS or NER tags.

See Section 7.3 for more information on NLP specific feature generation.

7.5 Batching

We provide a novel method to vectorize and batch data without wasteful pre-processing of each feature and embedding combination. Other similar frameworks pre-process data in an intermediate form only once before training. However, this leads to a brittle and difficult to reproduce dataset of ad-hoc text processing scripts that are challenging to re-execute, and thus, reproduce performance metrics.

Our framework addresses this with an organized intermediate file scheme and partitioned feature set so the input data is vectorized only once efficiently using a multi-processing pipeline. The output format of this process allows for quick feature swapping and hyper-parameter tuning for re-training. It leverages the fact that mini-batches are independent and fit nicely as independent units of work by segmenting datasets into smaller chunks, vectorizing each chunk in parallel sub-processes, and creating one or more batches independently across each sub processes.

This process by which data is written to the file system in a format that is fast to reassemble is called *batch encoding* (see Figure 37), and accomplished by a) splitting sentences and/or tokens into equal size “chunks” units of work, b) parsing natural language features from chunks across multiple processes, and c) vectorizing each chunk as tensor data in separate files by feature.

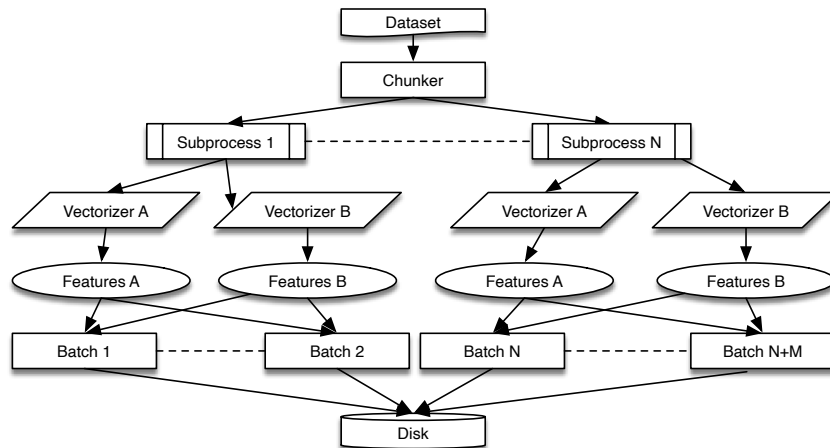


Figure 37: **Batch Process.** The mini-batch creation process splits the work done by vectorizers (two in this example) and writes constituent feature files to disk in sub processes.

After batch encoding is complete, the model is ready to be trained from data obtained from a *batch decoding* step, which is accomplished by: a) choosing a feature set for a training run, b) reassembling features by mini-batch (see Figure 38), c) decode each mini-batch into a tensor (see Figure 39), and d) load, cache and copy tensors to the GPU.

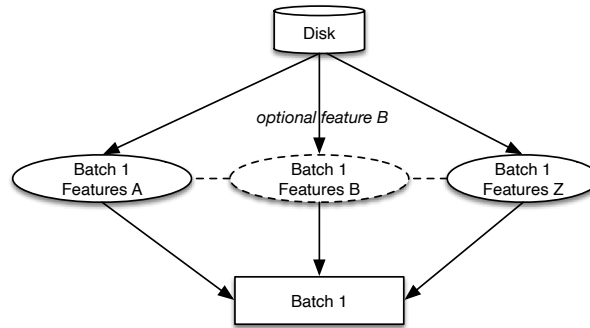


Figure 38: **Batch Reassembly.** Batch reassembly loads mini-batch tensors from only files containing features for the current run.

Reassembling mini-batches by feature greatly reduces load time and memory space, which speeds up model training (see Section 7.8) and ameliorates issues of complex models. This speedup is most apparent when comparing the time taken to fine-tune BERT with using the transformer’s static embeddings (see Section 7.8). In the case of the former, large data files with output tensors are read as wordpieces embeddings (Wu et al., 2016). The train, validation and test cycle is faster for other vectorized linguistic data such as spaCy features as well.

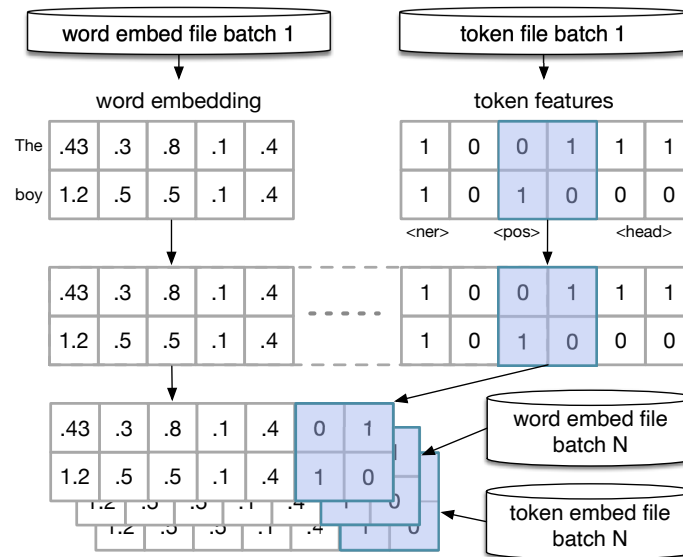


Figure 39: **Batch Decoding.** Batch decoding “stitches” mini-batches together from only files containing features for the current run.

7.6 Technology Stack

Each “layer” of the stack builds on more general libraries to reduce the installation footprint based on the needs of each use case. Each library contains the requirements for dependent third-party and lower tier packages. The framework consists of the following libraries (see Figure 40):

- `zensols.deepnlp`: Contains language vectorization, such as word embeddings (see Section 7.3), part of speech tags, named entity recognition, and head dependencies (McDonald et al., 2005).

- `zensols.deeplearn`: General purpose DL API that provides mini-batching, vectorization, and training, validation and testing of a model. The package provides general layers, many of which, are utilized as sub layers by `zensols.deepnlp`.
- `zensols.nlp.parse`: Parses natural language text using spaCy, generates and vectorizes features in an object graph that is easy to use and fast to serialize.
- `zensols.install`: An API to download, install and uncompress online web resources such as models and datasets. File system paths are then provided to the project in the same configuration set.
- `zensols.util`: Utility library command line parsing, persistence and a Java Spring like inversion of control (Mattsson, 1996) configuration system.

7.7 Execution

The framework provides both a command line and a Jupyter notebook interface to train, test and predict. A “glue” API is used to make a Python `dataclass`¹ class a dynamically generated command line with help usage message documentation. A set of default application classes are available with the framework, but they can be extended to include project specific workflows. The default application set provides interactive early stopping or epoch resetting during training.

The command line and Jupyter notebook use a common interface to the model itself, which is conducive as an entry point to larger projects or simple run scripts. Both interfaces provide

¹<https://docs.python.org/3/library/dataclasses.html>

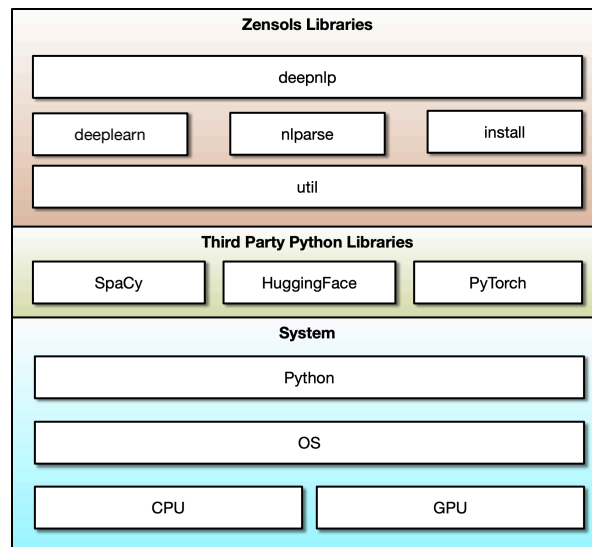


Figure 40: **Framework Library Stack.** The Python Library Stack

evaluation metrics such as recall, precision, F1, ROC, accuracy and several others and training and validation loss plots (see Figure 41). The command line writes the loss plot for each epoch and Jupyter notebook renders it in an inline cell.

A debugging mode that outputs a step of the model training with batch composition, layer names, dimension calculations, using the logging system, which is filterable by module or component is also available.

Results are organized by each run and carry a common file system structured named by either what is provided in the configuration or by model name. This directory structure contains the full model with all configuration, the PyTorch model, and results provided as human readable indented text, JSON and binary formats.

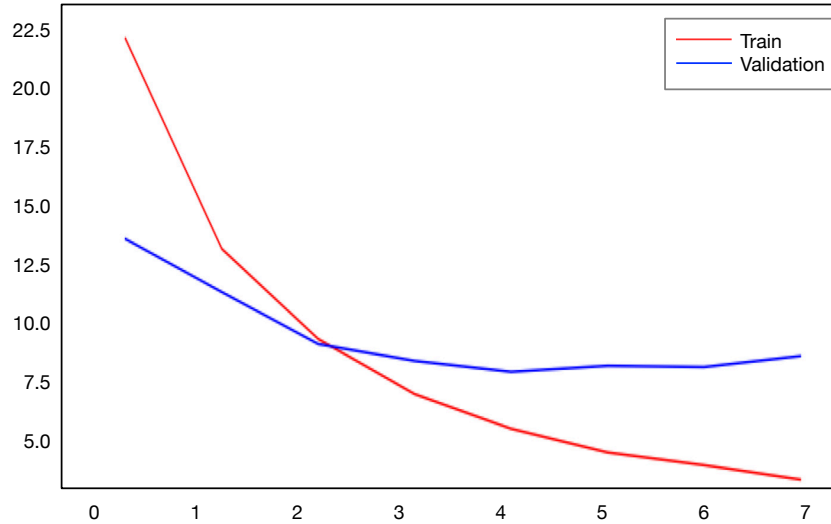


Figure 41: **Validation and Training Loss.** Plot from the NER Token Classification Application using RoBERTa embeddings. (Graph selected for illustrative purposes rather than performance.)

7.8 Runtime Analysis

Runtime analysis was performed for parsing, feature vectorization (see Section 7.4), batching (see Section 7.5), training and testing three different types of models using a Nvidia TITAN RTX graphics processor on an Intel 3.6GHz CPU summarized in Table XXIV. The time to train, validate and test was measured using the following criteria:

- Model: the model trained and evaluated.
- New Batch: whether or not the mini-batches were (re)created (see Section 7.5).

- Cached: whether or not the mini-batches were cached in GPU memory.¹

Since obtaining fast results allows for more experimentation with a variety of feature sets, embeddings, and NN architectures, the table shows each dataset’s quickest end-to-end completion with all time duration in hours, minutes and seconds.

The datasets used in the runtime analysis include: the CoNNL 2003 (Tjong Kim Sang and De Meulder, 2003) for named entity recognition, the movie review corpus (Pang and Lee, 2005; Socher et al., 2013) for sentiment classification, and the clickbate corpus (Chakraborty et al., 2016) for text classification. Batched feature sets were uniform for each dataset and included:

token: spaCy POS and named entity type

embedding: BERT, GloVE 50D, GloVE 300D, word2vec 300D

document: spaCy named entity counts, head depth, and statistics (such as sentence, token, character counts and averages)

Not surprisingly the GloVE model shows the best improvements since the word embeddings are non-contextual, and thus, need not be recomputed per sentence. The framework shows significant processing time improvements in the movie review sentiment dataset with a 2.6X speedup using only batching (see Section 7.5), and 3X using GPU caching. Interesting, the clickbate dataset shows a converse relationship with a 1.5X batching speed up, but a 4.7X

¹The framework offers GPU caching, CPU caching, and iterative buffering of mini-batches.

Data	Model	New Batch	Cached	Duration
N	BERT			01:00:20
N	BERT		✓	01:00:15
N	BERT	✓		01:04:43
N	BERT	✓	✓	01:04:31
N	GLoVE			00:29:45
N	GLoVE		✓	00:24:15
N	GLoVE	✓		00:34:04
N	GLoVE	✓	✓	00:28:23
M	BERT			00:19:11
M	BERT		✓	00:19:10
M	BERT	✓		00:21:14
M	BERT	✓	✓	00:21:41
M	GLoVE			00:01:55
M	GLoVE		✓	00:00:38
M	GLoVE	✓		00:05:01
M	GLoVE	✓	✓	00:03:41
C	BERT			00:03:54
C	BERT		✓	00:03:54
C	BERT	✓		00:05:44
C	BERT	✓	✓	00:05:45
C	GLoVE			00:03:52
C	GLoVE		✓	00:00:49
C	GLoVE	✓		00:05:43
C	GLoVE	✓	✓	00:02:40

TABLE XXIV: **DeepZensols Efficiency Benchmarks.** Efficiency benchmarks showing each dataset with N=named entity recognition, M=movie review sentiment, and C=clickbate.

speed up with scenario GPU caching. While BERT saw a less drastic speed up, it did finish almost 10 minutes faster under the NER model.

The framework shows significant processing time improvements in the clickbate corpus, which was achieved by reusing the vectorized feature batch tensors on disk and by caching those tensors in the GPU. Furthermore, the novel batching approach (see Section 7.5) saves

one minute and 40 seconds (3X faster) with a total five minutes and 40 seconds (28X faster) when buffering mini-batches in the GPU.

7.9 Conclusion

The DeepZensols framework is a viable solution to easily create NLP specific models with APIs and analysis tools to produce consistent results. Such frameworks create the types of models that give confidence and legitimacy by providing a way to produce reliable reproducible results, which is ideal for the summarization pipeline (see Chapter 1) Runtime analysis shows the framework offers significant processing time savings compared to systems that do not provide feature caching with stable results.

CHAPTER 8

CONCLUSIONS AND FUTURE WORK

The goal of this work is to find new ways to ground summaries using state of the art deep learning (DL) methods. The output of this work is several methods and annotation datasets that have been shown to successfully generate faithful and traceable discharge summaries.

8.1 Contributions

The contributions of this work include annotated datasets and summarization methods. One dataset is 2,002 annotated medical notes from the Medical Information Mart for Intensive Care III (MIMIC-III) corpus for clinical sections in both lexical demarcations and section type (see Chapter 3). This project also includes a model that automatically sections clinical text with an F1 metric of 96%, and an analysis by section that shows the orthogonality of medicine vs. surgical medical concepts.

Another contributed dataset includes 51 fully physician-annotated MIMIC-III admissions of semantically similar, or copy and pasted text spans, between electronic health record (EHR) notes and respective (per admission) discharge summaries. This work also includes a provenance of data study that explores how data flows from previously written EHR note spans of text (by note category and section type) to the discharge summary (see Chapter 4). This work also includes a novel unsupervised span-matching method to automatically create match spans in any domain (see Section 4.3).

The alignment method has shown to be successful for finding summarized content in news articles (see Chapter 5) and leverages the performant and well-established max flow algorithm, which has been seldom used in natural language processing. This work provides a method for aligning graphs useful for explainable summarization and uncovering the extent of overlap between source and summarized text. The work gives a new summarization specialized score to measure the summarized overlap and a metric for calculating similarity between multi-sentence AMR (abstract meaning representation) graphs, which were not previously available.

A general purpose NLP framework (see Chapter 7) with extensions for each module presented is another contribution. This framework provides a mechanism for parsing and vectorizing natural language text into features, then provides a method of creating reproducible results for DL NLP models. The clinical domain, alignment and summarization modules were implemented as extensions of this framework, and then automated as an end-to-end pipeline.

The most important contribution of this work is the method and implementation of a pipeline to automatically generate discharge summaries (see Chapter 6). The summarization pipeline takes previously written discharge summaries written by physicians and uses it as training data. This model is trained by first using the MedSecId model to section the notes, then uses Component ALignment for Abstract Meaning Representation (CALAMR) to align parsed AMR graphs. Finally, CALAMR alignments are used to train the source section model model to copy sentences from the note antecedents to the discharge summary as an extractive summarization task. Given the lack of any other peer-reviewed work on this topic, I believe this

is a beneficial first implementation of a complete end-to-end discharge summarization. However, more work is needed to improve the efficiency and performance of the results.

8.2 Future Work

The summarizer is able to generate a comprehensive list of sections when generating discharge summaries. These sections are listed (see Table XIX). However, certain sections have a sparse representation in the discharge summaries reviewed in the informal analysis. Additional error analysis in CALAMR alignments could be helpful in finding where and why more of these sections are not aligned. Loosening the sentence matching algorithm constraints might be helpful in generating more data. The most salient and important discharge summary section to summarize is the *Brief Hospital Course* (Adams et al., 2021). This section was also summarized by Damm et al. (2024), which used the MedSecId (see Chapter 3) model in their pipeline. Other important sections that are important and useful in a generated summary include (in order) *History of Previous Illness*, *Assessment and Plan*, *Past Medical History* and *Past Surgical History*.

Perhaps a more straightforward effort to include more summarized sections would be to adjust the sentence matching algorithm minimum sent flow (μ_s) hyperparameter (see Section 8.2.2). There is work left to be done to other models in the pipeline as well. Hyperparameter tuning the MedSecId models to improve how clinical notes are sectioned may help in reducing propagation error in the pipeline. Another obvious opportunity to improve performance is to concatenate *cui2vec* embeddings in the input layer as described in Section 3.2.3. Using pseudo or synthetic tokens in place of the MIMIC-III training data could improve perfor-

mance and shed light on how models learn with more realistic data. Implementation changes to increase the performance of CALAMR is definitely necessary given the time it takes to align large admissions in the UI Health dataset corpus.

To improve the current extractive summarization, the DSProv (see Chapter 4) annotations could be used as a validation and test set for the extractive method. One of the original goals of the data was to use it tune the CALAMR hyperparameters to maximize the alignments via the note matches. In terms of admissions, the dataset is relatively small (51), but the dataset’s 569 note matches could be used for training a supervised model to help in additional human annotation. However, the most important future work is an abstractive graph-based summarization method.

8.2.1 Abstractive Summarization

Switching from extractive to abstractive summarization using graph neural networks (GNNs) is the largest of the recommendations, but has the greatest potential. True abstractive summarization could be accomplished using the CALAMR alignments in combination with a GNN. A modified version of a deep auto-regressive model that generates graphs (GraphRNN) would be used with a model to predict the summary graph (You et al., 2018). However, this would require modifications to the published algorithm and would require more training examples than those used in the experiments described in Chapter 6. The need for more alignment data motivates the recommended optimizations of the CALAMR method (see Section 8.2.2).

Leveraging the source component might mitigate the need for a large training dataset since the model would have more prior data from which to draw. The modifications to the GraphRNN

would exploit the source component rather than just sampling new graphs from a probability distribution over the training data per the GraphRNN method. The neural network (NN) architecture would also need to be expanded to accommodate local network neighborhood embeddings (in addition to the adjacency matrix). However, the strategy of utilizing a recurrent neural network (RNN) cell for each iteration of the learning algorithm as a node traversal, and then growing the summary component for each iteration during testing, remains the same. A long-short term memory (LSTM) (Hochreiter and Schmidhuber, 1997) would be a good choice for the RNN for its ability to remember long node distances across the graph. Gradient clipping should be considered to alleviate issues of LSTM exploding gradients (Bengio et al., 1994).

8.2.1.1 Training

The model would learn network neighborhoods in the source as a prior to predict the summary graph across the alignments of the bipartite graph. The sequence based model's final layer model output (called the prediction head) would be used to classify the next AMR concept or attribute node based on previously predicted nodes. Another model would add edges between the new predicted node and existing nodes that were previously added to the graph. As with the GraphRNN, the edge predictions are folded into each step of the node prediction step as shown in Figure 42.

Each concept or attribute node embedding formulated in Section 5.4, and shown in Figure 16, is used as input layers to the GraphRNN node model (see Figure 43b). Similarly, each role edge's embeddings are used as input to the GraphRNN edge model. Before training begins alignment graphs are created (see Chapter 5) and given to the summarizer for each

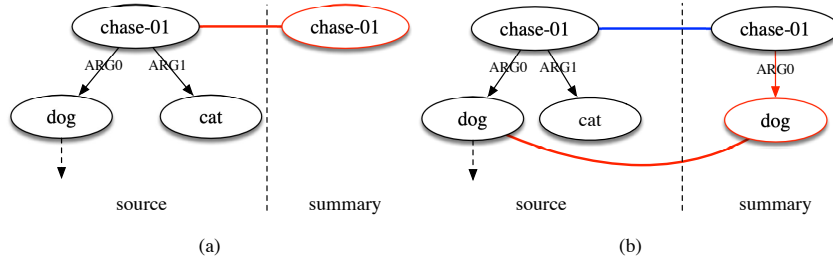


Figure 42: **Summarization Iteration.** Two iterations first classifying node `chase-01` and then `dog` with role `ARG0` with updates in red.

admission (or document) as input. The model would take minibatches of node or edge features and a label in typical GNN fashion, and the training algorithm closely resembles that of the GraphRNN (You et al., 2018).

The training algorithm identifies a target node that iterates as a breadth-first search through the source component. The GNN prediction head identifies a node in the summary component to learn or predict with the node model. The GNN prediction head identifies two nodes to connect for the edge model.

The training algorithm follows:

1. Configure the GNN models. Note that both models use node and edge embeddings.
 - (a) The node model's local network neighborhood embeddings of the source component and aligned edges to the summary component are added as successive layers of the NN from the current prediction head.

Each l^{th} layer has embeddings of the k^{th} order neighbor set nodes k hops from the predicted node \mathbf{n} such that layer $\mathcal{L} = \mathcal{U}(\mathbf{n}, l)$. See Figure 43 for the network neighborhood with a three layer ($k = 3$) GNN. The layer and network neighborhood depth is set with hyperparameter graph neural network layer depth (δ) a priori.

- (b) The edge model's model is created in the same was as the node model in step 1a.
2. The node model's RNN is initialized the **start** state per the GraphRNN training algorithm. The same is done for the edge model's RNN cell.
 3. Set the target node to the root node of the source component.
 4. Train the models using the following steps for a single back propagation.
 - (a) Set input layer of the GNN to the network neighborhood embeddings. The node embedding is used as the prediction head's output layer. Type of node (concept vs. attribute) could be modeled another output neuron. Regardless, the prediction would have to be a multilabel classification.

A classification over all PropBank role entries could be used for predicating element concept nodes (see Section 5.1.1). However, a vocabulary prediction would be needed for attribute nodes, light verbs, etc.

 - (b) Add the node to the adjacency matrix. Because the graph could be large, only the last N previously traversed target nodes are tracked (You et al., 2018).
 - (c) Update the node model's parameters.

- (d) Populate the edge model’s input layer with the edge embeddings of the network neighborhood with the same hyperparameters for network size as done in step 4a.

The label of the edge model is the role edge edge and two nodes to connect as sequence identifiers of nodes for the preview N nodes. Add the sequence identifiers to the adjacency matrix to be used as parameters of the model.

- (e) Repeat step 4d for all nodes in the summary component.

- (f) Update the edge model’s parameters.

5. Traverse to the next node in the breadth-first-search and set the target node to it.
6. Repeat steps 4 – 5 across all nodes in the source component.
7. Add the **end** state to the node model’s LSTM (You et al., 2018).

8.2.1.2 Testing

Testing follows a similar pattern as the training process:

1. Set the target node to the root of the source component.
2. Populate the node model’s input GNN layer to the network neighborhood of the target node in the source component. This network neighborhood includes some portion of the summary component (after the root node is created). An example of this prediction is the red **chase-01** concept node shown in (see Figure 42a).
3. Given the network neighborhood, use the node model to predict a new node in the summary component.

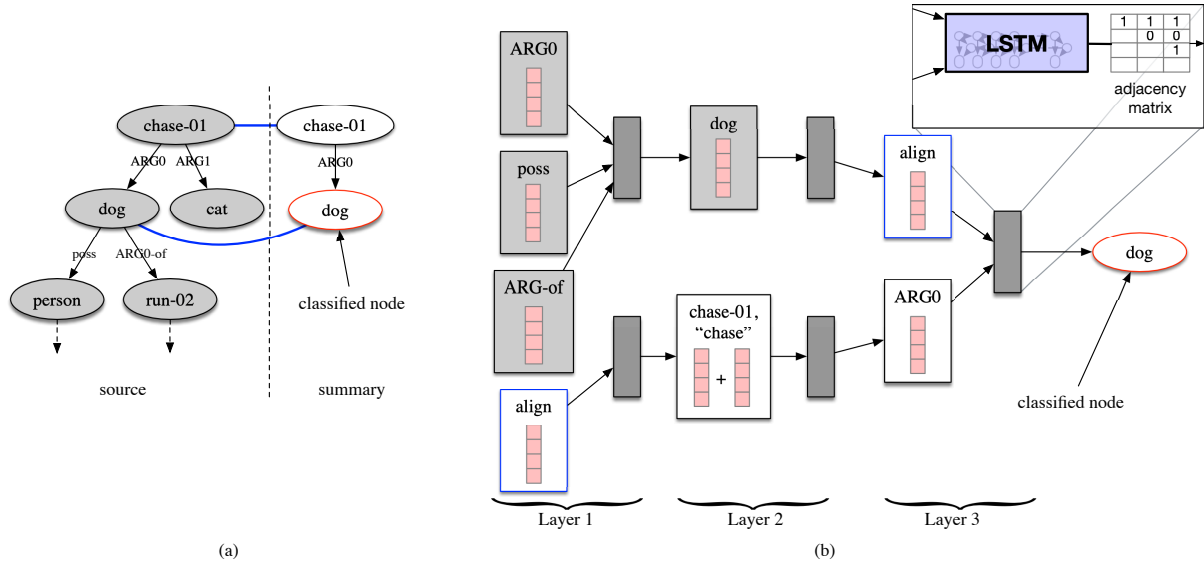


Figure 43: **Prediction with the Summarization Network.** Left (a) the AMR graph predicting the next node **dog**. Right (b) the GNN created from the AMR graph in (a). The source component nodes are shown in light gray, the predicted node in red, alignment edge and embedding in blue, and dark gray boxes are pooling operators.

4. Do the same input layer population as in step 2, but on the edge model.
5. Given the network neighborhood and the adjacency matrix, predict a role edge and two nodes in the summary component to connect. Add the new edge to the adjacency matrix.
An example of this prediction is the red alignment edge shown in (see Figure 42b).
6. Repeat step 5 for all nodes in the summary component.
7. Repeat steps 2 — 6 until the node model's LSTM **end** state is reached.

8.2.1.3 Generation

After the summary component is created, the discharge summary is generated. This is done by iterating over each section node in order (retained by the LSTM), which are children nodes of the summary component root. For each sentence node in each section use the sentence AMR root as input to a AMR graph-to-text model to generate the natural language sentence.

Tense can be added to the output by including part of speech (POS) tags in the nodes using the text-to-graph aligned tokens (TATs). Attaching the POS tags as nodes is currently supported by the framework that creates the admission graphs. However, the AMR graph-to-text model would need to be retrained with the POS nodes added for each training instance.

8.2.2 Model Tuning and Optimization

Currently, the process takes too long and too much computing power. The larger admissions take hours just to align in the ACER (covered entity) environment. Optimizing the code will speed things up. But to really improve performance, it would be ideal to find an implementation of the max flow algorithm that runs on the GPU rather than the CPU as the C/C++ library can take up to 20 minutes on a single iteration for larger admissions.

Adjusting the CALAMR hyperparameters could be even more impactful as shown in Section 6.3.1. Tuning these hyperparameters has the potential to create better aligned content and would improve training data for the source section model. Bayesian hyperparameter optimization (Bergstra et al., 2013) could be used to find the best alignment with the physician-annotated admissions in the DSProv dataset (see Section 4.2), and would very likely also produce better results.

To facilitate a faster alignment for a tighter hyperparameter update, align, and test loop, the alignment capacity calculation would need to be optimized as a matrix multiplication operation on the GPU from its current per node and node neighborhood CPU computation. There are many nonlinearities, such as the sigmoid sentence skew, for sentence dampening calculations (see Section 5.5). Optimizing over network neighborhoods further complicate this parallel computation since k^{th} order neighbor set nodes would have to be indexed (see Section 5.4.7).

In addition to more hyperparameter tuning in the source section model, a new model to predict sentence order would likely improve the quality of summarization output given that the medical professional feedback noted sentences out of order.

8.2.3 Application to Large Language Models

There is a strong likelihood that the alignment module (CALAMR) of this dissertation can be used to detect and reduce (or eliminate) hallucination in large language models (LLMs). For example, Figure 27 shows the final flow of a summary graph on the right (green graph). The sentence node “A plane flew over.” has no flow via the role edges because the method found no matching or overlapping content in the source.

As explained in Chapter 6, clinical notes can number in the hundreds per admission. For such large amounts of data, utilization of LLMs for summarization is currently not feasibly given text length and GPU memory constraints. However, for smaller admissions or other less memory-demanding clinical tasks, such as that of radiology report summarization, there is potential to use LLMs for summarization, and then “smooth” out the summary using CALAMR.

Specifically, this process would involve sectioning the text and then summarizing each section in turn. The summarized output of this step, along with the source, would then be aligned with CALAMR. As discussed in Section 8.2.1, the zero flow role edges would then be removed from the source graph, then a AMR graph-to-text is used to generate the text of the summary. Because none of these processes are specific to the clinical domain, any domain source could be summarized using LLMs with hallucination removed in this fashion. As explained in Chapter 6, these summaries would be faithful to the source and traceable through the CALAMR alignments.

APPENDICES

Appendix A

DEFINITIONS

- abstractive** A method that generates unique text not always found in the source text. p. 1
- admission graph** A graph that forms the semantic meaning of a hospital admission. p. 129
- alignment graph** A flow network bipartite graph whose component connecting edges serve to align nodes for the purpose of calculating a match score or to create training data for summarization. p. 80
- alignment graph algorithm** The algorithm that uses the alignment graph construction as a flow network to reduce the capacities. p. 11
- alignment edge** An edge that joins a concept node pair across the source and summary components. p. 11
- alignment graph construction** The portion of the graph alignment algorithm that constructs a flow network from the source and summary abstract meaning representation graphs. p. 101
- AMR** abstract meaning representation: A semantic representation language that describes the abstract meaning of a sentence as an acyclic graph. AMR captures “who is doing what to whom” in a sentence (Banarescu et al., 2013) and can also be represented in the context free notation of Penman (Kasper, 1989). p. 2
- AMR graph-to-text** A model that generates natural language from abstract meaning representation Penman notation. p. 145
- automatic summarization** A computational process to shorten natural language text by reproducing only the salient information from the source text. p. 1
- BERTScore** Bidirectional Encoder Representations from Transformers Score: An evaluation metric for scoring generated text. This quality is assessed by comparing reference and generated tokens’ similarity based on the contextual embeddings of a BERT family model. p. 19
- Bleu** Bilingual Evaluation Understudy: A measure that assess the quality of machine generated translations. This quality is assessed by comparing human translated text. p. 60
- Calamr** (Component ALignment for Abstract Meaning Representation) An Abstract Meaning Representation Alignment Method p. 68
- aggregate alignment portion** Harmonic mean of the alignment node portion across both components. p. 113
- aggregate flow** Harmonic mean of the flow across both components. p. 113
- summary aligned portion** The portion of nodes in the summary component that have at least one alignment. p. 113

Appendix A (Continued)

- summary root flow** The value of flow exiting the summary node to the sink. p. 112
- source aligned portion** The portion of nodes in the source component that have at least one alignment. p. 112
- source root flow** The value of flow exiting the source node to the sink. p. 112
- edge node child** The node at the target end of a directed edge e in the alignment graph. p. 89
- Coreference Resolution** The task of finding all mentions referring to the same named entity. p. 101
- CUI** (concept unique identifier) a unique identifier within the SNOMED-CT subset of the UMLS p. 44
- cui2vec*** A unique concept identifier embedding trained from biomedical text. p. 44
- node descendant** The descendants of v is defined as all paths from that node (grandchildren) to the terminal leaf nodes. p. 97
- document node** A node in the alignment graph that represent some aspect of the document rather than AMR such as paragraphs, sections, documents etc. p. 10
- discharge summary** A clinical note that describes a patient’s medical history, stay at a hospital, and the care they received. p. xv
- DSProv** A automatic discharge summarization feasibility study. p. 51
- DSProv dataset** A dataset of annotations that associates note antecedents with discharge summaries for the DSProv study. p. 56
- encoder-decoder** A deep learning architecture that encodes state in an intermediate network layer and useful when inputs and outputs have varying sizes. p. 18
- extractive** A method that copies selected sentences from the source text verbatim to the summary. p. 1
- faithful** How accurate the summary is. Refers to how accurate the summary is with respect to the source text whence it originates. p. 1
- flow network** A graph, or capacitance network, that associates a capacity and a flow with each edge of a graph. p. 10
- capacity** The upper bound constraint on the flow of material through an edge of a flow network graph. p. 27
- flow** A value of the amount of flow through a flow network. p. 76
- max flow** The maximum amount of flow available to traverse an s-t flow network given the capacities of the network. p. 6
- flow per node** The flow through an edge divided by the number descendants starting the node incident that’s closer to the leafs to the terminal leaf nodes. p. 108
- GatorTron** A pretrained embedding model trained on more than 90 billion words de-identified clinical note text. p. 142

Appendix A (Continued)

- GNN** (graph neural network) A deep neural network representing data as graphs. p. 23
- hallucination** Erroneous and nonfactual text automatically generated by a model. p. 12
- hallucination rate** The rate at which the output generates hallucinated text. p. 20
- hyperparameter** Any setting, or a “knob for tweaking”, that remains constant during both training and testing of the machine learning algorithm whose purpose is to increase better performance. p. 67
- 2010 i2b2** An i2b2 (Informatics for Integrating Biology and the Bedside) challenge on concepts, assertions, and relations in clinical text. p. 24
- knowledge graph** A graph representing a knowledge base of related entities. p. 29
- k^{th} order neighbor set** The subgraph induced as a set of vertexes that are at exactly k hops from a node. p. 91
- Levenshtein edit distance** A metric for measuring the sequence of characters in strings. p. 60
- MDS** (multi-document summarization) A task consisting of automatically generating a summary from multiple documents instead of one. p. 17
- MedLINE** An online life sciences bibliographic database containing journal references. p. 24
- MedSecId** A publicly available set of 2,002 fully annotated medical notes from the MIMIC-III. p. v
- MIMIC-III** (Medical Information Mart for Intensive Care III) a large freely accessible hospital database of ICU data from the Beth Israel Deaconess Medical Center in Boston, Massachusetts p. 4
- note antecedent** A clinical medical note containing the original information that contributed to the discharge summary. p. 53
- NER** (named entity recognition) Named entities are proper nouns or anything that can be named, such as organizations, a person’s name, locations, etc. p. 18
- network neighborhood** The subgraph induced as a set of vertexes adjacent to, and within N hops from the target node. p. 82
- n -gram** A sequence of words or symbols. p. 19
- note match** A “tie” between two spans of semantically similar or copied text segments between a note antecedent and a discharge summary. p. 54
- nominalization** Using the root verb form of a noun in place of a noun. p. 74
- edge node parent** The node at the source end of a directed edge e in the alignment graph. p. 108
- function tag** The identifier of a semantic function of a role, such as location, cause, time, purpose. p. 74
- predicate** A PropBank entry of a verb and its syntactic arguments that describes an event. p. 75

Appendix A (Continued)

- role** The role played by a sense of a predicate’s verb’s grammatical subject or object. p. 74
- roleset** A grouping of roles that make up a particular sense of a verb belonging to a predicate. p. 74
- Penman** A flat text representation widely used for abstract meaning representation graphs. p. 13
- predicating element** Verb or “verb-like” main event or action core to an AMR or AMR subtree. Examples include verbs (**chase-01**), adjectives (**attract-01**) and nominalizations (the root **destroy-01** is used in place of “destruction”). p. 74
- PropBank** Is a verb-oriented lexicon database consisting of its word sense as a “frame”. p. 21
- reentrancy** Nodes with more than one parent, or duplicate variable name in Penman notation. p. 102
- reference summary** The human written summary from which to learn. p. 20
- reduced graph** The summary to source aligned graph with 0-flow edges (and their orphaned nodes) removed. p. 136
- role edge** The edge that connects each abstract meaning representation graph (source and summary) as a disconnected component. p. 11
- Rouge** Recall-Oriented Understudy for Gisting Evaluation: A measure that assesses the quality of generated summaries. This quality is assessed by comparing to summaries generated by humans using overlapping n -grams counts or word pairs. p. 19
- SBERT** (Sentence-BERT) A siamese network model that captures sentinel semantic similarity. p. 71
- SDS** (single-document summarization) A task consisting of automatically generated a summary from a single document. p. 18
- section type** The identifier of a clinical note section, such as *History of Present Illness*. p. 2
- sentence matching algorithm** An algorithm that matches source with summary sentence candidates. p. 135
- sentence skew** Scales a sentence node and attribute semantic similarity by the linear factor sentence dampen (γ). p. 96
- Smatch** An evaluation metric for abstract meaning representation similarity that uses a greedy feature overlap method. p. 30
- summarization algorithm** The algorithm that uses the alignment graph algorithm output to learn to summarize using a graph neural network. p. 11
- summary component** The abstract meaning representation graph generated from the reference summary. p. 6
- summarizer** The algorithm and trained model used to summarize AMR graphs. p. 172
- SimpleNLG** A library used to generate natural language using a realizer. p. 23

Appendix A (Continued)

- SOAP** Subjective, Objective, Assessment and Plan. A structured format for healthcare workers to document encounters and a clinical cognitive framework. p. 26
- source component** The abstract meaning representation graph generated from the source text. p. 6
- source section dataset** The dataset used by the source section model created by aligning admissions. p. 130
- source section model** The model that classifies note antecedent sentences with a discharge summary section, or no section. p. 128
- source text** The document or text to be summarized. p. 1
- system summary** The machine generated summary. p. 20
- traceable** If the summary can be traced back to its source content. A summary is traceable if every sentence in it can be traced back from the source it was responsible. p. 1
- TAT** (text-to-graph aligned token) an alignment from a token’s index in the sentence to a node in the abstract meaning representation graph p. 81
- UI Health dataset** An IRB approved (protocol #2024-0109) private dataset of 11,001 admissions and 607,872 notes, which include daily progress, radiology, ECG and a variety of other notes from the University of Illinois Chicago hospital. p. 13
- Wlk** Weisfeiler-Leman Bamboo. p. 115
- wordpiece** Token sub-units with associated vectors and provided by the model’s tokenizer. p. 66

Appendix B

HYPERPARAMETERS

- τ_α The lowest value threshold an alignment capacity can be before being set to zero. (`alignment.align_min_capacity_cutoff` in the configuration.) p. 109
- δ The number of layers or node depth used in the GNN. (`graph_network.layer_depth` in the configuration.) p. 174
- Λ an array of weights, each of which is used to scale the embeddings of the k^{th} order neighbor set neighbors (`capacity_calculator.neighbor_embedding_weights` in the configuration.) p. 91
- τ transition and compression parameters to sigmoid function that adjusts the capacity values of alignment edges up or down (`capacity_calculator.neighbor_skew` in the configuration.) p. 95
- α_t This is used to scale the document positional embedding component. (`document_matcher.position_scale` in the configuration.) p. 67
- ρ_c The embedding weight of a role node's embedding. (`capacity_calculator.concept_embedding_role_weights.neighbor` in the configuration.) p. 89
- ρ_r The embedding weight of a role embedding. (`capacity_calculator.concept_embedding_role_weights.role` in the configuration.) p. 89
- τ_p The lowest value an alignment capacity can be before its parent node's descendants role edges are set to zero. (`align.role_min_flow_cutoff` in the configuration.) p. 110
- γ The slope for the linear dampening of nodes under a sentence by sentence cosine similarity, a number in range $[0, 1]$. The higher the value the lower the sentence similarity, which leads to lower concept and attribute node similarities. (`capacity_calculator.sentence_dampen` in the configuration.) p. 96
- μ_s The summary sentence flow distribution threshold cut-off for adding sentence matches in the source/summary match algorithm. (`calsum_sent_matcher.min_sent_flow` in the configuration.) p. 137
- τ_δ The threshold cut-off for creating component alignment edges during graph creation. (`graph_attrib_context.similarity_threshold` in the configuration.) p. 104

Appendix C

GOLD DISCHARGE SUMMARY

Discharge Summary by [**Doctor First Name**] [**Doctor Last Name**], MD at [**Date**] 6:00 AM

Author: [**Doctor First Name**] [**Doctor Last Name**], MD Service: Neuro Critical Care Author Type: Resident
 Filed: [**Date**] 6:09 PM Date of Service: [**Date**] 6:00 AM Status: Attested
 Editor: [**Doctor First Name**] [**Doctor Last Name**], MD (Resident) Cosigner: [**Doctor First Name**] [**Doctor Last Name**], MD at [**Date**] 12:45 PM
 Attestation signed by [**Doctor First Name**] [**Doctor Last Name**], MD at [**Date**] 12:45 PM Stroke attending: I have reviewed the above discharge summary and agree with the assessment.

Pt eloped before I could staff the pt.

UNIVERSITY OF ILLINOIS HOSPITAL AND CLINICS
 Discharge Summary

Patient: [**First Name**] [**Last Name**]
 Admission Date: [**Date**]
 Discharge Date: [**Date**]
 Discharge Disposition: Left Against Medical Advice
 Discharge Service: Stroke
 Discharge Attending: [**Doctor First Name**] [**Doctor Last Name**], MD
 Primary Diagnosis: Acute R medial temporal and internal capsule/thalamic stroke

Other Active Diagnoses
 Diagnosis Date Noted POA
 - Troponin level elevated [**Date**] Yes
 Priority: High
 - Stroke (CMS/HCC) [**Date**] Yes

Hospital Course

HPI: [**First Name**] [**Last Name**] 47 y.o. PMH HTN, HLD, previous stroke [**Date**] (R-MCA watershed), presenting to ED with L sided weakness, LLE numbness, and blurry vision b/l. Patient reports symptoms started acutely on Friday [**Date**] around 4pm while driving causing him to have to pull over. He decided to try to sleep it off. After waking up the following morning with no improvement, he went to [**Hospital Name**] ED where he was seen by neurology, but left AMA as he felt he was being asked the same questions repeatedly and nothing was getting done. CTH at [**Hospital Name**] was without ICH, reportedly showed wedge shaped hypodensity in frontal lobe

Appendix C (Continued)

likely from chronic infarct. The patient reports that the left sided weakness has improved somewhat today, but he still endorses b/l blurry vision with constant white floaters, as well as numbness/tingling in his LLE. NIHSS is 2 (LUQ quadrantopia, LLE numbness). CTA Hmultifocal narrowing of b/l ACAs as well as R M1 focal narrowing. B w/o acute ischemia in medial R temporal lobe involving the posterior limb of R internal capsule as well as possibly the thalamus, and redemonstrating old R frontal ACA-MCA watershed infarct and old L occipital cortical infarct. Patient to be admitted to stepdown under stroke service.

Patient with poorly controlled HTN and HLD, not taking the past 5 months she says he had trouble getting primary care appointment/prescription renewals. SBP gets up to 200s per patient. Labs today significant for Troponin of 0.77-0.62. EKG showing T inversions in V5 and V6. Cardiology consulted. Patient denying CP at this time. Cardio recommending trending EKG/trop until downtrend and ordering Echo.

Patient appeared to have left before he was evaluated [**Date**] AM.

Pertinent Physical Exam At Time of Discharge

Physical Exam

PATIENT NOT EXAMINED PRIOR TO DISCHARGE

Test Results Pending At Discharge

Discharge Medications

No medications have been prescribed.

Issues Requiring Follow-Up

- Patient not evaluated prior to discharge

Outpatient Follow-Up Appointments

No future appointments.

Referrals

No orders of the defined types were placed in this encounter.

Completed Consults:

Consults Ordered This Encounter

Procedures

- Inpatient consult to General Neurology

- Inpatient consult to Cardiology

Figure 44: **Gold Discharge Summary.** The physician hand written gold UI Health de-identified discharge summary. White space that spanned longer than three lines was reduced to two.

Appendix D

NOTATION

The mathematical notation follows that of “Short Guide to Typesetting Math in NLP Papers” (Dyer et al., 2017):

\mathcal{G}, \mathcal{E} : upper case for graphs and sets.

$\mathcal{P}(v)$: upper case curly letters select nodes or edges in a graph.

e as in $e \in E$: lower case letters for variables (edge e is in the set of edges E).

\mathbf{e} : lower case bold for vectors.

\mathbf{s} : lower case bold italic letter for structured elements such as sentences.

α : lower case bold Greek letters for hyperparameters.

count: functions in lowercase latex characters.

\triangleq : a definition

Appendix E

COPYRIGHT NOTICIES

International Committee on Computational Linguistics

Paul Landes, Kunal Patel, Sean S. Huang, Adam Webb, Barbara Di Eugenio, and Cornelia Caragea. “A New Public Corpus for Clinical Section Identification: MedSecId”. In Proceedings of the 29th International Conference on Computational Linguistics (COLING, 2022).

Copyright of each paper stays with the respective authors (or their employers).

ISSN 2951-2093

Association for Computing Linguistics (ACL)

Paul Landes, Aaron Chaise, Kunal Patel, Sean Huang, and Barbara Di Eugenio. “Hospital Discharge Summarization Data Provenance”. In the 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks (BioNLP 2023).

Paul Landes, Barbara Di Eugenio, and Cornelia Caragea. “DeepZensols: A Deep Learning Natural Language Processing Framework for Experimentation and Reproducibility”. In the Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023).

ACL materials are Copyright © 1963–2024 ACL; other materials are copyrighted by their respective copyright holders. Materials prior to 2016 here are licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 3.0 International License. Permission is granted to make copies for the purposes of teaching and research. Materials published in or after 2016 are licensed on a Creative Commons Attribution 4.0 International License.



ELRA Language Resource Association

Paul Landes and Barbara Di Eugenio. “CALAMR: Component ALignment for Abstract Meaning Representation”. In the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING, 2024).

Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)

Copyright ELRA Language Resources Association (ELRA), 2024
These proceedings are licensed under a Creative Commons Attribution-NonCommercial 4.0 International License (CC BY-NC 4.0)

ISBN 978-2-493814-10-4
ISSN 2951-2093 (COLING); 2522-2686 (LREC)

CITED LITERATURE

- Acharya, S.: Generating Personalized Hospital-Stay Summaries for Patients. thesis, University of Illinois at Chicago, Chicago, IL, 2019.
- Acharya, S., Di Eugenio, B., Boyd, A., Cameron, R., Dunn Lopez, K., Martyn-Nemeth, P., Dickens, C., and Ardati, A.: Towards Generating Personalized Hospitalization Summaries. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop, pages 74–82. Association for Computational Linguistics, 2018.
- Adams, G., Alsentzer, E., Ketenci, M., Zucker, J., and Elhadad, N.: What’s in a Summary? Laying the Groundwork for Advances in Hospital-Course Summarization. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4794–4811, 2021.
- Alberti, M., Pondenkandath, V., Würsch, M., Ingold, R., and Liwicki, M.: DeepDIVA: A Highly-Functional Python Framework for Reproducible Experiments. In 2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR), pages 423–428, 2018.
- Alsentzer, E. and Kim, A.: Extractive Summarization of EHR Discharge Notes. arXiv: 1810.12085 (Only available as arXiv preprint), 2018.
- Alsentzer, E., Murphy, J., Boag, W., Weng, W.-H., Jindi, D., Naumann, T., and McDermott, M.: Publicly Available Clinical BERT Embeddings. In Proceedings of the 2nd Clinical Natural Language Processing Workshop, pages 72–78. Association for Computational Linguistics, 2019.
- Badjatiya, P., Kurisinkel, L. J., Gupta, M., and Varma, V.: Attention-Based Neural Text Segmentation. In Advances in Information Retrieval, eds. G. Pasi, B. Piwowarski, L. Az-zopardi, and A. Hanbury, Lecture Notes in Computer Science, pages 180–193. Springer International Publishing, 2018.
- Banarescu, L., Bonial, C., Cai, S., Georgescu, M., Griffitt, K., Hermjakob, U., Knight, K., Koehn, P., Palmer, M., and Schneider, N.: Abstract Meaning Representation for Sem-banking. In Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse, pages 178–186. Association for Computational Linguistics, 2013.

- Banarescu, L., Bonial, C., Cai, S., Georgescu, M., Griffitt, K., Hermjakob, U., Knight, K., Koehn, P., Palmer, M., and Schneider, N.: Abstract Meaning Representation (AMR) 1.2.6 Specification. (Only available online), 2024.
- Barnes, S. L., Robinson, B. R. H., Richards, J. T., Zimmerman, C. E., Pritts, T. A., Tsuei, B. J., Butler, K. L., Muskat, P. C., Davis, K., and Johannigman, J. A.: The devil is in the details: Maximizing revenue for daily trauma care. Surgery, 144(4):670–676, 2008.
- Bazaraa, M. S., Jarvis, J. J., and Sherali, H. D.: Linear Programming and Network Flows. John Wiley & Sons, 2010.
- Beltagy, I., Peters, M. E., and Cohan, A.: Longformer: The Long-Document Transformer. arXiv: 2004.05150 (Only available as arXiv preprint), 2020.
- Bengio, Y., Simard, P., and Frasconi, P.: Learning long-term dependencies with gradient descent is difficult. IEEE Transactions on Neural Networks, 5(2):157–166, 1994.
- Berger, A. L., Pietra, V. J. D., and Pietra, S. A. D.: A maximum entropy approach to natural language processing. Computational Linguistics, 22(1):39–71, 1996.
- Bergstra, J., Yamins, D., and Cox, D.: Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures. In Proceedings of the 30th International Conference on Machine Learning, pages 115–123. PMLR, 2013.
- Bevilacqua, M., Blloshmi, R., and Navigli, R.: One SPRING to Rule Them Both: Symmetric AMR Semantic Parsing and Generation without a Complex Pipeline. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 12564–12573, 2021.
- Blei, D. M., Ng, A. Y., and Jordan, M. I.: Latent Dirichlet Allocation. Journal of Machine Learning Research, 3:993–1022, 2003.
- Blloshmi, R., Tripodi, R., and Navigli, R.: XL-AMR: Enabling Cross-Lingual AMR Parsing with Transfer Learning Techniques. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 2487–2500. Association for Computational Linguistics, 2020.
- Blodgett, A. and Schneider, N.: Probabilistic, Structure-Aware Algorithms for Improved Variety, Accuracy, and Coverage of AMR Alignments. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and

the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3310–3321. Association for Computational Linguistics, 2021.

- Bodenreider, O.: The Unified Medical Language System (UMLS): Integrating biomedical terminology. Nucleic Acids Research, 32:D267–D270, 2004.
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T.: Enriching Word Vectors with Subword Information. Transactions of the Association for Computational Linguistics, 5:135–146, 2017.
- Bonial, C., Donatelli, L., Abrams, M., Lukin, S. M., Tratz, S., Marge, M., Artstein, R., Traum, D., and Voss, C.: Dialogue-AMR: Abstract Meaning Representation for Dialogue. In Proceedings of the 12th Language Resources and Evaluation Conference, pages 684–695. European Language Resources Association, 2020.
- Bonial, C., Lukin, S. M., Doughty, D., Hill, S., and Voss, C.: InfoForager: Leveraging Semantic Search with AMR for COVID-19 Research. In Proceedings of the Second International Workshop on Designing Meaning Representations, pages 67–77. Association for Computational Linguistics, 2020.
- Boyd, A. D., Dunn Lopez, K., Lugaresi, C., Macieira, T., Sousa, V., Acharya, S., Balasubramanian, A., Roussi, K., Keenan, G. M., Lussier, Y. A., Li, J. J., Burton, M., and Di Eugenio, B.: Physician nurse care: A new use of UMLS to measure professional contribution: Are we talking about the same patient a new graph matching algorithm? International Journal of Medical Informatics, 113:63–71, 2018.
- Cai, D. and Lam, W.: AMR Parsing via Graph-Sequence Iterative Inference. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1290–1301. Association for Computational Linguistics, 2020.
- Cai, J. Z., Wright-Bettner, K., Palmer, M., Savova, G. K., and Martin, J. H.: Adapting Abstract Meaning Representation Parsing to the Clinical Narrative – the SPRING THYME parser. arXiv: 2405.09153 (Only available as arXiv preprint), 2024.
- Cai, S. and Knight, K.: Smatch: An Evaluation Metric for Semantic Feature Structures. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 748–752. Association for Computational Linguistics, 2013.

- Chakraborty, A., Paranjape, B., Kakarla, S., and Ganguly, N.: Stop clickbait: Detecting and preventing clickbaits in online news media. In Proceedings of the 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM '16, pages 9–16. IEEE Press, 2016.
- Chaturvedi, R., ., S., Dhani, J. S., Joshi, A., Khanna, A., Tomar, N., Duari, S., Khurana, A., and Bhatnagar, V.: Divide and Conquer: From Complexity to Simplicity for Lay Summarization. In Proceedings of the First Workshop on Scholarly Document Processing, pages 344–355. Association for Computational Linguistics, 2020.
- Chen, L., Kyng, R., Liu, Y. P., Peng, R., Gutenberg, M. P., and Sachdeva, S.: Maximum Flow and Minimum-Cost Flow in Almost-Linear Time. In 2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS), pages 612–623, 2022.
- Chu, X., Fan, X., Yao, D., Zhu, Z., Huang, J., and Bi, J.: Cross-Network Embedding for Multi-Network Alignment. In The World Wide Web Conference on - WWW '19, pages 273–284. ACM Press, 2019.
- Cohen, R., Elhadad, M., and Elhadad, N.: Redundancy in electronic health record corpora: Analysis, impact on text mining performance and mitigation strategies. BMC bioinformatics, 14:10, 2013.
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C.: Introduction to Algorithms. MIT Press and McGraw-Hill, second edition, 2001.
- Damm, H., Pakull, T. M. G., Eryılmaz, B., Becker, H., Idrissi-Yaghir, A., Schäfer, H., Schultenkämper, S., and Friedrich, C. M.: WisPerMed at "Discharge Me!": Advancing Text Generation in Healthcare with Large Language Models, Dynamic Expert Selection, and Priming Techniques on MIMIC-IV. arXiv: 2405.11255 (Only available as arXiv preprint), 2024.
- Dang, H. T. and Owczarzak, K.: Overview of the TAC 2008 Update Summarization Task. Proceedings of Text Analysis Conference (TAC), 2008.
- de Bruijn, B., Cherry, C., Kiritchenko, S., Martin, J., and Zhu, X.: NRC at i2b2: One challenge, three practical tasks, nine statistical systems, hundreds of clinical records, millions of useful features. In Proceedings of the 2010 I2b2/VA Workshop on Challenges in Natural Language Processing for Clinical Data, 2010.

- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, 2019.
- Dohare, S., Karnick, H., and Gupta, V.: Text Summarization using Abstract Meaning Representation. arXiv: 1706.01678 (Only available as arXiv preprint), 2017.
- Dyer, C., Gimpel, K., and Smith, N. A.: A Short Guide to Typesetting Math in NLP Papers. (Only available on the web), 2017.
- Faria, D., Pesquita, C., Santos, E., Palmonari, M., Cruz, I. F., and Couto, F. M.: The AgreementMakerLight Ontology Matching System. In On the Move to Meaningful Internet Systems: OTM 2013 Conferences, eds. R. Meersman, H. Panetto, T. Dillon, J. Eder, Z. Bellahsene, N. Ritter, P. De Leenheer, and D. Dou, Lecture Notes in Computer Science, pages 527–541. Springer, 2013.
- Flanigan, J., Thomson, S., Carbonell, J., Dyer, C., and Smith, N. A.: A Discriminative Graph-Based Parser for the Abstract Meaning Representation. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1426–1436. Association for Computational Linguistics, 2014.
- Ford, L. R. and Fulkerson, D. R.: Flows in networks. In Flows in Networks, Princeton Landmarks in Mathematics and Physics, page 212. Princeton University Press, 1962.
- Fossati, D., Ghidoni, G., Di Eugenio, B., Cruz, I., Xiao, H., and Subba, R.: The problem of ontology alignment on the Web: A first report. In Proceedings of the 2nd International Workshop on Web as Corpus, 2006.
- Fu, Q., Song, L., Du, W., and Zhang, Y.: End-to-End AMR Coreference Resolution. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4204–4214. Association for Computational Linguistics, 2021.
- Galitsky, B.: Summarized Logical Forms Based on Abstract Meaning Representation and Discourse Trees. In Artificial Intelligence for Customer Relationship Management: Keeping Customers Informed, ed. B. Galitsky, Human-Computer Interaction Series, pages 151–191. Springer International Publishing, 2020.

- Gallagher, B. D., Nematollahi, S., Park, H., and Kurra, S.: Comparing Students' Clinical Grades to Scores on a Standardized Patient Note-Writing Task. Journal of General Internal Medicine, 35(11):3243–3247, 2020.
- Gao, Y., Dligach, D., Miller, T., Xu, D., Churpek, M. M. M., and Afshar, M.: Summarizing Patients' Problems from Hospital Progress Notes Using Pre-trained Sequence-to-Sequence Models. In Proceedings of the 29th International Conference on Computational Linguistics, pages 2979–2991. International Committee on Computational Linguistics, 2022.
- Gao, Y., Liu, Y. P., and Peng, R.: Fully Dynamic Electrical Flows: Sparse Maxflow Faster Than Goldberg-Rao. In 2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS), pages 516–527, 2022.
- Gardner, M., Grus, J., Neumann, M., Tafjord, O., Dasigi, P., Liu, N. F., Peters, M. E., Schmitz, M., and Zettlemoyer, L.: AllenNLP: A Deep Semantic Natural Language Processing Platform. In Proceedings of Workshop for NLP Open Source Software (NLP-OSS), pages 1–6, 2018.
- Gatt, A. and Reiter, E.: SimpleNLG: A realisation engine for practical applications. In Proceedings of the 12th European Workshop on Natural Language Generation - ENLG '09, pages 90–93. Association for Computational Linguistics, 2009.
- Gaussier, E.: Flow Network Models for Word Alignment and Terminology Extraction from Bilingual Corpora. In 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 1, pages 444–450. Association for Computational Linguistics, 1998.
- Goldberg, A. V. and Tarjan, R. E.: A new approach to the maximum-flow problem. Journal of the ACM, 35(4):921–940, 1988.
- Goyal, P. and Ferrara, E.: Graph embedding techniques, applications, and performance: A survey. Knowledge-Based Systems, 151:78–94, 2018.
- Graves, A. and Schmidhuber, J.: Framewise phoneme classification with bidirectional LSTM networks. In Proceedings. 2005 IEEE International Joint Conference on Neural Networks, volume 4, pages 2047–2052. IEEE, 2005.

- Gu, X., Mao, Y., Han, J., Liu, J., Wu, Y., Yu, C., Finnie, D., Yu, H., Zhai, J., and Zukoski, N.: Generating Representative Headlines for News Stories. In Proceedings of The Web Conference 2020, WWW '20, pages 1773–1784. Association for Computing Machinery, 2020.
- Guo, J., Fan, Y., Ai, Q., and Croft, W. B.: Semantic Matching by Non-Linear Word Transportation for Information Retrieval. In Proceedings of the 25th ACM International on Conference on Information and Knowledge Management, CIKM '16, pages 701–710. ACM Press, 2016.
- Gupta, A., Kaur, M., Bajaj, A., and Khanna, A.: Entailment and Spectral Clustering based Single and Multiple Document Summarization. International Journal of Intelligent Systems and Applications, 11(4):39–51, 2019.
- Habernal, I. and Gurevych, I.: Argumentation Mining in User-Generated Web Discourse. Computational Linguistics, 43(1):125–179, 2016.
- Harris, T. and Ross, F.: Fundamentals of a method for evaluating rail net capacities. Rand Corporation, Santa Monica, CA, 1955.
- Heidary, G., Zamanifar, K., Nematbakhsh, N., and Mardukhi, F.: How to discover a semantic Web service by knowing its functionality parameters. In 2010 2nd International Conference on Software Technology and Engineering, volume 2, pages V2–194–V2–198, 2010.
- Heinecke, J. and Shimorina, A.: Multilingual Abstract Meaning Representation for Celtic Languages. In Proceedings of the 4th Celtic Language Technology Workshop within LREC2022, pages 1–6. European Language Resources Association, 2022.
- Hirohata, K., Okazaki, N., Ananiadou, S., and Ishizuka, M.: Identifying Sections in Scientific Abstracts using Conditional Random Fields. In Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-I, 2008.
- Hirsch, J. S., Tanenbaum, J. S., Lipsky Gorman, S., Liu, C., Schmitz, E., Hashorva, D., Ervits, A., Vawdrey, D., Sturm, M., and Elhadad, N.: HARVEST, a longitudinal patient record summarizer. Journal of the American Medical Informatics Association, 22(2):263–274, 2015.
- Hochreiter, S. and Schmidhuber, J.: Long short-term memory. Neural computation, 9(8):1735–1780, 1997.

- Hopcroft, J. E. and Karp, R. M.: An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs. SIAM Journal on Computing, 2(4):225–231, 1973.
- Irving, R. W., Leather, P., and Gusfield, D.: An efficient algorithm for the “optimal” stable marriage. Journal of the ACM, 34(3):532–543, 1987.
- Ivanov, P., Bichsel, B., Mustafa, H., Kahles, A., Rättsch, G., and Vechev, M.: AStarix: Fast and Optimal Sequence-to-Graph Alignment. In Research in Computational Molecular Biology: 24th Annual International Conference, RECOMB 2020, Padua, Italy, May 10–13, 2020, Proceedings, pages 104–119. Springer-Verlag, 2020.
- Jin, H. and Wan, X.: Abstractive Multi-Document Summarization via Joint Learning with Single-Document Summarization. In Findings of the Association for Computational Linguistics: EMNLP 2020, pages 2545–2554. Association for Computational Linguistics, 2020.
- Johnson, A. E. W., Pollard, T. J., Shen, L., Lehman, L.-w. H., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Celi, L. A., and Mark, R. G.: MIMIC-III, a freely accessible critical care database. Scientific Data, 3(1):1–9, 2016.
- Kasper, R. T.: A Flexible Interface for Linking Applications to Penman’s Sentence Generator. In Speech and Natural Language: Proceedings of a Workshop Held at Philadelphia, Pennsylvania, February 21-23, 1989, 1989.
- Kingsbury, P. and Palmer, M.: From TreeBank to PropBank. In Proceedings of the Third International Conference on Language Resources and Evaluation (LREC’02). European Language Resources Association (ELRA), 2002.
- Kleinberg, J. and Tardos, E.: Algorithm Design. Addison-Wesley Longman Publishing Co., Inc., 2005.
- Klie, J.-C., Bugert, M., Boullosa, B., Eckart de Castilho, R., and Gurevych, I.: The INCEpTION Platform: Machine-Assisted and Knowledge-Oriented Interactive Annotation. In Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations, pages 5–9. Association for Computational Linguistics, 2018.
- Knight, K., Badarau, B., Baranescu, L., Bonial, C., Bardocz, M., Griffitt, K., Hermjakob, U., Marcu, D., Palmer, M., O’Gorman, T., and Schneider, N.: Abstract Meaning Representation (AMR) Annotation Release 3.0, 2021.

- Knight, K. and Marcu, D.: Statistics-Based Summarization - Step One: Sentence Compression. In Proceedings of the Seventeenth National Conference on Artificial Intelligence and Twelfth Conference on Innovative Applications of Artificial Intelligence, pages 703–710. AAAI Press, 2000.
- Knight, K. and Marcu, D.: Summarization beyond sentence extraction: A probabilistic approach to sentence compression. Artificial Intelligence, 139(1):91–107, 2002.
- Koutra, D., Tong, H., and Lubensky, D.: BIG-ALIGN: Fast Bipartite Graph Alignment. In 2013 IEEE 13th International Conference on Data Mining, pages 389–398, 2013.
- Kraljevic, Z., Searle, T., Shek, A., Roguski, L., Noor, K., Bean, D., Mascio, A., Zhu, L., Folarin, A. A., Roberts, A., Bendayan, R., Richardson, M. P., Stewart, R., Shah, A. D., Wong, W. K., Ibrahim, Z., Teo, J. T., and Dobson, R. J. B.: Multi-domain clinical natural language processing with MedCAT: The Medical Concept Annotation Toolkit. Artificial Intelligence in Medicine, 117:102083, 2021.
- Krauss, T., McCollum, J., Pendery, C., Litwin, S., and Michaels, A. J.: Solving the max-flow problem on a quantum annealing computer. IEEE Transactions on Quantum Engineering, 1:1–10, 2020.
- Krippendorff, K.: Reliability in Content Analysis. Human Communication Research, 30(3):411–433, 2004.
- Krippendorff, K.: Agreement and Information in the Reliability of Coding. Communication Methods and Measures, 5(2):93–112, 2011.
- Krishna, K., Khosla, S., Bigham, J., and Lipton, Z. C.: Generating SOAP Notes from Doctor-Patient Conversations Using Modular Summarization Techniques. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4958–4972. Association for Computational Linguistics, 2021.
- Kuchaiev, O., Milenković, T., Memišević, V., Hayes, W., and Pržulj, N.: Topological network alignment uncovers biological function and phylogeny. Journal of The Royal Society Interface, 7(50):1341–1354, 2010.
- Kupiec, J., Pedersen, J., and Chen, F.: A trainable document summarizer. In Proceedings of the 18th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval - SIGIR '95, pages 68–73. ACM Press, 1995.

- Kusner, M. J., Sun, Y., Kolkin, N. I., and Weinberger, K. Q.: From Word Embeddings to Document Distances. In Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15, pages 957–966. JMLR.org, 2015.
- Landes, P., Chaise, A., Patel, K., Huang, S., and Di Eugenio, B.: Hospital Discharge Summarization Data Provenance. In The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks, pages 439–448. Association for Computational Linguistics, 2023.
- Landes, P. and Di Eugenio, B.: CALAMR: Component ALignment for Abstract Meaning Representation. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), eds. N. Calzolari, M.-Y. Kan, V. Hoste, A. Lenci, S. Sakti, and N. Xue, pages 2622–2637. ELRA and ICCL, 2024.
- Landes, P., Di Eugenio, B., and Caragea, C.: DeepZensols: A Deep Learning Natural Language Processing Framework for Experimentation and Reproducibility. In Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023), eds. L. Tan, D. Milajevs, G. Chauhan, J. Gwinnup, and E. Rippeth, pages 141–146. Association for Computational Linguistics, 2023.
- Landes, P., Patel, K., Huang, S. S., Webb, A., Di Eugenio, B., and Caragea, C.: A New Public Corpus for Clinical Section Identification: MedSecId. In Proceedings of the 29th International Conference on Computational Linguistics, pages 3709–3721. International Committee on Computational Linguistics, 2022.
- Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., and Kang, J.: BioBERT: A pre-trained biomedical language representation model for biomedical text mining. Bioinformatics, 36(4):1234–1240, 2020.
- Levenshtein, V. I.: Binary codes capable of correcting deletions, insertions, and reversals. In Soviet Physics Doklady, volume 10, pages 707–710, 1966.
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., and Zettlemoyer, L.: BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7871–7880. Association for Computational Linguistics, 2020.

- Li, C., Liu, Y., Liu, F., Zhao, L., and Weng, F.: Improving Multi-documents Summarization by Sentence Compression based on Expanded Constituent Parse Trees. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 691–701. Association for Computational Linguistics, 2014.
- Liao, K., Lebanoff, L., and Liu, F.: Abstract Meaning Representation for Multi-Document Summarization. In Proceedings of the 27th International Conference on Computational Linguistics, pages 1178–1190. Association for Computational Linguistics, 2018.
- Likert, R.: A technique for the measurement of attitudes. Archives of Psychology, 22 140:55–55, 1932.
- Lim, J., Oh, D., Jang, Y., Yang, K., and Lim, H.: I Know What You Asked: Graph Path Learning using AMR for Commonsense Reasoning. In Proceedings of the 28th International Conference on Computational Linguistics, pages 2459–2471. International Committee on Computational Linguistics, 2020.
- Lin, C.-Y.: ROUGE: A Package for Automatic Evaluation of Summaries. In Text Summarization Branches Out, pages 74–81. Association for Computational Linguistics, 2004.
- Lin, C.-Y. and Hovy, E.: Automatic evaluation of summaries using N-gram co-occurrence statistics. In Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1, NAACL '03, pages 71–78. Association for Computational Linguistics, 2003.
- Liu, F., Shareghi, E., Meng, Z., Basaldella, M., and Collier, N.: Self-Alignment Pre-training for Biomedical Entity Representations. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4228–4238. Association for Computational Linguistics, 2021.
- Liu, F., Flanigan, J., Thomson, S., Sadeh, N., and Smith, N. A.: Toward Abstractive Summarization Using Semantic Representations. In Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1077–1086. Association for Computational Linguistics, 2015.

- Loper, E., Yi, S.-T., and Palmer, M.: Combining Lexical Resources: Mapping Between PropBank and VerbNet. In Proceedings of the 7th International Workshop on Computational Semantics, 2007.
- Luhn, H. P.: The Automatic Creation of Literature Abstracts. IBM Journal of Research and Development, 2(2):159–165, 1958.
- Lyu, C. and Titov, I.: AMR Parsing as Graph Prediction with Latent Alignment. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 397–407. Association for Computational Linguistics, 2018.
- Magnanti, T. L. and Wolsey, L. A.: Optimal trees. In Handbooks in Operations Research and Management Science, volume 7 of Network Models, pages 503–615. Elsevier, 1995.
- Malod-Dognin, N. and Pržulj, N.: L-GRAAL: Lagrangian graphlet-based network aligner. Bioinformatics, 31(13):2182–2189, 2015.
- Manning, E., Wein, S., and Schneider, N.: A Human Evaluation of AMR-to-English Generation Systems. In Proceedings of the 28th International Conference on Computational Linguistics, pages 4773–4786. International Committee on Computational Linguistics, 2020.
- Mattsson, M.: Object-Oriented Frameworks. mathesis, Lund University, Ronneby, Sweden, 1996.
- Maynez, J., Narayan, S., Bohnet, B., and McDonald, R.: On Faithfulness and Factuality in Abstractive Summarization. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1906–1919. Association for Computational Linguistics, 2020.
- McDonald, R., Crammer, K., and Pereira, F.: Online Large-margin Training of Dependency Parsers. In Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics, ACL '05, pages 91–98. Association for Computational Linguistics, 2005.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J.: Efficient estimation of word representations in vector space. arXiv: 1301.3781 (Only available as arXiv preprint), 2013.

- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J.: Distributed Representations of Words and Phrases and their Compositionality. In Advances in Neural Information Processing Systems, volume 26. Curran Associates, Inc., 2013.
- Miller, D.: Leveraging BERT for Extractive Text Summarization on Lectures. arXiv: 1906.04165 (Only available as arXiv preprint), 2019.
- Molla, D. and Santiago-Martinez, M. E.: Development of a Corpus for Evidence Based Medicine Summarisation. In Proceedings of the Australasian Language Technology Association Workshop 2011, pages 86–94, 2011.
- Naseem, T., Blodgett, A., Kumaravel, S., O’Gorman, T., Lee, Y.-S., Flanigan, J., Astudillo, R., Florian, R., Roukos, S., and Schneider, N.: DocAMR: Multi-Sentence AMR Representation and Evaluation. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3496–3505. Association for Computational Linguistics, 2022.
- Ni, C.-C., Lin, Y.-Y., Gao, J., and Gu, X.: Network Alignment by Discrete Ollivier-Ricci Flow. In Graph Drawing and Network Visualization, eds. T. Biedl and A. Kerren, Lecture Notes in Computer Science, pages 447–462. Springer International Publishing, 2018.
- Ning, Q., Wu, H., Dasigi, P., Dua, D., Gardner, M., Iv, R. L. L., Marasović, A., and Nie, Z.: Easy, Reproducible and Quality-Controlled Data Collection with CROWDAQ. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 127–134, 2020.
- Opitz, J., Daza, A., and Frank, A.: Weisfeiler-Leman in the Bamboo: Novel AMR Graph Metrics and a Benchmark for AMR Graph Similarity. Transactions of the Association for Computational Linguistics, 9:1425–1441, 2021.
- Opitz, J. and Frank, A.: Better Smatch = Better Parser? AMR evaluation is not so simple anymore. In Proceedings of the 3rd Workshop on Evaluation and Comparison of NLP Systems, pages 32–43. Association for Computational Linguistics, 2022.
- Palmer, M., Gildea, D., and Kingsbury, P.: The Proposition Bank: An Annotated Corpus of Semantic Roles. Computational Linguistics, 31(1):71–106, 2005.
- Pang, B. and Lee, L.: Seeing Stars: Exploiting Class Relationships for Sentiment Categorization with Respect to Rating Scales. In Proceedings of the 43rd Annual Meeting of

the Association for Computational Linguistics (ACL'05), pages 115–124. Association for Computational Linguistics, 2005.

Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J.: Bleu: A Method for Automatic Evaluation of Machine Translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318. Association for Computational Linguistics, 2002.

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S.: PyTorch: An Imperative Style, High-Performance Deep Learning Library. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019.

Pei, S., Yu, L., Yu, G., and Zhang, X.: Graph Alignment with Noisy Supervision. In Proceedings of the ACM Web Conference 2022, pages 1104–1114. ACM, 2022.

Pele, O. and Werman, M.: Fast and robust Earth Mover’s Distances. In 2009 IEEE 12th International Conference on Computer Vision, pages 460–467, 2009.

Pennington, J., Socher, R., and Manning, C.: Glove: Global Vectors for Word Representation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1532–1543. Association for Computational Linguistics, 2014.

Pham, H. V., Qian, S., Wang, J., Lutellier, T., Rosenthal, J., Tan, L., Yu, Y., and Nagappan, N.: Problems and opportunities in training deep learning software systems: An analysis of variance. In Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering, pages 771–783. Association for Computing Machinery, 2020.

Pivovarov, R. and Elhadad, N.: Automated methods for the summarization of electronic health records. Journal of the American Medical Informatics Association, 22(5):938–947, 2015.

Pomares-Quimbaya, A., Kreuzthaler, M., and Schulz, S.: Current approaches to identify sections within clinical narratives from electronic health records: A systematic review. BMC Medical Research Methodology, 19(1):155, 2019.

- Portet, F., Reiter, E., Gatt, A., Hunter, J., Sripada, S., Freer, Y., and Sykes, C.: Automatic generation of textual summaries from neonatal intensive care data. Artificial Intelligence, 173(7):789–816, 2009.
- Powsner, S. M. and Tufte, E. R.: Summarizing clinical psychiatric data. Psychiatric Services, 48(11):1458–1461, 1997.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J.: Exploring the limits of transfer learning with a unified text-to-text transformer. The Journal of Machine Learning Research, 21(1):140:5485–140:5551, 2020.
- Ramshaw, L. and Marcus, M.: Text Chunking using Transformation-Based Learning. In Third Workshop on Very Large Corpora. ACL, 1995.
- Ranjitha, N. S. and Kallimani, J. S.: Abstractive multi-document summarization. In 2017 International Conference on Advances in Computing, Communications and Informatics (ICACCI), pages 1690–1694, 2017.
- Rao, S., Marcu, D., Knight, K., and Daumé III, H.: Biomedical Event Extraction using Abstract Meaning Representation. In BioNLP 2017, pages 126–135. Association for Computational Linguistics, 2017.
- Rautiainen, M. and Marschall, T.: GraphAligner: Rapid and versatile sequence-to-graph alignment. Genome Biology, 21(1):253, 2020.
- Reimers, N. and Gurevych, I.: Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992. Association for Computational Linguistics, 2019.
- Reiter, E.: An Architecture for Data-to-Text Systems. In Proceedings of the Eleventh European Workshop on Natural Language Generation (ENLG 07), ed. S. Busemann, pages 97–104. DFKI GmbH, 2007.
- Sadoughi, N., Finley, G. P., Edwards, E., Robinson, A., Korenevsky, M., Brenndoerfer, M., Axtmann, N., Miller, M., and Suendermann-Oeft, D.: Detecting Section Boundaries in Medical Dictations: Toward Real-Time Conversion of Medical Dictations to Clinical Reports. In Proceedings of the 20th International Conference On Speech And Computer,

eds. A. Karpov, O. Jokisch, and R. Potapova, *Lecture Notes in Computer Science*, pages 563–573. Springer International Publishing, 2018.

Schluter, N.: The limits of automatic summarisation according to ROUGE. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 41–45. Association for Computational Linguistics, 2017.

Schrijver, A.: On the history of the transportation and maximum flow problems. *Mathematical Programming*, 91(3):437–445, 2002.

Schwartz, B. L.: Possible winners in partially completed tournaments. *SIAM Review*, 8(3):302–308, 1966.

See, A., Liu, P. J., and Manning, C. D.: Get To The Point: Summarization with Pointer-Generator Networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083. Association for Computational Linguistics, 2017.

Segura-Bedmar, I., Martínez, P., and Herrero-Zazo, M.: SemEval-2013 Task 9 : Extraction of Drug-Drug Interactions from Biomedical Texts (DDIExtraction 2013). In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pages 341–350. Association for Computational Linguistics, 2013.

Shing, H.-C., Shivade, C., Pourdamghani, N., Resnik, P., Oard, D., and Bhatia, P.: Towards Clinical Encounter Summarization: Learning to Compose Discharge Summaries from Prior Notes. arXiv: 2104.13498 (Only available as arXiv preprint), 2021.

Shivade, C., Malewadkar, P., Fosler-Lussier, E., and Lai, A. M.: Comparison of UMLS terminologies to identify risk of heart disease using clinical notes. *Journal of Biomedical Informatics*, 58 Suppl:S103–S110, 2015.

Sinsky, C., Colligan, L., Li, L., Prgomet, M., Reynolds, S., Goeders, L., Westbrook, J., Tutty, M., and Blike, G.: Allocation of Physician Time in Ambulatory Practice: A Time and Motion Study in 4 Specialties. *Annals of Internal Medicine*, 165(11):753, 2016.

Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A., and Potts, C.: Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. In *Proceedings of*

the 2013 Conference on Empirical Methods in Natural Language Processing, pages 1631–1642. Association for Computational Linguistics, 2013.

Sparck Jones, K.: A statistical interpretation of term specificity and its application in retrieval. Journal of documentation, 28(1):11–21, 1972.

Steinkamp, J., Kantrowitz, J. J., and Airan-Javia, S.: Prevalence and Sources of Duplicate Information in the Electronic Medical Record. JAMA Network Open, 5(9):e2233348, 2022.

Szubert, I., Damonte, M., Cohen, S. B., and Steedman, M.: The Role of Reentrancies in Abstract Meaning Representation Parsing. In Findings of the Association for Computational Linguistics: EMNLP 2020, pages 2198–2207. Association for Computational Linguistics, 2020.

Thadani, K. and McKeown, K.: Sentence Compression with Joint Structural Inference. In Proceedings of the Seventeenth Conference on Computational Natural Language Learning, pages 65–74. Association for Computational Linguistics, 2013.

Tjong Kim Sang, E. F. and De Meulder, F.: Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition. In Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003, pages 142–147, 2003.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G.: LLaMA: Open and Efficient Foundation Language Models. arXiv: 2302.13971 (Only available as arXiv preprint), 2023.

Trung, H. T., Toan, N. T., Vinh, T. V., Dat, H. T., Thang, D. C., Hung, N. Q. V., and Sattar, A.: A comparative study on network alignment techniques. Expert Systems with Applications, 140:112883, 2020.

Uzuner, O., South, B. R., Shen, S., and DuVall, S. L.: 2010 i2b2/VA challenge on concepts, assertions, and relations in clinical text. Journal of the American Medical Informatics Association : JAMIA, 18(5):552–556, 2011.

Vilca, G. C. V. and Cabezudo, M. A. S.: A Study of Abstractive Summarization Using Semantic Representations and Discourse Level Information. In Text, Speech, and Dialogue, eds. K.

Ekšteín and V. Matoušek, *Lecture Notes in Computer Science*, pages 482–490. Springer International Publishing, 2017.

Wang, D., Liu, P., Zheng, Y., Qiu, X., and Huang, X.: Heterogeneous Graph Neural Networks for Extractive Document Summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6209–6219. Association for Computational Linguistics, 2020.

Wang, Y., Liu, S., Rastegar-Mojarad, M., Wang, L., Shen, F., Liu, F., and Liu, H.: Dependency and AMR Embeddings for Drug-Drug Interaction Extraction from Biomedical Literature. In *Proceedings of the 8th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, pages 36–43. ACM, 2017.

Wei, Y.: Document summarization method based on heterogeneous graph. In *2012 9th International Conference on Fuzzy Systems and Knowledge Discovery*, pages 1285–1289, 2012.

Wiseman, S., Shieber, S., and Rush, A.: Challenges in Data-to-Document Generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263. Association for Computational Linguistics, 2017.

Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Kaiser, L., Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M., and Dean, J.: Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. arXiv: 1609.08144 (Only available as arXiv preprint), 2016.

Xu, S., Pang, L., Shen, H., Cheng, X., and Chua, T.-S.: Search-in-the-Chain: Interactively Enhancing Large Language Models with Search for Knowledge-intensive Tasks. In *Proceedings of the ACM on Web Conference 2024, WWW ’24*, pages 1362–1373. Association for Computing Machinery, 2024.

Yang, X., Chen, A., PourNejatian, N., Shin, H. C., Smith, K. E., Parisien, C., Compas, C., Martin, C., Flores, M. G., Zhang, Y., Magoc, T., Harle, C. A., Lipori, G., Mitchell, D. A., Hogan, W. R., Shenkman, E. A., Bian, J., and Wu, Y.: GatorTron: A Large Clinical Language Model to Unlock Patient Information from Unstructured Electronic Health Records. arXiv: 2203.03540 (Only available as arXiv preprint), 2022.

- You, J., Ying, R., Ren, X., Hamilton, W., and Leskovec, J.: GraphRNN: Generating Realistic Graphs with Deep Auto-regressive Models. In Proceedings of the 35th International Conference on Machine Learning, pages 5708–5717. PMLR, 2018.
- Zayyan, M.: Objective Structured Clinical Examination: The Assessment of Choice. Oman Medical Journal, 26(4):219–222, 2011.
- Zhang, S. and Tong, H.: FINAL: Fast Attributed Network Alignment. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16, pages 1345–1354. Association for Computing Machinery, 2016.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., and Artzi, Y.: BERTScore: Evaluating Text Generation with BERT. In Proceedings of the 8th International Conference on Learning Representations, 2020.
- Zhang, Y., Merck, D., Tsai, E., Manning, C. D., and Langlotz, C.: Optimizing the Factual Correctness of a Summary: A Study of Summarizing Radiology Reports. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 5108–5120. Association for Computational Linguistics, 2020.
- Zhang, Z., Parulian, N., Ji, H., Elsayed, A., Myers, S., and Palmer, M.: Fine-grained Information Extraction from Biomedical Literature based on Knowledge-enriched Abstract Meaning Representation. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6261–6270. Association for Computational Linguistics, 2021.
- Zhou, Q., Yang, N., Wei, F., Huang, S., Zhou, M., and Zhao, T.: Neural Document Summarization by Jointly Learning to Score and Select Sentences. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 654–663. Association for Computational Linguistics, 2018.

VITA

PAUL LANDES

Personal Summary

- Specialization in natural language processing (NLP), machine learning (ML) algorithms and models, meaning representation, and automatic summarization.
- More than nine years of research experience in academia and industry with a focus on NLP and deep learning problems in text classification, summarization and clinical domain.

Education

- 2013-2017 Master of Science, Computer Science
- 1995–1999 Bachelor of Science, Computer Science University of North Texas Denton

Experience

- 2023 Research Internship Argonne National Labs Lemont, IL
- 2020–2022 Research Assistant University of Illinois Chicago, IL
- 2019–2020 AI/Machine Learning Engineer Blue Cross and Blue Shield Chicago, IL
- 2018–2019 NLP/Machine Learning Consultant United Airlines Chicago
- 2016–2019 Senior R&D Scientist OpinionLab, A Verint Company Chicago
- 2015–2016 NLP Data Scientist ChoreHat Chicago

Publications

- Paul Landes and Barbara Di Eugenio. “CALAMR: Component ALignment for Abstract Meaning Representation”. In the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING, 2024).
- Krenare Nuci, Paul Landes, and Barbara Di Eugenio. “RoBERTa Low Resource Fine Tuning for Sentiment Analysis in Albanian”. In the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024).

- Paul Landes, Barbara Di Eugenio, and Cornelia Caragea. “DeepZensols: A Deep Learning Natural Language Processing Framework for Experimentation and Reproducibility”. In the Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023).
- Paul Landes, Aaron Chaise, Kunal Patel, Sean Huang, and Barbara Di Eugenio. “Hospital Discharge Summarization Data Provenance”. In the 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks (BioNLP 2023).
- Paul Landes, Kunal Patel, Sean S. Huang, Adam Webb, Barbara Di Eugenio, and Cornelia Caragea. “A New Public Corpus for Clinical Section Identification: MedSecId”. In Proceedings of the 29th International Conference on Computational Linguistics (COLING, 2022).
- Paul Landes and Barbara Di Eugenio. “Supervised Approach To The Interpretation Of Imperative To-Do List”. arXiv: 1806.07999.

Patents

- Paul Landes. “Automatic Corpora Annotation”. U.S. pat. req. 2021/0263971A1. Health Care Service Corporation.
- Paul Landes. “Structural Representation and Facilitation of Manipulation Thereof via Implicit Vertex Relationships”. U.S. pat. 10,936,762B2. Navteq. Jan. 26, 2016.
- Nimesh Patel, Greg Blew, and Paul Landes. “Dynamic Management, Assembly, and Presentation of Web-Based Content”. U.S. pat. 9,870,451B1. EmmiSolutions. Jan. 16, 2018.

Books

- Knute Axelson, Mary Bellino, Dave Harper, Dave Iffland, Dan Ignat, Greg Kick, Paul Landes, Paul Millar, Sonny Mounicou, and Munger Greg. “Coding: The Handbook for Information Technology”. BlueLine Press, 2005-10-03, 2005. 536 pp. isbn: 0-9766279-0-6.

Program Committee Member or Reviewer

- ACL Association for Computational Linguistics Rolling Review
- Workshop on Designing Meaning Representations (DRM)
- COLING International Committee on Computational Linguistics
- AIED Conference on Artificial Intelligence in Education