

# External Knowledge-Driven Argument Mining: Leveraging Attention-Enhanced Multi-Network Models

Anonymous ACL submission

## Abstract

Argument mining (AM) involves the identification of argument relations (AR) between Argumentative Discourse Units (ADUs). The essence of ARs among ADUs is context-dependent and lies in maintaining a coherent flow of ideas, often centered around the relations between discussed entities, topics, themes or concepts. However, these relations are not always explicitly stated; rather, inferred from implicit chains of reasoning connecting the concepts addressed in the ADUs. While humans can infer such background knowledge, machines face challenges when the contextual cues are not explicitly provided. This paper leverages external resources, including WordNet, ConceptNet, and Wikipedia to identify semantic paths (knowledge paths) connecting the concepts discussed in the ADUs to obtain the implicit chains of reasoning. To effectively leverage these paths for AR prediction, we propose attention-based Multi-Network architectures. Various configurations of the architecture are evaluated on the external resources, and the configuration using Wikipedia achieves a new state-of-the-art performance with F-scores of 0.85, 0.84, 0.70, and 87, respectively, on four diverse datasets.

## 1 Introduction

Argument mining involves identifying the argumentative structure within a text. It includes segmenting arguments into Argumentative Discourse Units (ADUs) (Peldszus and Stede, 2015a), distinguishing argumentative units from non-argumentative ones, classifying ADUs, labeling argument relation (AR) between ADUs, and identifying argument schemes (Persing and Ng, 2016; Stab and Gurevych, 2017; Lawrence and Reed, 2020). This study focuses on classifying the AR between ADUs into three categories: Inference (RA) (when one ADU supports the other), Conflict (CA) (when one ADU attacks the other), and None (when there is no AR).

The nature of AR is inherently context-dependent (Potash et al., 2017; Habernal et al., 2017; Choi and Lee, 2018; Rinott et al., 2015), relying on maintaining a coherent flow of interconnected ideas. This cohesion is often centered around the connections between the discussed entities, topics, themes or concepts, commonly referred to as **Local coherence** (Foltz et al., 1998; Marcu, 2000). Local coherence facilitates smooth idea transitions between ADUs by recognising inherent regularities in entity distribution. Similarly, other entity-based theories of discourse (Givón, 1987; Prince, 1981) and **Centering Theory** (Grosz et al., 1995) propose that these regularities contribute to the coherence of discourse by guiding the organisation of ideas around salient entities. Following a similar framework, aspect-based argument mining techniques use the relationships between the concepts discussed in ADUs, to identify argument structures (Misra et al., 2017; Dragoni et al., 2018; Gemechu and Reed, 2019; Trautmann, 2020). Yet, the contexts required to link these concepts are not always explicit and are often inferred from background knowledge.

Pre-trained large language models (LLMs) have transformed NLP, moving from traditional feature engineering to data-driven approaches. Studies indicate that these models implicitly capture various types of knowledge, including relational, commonsense, and structural linguistic knowledge, within their parameters (Petroni et al., 2019; Goldberg, 2019; Safavi and Koutra, 2021; AIKhamissi et al., 2022). While excelling in various NLP tasks, their ability to encode the necessary background knowledge for identifying ARs remains uncertain (Kassner and Schütze, 2019). For example, Polu et al. (2022) revealed their limitations in chaining multiple steps of complex logical reasoning, while Merrill et al. (2021) demonstrated they fail to comprehend the semantics behind commonsense reasoning tasks. This limitation is critical in AR iden-

043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083

tification, as linking ADUs relies on the implicit chain of reasoning, often inferred from the chain of relations between the concepts discussed in the ADUs. This highlights the need for supplementary contextual information from external sources to establish these connections.

Consider the ADUs from the 2016 presidential election debate corpus (Visser et al., 2019) in Table 1. Identifying the AR between (1) and (2) relies on recognising the relationship between “NAFTA agreement” and “USA”, whereas for (4) and (5), it requires understanding “building electric grid” is an “economic activity”. While these connections are straightforward for human experts, computers face challenges as such interconnections are often implicitly inferred. For example, the AR between (3) and (4) is direct as the relation between the concepts mentioned in the respective ADUs can be obtained from an ontology (Miller, 1995; Speer et al., 2017) (“Electric grid; grid” is directly related to “power; electrical power” in WordNet (Miller, 1995)) or by comparing their embeddings (Pilehvar et al., 2013; Le and Mikolov, 2014; Reimers and Gurevych, 2019). However, identifying the AR between (4) and (5) is challenging since the path linking “electric grid” to “economic activity” is missing in existing knowledge resources including WordNet (Miller, 1995) or ConceptNet (Speer et al., 2017) or DBpedia<sup>1</sup>. However, the concepts are indirectly linked in Wikipedia through a chain of concepts interlinked using a set of semantic relation types: “economic activity” *involves* “innovation” *which constitutes* developing “clean energy” *transmitted by* “electric grid”. This study aims to identify and leverage the chain of such semantic relations between the concepts, to capture implicit referential information between ADUs (Asher and Lascarides, 2003) and use it for AR prediction.

No	ADUs
1	[USA] <sub>C</sub> [is in deep trouble] <sub>OC</sub>
2	[NAFTA agreement] <sub>C</sub> [is defective] <sub>OC</sub>
3	[We] <sub>C</sub> [can have] <sub>OC</sub> [clean energy] <sub>A</sub>
4	[We] <sub>C</sub> [can build] <sub>OC</sub> a new modern [electric grid] <sub>A</sub>
5	[This] <sub>C</sub> [is a lot of] <sub>OC</sub> new [economic activity] <sub>A</sub>

Table 1: Examples from 2016 presidential election debate corpus (Visser et al., 2019) to illustrate the relation between the functional components of ADUs. *C* represents the theme of the sentence, *A* represents the aspects specialising the theme, while the opinion on *C* is represented by *OC*.

<sup>1</sup><https://wiki.dbpedia.org/>

Leveraging knowledge from external resources has been shown to improve performance in AM (Kobbe et al., 2019; Botschen et al., 2018; Fromm et al., 2019; Plenz et al., 2023) and related tasks, such as semantic plausibility (Wang et al., 2018), identifying inferences (Chen et al., 2017), and determining entailment (Glockner et al., 2018). However, existing studies on AR prediction exclusively utilise structured knowledge bases and overlook semi-structured resources like Wikipedia, which contains over 6,805,837 articles (as of April 1, 2024), offering richer connections through hyperlinks embedded within articles. Moreover, these methods rely on entities, events, and factual information sourced from structured databases, limiting their applicability to specific domains. In contrast, using generic semantic relation types that encode AR ensures adaptability across domains (refer to Table 7 for examples of such relation types). Furthermore, they lack effective method for integrating the external information into model architectures, relying instead on conventional feature engineering techniques. For instance, Kobbe et al. (2019) leverage features derived from graph representations of the resources, including the inter-concept distances within the graph. Similarly, Plenz et al. (2023) employ semantic similarity to determine the relevance of external knowledge, in conjunction with traditional features derived from the graph representation of the resources.

In this paper, we propose traversing Wikipedia, WordNet, and ConceptNet to find semantic paths linking concepts mentioned in ADU pairs. ARs between ADUs are identified by leveraging these paths using attention-based Multi-Network architectures. To establish a benchmark, we evaluate LLMs across various configurations, comparing the knowledge obtained from external resources with that inherent in LLMs. The evaluation demonstrated that integrating external resources consistently enhances performance, achieving new state-of-the-art results. Additionally, we assess the effectiveness of the attention-based Multi-Network architecture in leveraging external knowledge, consistently demonstrating its superiority over the standard linear classification baseline. The contribution of this paper is four-fold: (a) the utilisation of both structured and semi-structured external resources for AR prediction, (b) architecture for effectively leveraging external knowledge, (c) features adaptable across domains, and (d) the state-of-the-art performance.

## 2 Related Works

In the literature, AM has been approached using various configurations, including dependency parsing (Peldszus and Stede, 2015b), discourse parsing (Muller et al., 2012), sequence tagging (Eger et al., 2017; Mayer et al., 2020), and sequence classification configurations (Reimers et al., 2019; Ruiz-Dolz et al., 2021; Mayer et al., 2020). Various works tackle specific AM tasks. Some focus exclusively on argument segmentation (Chernodub et al., 2019; Ajjour et al., 2017), while others start with segmented data and focus solely on AR identification (Potash et al., 2016; Gemechu and Reed, 2019; Ruiz-Dolz et al., 2021). Potash et al. (2016) train an encoder-decoder (Sutskever et al., 2014) with attention mechanism (Bahdanau et al., 2014) to identify AR. Gemechu and Reed (2019) decompose ADUs into fine-grained components and use classifiers to predict AR based on the relations between these components. Chakrabarty et al. (2020) identify argument components and ARs within both inter-turn and intra-turn interactions in dialogues. They classify ARs as a binary prediction, determining only the presence of a relation without specifying its type. Their findings indicate that using distant-labeled data and integrating discourse relations from Rhetorical Structure Theory (Mann and Thompson, 1988) improve performance.

End-to-end AM approaches address multiple AM tasks, simultaneously. Persing and Ng (2016) and Stab and Gurevych (2017) adopt a pipeline architecture and train separate models for each subtask to then utilise an Integer Linear Programming (ILP) model to encode global constraints. Eger et al. (2017) propose a neural end-to-end approach, framing the task in various configurations including dependency parsing and token-based sequence tagging. They also employ a multi-task setup to leverage the dependencies between AM tasks, including component identification and AR prediction. Their best-performing configuration achieves an F1-score of 0.51 for AR identification on the AAEC dataset. Peldszus and Stede (2016) aim to map RST trees to argumentation structures (Taboada and Mann, 2006) using sub-graph matching and an evidence graph model. They evaluate various features of their system on the AMT dataset and achieve an overall F-measure of 0.76 in identifying ARs. Similarly, Morio et al. (2022) introduce an end-to-end cross-corpus training strategy that facilitate information transfer between datasets.

Mayer et al. (2020) address argument component and relation identification on a dataset comprising various disease treatments. The approach involves combining static and dynamic embeddings using various configurations of RNN and CRFs. They demonstrate the efficacy of specialised LLMs like SciBERT (Beltagy et al., 2019), highlighting their relevance in medical domain adaptations. However, most of these works rely on the information explicitly provided in the argument alone.

Recent AM works fine-tuned LLMs in sequence classification fashion (Reimers et al., 2019; Ruiz-Dolz et al., 2021). Studies show that such LLMs implicitly capture relational, commonsense, and structural linguistic knowledge (Petroni et al., 2019; Goldberg, 2019; Safavi and Koutra, 2021; Alkhamissi et al., 2022). Despite their significant performance, the ability of LLMs to encode the requisite background knowledge for identifying ARs remains uncertain, raising concerns about relying solely on LLMs for this task (Kassner and Schütze, 2019). For instance, Polu et al. (2022) exposed their limitations in complex logical reasoning, while Merrill et al. (2021) showed they struggle in comprehending the semantics of commonsense reasoning tasks.

The works most related to ours are those of Kobbe et al. (2019) and Plenz et al. (2023), as they also leverage external knowledge bases to identify AR. However, their methodologies differ significantly from ours. Firstly, they rely on structured knowledge bases with predefined relation types, while we also use semi-structured resources like Wikipedia that cover diverse relations. Furthermore, they struggle to effectively integrate external information into model architectures, relying instead on conventional feature engineering techniques that exploit structural features obtained from sub-graph extracted from external knowledge bases. For instance, Kobbe et al. (2019) use features like the frequency of relations existing between ADUs. Similarly, Plenz et al. (2023) leverage the similarity between external knowledge and ADUs to identify relevant sub-knowledge graphs and exploit the sub-graph to extract categorical features, such as the number of shared concepts between ADU pairs and the path lengths between the concepts. Additionally, the formalisation of the “concepts” used for alignment with external resources is vague, relying on arbitrary entity mentioned in the ADUs. Moreover, their approach for AR identification has not been evaluated.



277	<b>3 Methodology</b>	
278	<b>3.1 Data</b>	
279	We use four corpora. The first is AAEC (Stab and	
280	Gurevych, 2017) which has a total of 402 argu-	
281	ments. ADUs under each argument are labelled	
282	as premise, claim or major claim. It has 147,271	
283	tokens, 6,089 ADUs and 5335 ARs (4841 support	
284	and 497 attack).	
285	The second corpus is the Argumentative Micro	
286	Text (AMT) (Peldszus and Stede, 2013) which is a	
287	collection of 112 short texts collected from human	
288	subjects in German translated into English. It is	
289	annotated following the argumentation structure	
290	outlined by MicroTextAnnotation. The structure	
291	consists of a central claim, and supporting ADUs.	
292	It has a total of 8007 tokens, 576 ADUs and 443	
293	ARs (272 support and 171 attack).	
294	The third corpus is part of the US 2016 presi-	
295	dential election debate corpus (US2016) (Visser	
296	et al., 2019) which is annotated based on Inference	
297	Anchoring Theory (IAT) (Budzynska and Reed,	
298	2011). Argument components are referred to as	
299	propositions, with the relations between them an-	
300	notated as default inference for support and default	
301	conflict for attack. The corpus has a total of 15805	
302	tokens, 1473 ADUs and 584 ARs (505 support and	
303	79 attack).	
304	The fourth corpus is the AbstrCT corpus	
305	(Mayer et al., 2020) which consists of abstracts	
306	extracted from the MEDLINE database. Argument	
307	components are categorised into major claim,	
308	claim, and evidence components, and the relations	
309	between them are categorised into support, attack,	
310	and partial-attack. The corpus consists of 100,253	
311	tokens, 4,679 ADUs, and 2,634 ARs, including	
312	344 attack relations (combining attack and partial-	
313	attack relations) and 2,290 support relations.	
314	As described above, argument components are	
315	annotated non-uniformly across datasets, based on	
316	the underlying theoretical framework. For exam-	
317	ple, in AAEC, argument components are annotated	
318	as premises, claims, and major claims. However,	
319	in US2016, the components are not explicitly cate-	
320	gorised, but the premise-conclusion notion can be	
321	inferred from the direction of the AR. As our cur-	
322	rent objective does not involve classifying the com-	
323	ponents or the direction of the relation, we focus on	
324	the AR existing between the components without	
325	classifying the categories of the components (into	
326	claim/conclusion/major-claim, premise/evidence).	
	<b>3.2 External Knowledge Alignment and</b>	<b>327</b>
	<b>Paths Extraction</b>	<b>328</b>
	Each ADU is annotated into its four functional	329
	components, following the framework proposed	330
	by Gemechu and Reed (Gemechu and Reed, 2019)	331
	(see Appendix A.3 for more details). These com-	332
	ponents are used to align the respective ADUs with	333
	the external resources. The functional components	334
	consist of target concepts ( <i>C</i> ), aspects ( <i>A</i> ), opin-	335
	ions on <i>C</i> ( <i>OC</i> ), and opinions on <i>A</i> ( <i>OA</i> ). ( <i>C</i> ) refers	336
	to the set of concepts related to the ADUs’s topic,	337
	while ( <i>A</i> ) refers to the set of concepts further speci-	338
	fying that topic (examples provided in Table 1). In	339
	this study, we focus on ( <i>C</i> ) and ( <i>A</i> ), which repre-	340
	sent the topics and aspects addressed by the ADUs.	341
	The statistics of these components can be found in	342
	Table 4 in the Appendix.	343
	To extract relevant external knowledge, we	344
	align these components with two ontological res-	345
	ources—WordNet (Miller, 1995) and ConceptNet	346
	(Speer et al., 2017)—as well as a semi-structured	347
	resource, Wikipedia. The detailed alignment pro-	348
	cess is described in Sections 3.2.1 to 3.2.2.	349
	<b>3.2.1 Ontology as External Source</b>	<b>350</b>
	We traverse WordNet (Miller, 1995) and Concept-	351
	Net (Speer et al., 2017) Synset hierarchies and	352
	align the components of ADUs with the Synsets,	353
	to identify the chain (path) of Synsets that connects	354
	the components. The alignment relies on cosine	355
	similarity between the embeddings of the compo-	356
	nents and Synsets, determined by the cosine simi-	357
	larity threshold $\beta$ . Sentence-transformer (Reimers	358
	and Gurevych, 2019) is utilised to identify the em-	359
	beddings. For more details on the embeddings and	360
	similarity threshold, check Appendix A.3.2.	361
	By treating the ontology as a graph, with Synsets	362
	as nodes and relation types as edges, we begin the	363
	search with one of the components and traverse the	364
	knowledge graphs until either the other component	365
	is found or the search depth reaches the threshold	366
	$\alpha = 5$ . For more details on setting the value of	367
	$\alpha$ , refer to Appendix A.3.3. If the search is suc-	368
	cessful, we concatenate the Synsets and the type	369
	of semantic relation between, otherwise return the	370
	concatenation of both components with the con-	371
	stant string “None” in between. We use relation	372
	types with frequency higher than $m=3$ to form the	373
	paths (see Appendix A.3.5 for more information	374
	about the relation filtering process).	375

### 3.2.2 Wikipedia as External Source

We also traverse Wikipedia to identify the chain (path) of Wikipedia pages linking the functional components of ADUs. For any pair of components (e.g.,  $C_1$ ,  $A_2$  or  $C_1$ ,  $C_2$  or  $A_1$ ,  $A_2$ ) associated with a pair of ADUs ( $p_1$ ,  $p_2$ ), the initial step involves aligning these components with corresponding Wikipedia pages. This alignment is achieved by computing the similarity between the embeddings (Reimers and Gurevych, 2019) of the Wikipedia page titles and the components.

Viewing Wikipedia as a graph (with pages as nodes and hyperlinks as edges), we begin a breadth-first search from the Wikipedia page of one concept ( $c_1$ ), continuing until we locate the second concept ( $c_2$ ) or reach a depth threshold,  $\alpha = 5$ . During this search, we record sentences ( $S$ ) containing Wikipedia page titles of the current page ( $hl_1$ ) and the hyperlinks leading to the next Wikipedia page ( $hl_2$ ) along the path. These sentences contribute to the formation of a tuple:  $\langle hl_1, hl_2, keywords \rangle$ , where the keywords represent the semantic relation type linking  $hl_1$  and  $hl_2$  within the sentences.

We utilise semantic role labeling (SRL) to identify the keywords that connect  $hl_1$  and  $hl_2$  within the sentences ( $S$ ) containing the hyperlinks. The SRL tool from AllenNLP<sup>2</sup> is used for this purpose. The process involves extracting subject-predicate structures that link  $hl_1$  and  $hl_2$  in the sentences involving the hyperlinks, followed by identifying phrases that connect them across the semantic roles assigned (see Appendix A.3.4). Top  $m$  most frequent relations are selected to construct the paths.

### 3.3 Model

We propose attention-based Multi-Network to leverage the information obtained from external resources for AR prediction (Section 3.3.1). Section 3.3.2 presents baseline models that utilise LLMs alone as sources of background knowledge.

#### 3.3.1 Attention-Based Multi-Network

We investigate two attention-based Multi-Network configurations, namely Siamese and Triplet (Schroff et al., 2015) networks. Initially, we utilise the Siamese network involving two sub-networks, where one sub-network encodes the concatenation of both ADUs together while the other encodes the external information. Furthermore, we examine Triplet network, which uses three sub-network

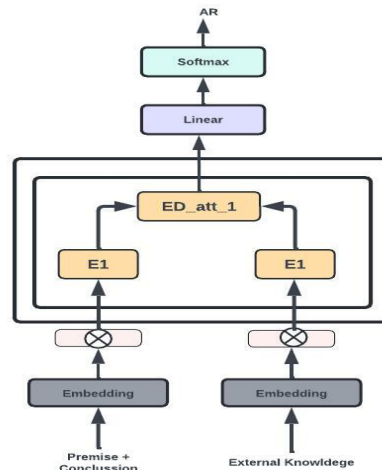


Figure 1: Siamese-networked with attention layers.

to encodes each of the ADUs and the external resources separately.

**Siamese Network Architecture with Attention.** In this setup, given the two sub-networks (**E1** and **E2**) in Siamese network, **E1** processes the concatenation of the pair of ADUs (premise and conclusion), while **E2** handles the concatenation of the information from external resources. Cross attention layer (**ED-att-1**) is applied to the outputs of these sub-networks for attending to the external resources relevant to the premise and conclusion (see Figure 1). Accordingly, the output of (**E1**) serves as the query, while the output of (**E2**) is used as keys and values, enabling to query the external information relevant to the premise and conclusion. It employs multi-head attention  $h$ , where each head  $j$  computes scaled dot-product attention using query  $Q^j$ , key  $K^j$ , and value  $V^j$  matrices, which are linear transformations of the input hidden state  $h_i$ . The final attention weight  $e_i$  is obtained by concatenating over all attention heads. The resulting attention weights are then multiplied with the output of **E1**, and passed through a fully connected classification layer, for AR classification. This fusion allows the model to integrate the original representations of the premise and conclusion with the extracted external information (see detailed model parameters in the Appendix A.5).

**Triplet Network Architecture with Attention.** In contrast to the Siamese Architecture, the Triplet Network Architecture consists of three sub-networks: **E1**, **E2**, and **E3** (see Figure 3 in appendix 3.3.1). Sub-networks **E1** and **E2** encode the premise and conclusion, respectively, while **E3**

<sup>2</sup> [https://docs.allennlp.org/v0.9.0/api/allennlp.models.semantic\\_role\\_labeler.html](https://docs.allennlp.org/v0.9.0/api/allennlp.models.semantic_role_labeler.html)

encodes the external knowledge connecting them. Two cross-attention layers are introduced (**ED-att-1** and **ED-att-2**). **ED-att-1** focuses on the relation between the premise and conclusion, where the output of **E1** serves as queries and the output of **E2** is used as keys and values. On the other hand, **ED-att-2** attends to the external knowledge relevant to the premise and conclusion. Specifically, the output of **ED-att-1** acts as the query, while the output of **E3** is used as keys and values. Similar to the Siamese architecture, we combine the output of the two attention layers for classification. The rationale behind this approach is that **ED-att-1** encodes the relation between the premise and conclusion, while **ED-att-2** encodes the relevant external resource, enabling the model to effectively leverage both the relationship between the premise and conclusion and the relevant external knowledge for AR classification.

### 3.3.2 LLMs as Baseline Models

We establish LLMs without external resources as baseline models under two configurations: few-shot and fully fine-tuning configurations. We evaluate these baselines against configurations that leverage external knowledge sources to enhance the performance of LLMs.

**Zero-shot setup:** We prompt GPT-4<sup>3</sup>, a generative LLM, to perform two tasks: (a) predicting ARs for comparative analysis against models utilising external resources, and (b) generating paths between ADU components for comparison with models using paths derived from ontology and Wikipedia. Accordingly, GPT-4-generated paths are used as external knowledge to train the Multi-Network configuration for AR classification. This enables a direct comparison between GPT-4-generated paths and those obtained from other external knowledge sources. The experimental setup for prompting GPT-4 is provided in A.4.

**Fine-tuning setup:** We also fine-tuned BERT (Devlin et al., 2018) using various configurations for comparison. Initially, we use the vanilla sequence classification setup ( $SM \odot V \oplus bert$ ), where the concatenation of ADUs is presented as an input. Furthermore, we fine-tune BERT within Siamese architectures, both with ( $SM \odot A \oplus bert$ ) and without attention layers ( $SM \odot V \oplus bert$ ). See A.1 for the details of model configuration and experimental setups.

<sup>3</sup><https://openai.com/chatgpt>

## 4 Experiments

### 4.1 Experimental Setup

The dataset is randomly partitioned, with 70%, 10%, and 20% allocation for training, validation, and testing respectively, ensuring uniformity throughout the dataset. Refer to Table 3 in the Appendix for the breakdown of ARs across the datasets. Results represent the average of three runs using different random seeds. Precision (P), recall (R), and F-measure (F) are computed, and macro-averaged P, R, and F are reported for the test dataset (more experimental setup provided in Appendix A.1). The datasets and code utilised in our experiments are available for public access at ANONYMISED\_URL.

### 4.2 Model Configurations

We evaluate various configurations of approaches leveraging the two ontological resources (WordNet, and ConceptNet) and Wikipedia across the four datasets. These configurations encompass three Triplet network architectures:  $TL \odot A \oplus wp$  for Wikipedia,  $TL \odot A \oplus wn$  for WordNet, and  $TL \odot A \oplus cn$  for ConceptNet. Similarly, three Siamese network architectures are evaluated across these ontological resources:  $SM \odot A \oplus wp$ ,  $SM \odot A \oplus wn$ , and  $SM \odot A \oplus cn$ .

Furthermore, to evaluate the attention layers' impact on external resources, we compare Triplet and Siamese architectures without attention layers across the three external resources, totaling six configurations:  $TL \odot V \oplus wp$ ,  $TL \odot V \oplus wn$ ,  $TL \odot V \oplus cn$ ,  $SM \odot V \oplus wp$ ,  $SM \odot V \oplus wn$ , and  $SM \odot V \oplus cn$ . Finally, we evaluate the Triplet architecture on GPT-4 generated paths ( $TL \odot V \oplus gpt$ ,  $TL \odot A \oplus gpt$ ).

### 4.3 Results and Discussions

The evaluation results depicted in Table 2 revealed clear trends in performance. Particularly, the influence of model architecture and the incorporation of external knowledge on AR prediction. This is evidenced by the performance improvement observed in configurations with such integration compared to those without.

Models incorporating external resources outperformed those lacking such integration, indicating the importance of leveraging additional knowledge sources for AR identification. This led to a notable enhancement, surpassing the baseline by over 5.4% in F-measure. For example, the Siamese architecture leveraging Wikipedia achieved an average F-

Configs	AAEC			AMT			US2016			AbstRCT		
	P	R	F	P	R	F	P	R	F	P	R	F
<b>Comparison</b>												
P2016	n/a	n/a	77	n/a	n/a	74	n/a	n/a	n/a	n/a	n/a	n/a
K2019	n/a	n/a	59	n/a	n/a	67	n/a	n/a	n/a	n/a	n/a	n/a
PS2016	n/a	n/a	n/a	n/a	n/a	76	n/a	n/a	n/a	n/a	n/a	n/a
E2017	n/a	n/a	51	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
GPT-4	63±2	48±2	55±2	60±2	47±2	52±2	58±1	43±2	50±1	69±3	58±2	63±2
GR2019	81	74	77	<b>88</b>	66	75	63	61	62	n/a	n/a	n/a
M2020	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	62	n/a	n/a	69
<b>LLMs as KB</b>												
SC⊙V⊕bert	78±0.3	73±0.2	75±0.1	79±0.4	67±0.1	72±0.1	56±0.4	64±0.2	60±0.2	84±0	82±0	83±0
SM⊙V⊕bert	77±0.1	72±0.1	74±0.1	80±0.9	65±0.2	72±0.5	55±0.2	63±0.1	59±0.1	82±0	82±0	82±0
SM⊙A⊕bert	80±0.2	73±0.3	76±0.2	80±0.1	68±0.3	74±0.2	57±0.1	64±0.2	60±0.1	85±0	83±0	84±0
<b>No Att + Ext</b>												
TL⊙V⊕gpt	77±2	84±2	80±2	74±3	81±2	77±3	54±4	76±3	64±3	72±4	<b>87±3</b>	80±4
SM⊙V⊕wn	84±0	79±0.2	81±0.1	82±0.4	73±0.3	77±0.3	62±0.1	69±0.1	65±0.1	82±0	82±0	82±0
SM⊙V⊕cn	83±0.3	76±0.1	80±0.2	82±0.1	72±0.2	77±0.1	61±0.2	71±0.2	66±0.2	84±0	85±0.1	85±0.1
SM⊙V⊕wp	82±0.2	82±0.1	82±0.1	84±0.2	76±0.3	80±0.2	63±0.3	71±0.6	67±0.3	85±0	85±0	85±0
TL⊙V⊕wn	83±0.1	79±0.2	81±0.1	<b>84±0.1</b>	75±0.1	80±0.1	61±0	70±0	65±0	83±0.1	82±0.1	82±0.1
TL⊙V⊕cn	83±0.1	80±0.2	82±0.1	84±0.1	76±0.1	80±0.1	61±0.1	71±0.1	66±0	85±0	85±0	85±0
TL⊙V⊕wp	84±0	80±0	82±0	82±0.1	76±0.1	79±0.1	64±0.1	70±0.1	67±0.1	86±0.1	85±0	86±0.1
<b>Att + Ext</b>												
TL⊙A⊕gpt	77±3	<b>84±2</b>	80±3	71±4	<b>85±4</b>	77±3	56±3	72±3	63±3	73±2	85±4	79±3
SM⊙A⊕wn	84±0.1	81±0.1	82±0.1	83±0.1	79±0.1	81±0.1	62±0.1	72±0.1	67±0.1	83±0	82±0	83±0
SM⊙A⊕cn	83±0.3	81±0.2	82±0.2	81±0.2	82±0.2	81±0.2	65±0.1	72±0.1	68±0.1	85±0.1	85±0	85±0.1
SM⊙A⊕wp	85±0.2	81±0.1	83±0.1	82±0.2	83±0.2	82±0.1	65±0.2	73±0.2	69±0.2	85±0.1	86±0	86±0
TL⊙A⊕wn	85±0.1	82±0.1	84±0.1	83±0.2	84±0.2	<b>84±0.1</b>	64±0.2	72±0.3	68±0.2	83±0.1	83±0	83±0
TL⊙A⊕cn	84±0.1	82±0.2	83±0.1	82±0.1	84±0.1	83±0.1	65±0.2	73±0.3	69±0.2	86±0	85±0	86±0
TL⊙A⊕wp	<b>86±0.1</b>	83±0.2	<b>85±0.2</b>	83±0.1	<b>85±0</b>	<b>84±0</b>	<b>66±0.1</b>	75±0.1	<b>70±0.1</b>	<b>87±0.1</b>	86±0	<b>87±0</b>

Table 2: Performance of our models and the comparison systems including, (Potash et al., 2016) (P2016), (Eger et al., 2017) (E2017), (Peldszus and Stede, 2016) (PS16), (Kobbe et al., 2019) (K2019), (OpenAI, 2023) (GPT-4), (Gemechu and Reed, 2019) (GR2019), (Mayer et al., 2020) (M2020) across the four datasets. The reported results have been averaged from 3 randomly initialised sequential runs. The table is divided into subsections: Comparison approaches; LLM-alone; non-attention with external sources; attention-based with external resources.

measure of 80% across datasets, whereas its counterpart, lacking the external resource, achieved 74%. This finding aligns with previous research demonstrating that while LLMs tend to encode world knowledge, LLMs alone may not fully present the depth and specificity of knowledge required for certain tasks, such as AR identification involving structured and chained reasoning (Kassner and Schütze, 2019; Polu et al., 2022; Merrill et al., 2021). Likewise, models equipped with attention mechanisms consistently surpassed those without, demonstrating an average increase in F-measure of over 2% across diverse configurations. Notably, Triplet Network architecture with attention mechanism leveraging Wikipedia as an external knowledge source, attained an average F-measure of 81% across the datasets. This represents a new state-of-the-art performance in AR identification, showcasing the effectiveness of the architecture in integrating external knowledge.

We also compare our approach to other related works including Potash et al.’s (2016), Eger et al.’s (2017), Peldszus and Stede’s (2016), Kobbe et al.’s (Kobbe et al., 2019), OpenAI’s GPT-4 (OpenAI, 2023), Gemechu and Reed’s (2019; 2023) and Mayer et al.’s (2020) work. Please note that direct comparisons with some of these works need additional contextual nuance in interpretation due to variations in task setup and complexities. For instance, the works of Eger et al. (2017) and Mayer et al. (2020) involve argument segmentation in addition to AR identification as an end-to-end task. In our case, the goal is to identify AR based on correct segments in the gold datasets. Similarly, Plenz et al. (Plenz et al., 2023) evaluate their approach on several AM tasks, including ValNov Shared Task (Heinisch et al., 2022), which involves assessing the validity and novelty of a conclusion given a premise—a task closely related to AR prediction. They report an F1 score of 70.69% for this task.



As can be seen from Table 2, our approach outperforms the comparison systems, including OpenAI’s GPT-4 (OpenAI, 2023) across the datasets.

**Model Architecture Influence.** As shown in Table 2, incorporating attention layers into Multi-Network architectures, brought clear benefits. Multi-network configurations with attention mechanisms outperformed the vanilla sequence classification setup, both with and without external knowledge, achieving an average F1 gains of 6.4% and 1%, respectively. Attention-based configurations leveraging external resources consistently outperform their counterparts without attention, yielding an average F1-score improvement of 2%. The attention-based Triplet architecture outperformed their counterpart Siamese architecture, with an average performance increase of 1.2% in leveraging external knowledge. It is noteworthy that in the absence of attention and external resources, multi-network configurations ( $SM \odot V \oplus bert$ ) underperform as compared to the vanilla sequence classification approach ( $SC \odot V \oplus bert$ ).

This highlights the efficacy of attention-based Multi-Network architectures in leveraging external resources for AR prediction, contrasting with standard sequence classification setups. Additionally, the performance advantage of Triplet architecture over Siamese architecture can be attributed to its design, enabling each sub-network to focus on learning two levels of alignment: between the premise and conclusion, and between the external resource and the premise-conclusion pair. To explore whether the performance gap solely stems from the additional parameters in the attention layer, we introduced extra linear layers to the Multi-Network architecture (without attention layers) and observed no change in performance despite the additional layers. However, attention analysis is required to substantiate this claim.

**External knowledge influence.** Wikipedia-based models outperformed the baselines and ontology-based models across all four datasets. The attention-based Triplet-network on Wikipedia ( $TL \odot A \oplus wp$ ) achieved an F-measure of 0.85, 0.84, 0.70 and 0.87 in identifying AR on AAEC, AMT, US2016, and AbstRCT respectively. Upon analysis of the paths connecting the components of ADUs, we found that 37% of concepts not present in ontological resources are connected in Wikipedia, while only 7% of concepts absent in Wikipedia are covered by ontological resources. For further details, please refer to Appendix A.3.6.

This disparity can be attributed to Wikipedia’s rich network of hyperlinks connecting pages using diverse relations, unlike ontological resources that only connect Synsets based on predefined sets of semantic relations.

Models trained on GPT-4 generated paths outperformed those without external knowledge, aligning with other works leveraging LLM-generated commonsense knowledge (Bansal et al., 2022). However, despite exhibiting higher accuracy, they still demonstrated lower precision compared to the approaches used the external knowledge sources. The observed high recall and low precision can be attributed to the models’ inclination to identify unintended paths between concepts. Twenty errors were randomly selected for analysis, with two human annotators collaboratively examining the paths. Of these, 14 errors were considered contextually irrelevant, despite the logical coherence evident in the generated paths. These paths introduce chains of thought that diverge from the original argument, as previously noted by other studies that rely on LLMs to generate commonsense knowledge (Levy et al., 2022). An error analysis can be found in Appendix A.3.6.

## 5 Conclusion

Our exploration of various model configurations underscored the importance of external resources and multi-network architecture with attention mechanisms in AR prediction. Models augmented with external resources consistently outperform those relying solely on LLMs. This emphasises the necessity of leveraging supplementary knowledge sources to enrich LLMs for AR prediction. Furthermore, multi-network architectures with attention mechanisms, notably the attention-based Triplet Network architecture, demonstrates superiority across all configurations. Further work is required to delve deeper into attention analysis, to shed-light on its role in encouraging the model to focus in aligning the premise with the conclusion, as well as in linking the premise-conclusion pair with external knowledge. While configurations leveraging Wikipedia outperformed those using other resources, more work is required to evaluate the quality of keywords representing semantic relations between concepts identified from Wikipedia against the standard semantic relation types in ontologies. Furthermore, alternative methods for extracting these keywords should be explored.



## 698 Limitations

699 Although our work presents promising advance-  
700 ments, it also entails the following limitations.

701 **Cross-Domain Evaluation.** Robust evaluation  
702 involving cross-domain evaluation, where models  
703 are trained on one domain and evaluated on a new  
704 domain, is essential for uncovering the robustness  
705 of the proposed approaches. While our evalua-  
706 tion has primarily focused on specific domains or  
707 datasets, cross-domain evaluation can provide in-  
708 sights into the generalisability and adaptability of  
709 the models across diverse domains and real-world  
710 applications.

711 **External Knowledge Alignment and Rela-**  
712 **tion Identification.** More work is required in  
713 aligning the concepts with external resources,  
714 particularly in disambiguating the senses of the  
715 Synsets and Wikipedia page titles. Our current  
716 approach relies on simple similarity measures be-  
717 tween the embeddings of glosses of the resources  
718 and the components, which may lead to missing  
719 alignments and incorrect alignment. Improving  
720 the alignment procedure to account for semantic  
721 ambiguity and variability in external resources is  
722 crucial for enhancing the effectiveness of the pro-  
723 posed approach. Additionally, sophisticated tech-  
724 niques are needed to identify the semantic relation  
725 types existing between Wikipedia hyperlinks. Un-  
726 like ontologies, Wikipedia does not encode explicit  
727 semantic relation types between hyperlinks. There-  
728 fore, developing robust method to identify seman-  
729 tic relations from Wikipedia articles can improve  
730 the quality and relevance of external knowledge  
731 integration in AR prediction.

732 **Interpretability and Explainability.** The ex-  
733 planations provided regarding the performance of  
734 the architectures and external resources are based  
735 on the analysis of empirical results. While empir-  
736 ical analysis is valuable for understanding model  
737 behavior, additional techniques beyond the results  
738 themselves can provide deeper insights into model  
739 performance. Exploring techniques such as model  
740 visualisation, attention mechanisms analysis, and  
741 interpretability methods like LIME (Local Inter-  
742 pretable Model-Agnostic Explanations) (Ribeiro  
743 et al., 2016) or SHAP (SHapley Additive exPlan-  
744 ations) (Lundberg and Lee, 2017) can help uncover  
745 the underlying reasons behind model decisions and  
746 configurations. Complementing empirical analysis  
747 with interpretability techniques can allow a more  
748 comprehensive understanding of model behavior.

## References 749

- 750 Yamen Ajjour, Wei-Fan Chen, Johannes Kiesel, Hen-  
751 ning Wachsmuth, and Benno Stein. 2017. Unit seg-  
752 mentation of argumentative texts. In *Proceedings of*  
753 *the 4th Workshop on Argument Mining*, pages 118–  
754 128.
- 755 Badr AlKhamissi, Millicent Li, Asli Celikyilmaz,  
756 Mona Diab, and Marjan Ghazvininejad. 2022. A  
757 review on language models as knowledge bases.  
758 *arXiv preprint arXiv:2204.06031*.
- 759 Nicholas Asher and Alex Lascarides. 2003. *Logics of*  
760 *conversation*. Cambridge University Press.
- 761 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-  
762 gio. 2014. Neural machine translation by jointly  
763 learning to align and translate. *arXiv preprint*  
764 *arXiv:1409.0473*.
- 765 Rachit Bansal, Milan Aggarwal, Sumit Bhatia, Ji-  
766 vat Neet Kaur, and Balaji Krishnamurthy. 2022.  
767 Cose-co: Text conditioned generative commonsense  
768 contextualizer. *arXiv preprint arXiv:2206.05706*.
- 769 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. Scib-  
770 ert: A pretrained language model for scientific text.  
771 *arXiv preprint arXiv:1903.10676*.
- 772 Teresa Botschen, Daniil Sorokin, and Iryna Gurevych.  
773 2018. Frame-and entity-based knowledge for  
774 common-sense argumentative reasoning. In *Pro-*  
775 *ceedings of the 5th Workshop on Argument Mining*,  
776 pages 90–96.
- 777 Katarzyna Budzynska and Chris Reed. 2011. Whence  
778 inference. *University of Dundee Technical Report*.
- 779 Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-  
780 Gazpio, and Lucia Specia. 2017. Semeval-2017  
781 task 1: Semantic textual similarity-multilingual and  
782 cross-lingual focused evaluation. *arXiv preprint*  
783 *arXiv:1708.00055*.
- 784 Tuhin Chakrabarty, Christopher Hidey, Smaranda  
785 Muresan, Kathy McKeown, and Alyssa Hwang.  
786 2020. Ampersand: Argument mining for  
787 persuasive online discussions. *arXiv preprint*  
788 *arXiv:2004.14677*.
- 789 Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Diana  
790 Inkpen, and Si Wei. 2017. Neural natural language  
791 inference models enhanced with external knowl-  
792 edge. *arXiv preprint arXiv:1711.04289*.
- 793 Artem Chernodub, Oleksiy Oliynyk, Philipp Heidenre-  
794 ich, Alexander Bondarenko, Matthias Hagen, Chris  
795 Biemann, and Alexander Panchenko. 2019. Targer:  
796 Neural argument mining at your fingertips. In *Pro-*  
797 *ceedings of the 57th Annual Meeting of the Associa-*  
798 *tion for Computational Linguistics: System Demon-*  
799 *strations*, pages 195–200.

800	HongSeok Choi and Hyunju Lee. 2018. Gist at semeval-2018 task 12: A network transferring inference knowledge to argument reasoning comprehension task. In <i>Proceedings of The 12th International Workshop on Semantic Evaluation</i> , pages 773–777.	853
801		854
802		855
803		856
804		
805	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> .	857
806		858
807		859
808		
809	Mauro Dragoni, Celia da Costa Pereira, Andrea GB Tettamanzi, and Serena Villata. 2018. Combining argumentation and aspect-based opinion mining: the smack system. <i>AI Communications</i> , 31(1):75–95.	860
810		861
811		862
812		863
813		864
814	Steffen Eger, Johannes Daxenberger, and Iryna Gurevych. 2017. Neural end-to-end learning for computational argumentation mining. <i>arXiv preprint arXiv:1704.06104</i> .	866
815		867
816		868
817		
818	Peter W Foltz, Walter Kintsch, and Thomas K Landauer. 1998. The measurement of textual coherence with latent semantic analysis. <i>Discourse processes</i> , 25(2-3):285–307.	869
819		870
820		871
821		872
822	Michael Fromm, Evgeniy Faerman, and Thomas Seidl. 2019. Tacam: Topic and context aware argument mining. <i>arXiv preprint arXiv:1906.00923</i> .	873
823		874
824		875
825	Debela Gemechu and Chris Reed. 2019. Decompositional argument mining: A general purpose approach for argument graph construction. In <i>Proceedings of the 57st Annual Meeting of the Association for Computational Linguistics</i> , pages 1341–1351.	876
826		877
827		
828		878
829		879
830		880
831	Talmy Givón. 1987. Beyond foreground and background. <i>Coherence and grounding in discourse</i> , 11:175–188.	881
832		882
833		883
834	Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. Breaking nli systems with sentences that require simple lexical inferences. <i>arXiv preprint arXiv:1805.02266</i> .	884
835		
836		885
837		886
838	Yoav Goldberg. 2019. Assessing bert’s syntactic abilities. <i>arXiv preprint arXiv:1901.05287</i> .	887
839		888
840	Barbara J Grosz, Aravind K Joshi, and Scott Weinstein. 1995. Centering: A framework for modelling the local coherence of discourse.	889
841		890
842		891
843	Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2017. The argument reasoning comprehension task: Identification and reconstruction of implicit warrants. <i>arXiv preprint arXiv:1708.01425</i> .	892
844		893
845		894
846		895
847		896
848	Philipp Heinisch, Anette Frank, Juri Opitz, Moritz Plenz, and Philipp Cimiano. 2022. Overview of the 2022 validity and novelty prediction shared task. In <i>Proceedings of the 9th Workshop on Argument Mining</i> , pages 84–94.	897
849		898
850		899
851		900
852		901
		902
		903
	Nora Kassner and Hinrich Schütze. 2019. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly. <i>arXiv preprint arXiv:1911.03343</i> .	
	Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	
	Jonathan Kobbe, Juri Opitz, Maria Becker, Ioana Hulpus, Heiner Stuckenschmidt, and Anette Frank. 2019. Exploiting background knowledge for argumentative relation classification. In <i>2nd Conference on Language, Data and Knowledge (LDK 2019)</i> . Schloss Dagstuhl-Leibniz-Zentrum für Informatik.	
	John Lawrence and Chris Reed. 2020. Argument mining: A survey. <i>Computational Linguistics</i> , 45(4):765–818.	
	Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In <i>International Conference on Machine Learning</i> , pages 1188–1196.	
	Sharon Levy, Emily Allaway, Melanie Subbiah, Lydia Chilton, Desmond Patton, Kathleen McKeown, and William Yang Wang. 2022. Safetext: A benchmark for exploring physical safety in language models. <i>arXiv preprint arXiv:2210.10045</i> .	
	Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. <i>Advances in neural information processing systems</i> , 30.	
	William C Mann and Sandra A Thompson. 1988. Rhetorical structure theory: Toward a functional theory of text organization. <i>Text-interdisciplinary Journal for the Study of Discourse</i> , 8(3):243–281.	
	Daniel Marcu. 2000. <i>The theory and practice of discourse parsing and summarization</i> . MIT press.	
	Tobias Mayer, Elena Cabrio, and Serena Villata. 2020. Transformer-based argument mining for healthcare applications. In <i>ECAI 2020</i> , pages 2108–2115. IOS Press.	
	William Merrill, Yoav Goldberg, Roy Schwartz, and Noah A Smith. 2021. Provable limitations of acquiring meaning from ungrounded form: What will future language models understand? <i>Transactions of the Association for Computational Linguistics</i> , 9:1047–1060.	
	George A Miller. 1995. Wordnet: a lexical database for english. <i>Communications of the ACM</i> , 38(11):39–41.	
	Amita Misra, Pranav Anand, Jean E Fox Tree, and Marilyn Walker. 2017. Using summarization to discover argument facets in online ideological dialog. <i>arXiv preprint arXiv:1709.00662</i> .	

904	Gaku Morio, Hiroaki Ozaki, Terufumi Morishita, and Kohsuke Yanai. 2022. End-to-end argument mining with cross-corpora multi-task learning. <i>Transactions of the Association for Computational Linguistics</i> , 10:639–658.	Peter Potash, Robin Bhattacharya, and Anna Rumshisky. 2017. Length, interchangeability, and external knowledge: Observations from predicting argument convincingsness. In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 342–351.	959
905			960
906			961
907			962
908			963
909	Philippe Muller, Stergos Afantenos, Pascal Denis, and Nicholas Asher. 2012. Constrained decoding for text-level discourse parsing. In <i>Proceedings of COLING 2012</i> , pages 1883–1900.	Peter Potash, Alexey Romanov, and Anna Rumshisky. 2016. Here’s my point: Joint pointer architecture for argument mining. <i>arXiv preprint arXiv:1612.08994</i> .	964
910			965
911			966
912			967
913	R OpenAI. 2023. Gpt-4 technical report. arxiv 2303.08774. <i>View in Article</i> , 2:13.		968
914			969
915	Andreas Peldszus and Manfred Stede. 2013. Ranking the annotators: An agreement study on argumentation structure. In <i>Proceedings of the 7th linguistic annotation workshop and interoperability with discourse</i> , pages 196–204.	Ellen F Prince. 1981. Toward a taxonomy of given-new information. <i>Radical pragmatics</i> .	970
916			971
917			972
918		Nils Reimers and Iryna Gurevych. 2019. Sentencebert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .	973
919			974
920	Andreas Peldszus and Manfred Stede. 2015a. Joint prediction in mst-style discourse parsing for argumentation mining. In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 938–948.	Nils Reimers, Benjamin Schiller, Tilman Beck, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2019. Classification and clustering of arguments with contextualized word embeddings. <i>arXiv preprint arXiv:1906.09821</i> .	975
921			976
922			977
923			978
924			979
925	Andreas Peldszus and Manfred Stede. 2015b. Joint prediction in mst-style discourse parsing for argumentation mining. In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 938–948.	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " why should i trust you?" explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pages 1135–1144.	980
926			981
927			982
928			983
929			984
930	Andreas Peldszus and Manfred Stede. 2016. Rhetorical structure and argumentation structure in monologue text. In <i>Proceedings of the Third Workshop on Argument Mining</i> , pages 103–112.	Ruty Rinott, Lena Dankin, Carlos Alzate Perez, Mitesh M Khapra, Ehud Aharoni, and Noam Slonim. 2015. Show me your evidence-an automatic method for context dependent evidence detection. In <i>Proceedings of the 2015 conference on empirical methods in natural language processing</i> , pages 440–450.	985
931			986
932			987
933			988
934	Isaac Persing and Vincent Ng. 2016. End-to-end argumentation mining in student essays. In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1384–1394.	Ramón Ruiz-Dolz, José Alemany, Stella M Heras Barberá, and Ana García-Fornes. 2021. Transformer-based models for automatic identification of argument relations: A cross-domain evaluation. <i>IEEE Intelligent Systems</i> , 36(6):62–70.	989
935			990
936			991
937			992
938			993
939	Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? <i>arXiv preprint arXiv:1909.01066</i> .	Tara Safavi and Danai Koutra. 2021. Relational world knowledge representation in contextual language models: A review. <i>arXiv preprint arXiv:2104.05837</i> .	994
940			995
941			996
942			997
943	Mohammad Taher Pilehvar, David Jurgens, and Roberto Navigli. 2013. Align, disambiguate and walk: A unified approach for measuring semantic similarity. In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics</i> , pages 1341–1351.	Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 815–823.	998
944			999
945			1000
946			1001
947			1002
948			1003
949	Moritz Pleniz, Juri Opitz, Philipp Heinsch, Philipp Cimiano, and Anette Frank. 2023. Similarity-weighted construction of contextualized commonsense knowledge graphs for knowledge-intensive argumentation tasks. <i>arXiv preprint arXiv:2305.08495</i> .	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 31.	1004
950			1005
951			1006
952			1007
953			1008
954			1009
955	Stanislas Polu, Jesse Michael Han, Kunhao Zheng, Mantas Baksys, Igor Babuschkin, and Ilya Sutskever. 2022. Formal mathematics statement curriculum learning. <i>arXiv preprint arXiv:2202.01344</i> .	Christian Stab and Iryna Gurevych. 2017. Parsing argumentation structures in persuasive essays. <i>Computational Linguistics</i> , 43(3):619–659.	1010
956			1011
957			1012
958			1013



Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.

Maitte Taboada and William Mann. 2006. Rhetorical structure theory: Looking back and moving ahead. *Discourse studies*, 8(3):423–459.

Dietrich Trautmann. 2020. Aspect-based argument mining. *arXiv preprint arXiv:2011.00633*.

Jacky Visser, Barbara Konat, Rory Duthie, Marcin Koszowy, Katarzyna Budzynska, and Chris Reed. 2019. Argumentation in the 2016 us presidential elections: annotated corpora of television debates and social media reaction. *Language Resources and Evaluation*, pages 1–32.

Su Wang, Greg Durrett, and Katrin Erk. 2018. Modeling semantic plausibility by injecting world knowledge. *arXiv preprint arXiv:1804.00619*.

## A Appendix

We provide additional details regarding the methodology and experimental setups used in our study.

### A.1 Experiment Setup

#### A.1.1 Training Procedure

**Hyper-parameters:** We employ Adam optimisation (Kingma and Ba, 2014) to minimise the cost function. The learning rate is set to  $2e^{-5}$  with a batch size of 16. Categorical cross-entropy loss was used as the loss function.

**Gradient Clipping:** To prevent exploding gradients during training, we apply gradient clipping. We use a maximum gradient norm (max\_grad\_norm) parameter set to 1.0 to determine the threshold for gradient clipping.

**Warm-up and Learning Rate Schedule:** We employed a linear warm-up strategy for the learning rate. The number of warm-up steps is set to 10% of the total training steps. Following the warm-up phase, the learning rate schedule is determined by a lambda function. This function linearly increases the learning rate during the warm-up phase and decreases it linearly thereafter.

**Early Stopping:** We implement early stopping to prevent overfitting and to determine the optimal number of epochs. This technique involves continuously monitoring the loss and F-score on the validation set throughout training. If there is a sustained degradation in performance over consecutive epochs, training is terminated to prevent the model from being influenced by noise present in the training data.

Dataset	Training		Validation		Test	
	RA	CA	RA	CA	RA	CA
AAEC	4235	411	605	59	1210	117
US2016	353	55	51	8	101	16
MTC	190	120	27	17	55	34
AbstRCT	1603	241	229	34	458	69

Table 3: Distribution of support and attack relations in the train, validation, and test splits across the datasets.

#### A.1.2 Input Setup

For the baseline sequence classification configurations, we concatenate the premise to the conclusion using a special token [SEP]. In the Siamese architecture, one of the sub-networks takes the concatenation of the premise and conclusion based on the special token [SEP], while the other takes the concatenation of the paths. The paths are concatenated using the special token [SEP].

The number and length of the paths between the components of the ADUs vary, with some ADUs not involving any path at all. For ADU pairs involving a large number of paths exceeding the maximum sequence length, we concatenate the paths until the maximum sequence length is reached. In such cases, we sort the paths based on their frequency. The concatenation process starts from the most frequent paths until the maximum sequence length is reached.

#### A.1.3 Fully Fine-tuned Baseline LLM Configuration

For the fully fine-tuned baseline LLM configuration, we utilise the HuggingFace implementation of BERT for sequence classification (bert-base-uncased<sup>4</sup>). We experimented with two variants of BERT: bert-base-cased and bert-large-uncased. Our experiments revealed that bert-base-uncased consistently provided better performance compared to bert-base-cased. In the baseline Siamese architecture, each sub-network independently encodes the ADUs.

### A.2 External Knowledge Extraction

#### A.3 ADU Decomposition

To identify the functional components (C and A) from ADUs, we adopt a sequence labeling approach following the methodology outlined by Gemechu and Reed (2019). Unlike Gemechu and

<sup>4</sup><https://huggingface.co/google-bert/bert-base-uncased>

Dataset	Total C	Total A	Unique C	Unique A
AAEC	9875	6789	5634	4356
US2016	3225	1737	1854	1566
MTC	870	589	756	470
AbstRCT	7343	6432	5554	4546

Table 4: Distribution of target concepts (C) and aspects (A) across the datasets.

Reed (2019) method, which employs a convolutional neural network (CNN), we fine-tune BERT for token classification using their dataset annotated with the BIOES-style sequence labeling scheme, outperforming their top-performing method by 3% and achieving a macro F-score of 0.784. We utilise the HuggingFace implementation of BERT (bert-base-uncased<sup>5</sup>). The inputs are padded to 256 maximum size. We use the train-test split in the original dataset. Training is conducted over 6 epochs, and evaluation is reported as the average performance over 3 runs of the experiment on the test dataset. Using the fine-tuned model, we identify the functional components of ADUs, and the distribution of these components is presented in Table 4.

### A.3.1 Alignment of Ontologies and Wikipedia

For aligning ontologies and Wikipedia with the components of ADUs, the cosine similarity between the embeddings of the components and the Synsets of the ontologies or the corresponding Wikipedia page title is used. Additionally, we utilise the similarity between the concepts and the gloss texts of the respective sources for disambiguating senses, for concepts involving multiple senses.

### A.3.2 Similarity Threshold

We leverage embeddings derived from Sentence-transformers, particularly the *all-roberta-large-v1*<sup>6</sup> variant, for determining similarity. We set a similarity threshold of  $\beta = 0.80$  based on experimental comparisons of similarity scores between related and unrelated text pairs in the STSB dataset<sup>7</sup>.

The dataset is originally annotated on a scale of 0-5 based on the degree of similarity. We transform the original 5-class labels into binary labels,

<sup>5</sup><https://huggingface.co/google-bert/bert-base-uncased>

<sup>6</sup><https://huggingface.co/sentence-transformers/all-roberta-large-v1>

<sup>7</sup><https://huggingface.co/datasets/nyu-ml/glue/viewer/stsb/train>

where labels below 4 are considered unrelated, and labels 4 and above are deemed related. In the original annotation rubric provided by SemEval-2017 (Cer et al., 2017), label 3 indicates sentences that are roughly equivalent, but some important information differs. However, we found that this definition allows for a certain degree of looseness in similarity assessment. Consequently, to impose a stricter criterion for similarity, we decided to raise the threshold from label 3 to label 4. To this end, we calculate the similarity between the sentence pairs in the training dataset and select the threshold yielding the highest F1-score. We compute F1-scores at 20 similarity threshold points (ranging from 0 to 1 with increments of 0.05), as outlined in Algorithm 1.

---

#### Algorithm 1 Find Optimal Similarity Threshold

---

**Require:** List of sentence pairs  $(s_1, s_2)$

**Ensure:** Threshold

```

best_threshold ← min_thr
max_f_score ← 0
for thr ← min_thr to max_thr by thr_step do
    tp ← 0
    fp ← 0
    fn ← 0
    for each sentence pair  $(s_1, s_2)$  in data do
        sim_score ← sim score( $s_1, s_2$ )
        if similarity_score ≥ thr then
            if pair is similar then
                tp ← tp + 1
            else
                fp ← fp + 1
            end if
        else
            if pair is dissimilar then
                tn ← tn + 1
            else
                fn ← fn + 1
            end if
        end if
    end for
    precision ←  $\frac{tp}{tp+fp}$ 
    recall ←  $\frac{tp}{tp+fn}$ 
    f1_score ←  $2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$ 
    if f1_score > max_f_score then
        max_f_score ← f1_score
        best_threshold ← thr
    end if
end for
return best_threshold

```

---

### 1155 A.3.3 Search Depth Threshold

1156 To estimate the optimal depth threshold for navigat-  
1157 ing through the knowledge graphs, we employ the  
1158 following procedure: we randomly select 20 pairs  
1159 of concepts and initiate a complete search from one  
1160 concept to identify paths leading to the other. This  
1161 provides a total of 728 paths with various depths  
1162 from the three resources. Human annotators then  
1163 evaluate the relevance of the retrieved paths based  
1164 on a binary value indicating if the path is relevant  
1165 to the given AR or not. The cumulative F1-score  
1166 at each depth is computed based on the total num-  
1167 ber of relevant paths retrieved up to that depth.  
1168 The depth with the highest cumulative F-score is  
1169 chosen as the optimal threshold. Accordingly, the  
1170 threshold of  $\alpha = 5$  yielded the highest score.

### 1171 A.3.4 Extracting keywords encoding 1172 semantic relation types from 1173 Wikipedia.

1174 The AllenNLP semantic role labeling (SRL)<sup>8</sup> is  
1175 used to parse sentences and assign semantic roles  
1176 to each word. This enables to extract phrases  
1177 linking the concepts of interest along the subject-  
1178 predicate structure of the sentences. To mention, if  
1179 one concept is identified as the *agent* and another  
1180 as the *patient*, the phrase denoting the action per-  
1181 formed by the agent on the patient is used as the  
1182 relation type between them.

1183 Consider the concepts *exercise* and *cardiovascu-*  
1184 *lar diseases* in the sentence:

1185 *According to the American Heart Associ-*  
1186 *ation, exercise reduces the risk of cardio-*  
1187 *vascular diseases, including heart attack*  
1188 *and stroke.*

1189 Below is the output of SRL for this sen-  
1190 tence (the concepts are highlighted in light  
1191 blue while the keywords representing the re-  
1192 lation type are highlighted in red): {'verbs':  
1193 [{'verb': 'According', 'description':  
1194 '[V:According] to the American Heart  
1195 Association , exercise reduces the risk  
1196 of cardiovascular diseases , including  
1197 heart attack and stroke', 'tags': ['B-V',  
1198 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O',  
1199 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O',  
1200 'O', 'O', 'O']}], {'verb': 'reduces',  
1201 'description': '[ARGM-ADV: According

<sup>8</sup> [https://docs.allennlp.org/v0.9.0/api/allennlp.models.semantic\\_role\\_labeler.html](https://docs.allennlp.org/v0.9.0/api/allennlp.models.semantic_role_labeler.html)

1202 to the American Heart Association] ,  
1203 [ARG0: **exercise**] [V: **reduces**] [ARG1:  
1204 the risk of **cardiovascular diseases** ,  
1205 including heart attack and stroke'] ,  
1206 'tags': ['B-ARGM-ADV', 'I-ARGM-ADV',  
1207 'I-ARGM-ADV', 'I-ARGM-ADV', 'I-ARGM-ADV',  
1208 'I-ARGM-ADV', 'O', 'B-ARG0', 'B-V',  
1209 'B-ARG1', 'I-ARG1', 'I-ARG1', 'I-ARG1',  
1210 'I-ARG1', 'I-ARG1', 'I-ARG1', 'I-ARG1',  
1211 'I-ARG1', 'I-ARG1', 'I-ARG1']}, {'verb':  
1212 'including', 'description': 'According  
1213 to the American Heart Association  
1214 , exercise reduces the risk of [ARG2:  
1215 cardiovascular diseases] , [V: including]  
1216 [ARG1: heart attack and stroke'] ,  
1217 'tags': ['O', 'O', 'O', 'O', 'O',  
1218 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O',  
1219 'B-ARG2', 'I-ARG2', 'O', 'B-V', 'B-ARG1',  
1220 'I-ARG1', 'I-ARG1', 'I-ARG1']},  
1221 'words': ['According', 'to', 'the',  
1222 'American', 'Heart', 'Association', ',',  
1223 'exercise', 'reduces', 'the', 'risk',  
1224 'of', 'cardiovascular', 'diseases', ',',  
1225 'including', 'heart', 'attack', 'and',  
1226 'stroke']}]

1227 We navigate through the SRL output to identify  
1228 the predicate-argument structures connecting both  
1229 concepts (*exercise* and *cardiovascular diseases* in  
1230 this case). We then use predefined rules to extract  
1231 keywords encoding the semantic relations existing  
1232 between the concepts. To mention, if one concept  
1233 is part of *ARG0* and the other being part of *ARG1*,  
1234 the predicate term is used as the relation type. In  
1235 the example output above, the predicate term repre-  
1236 senting the semantic relation type is *reduces*. More  
1237 examples are provided below. The pair of concepts  
1238 are highlighted in light blue and the relation type  
1239 highlighted in red:

#### 1240 1. Concept Pair: Exercise, Cardiovascular 1241 Diseases

- 1242 • **Semantic Relation:** **increase**
- 1243 • **Sentence:** "Low levels of physical exer-  
1244 cise increase the risk of cardiovascular  
1245 diseases mortality."
- 1246 • **Predicate structure:** [ARG0: Low lev-  
1247 els of **physical exercise**] [V: **increase**]  
1248 [ARG1: the risk of **cardiovascular dis-**  
1249 **eases** mortality].

#### 1250 2. Concept Pair: Exercise, Cardiovascular 1251 Profiles



- **Semantic Relation:** **leads**
- **Sentence:** "Studies have shown that since heart disease is the leading cause of death in women, regular exercise in aging women leads to healthier cardiovascular profiles."
- **Predicate structure:** Studies have shown that [ARGM-CAU: since heart disease is the leading cause of death in women], [ARG0: regular **exercise** in aging women] [V: **leads**] [ARG2: to healthier **cardiovascular profiles**].

### 3. Concept Pair: Innovation, Economy

- **Semantic Relation:** **is**
- **Sentence:** "Given the noticeable effects on efficiency, quality of life, and productive growth, innovation is a key factor in society and economy."
- **Predicate structure:** [ARGM-ADV: Given the noticeable effects on efficiency , quality of life , and productive growth], [ARG1: **innovation**] [V: **is**] [ARG2: a key factor in society and **economy**]

### 4. Concept Pair: Sustainable Energy, Renewable Energy

- **Semantic Relation:** **involves**
- **Sentence:** "Sustainable energy involves increasing production of renewable energy, making safe energy universally available, and energy conservation."
- **Predicate structure:** [ARG2: **Sustainable energy**] [V: **involves**] [ARG1: increasing production of **renewable energy** , making safe energy universally available, and energy conservation]

#### A.3.5 Filtering semantic relations.

A total of 7959 unique relation types are extracted. Please note that similar relation types like “leads to”, “leads” and “can lead to” are counted as different relation types, as we only consider surface-level counts. To exclude arbitrary paths between concepts only relation types with a frequency greater than “ $m=3$ ” are considered. This yields a total of 1488 unique relation types. However, as can be seen in Table 7, manual analysis revealed similarities among certain tuples; for example, the relation

type “influences” is similar to other relations like “contributes to”, “leads to”, and “results in”.

Some concepts are directly related through single relation type (one-hop path), while others are indirectly connected via paths involving multiple relation types (multi-hop). See examples in Table 6. The length of these paths ranges from 1 (indicating direct links between concepts) to 5 (the maximum search depth), with an average path length of 1.9.

#### A.3.6 External Resource Evaluations

**Ontology and Wikipedia:** We analyse the three resources to showcase their contributions in terms of coverage and the quality of connections.

**Coverage.** The aim is to show the proportion of pairs connected exclusively by one resource but not by others. To this end, we randomly select 500 unconnected pairs from each resource and generate a heatmap illustrating the ratio of pairs exclusively connected by each resource compared to the others to identify which resource is most effective in covering concepts absent in others. On average, Wikipedia covers 37% of pairs unconnected in both WordNet and ConceptNet, while only 7% of the concepts missing in Wikipedia are covered by both WordNet and ConceptNet. Please note that pairs of concepts connected by relation types occurring less than three times are considered unconnected.

**Connection quality.** We further analyse the quality of the paths by ranking component pairs based on the number of paths linking them from each respective resource. From this ranking, we select the top 25 most connected and 25 least connected pairs from each resource for detailed evaluation. Two annotators independently rate the relevance of these paths by assigning binary labels, reflecting their subjective assessments of the paths’ pertinence to the AR between the ADUs. The evaluation reveals that Wikipedia is the top-rated source for both well-connected and least connected paths, followed by ConceptNet.

**GPT-generated paths:** As shown in Table 2, configuration utilising GPT-generated paths show higher accuracy but lower precision. Of the total errors observed, 79% are identified as false positives for approaches using GPT-generated paths in predicting AR, while the average false positive rate for the other three external resources is 53%. To further investigate, we randomly select 20 errors and engage two human annotators to jointly analyse the paths connecting the pair of ADUs.

Out of the 20 errors, the paths for the 14 of the

ADUs	Components, Chain of Hyperlinks
ADU1 - Trump tax cut is the biggest since Ronald Reagan; ADU2 - It will create tremendous numbers of new jobs	<b>Chain of Hyperlinks for the components, Tax (C), Jobs (A):</b> - Job → Working hour system → Income tax → Tax - Job → Labor economics → Economic policy → Tax - Job → Unemployed → Tariff → Tax
ADU1 - Clinton is going to approve one of the biggest tax cut in history; ADU2 - Clinton is going to drive business out	<b>Chain of Hyperlinks for the components, Business (C), Tax cut (A):</b> - Business → Adam Smith → Economic theory → Tax cut - Business → Adam Smith → Neoliberalism → Tax cut - Business → Corporate tax → Effect of taxes and subsidies on price → Tax cut

Table 5: Examples showing the connection between ADUs via the chain of hyperlinks linking their components.

errors are categorised as contextually irrelevant for the ADU pairs. The primary reason cited by the annotators for the irrelevant paths indicates that while the generated paths make logical sense and provide valid lines of reasoning between the ADUs, there were no AR between these ADUs as originally annotated in the dataset.

For example, consider the pair of ADUs "*Researches into humanities and art still need large amount of money*" and "*a government should spare effort on young children education as well as universities*", taken from the argument graph depicted in Figure 2 (taken from AAEC dataset). GPT identified the following semantic relation paths linking the concepts "*money*" and "*young children education*":

- **money** → *facilitates* → **technology adoption** → *enables* → **digital literacy programs** → *encourages* → **young children education**
- **money** → *stimulates* → **philanthropic endeavors** → *cultivates* → **community partnerships** → *fosters* → **early childhood learning opportunities**
- **money** → *fuels* → **economic growth** → *stimulates* → **job creation** → *expands access to* → **early childhood education**
- **money** → *drives* → **philanthropic activities** → *funds* → **charitable organisations** → *supports* → **early childhood education initiatives**
- **money** → *empowers* → **local communities** → *cultivates* → **community engagement** → *enhances* → **early childhood learning environments**
- **money** → *encourages* → **resource allocation** → *drives* → **research and development** → *inspires* → **pedagogical advancements**

Despite these two ADUs not being linked by AR in the gold dataset, the paths between the concepts they address mimic the paths typically associated with ADUs involving AR. However, the reasoning conveyed by these paths is categorised as unintended, as they involve reasoning diverging from the original argument, and the AR between these ADUs is absent in the gold dataset.

The same applies to the paths identified for the concepts *money* and *future* addressed by the pair of ADUs: "*Researches into humanities and art still need large amounts of money*" and "*both are crucial on the way to a brighter future*".

- **money** → *allows for* → **travel experiences** → *impacts* → **cultural enrichment** → *shapes* → **future memories**
- **money** → *allows for* → **excessive spending** → *impacts* → **short-term pleasure** → *shapes* → **future goals**
- **money** → *initiates* → **investment opportunities** → *promotes* → **financial stability** → *contributes to* → **future security**
- **money** → *used for* → **educational funding** → *influences* → **career advancement** → *impacts* → **future**
- **money** → *is used for* → **investment in property** → *helps in* → **wealth accumulation** → *contributes to* → **future**
- **money** → *is used for* → **infrastructure development** → *helps in* → **urban planning** → *contributes to* → **future city growth**
- **money** → *allows for* → **business expansion** → *impacts* → **economic prosperity** → *shapes* → **future**
- **money** → *leads to* → **business expansion** → *is linked to* → **economic growth** → *impacts* → **future prosperity**

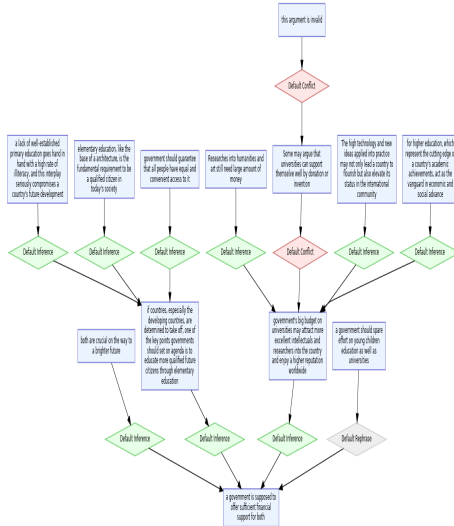


Figure 2: Example argument graph.

- **money** → *is essential for* → **scientific research** → *contributes to* → **technological advancement** → *shapes* → **future innovation**

## A.4 GPT for Path Generation and AR Prediction

### A.4.1 Experimental Settings

We utilise the chat completion configuration of ChatGPT-4 for two tasks: (a) generating the chain of semantic relation between ADU components, and (b) predicting AR.

- Configurations:** We use GPT-4 based on gpt-3.5-turbo-instruct. We set a maximum token limit of 2048, a temperature of 0.7, a top-p probability of 0.9.
- Prompts Strategy:** We explored two strategies: zero-shot and few-shot prompts. In the zero-shot setting, only instruction based prompts without examples are used. Based on the insightful recommendation from the reviews, we also try few-shot setup, where specific examples are provided as part of the instruction. Interestingly, our analysis revealed that the example-based experiment achieved a 1.3%, 2.1% higher score compared to the zero-shot prompt in the AR prediction and path generation, respectively. As a result, our experiment is based on example-based prompting. We create prompt templates that include instructions and two examples randomly selected from a list of examples. These examples consist of ADU pairs, concept pairs

identified from the ADUs, and paths obtained from three external resources. The placeholder variables in the template are replaced with the ADUs, concepts, and paths.

**Prompt Design for Path Generation.** GPT-4 is tasked with generating paths between components of ADUs using the following template:

You are a model trained to identify chains of semantic relations between a pair of concepts derived from two sentences (ADU1 and ADU2). Given concepts c1 and c2 extracted from ADU1 and ADU2 respectively, your goal is to identify chains of semantic relation types connecting these concepts. These relations may include meronymy, hypernymy, hyponymy, cause-effect, or any other valid semantic relation. Concepts are often indirectly linked via intermediate concepts and their relations. Include both direct and indirect paths between the concepts whenever possible, using only the context provided by the ADU pairs. Provide up to 10 paths if possible; otherwise, return an empty list. Each relation type should be represented as a tuple in the format (concept1, relation type, concept2). For indirect paths involving multiple tuples, return them as a list of tuples. Example 1: between the concepts "USA" and "NAFTA" identified from the pair of ADUs "USA is in deep trouble" and "NAFTA agreement is defective", a valid list of paths could be, [{"USA, part-of, NAFTA"}], [{"USA, has, trade deal"}, {"trade deal, instance of, NAFTA"}]. Example 2: between the concepts {c1} and {c2} identified from the pair of ADUs {ADU1} and {ADU2}, the list of paths should include, {list\_path}. Provide your answer as a python list.

Note: In Example 1, we show an actual example, but it should be a placeholder variable in the prompt template, as shown in Example 2.

**Prompt Design for Zero-Shot AR Prediction:** We prompt GPT-4 to classify the relationship between the ADUs as supporting, contradicting, or



1506 having no clear AR using the following prompt  
 1507 template.

1508 You are a 3-class classifier model tasked with  
 1509 assigning a label to the argument  
 1510 relation between two argument units  
 1511 (argument 1 and argument 2).

1512 Classify the following pair of arguments,  
 1513 argument 1: {ADU\_1}  
 1514 argument 2: {ADU\_2},  
 1515 into:

1516 "support" (if argument 1 supports  
 1517 argument 2),  
 1518 "contradict" (if argument 1 attacks  
 1519 argument 2),  
 1520 and "None" (if no argument relation exists  
 1521 between argument 1 and argument 2).

1522 Please enter:

1523 1 - for support,  
 1524 2 - for contradict,  
 1525 0 - for None relation.

1526 Examples from each argument  
 1527 relation types are provided below:

1528 Example 1: the argument relation between  
 1529 the argument "people feel, when they have  
 1530 been voicing opinions on different matters,  
 1531 that they have been not listened to", and  
 1532 the argument "people  
 1533 feel that they have been treated  
 1534 disrespectfully on all sides of the  
 1535 different arguments and disputes going on"  
 1536 is support, and hence prediction label is 1.

1537 Example 2: The argument relation between  
 1538 "there would be no non-tariff barriers  
 1539 with the deal done with the EU" and  
 1540 the argument "there are lots of  
 1541 non-tariff barriers  
 1542 with the deal done with the EU"  
 1543 is contradiction, and  
 1544 hence prediction label is 2.

1545 Note: We use the actual examples to show sup-  
 1546 port and contradiction relations, which should be a  
 1547 placeholder variable in the final prompt template.

## 1548 A.5 Multi-Network Architectures

1549 The encoder blocks within the multi-networks are  
 1550 constructed using the HuggingFace implementa-  
 1551 tion of BERT (bert-base-uncased)<sup>9</sup>. In all con-  
 1552 figurations, we employ 8 attention heads to align  
 1553 with the standard transformers implementations.

<sup>9</sup><https://huggingface.co/google-bert/bert-base-uncased>

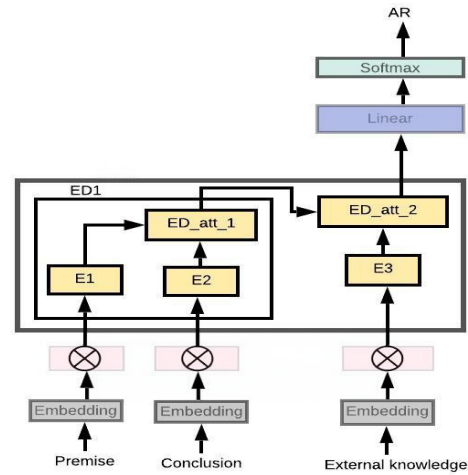


Figure 3: Triplet-networked with attention layers.

### A.5.1 Attention Mechanisms in Multi-Network Architectures

The Triplet Network architecture is aimed to encode the individual components of ADUs as well as the external knowledge paths connecting them. The architecture consists of three sub-networks, each focusing on a different aspect of the input:

- **Sub-Network 1:** Encodes the premise.
- **Sub-Network 2:** Encodes the conclusion.
- **Sub-Network 3:** Encodes the paths between components of ADUs.

Two attention layers are used to attend to the alignment between the inputs (premise, conclusion and external knowledge).

1. **First Attention Layer:** This layer attends to the alignment between the premise and conclusion based on the outputs of Sub-Networks 1 and 2, respectively.
2. **Second Attention Layer:** Building upon the output of the first attention layer, this layer aligns the information from the first attention layer with the external knowledge provided by Sub-Network 3.

Finally, the outputs of the attention layers are averaged to obtain a unified representation of the input, which is then passed through a linear classification layer to predict AR. We experiment with two configurations for representing the ADUs and the external knowledge as input to the attention

1583 layer: (a) using the final output of the [CLS] token  
1584 and (b) using the mean of the last hidden layer of  
1585 all tokens from BERT's output. Consistently, the  
1586 mean of the last hidden layer of all tokens yields  
1587 superior performance compared to the [CLS] to-  
1588 ken.

Path	Path	Path	Path
related to → leads to → affects	related to → related to	synonym	involves
related to	affects → associated with → impacts	synonym → related to	causes
has	is related to	leads to	is a → related to → related to
impacts	related to → involves	contains	is a → involves
is a → is a → is a	part of	related to → part of	causes → related to
influences → affects	antonym → hyponym → hyponym	associated with	involves → related to
leads to → results in	entails → entails	is a → has	can lead to
implies	related to → related to → related to	affects → influences	causes → leads to
has → includes	part of → includes	related to → includes	related to → related to → causes
supports	synonym → hypernym → hyponym	entails → involves	related to → entails
causes → affects	is a → belongs	is associated with → involves	regulate
related to → impacts	can result in	is an umbrella term → is related to	leads to → involves
found in	administers → impacts → involves	affiliated with → associated with	aids → helps → helps
developed through → facilitated by → leads to	discussed in → is a → lead to	are → show → can lead to → can result in	assessment of → measure of → related to
associated with → is a type of → can be	be used for → have quality	can be obtained → is extended for → is a type of	occur in → experiencing → necessitate
can provide → may lead to → changes	common in → generally involves	includes → example of	determines → affects
empowerment through → instance of	entails → sustain	entails → includes → involves	entails → is a → is a → is a
entails → is required for → can lead to → can result in	experienced → includes	often favours → which stems from	fosters → crucial for
give → way to	impacts → evaluates	influences → lead to → affect	influences → are reflected by
influences → is achieved by	influences → importance of → includes	involves → involvement of → can come under	involves → is represented by
involves → brings → used for	is a factor in → generates → can include	symptom of → includes → has code	is a type of → may require → is associated with
is a → is delivered through → facilitates	is essential to → has an impact → results in	is important for → used in → opportunity for	is involved in → has phase → is type for
is often accompanied by → is similar to	is often associated with → has effects on → are linked to	related → shapes → contribute to → are crucial for	required → necessary for
supported by → promotes → reduces → are important	is the goal of → can include	is type of → can involve → is related to	is a → involves → relieves
live in → has	may bring → followed by → result in	may lead to → requires → found	necessitates → involves → category of
offer → facilitate → contributes to → aids in	opposite of → causes → leads to	organised by → hold	participates in → can involve
provide → attend	provides challenge in	provides → enable	provides → includes → develops → can lead to
refers to → impacts → affects	related to → can lead to → results in → results in	related to → improves → essential → crucial for	related to → indicates → compared to
leads to → is likely to → often achieved by	represents → causes	require → achieved by → help	shares → comprises of
duration → has value	used for → associated with → part of	convey → interpreted by → part of → makes up	requires → causes → leads → necessary

Table 6: Examples of semantic relation paths.



Relation	Relation	Relation	Relation	Relation	Relation
related to	involves	hyponym	antonym	synonym	is a
has	results in	affects	related term	leads to	can lead to
causes	entails	associated with	part of	includes	hypernym
influences	impacts	is related to	contributes to	include	is a type of
instance of	can result in	requires	related	connected to	contains
have	require	can be	involve	used for	implies
consists of	versus	are	lead to	greater than	affect
influence	entailment	type of	causes desire	linked to	cause effect
opposite of	relates to	is essential for	is similar to	impact	may lead to
supports	provides	can involve	is crucial for	result in	is part of
cause	essential for	may result in	symptom of	is a form of	comparison
facilitates	enhances	motivated by	contribute to	can cause	similar
used in	experience	is important for	enables	influenced by	drives
at location	provide	are part of	percentage	may involve	comprises
synonym of	opposite	indicates	describes	attribute	attend
refers to	is	can include	determines	promotes	has instance
use	participate in	entails action	treated with	utilises	measured by
shapes	pertains to	is connected to	necessitates	encourages	improves
antonym of	is used in	similar to	measures	is used for	represents
chain map	offers	is influenced by	treat	may cause	of
negation	has context	shape	consist of	has property	example of
motivates	are associated with	equals	can affect	location	has quality
enhance	relate to	affected by	undergo	may include	contributes to
belongs to	can influence	found in	addresses	impact on	create
seek	possess	increases	can impact	receive	compares
opposes	member of	feature	subset of	concerns	is required for
derived from	is a part of	has attribute	resulted in	comprise	is equivalent to
treatment for	used by	activity	treatment	regulates	correlates with
enable	produces	is necessary for	triggers	target	ensures
inspires	correlated with	impacted by	inspire	helps in	has duration
need	is a factor in	component of	is known for	modifies	related term of
measure	treated by	is less than	characteristic of	has numeric value	covers
is about	is a symptom of	employs	entail	located in	has part
located near	shows	are crucial for	focuses on	engage in	depends on
pursue	is needed for	brings about	motivated by goal	cause of	can have
can	associated with	are related to	range	involved in	utilise
targets	means	attribute of	benefits	characterised by	measured by
spouse of	is a measure of	side effect of	comparative of	can be influenced by	is vital for
has member	occurs in	evaluate	implement	allows for	has symptom
is equal to	less than	belong to	linked to	involves	encourage
fosters	component of	known for	capable of	is key to	helps
interact with	drive	constitute	relies on	comprises of	meronym
defines	generates	correlate with	determine	has subevent	represent
is an umbrella term	compare	has prerequisite	facilitate	desires	percentage of
a type of	acquired through	address	addressed	advocates for	agent
agent of	are important for	aligns with	aid in	are used in	assesses
attract	belongs to group	belong to	boosts	achieved through	be found in
can be represented	can contribute to	can create	can enhance	can develop into	can lead to
can occur in	can require	can stimulate	capability of	category of	caused by
combined with	complication of	concept in	connects	conceptually related to	deals with
essential for	establish	evaluated by	examines	example of	exhibit
has percent	has range	has impact on	have activity level	helps to get more	helps gain
holonym	impacts result in	implemented by	imposes	indicate	induces
inhibits	is a medication	is a metric for	is a side effect of	is a source of	is a subclass of
is a symptom of	is a unit of time	is a way to	is beneficial for	is critical for	is defined by
is fundamental for	is funded by	is greater than	is key to	is opposite of	is perceived as
is quantified by	is significant for	has value	is the target of	is treated with	is used to assess
is treated by	lack of	lacks	live in	location of	made by some
made of	manifests as	negatively impacts	numerical value	often involve	outcome of
percentage value	play a role in	plays a role in	politician	possesses	prevents
process of	produce	promoted by	provide access to	provided by	qualifies
quantity	reduces	reflect	reflects	regulated by	rely on
restrict	results in state	show	stimulates	studies	suggests
superlative	to be gained by	tool for	treats	treatment includes	treated by
treatment involves	treatment with	trigger	utilised for	increased expression of	yield

Table 7: Examples of semantic relation types. We normalised (lower cased and expanded relation types like IsA, RelatedTo, HasProperty) the relation types for consistency across the external resources.