

Cross-Domain Argument Mining Midterm Report

Oscar Morris

January 2025

1 Introduction

Argument Mining (AM) is the automatic identification and extraction of argumentative structures from natural language discourse [1]. In order to achieve this, the discourse can first be segmented into Elementary Discourse Units (EDUs), atomic units of the discourse. The EDUs can then be classified into argument and non-argument, the argumentative EDUs can be termed Argumentative Discourse Units (ADUs). Following this, clausal properties can be determined for each ADU (e.g. for an ADU X, is X evidence? or is X a premise?). Finally, the relations between ADUs can be determined (e.g. does one ADU support or attack another?). It is with this final step that this research is concerned.

The task of identifying the relational properties between already identified ADUs is known as Argument Relation Identification (ARI). ARI concerns itself both with identifying whether a relation exists between ADUs, and also classifying the type of relation, (typically either support or attack). This typically becomes a 3-class classification problem, concerning no relation, support and attack [2].

There are, however, varying opinions on whether having only two relational classes is enough, Inference Anchoring Theory (IAT) presents three different classes [3], [4]. In IAT, the ‘support’ class is split into ‘inference’ or RA and ‘rephrase’ or MA. The ‘attack’ relation is also termed ‘conflict’ or CA. Using the IAT relation classes produces a 4-class classification problem.

Transformer models [5] have significantly improved performance in many tasks across natural language processing, including AM [6], [7]. Subsequent task-agnostic pretraining approaches have allowed these models to be fine-tuned on a wide variety of tasks. As the RoBERTa model [8] has been shown to perform well in the ARI task [6] it is used for this research to allow comparison with previous work.

Multimodal models have also proven useful in audio-based natural language tasks, where transcripts of the spoken word are also available. Generally it has appeared that combining both audio and textual features has seen improved performance over unimodal techniques [9], [10].

2 Datasets

All datasets used in this project are available as corpora on AIFdb¹. Using consistently annotated Argument Interchange Format (AIF) data allows many different datasets to be used and tested. The AIF Format [11] allows the annotation of argument data across all AM tasks, providing a platform for many different kinds of research.

Throughout the project two primary corpora have been considered: QT30 [12], a corpus consisting of 30 AIF annotated Question Time episodes, and a corpus of 9 AIF annotated Moral Maze episodes available on AIFdb.

¹<https://corpora.aifdb.org/>

2.1 Preprocessing

2.1.1 Argument Data

In order to use AIF data efficiently for ARI, it is useful to perform some preprocessing. This process produces a graph, where each node contains a locution, its related proposition, and the proposition’s AIF identifier. This identifier corresponds to the audio data, allowing it to be easily loaded when required. Each edge in this graph corresponds to a relation between the propositions, one of RA (inference), MA (rephrase) or CA (conflict).

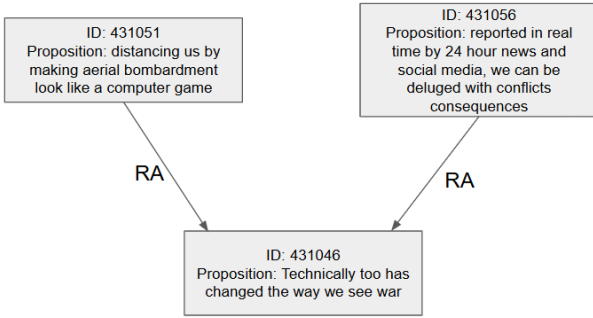


Figure 1: Example sub-graph.

Figure 1 shows an example sub-graph from the larger argument graph. Each node is truncated for brevity and only shows the node’s ID, and the proposition. This sub-graph is taken from the Moral Maze episode on the 75th Anniversary of D-Day.

An example of the JSON structure used to store the argument data is shown in Listing 1.

An array of these JSON objects can then be used to create the node pairs required for the training and evaluation of the model.

2.1.2 Audio Data

The audio data first had to be downsampled from 44.1kHz to the 16kHz which is best accepted by the Wav2Vec2 transformer [13] among many others. This can easily be achieved using FFmpeg². In the case of QT30, first audio had to be extracted from the video,

²<https://ffmpeg.org/>

Listing 1 Example JSON object corresponding to a Node.

```

1 {
2   "id": 433407,
3   "locution": "Matthew Taylor : she
   ↳ answered questions about
   ↳ norms and structures by
   ↳ talking about beliefs and
   ↳ campaigns and I think beliefs
   ↳ are different to norms and I
   ↳ think campaigns are different
   ↳ to social structures",
4   "proposition": "Nancy Sherman
   ↳ answered questions about
   ↳ norms and structures by
   ↳ talking about beliefs and
   ↳ campaigns and beliefs are
   ↳ different to norms and
   ↳ campaigns are different to
   ↳ social structures",
5   "relations": [
6     {
7       "type": "CA",
8       "to_node_id": 433393
9     },
10    {
11      "type": "RA",
12      "to_node_id": 433416
13    }
14  ],
15 }
```

and collapsed into a mono track before it could be downsampled, this was also easily achieved with FFmpeg.

Next, start and end times for each locution in the argument graph need to be found, to allow the audio to be split per-locution (and therefore per node). This can be achieved using PyTorch’s forced alignment api³.

Initially this was achieved by aligning each word in the transcript of the episode, producing start and end times for each word. A search

³<https://pytorch.org/audio/>

can then be performed through this data to find the required locution. While this technique initially produced promising results, it was not robust enough to allow for errors in the transcripts or the crosstalk common in debates.

To solve this problem, the PyTorch forced alignment api is able to take wildcard tokens as input, therefore, each locution can be searched for individually. To achieve this, the partial transcript used as input to the forced aligner took the following form: `* {locution} *`.

Using this system allows the forced aligner to work well through crosstalk (since each locution’s alignments are searched for independently of all others), and qualitatively seems to be more resilient to errors. Error resilience is helped since errors are less common in the locution texts as opposed to the transcripts. Using this system also allowed for confidence scores to be collected for analysis. In this section general analysis across all corpora is performed, with corpus specific analysis in the relevant section.

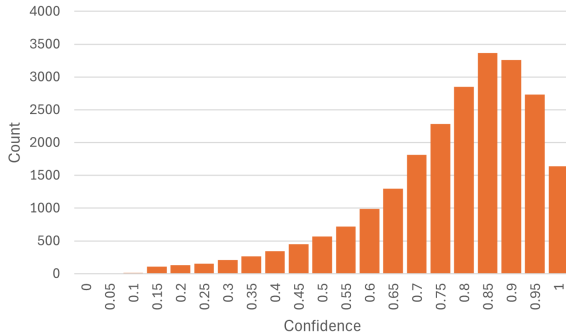


Figure 2: Confidence distribution across all corpora.

Figure 2 shows the distribution of confidence scores across both the QT30 corpus and all Moral Maze corpora. This distribution shows that the system can relatively confidently align the majority of locutions, with only approx. 8% of locutions with a confidence score less than 0.50.

In order to further analyse the performance

of this system, locutions were selected at random and qualitatively analysed. Throughout this process, all locutions appeared correct, however, it was very challenging to accurately determine the accuracy of the system on locutions with confidence scores < 0.2 . This shows that this method of aligning locutions with their corresponding audio is accurate.

2.1.3 Pair Creation

Finally, a set of node pairs and their relations can be generated in order to train a neural network. For related nodes this can be done trivially in that for each relation, the corresponding pair of nodes can be added to the set. When sampling unrelated nodes, however, things are more complex.

For this project, Short Context Sampling (SCS) is used as presented in [14]. Given the episodic structure of both QT30 and the Moral Maze corpora, a short context can be defined as the episode. This also allows the model to learn in a more realistic environment. Since the vast majority of node pairs have no relation, a number equal to that of Inference relations are generated.

2.2 QT30

The QT30 argument corpus [12] contains transcripts and argument annotations for 30 episodes of the BBC’s Question Time, a series of televised topical debates across the United Kingdom. All episodes aired in 2020 and 2021. The corpus is split into 30 subcorpora, each spanning a single episode. This allows analysis of each episode individually, or combined as a single corpus.

Table 1 shows the distribution of each type of relation across QT30. Inference and Rephrase relations make up a total of 91.5% of the dataset, with Conflict relations being significantly less common, only making up 8.5% of the dataset. It is obvious that this is an unbalanced dataset, which will have to be considered

Table 1: Disribution of propositional relations in QT30.

| Relation Type | Count | Proportion (%) |
|---------------|--------|----------------|
| Inference | 5,761 | 51.4% |
| Conflict | 947 | 8.5% |
| Rephrase | 4,496 | 40.1% |
| Total | 11,204 | 100% |

during training.

Table 2: Mean confidence scores (μ) and standard deviation of confidence scores (σ) across each QT30 subcorpus.

| Corpus Name | μ | σ |
|------------------------|-------------|-------------|
| 28May2020 | 0.76 | 0.16 |
| 4June2020 | 0.72 | 0.17 |
| 18June2020 | 0.76 | 0.16 |
| 30July2020 | 0.75 | 0.16 |
| 2September2020 | 0.78 | 0.15 |
| 22October2020 | 0.76 | 0.16 |
| 5November2020 | 0.77 | 0.17 |
| 19November2020 | 0.74 | 0.19 |
| 10December2020 | 0.77 | 0.16 |
| 14January2021 | 0.74 | 0.17 |
| 28January2021 | 0.70 | 0.17 |
| 18February2021 | 0.75 | 0.16 |
| 4March2021 | 0.76 | 0.16 |
| 18March2021 | 0.75 | 0.17 |
| 15April2021 | 0.75 | 0.15 |
| 29April2021 | 0.70 | 0.19 |
| 20May2021 | 0.76 | 0.17 |
| 27May2021 | 0.79 | 0.15 |
| 10June2021 | 0.75 | 0.17 |
| 24June2021 | 0.74 | 0.17 |
| 8July2021 | 0.72 | 0.17 |
| 22July2021 | 0.44 | 0.28 |
| 5August2021 | 0.76 | 0.17 |
| 19August2021 | 0.77 | 0.17 |
| 2September2021 | 0.78 | 0.16 |
| 16September2021 | 0.77 | 0.16 |
| 30September2021 | 0.24 | 0.08 |
| 14October2021 | 0.75 | 0.19 |
| 28October2021 | 0.75 | 0.17 |
| 11November2021 | 0.78 | 0.16 |
| QT30 | 0.75 | 0.17 |

Table 3: Disribution of propositional relations in QT30-MM.

| Relation Type | Count | Proportion (%) |
|---------------|--------|----------------|
| Inference | 5,740 | 51% |
| Conflict | 937 | 8.4% |
| Rephrase | 4,479 | 40.6% |
| Total | 11,156 | 100% |

Table 2 shows the mean and standard deviation of the confidence scores across each of the QT30 subcorpora. The lowest two mean scores are shown in bold. Manually analysing samples in these episodes indicates a high error rate in the alignment of locutions. Because of this high error rate, it was decided to exclude these episodes from the corpus used for training. The excluded episodes are: 22July2021 and 30September2021. The rest of the episodes from QT30 will form its multimodal subcorpus (QT30-MM).

Similarly to the complete QT30 corpus, Inference and Rephrase make up the vast majority of the QT30-MM dataset, with the proportion of Conflict relations decreasing to 8.4%.

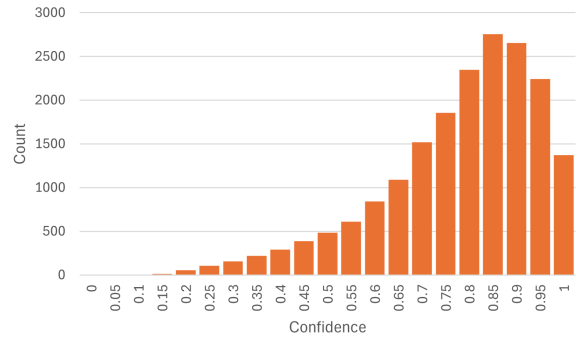


Figure 3: Confidence distribution across QT30-MM.

As can be seen in Figure 3 the distribution of confidence scores closely matches that shown in Figure 2. This still indicates that the audio alignments are calculated with high accuracy.

Table 4: Distribution of propositional relations after sampling non-related nodes.

| Relation Type | Count | Proportion (%) |
|---------------|--------|----------------|
| None | 5,470 | 34% |
| Inference | 5,470 | 34% |
| Conflict | 937 | 6% |
| Rephrase | 4,479 | 27% |
| Total | 16,896 | 100% |

2.3 Moral Maze

Similar to Question Time, the BBC’s Moral Maze is a series of radio broadcast debates, with each episode focusing on a certain topic. Seven different Moral Maze episodes have been AIF annotated and made available on AIFdb. It is therefore these seven episodes, released from 2012 to 2019, which this project considers. Each episode focuses on a very different domain which allows for a robust, cross-domain analysis of any models trained on another corpus (e.g. QT30). The Moral Maze corpus contains data from nine different episodes: Banking (B), Empire (E), Money (M), Problem (P), Syria (S), Green Belt (G), D-Day (D) and Hypocrisy (H). Each episode consists of a debate focusing on a different topic, and hence has a different distribution of classes.

Table 5 shows the distribution of propositional relations across the Moral Maze corpus, after non-related pairs have been sampled. Comparing the corpus to QT30, a significantly lower proportion of the corpus is made up of Rephrase relations. It is possible that the differing formats of the debates has an impact here.

In Table 6 the mean and standard deviation of audio alignment confidence scores are compared across subcorpora. Generally the results match what is expected and are similar to those in QT30, the system does achieve unusually low scores when considering the D-Day subcorpus, the reason for this is unclear, however, manually analysing both random and low-confidence samples indicates they are gen-

erally correct and so the subcorpus can be used for the cross-domain evaluation.

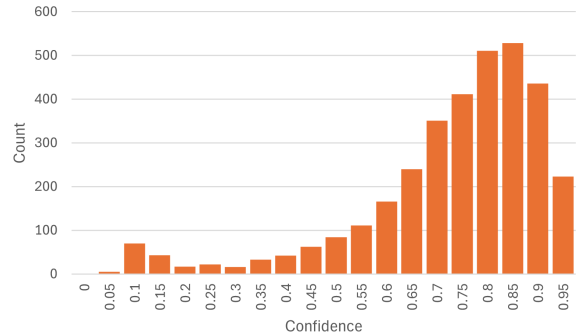


Figure 4: Confidence distribution across all Moral Maze subcorpora.

Figure 4 shows the distribution of audio alignment confidence scores across all Moral Maze episodes. This follows the expected pattern as shown in Figures 2 and 3. The secondary peak around a confidence of 0.10 is caused by the D-Day subcorpus.

3 Models

There are many different approaches when considering multimodal models. Generally they can be split into two categories: early fusion and late fusion. Early fusion models initially combine the features from each modality together, before being used as input to a single transformer model. Late fusion models use multiple transformers, one specialised in each modality. After being fed through each transformer, the hidden vectors are combined before being fed into the model’s head. Initially, only late fusion models have been considered in this research.

This late fusion model uses RoBERTa-base [8] as its text transformer. To process the audio data, the model uses the Wav2Vec2-base audio transformer [13], having been pre-trained on 960 hours of dialogue⁴. The base models are

⁴<https://huggingface.co/facebook/wav2vec2-base-960h>

Table 5: Distribution of propositional relations across the Moral Maze corpus.

| Relation Type | None | Inference | Conflict | Rephrase | Total |
|---------------|-------------|-------------|----------|----------|-------|
| B | 132 (45%) | 132 (45%) | 24 (8%) | 3 (1%) | 291 |
| E | 151 (41%) | 151 (41%) | 39 (11%) | 25 (7%) | 366 |
| M | 255 (43%) | 255 (43%) | 29 (5%) | 58 (10%) | 597 |
| P | 236 (42%) | 236 (42%) | 40 (7%) | 45 (8%) | 557 |
| S | 181 (42%) | 181 (42%) | 63 (14%) | 10 (2%) | 435 |
| G | 301 (41%) | 301 (41%) | 46 (6%) | 93 (13%) | 741 |
| D | 72 (40%) | 72 (40%) | 7 (4%) | 28 (16%) | 179 |
| H | 207 (43%) | 207 (43%) | 23 (5%) | 43 (9%) | 480 |
| Total | 1,535 (42%) | 1,535 (42%) | 271 (7%) | 305 (8%) | 3,646 |

Table 6: Mean confidence scores (μ) and Standard Deviation of confidence scores (σ) across Moral Maze subcorpora.

| Subcorpus | μ | σ |
|-----------|-------------|-------------|
| B | 0.79 | 0.14 |
| E | 0.76 | 0.15 |
| M | 0.78 | 0.14 |
| P | 0.80 | 0.15 |
| S | 0.74 | 0.16 |
| G | 0.80 | 0.14 |
| D | 0.50 | 0.34 |
| H | 0.74 | 0.16 |

used currently in order to increase the speed at which experiments can be conducted, however, they are used with a view to transitioning to their large variants towards the end of the project.

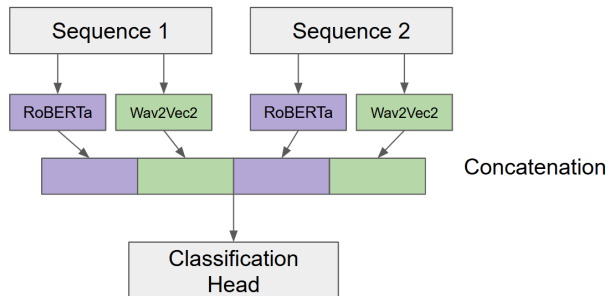


Figure 5: Concatenation model with late fusion.

Figure 5 provides a visualisation of the primary model used at this point in the project. Similar

models have been shown to achieve good multi-modal performance while also being relatively simple to implement [9], [10]. It is for these reasons that this model has been implemented using the Huggingface⁵ and PyTorch⁶ python modules.

Some preliminary experiments have also been conducted using text-only models in an attempt to replicate results produced in [2]. For this the RoBERTa-base model is used and each sentence is concatenated before tokenisation, delimited by RoBERTa’s special purpose token `</s>`.

4 Future Work

Throughout the rest of the project the following aims exist:

- A performance analysis of the models presented in Section 3 and other both early and late fusion strategies (e.g. swapping concatenation for a Hadamard product or cross-attention layer). Other pooling strategies could also be analysed.
- Implementation and analysis of many different models, including different transformer models (e.g. HuBERT [15]).
- Extension of MAMKit datasets to allow the simple implementation of QT30-MM

⁵<https://huggingface.co/>

⁶<https://pytorch.org/>

and the Moral Maze corpora created in this project by others.

It is likely that these aims will be modified as the project continues given its dynamic nature, as additional results are produced they will be interpreted and a way forward produced.

References

- [1] J. Lawrence and C. Reed, “Argument Mining: A Survey,” *Computational Linguistics*, vol. 45, no. 4, pp. 765–818, Jan. 2020, doi: 10.1162/coli_a_00364.
- [2] D. Gemechu, R. Ruiz-Dolz, and C. Reed, “ARIES: A General Benchmark for Argument Relation Identification,” in *Proceedings of the 11th Workshop on Argument Mining (ArgMining 2024)*, Aug. 2024, pp. 1–14. doi: 10.18653/v1/2024.argmining-1.1.
- [3] K. Budzynska *et al.*, “Towards Argument Mining from Dialogue,” in *Computational Models of Argument*, IOS Press, 2014, pp. 185–196. doi: 10.3233/978-1-61499-436-7-185.
- [4] K. Budzynska, M. Janier, C. Reed, P. Saint-Dizier, M. Stede, and O. Yakorska, “A Model for Processing Illocutionary Structures and Argumentation in Debates,” in *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, May 2014, pp. 917–924.
- [5] A. Vaswani *et al.*, “Attention Is All You Need,” p. 11, 2017.
- [6] R. Ruiz-Dolz, J. Alemany, S. M. H. Barberá, and A. García-Fornes, “Transformer-Based Models for Automatic Identification of Argument Relations: A Cross-Domain Evaluation,” *IEEE Intelligent Systems*, vol. 36, no. 6, pp. 62–70, Nov. 2021, doi: 10.1109/MIS.2021.3073993.
- [7] Y. Wu, Y. Zhou, B. Xu, W. Wang, and Y. Song, “KnowComp at DialAM-2024: Fine-tuning Pre-trained Language Models for Dialogical Argument Mining with Inference Anchoring Theory,” in *Proceedings of the 11th Workshop on Argument Mining (ArgMining 2024)*, Aug. 2024, pp. 103–109. doi: 10.18653/v1/2024.argmining-1.10.
- [8] Y. Liu *et al.*, “RoBERTa: A Robustly Optimized BERT Pretraining Approach.” arXiv, Jul. 2019. doi: 10.48550/arXiv.1907.11692.
- [9] E. Mancini, F. Ruggeri, A. Galassi, and P. Torroni, “Multimodal Argument Mining: A Case Study in Political Debates,” in *Proceedings of the 9th Workshop on Argument Mining*, Oct. 2022, pp. 158–170.
- [10] E. Mancini, F. Ruggeri, S. Colamonaco, A. Zecca, S. Marro, and P. Torroni, “MAMKit: A Comprehensive Multimodal Argument Mining Toolkit,” in *Proceedings of the 11th Workshop on Argument Mining (ArgMining 2024)*, Aug. 2024, pp. 69–82. doi: 10.18653/v1/2024.argmining-1.7.
- [11] C. Chesñevar *et al.*, “Towards an argument interchange format,” *The Knowledge Engineering Review*, vol. 21, no. 4, pp. 293–316, Dec. 2006, doi: 10.1017/S0269888906001044.
- [12] A. Hautli-Janisz, Z. Kikteva, W. Siskou, K. Gorska, R. Becker, and C. Reed, “QT30: A Corpus of Argument and Conflict in Broadcast Debate,” in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, Jun. 2022, pp. 3291–3300.
- [13] A. Baeviski, H. Zhou, A. Mohamed, and M. Auli, “Wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations.” arXiv, Oct. 2020. Accessed: Oct. 16, 2024. [Online]. Available: <https://arxiv.org/abs/2006.11477>

- [14] R. Ruiz-Dolz, D. Gemechu, Z. Kik-teva, and C. Reed, “Looking at the Unseen: Effective Sampling of Non-Related Propositions for Argument Mining,” in *Proceedings of the 31st International Conference on Computational Linguistics*, Jan. 2025, pp. 2131–2143.
- [15] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, “HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units.” arXiv, Jun. 2021. doi: 10.48550/arXiv.2106.07447.