

ArguNet: Exploiting Dialogue Acts and Context to Identify Argumentative Relations in Online Debates

Anonymous ACL submission

Abstract

Argumentative Relation Classification is the task of determining the relationship between two arguments within the context of an argumentative dialogue. Existing models in the literature rely on a combination of lexical features and pre-trained Large Language Models (LLMs) to tackle this task; while this approach is somewhat effective, it fails to take into account the importance of pragmatic features such as the illocutionary force of the argument or the structure of previous utterances in the discussion. In this work, we introduce ArguNet, a new model for Argumentative Relations Classification which relies on a combination of Dialogue Acts and Dialogue Context to obtain a more nuanced understanding of an argument’s stance. We show that our model achieves state-of-the-art results on the Kialo benchmark test set, and provide evidence of its robustness in an open-domain scenario.

1 Introduction

Argumentative Dialogues are discussions between two or more parties involving an opinionated topic, i.e. any topic which may divide the interlocutors into a number of conflicting opinions. These discussions are usually different from ordinary conversations, in that the speakers’ goal is usually to convince their interlocutor of their own point of view by defending their own stance and attacking their opponent’s arguments. Figure 1 shows an example of a debate from the Kialo online debate platform. A key aspect in the study of Argumentative Dialogues is identifying the relationship between an argument step in the discussion and preceding argument steps introduced by other speakers; this task is commonly referred to as *Argumentative Relation Classification* (Stab and Gurevych, 2014), or sometimes *Argument Polarity Prediction* (Cayrol and Lagasque-Schiex, 2005) when it only involves a binary classification between two possible relations.

In this work, we will use the term **Argumentative Relation Classification**, to avoid any confusion with similar tasks such as *Sentiment Analysis* or *Stance Classification*.

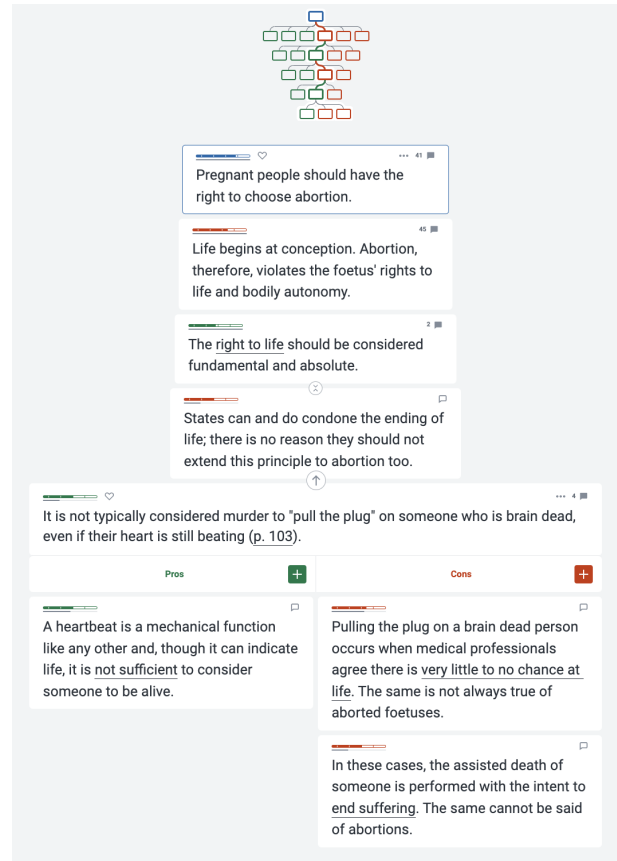


Figure 1: An example of a debate from the Kialo online debate platform. Green nodes agree with the original thesis (in blue), while red nodes disagree with it. Nodes are annotated with the *argumentative move* that they perform on their parent node in the graph (i.e. *Support* or *Attack*). Users annotate their own stance towards the thesis, as well as their argumentative move towards the node they are interacting with.

Existing works in the literature that aim at solving this task usually rely on either hand-crafted syntactic and lexical features (Stab and Gurevych,

2014; Lenz et al., 2020), pre-trained Large Language Models (LLMs) (Agarwal et al., 2022) or both (Cocarascu et al., 2020). While these models are becoming increasingly accurate, there are some shortcomings in their approach: they often ignore any non-lexical aspect of the dialogue, which hinders their capability to correctly understand the conversation. They have limited understanding of the surrounding context of the argument, and struggle to take long-term dependencies into account. Finally, they are often tested in a domain-specific scenario in which the system learns to predict relations between arguments that belong in the same dataset it was trained on; this makes it hard to correctly assess their capability to adapt to unseen conversations, which is crucial for practical applications such as the development of Automated Dialogue Agents.

In this work, we explore the hypothesis that contextual information and pragmatic features (such as Dialogue Act Tags) can be highly beneficial in increasing the accuracy of Argumentative Relation Classification models. We also aim at analysing how much existing models can generalise to entirely unseen topics of discussion, and how these features can help a model become less dependent on its training domain. There is evidence in the literature that Dialogue Act Tags may be used as a feature to improve a model’s understanding of the argumentative structure of a debate (Petukhova et al., 2016; Budzyska et al., 2014). There is also evidence that contextual information is highly beneficial for Argument Mining tasks and, more specifically, to increase the accuracy of Argumentative Relation Classification models (Agarwal et al., 2022).

We build on this existing evidence and introduce **ArguNet**, a novel neural architecture for Argumentative Relation Classification that relies on a combination of Dialogue Acts and a specialised encoding of the previous nodes in the debate. ArguNet uses ISO 24617-2 Dialogue Acts (DAs) annotated with the DASHNet architecture (Mezza et al., 2022) to enrich the input utterances with additional syntactic and pragmatic information. BERT (Devlin et al., 2018) is used to encode the enriched input utterances into dense sentence embeddings, with the addition of Utterance Manipulation Strategies from (Whang et al., 2021) to further increase the effectiveness of the contextual

embeddings from BERT. Our approach is trained and tested on data from the Kialo online debate platform, a high-quality, publicly-available source of conversations annotated with argumentative relations. We use the same Kialo scrape introduced by (Agarwal et al., 2022); however, instead of shuffling the arguments and dividing them in a training and test split, we split at the debate level, so that arguments from the same debate will not appear in different splits. This is done to test the hypothesis that existing models identify lexical information in the training debates and are able to use this information when tested on arguments from the same debates. We also sampled an additional, smaller collection of Kialo debates called *KialoAbortion* that involve discussions on reproductive rights, which we use to further test our hypothesis that Argumentative Relation classification is highly sensitive to the topic of the classified arguments.

In our experimental section, we provide evidence that the ArguNet architecture achieves state-of-the-art results on the Kialo dataset; we also provide evidence that our model outperforms existing models in the literature when tested on debates from the *KialoAbortion* test set, which shows how ArguNet can generalise to unseen domains better than existing architectures.

2 Related Work

The formal study of argumentative discussions is known in the literature as *Argumentation Theory* (van Eemeren et al., 1996) and it has been the subject of interest of various disciplines, including logic, rhetoric and philosophy. (Walton, 2009) divides argumentative study into four separate tasks: *identification*, which involves identifying Argumentative Dialogue Units (ADUs) in a dialogue and inserting them into a pre-determined *argumentation scheme*; *analysis*, which deals with identifying premises and conclusion of each argument; *evaluation*, which involves assessing an argument’s quality and persuasive power; and *invention*, which involves the creation of novel arguments for the debate. In this work we will focus on the task of *identification*, and propose a method to insert pre-constructed ADUs in an argumentation scheme.

The identification of a logical structure for

reasoning goes back to the seminal works by (Pollock, 1987) and (Nute, 1988), which introduced *Defeasible Logic*, a formalism in which *conclusions* are supported by *premises* that may no longer be justified when additional premises are introduced. (Dung, 1995) introduced an abstract theory of *Acceptability of Arguments* in which arguments are seen as a set of logical statements, and each argument can be *accepted* or *defeated* depending on whether it clashes with other arguments. (Prakken, 2010) elaborated on this theory and presented a framework for structured arguments in which arguments can be supported with premises that justify their validity, and other arguments can attack the speaker’s viewpoint by either attacking the argument directly, or one of its premises. (Cabrio and Villata, 2012) combine textual entailment and argumentation graph into a unified framework that aims at automatically detecting accepted and defeated arguments based on the entailment between them. (Lenz et al., 2020) adopted this scheme in their study on Argumentative Relation Classification on the Kialo corpus, and defined *Default Inference* and *Default Conflict* relations between arguments that support and attack each other respectively. The scheme was adopted by (Fabbri et al., 2021), who use Natural Language Inference models to directly compute Argumentative Relations. This approach, however, does not distinguish between the semantic problem of determining logical relations between argument steps and the pragmatic problem of determining dialogue moves in a sequence of contributions in a debate.

(Rosenfeld and Kraus, 2016) introduced a graph-like scheme for argumentative moves in a debate called the *Bipolar Argumentation Graph* (BAG), in which claims are represented as nodes in a weighted graph, and can be supported by other claims or *premises* that can either *Support* or *Attack* each other. Figure 2 shows the BAG scheme as illustrated in (Rach et al., 2019), one of the works that adopt it. As the Kialo dataset uses a graph-like structure that resembles a BAG, we will sometimes use their terminology in this work, particularly when referring to the argumentative moves between argument nodes.

Automatic annotation of argumentation schemes through Machine Learning algorithms has been studied extensively in recent years. One of

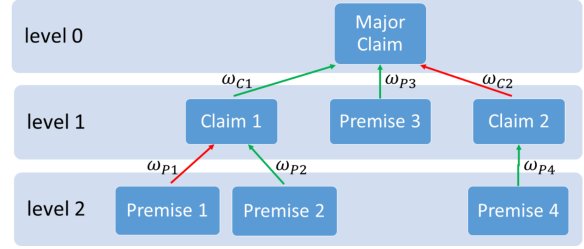


Figure 2: The Bipolar Argument Graph argumentation framework, as illustrated in (Aicher et al., 2021).

the earliest examples of a formal approach to Argumentative Relation Classification is (Cabrio and Villata, 2012), which proposes an approach based on Textual Entailment. (Naderi and Hirst, 2016) uses a combination of Skip-Thought Vectors and Cosine Similarity to predict argumentative relations in parliamentary debates; their work is one of the earliest that takes advantage of pre-trained word embeddings for this task. (Cocarascu and Toni, 2017) propose a neural architecture based on Long-Short Term Memory cells to annotate a multi-topic corpus which included debates on movies, technology and politics; they formulate the problem as a three-way classification problem between the classes *Attack*, *Support* and *Neither*. (Cocarascu et al., 2020) proposed a set of strong baselines for argumentative relation prediction in a dataset-independent setting, which included an attention-based model and an autoencoder. Their emphasis on dataset-independent classification is highly relevant to our work; however, they do not analyse the difference between in-domain and out-of-domain accuracy for their model and they do not provide details on how they split their data when separating training and test sets.

Recently, (Agarwal et al., 2022) proposed GraphNLI, a graph-based neural architecture that uses graph walking techniques to obtain contextual information, which is then encoded with RoBERTa embeddings (Liu et al., 2019). Their model was a source of inspiration for our work, as it shares our reliance on context encoding for Argumentative Relations Classification; however, their approach does not use pragmatic features like Dialogue Acts, and it also uses weighted averaging for embeddings rather than relying on a structured approach for context encoding, which we argue is less effective when trying to capture contextual information. Finally, some of their graph-walking

techniques rely on visiting neighbouring nodes or even future nodes in the discussion, which is not suitable for a real-life application like an interactive debater or a Spoken Dialogue System.

The idea of adopting Dialogue Acts as input features for Argument Mining systems has been investigated before in the literature. (Fouqueré and Quatrini, 2013) proposed a unified framework for argumentative analysis and inference which used Dialogue Acts as part of the argumentation scheme, and used it to annotate a discussion from (Prakken, 2008). (Budziyska et al., 2014) introduced Inference Anchoring Theory (IAT), a framework designed to model arguments via a combination of argumentative moves and the dialogue acts associated with them. Both of these works utilised Dialogue Act schemes that are difficult to adopt due to the scarcity of annotated data. (Petukhova et al., 2016) use ISO 24617-2 Dialogue Acts as part of a model designed to understand the argumentative behaviour of participants in a debate in order to predict its outcome. Their study provides some useful insights on how ISO 24617-2 Dialogue Acts can be used to model an argumentative discussion; however, their model is limited by the use of outdated Machine Learning methods for the task and was only tested on a limited number of debates. In our work, we also decided to adopt ISO 24617-2 Dialogue Acts due to their flexible, multi-dimensional and domain-independent taxonomy; we rely on the DASHNet model from (Mezza et al., 2022) which achieved state-of-the-art accuracy on various benchmark test sets. Their work also provides a detailed overview of previous works in Dialogue Act annotation and DA taxonomy design.

3 Methodology

3.1 Task Definition

Given an argument contribution A_j , which can be comprised of one or more sentences, and the list of nodes P_j connecting it to the thesis node T , which we will call the *Context* of the argument, we define **Argumentative Relation Classification** as the task of automatically identifying the weight of the edge $E_j = (A_j, A_{j-1})$, which represents the argumentative relation $R_{j,j-1}$ between A_j and its preceding node in the debate, A_{j-1} . We modeled this task as a Statistical Machine Learning model and designed a Neural Network architecture called ArguNet, which we describe in Section 3.3.

3.2 Data

For this study, we chose to work with data from the Kialo online debate platform¹. We have decided to use Kialo because it is a highly-curated platform with moderated arguments and a vote system for posts, which minimizes the amount of noise, ad hominem attacks and other irrelevant information in the arguments. Moreover, as the dataset is moderated, it is free of identifiable information about individuals or offensive content. Kialo debates are organised in a weighted graph-like structure: nodes in the graph represent individual, fully-formed arguments from a single participant in the debate and are called *Contributions*. Contributions are linked together with weighted edges, with the weights representing the *Argumentative Relation* between the two contributions linked by the edge. Every debate graph forms a tree-like structure, with the thesis being debated as the root node of the tree; dialogues have multiple participants, and the participants construct the tree structure collectively as they debate.

We use a scrape of Kialo introduced in (Agarwal et al., 2022), which we refer to as *KialoDataset*. This is a complete scrape of the website as of January 2020, and contains about 1,400 debates in total. We also collected our own scrape of the website, which we refer to as *KialoAbortion*, focusing on a specific topic; we chose to focus on *Reproductive Rights*, as this is a very popular and polarising debate topic at the time of writing. We collected 40 debates related to the topic via a combination of keyword extraction and manual filtering. Table 1 contains some quantitative information on our splits.

Dataset	# of Debates	# of contributions
KialoDataset (train)	1,051	231,945
KialoDataset (test)	278	53,699
KialoDataset (valid)	141	25,594
KialoAbortion (train)	27	8,970
KialoAbortion (test)	13	1,614

Table 1: Quantitative information about our data splits.

We divided our datasets into a 75% train split, a 15% test split and a 10% validation split. Experi-

¹<https://www.kialo.com/>

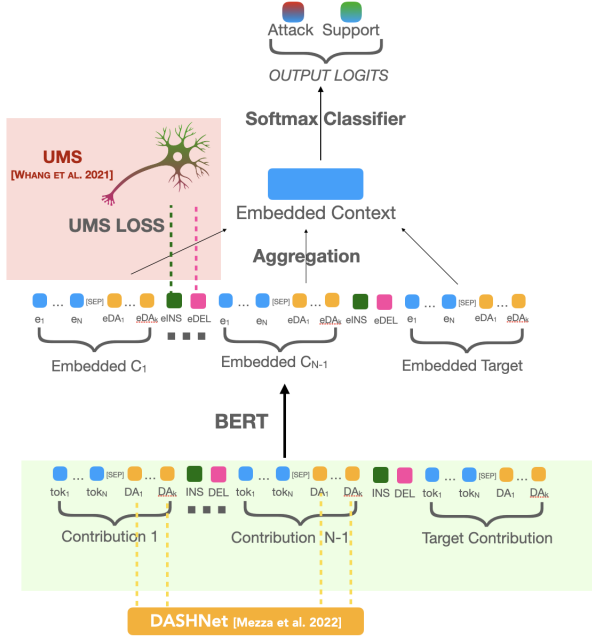


Figure 3: The ArguNet architecture.

ments in the literature sometimes split the debates without preserving their integrity; this *Single Contribution* splitting strategy produces splits which may contain argument contributions from the same debates. In contrast to that approach, we adopt a *Whole Debate* splitting strategy and split our data at the debate level, meaning that each split contains whole debates and contributions from the same debate do not appear in different splits.

3.3 Model

In this section we will outline the details of the **ArguNet** model, our neural architecture designed for *Argumentative Relation Classification*. Figure 3 provides an overview of the architecture. ArguNet is a transformer-based architecture with a few enhancements designed to increase its accuracy when dealing with argumentative data. It uses BERT (Devlin et al., 2018) to produce dense embeddings of each token in the input arguments. In order to increase the model’s ability to correctly understand each argument’s underlying meaning, we enhanced the input of ArguNet with ISO 24617-2 Dialogue Act Tags extracted with the DASHNet architecture (Mezza et al., 2022). We chose the DASHNet classifier because of its multidimensional and open-domain nature, which suits our use case very well; moreover, the model uses data from the Internet Argument Corpus (Abbott et al., 2016; Walker et al., 2012), which is similar in nature and scope to the Kialo data.

ArguNet also uses Utterance Manipulation Strategies (UMS) from (Whang et al., 2021) to obtain a better encoding of the context of the arguments to classify: special "[INS]" and "[DEL]" tokens are randomly inserted in the input and the corresponding utterance is either removed (in the case of "[DEL]") or erroneously inserted in the wrong spot (in the case of "[INS]"). The network has separate loss functions that control its learning of the correct UMS tags; this is combined with the classification loss from the final Softmax classifier, and the losses are averaged together to produce the final loss of the network. These strategies were originally introduced within the scope of Response Selection models to enhance their understanding of the history of the discussion; however, we chose to use them as the underlying principle of UMS should also apply to our task.

Our input is an argument contribution $A_N = T_1, \dots, T_N$, where T_i is the i -th token of the contribution, together with its context $C_{A_N} = A_{N-1} \dots A_{N-k}$, where k is the context window size of our model. We keep the window size at 5, following evidence in the literature that this is the optimal amount of context for an Argumentative Relation Classification model (Agarwal et al., 2022). We also only utilise argument contributions that directly preceded the target contribution in the debate, as opposed to alternative branches in the graph or future arguments in the discussion; this is done to make our model suitable for a real-life application in which future arguments may not be available for the analysis. Our data is pre-annotated with the DASHNet architecture to obtain a DA-enriched argument contribution $A_N = T_1, \dots, T_N, [SEP], DA_1, \dots, DA_M$. Each contribution in the context is also annotated with its DA tags. Note that DAs extracted from the DASHNet model are multi-dimensional, therefore there may be multiple tags for a single contribution. The input is then reshaped to utilize Utterance Manipulation Strategies (UMS), similarly to the UMS-ResSel model introduced in (Whang et al., 2021). We only utilize *Insertion* and *Deletion* strategies, as we found in our experiments that the *Search* strategies did not impact the accuracy of the resulting model when the other two strategies were present. For the insertion strategies, a target argument contribution in the context is randomly removed from its original position and placed at

the end of the context window. Special [INS] token are placed before each contribution in the context to encode whether the target contribution should be placed in that position. Target values for the [INS] tokens are 1 for the position in which the target argument contribution originally belonged, and 0 for all other tokens. For the deletion strategies, a random outlier contribution from a different context window is randomly placed in a random place in the context. Special [DEL] tokens are placed before each argument contribution in the context to encode whether that contribution is the outlier argument or not.

The input is concatenated with its UMS-enhanced context and they are all passed to the BERT model, which produces embeddings for each token in the input (including the DA tags and the UMS tokens). A binary cross-entropy loss function is applied to the UMS tokens to determine whether the network correctly guessed the positions of the argument contributions in the context. The tokens are then stacked together to produce a dense input representation which is then fed to a Softmax Classifier similar to the one used in Sentence-BERT (Reimers and Gurevych, 2019). The final loss of the model is the sum of the classification loss and the UMS losses.

4 Experiments and Results

In this section, we illustrate the results of our experimental study. We ran two sets of studies for these experiments: the first one was aimed at assessing ArguNet’s accuracy when trying to determine the Argumentative Relation between two argument contributions, and compare it to existing methods in the literature, while the second one aimed at measuring how much our model and existing models rely on domain-specific lexical information in order to produce their prediction. Both these sets of experiments involved the same models:

- **Majority Baseline:** this is just the frequency of the most likely argumentative move in the dataset. As this is a binary classification task with a reasonably balanced dataset, the majority baseline sits around 50% for both our test sets
- **ReCAP:** this is a model trained and tested on the Kialo corpus, originally introduced in (Lenz et al., 2020) as part of a larger study

on argument mining pipelines to transform textual arguments into argument graphs. The authors trained various machine learning models to predict the relation type between Kialo posts. We report results for their XG Boosting model, which is the most accurate based on our replication study, and used the code released by the authors under the Apache License, which was suitable for the purpose of this work.

- **BERT-1:** this is the result of fine-tuning the BERT model on the Kialo dataset, using a single argument contribution as the context window ($k = 1$). A softmax classifier is applied to the output BERT embeddings to produce the final output.
- **BERT-5:** this model is the same as BERT-1, but the context window length is increased to 5.
- **UMS:** this is a slight variation of the UMS-ResSel model introduced in (Whang et al., 2021) and originally intended for the task of Response Selection. We used the code released by the authors under the CC BY-SA license, which was suitable for the purpose of this work. We changed the final layer of the network to predict binary logits instead of ranks for response candidates.
- **GraphNLI:** this is the GraphNLI model as presented in (Agarwal et al., 2022). We used the code released by the authors under the MIT license, which was suitable for the purpose of this work. We use the weighted sum average method for aggregation, as it is the one that achieves the highest accuracy, and use the *Weighted root-seeking path* with a context length of 5. As described in Section 3.2, we altered the training and test splits of the Kialo dataset to keep debates intact, rather than randomly shuffling the argument contributions and splitting them; because of this, while we were able to replicate the authors’ results with their settings, our results when evaluating this model are different from the ones they reported.
- **GraphNLI-DA:** this is the same as GraphNLI, but we altered the input to also use Dialogue Act (DA) tags. We used the "<s/>" token to separate the input argument contributions

from their DA tags, and added the list of DASHNet DA tags as extra special tokens to the RoBERTa tokenizer that the authors use.

- **ArguNet:** our model as described in section 3.3.

In our initial plans we were aiming at combining the GraphNLI architecture with UMS tags; however, we decided against it as the two approaches were not compatible from an architectural or theoretical perspective. We also did not implement alternative versions of UMS, GraphNLI and ArguNet with context length equal to 1: as these models are built around using contextual information in their architecture, removing such information would have made these architecture meaningless without a history of previous utterances.

4.1 Implementation Details

We trained our models on Google Colab, using an NVIDIA A100 GPU with the "High RAM" setting. Training of our models took a total of roughly 400 GPU Hours, which includes all the re-trainings we had to do for our various experiments. We trained the UMS and ArguNet models for 20 epochs, but implemented early stopping with a patience of 3 (most models finished training between epochs 8 and 12). We use a Dropout rate of 0.8 for the final classification layer, a learning rate of 3e-05 and AdamW optimiser with epsilon value of 1e-8. We used BERT with 12 hidden layers, and an embedding dimension of 768, with a Dropout rate for its attention layer of 0.1. We validated all of these hyperparameters using the validation set of the *KialoDataset*. The original UMS-ResSel paper experimented with different weights for the UMS losses and the classification loss; we tried altering these values, but found that the best performing model was the one combining all losses with a ratio of 1.0. For the GraphNLI model, we maintained the original settings as detailed in the original paper (Agarwal et al., 2022), including hyper-parameters and number of epochs.

4.2 Argumentative Relation Classification

We trained various models from the literature on the combined train splits of the *KialoDataset* and *KialoAbortion* datasets, and measured their results to the ones obtained by the ArguNet model. We use accuracy as a metric and test on both the *KialoDataset* and *KialoAbortion* test sets separately. All the models were trained and

tested on the same data, and were trained with the *Whole Debate* splitting strategy (i.e. contributions from the same debate are kept in the same split). Because of this reason, some of the results we obtained are slightly different from the ones reported by the original authors of the respective papers. Table 3 shows the results of our study.

The results confirm our hypothesis that contextual information is highly beneficial for Argumentative Relation Classification: the BERT-5 model shows a significant improvement when compared to BERT-1, especially when tested on the *KialoAbortion* dataset, where it shows a 4.1% increase in accuracy. The results also show that an unstructured encoding of the context is less effective than a specialised encoding, as the model based on Utterance Manipulation Strategies (UMS) outperforms the BERT-5 model on both the *KialoDataset* and *KialoAbortion* corpora, with a 1.5% and 0.5% accuracy increase respectively. Moreover, the Dialogue Act feature appears to be highly beneficial to the classification for both the GraphNLI-DA and ArguNet models; this is particularly evident in the GraphNLI-DA model, which exhibits a 2.2% increase in accuracy on the *KialoAbortion* test set, and a 1.7% increase on the *KialoDataset* corpus when compared to the base GraphNLI model. This follows our hypothesis that Dialogue Act Tags provide an input signal that correlates with Argumentative Relation types. The DASHNet model uses data from the Internet Argument Corpus V2 (IAC) (Abbott et al., 2016; Walker et al., 2012) ; as this corpus contains argumentative discussions that are similar in scope and style to those found in Kialo, this may also have helped the classification further.

4.3 In-domain vs Out-of-domain accuracy

One of the main hypotheses that led to the design of the ArguNet architecture is that existing models in the literature largely rely on lexical information from their training corpora, which makes them less accurate when annotating debates on entirely unseen topics. In order to test this hypothesis, we compared the results of our implemented models when trained with and without the *KialoAbortion* training data. We used accuracy on the *KialoAbortion* benchmark test set as a metric. Table 2 shows the results of this study.

Results indicate that the ArguNet architecture

Model	Accuracy (KialoDataset)	Accuracy (KialoAbortion)
Majority Baseline	54.7%	54.5%
ReCAP (Lenz et al., 2020)	66.8 %	64.1%
BERT-1 (Devlin et al., 2018)	79.7%	74.4%
BERT-5	80.2%	78.5%
GraphNLI (Agarwal et al., 2022)	79.9%	78.9%
UMS (Whang et al., 2021)	80.7%	80.0%
GraphNLI-DA	81.6%	81.1%
ArguNet (our model)	82.1%	81.6%

Table 2: Argumentative Relation Classification results for the ArguNet model, compared with other models in the literature. We replicated all models for this work, and managed to replicate the original authors’ results.

outperforms existing approaches in the literature on both the in-domain and out-of-domain settings, while maintaining a relatively low difference in accuracy when trained with and without in-domain data. Moreover, it is noteworthy that models that utilise contextual information and other non-lexical features seem to be less prone to accuracy loss when trained without in-domain data; for example, the BERT-5 model has a 1.7% accuracy loss when trained without in-domain data, while the BERT-1 model has a significant 5.3% accuracy loss. More sophisticated models like GraphNLI, UMS or ArguNet which use contextual and Dialogue Act features have even lower differences in accuracy. This appears to validate our hypothesis that models that rely solely or mainly on lexical features are more prone to committing annotation errors when compared to models that adopt a more sophisticated encoding of the input.

5 Conclusions

In this work, we introduced a neural architecture called ArguNet which is optimised for the analysis of Argumentative Relations between argument contributions in online debates. We showed how it achieves state-of-the-art results when tested on the Kialo dataset of online debates, and provided evidence that its defining features, namely the use of Dialogue Acts and well-structured encoding of the

Model	OOD training	In-domain training	difference (%)
ReCAP (Lenz et al., 2020)	62.3 %	64.1%	1.8%
BERT-1 (Devlin et al., 2018)	72.3%	74.4%	2.1%
BERT-5	77.3%	78.5%	1.2%
GraphNLI (Agarwal et al., 2022)	78.8%	79.9%	1.1%
UMS (Whang et al., 2021)	79.4%	80.0%	0.6%
GraphNLI-DA	80.2%	81.1%	0.9%
ArguNet (our model)	80.9%	81.6%	0.7%

Table 3: Difference in accuracy between our implemented models when trained with and without in-domain data. All models were tested on the KialoAbortion test set.

context of the conversation, are highly beneficial for the task at hand. Finally, we showed how its architecture is more robust to out-of-domain classification when compared to existing approaches in the literature, and provided a comparison between in-domain and out-of-domain performance for all of our baselines.

5.1 Limitations and Future Work

The ArguNet architecture currently uses Dialogue Acts as an input feature that is annotated offline with the DASHNet Dialogue Act Tagger. We plan to refine its architecture to incorporate the DA tagging as part of its internal annotation, so that information about the argumentative relations can backtrack through the weights of the DA tagger and further increase its accuracy. The model is currently unable to determine the position of an input contribution in the graph, as it only predicts the Argumentative Relations between contributions. Argumentative Relations are a useful input feature that could be used for a variety of different tasks; we plan to experiment with how they can be used as part of a Dialogue Agent that can converse about opinionated topics and use the ArguNet architecture to understand a user’s stance on a topic and provide meaningful and relevant counter-arguments.

References

- Rob Abbott, Brian Ecker, Pranav Anand, and Marilyn Walker. 2016. Internet Argument Corpus 2.0: An SQL Schema for Dialogic Social Media and the Corpora to Go with It. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 4445–4452.
- Vibhor Agarwal, Sagar Joglekar, Anthony P Young, and Nishanth Sastry. 2022. GraphNLI: A Graph-based Natural Language Inference Model for Polarity Prediction in Online Debates. In *Proceedings of the ACM Web Conference 2022*, pages 2729–2737.
- Annalena Aicher, Niklas Rach, Wolfgang Minker, and Stefan Ultes. 2021. Opinion Building Based on the Argumentative Dialogue System BEA. In *Proceedings of the 10th International Workshop on Spoken Dialogue Systems*, pages 307–318.
- Kasia Budsziyska, Mathilde Janier, Chris Reed, Patrick Saint-Dizier, Manfred Stede, and Olena Yakorska. 2014. A Model for Processing Illocutionary Structures and Argumentation in Debates. In *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC 2014)*, pages 917–924.
- Elena Cabrio and Serena Villata. 2012. Combining Textual Entailment and Argumentation Theory for supporting Online Debates Interactions. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 208–212.
- Claudette Cayrol and Marie-Christine Lagasquie-Schiex. 2005. On the Acceptability of Arguments in Bipolar Argumentation Frameworks. In *Proceedings of the European Conference on Symbolic and Quantitative Approaches to Reasoning and Uncertainty*, pages 378–389.
- Oana Cocarascu, Elena Cabrio, Serena Villata, and Francesca Toni. 2020. Dataset Independent Baselines for Relation Prediction in Argument Mining. In *Proceedings of the 8th International Conference on Computational Models of Argument*, pages 45–52.
- Oana Cocarascu and Francesca Toni. 2017. Identifying Attack and Support Argumentative Relations using Deep Learning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1374–1379.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*.
- Phan Minh Dung. 1995. On the Acceptability of Arguments and its Fundamental Role in Nonmonotonic Reasoning, Logic Programming and N-Person Games. *Artificial Intelligence*, 77:321–357.
- Alexander Richard Fabbri, Faiaz Rahman, Imad Rizvi, Borui Wang, Haoran Li, Yashar Mehdad, and Dragomir Radev. 2021. ConvoSumm: Conversation Summarization Benchmark and Improved Abstractive Summarization with Argument Mining. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6866–6880.
- Christophe Fouqueré and Myriam Quatrini. 2013. Argumentation and Inference: A Unified Approach. *Baltic International Yearbook of Cognition, Logic and Communication*, 8(1):4.
- Mirko Lenz, Premtim Sahitaj, Sean Kallenberg, Christopher Coors, Lorik Dumani, Ralf Schenkel, and Ralph Bergmann. 2020. Towards an Argument Mining Pipeline Transforming Texts to Argument Graphs. In *Computational Models of Argument*, pages 263–270. IOS Press.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692*.
- Stefano Mezza, Wayne Wobcke, and Alan Blair. 2022. A Multi-Dimensional, Cross-Domain and Hierarchy-Aware Neural Architecture for ISO-Standard Dialogue Act Tagging. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 542–552.
- Nona Naderi and Graeme Hirst. 2016. Argumentation mining in parliamentary discourse. In *Principles and Practice of Multi-Agent Systems*, pages 16–25, Cham. Springer.
- Donald Nute. 1988. Defeasible reasoning and decision support systems. *Decision support systems*, 4(1):97–110.
- Volha Petukhova, Andrei Malchanau, and Harry Bunt. 2016. Modelling argumentative behaviour in parliamentary debates: Data collection, analysis and test case. In *Principles and Practice of Multi-Agent Systems*, pages 26–46, Cham. Springer.
- John L Pollock. 1987. Defeasible reasoning. *Cognitive science*, 11(4):481–518.
- Henry Prakken. 2008. A Formal Model of Adjudication Dialogues. *Artificial Intelligence and Law*, 16:305–328.
- Henry Prakken. 2010. An Abstract Framework for Argumentation with Structured Arguments. *Argument & Computation*, 1(2):93–124.
- Niklas Rach, Klaus Weber, Annalena Aicher, Florian Lingenfelder, Elisabeth André, and Wolfgang Minker. 2019. Emotion Recognition based Preference Modelling in Argumentative Dialogue Systems. In *2019*

774 *IEEE International Conference on Pervasive Comput-*
775 *ing and Communications Workshops (PerCom Work-*
776 *shops)*, pages 838–843.

777 Nils Reimers and Iryna Gurevych. 2019. Sentence-
778 BERT: Sentence Embeddings using Siamese BERT
779 Networks. *arXiv preprint arXiv:1908.10084*.

780 Ariel Rosenfeld and Sarit Kraus. 2016. Strategical Argu-
781 mentative Agent for Human Persuasion. In *Proceed-*
782 *ings of the 22nd European Conference on Artificial*
783 *Intelligence*, pages 320–328.

784 Christian Stab and Iryna Gurevych. 2014. Identifying
785 Argumentative Discourse Structures in Persuasive
786 Essays. In *Proceedings of the 2014 Conference on*
787 *Empirical Methods in Natural Language Processing*
788 *(EMNLP)*, pages 46–56.

789 Frans H van Eemeren, Rob Grootendorst, A Fran-
790 cisca Snoeck Henkemans, J Anthony Blair, Ralph H
791 Johnson, Erik CW Krabbe, Christian Plantin, Dou-
792 glas N Walton, Charles A Willard, John Woods, et al.
793 1996. *Fundamentals of Argumentation Theory: A*
794 *Handbook of Historical Backgrounds and Contempo-*
795 *rary Developments*.

796 Marilyn Walker, Jean E Fox Tree, Pranav Anand, Rob
797 Abbott, and Joseph King. 2012. A Corpus for Re-
798 search on Deliberation and Debate. In *Proceedings*
799 *of the Eighth International Conference on Language*
800 *Resources and Evaluation (LREC’12)*, pages 812–
801 817.

802 Douglas Walton. 2009. *Argumentation Theory: A Very*
803 *Short Introduction*, pages 1–22. Springer, Boston,
804 MA.

805 Taesun Whang, Dongyub Lee, Dongsuk Oh, Chanhee
806 Lee, Kijong Han, Dong-hun Lee, and Saebyeok Lee.
807 2021. Do Response Selection Models Really Know
808 What’s Next? Utterance Manipulation Strategies for
809 Multi-Turn Response Selection. In *Proceedings of*
810 *the AAAI Conference on Artificial Intelligence*, 35,
811 pages 14041–14049.