

Alma Mater Studiorum Università di Bologna
Archivio istituzionale della ricerca

Multi-Task Attentive Residual Networks for Argument Mining

This is the submitted version (pre peer-review, preprint) of the following publication:

Published Version:

Andrea Galassi, Marco Lippi, Paolo Torroni (2023). Multi-Task Attentive Residual Networks for Argument Mining. IEEE/ACM TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING, 31, 1877-1892 [10.1109/TASLP.2023.3275040].

Availability:

This version is available at: <https://hdl.handle.net/11585/816442> since: 2023-06-01

Published:

DOI: <http://doi.org/10.1109/TASLP.2023.3275040>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

Multi-Task Attentive Residual Networks for Argument Mining

Andrea Galassi, Marco Lippi, and Paolo Torroni

Abstract—We explore the use of residual networks and neural attention for argument mining and in particular link prediction. The method we propose makes no assumptions on document or argument structure. We propose a residual architecture that exploits attention, multi-task learning, and makes use of ensemble. We evaluate it on a challenging data set consisting of user-generated comments, as well as on two other datasets consisting of scientific publications. On the user-generated content dataset, our model outperforms state-of-the-art methods that rely on domain knowledge. On the scientific literature datasets it achieves results comparable to those yielded by BERT-based approaches but with a much smaller model size.

Index Terms—Argument mining, residual networks, neural attention, multi-task learning, ensemble learning, natural language processing.

I. INTRODUCTION

ARGUMENT mining (AM) is an emerging research area in natural language processing (NLP) which aims to extract arguments from text collections [1]. AM consists of several different tasks, that include argument detection, stance classification, topic-based argumentative content retrieval, and many others [2]. In this work we focus on the challenging problem of assembling the structure of the argumentation graph of a given input document. Such problem comprises the detection of both argument components, and relations (or links) amongst them, and is thus one of the most difficult steps for AM systems.

While there is no unique definition of an argument, one of the most popular ones was proposed by Douglas Walton [3], who defined an argument as the collection of three parts: (i) a *claim*, or assertion, about a given topic; (ii) a set of *premises* supporting the claim; (iii) the inference between the premises and the claim. Relations between arguments, or argument components, typically consist of either *support* or *attack* links.

AM approaches are very often tailored to specific corpora or genres [4]–[6], with solutions that are seldom general enough to be directly applied to different data sets without the need of any adaptation. It is very often the case that AM systems build

upon sets of handcrafted features which encode information about the underlying argument model, the genre, or the topic of interest. These approaches typically make some assumptions on the argumentative structure of the given input document, thus constraining the resulting argument graph.

We propose a general-purpose neural architecture that is domain-agnostic, and that does not rely on specific genre- or topic-dependent features. The model exploits neural attention and multi-task learning, jointly addressing the problems of identifying the category of argument components, and predicting the relations among them. Experimental results conducted on a variety of different corpora show that the model is robust and achieves good performance across the considered data sets. They also suggest that, when background information about the structure of annotations in a corpus is given, ad-hoc approaches may yield better performance.

Our main contributions are:

- A novel approach to AM, which extends our previous work [7] by introducing an attention module and ensemble learning. Such a model performs multiple AM tasks at the same time, and does not rely on ad-hoc features or rich contextual information, but only on GloVe embeddings and on a widely applicable notion of distance.
- An analytical evaluation of the contribution of each added module through an ablation study and a validation of our model on a challenging corpus, indicating that our proposed model improves state-of-the-art results in all the tasks we address.
- A set of experiments designed to assess generality, whereby we test our approach on three additional corpora that vary in domain, style of writing, formatting, length, and annotation model.

With respect to our previous work [7], this paper extends the neural architecture with attention and ensemble learning, and presents a more thorough and extensive experimental evaluation, offering comparisons with state-of-the-art systems across four different argument mining corpora. All the code used in our experiments is publicly available.¹

The paper is organized as follows. We present related work in Section II. Section III introduces our architectures. Section IV describes the data used for evaluation. Section V describes our experimental setting, whereas results are presented and discussed in Section VI. Section VII concludes.

¹<https://github.com/AGalassi/StructurePrediction18>.

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

This work was supported by the European Union’s Justice programme, project “Analytics for DEcision of LEGal Cases”, under Grant 101007420. (Corresponding author: Andrea Galassi.)

A. Galassi and P. Torroni are with the Department of Computer Science and Engineering (DISI), University of Bologna, Bologna 40126, Italy (e-mail: a.galassi@unibo.it; paolo.torroni@unibo.it).

M. Lippi is with the Department of Sciences and Methods for Engineering (DISMI), University of Modena and Reggio Emilia, Modena 42122, Italy (e-mail: marco.lippi@unimore.it).

II. RELATED WORK

The adoption of deep learning approaches in AM is relatively recent, compared to other areas of NLP. That is probably a consequence of a lack of large AM corpora, considering the complexity and peculiarities of the tasks at hand. Indeed, the annotation of large corpora for AM system evaluation and training proved to be challenging, as demonstrated by relatively low IAA indicators and several unsatisfactory attempts at crowdsourcing annotations. That is especially true for some genres like user-generated content [8]. Reasons for that are the nature of the task, which is intellectually demanding, and the lack of a unified argument model, as “arguments” may take very different shapes in different genres, also leading to a trade-off between the expressiveness of the argument model and the complexity of the annotation process and availability of relevant data points, often resolved in favor of simple argument models [1]. Earlier research mainly focused on the definition of features for specific genres or even for specific corpora. The differences between corpora, both regarding the domain and the theoretical framework followed during the annotation process, force researchers to test a model on the same corpora on which it was trained, and to the best of our knowledge, transfer learning approaches have not seen wide experimentation. These two elements lead to the common practice to define a method or a model and validate it only on a single corpus or on a few corpora [1].

A. Multi-task Learning and Joint Learning for AM

Since AM includes many subtasks that are strongly inter-related, a recent trend of this research field is to address many of them at the same time using multi-task or joint learning techniques. The aim of such approaches is to transfer knowledge from the auxiliary tasks to the main one, or to obtain coherent results on multiple tasks performed at once.

Stab and Gurevych [4] jointly address component classification and link prediction on persuasive essays, using Integer Linear Programming and a rich set of specific features, such as lexical, structural, and contextual information. Various neural architectures are tested in [9], including the deep biLSTM multi-task learning (MTL) setting of [10], using sub-tasks as auxiliary tasks. They conclude that neural networks can outperform feature-based techniques in argument mining tasks. Schulz et al. [11] investigate MTL settings addressing component detection on five datasets as five different tasks. Their architecture is composed of a CRF layer on top of a biLSTM, whose recurrent layers are shared across the tasks. They obtain positive results, and the MTL setting shows to be beneficial especially for small datasets, even if the auxiliary AM tasks involve different domains and even different component classes. Lauscher et al. [12] analyze an MTL setting where rhetorical classification tasks are performed along with component detection. They use a hierarchical attention-based model so to perform both word-level and sentence-level tasks with the same neural architecture. The results show improvements in the rhetorical tasks, but not in AM.

In [13], a structured learning framework based on factor graphs is used to jointly classify all the propositions in a

document and determine which ones are linked together. The models heavily rely on a priori knowledge, encoded as factors and constraints, designed so to enforce adherence to the desired argumentation structure, according to the argument model and domain characteristics. The authors discuss experiments with six different models, which differ by complexity and by how they model the factors, using RNNs and SVMs. Their best result is obtained by using the same set of features used in [4], resulting in a total feature size of around 7,000 for propositions and 2,100 for links.

B. Neural Attention for AM

Neural attention is a mechanism widely used in NLP to improve performance and interpretability of neural networks, and it is the core of many NLP architectures like RNNsearch [14], Pointer Networks [15], and Transformer [16]. Given an input sequence, and possibly a query element, attention consists in the computation of a set of weights that represent the importance of each element of the sequence, which can be further used to create a compact representation of such an input. There are many different ways to compute such weights. A taxonomy of attention models is proposed in our survey [17].

Among the AM systems that use neural attention, the one used in [18] integrate hierarchical attention and biGRU for the analysis of the quality of the argument, the one in [19] use attention to integrate sentiment lexicon, while in other works [20]–[22] attention modules are stacked on top of recurrent layers. The use of Pointer Networks for AM has also been investigated [23]. Transformer-based approaches in AM use language representation models such as BERT [24] and ELMO [25] to create contextualized word embeddings. Specifically, Reimers et al. [26] address component classification and argument clustering, a related task whose aim is to identify similar arguments. Similarly, Lugini and Litman [27] use BERT embeddings alongside other contextual information to perform component classification, and Wang et al. [28] use them to train a different model for each type of component. Trautmann et al. [29] use pre-trained BERT models to perform word-level classification of the stance of components regarding a given topic, while Poudyal et al. [30] use RoBERTa [31], an improved version of the original BERT.

Mayer et al. [32] present and conduct extensive experimentation on the AbstRCT corpus, addressing four AM subtasks with a pipeline scheme. They analyze the impact of various BERT models, which are pre-trained on other corpora and then fine-tuned on the corpus at hand. Segmentation and component classification are performed as sequence tagging with BIO scheme. Link prediction and relation classification follow, taking into account all the pairs of components obtained in the first step and classifying their relations as attack, support, or non-existing. Their architecture is based on bi-directional transformers followed by a softmax layer and various encoders. Their approach is completely distance-independent, but since they compare every possible pair of components, the size of the dataset grows quadratically with the number of components in the document, which makes it hardly scalable to large documents. Another approach,

consisting of predicting at most one related component for each component, and then classifying their relation, has been tested but yields worse results. The architectures that yield the best results are BioBERT [33], which is pre-trained on a large-scale biomedical corpus, SciBERT [34], which is pre-trained on scientific articles of various nature, and RoBERTa.

III. MODEL

The design of our model is inspired by the great success that residual networks [35] have obtained across many different tasks related to NLP [16], [36], [37]. The core idea behind residual networks is to create shortcuts that link neurons belonging to distant layers, whereas standard feed-forward networks typically link neurons belonging to subsequent layers only. This kind of architecture usually results in a more efficient training phase, allowing to train networks with considerably more layers, reducing the overall computational footprint.

The architecture we propose makes use of the dense residual network model, along with a Long Short-Term Memory (LSTM) network [38], and an attention module [17]. The network is trained to jointly perform three argument mining sub-tasks: argument component classification, link prediction, and relation classification.

More specifically, our approach operates on sentence pairs, does not rely on document-level global optimization, and does not enforce model constraints induced, for example, by domain- or genre-specific background knowledge. This makes our approach amenable to a possible integration within more complex and sophisticated systems.

We performed model selection and hyper-parameter tuning on a single corpus (CDCP, see Section IV) and we collected results on validation data in order to tune the whole architecture. There are two reasons for this choice: one the one hand, we aim to show the robustness of the approach across different corpora, while on the other we believe it is important to limit the footprint of these experiments – an issue that is receiving a growing attention in the community [39].

A. Model description

In order to achieve a general method which may be applicable in any domain, our method does not rely on a specific argument model, but rather it reasons in terms of abstract entities, such as argumentative components and links among them. We instantiate such abstract entities into concrete categories given by annotations, such as claims and premises, supports and attacks, as soon as we apply the method to a specific corpus whose annotations follow a concrete argument model.

The detection of argumentative content in text is one typical stage of AM systems [1]. Other works only focus on AM tasks that assume that argumentative components and their boundaries are already identified in the data. Such is the case with Niculae et al. [13], whose CDCP dataset only consists of argumentative elements, and with others [20], [27], [40] who simply ignore the non-argumentative elements of the input text. Accordingly, we define a *document* D as a sequence

of *argumentative components* and disregard the rest of the input text. An argumentative component in turn is a sequence of *tokens*, i.e., words and punctuation marks, representing an argument, or part thereof. The labeling of components is induced by the chosen argument model. Such a labeling associates each component with the corresponding category P of the argument component it contains. For this reason, we will use the terms component, sentence, and proposition as equivalent, and implying them as being argumentative by assumption.

Given two argumentative components a and b belonging to the same document, we represent a directed relation from the former (*source*) to the latter (*target*) as $a \rightarrow b$. Reflexive relations ($a \rightarrow a$) are not allowed.² Any pair of components is characterized by four labels: the types of the two components (P_a and P_b), the Boolean *link label* $L_{a \rightarrow b}$, and *relation (type) label* ($R_{a \rightarrow b}$). The link label indicates the presence of a link, and is therefore *true* if there exists a directed link from a to b , and *false* otherwise. The relation label instead contains information on the nature of the link connecting a and b . It represents the relationship between the two components, according to the links that connect a to b or b to a . Its domain is composed, according to the underlying argument model, not only by all the possible link types, but also by their opposite types (e.g., *attack* and *attackedBy*), as well as by a special category, *None*, meaning no link in either direction. One reason to introduce opposite relation types is to mitigate the unbalance caused by limited amount of instances each relation type typically has, if compared with the number of instances belonging to the *None* class. Likewise, we speculate that the introduction of additional labels may contribute positively to the optimization process. We shall remark that opposite relation labels are exploited during training, but they are discarded in the test phase, where they are simply substituted with the *None* label, consistently with previous work.

We use a multi-objective learning setting where multiple tasks are performed jointly for each possible input pair of components (a, b) belonging to the same document D . Our main focus is the identification of the link label $L_{a \rightarrow b}$ for each possible input pair of propositions (a, b) belonging to the same document D . Our first objective is thus a *link prediction* task, which can be considered as a sub-task of argument structure prediction. A second objective is the *classification* of the two components,³ and our final objective is the classification of the relationship between such components, i.e., the prediction of labels P_a , P_b , $R_{a \rightarrow b}$. A common issue in the classification of pairs of document components is the fact that pairs grow quadratically with the number of components, causing a large imbalance against the negative class [30], [32]. One way of dealing with that issue is to limit the possible pairs by setting a maximum distance, thus obtaining a number of pairs proportional to the number of components. Such a distance is a hyper-parameter, and as such it may be empirically determined [30].

²We will partially consider reflexive relations for the UKP dataset for a specific reason explained in Section V.

³Since we examine only argumentative propositions, we do not consider the non-argumentative class for component classification.

B. Embeddings and features

Faithful to the main purpose of this work, of evaluating the effectiveness of deep residual networks and attention for AM without resorting to domain- or genre-specific information, our system relies on a minimal set of widely applicable features.

Words are encoded using pre-trained GloVe embeddings [41] of size 300. Input sequences are zero-padded to the length of the longest sequence in the datasets (henceforth T). Out-of-vocabulary terms are handled by creating random embeddings.

In our previous work, we empirically assessed how the distance between two components may be a relevant feature for argument mining in the CDCP corpus [7]. The same observation has been recently made also with reference to other corpora [9], [30], [32]. Similarly to what has been done in [9], we define the number of argumentative components separating source and target as *argumentative distance*, using the positive sign when the source precedes the target, and the negative sign otherwise. Inspired by works in other domains [42]–[44], we encode such a scalar number in a 10-bit array, using the first 5 bits for those cases where the source precedes the target, and the other 5 bits for the opposite case. The number of consecutive “1” values encodes the value of the distance, with a maximum value of 5. For example, if the argumentative distance is -3 , the encoding is 00111 00000; if the argumentative distance is 2, the encoding is 00000 11000.

C. The RESARG architecture

We use our own previous system [7] as a baseline. We refer to it as RESARG. Its architecture, depicted in Figure 1a, is based on residual networks [35] and comprises the following macro blocks:

- two deep embedders, one for sources and one for targets, that manipulate token embeddings;
- a dense encoding layer that reduces the dimensionality of the features;
- a biLSTM that processes the sequences;
- a residual network;
- the final-stage classifiers.

The purpose of the deep embedders is to fine-tune the pre-trained embeddings, a common procedure in deep learning-based NLP solutions [45] whose usefulness was confirmed by preliminary experiments. Each embedder is composed of a single residual block consisting of four pre-activated time-distributed dense layers. Accordingly, each layer applies the same transformation to each embedding, regardless of their position inside the sentence. All the layers have 50 neurons, except for the last one, which has 300 neurons.

The dense encoding layer is necessary in order to obtain an LSTM with fewer parameters, thus reducing the time needed for training, and limiting overfitting. It applies a time-distributed dense layer, which reduces the embedding size to 50, and a time average-pooling layer [46], which reduces the sequence size by a factor of 10.

The resulting sequences are then given as input to the same bidirectional LSTM, producing a single representation of size

50 for each component. Thus, for each proposition, T embeddings of size 300 are transformed first into T embeddings of size 50, then into $T/10$ embeddings of size 50, and finally in a single feature of size 50.

Source and target are processed in parallel in the first three blocks, then concatenated together, along with the encoding of the distance, and given as input to the final residual network. The first level of the final residual network is a dense encoding layer with 20 neurons, while the residual block is composed of a layer with 5 neurons and one with 20 neurons. The outputs of the first and the last layers of the residual networks are summed up and provided as input to the classifiers.

The final stage of RESARG are three independent softmax classifiers used to predict the source, the target, and the relation labels. Each classifier, which predicts a label for a dedicated task, contributes simultaneously to our learning model. The link classifier is obtained by summing the relevant scores produced by the relation classifier, aggregating the probability assigned to the relation labels into a single link label.

All the dense layers use the rectifier activation function [47], and they randomly initialize weights with He initialization [48]. The application of all non-linearity functions is preceded by batch-normalization layers [49] and by dropout layers [50], with probability $p = 0.1$.

D. The RESATTARG architecture

Motivated by the remarkable results obtained by attention-based architectures in NLP tasks, we have extended RESARG by including a neural attention block after the bi-LSTM module. To better exploit the new attention module, we removed the time pooling layer from the dense encoding block, so as to avoid loss of information along the temporal axis, and to maintain the whole output sequence from the LSTM. Therefore, in this new model, the input and the output of the LSTM module have size $(T, 50)$. The resulting architecture, named RESATTARG, is depicted in Figure 1b.

The attention module is implemented as *coarse-grained parallel co-attention* [17], so as to consider both components at the same time while computing attention on each of them. Our method consists of exploiting the average embedding of one proposition as a query element while computing attention on the other, similarly to what has been done in [51]. Specifically, calling K_s and K_t the outputs of the bi-LSTM obtained from, respectively, the processing of the source and the target propositions, we compute the (masked) average of K_t , obtaining a single embedding g_t of size 50 (Eq. 1). This embedding is used as query element to compute *additive soft attention* [17] on K_s , obtaining a single source context vector c_s of size 50. The details of this process are described in Equations 2, 3, and 4, where the matrices W_1 , W_2 and the vectors b , w_3 are learnable parameters. An equivalent symmetric procedure is used to compute attention on K_t so as to obtain c_t . The output of this block are two embeddings of size 50, as in our previous architecture.

$$g_t = \text{masked_avg}(K_t) \quad (1)$$

$$e_s = w_3^T \text{relu}(W_1 K_s + W_2 g_t + b) \quad (2)$$

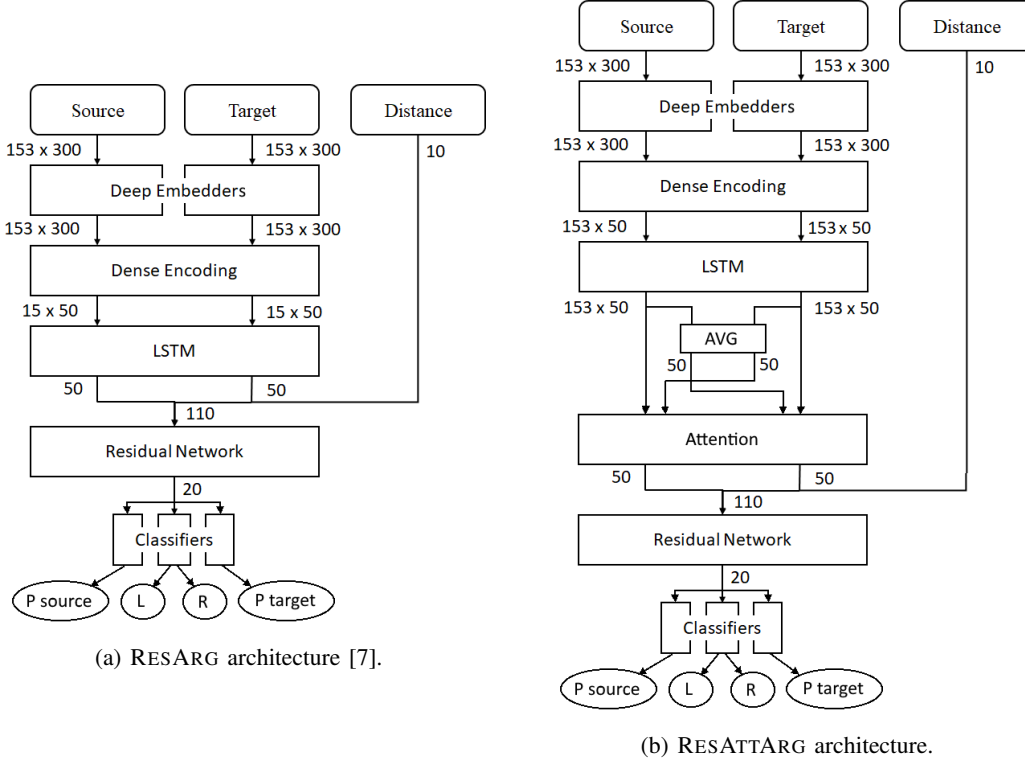


Fig. 1: A block diagram of the proposed architectures. The figure shows, next to each arrow, the dimensionality of the data involved (using the CDCP temporal size $T = 153$), so as to clarify the size of the inputs and the outputs of each block.

$$\mathbf{a}_s = \text{softmax}(\mathbf{e}_s) \quad (3)$$

$$\mathbf{c}_s = \sum_{i=1}^T \mathbf{k}_{si} \mathbf{a}_{si} \quad (4)$$

The resulting architecture has close to 5.5M parameters, 140,000 of which are trainable. If compared with other state-of-the-art neural architectures, such as BERT_{BASE} and its 110M parameters, RESATTARG is considerably smaller, and accordingly it is less computationally demanding.

E. Repeated trainings and ensemble learning

Since the training of neural models is non-deterministic, the results of a single training procedure are influenced by the random seed that is used, thus they may not be reliable or reproducible [52], [53]. Such problem also affects our previous results [7], since they were obtained from a single training experiment.

We have decided to replicate that experiment by repeating the training procedure 10 times, with different seeds, obtaining 10 trained neural networks for each configuration. We will evaluate our models in two different ways. At first, we will consider the average of the scores obtained by every single network for each metric. Then, we evaluate the predictions obtained using all the 10 models in ensemble voting.

In our ensemble setting the class of each entity is assigned as the class voted by the majority of the networks. This technique is similar to the concept of bootstrap aggregating, also known as bagging [54]. However, while in standard

bagging each model is trained on a random sample of the training set, here we train all the models on the same training set, since stochastic elements are already present in the training procedure itself. We have chosen this ensemble method for the sake of simplicity, but more advanced techniques do exist and may yield better results [55].

IV. CORPORA

We validate our approach on 4 corpora differing from each other in various dimensions: the domain of the documents, their average length, the formatting, and the argumentative model followed for the annotations.

A. Cornell eRulemaking Corpus (CDCP)

The Cornell eRulemaking Corpus (CDCP) [13], [56] consists of user-generated documents in which specific regulations are discussed. The authors have collected user comments from an eRulemaking website on the topic of Consumer Debt Collection Practices rule. The corpus contains 731 user comments, for a total of about 4,700 components, all considered to be argumentative.

As typical of user-generated data, the comments are not structured, and often present grammatical errors, typos, and do not follow usual writing conventions (such as the blank space after the period mark). This complicates pre-processing, since most of the off-the-shelf tools turn out to be inaccurate even in simple tasks such as tokenization.

Annotations follow the argument model proposed in [57], where links are constrained to form directed graphs. The

corpus is suitable both for component and relation classification, since it presents 5 classes of propositions and two types of links. We will use the version of CDCP without nested proposition and guaranteed transitive closure [13].

The components are addressed as propositions, and they consist of a sentence or a clause. Propositions are divided into POLICY (17%), VALUE (45%), FACT (16%), TESTIMONY (21%) and REFERENCE (1%). Only 3% of more than 43,000 possible proposition pairs are linked; almost all links are labeled as REASON (97%), whereas only a few are labeled as EVIDENCE (3%).

The unstructured nature of documents, the strong unbalance between the classes, and the presence of noise make the corpus particularly challenging for all the subtasks of argument mining, especially those that involve the relationships between components.

B. *AbstRCT*

The AbstRCT Corpus [32] extends previous work [58], and consists of abstracts of scientific papers regarding randomized control trials for the treatment of specific diseases (i.e., neoplasm, glaucoma, hypertension, hepatitis b, diabetes). The final corpus contains 659 abstracts, for a total of about 4,000 argumentative components. AbstRCT is divided into three parts: neoplasm, glaucoma, and mixed. The first one contains 500 abstracts about neoplasm, divided into train (350), test (100), and validation (50) splits. The remaining two are designed to be test sets. One contains 100 abstracts for glaucoma, the other 20 abstracts for each disease⁴.

Components are labeled as EVIDENCE (2,808) and CLAIM (1,390), while relations are labeled as SUPPORT (2,259) and ATTACK (342).⁵ About 10% of about 25,000 possible component pairs have a labeled relationship. The argumentative model chosen for annotation enforces only one constraint: claims can have an outgoing link only to other claims.

With respect of CDCP, this corpus is less noisy and the distribution of the classes is more balanced. We have chosen this as a benchmark to demonstrate that our approach is independent of the domain and of the argument model.

C. *Doctor Inventor Argumentative Corpus (DrInventor)*

The Doctor Inventor Argumentative Corpus (DrInventor) [59] is the result of an extension of the Doctor Inventor corpus [60], which includes an annotation layer containing argumentative components and relations. DrInventor consists of 40 scientific publications from computer graphics, which contain about 12,000 argumentative component labels, as well as annotations for other tasks.

The classes of argumentative components are DATA (4,093), OWN CLAIM (5,445), and BACKGROUND CLAIM (2,751). The former two are related to the concepts of premises and claims, while the latter is something in between, since it is a

claim related to some background knowledge, such as that made by another author in a previous work. The relation classes are SUPPORTS (5,790), CONTRADICTS (696), and SEMANTICALLY SAME (44), since it is common practice in scientific publications to re-iterate the same claim (or more rarely the same data) multiple times.

Since DrInventor includes documents where the structure of the discourse is complex, and data are often presented along with claims, it makes argument mining more challenging: in more than 1,000 cases some components are split into multiple text sequences, located in non-contiguous parts of the documents. This phenomenon mostly concerns claims, but data are affected too, in fewer cases. This introduces the difficulty of recognizing different segments of the documents as part of a single component and makes link prediction more difficult to address through non-pipeline approaches.

The unbalanced distribution between the three classes and the presence of split components makes this corpus quite challenging for link prediction, a difficulty which is highlighted also by the low inter-annotator agreement reported in the original paper.

D. *Persuasive Essays Corpus (UKP)*

The Persuasive Essays Corpus (UKP) [4] consist of 402 documents coming from an online community where users post essays and other material, provide feedback, and advise each other. The dataset is divided into a test split of 80 essays and a training split with the remaining documents.

UKP defines three classes of argumentative components: MAJOR CLAIM (751), CLAIM (1,506), and PREMISE (3,832). Premises may be linked to CLAIMS through relations of SUPPORT (3,613) or ATTACK (219). MAJOR CLAIMS are not linked to other components.⁶ While these classes of components are not unlike those used in other datasets, the argumentation model instead is tailored on this specific corpus. It amounts to an argument graph consisting of trees, each rooted on a claim. Each tree can only include components belonging to the same paragraph. Claims do not have outgoing relationships because only premises can descend from claims. Each premises has exactly one outgoing relation. Finally, the structure of the argumentation follows domain-specific conventions. For example, in most cases, the MAJOR CLAIM is in the first or the last paragraph of the document, and more often than not it is the only argumentative component in its paragraph.

Thanks to the highly constrained nature of the UKP data, we expect methods based on the document structure to have an edge over general, structure-agnostic methods such as the one we propose. However, we believe it is important to include also this type of data in the present study, in order to evaluate our approach in as many and diverse scenarios as possible, and better understand its advantages as well as its limitations.

⁴Glaucoma and neoplasm documents of the mixed set are present also in the respective test set.

⁵The corpus allows also the distinction between CLAIM/MAJOR CLAIM and ATTACK/PARTIAL ATTACK. For the sake of consistency with previous works, this detail will not be considered.

⁶In fact, each component is linked to a MAJOR CLAIM via an attribute called *stance*. Therefore, one could use this dataset for stance detection, by creating explicit relationships toward the major claims. However, that would be outside the scope of this work.

V. EXPERIMENTAL SETTING

We consider a multi-task formulation for our learning problem. The loss function is given by the weighted sum of four different components: the categorical cross-entropy on three labels (source and target categories, link relation category) and an L_2 regularization on the network parameters.

We initially evaluate our new architecture against our previous model and the structured learning approach of [13] on CDCP, presenting an ablation study of the new components we have introduced. Then, we extend the evaluation to other three data sets, for which we compare our approach against the state-of-the-art.

In our approach each component is involved in many pairs, both as a source and as a target, and accordingly it is classified multiple times by the same network. The label will be assigned by the model by considering the average probability computed by the ensemble for each class, and by thus choosing the class with the highest score. Alternative approaches could be to assign the class that results to be the most probable in most of the cases, thus relying on a majority vote. A further option could be to simply consider the label with the highest confidence. However, the latter procedure might be more sensitive to outliers, because the misclassification of a component in just one pair would lead to the final misclassification of the component, regardless of all the other pairs. A deeper analysis of different techniques to address these issues is left to future research.

A. Data Preparation

Our architecture allows us to use the CDCP and AbstRCT datasets directly, without need for further pre-processing.

For what concerns DrInventor, instead, specific data pre-processing is needed to address two aspects of this dataset: the presence of lengthy documents and split components. Lengthy documents make it inconvenient to consider all the possible pairs of argumentative units. Doing so would not only be infeasible with regular computational resources, but it would also yield an extremely unbalanced dataset for link prediction, with less than 1% of pairs linked. We thus filtered out all the pairs that did not appear in the same section of the document, and whose argumentative distance is included between -10 and +10. A second peculiarity of this dataset is the presence of components that include non-argumentative material. These “split components” are made of two sequences x and y separated by a third, non-argumentative sequence z . In those cases, we split x and y into two unrelated components, and attributed them the same label, the same links, and the same argumentative relations with the other components. The resulting dataset consists of about 8,700 links out of 240,000 possible pairs, which amount roughly to 3.6%. Among these links, SUPPORTS amount to 89%, CONTRADICTS to 10%, and the remaining 1% are SEMANTICALLY SAME relations.

Regarding UKP, like others did before us [4], [13], we also consider exclusively pairs of components that belong to the same paragraph. However, many paragraphs contain only a single component. That is the case, for instance, with about 400 paragraphs containing a single major claim. In order to

include also them in our pair-based classification method, we decided to introduce “self pairs” into our dataset, which are instances where the same component acts both as source and target. This significantly increases the number of pairs (from 22,000 to 28,000). So, to improve optimization and enable a comparison with previous approaches, we did not consider these pairs for link prediction and relation classification in validation and testing.

B. Comparison with other methods

Not all the approaches to AM are easily compared against one another. This is the case, for example, of approaches that perform only few tasks versus end-to-end systems, or pipeline versus joint learning approaches. Since we perform component classification on propositions or sentences, to make our results comparable with architectures that perform it token-wise, we split each classified component into tokens that share the same label, and compute the evaluation of token-wise classification. Since the tokenization method may not be the same one used by other approaches, the final results may not be perfectly comparable, but we believe that this minor difference will not introduce appreciable errors.

We shall also remark that in our approach we consider argumentative components as already selected and perfectly bounded, therefore we perform component classification only between argumentative classes and we do not consider the “non-argumentative” class as a possibility. This makes our figures incomparable against those obtained by architectures that address both component identification and classification at once, such as [32], since they include “non-argumentative” among the possible classes and thus address a harder problem. A similar consideration holds regarding the pipeline approaches that perform evaluation of each step based on the result of the previous one instead of using the gold standard. In this case, the errors introduced by early steps introduce noise which may affect the evaluation of subsequent steps. It is once again the case of [32], where errors obtained during the first step may introduce noise in the link prediction/relation classification tasks. We could not find a solution to this problem, but we argue that, nevertheless, a qualitative evaluation of our method can still benefit from a comparison with these other approaches.

C. Optimization

We shall remark that the hyper-parameters of the architecture and of the learning model have been tuned on the validation set of CDCP. It is also important to highlight that we use the same set of hyper-parameters in all the experiments. Our purpose is to test whether our approach can yield satisfactory results across different and heterogeneous corpora without the need of re-tuning, and therefore limiting its cost and its environmental impact [39]. Nonetheless, we are aware that performing a specific calibration for each corpus would probably improve our results. We use the Adam optimizer [61] with parameters $b_1 = 0.9$ and $b_2 = 0.9999$, applying proportional decay of the initial learning rate $\alpha_0 = 5 \times 10^{-3}$. The weights of the four components of the loss function are

set to 1 for the cross entropy of source and target, 10 for the cross entropy of relation, and 10^{-4} for the regularization. The training was early-stopped after 100 epochs with no improvements on the F_1 score of the Link class computed over validation data, except for DrInventor, where we early-stopped after 20 epochs of patience due to the dataset’s size and much heavier computational footprint.

VI. RESULTS AND DISCUSSION

In this section we will present the experimental results obtained for each corpus. We will compare the RESATTARG and RESARG architectures, thus assessing the impact of the attention module, and we will also analyze the performance gain introduced by the ensemble approach.

A. CDCP

We used the same validation set as in our earlier work [7], which was created by randomly selecting documents from the original training split with 10% probability. We used the remaining documents as training data and the original test split as is. To provide a summary evaluation, following [13], we measured the performance of the models by computing the F_1 score for links, propositions, and the average between the two. More specifically, for the links we measured the F_1 of the positive classes, whereas for the propositions we used the score of each class and then we computed the macro-average. We also reported the F_1 score for each relation class, alongside their macro-average. The NONE class of relation classification corresponds to the negative class of link prediction.

A first question we address was whether our results with a single model were solid or to what extent influenced by the non-deterministic nature of the training procedure. We compared our baseline model with the average scores obtained by 10 networks, with our ensemble setting, and against the structured approach used in [13]. The results are shown in Table I. The average computed over the 10 networks leads to a worse performance on Link prediction with respect to our previous results. This difference suggests that our previous results were the result of a “lucky” training. Nonetheless, the average score on the two tasks remains similar (between 47 and 48), just a few points below the state of the art. The ensemble approach substantially improves the results, outperforming the structured learning approach on both tasks. The results on link prediction are still below those obtained in the first experiment, if only by less than 1%.

Introducing the attention module in the architecture leads to appreciable improvements for both the average and the ensemble approach. In particular, the latter’s performance marks a new state-of-the-art for this corpus, even for the relation classification task. As far as relation label prediction, our approaches fail to predict the EVIDENCE relation. It is a negative result, but hardly surprising, since EVIDENCE is a rather rare class in this dataset (less than 1% of all relations).

To estimate the agreement among the networks in the ensemble architecture, and have a measure of the architecture’s robustness against the implicit randomness of the training procedure, we have computed Krippendorff’s alpha [62] for

the three tasks. We obtained $\alpha = 0.68$ for component classification, and $\alpha = 0.50$ for both link prediction and relation classification. These values are similar to the IAA obtained by the authors of the corpus, and confirm the difficulty of the link prediction task.

Figure 2 shows confusion matrices for component classification on CDCP. Unsurprisingly, the most common mistake regards the prediction of facts as values – VALUE being by far the largest class in the corpus, and so affected by many false positives. Such an ambiguity between the two classes has also been reported during the annotation process.

Interestingly, the confusion matrices of the structured approach and of our methods are quite similar. We speculate that our networks may have learned a behavior similar to that produced by the structured approach, with no need to receive any of the constraints or information regarding the argumentative structure that are instead injected in the structured approach.

B. AbstRCT

For what concerns AbstRCT, we compare our architectures against the best methods presented by its authors [32], whose results are reported in the first rows of Tables II and III. We trained and validated our model on the respective splits of the Neoplasm dataset, using the remainder of the dataset for testing. For reasons we already explained, the approach presented by Mayer et al. [32] is not directly comparable with ours, therefore the comparison can only be qualitative. To ease comparison with future approaches, we report in Table IV some additional details on our results.

As for component classification, RESATTARG with ensemble yields the best result, performing comparably with the state of the art. Our approaches obtain better scores for EVIDENCE than CLAIM on all datasets. Similarly to the Transformer-based approaches, our architectures perform better on the mixed test set than on the neoplasm one. We yield better results on all datasets for what concerns the micro f_1 score. However, for what concerns macro F_1 , although our architecture improves the previous approaches on Neoplasm, it is outperformed by BioBERT on Glaucoma and Mixed. In relation classification, RESATTARG with ensemble outperforms all the other models on Neoplasm and Glaucoma, and it performs about 1% worse than SciBERT on Mixed. It is interesting to notice that in this task BioBERT is largely outperformed by our approach. Almost all the metrics confirm that the introduction of attention and ensemble improve our architectures. The agreement between the networks RESATTARG is very high for token-wise component classification in each dataset (α between 0.81 and 0.83), and lower but still acceptable for the other two tasks ($\alpha = 0.67$ on neoplasm and $\alpha = 0.62$ for the other two).

These good results indicate that our method may be a valuable approach with well-structured corpora. Moreover, such results are attained without resorting to contextual embeddings or pre-training on domain-related corpora, but by only relying on non-contextual, general-purpose embeddings.

TABLE I: Results of the experiments involving multiple trainings of the same models and use of attention on CDCP. From left to right: our previous result [7] obtained with a single training of RESARG, the average scores of the same architecture trained 10 times, the scores of the ensemble learning setting of the same model, the average and the ensemble scores of the new attention-based architecture RESATTARG, and the best results of structured approaches based on SVM and RNN. For each class, the number of instances is reported in parenthesis. F_1 and macro F_1 percentage scores are reported.

	Single [7]	RESARG Average	Ensemble	RESATTARG Average	RESATTARG Ensemble	Structured [13] SVM	RNN
Average (Link and Components)	47.28	47.75	52.14	51.57	54.22	50	43.5
Link (272)	29.29	24.99	28.76	27.40	29.73	26.7	14.6
Components (973)	65.28	70.51	75.53	75.75	78.71	73.50	72.7
VALUE (491)	72.19	72.30	75.37	77.84	80.37	76.4	73.70
POLICY (153)	74.36	75.39	79.60	80.09	82.55	77.3	76.8
TESTIMONY (204)	72.86	73.46	76.33	76.42	81.19	71.7	75.8
FACT (124)	40.31	41.39	46.37	44.39	49.42	42.5	42.2
REFERENCE (1)	66.67	90.00	100	100	100	100	100
Relation (9,484)	-	41.78	42.25	42.69	42.95	-	-
REASON (265)	30.02	25.10	28.57	27.88	30.56	-	-
EVIDENCE (7)	0	2.50	0	2.22	0	-	-
NONE (9,212)	-	97.76	98.21	98.03	98.32	-	-

Structured SVM full	Predicted				
	P	F	T	V	R
True	P 0.76	0.05	0.04	0.16	0.00
	F 0.04	0.44	0.10	0.42	0.00
	T 0.01	0.06	0.72	0.21	0.00
	V 0.05	0.11	0.08	0.76	0.00
	R 0.00	0.00	0.00	0.00	1.00

ResArg Single	Predicted				
	P	F	T	V	R
True	P 0.76	0.06	0.01	0.17	0.00
	F 0.06	0.42	0.08	0.44	0.00
	T 0.00	0.06	0.75	0.18	0.00
	V 0.07	0.12	0.10	0.70	0.00
	R 0.00	0.00	0.00	0.00	1.00

ResArg Ensemble	Predicted				
	P	F	T	V	R
True	P 0.78	0.07	0.01	0.14	0.00
	F 0.03	0.54	0.06	0.36	0.00
	T 0.00	0.07	0.77	0.15	0.00
	V 0.04	0.15	0.09	0.72	0.00
	R 0.00	0.00	0.00	0.00	1.00

ResAttArg Ensemble	Predicted				
	P	F	T	V	R
True	P 0.80	0.08	0.01	0.11	0.00
	F 0.02	0.52	0.05	0.41	0.00
	T 0.00	0.05	0.80	0.14	0.00
	V 0.04	0.10	0.06	0.80	0.00
	R 0.00	0.00	0.00	0.00	1.00

Fig. 2: Confusion matrices for component classification on CDCP. From left to right: Structured Learning approach, our previous result with RESARG [7], RESARG used in ensemble fashion, and RESATTARG used in ensemble fashion.

C. DrInventor

To the best of our knowledge, the only approach tested on this corpus is the architecture for token-wise component classification used by Lauscher et al. [12], which makes use of GloVe embeddings and a Bi-LSTM followed by a feed-forward neural network with a single hidden layer as classifier. We thus consider such an approach as a baseline. Like Lauscher et al., we reserved 30% of the documents of the DrInventor corpus as test set, and 20% of the remaining part as validation set. It is worth remarking that for the tasks of link prediction and relation classification we are considering a limited number of pairs.

Tables V and VI includes a detailed report of our performance on the dataset. We outperform the baseline by a wide margin. Moreover, we address two additional tasks, link prediction and relation classification, thus offering a benchmark for future work. These results confirm once more that attention and ensemble together give a crucial contribution to the classifier.

Differently from previous experiments, the agreement between the networks RESATTARG is similar for all the tasks, with only $\alpha = 0.56$ for component classification and $\alpha = 0.60$ for the remaining tasks. The agreement for Component Classification is lower than on the previous datasets and may suggest that this dataset is more challenging.

Our model is incapable of classifying the SEMANTICALLY SAME relation and has difficulties also with CON-

TRADICTS. That is hardly surprising, if we consider that these are the two least represented classes in this dataset. It is less straightforward to understand why the model is better at classifying BACKGROUND CLAIM rather than DATA, even if the latter are more represented than the former. We speculate it may be related to the fact that in some instances data may amount to citations or text other than proper sentences.

D. UKP

UKP comes with two strong baselines: the ILP joint model proposed by the authors of the dataset [4] and Niculae et al.'s structured learning approach [13]. We compared based on the original test split of the dataset, using about 10% of the documents of the training split as validation split. Apparently our approach is largely below the baselines, with a difference in F_1 scores between 20% and 30%. The agreement between the networks is also low, with $\alpha = 0.57$ for component classification and $\alpha = 0.38$ for link prediction, assessing them as nearly acceptable for the first task but unreliable for the others. We interpret these results as an indication that a general-purpose and domain-agnostic architecture such as ours struggles with datasets characterized by strong regularities that can be best exploited by domain-oriented approaches.

VII. CONCLUSIONS

In this paper we presented RESATTARG, a new neural architecture for argument mining based on residual networks,

TABLE II: Results of component classification on AbstrCT. We report the F_1 score related to the micro average, the macro average, the EVIDENCE class, and the CLAIM class obtained on the 3 test sets. Our approach is not directly comparable with component classification of [32] and the comparison must be considered qualitative.

Level	Approach	Neoplasm				Glaucoma				Mixed			
		f_1	F_1	E	C	f_1	F_1	E	C	f_1	F_1	E	C
Token	Transformer-based [32]												
	BioBERT+GRU+CRF	90	84	90	87	92	91	91	93	92	91	92	91
	SciBERT+GRU+CRF	90	87	92	88	91	89	91	93	91	88	93	90
	Our approach												
	RESARG (avg)	90.66	88.20	93.58	82.81	91.84	87.49	94.86	80.11	91.25	87.79	94.29	81.29
	RESARG+Ensemble	90.75	88.10	93.72	82.47	92.50	88.48	95.28	81.68	91.61	88.21	94.54	91.61
	RESATTARG (avg)	90.80	88.60	93.59	83.61	92.02	88.02	94.93	81.11	91.58	88.72	94.40	83.04
Component	RESATTARG+Ensemble	92.12	90.14	94.56	85.72	92.92	89.35	95.52	83.19	92.79	90.26	95.23	85.30
	Our approach												
	RESARG (avg)	87.42	86.18	90.31	82.04	88.08	85.53	91.59	79.48	88.20	86.74	91.13	82.35
	RESARG+Ensemble	87.76	86.38	90.71	82.05	89.39	87.13	92.53	81.74	89.00	87.59	91.77	83.42
	RESATTARG (avg)	87.32	86.19	90.11	82.27	88.50	86.26	91.79	80.72	88.65	87.51	91.27	83.74
	RESATTARG+Ensemble	88.92	87.87	91.44	84.30	89.73	87.71	92.69	86.54	90.67	89.70	92.86	82.72

TABLE III: Results of relation classification on AbstrCT. We report the macro averaged F_1 score obtained on the 3 test sets.

Approach	Neoplasm	Glaucoma	Mixed
Transformer-based [32]			
BioBERT	64	58	61
SciBERT	68	62	69
RoBERTa	67	66	67
Our approach			
RESARG (avg)	59.15	57.23	60.31
RESARG+Ensemble	63.16	61.86	68.35
RESATTARG (avg)	66.49	62.68	63.47
RESATTARG+Ensemble	70.92	68.40	67.66

TABLE IV: Results of RESATTARG with Ensemble for Link Prediction and Relation Classification on AbstrCT. The “Link” column refers to the F_1 score of the positive class. The “Relation” column reports the result of the macro F_1 score. The other columns report the F_1 score of the respective classes.

Test set	Link	Relation	SUPPORT	ATTACK	NONE
Neoplasm	54.43	70.92	52.77	65.38	94.54
Glaucoma	55.23	68.40	54.73	56.00	94.36
Mixed	51.20	67.66	49.62	59.09	94.21

multi-task learning, neural attention, and ensemble learning. Our approach does not rely on domain-tailored features or encodings. On the contrary, it only uses general-purpose embeddings and a broadly applicable distance feature, making it suitable for any domain and argumentative model. Moreover, RESATTARG is considerably smaller than other state-of-the-art approaches, making it less expensive to train, and more sustainable from an environmental perspective [39].

RESATTARG outperforms state-of-the-art architectures on a variety of tasks and datasets. We conducted ablation studies to demonstrate that the attention module and the ensemble learning addition give a positive contribution, improving on our previous architecture [7]. The use of ensemble also in-

TABLE V: Results of component classification on DrInventor. “Average” is the macro F_1 score, the remaining columns report the F_1 score of the classes OWN CLAIM, BACKGROUND CLAIM, DATA.

Level	Approach	Average	O C	B C	D
Token	Bi-LSTM [12]	44.70	-	-	-
	RESARG (avg)	58.27	72.31	51.36	51.14
	RESARG+Ens	61.16	75.77	54.24	53.48
	RESATTARG (avg)	61.77	73.66	57.70	53.95
	RESATTARG+Ens	65.71	78.03	61.58	57.53
Component	RESARG (avg)	60.62	65.65	45.63	70.56
	RESARG+Ens	62.97	68.97	48.03	71.90
	RESATTARG (avg)	66.19	68.61	56.07	73.89
	RESATTARG+Ens	69.64	73.13	59.63	76.15

TABLE VI: Results of RESATTARG with Ensemble for Link Prediction, and Relation Classification on DrInventor. The “Link” column refers to the F_1 score of the positive class. The “Relation” column reports the result of the macro F_1 score. The other columns report the F_1 score of the respective classes.

Link	Relation	SUPP	CON	SEM SAME	NONE
43.66	37.72	45.90	6.61	0	98.37

creases the robustness of this approach against the intrinsic randomness of neural architecture training.

The main limitations of RESATTARG are its seemingly poor performance on datasets characterized by strong regularities and its limited scalability to large documents. To fill this gap, we believe that neural-symbolic approaches [63] may enable a systematic and modular integration of background knowledge. Such a knowledge would contribute during the optimization process, so as to influence and improve the training, without compromising the generality of the neural architecture. For what concerns scalability, we have addressed this problem by limiting the range of argumentative relationships using a fixed-

size window, with the drawback of imposing a constraint on the model of the argument. Alternatively to our pair-based approach, other authors have proposed multiple-choice classifiers [32], pointer networks [23], and sequence labelling [9]. Such methods should scale better, but they enforce a constraint on the argument model as well, imposing that any component can have only one outgoing relationship, which makes them unsuitable to some corpora. While successful approaches to argument retrieval have been recently published [64], scalability is still an open challenge for AM systems aiming to reconstruct the argumentation structure.

REFERENCES

- [1] M. Lippi and P. Torroni, "Argumentation mining: State of the art and emerging trends," *ACM Trans. Internet Technol.*, vol. 16, no. 2, pp. 10:1–10:25, Mar. 2016. [Online]. Available: <http://doi.acm.org/10.1145/2850417>
- [2] J. Lawrence and C. Reed, "Argument mining: A survey," *Computational Linguistics*, vol. 45, no. 4, pp. 765–818, 2020.
- [3] D. Walton, "Argumentation theory: A very short introduction," in *Argumentation in Artificial Intelligence*, G. Simari and I. Rahwan, Eds. Springer US, 2009, pp. 1–22.
- [4] C. Stab and I. Gurevych, "Parsing argumentation structures in persuasive essays," *Comput. Linguistics*, vol. 43, no. 3, pp. 619–659, 2017.
- [5] I. Persing and V. Ng, "End-to-end argumentation mining in student essays," in *HLT-NAACL*. The Association for Computational Linguistics, 2016, pp. 1384–1394.
- [6] N. Slonim, E. Aharoni, C. Alzate, R. Bar-Haim, Y. Bilu, L. Dankin, I. Eiron, D. Hershovich, S. Hummel, M. M. Khapra, T. Lavee, R. Levy, P. Matchen, A. Polnarov, V. C. Raykar, R. Rinott, A. Saha, N. Zwerdling, D. Konopnicki, and D. Gutfreund, "Claims on demand - An initial demonstration of a system for automatic detection and polarity identification of context dependent claims in massive corpora," in *COLING 2014, Dublin, Ireland*, L. Tounsi and R. Rak, Eds. ACL, 2014, pp. 6–9. [Online]. Available: <http://www.aclweb.org/anthology/C14-2002>
- [7] A. Galassi, M. Lippi, and P. Torroni, "Argumentative link prediction using residual networks and multi-objective learning," in *Proceedings of the 5th Workshop on Argument Mining, ArgMining@EMNLP*. Association for Computational Linguistics, 2018, pp. 1–10. [Online]. Available: <https://doi.org/10.18653/v1/w18-5201>
- [8] I. Habernal and I. Gurevych, "Argumentation mining in user-generated web discourse," *Computational Linguistics*, vol. 43, no. 1, pp. 125–179, 2017. [Online]. Available: <http://aclweb.org/anthology/J17-1004>
- [9] S. Eger, J. Daxenberger, and I. Gurevych, "Neural end-to-end learning for computational argumentation mining," in *ACL (1)*. Association for Computational Linguistics, 2017, pp. 11–22.
- [10] A. Søgaard and Y. Goldberg, "Deep multi-task learning with low level tasks supervised at lower layers," in *ACL (2)*. The Association for Computer Linguistics, 2016.
- [11] C. Schulz, S. Eger, J. Daxenberger, T. Kahse, and I. Gurevych, "Multi-task learning for argumentation mining in low-resource settings," in *NAACL-HLT (2)*. Association for Computational Linguistics, 2018, pp. 35–41.
- [12] A. Lauscher, G. Glavas, S. P. Ponzetto, and K. Eckert, "Investigating the role of argumentation in the rhetorical analysis of scientific publications with neural multi-task learning models," in *EMNLP*. Association for Computational Linguistics, 2018, pp. 3326–3338.
- [13] V. Niculae, J. Park, and C. Cardie, "Argument mining with structured svms and rnns," in *ACL (1)*. Association for Computational Linguistics, 2017, pp. 985–995.
- [14] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *ICLR*, 2015.
- [15] O. Vinyals, M. Fortunato, and N. Jaitly, "Pointer networks," in *NIPS*, 2015, pp. 2692–2700.
- [16] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *NIPS*. Curran Associates, Inc., 2017, pp. 5998–6008.
- [17] A. Galassi, M. Lippi, and P. Torroni, "Attention in natural language processing," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–18, 2020. [Online]. Available: <https://doi.org/10.1109/TNNLS.2020.3019893>
- [18] D. Suhartono, A. P. Gema, S. Winton, T. David, M. I. Fanany, and A. M. Arymurthy, "Argument annotation and analysis using deep learning with attention mechanism in bahasa indonesia," *J. Big Data*, vol. 7, no. 1, p. 90, 2020. [Online]. Available: <https://doi.org/10.1186/s40537-020-00364-z>
- [19] J. Lin, K. Y. Huang, H. Huang, and H. Chen, "Lexicon guided attentive neural network model for argument mining," in *ArgMining@ACL*. Association for Computational Linguistics, 2019, pp. 67–73.
- [20] C. Stab, T. Miller, B. Schiller, P. Rai, and I. Gurevych, "Cross-topic argument mining from heterogeneous sources," in *EMNLP*. ACL, 2018, pp. 3664–3674.
- [21] M. Spliethöfer, J. Klaff, and H. Heuer, "Is it worth the attention? a comparative evaluation of attention layers for argument unit segmentation," in *ArgMining@ACL*. Florence, Italy: Association for Computational Linguistics, Aug. 2019, pp. 74–82.
- [22] I. Habernal and I. Gurevych, "What makes a convincing argument? empirical analysis and detecting attributes of convincings in web argumentation," in *EMNLP*. ACL, 2016, pp. 1214–1223.
- [23] P. Potash, A. Romanov, and A. Rumshisky, "Here's my point: Joint pointer architecture for argument mining," in *EMNLP*. Association for Computational Linguistics, 2017, pp. 1364–1373.
- [24] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *NAACL-HLT (1)*. ACL, 2019, pp. 4171–4186.
- [25] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," in *NAACL-HLT*. Association for Computational Linguistics, 2018, pp. 2227–2237.
- [26] N. Reimers, B. Schiller, T. Beck, J. Daxenberger, C. Stab, and I. Gurevych, "Classification and clustering of arguments with contextualized word embeddings," in *ACL (1)*. Association for Computational Linguistics, 2019, pp. 567–578.
- [27] L. Lugini and D. J. Litman, "Contextual argument component classification for class discussions," in *COLING*. International Committee on Computational Linguistics, 2020, pp. 1475–1480.
- [28] H. Wang, Z. Huang, Y. Dou, and Y. Hong, "Argumentation mining on essays at multi scales," in *COLING*. International Committee on Computational Linguistics, 2020, pp. 5480–5493.
- [29] D. Trautmann, J. Daxenberger, C. Stab, H. Schütze, and I. Gurevych, "Fine-grained argument unit recognition and classification," in *AAAI*. AAAI Press, 2020, pp. 9048–9056.
- [30] P. Poudyal, J. Savelka, A. Ieven, M. F. Moens, T. Goncalves, and P. Quaresma, "ECHR: Legal corpus for argument mining," in *ArgMining*. Online: Association for Computational Linguistics, Dec. 2020, pp. 67–75.
- [31] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized BERT pretraining approach," *CoRR*, vol. abs/1907.11692, 2019. [Online]. Available: <http://arxiv.org/abs/1907.11692>
- [32] T. Mayer, E. Cabrio, and S. Villata, "Transformer-based argument mining for healthcare applications," in *ECAI 2020*, ser. Frontiers in Artificial Intelligence and Applications, vol. 325. IOS Press, 2020, pp. 2108–2115. [Online]. Available: <https://doi.org/10.3233/FAIA200334>
- [33] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinform.*, vol. 36, no. 4, pp. 1234–1240, 2020. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btz682>
- [34] I. Beltagy, K. Lo, and A. Cohan, "Scibert: A pretrained language model for scientific text," in *EMNLP*, K. Inui, J. Jiang, V. Ng, and X. Wan, Eds. Association for Computational Linguistics, 2019, pp. 3613–3618. [Online]. Available: <https://doi.org/10.18653/v1/D19-1371>
- [35] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*. IEEE Computer Society, 2016, pp. 770–778.
- [36] A. Conneau, H. Schwenk, L. Barrault, and Y. LeCun, "Very deep convolutional networks for text classification," in *EACL (1)*. Association for Computational Linguistics, 2017, pp. 1107–1116.
- [37] Y. Y. Huang and W. Y. Wang, "Deep residual learning for weakly-supervised relation extraction," in *EMNLP*. Association for Computational Linguistics, 2017, pp. 1803–1807.
- [38] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [39] E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning in NLP," in *ACL (1)*. Association for Computational Linguistics, 2019, pp. 3645–3650.
- [40] A. Peldszus and M. Stede, "An annotated corpus of argumentative microtexts," in *Argumentation and Reasoned Action: Proceedings of the 1st European Conference on Argumentation, Lisbon*, vol. 2, 2015, pp. 801–815.

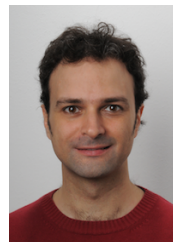
- [41] J. Pennington, R. Socher, and C. Manning, “Glove: Global vectors for word representation,” in *EMNLP*. Association for Computational Linguistics, 2014, pp. 1532–1543.
- [42] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, “Mastering the game of go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [43] F. Chesani, A. Galassi, M. Lippi, and P. Mello, “Can deep networks learn to play by the rules? A case study on nine men’s morris,” *IEEE Transactions on Games*, vol. 10, no. 4, pp. 344–353, 2018. [Online]. Available: <https://doi.org/10.1109/TG.2018.2804039>
- [44] A. Galassi, M. Lombardi, P. Mello, and M. Milano, “Model agnostic solution of CSPs via deep learning: A preliminary study,” in *CPAIOR*, ser. Lecture Notes in Computer Science, vol. 10848. Springer, 2018, pp. 254–262. [Online]. Available: https://doi.org/10.1007/978-3-319-93031-2_18
- [45] D. Jurafsky and J. H. Martin, *Speech and Language Processing (3rd Ed. draft)*, 2021. [Online]. Available: <https://web.stanford.edu/~jurafsky/slp3/>
- [46] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa, “Natural language processing (almost) from scratch,” *Journal of Machine Learning Research*, vol. 12, no. Aug, pp. 2493–2537, 2011.
- [47] X. Glorot, A. Bordes, and Y. Bengio, “Deep sparse rectifier neural networks,” in *AISTATS*, ser. JMLR Proceedings, vol. 15. JMLR.org, 2011, pp. 315–323.
- [48] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *ICCV*. IEEE Computer Society, 2015, pp. 1026–1034.
- [49] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *ICML*, ser. JMLR Workshop and Conference Proceedings, vol. 37. JMLR.org, 2015, pp. 448–456.
- [50] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [51] D. Ma, S. Li, X. Zhang, and H. Wang, “Interactive attention networks for aspect-level sentiment classification,” in *IJCAI*. IJCAI.org, 2017, pp. 4068–4074.
- [52] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016. [Online]. Available: <http://www.deeplearningbook.org>
- [53] N. Reimers and I. Gurevych, “Reporting score distributions makes a difference: Performance study of lstm-networks for sequence tagging,” in *EMNLP*. Association for Computational Linguistics, 2017, pp. 338–348.
- [54] L. Breiman, “Bagging predictors,” *Mach. Learn.*, vol. 24, no. 2, pp. 123–140, 1996. [Online]. Available: <https://doi.org/10.1007/BF00058655>
- [55] X. Dong, Z. Yu, W. Cao, Y. Shi, and Q. Ma, “A survey on ensemble learning,” *Frontiers Comput. Sci.*, vol. 14, no. 2, pp. 241–258, 2020.
- [56] J. Park and C. Cardie, “A corpus of erulemaking user comments for measuring evaluability of arguments,” in *LREC*. European Language Resources Association (ELRA), 2018.
- [57] J. Park, C. Blake, and C. Cardie, “Toward machine-assisted participation in erulemaking: an argumentation model of evaluability,” in *ICAIL*. ACM, 2015, pp. 206–210.
- [58] T. Mayer, E. Cabrio, M. Lippi, P. Torroni, and S. Villata, “Argument mining on clinical trials,” in *COMMA*, ser. Frontiers in Artificial Intelligence and Applications, vol. 305. IOS Press, 2018, pp. 137–148.
- [59] A. Lauscher, G. Glavas, and S. P. Ponzetto, “An argument-annotated corpus of scientific publications,” in *ArgMining@EMNLP*. Association for Computational Linguistics, 2018, pp. 40–46.
- [60] B. Fisas, F. Ronzano, and H. Saggion, “A multi-layered annotated corpus of scientific papers,” in *LREC*. European Language Resources Association (ELRA), 2016.
- [61] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *ICLR (Poster)*, 2015.
- [62] K. Krippendorff, *Content analysis: An introduction to its methodology*. Sage publications, 2018.
- [63] A. Galassi, K. Kersting, M. Lippi, X. Shao, and P. Torroni, “Neural-symbolic argumentation mining: An argument in favor of deep learning and reasoning,” *Frontiers in Big Data*, vol. 2, p. 52, 2019. [Online]. Available: <https://doi.org/10.3389/fdata.2019.00052>
- [64] L. Ein-Dor, E. Shnarch, L. Dankin, A. Halfon, B. Sznajder, A. Gera, C. Alzate, M. Gleize, L. Choshen, Y. Hou, Y. Bilu, R. Aharonov, and N. Slonim, “Corpus wide argument mining - A working solution,” in *AAAI*. AAAI Press, 2020, pp. 7683–7691.



Andrea Galassi holds a Master Degree in Computer Engineering, and is a PhD student and Research Fellow at the Department of Computer Science and Engineering at the University of Bologna. His research activity concerns artificial intelligence and machine learning, focusing on argumentation mining and related NLP tasks. Other research interests involve deep learning applications to games and CSPs.



Marco Lippi is associate professor at the Department of Sciences and Methods for Engineering, University of Modena and Reggio Emilia. He previously held positions at the Universities of Florence, Siena and Bologna, and he was visiting scholar at UPMC, Paris. His work focuses on machine learning and artificial intelligence, with applications to several areas, including argumentation mining, legal informatics, and medicine. In 2012 he was awarded the “E. Caianiello” prize for the best Italian PhD thesis in the field of neural networks.



Paolo Torroni is an associate professor with the University of Bologna since 2015. His main research focus is artificial intelligence. He edited over 20 books and special issues and authored over 150 articles in computational logics, multi-agent systems, argumentation, and natural language processing. He serves as an associate editor for *Fundamenta Informaticae* and *Intelligenza Artificiale*.