

# Transferring Confluent Knowledge to Argument Mining

João A. Rodrigues and António Branco

University of Lisbon

NLX – Natural Language and Speech Group, Department of Informatics

Faculdade de Ciências, Campo Grande, 1749-016 Lisboa, Portugal

{jarodrigues, ambranco}@fc.ul.pt

## Abstract

Relevant to all application domains where it is important to get at the reasons underlying sentiments and decisions, argument mining seeks to obtain structured arguments from unstructured text and has been addressed by approaches typically involving some feature and/or neural architecture engineering.

By adopting a transfer learning methodology, and by means of a systematic study with a wide range of knowledge sources promisingly suitable to leverage argument mining, the aim of this paper is to empirically assess the potential of transferring such knowledge learned with confluent tasks.

By adopting a lean approach that dispenses with heavier feature and model engineering, this study permitted both to gain novel empirically based insights into the argument mining task and to establish new state of the art levels of performance for its three main sub-tasks, viz. identification of argument components, classification of the components, and determination of the relation among them.

## 1 Introduction

Argument mining is a Natural Language Processing (NLP) task consisting in taking unstructured text as input and returning it annotated such that each portion occurring in it that is an argument is properly delimited and analysed (Schneider et al., 2013; Peldszus and Stede, 2013; Lippi and Torroni, 2016; Habernal and Gurevych, 2017; Wachsmuth et al., 2017; Stede and Schneider, 2018; Lawrence and Reed, 2020). Argument mining relates to the high-level human capacity of reasoning (Walton et al., 2005), it is at the core of social interaction concerned with persuasion (Mercier and Sperber, 2017), and it is of utmost importance to enhance applications across different domains that aim at enhancing their services beyond mere sentiment analysis, on the basis of the reasons uncovered for

the associated sentiments and decisions (Habernal et al., 2014).

Argument mining has been decomposed into a number of sub-tasks. While the number and profiling of these tasks depends on the theoretical approach adopted to analyse arguments (Van Eemeren et al., 2019), they typically involve some sort of delimitation of the text segments conveying argument components, the classification of the roles of these components (e.g. premises, conclusions, etc.), and the classification of the type of relation among those components (e.g. support, attack, etc.) (Lawrence and Reed, 2020).

These sub-tasks and their eventual pipeline in argument mining have been addressed by means of supervised deep learning approaches that involve some degree of neural architecture engineering (Eger et al., 2017; Potash et al., 2017; Nguyen and Litman, 2016) a.o. Recently, first attempts to approach argument mining with Transformers have been reported in the literature (Wang et al., 2020; Rodrigues et al., 2020a) a.o., though at an exploratory level that leaves much of its strength still untapped.

This has been combined with experimentation with transfer learning (Caruana, 1997; Ruder, 2019). Given its complexity, and the associated difficulty in producing gold labelled data, argument mining is a task with a scarcity of data sets that are needed to support supervised learning approaches. Enhancing the argument mining task by transferring knowledge elicited when solving other natural language processing tasks is thus a promising approach to alleviate such scarceness. This has been tried in the literature (Mohammad et al., 2016; Stab et al., 2018; Choi and Lee, 2018; Habernal et al., 2018; Rodrigues and Branco, 2020) a.o., though at a haphazard level that leaves still much of its potential to be studied.

For humans, argumentation is a high level cognitive task that goes together with a number of other

capacities relating to linguistic syntactic and semantic processing, to entailment and paraphrasing, to question answering and language comprehension, to reasoning, to common sense, etc. (Lawrence and Reed, 2020; Lauscher et al., 2021). Interestingly, there is now available in the literature a wide range of data sets and respective NLP tasks that permit to address a wide range of these different dimensions and use them as auxiliary sources of knowledge in transfer learning approaches to argument mining (Wang et al., 2018, 2019a) a.o.

In this context, our goal is to empirically assess the potential of transfer learning to support argument mining by means of a systematic study with a wide range of possible sources of related tasks and knowledge possibly suitable to be transferred. In this paper we report on the findings of exploring a vast experimental space that results from: performing sequential single-step transfer learning from over 40 auxiliary tasks to each one of three main sub-tasks of argument mining (Stab and Gurevych, 2014, 2017) during the fine-tuning phase (Section 4); further explore the source tasks that supported the best single-step transfer learning by experimenting with ways of possibly combining them in multi-step transfer learning processes, and further explore these tasks in a multi-task transfer learning setting (Section 5). This is preceded by an overview of related work (Section 2) and by the presentation of the experimental setup adopted (Section 3).

By undertaking this study, not only new state-of-the-art results were achieved for argument mining, as also new empirically based insights were gained on how this task can be enhanced, showing the effectiveness of transfer learning to leverage argument mining and to alleviate its data scarcity when combined with a lean approach that dispenses with heavier feature and model engineering.

## 2 Related work

Transfer learning is a technique in machine learning that leverages knowledge from other, so called source tasks to improve the learning of a target task (Caruana, 1997), being a methodology to alleviate the lack of labelled data for the latter (Ruder, 2019).

### 2.1 Transfer learning for argument mining

Four families of approaches of transfer learning for argument mining have been reported in the

literature: (i) transfer learning across discourse domains for the same argument mining sub-task; (ii) cross-lingual transfer learning for a given sub-task; (iii) multi-task learning among argument mining sub-tasks; and (iv) sequential transfer learning from sources tasks that are not argument mining sub-tasks. A brief overview follows below.

Several papers have applied transfer learning with a **domain adaptation** approach for identifying components and clausal properties (Al-Khatib et al., 2016; Ajjour et al., 2017; Daxenberger et al., 2017). Typically, a model is trained with data sets from various discourse domains and is evaluated over each domain.

**Cross-lingual** transfer learning for argument mining (Aker and Zhang, 2017; Sliwa et al., 2018; Eger et al., 2018; Rocha et al., 2018) is mainly performed through direct transfer (McDonald et al., 2011) or projection (David et al., 2001) techniques. Direct transfer techniques train a model with the source language data that initializes a new model for a target language, typically with less to no data. Projection techniques resort to mapping the same labels from the source language data set to a target language data set by resorting to parallel corpora.

The argument mining pipeline has been addressed also with transfer learning by **multi-task** and **sequential** approaches (Cabrio and Villata, 2013; Peldszus and Stede, 2015; Eger et al., 2017; Potash et al., 2017; Niculae et al., 2017; Galassi et al., 2018; Schulz et al., 2018; Mensonides et al., 2019; Chakrabarty et al., 2019; Accuosto and Saggion, 2019; Cheng et al., 2020). Most proposals train models pipelining the sub-tasks in some way.

Transfer learning from **related tasks** has also been shown to improve the performance of argument mining sub-tasks. (Stab et al., 2018) transferred shared knowledge from two different tasks: a stance detection task (Mohammad et al., 2016) and a topic identification task. (Choi and Lee, 2018), in turn, transferred knowledge from the Argument Reasoning Comprehension Task (Habernal et al., 2018) for a clausal classification sub-task.

### 2.2 Main sub-tasks

To proceed with our systematic study of transfer learning for argument mining on a mainstream pipeline of sub-tasks (Lawrence and Reed, 2020), which includes identifying argument components, classifying their clausal roles and determining the relational properties among them, we resorted to

the AAEC corpus (Stab and Gurevych, 2014, 2017), a collection of annotated essays in English, which has been subject to various studies. An example from this data set is displayed in Figure 1.

In order to further support this option, it is worth noting that there is not in the literature a set of commonly agreed standard argument mining sub-tasks and that persuasive arguments, contained in the AAEC corpus, are by no means peripheral to argumentation, which is ultimately about persuasion. It is also worth noting that, while in NLP in general, it is always better to have more data sets/tasks for evaluation, the empirical study in this paper builds on a strong series of recent investigations that are based on one of the few data sets for argument mining, the AAEC, that given its quality and volume, has permitted comparison of results and the objective assessment of possible advances.

Title: Children should grow up in a big city!  
 Essay: It's certainly better for children to grow up in a big city<sup>1</sup>. Of course you need to choose a good neighborhood. I hold this belief because of two main reasons, academic and social reasons<sup>2</sup>.  
 Some people think that if a child grows up in a big city they will be all day at home at the computer or at the video-game<sup>3</sup>, but this is not true if you live in a neighborhood with other people about your age as I did<sup>4</sup>. My friends and I used to play soccer, bike, climb trees and do a lot of other stuff every day<sup>5</sup>. We did play video-games, but that wasn't our main activity<sup>6</sup>. In a big city there are more kinds of people and more things to do<sup>7</sup>.  
 I have a friend that grew up in the countryside<sup>8</sup>. He said that he had to study a lot to pass the test to enter the university<sup>9</sup>. This is another downside of growing up in the countryside. In a big city you have more qualified teachers and a better access to technology<sup>10</sup>.  
 Growing up in the countryside is not such a good experience<sup>11</sup>, you won't know a lot of people, there are gossips everywhere, and your life will be really limited<sup>12</sup>. If someday I have children, I'm absolutely sure that they will grow up in a good neighborhood of a big city and they will be very happy about it<sup>13</sup>.  
 Labels: Major Claim / Claim / Premise  
 Relations: Support (3→4; 13→11; 12→11) Attacks (7→6; 8→6; 9→6; 10→6)

Figure 1: Example of a labelled essay in AAEC.

The AAEC corpus integrates the annotation of every sub-task in the argument mining pipeline into a single data set. It contains 402 manually annotated essays,<sup>1</sup> in English, with 7,116 sentences over 1,833 paragraphs spanning 147,271 tokens.

It adopts an argument structure model in the form of a tree composed of major claim (in the root node, as the author's standpoint on the argument topic), claims and premises. Individual paragraphs of the essay include arguments that may be linked or not-linked (via relational properties) to the author's major claim. Both "support" and "attack" relations are taken into account.

The annotation of text segments with argument components resorted to an IOB tagging scheme (Ramshaw and Marcus, 1999). The beginning of an argument component is tagged with *Arg-B*, the following tokens in that component are

<sup>1</sup>80 essays, i.e 20% for testing, were annotated by three annotators and the remaining 322, for training, by an expert.

tagged with *Arg-I* and non-argumentative tokens with *O*. Identifying argument components consists of tagging each token with this IOB-tagset given a complete essay as a single input sequence. Identifying clausal properties consists of classifying spans of discourse with one of the three classes (major claim, claim and premise) given an entire essay as input. Following the literature, given the large imbalance between "support" and "attack" classes, identifying relational properties consists in classifying pairs of segments just as linked or not-linked. Statistics are displayed in Table 1.<sup>2</sup>

Task	Labels	Total	Train	Test
Comp.	Arg-B	11%	6,089	79%
	Arg-I	64%	93,618	80%
	O	25%	47,474	80%
Clausal	Major Cl	12%	751	80%
	Claim	25%	1,506	80%
	Premise	63%	3,832	79%
Relat.	Not-Link	82%	18,340	78%
	Linked	18%	3,832	79%

Table 1: For the tasks annotated in AAEC (rows), the number of instances for labels and data set split (columns) are indicated.

### 2.3 Literature on the AAEC tasks

Several papers on argument mining address the AAEC tasks, although none addresses all of them, except (Stab and Gurevych, 2017), which addressed each task with a feature-engineered SVM (components: 0.849 macro-F1; clausal: 0.773; relational: 0.736), and an Integer Linear Programming (ILP) algorithm (0.867, 0.826, 0.751 respectively), that is an ensemble of the SVM models supplemented by rules to ensure the correct tree structure. Table 2 presents the performance scores reported in the literature for the AAEC tasks.

Regarding the identification of argument **components** task: (Ajjour et al., 2017) implement a BiLSTM with extensive use of features and obtain 0.885 macro-F1. (Petasis, 2019) applies several types of neural networks for segmentation, with the top-performing model, a BiLSTM-CRF, obtaining 0.901 macro-F1. (Spliethöver et al., 2019) resorts to attention mechanisms with BiLSTMs for unit segmentation, with the top-performing model obtaining 0.87 weighted-F1. (Eger et al., 2017) apply different models, including multi-task learning experiments, and report 0.908 macro-F1 for the identification of components sub-task.

<sup>2</sup>Further descriptions of the data set and the framing of the tasks are provided in the Appendix A.

	Comp.	Clau.	Rel.
SVMs (Stab and Gurevych, 2017)	.849	.773	.736
ILP (Stab and Gurevych, 2017)	.867	.826	.751
S2S (Potash et al., 2017)		<b>.849</b>	<b>.767</b>
BL (Ajour et al., 2017)	.885		
BL (Eger et al., 2017)	<b>.908</b>		
BL (Spliethöver et al., 2019)	.870		
BL-CRF (Petasis, 2019)	.901		
BL-CRF (Schulz et al., 2018)		.606	
BL-CNN-CRF (Chernodub et al., 2019)		.471	
CNN-Seq. (Gemechu and Reed, 2019)		.790	
BERT (Wang et al., 2020)		.640	
LibLINEAR (Nguyen and Litman, 2016)			.753

Table 2: **Comparison of different performance scores in the literature on the AAEC tasks**, in macro-F1 (except weighted-F1 in (Spliethöver et al., 2019)), with the top results in bold, indicating the state-of-the-art (BL stands for BiLSTM). It should be noted that LibLINEAR uses the first version of the AAEC data set.

For the identification of **clausal** properties task: (Gemechu and Reed, 2019) obtain 0.79 macro-F1 for clausal properties linking premises and conclusions, taking into account the similarity of target concepts and aspects. (Chernodub et al., 2019) applied a framework for tagging arguments and their retrieval, including a BiLSTM-CNN-CRF sequence tagger. A micro-F1 of 0.645 was the top-performing performance in identifying clausal properties (0.471 macro-F1 is the reproduction in (Wang et al., 2020)). (Wang et al., 2020) propose a multi-scale mining model, resorting to several encoder-only Transformers (BERT) that mine different argumentation components at different textual levels, namely at the essay/paragraph/word-level. The top-performing model obtains 0.64 macro-F1 in identifying clausal properties. (Schulz et al., 2018) also apply a multi-task learning approach from different domains and argumentative structures, including AAEC, with a BiLSTM-CRF, obtaining 0.606 macro-F1 score.

Finally, as for **relational** properties: (Nguyen and Litman, 2016) obtain 0.753 macro-F1 combining different topic to window context features with a linear classifier (LibLINEAR). (Potash et al., 2017) report a 0.849 clausal and 0.767 relational macro-F1 using a joint pointer architecture (sequence-to-sequence model with attention), simultaneously addressing clausal and relational properties with several features.

### 3 Experimental space and settings

For the tasks that are the source of knowledge to be transferred to argument mining models, we resorted to a vast array of annotated data sets listed in Table

3. They cover different dimensions in terms of linguistic and cognitive processing:<sup>3</sup>

#### 3.1 Source tasks

**Syntax** - Information on syntax is typically included in structured machine learning algorithms that address the argument mining in a feature engineering approach. We included part-of-speech (POS) tagging, named entity recognition (NER) (Hu et al., 2020) and several other tasks regarding linguistic properties of sentences (Conneau and Kiela, 2018).

**Semantics** - Features from semantic similarity (SS) are widely used in argument mining literature. For example, (Boltužić and Šnajder, 2015) use SS to identify prominent arguments in online debates, and (Lawrence and Reed, 2015) use SS obtained from WordNet to identify the components of argumentation schemes. We included a diversity of SS data sets, from the context-sensitive similarity task Wic (Pilehvar and Camacho-Collados, 2019) to the large data set obtained from Quora Question Pairs (QQP) (Iyer et al., 2017).

**Grammaticality** - To address the widest spectrum of linguistic aspects, we included also tasks on determining the grammaticality of input sentences. Data sets such as the Corpus of Linguistic Acceptability (CoLA) (Warstadt et al., 2019) were used, that are challenging with regards this type of task.

**Sentiment** - Sentiment analysis has a certain proximity to argument mining, which adds an extra dimension to it by providing reasons for sentiments (Habernal et al., 2014). The Stanford Sentiment Treebank (SST) (Socher et al., 2013) was included.

**Reasoning & Comprehension** - Reasoning is at the core of argumentation given it is crucial in formulating and accepting or rejecting an argument. We included several related tasks, as for instance the AI2 Reasoning Challenge (ARC) (Clark et al., 2018) in the domain of grade-school science.

**Question Answering & Common sense** - Question Answering (QA) relates to argument mining given linguistic similarities between the Question/Answer and Claim/Premise pairs. Several QA

<sup>3</sup>We resorted also partly to PORTULAN CLARIN workbench consisting of language processing services: (Gomes et al., 2018; Branco et al., 2020; Barreto et al., 2006; Branco et al., 2010; Cruz et al., 2018; Veiga et al., 2011; Branco and Henriques, 2003; Branco and Silva, 2003; Branco et al., 2011, 2022; Branco and Nunes, 2012; Silva et al., 2010; Branco et al., 2014b; Silveira and Branco, 2012a,b; Branco and Costa, 2008; Branco et al., 2014a; Rodrigues et al., 2016; Branco and Silva, 2006; Rodrigues et al., 2020b; Costa and Branco, 2012; Santos et al., 2019; Neale et al., 2016; Miranda et al., 2011).

tasks were included that address common sense as this is closely related to argumentation, given that several implicit premises, tacit assumptions or inferences are to some extent regarded as common sense—for example, (Saint-Dizier, 2017) uses QA techniques for argument mining.

**Entailment & Paraphrase** - Although argument mining and Textual Entailment (TE) are different tasks, they are closely related given the similarity between specific entailment properties and argument clausal and relational properties. Works such as (Cabrio and Villata, 2012; Cocarascu and Toni, 2017) use models for TE to address argument relational properties. We included several TE tasks in different discourse domains, such as news and forums, with STSB (Cer et al., 2017), and science, with SciTAIL (Khot et al., 2018).

**Argument mining** - In addition to non argument mining tasks, we considered also as a source task for transfer learning the predecessor sub-task in the argument mining pipeline, that is the identification of components (for the clausal sub-task) and the clausal classification (for the relational sub-task).

### 3.2 Computational models

In order to explore the experimental space setup for our study, we resorted to the Transformer architecture (Vaswani et al., 2017), which became mainstream in NLP, surpassing several state-of-the-art results in a wide range of tasks of all sorts (Wang et al., 2018, 2019a). In contrast to most literature on argument mining, where structured feature engineering has been the favoured approach, a Transformer is a deep learning approach that obtains linguistic knowledge by transfer learning from a language modelling task.

In order to factorize out the impact of different possible models and obtain results that can be comparable across the different data points in our experimental space, we adopt the same type of model for all of them. Taking a look at a task closely related to argument mining, namely common sense reasoning, there are works in the literature (Branco et al., 2021) that, for this task, under comparable circumstance, have experimented with prominent exemplars of encoder-only, decoder-only, encoder-decoder, and neuro-symbolic types of Transformers, which found that RoBERTa (Liu et al., 2019) offers a clear advantage. Inspired by these results, we undertook an exploratory study, repeating the above experiments but now for sample cases of ar-

Task	#Train
<i>Syntax</i>	
PANX (Hu et al., 2020)	20K
UDPOS (Hu et al., 2020)	21K
Bigram Shift (Conneau and Kiela, 2018)	100K
Coord Inversion (Conneau and Kiela, 2018)	100K
Obj number (Conneau and Kiela, 2018)	100K
Odd Man Out (Conneau and Kiela, 2018)	100K
Past-Present (Conneau and Kiela, 2018)	100K
Sentence Length (Conneau and Kiela, 2018)	100K
Subj Number (Conneau and Kiela, 2018)	100K
Top Constituents (Conneau and Kiela, 2018)	100K
Tree Depth (Conneau and Kiela, 2018)	100K
Word Content (Conneau and Kiela, 2018)	100K
<i>Semantics</i>	
COPA (Roemmele et al., 2011)	400
WIC (Pilehvar and Camacho-Collados, 2019)	5.4K
STSB (Cer et al., 2017)	7K
QQP (Iyer et al., 2017)	364K
<i>Grammaticality</i>	
Coord (White et al., 2020)	458
Eos (White et al., 2020)	479
Definiteness (White et al., 2020)	508
Whwords (White et al., 2020)	585
CoLA (Warstadt et al., 2019)	8.5K
<i>Sentiment</i>	
SST (Socher et al., 2013)	67K
<i>Reasoning &amp; Comprehension</i>	
MULTIRC (Khashabi et al., 2018)	456
WNLI (Levesque et al., 2012)	635
ARC (Clark et al., 2018)	2.2K
ROPES (Lin et al., 2019)	10K
ANLI (Bhagavatula et al., 2020)	169.6K
FEVER (Nie et al., 2019)	208.3K
<i>Question Answering &amp; Common sense</i>	
WSC (Levesque et al., 2012)	554
CommonsenseQA (Talmor et al., 2019)	9.7K
QUAIL (Rogers et al., 2020)	10.2K
BoolQ (Clark et al., 2019)	16K
PIQA (Bisk et al., 2020)	16.1K
CosmosQA (Huang et al., 2019)	25K
HellaSwag (Zellers et al., 2019)	39.9K
MRQA (Fisch et al., 2019)	104K
QNLI (Wang et al., 2018)	105K
<i>Entailment/Paraphrase</i>	
CB (De Marneffe et al., 2019)	1.2K
RTE (Dagan et al., 2005)	2.5K
MRPC (Dolan and Brockett, 2005)	3.7K
SciTAIL (Khot et al., 2018)	27K
MNLI (Williams et al., 2018)	393K
<i>Argument mining</i>	
Components (Stab and Gurevych, 2017)	117k
Clausal (Stab and Gurevych, 2017)	4k

Table 3: **Data sets used for source tasks**, clustered by linguistic and cognitive dimensions.

gument mining from our experimental space and arrived at the same finding. Accordingly, and given also its accessible compute requirements and top performance in several NLP tasks, we adopted the off-the-shelf RoBERTa model, resorting to RoBERTa-large variant only when the RoBERTa-base was shown not to be enough to beat the SoTA.

We used the Jiant framework (Wang et al., 2019b; Phang et al., 2020) and Huggingface (Wolf

et al., 2020). The training objective for the pre-training model was the Mask Language Modelling (MLM), which randomly masks a word in a sentence and predicts it.

To identify argument components, a token classification head classifies the input sequence  $x_{1:N}$  (full essay) and gives a possible output  $y_{1:N}$  from a class set  $C$ . To identify clausal and relational properties, a sequence classification head classifies each input sequence  $x_{1:N}$  and gives a possible output  $y$  from a class set  $C$ .

### 3.3 Baselines

As for the baselines, we included the class majority, and the scores of a RoBERTa-base model fine-tuned for each AAEC task. We also included the SVMs and ILP model from (Stab and Gurevych, 2017) as a strong baseline.

### 3.4 Evaluation

For the evaluation of the transfer learning, we used the final result of each main sub-task in argument mining, which is the mean score of three runs. As in the original AAEC work and given that classes are unbalanced, for all tasks we used a macro-F1 averaging (Sokolova and Lapalme, 2009). We applied the Independent Samples  $t$ -Test regarding the RoBERTa baseline and different data points obtained in our experimental space to evaluate the statistical significance (Dror et al., 2018).

## 4 Single-step transfer

A first batch of experiments was concerned with single-step sequential transfer learning where the source tasks were those listed in Table 3.

Given the large number of data points in this experimental space, concessions were made considering the compute footprint, and we limited the hyper-parameter search by using the recommended values (Liu et al., 2019; Wolf et al., 2020).<sup>4</sup>

<sup>4</sup>Inspired by the STILT approach (Phang et al., 2018) we adopted the jiant toolkit (Pruksachatkun et al., 2020), an open source toolkit for transfer learning experiments.

For the fine-tuning of the target tasks, we performed a hyper-parameter search with three learning rates and three seeds on the target task development set, creating a total of 396 models.

The AAEC development set was extracted from 10% of the original training data, thus the training data consists of the remaining 90%. Based on the top-performing result obtained from the development set, hyper-parameters were determined for the test set. Further descriptions of hyper-parameterization together with all materials to reproduce the experiments are available at <https://github.com/nlx-group/transfer-am>.

	Comp.	Clausal	Relational
Human	.886	.868	.854
SoTA - Table 2	.908	.849	.767
<i>Baselines</i>			
RoBERTa no transfer	.916	.820	.727
ILP	.867	.826	.751
SVM	.849	.773	.736
Majority	.259	.257	.455
<i>Syntax</i>			
PANX	.917	.815	.756
UDPOS	.914	.804	.743
Bigram Shift	.912	.710	.743
Coord Inversion	.910	.696	.735
Obj number	.907	.715	.729
Odd Man Out	.914	.703	.752
Past-Present	.901	.713	.718
Sentence Length	.885	.652	.466
Subj Number	.913	.707	.746
Top Constituents	.896	.708	<b>.762*</b>
Tree Depth	.904	.674	.735
Word Content	.896	.713	.455
<i>Semantics</i>			
COPA	.919*	.823	.738
WIC	.918	.821	.744
STSB	.917	.805	.753
QQP	.911	.800	.746
<i>Grammaticality</i>			
Coord	.910	.722	.754*
Eos	.914	.712	.745
Definiteness	.914	.705	.755
Whwords	.915	.702	<b>.758</b>
CoLA	<b>.924</b>	.713	.752*
<i>Sentiment</i>			
SST	.916	.820	.747*
<i>Reasoning &amp; Compreh</i>			
MULTIRC	.919	<b>.831</b>	<b>.758</b>
WNLI	.913	.788	.455
ARC	<b>.921</b>	.820	<b>.758</b>
ROPES	.920	.806	.748
ANLI	.917	.807	.749
FEVER	.914	.814	.736
<i>QA &amp; Common sense</i>			
WSC	.919	.820	<b>.758</b>
CommonsenseQA	.916	.819	.755*
QUAIL	<b>.921</b>	.827	.755*
BoolQ	.916	<b>.837</b>	.742
PIQA	.914	.774	.455
CosmosQA	.917	.817	.745
HellaSwag	.916	.823	.746
MRQA	<b>.924</b>	.825	.750
QNLI	.916	.826	.751
<i>Entailment/Paraphrase</i>			
CB	<b>.923*</b>	.819	.734
RTE	.916	<b>.843*</b>	<b>.757</b>
MRPC	.916	.790	.746
SciTAIL	.919	.827	.751*
MNLI	.919	.812	.731
<i>Argument mining</i>			
Components		<b>.843</b>	.664
Clausal			.657

Table 4: Performance in macro-F1 scores on the main sub-tasks (columns) by different source tasks (rows). Top score underlined, top 3 scores in bold, average score in the same family of tasks in italics. All values found to be statistical significant ( $p$ -value  $< .05$ ) are noted with an \*

## 4.1 Results and Analysis

Table 4 shows the results from this first batch of experiments,<sup>5</sup> which support the following major empirical findings:

- **The Transformer with no transfer is a very strong baseline** (off-the-self RoBERTa-base fine-tuned to each AAEC task). It overcomes (with 0.916 in components) the SoTA (0.908) of one of the three main tasks, and has strong scores in the other two.

- **Transfer learning is effective to leverage argument mining.** This is supported by scores above the Transformer baseline: with 0.924 (against the baseline 0.916) in the components task; 0.843 (against 0.820) in the clausal task; and 0.762 (against 0.727) in the relational task.

- **Transfer learning with a Transformer is very competitive with respect to, or even surpass, the SoTA.** This is supported by a new SoTA of 0.924 in components (against 0.908), and by very good scores, 0.843 and 0.762, against respectively 0.849 and 0.767, in clausal and relational.

- **Source tasks whose overall cognitive complexity is high and closer to the argument mining task tend to be more successful in supporting effective transfer.** The overall trend is that better results are found with source tasks for Reasoning, Common sense and Entailment, as shown by the respective averages and the larger number of top scores therein. Interestingly, the top score of 0.762 for relational is obtained with a syntactic source task, that seeks to identify Top Constituents: this is of relevance for the relational task as this is about relating clausal segments, which are univocally associated with their top constituents.

- **A main sub-task can be a good source task to other sub-task for effective transfer.** This is supported by the top score 0.843 in the clausal task when the components was the source in transfer.

- **A larger size of a data set for a source task, in contrast to other sources tasks, does not necessarily lead to an enhanced performance of the transfer chain.** This is illustrated, for instance, by the case of RTE, with a small data set of only 2.5K, but with the top score for clausal.

## 5 Multi-step and multi-task transfer

A second batch of experiments was concerned with multi-step and multi-task transfer learning. The

<sup>5</sup>All scores were obtained with RoBERTa-base.

source tasks considered here were the ones with the best results in the previous batch of experiments with single-step transfer.

Hence, **two-step** transfer was experimented with, where the typical chain encompasses the transfer from the components task to the clausal task and from the latter to the relational task. But we experimented also with other two-step instances, where the initial source tasks in the chain, viz. RTE, CB and Top Constituents (TC), are none of the argument mining sub-tasks. Experiments with **three-step** transfer were also undertaken, where besides the main tasks, these other source tasks contributed to the chain.

Finally, besides sequential transfer, also **multi-task** transfer learning was experimented with, involving the three argument mining sub-tasks altogether, and also pairs including two of them. Motivated by these pairings of the sub-tasks, we returned to one-step methodology, and for the sake of completeness, we experimented also with every combination of two such sub-tasks.

	Comp.	Clausal	Relational
Human	.886	.868	.854
SoTA Table 2	.908	.849	.767
<i>Baselines</i>			
RoBERTa no transfer	.916	.820	.727
ILP	.867	.826	.751
SVM	.849	.773	.736
Majority	.259	.257	.455
<i>Sequential</i>			
Cl ⇒ Cp	.920		
Re ⇒ Cp	<b>.924</b>		
RTE ⇒ Cp	.916		
Re ⇒ Cl ⇒ Cp	.912		
CB ⇒ Re ⇒ Cp	.915		
Cp ⇒ Cl		.843*	
Re ⇒ Cl		.811	
RTE ⇒ Cl		.843*	
Re ⇒ Cp ⇒ Cl		.839	
RTE ⇒ Cp ⇒ Cl		<b>.888*</b>	
Cp ⇒ Re			.664
Cl ⇒ Re			.657
RTE ⇒ Re			.757
Cp ⇒ Cl ⇒ Re			.781*
RTE ⇒ Cp ⇒ Cl ⇒ Re			<b>.783*</b>
TC ⇒ Cp ⇒ Cl ⇒ Re			.761
<i>Multi-task</i>			
Cp ⇔ Cl	.915	.813	
Cp ⇔ Re	.911		.684
Cl ⇔ Re		.738	.714
Cp ⇔ Cl ⇔ Re	.906	.796	.757

Table 5: **Performance on the three main sub-tasks (columns) by different transfer learning source tasks and their chaining (rows)**, reported with macro-F1, with the top results in bold, indicating new state-of-the-art scores. Cp stands for Components, Cl for Clausal, Re for Relational and TC for Top Constituents.

## 5.1 Results and Analysis

Table 5 presents the results for this second batch of experiments,<sup>6</sup> supporting these major findings:

– **Sequential transfer is more effective than multi-task transfer.** This is supported by the overall stronger scores in sequential transfer experiments for similar clusters of tasks.

– **Multi-step transfer can be more effective than single-step.** This is supported by the results obtained for the relational task: with the best score to relational in all experimental space of 0.783, this result was supported by a three step transfer that leveraged the relational task with the knowledge from the other two main tasks, components and clausal, and from RTE; and it is supported also by the results obtained for the clausal task: with the best score in all experimental space of 0.888, this result was supported by a two step transfer that leveraged the clausal task with the knowledge from other two tasks, one from the entailment (RTE) and the other being another main task (components).

– **Source tasks that are sub-tasks in the argument mining pipeline are very successful in enhancing effective transfer.** This is supported by the results obtained with the transfer being organized along the default argument mining pipeline direction, with top or very close to the top second scores for the chains  $C_p \Rightarrow C_l$  and  $C_p \Rightarrow C_l \Rightarrow R_e$ , with 0.843 and 0.781, respectively. But this is supported by the results obtained with the transfer being organized also in different directions, like for instance, the best score to components in all experimental space, of 0.924, with  $R_e \Rightarrow C_p$ .

– **Source tasks with the best performance for a given main task in the single-step setting are very successful in enhancing multi-step effective transfer, specially for that main task.** This is supported by the results obtained with top or very close to the top second scores for the chains  $RTE \Rightarrow C_p$ , with 0.916 (over the SoTA 0.908 for components),  $RTE \Rightarrow C_p \Rightarrow C_l$ , with 0.888 (top score for clausal, and over its SoTA 0.849), and  $RTE \Rightarrow C_p \Rightarrow C_l \Rightarrow R_e$ , with 0.774 (over the SoTA 0.767 for relational).

– **Transfer learning in the setting of an off-the-self Transformer architecture renders new SoTA scores for the argument mining tasks.** This is supported by the scores of 0.924 for components (against 0.908 in previous SoTA), 0.888 for clausal

<sup>6</sup>All scores obtained with RoBERTa-base except clausal  $RTE \Rightarrow C_p \Rightarrow C_l$ .

(against 0.849), and 0.783 for relational (against 0.767).

## 6 Further analysis

No correlation was found between the task scores and the size of their training data. Using the coefficient of determination, .101/.002 and .001  $R^2$  is obtained for identifying argument components, clausal and relational properties, respectively.

We performed a manual analysis of the output on top-performing tasks in the single-step transfer (CB, RTE, QUAIL). We notice that shorter arguments tend to be incorrectly tagged as *O* (outside) while more extensive arguments tend to be incorrectly divided into two arguments; also, some discourse markers introducing arguments, as "there is clear evidence that" or "thus it is apparent that", tend to be wrongly labelled as the beginning and inside of an argument segment.

Transfer learning experiments on clausal properties follow the same error pattern as the baseline, with most errors emerging from labelling major claims as claims, claims as premises and premises as claims. For relation identification, linked arguments were identified with higher precision and recall than the baseline.

Transferring knowledge from argument mining sources was examined also by extending the language modelling phase. Despite above-baseline scores, no statistical significance was found.<sup>7</sup>

## 7 Conclusions and future work

The results in this paper were obtained from a large experimental space that permitted a systematic empirical study aimed at assessing the viability of transfer learning to leverage neural argument mining with confluent knowledge. Major findings and results are:

- The knowledge transfer enabled by the transfer learning from language processing tasks that are confluent to argument mining is an effective approach to improve neural argument mining.

- Sequential transfer learning appears as more effective than multi-task transfer, and multi-step sequential transfer can achieve better performance than single-step.

- Source language processing tasks more closely related to argument mining and to the higher-level cognitive capacities mobilized for argumentation tend to provide better support.

<sup>7</sup>More details can be found in Appendix B.



- New state of the art levels of performance were established for the three main sub-tasks in argument mining, namely identification of argument components, classification of components, and determination of the relation among them.

- State of the art was obtained with a lean Transformer-based neural approach that dispensed with heavier feature and model engineering.

- There is much room for further improvements of performance in argument mining given that the new state of the art advanced in the present paper was possible even when deployed on top of just an off-the-shelf Transformer model, viz. RoBERTa,

Concomitantly, these advances open the way to future work. On the side of the mere race for brute force improvement of the state of the art levels of performance, resorting to available Transformer language models that are larger and more powerful than RoBERTa, which was used here, can be easily explored.

On the side of empirically motivated improvements based on more thoughtful approaches, it is possible to explore carefully articulated chains of transfer with curriculum and meta-learning, and also hybrid deep learning and symbolic approaches aimed to solve transfer learning catastrophic forgetting among other issues.

## Acknowledgments

The research reported here was supported by the individual PhD grant SFRH/BD/129824/2017 from FCT—Fundação para a Ciência e Tecnologia and it was also supported partially by PORTULAN CLARIN—Research Infrastructure for the Science and Technology of Language, funded by Lisboa 2020, Alentejo 2020 and FCT—Fundação para a Ciência e Tecnologia under the grant PINFRA/22117/2016.

## References

Pablo Accuosto and Horacio Saggion. 2019. *Transferring knowledge from discourse to arguments: A case study with scientific abstracts*. In *Proceedings of the 6th Workshop on Argument Mining*, pages 41–51, Florence, Italy. Association for Computational Linguistics.

Yamen Ajjour, Wei-Fan Chen, Johannes Kiesel, Henning Wachsmuth, and Benno Stein. 2017. Unit Segmentation of Argumentative Texts. In *Proceedings of the 4th Workshop on Argument Mining*, pages 118–128.

Yamen Ajjour, Henning Wachsmuth, Johannes Kiesel, Martin Potthast, Matthias Hagen, and Benno Stein. 2019. *Data Acquisition for Argument Search: The args.me corpus*. In *42nd German Conference on Artificial Intelligence (KI 2019)*, pages 48–59, Berlin Heidelberg New York. Springer.

Ahmet Aker and Huangpan Zhang. 2017. Projection of argumentative corpora from source to target languages. In *Proceedings of the 4th Workshop on Argument Mining*, pages 67–72.

Khalid Al-Khatib, Henning Wachsmuth, Matthias Hagen, Jonas Köhler, and Benno Stein. 2016. Cross-Domain Mining of Argumentative Text Through Distant Supervision. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 1395–1404.

Florbela Barreto, António Branco, Eduardo Ferreira, Amália Mendes, Maria Fernanda Bacelar do Nascimento, Filipe Nunes, and João Ricardo Silva. 2006. Open resources and tools for the shallow processing of Portuguese: the TagShare project. In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC 2006)*, pages 1438–1443.

Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. 2020. Abductive commonsense reasoning. In *ICLR*.

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*.

Filip Boltužić and Jan Šnajder. 2015. Identifying prominent arguments in online debates using semantic textual similarity. In *Proceedings of the 2nd Workshop on Argumentation Mining*, pages 110–115.

António Branco, Sérgio Castro, João Silva, and Francisco Costa. 2011. CINTIL DepBank handbook: Design options for the representation of grammatical dependencies. Technical Report DI-FCUL-TR-2011-03, University of Lisbon.

António Branco, Francisco Costa, João Silva, Sara Silveira, Sérgio Castro, Mariana Avelãs, Clara Pinto, and João Graça. 2010. Developing a deep linguistic databank supporting a collection of treebanks: the CINTIL DeepGramBank. In *Proceedings of the 7th International Conference on Language Resources and Evaluation (LREC 2010)*.

António Branco and Filipe Nunes. 2012. Verb analysis in a highly inflective language with an MFF algorithm. In *Proceedings of the 11th International Conference on the Computational Processing of Portuguese (PROPOR)*, number 7243 in Lecture Notes in Artificial Intelligence, pages 1–11. Springer.

- António Branco, João Rodrigues, João Silva, Francisco Costa, and Rui Vaz. 2014a. Assessing automatic text classification for interactive language learning. In *Proceedings of the IEEE International Conference on Information Society (iSociety)*, pages 72–80.
- António Branco and Francisco Costa. 2008. A computational grammar for deep linguistic processing of portuguese: Lxgram.
- António Branco and Tiago Henriques. 2003. Aspects of verbal inflection and lemmatization: Generalizations and algorithms. In *Proceedings of XVIII Annual Meeting of the Portuguese Association of Linguistics (APL)*, pages 201–210.
- António Branco, Amália Mendes, Paulo Quaresma, Luís Gomes, João Silva, and Andrea Teixeira. 2020. Infrastructure for the science and technology of language PORTULAN CLARIN. In *Proceedings, 1st International Workshop on Language Technology Platforms (IWLTP 2020)*, pages 1–7. European Language Resources Association (ELRA).
- António Branco, João Rodrigues, Francisco Costa, João Ricardo Silva, and Rui Vaz. 2014b. Rolling out text categorization for language learning assessment supported by language technology. In *Lecture Notes in Artificial Intelligence*, 8775, pages 256–261.
- António Branco and João Silva. 2003. Contractions: Breaking the tokenization-tagging circularity. In *Lecture Notes in Artificial Intelligence* 2721, pages 167–170.
- António Branco and João Silva. 2006. A suite of shallow processing tools for Portuguese: LX-Suite. In *Proceedings of the 11th European Chapter of the Association for Computational Linguistics (EACL)*, pages 179–182.
- António Branco, João Ricardo Silva, Luís Gomes, and João António Rodrigues. 2022. [Universal grammatical dependencies for portuguese with cintil data, lx processing and clarin support](#). In *Proceedings of the Language Resources and Evaluation Conference*, pages 5617–5626, Marseille, France. European Language Resources Association.
- Ruben Branco, António Branco, João António Rodrigues, and João Ricardo Silva. 2021. [Shortcutted commonsense: Data spuriousness in deep learning of commonsense reasoning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1504–1521, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- 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.
- Elena Cabrio and Serena Villata. 2013. A natural language bipolar argumentation approach to support users in online debate interactions. *Argument & Computation*, 4(3):209–230.
- Rich Caruana. 1997. Multitask Learning. *Machine learning*, 28(1):41–75.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. [SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation](#). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics.
- Tuhin Chakrabarty, Christopher Hidey, and Kathleen McKeown. 2019. Imho fine-tuning improves claim detection. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 558–563.
- Liyong Cheng, Lidong Bing, Qian Yu, Wei Lu, and Luo Si. 2020. Argument pair extraction from peer review and rebuttal via multi-task learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7000–7011.
- Artem Chernodub, Oleksiy Oliynyk, Philipp Heidenreich, Alexander Bondarenko, Matthias Hagen, Chris Biemann, and Alexander Panchenko. 2019. Targer: Neural argument mining at your fingertips. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 195–200.
- HongSeok Choi and Hyunju Lee. 2018. GIST at SemEval-2018 Task 12: A Network Transferring Inference Knowledge to Argument Reasoning Comprehension Task. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 773–777.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. [BoolQ: Exploring the surprising difficulty of natural yes/no questions](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- 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*, pages 1374–1379.
- Alexis Conneau and Douwe Kiela. 2018. Senteval: An evaluation toolkit for universal sentence representations. *arXiv preprint arXiv:1803.05449*.
- Francisco Costa and António Branco. 2012. Aspectual type and temporal relation classification. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 266–275.
- André Ferreira Cruz, Gil Rocha, and Henrique Lopes Cardoso. 2018. [Exploring Spanish corpora for Portuguese coreference resolution](#). In *Fifth International Conference on Social Networks Analysis, Management and Security (SNAMS)*, pages 290–295.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Machine Learning Challenges Workshop*, pages 177–190. Springer.
- Yarowsky David, Ngai Grace, Wicentowski Richard, et al. 2001. Inducing multilingual text analysis tools via robust projection across aligned corpora. In *Proceedings of the First International Conference on Human Language Technology Research*, pages 1–8.
- Johannes Daxenberger, Steffen Eger, Ivan Habernal, Christian Stab, and Iryna Gurevych. 2017. What is the essence of a claim? cross-domain claim identification. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2055–2066.
- Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank: Investigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*, volume 23, pages 107–124.
- William B. Dolan and Chris Brockett. 2005. [Automatically constructing a corpus of sentential paraphrases](#). In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*.
- Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. [The hitchhiker’s guide to testing statistical significance in natural language processing](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1383–1392, Melbourne, Australia. Association for Computational Linguistics.
- Steffen Eger, Johannes Daxenberger, and Iryna Gurevych. 2017. Neural End-To-End Learning for Computational Argumentation Mining. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 11–22.
- Steffen Eger, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2018. Cross-lingual argumentation mining: Machine translation (and a bit of projection) is all you need! In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 831–844.
- Adam Fisch, Alon Talmor, Robin Jia, Minjoon Seo, Eun-sol Choi, and Danqi Chen. 2019. MRQA 2019 shared task: Evaluating generalization in reading comprehension. In *Proceedings of 2nd Machine Reading for Reading Comprehension (MRQA) Workshop at EMNLP*.
- Andrea Galassi, Marco Lippi, and Paolo Torrioni. 2018. [Argumentative link prediction using residual networks and multi-objective learning](#). In *Proceedings of the 5th Workshop on Argument Mining*, pages 1–10, Brussels, Belgium. Association for Computational Linguistics.
- Debela Gemechu and Chris Reed. 2019. Decompositional argument mining: A general purpose approach for argument graph construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 516–526.
- Luís Gomes, Frederico Apolónia, Ruben Branco, João Ricardo Silva, and António Branco. 2018. Setting up the PORTULAN / CLARIN repository. In *CLARIN Annual Conference (CLARIN2018)*, pages 108–11.
- Ivan Habernal, Judith Eckle-Kohler, and Iryna Gurevych. 2014. Argumentation Mining On the Web From Information Seeking Perspective. In *ArgNLP*.
- Ivan Habernal and Iryna Gurevych. 2017. Argumentation Mining in User-Generated Web Discourse. *Computational Linguistics*, 43:125–179.
- Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2018. SemEval-2018 task 12: The argument reasoning comprehension task. In *Proceedings of the 12th International Workshop on Semantic Evaluation*, pages 763–772.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In *International Conference on Machine Learning*, pages 4411–4421. PMLR.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2391–2401.
- Shankar Iyer, Nikhil Dandekar, Kornél Csernai, et al. 2017. First quora dataset release: Question pairs. *data. quora. com*.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: a challenge set for reading comprehension over multiple sentences. In *NAACL*.

- Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018. Scitail: A textual entailment dataset from science question answering. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Anne Lauscher, Henning Wachsmuth, Iryna Gurevych, and Goran Glavaš. 2021. Scientia potentia est—on the role of knowledge in computational argumentation. *arXiv preprint arXiv:2107.00281*.
- John Lawrence and Chris Reed. 2015. Combining Argument Mining Techniques. In *Proceedings of the 2nd Workshop on Argumentation Mining*, pages 127–136.
- John Lawrence and Chris Reed. 2020. Argument mining: A survey. *Computational Linguistics*, 45(4):765–818.
- Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*.
- Kevin Lin, Oyvind Tafjord, Peter Clark, and Matt Gardner. 2019. Reasoning over paragraph effects in situations. In *MRQA@EMNLP*.
- Marco Lippi and Paolo Torroni. 2016. Argumentation Mining: State of the Art and Emerging Trends. *ACM Transactions on Internet Technology (TOIT)*, 16(2):10.
- 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*.
- Ryan McDonald, Slav Petrov, and Keith Hall. 2011. Multi-source transfer of delexicalized dependency parsers. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 62–72.
- Jean-Christophe Menonides, Sébastien Harispe, Jacky Montmain, and Véronique Thireau. 2019. Automatic detection and classification of argument components using multi-task deep neural network. In *Proceedings of the 3rd International Conference on Natural Language and Speech Processing*, pages 25–33, Trento, Italy. Association for Computational Linguistics.
- Hugo Mercier and Dan Sperber. 2017. *The enigma of reason*. Harvard University Press.
- Nuno Miranda, Ricardo Raminhos, Pedro Seabra, Joao Sequeira, Teresa Gonçalves, and Paulo Quaresma. 2011. Named entity recognition using machine learning techniques. In *Proceedings of the 15th Portuguese Conference on Artificial Intelligence (EPIA)*, pages 818–831.
- Saif Mohammad, Svetlana Kiritchenko, Parinaz Sobhani, Xiaodan Zhu, and Colin Cherry. 2016. SemEval-2016 Task 6: Detecting Stance in Tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 31–41.
- Steven Neale, Oier Lopez de Lacalle Luís Gomes, Eneko Agirre, and António Branco. 2016. Word sense-aware machine translation: Including senses as contextual features for improved translation models. In *Proceedings of the 10th Language Resources and Evaluation Conference (LREC2016)*, pages 2777–2783, Portorož, Slovenia. European Language Resources Association.
- Huy Nguyen and Diane Litman. 2016. Context-aware argumentative relation mining. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1127–1137, Berlin, Germany. Association for Computational Linguistics.
- Vlad Niculae, Joonsuk Park, and Claire Cardie. 2017. Argument mining with structured svms and rnns. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 985–995.
- Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks. In *Association for the Advancement of Artificial Intelligence (AAAI)*.
- Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. 2019. Asynchronous pipelines for processing huge corpora on medium to low resource infrastructures. Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-7) 2019. Cardiff, 22nd July 2019, pages 9 – 16, Mannheim. Leibniz-Institut für Deutsche Sprache.
- Andreas Peldszus and Manfred Stede. 2013. From Argument Diagrams to Argumentation Mining in Texts: A Survey. *International Journal of Cognitive Informatics and Natural Intelligence (IJCINI)*, 7(1):1–31.
- Andreas Peldszus and Manfred Stede. 2015. Joint Prediction in MST-Style Discourse Parsing for Argumentation Mining. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 938–948.
- Georgios Petasis. 2019. Segmentation of argumentative texts with contextualised word representations. In *Proceedings of the 6th Workshop on Argument Mining*, pages 1–10.
- Jason Phang, Thibault Févry, and Samuel R Bowman. 2018. Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks. *arXiv preprint arXiv:1811.01088*.
- Jason Phang, Phil Yeres, Jesse Swanson, Haokun Liu, Ian F. Tenney, Phu Mon Htut, Clara Vania, Alex Wang, and Samuel R. Bowman. 2020. jiant 2.0: A software toolkit for research on general-purpose text understanding models. <http://jiant.info/>.

- Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1267–1273.
- Peter Potash, Alexey Romanov, and Anna Rumshisky. 2017. Here’S My Point: Joint Pointer Architecture for Argument Mining. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1364–1373.
- Yada Pruksachatkun, Phil Yeres, Haokun Liu, Jason Phang, Phu Mon Htut, Alex Wang, Ian Tenney, and Samuel Bowman. 2020. jiant: A software toolkit for research on general-purpose text understanding models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 109–117.
- Lance A Ramshaw and Mitchell P Marcus. 1999. Text chunking using transformation-based learning. In *Natural language processing using very large corpora*, pages 157–176. Springer.
- Gil Rocha, Christian Stab, Henrique Lopes Cardoso, and Iryna Gurevych. 2018. Cross-lingual argumentative relation identification: from English to Portuguese. In *Proceedings of the 5th Workshop on Argument Mining*, pages 144–154.
- João Rodrigues and António Branco. 2020. Argument identification in a language without labeled data. In *International Conference on Computational Processing of the Portuguese Language*, pages 335–345. Springer.
- Joao Rodrigues, Ruben Branco, João Silva, and António Branco. 2020a. Reproduction and revival of the argument reasoning comprehension task. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 5055–5064.
- João Rodrigues, Francisco Costa, João Silva, and António Branco. 2020b. Automatic syllabification of Portuguese. *Revista da Associação Portuguesa de Linguística*, 1.
- João Rodrigues, António Branco, Steven Neale, and João Silva. 2016. LX-DSemVectors: Distributional semantics models for the Portuguese language. In *Proceedings of the 12th International Conference on the Computational Processing of Portuguese (PROPOR’16)*, pages 259–270.
- Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *2011 AAAI Spring Symposium Series*.
- Anna Rogers, Olga Kovaleva, Matthew Downey, and Anna Rumshisky. 2020. Getting closer to ai complete question answering: A set of prerequisite real tasks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8722–8731.
- Sebastian Ruder. 2019. *Neural transfer learning for natural language processing*. Ph.D. thesis, NUI Galway.
- Patrick Saint-Dizier. 2017. [Using question-answering techniques to implement a knowledge-driven argument mining approach](#). In *Proceedings of the 4th Workshop on Argument Mining*, pages 85–90, Copenhagen, Denmark. Association for Computational Linguistics.
- Rodrigo Santos, João Silva, António Branco, and Deyi Xiong. 2019. The direct path may not be the best: Portuguese-Chinese neural machine translation. In *Proceedings of the 19th Portuguese Conference on Artificial Intelligence (EPIA)*, pages 757–768.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.
- Jodi Schneider, Tudor Groza, and Alexandre Passant. 2013. A Review of Argumentation for the Social Semantic Web. *Semantic Web*, 4(2):159–218.
- Claudia Schulz, Steffen Eger, Johannes Daxenberger, Tobias Kahse, and Iryna Gurevych. 2018. Multi-Task Learning for Argumentation Mining in Low-Resource Settings. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 35–41.
- João Silva, António Branco, Sérgio Castro, and Rúben Reis. 2010. Out-of-the-box robust parsing of Portuguese. In *Lecture Notes in Artificial Intelligence, 6001*, pages 75–85.
- Sara Silveira and António Branco. 2012a. Combining a double clustering approach with sentence simplification to produce highly informative multi-document summaries. In *Proceedings, IEEE IRI2012-International Conference on Information Reuse and Integration*, pages 482–489.
- Sara Silveira and António Branco. 2012b. Enhancing multi-document summaries with sentence simplification. In *Proceedings, ICAI2012-14th International Conference on Artificial Intelligence*, pages 742–748.
- Alfred Sliwa, Yuan Ma, Ruishen Liu, Niravkumar Borad, Seyedeh Ziyaei, Mina Ghobadi, Firas Sabbah, and Ahmet Aker. 2018. Multi-lingual argumentative corpora in english, turkish, greek, albanian, croatian, serbian, macedonian, bulgarian, romanian and arabic. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for

- semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.
- Marina Sokolova and Guy Lapalme. 2009. A systematic analysis of performance measures for classification tasks. *Information processing & management*, 45(4):427–437.
- Maximilian Spliethöver, Jonas Klaff, and Hendrik Heuer. 2019. Is it worth the attention? a comparative evaluation of attention layers for argument unit segmentation. In *Proceedings of the 6th Workshop on Argument Mining*, pages 74–82.
- Christian Stab and Iryna Gurevych. 2014. Annotating Argument Components and Relations in Persuasive Essays. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1501–1510.
- Christian Stab and Iryna Gurevych. 2017. Parsing Argumentation Structures in Persuasive Essays. *Computational Linguistics*, 43:619–659.
- Christian Stab, Tristan Miller, Benjamin Schiller, Pranav Rai, and Iryna Gurevych. 2018. Cross-topic argument mining from heterogeneous sources using attention-based neural networks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3664–3674.
- Manfred Stede and Jodi Schneider. 2018. Argumentation Mining. *Synthesis Lectures on Human Language Technologies*, 11(2):1–191.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158.
- Frans H Van Eemeren, Rob Grootendorst, and Tjark Kruiger. 2019. *Handbook of argumentation theory*. De Gruyter Mouton.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Arlindo Veiga, Sara Candeias, and Fernando Perdigão. 2011. Generating a pronunciation dictionary for European Portuguese using a joint-sequence model with embedded stress assignment. In *Proceedings of the 8th Brazilian Symposium in Information and Human Language Technology*.
- Henning Wachsmuth, Nona Naderi, Yufang Hou, Yonatan Bilu, Vinodkumar Prabhakaran, Tim Alberdingk Thijm, Graeme Hirst, and Benno Stein. 2017. Computational Argumentation Quality Assessment in Natural Language. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 176–187.
- Douglas Walton, Christopher W. Tindale, and David Zarefsky. 2005. *Critical Reasoning and Argumentation*. Cambridge University Press.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2019a. Superglue: a stickier benchmark for general-purpose language understanding systems. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 3266–3280.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355.
- Alex Wang, Ian F. Tenney, Yada Pruksachatkun, Phil Yeres, Jason Phang, Haokun Liu, Phu Mon Htut, Katherin Yu, Jan Hula, Patrick Xia, Raghu Pappagari, Shuning Jin, R. Thomas McCoy, Roma Patel, Yinghui Huang, Edouard Grave, Najoung Kim, Thibault Févry, Berlin Chen, Nikita Nangia, Anhad Mohananey, Katharina Kann, Shikha Bordia, Nicolas Patry, David Benton, Ellie Pavlick, and Samuel R. Bowman. 2019b. jiant 1.3: A software toolkit for research on general-purpose text understanding models. <http://jiant.info/>.
- Hao Wang, Zhen Huang, Yong Dou, and Yu Hong. 2020. Argumentation mining on essays at multi scales. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5480–5493.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. *Neural network acceptability judgments*. *Transactions of the Association for Computational Linguistics*, 7:625–641.
- Aaron Steven White, Elias Stengel-Eskin, Siddharth Vashishtha, Venkata Subrahmanyam Govindarajan, Dee Ann Reisinger, Tim Vieira, Keisuke Sakaguchi, Sheng Zhang, Francis Ferraro, Rachel Rudinger, et al. 2020. The universal decompositional semantics dataset and decomp toolkit. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 5698–5707.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. *A broad-coverage challenge corpus for sentence understanding through inference*. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800.

## A AAEC Data set

The AAEC corpus (Stab and Gurevych, 2014, 2017) is a freely available collection of annotated essays extensively adopted by the argument mining community. For the transformation of the original annotations to the different tasks data sets inputs/outputs, we followed the original work (Stab and Gurevych, 2017).

In the literature, the three tasks, argument component identification, clausal properties and relation properties, typically follow the original frame.

**Identifying argument components** is performed by tagging each token with their tag from the IOB-tagset given a complete essay as an input sequence.

**Identifying clausal properties** is performed by individually classifying a span of components of an argument with one of the three classes (major claim, claim and premise) given the entire essay in context. In this task, the IOB-tags are not provided, and the span of components consists of raw text. As input to the model, we separated the span of components and the full essay with a separator token, for example, *components\_span <S> essay </S>*.

**Identifying relational** properties, in turn, is performed by individually classifying two components spans as linked or not-linked among themselves, given the entire essay as context. In this task, no IOB-tagset or clausal properties are provided. The spans consist of raw text. As input to the model, we separated the spans with separator tokens, like in the previous sub-task.

The three tasks are handled separately during training. There was no overlap of the test sets with

the training or development data sets. In the literature, some papers use the entire essay while some others only the paragraph as context for determining the clausal properties and the relation properties. We followed the typical approach described above for all base models, that is, we used the entire essay as context for RoBERTa-base models and only the paragraph for the RoBERTa-large model, given the large memory footprint and time processing when providing the entire essay to a larger model.

The original AAEC corpus also includes a fourth task, namely, stance recognition, where relational properties are classified with stance attributes (*for* or *against*). In our experiments, we did not perform this fourth task nor used the extra information provided with these stance attributes for the relational properties task.

## B Transfer during language modelling

We experimented with transferring knowledge from argument mining related sources by extending the pre-train, language modelling phase, rather than expanding the fine-tuning phase (as in the first and second batch of experiments). We experimented with three argumentation-oriented data sets under the Masked Language Modelling objective: a self-supervised approach was thus adopted, with no further labelled data resorted to during training.

In a first experiment, we extended the model with a train set obtained from the Oscar corpus (Ortiz Suárez et al., 2019) by parsing 1M sentences containing argumentative discourse markers. We extracted all sentences that contained argumentative discourse markers from premise to conclusion and conclusion to premise in an equal distribution.

In a second experiment, we extended the model with an argumentation data set, the Args.me corpus (Ajjour et al., 2019), containing 350k arguments from forum debates. Thirdly, we extended the model with ATOMIC, a common sense knowledge base converted to raw text (Sap et al., 2019) containing 877k inferential relations.

Each model was trained with three randomly initialized runs, for three epochs, with a learning rate of 1e-05 and fine-tuned for each task. The results are in Table 6.

**Results:** Some performance scores of these models are higher than the respective RoBERTa baseline, also used in the first two batches, however without a statistically significant difference. This

	Components	Clausal	Relational
Baseline	.916	.820	.727
Arg. markers	.908	<b>.825</b>	.717
Args.me	.915	.725	<b>.757</b>
ATOMIC	<b>.917</b>	.787	.716

Table 6: **Performance of models** obtained by further pre-training with data related to argument mining.

may indicate that for this type of approach to leveraging argument mining to be as effective as the approach in the first two batches of experiments, the volume of unlabelled data related to argument mining possibly needs to be higher than the labelled data resorted to there by far more orders of magnitude.