

# BERTweet’s TACO Fiesta: Contrasting Flavors On The Path Of Inference And Information-Driven Argument Mining On Twitter

Anonymous ACL submission

## Abstract

Argument mining, dealing with the classification of text based on inference and information, denotes a challenging analytical task in the rich context of Twitter (now  $\mathbb{X}$ ), a key platform for online discourse and exchange. Thereby, Twitter offers a diverse repository of short messages bearing on both of these elements. For text classification, transformer approaches, particularly BERT, offer state-of-the-art solutions. Our study delves into optimizing the embeddings of the understudied BERTweet transformer for argument mining on Twitter and broader generalization across topics. We explore the impact of pre-classification fine-tuning by aligning similar manifestations of inference and information while contrasting dissimilar instances. Using the TACO dataset, our approach augments tweets for optimizing BERTweet in a Siamese network, strongly improving classification and cross-topic generalization compared to standard methods. Overall, we contribute the transformer WRAPresentations and classifier WRAP, scoring 86.62% F1 for inference detection, 86.30% for information recognition, and 75.29% across four combinations of these elements, to enhance inference and information-driven argument mining on Twitter.

## 1 Introduction

Twitter (now  $\mathbb{X}$ ) is a global hub for opinions, news and information and serves as a primary data source for research, which had already recognized the value of its user-generated content prior to its transition to  $\mathbb{X}$  (Kwak et al., 2010; Boyd et al., 2010; Castillo et al., 2011).

Argument Mining is primarily about text classification, considering the structure of arguments, encompassing both informative and inferential elements (Palau and Moens, 2009; Peldszus and Stede, 2013; Lawrence and Reed, 2019).

For text classification, the pre-trained transformer BERT (Devlin et al., 2019) and its numerous domain-specific derivatives, such as BERTweet

(Nguyen et al., 2020), achieve state-of-the-art performance (Houlsby et al., 2019; Sun et al., 2019) with a soft-max classification head added as additional layers. During the fine-tuning process, such transformers are used to generate universal text representations serving as contextualized language features to inform the head, which in turn are further specialized for the actual downstream task.

Thereby, the field of argument mining has also witnessed the benefits of transformer models like BERT for cross-topic classification (Bhatti et al., 2021; Thorn Jakobsen et al., 2021) and argument similarity (Reimers and Gurevych, 2019; Reimers et al., 2019; Thakur et al., 2021) on the AFS (Misra et al., 2016), UKP (Stab et al., 2018), and IBM-Debater (Shnarch et al., 2018) corpus.

Besides the common methods of adjusting the in-task performance through parameter tweaks (Lan et al., 2019; You et al., 2019) or incorporating augmentations (Kaushik et al., 2019; Anaby-Tavor et al., 2019; Feng et al., 2021; Thakur et al., 2021), multi-task learning is recommended as an additional fine-tuning strategy (Sun et al., 2019; Stab et al., 2018). Thereby, multi-task learning denotes a prior phase of fine-tuning representations on auxiliary tasks such as clustering or semantic similarity before proceeding to the actual classification step and is argued to effectively reduce a model’s sensitivity to spurious correlations (Liu et al., 2019; Tu et al., 2020), which in turn is key to cross-topic argument mining (Thorn Jakobsen et al., 2021).

We believe that acquiring robust and meaningful representations, in the sense of perceiving the constituent elements of arguments, prior to classification is particularly useful for the nuanced task of argument mining when applied to diverse topics.

Generalizability in terms of cross-topic classification is crucial for practical argument mining in realistic scenarios, both in general research (Daxenberger et al., 2017; Stab et al., 2018) and specifically on Twitter (Schaefer and Stede, 2021), neces-

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

084 sitating models to focus on argument components  
085 while avoiding reliance on spurious correlations  
086 like topic words (Thorn Jakobsen et al., 2021).

087 In this paper, we pioneer the optimization of the  
088 understudied transformer BERTweet for argument  
089 mining on Twitter. Thereby, we refine its linguistic  
090 knowledge of tweets within the embedding space,  
091 specializing BERTweet to better encode inference  
092 and information across diverse topics.

093 Utilizing the TACO dataset (Feger and Di-  
094 etze, 2023), offering initial baseline evaluations of  
095 BERTweet for argument mining on Twitter, we op-  
096 timize the model’s representation layers in a multi-  
097 task approach by accentuating the contrast between  
098 inference and information while centering similar  
099 manifestations before the actual classification step.

100 We achieve this by configuring a Siamese  
101 BERTweet network using SBERT (Reimers and  
102 Gurevych, 2019). Applying contrastive loss (Had-  
103 sell et al., 2006) and text augmentation tech-  
104 niques (Wei and Zou, 2019), this network teaches  
105 BERTweet to cluster tweet embeddings according  
106 to their respective roles in argument mining, that  
107 is, to generally encode the presence or absence of  
108 both inference and information in the final repre-  
109 sentations for classification.

110 Utilizing BERTweet’s enhanced embeddings, it  
111 excels in both closed and cross-topic argument min-  
112 ing on Twitter, outperforming standard methods  
113 (Schaefer and Stede, 2021) in this domain.

114 Towards inference and information-driven argu-  
115 ment mining on Twitter, we contribute<sup>1</sup>:

- 116 • A pre-classification fine-tuning approach for  
117 BERTweet, enhancing its capacity to encode  
118 information and inference for closed and  
119 cross-topic argument mining on Twitter.
- 120 • An augmentation strategy to reduce spurious  
121 entity and topic signals while increasing sen-  
122 tence variability in tweets.
- 123 • WRAPresentations<sup>2</sup>, an enhanced BERTweet  
124 embedding model driven by inference and in-  
125 formation, achieved through contrastive opti-  
126 mization on augmented TACO tweets.
- 127 • WRAP<sup>3</sup>, our tweet argument classifier leverag-  
128 ing WRAPresentations for argument mining  
129 across diverse topics on Twitter.

<sup>1</sup>Code: [anonymous.4open.science/r/TACO-Fiesta](https://anonymous.4open.science/r/TACO-Fiesta)

<sup>2</sup>[huggingface.co/TomatenMarc/WRAPresentations](https://huggingface.co/TomatenMarc/WRAPresentations)

<sup>3</sup>[huggingface.co/TomatenMarc/WRAP](https://huggingface.co/TomatenMarc/WRAP)

## 2 Twitter Arguments from Conversations

Our primary dataset, TACO (Feger and Dietze, 2023), encompasses 1,734 tweets from 200 en-  
tire conversations spanning six topics: #Abortion  
(25.9%), #Brexit (29.0%), #GOT (11.0%), #LOTR-  
ROP (12.1%), #SquidGame (12.7%), and #Twitter-  
Takeover (9.3%). So far, it stands as the sole  
publicly available labeled tweet dataset tailored  
for conversation-level inference and information  
extraction, strategically addressing reply-pattern  
nuances inherent to their conversational contexts.

Annotations were conducted by six experts ac-  
cording to the Cambridge Dictionary<sup>4</sup> definitions,  
differentiating *inference* as *a guess that you make  
or an opinion that you form based on the informa-  
tion that you have* and *information* as *facts or de-  
tails about a person, company, product, etc..* With  
a robust agreement of 0.718 Krippendorff’s  $\alpha$ , four  
classes emerged of these elements: *Reason* (infer-  
ence and information), *Statement* (inference with-  
out information), *Notification* (information without  
inference), and *None* (neither element).

Table 1 details the class distribution of TACO.

Reason	Statement	Notification	None
581 (33.50%)	284 (16.38%)	500 (28.84%)	369 (21.28%)

Table 1: The class distribution of tweets in TACO.

On TACO, Vanilla BERTweet serves as the best  
performing baseline, excelling with 74.45% F1 for  
Reason, 56.66% F1 for Statement, 78.30% F1 for  
Notification, and 80.56% F1 for None after fine-  
tuning on these classes (Feger and Dietze, 2023).

## 3 Inference and Information-Driven Representations for Mining Arguments

In text classification, transformers like BERTweet  
use the final hidden state of the first token  $[CLS]$   
as the sequence representation. Classification in-  
volves a soft-max classifier added as an extension  
after the final representation layer, determining la-  
bel probabilities for a tweet  $t$  by evaluating the  
likelihood of assigning a label  $y$  as:

$$p(y|h) = \text{softmax}(\hat{W}h) \quad (1)$$

where,  $\hat{W}$  signifies the task-specific weights  
of the classification head, and  $h$  represents the  
final representation of  $t$  obtained with the trans-  
former. Achieved through pooling an entire se-  
quence representation via  $[CLS]$ ,  $h$  is expressed as

<sup>4</sup>[dictionary.cambridge.org](https://dictionary.cambridge.org)

$G_W(t) = h$ , where the transformer is considered an independent function  $G_W(t)$  with its distinct weights  $W$ , taking  $t$  as input. For the specific classification task, both  $\hat{W}$  and  $W$  are jointly fine-tuned by maximizing the log-probability of the correct label, where  $h$  implicitly undergoes optimization.

For optimizing class assignments on TACO, we emphasize the impact of specializing  $h$  for encoding inference and information before classification.

Hence, we consider the pre-classification specialization of an embedding  $h$  as a contrastive problem of semantic similarity, where tweets with similar expressions of the text dimensions inference and information are brought closer together, while those lacking in similarity are positioned farther apart.

### 3.1 Embedding Inference and Information

We measure the semantic similarity between two tweet representations, denoted as  $h_1$  and  $h_2$ , using cosine distance:

$$D(h_1, h_2) = 1 - \cos(h_1, h_2) \in [0, 2] \quad (2)$$

a standard metric (Mikolov et al., 2013; Kim, 2014; Tai et al., 2015; Chen and He, 2020) for assessing text vector similarity.  $D(h_1, h_2)$  reflects complete equivalence at 0, orthogonality at 1, and absolute dissimilarity at 2. Mainly defined by the cosine similarity  $\cos(h_1, h_2) \in [-1, 1]$ , where  $-1$  represents complete dissimilarity, 1 indicates equivalence, and values closer to 0 suggest orthogonality, this distance is length-independent and primarily influenced by the angle between two embeddings.

Building on this circumstance, we assume that the actual representation  $h$  of a tweet can be normalized and lies on the  $n$ -sphere:

$$S(n) = \{h \in \mathbb{R}^{n+1} : \|h\| = 1\} \quad (3)$$

Transferred to the Cartesian nature of arguments  $h = \langle inference, information \rangle$ , we consider their representations to live on the unit sphere  $h \in S(1)$  (Wang and Isola, 2020; Khosla et al., 2020; Chen and He, 2020). In  $h$ , 1 signifies full presence, and  $-1$  implies total absence of a component. Consequently, an ideal class center on the unit sphere heads towards the pole  $\langle 1, 1 \rangle$  for Reason,  $\langle 1, -1 \rangle$  for Statement,  $\langle -1, 1 \rangle$  for Notification, and  $\langle -1, -1 \rangle$  for None. A breakdown of this is shown in the upper part of Figure 1, acknowledging the realistic expectation that the actual embeddings may differ from the ideals while the objective is to get them closer to them.

### 3.2 Contrastive Siamese Network

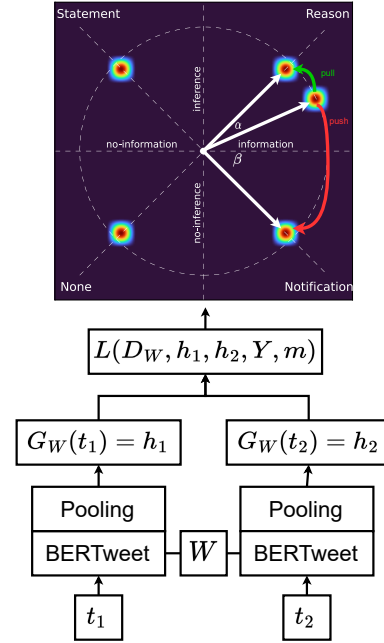


Figure 1: Visualization of the employed Siamese BERTweet architecture, with parameterized cosine distance  $D_W(h_1, h_2)$  and contrastive loss  $L(D_W, h_1, h_2, Y, m)$ . Atop this architecture, the Cartesian embedding space for an argument representation  $h = \langle inference, information \rangle$  is presented as target.

To address semantic similarity, a prevalent strategy involves enhancing representations through learning a metric (Chopra et al., 2005; Xing et al., 2002; Hadsell et al., 2006). Precisely, metric learning entails the implicit acquisition of a metric  $D_W(h_1, h_2)$  parameterized by the weights  $W$  of the representation model  $G_W$  (Chopra et al., 2005).

We seek to find  $W$  such that the target metric:

$$D_W(t_1, t_2) = 1 - \cos(G_W(t_1), G_W(t_2)) \quad (4)$$

is smaller if  $t_1, t_2$  are semantically similar, and higher if not.

By utilizing the identical embedding function  $G_W(t)$  (BERTweet) with shared weights  $W$  to learn the metric, our architecture is referred to as a Siamese network (Bromley et al., 1993; Chopra et al., 2005). Similar and dissimilar tweet pairs are provided as input to this network. To update the weights and optimize the network's performance, a loss function is applied on top of this architecture.

To attain the goal of increasing the differentiation between similar and dissimilar pairs, it is suggested to employ the contrastive loss (Chopra et al., 2005; Hadsell et al., 2006):

$$L(D_W, h_1, h_2, Y, m) = (Y) \frac{1}{2} D_w(h_1, h_2)^2 + (1 - Y) \frac{1}{2} \{ \max(0, m - D_w(h_1, h_2)) \}^2 \quad (5)$$

where,  $h_1, h_2$  are two representations ( $G_W(t_i) = h_i$ ) of different tweets  $t_1, t_2$  to be optimized given  $D_W(h_1, h_2)$  as metric.  $Y$  denotes the binary label indicating if  $t_1, t_2$  are similar ( $Y = 1$ ) or contrasting ( $Y = 0$ ). Furthermore, a margin value  $m > 0$  is introduced as the minimal distance between two contrasting tweets.

When establishing  $m$ , our objective was to set  $D_W(h_1, h_2)$  in a way that maximizes contrast between dissimilar pairs while avoiding over-estimation of their true distance. Focusing on  $D_W(h_1, h_2) \in [0, 1]$ , representing positive similarity, we selected  $m = 0.5$ . This choice intuitively represents the minimum threshold for high similarity, yielding optimal results in our study.

With  $m = 0.5$  we ensure that even if a representation closely matches an ideal center but is labeled as dissimilar, the optimized representation pushes  $60^\circ$  away and into an adjacent quadrant.

### 3.3 Augmentation of TACO

In the initial phase of processing TACO data, we generated a unique copy for each tweet through augmentation, denoted as A-TACO. Employing EDA (Easy Data Augmentation) techniques (Wei and Zou, 2019) of (1) synonym replacement, random (2) insertion, (3) swap, and (4) deletion, this procedure segregates our total ground truth into A-TACO, for optimization the embedding space of BERTweet prior to classification, and TACO, designated for fine-tuning and evaluating classifiers.

Maintaining independence between optimization and evaluation data is crucial to avoid spurious correlations (Thorn Jakobsen et al., 2021) and ensure that the data includes essential signals for class representations, thus enabling broad generalization across varying sentence structures and cross-topic evaluations for classifiers.

Following technique (1), we utilized spaCy<sup>5</sup> to identify all entities and specific words related to the six topics in the TACO dataset. Subsequently, we replace these words with the  $[MASK]$  token, a placeholder commonly used by BERT-like models, including BERTweet, for predicting missing words.

<sup>5</sup><https://spacy.io>

In particular, we utilized BERTweet as a fill-mask model to generate new tokens for those masked in the input sequence (Kumar et al., 2020).

To introduce word choice variability while minimizing semantic changes, random replacement of 10-90% of all words is applied using techniques (2-4). The optimal coherence, indicated by an average cosine distance of  $\sim 0.08$  between the  $[CLS]$  tokens of tweets and augmentations, is observed at a 10% replacement rate, maintaining overall semantic consistency while increasing sentence variability. Again, step 1 is applied to avoid reintroducing topic words. An example of the resulting tweet augmentation is shown in Table 2.

TACO	Elon Musk ready with 'Plan B' if Twitter rejects his offer Read @USER Story   HTTPURL #ElonMusk #ElonMuskTwitter #TwitterTakeover HTTPURL
A-TACO	Wenger ready with 'Plan B' as Wenger rejects his offer - HTTPURL via @USER

Table 2: An augmented Notification demonstrating entity replacement, topic word removal related to #TwitterTakeover, and altered sentence structure for enhanced anonymity, particularly at the end.

## 4 Experimental Setup

This section outlines the protocols used for evaluating and optimizing BERTweet’s embedding space with A-TACO and follow-up classification on TACO. Our primary objective is to acquire enhanced semantic similarity, with a specific emphasis on overall F1, while considering recall for generalizability to unseen topics.

### 4.1 Models

In our approach, it is important to differentiate between the pre-classification fine-tuning for specializing embeddings and their subsequent fine-tuning tailored for mining arguments on TACO.

For both tasks, we utilize the Vanilla BERTweet model<sup>6</sup>, with 12 transformer blocks and 12 self-attention heads processing sequences of up to 128 tokens, consistent with the baseline evaluation model of TACO (Feger and Dietze, 2023).

The embedding model BERTweet, enhanced through the application of contrastive loss within the Siamese network using A-TACO, is referred to as WRAPresentations.

Distinct from our multi-task approach, we introduce Augmented BERTweet, which undergoes

<sup>6</sup>[huggingface.co/vinai/bertweet-base](https://huggingface.co/vinai/bertweet-base)

pre-classification fine-tuning using the same tweets of A-TACO as WRAPresentations but directly optimizes  $p(y|h)$  through standard cross-entropy loss.

For classification on TACO, we utilize TF-IDF representations, where word frequency is widely recognized as a feature in strong baselines for argument mining on Twitter, which are Support Vector Machine (SVM) (Addawood and Bashir, 2016), Logistic Regression (LR) (Bosc et al., 2016; Dusmanu et al., 2017), and Random Forest (RF) (Dusmanu et al., 2017). These models go beyond considering individual words by incorporating tweet-related features like emoji, URL, and hashtag frequencies. Despite this, their potential for cross-topic generalizability remains unexplored.

For each classifier, we evaluate the average class length for classification to examine linguistic feature acquisition.

## 4.2 Pre-Classification Fine-Tuning

To enhance BERTweet’s embeddings, we chose TACO’s golden tweets with flawless annotation agreement, accounting for 70.30% of all tweets, with class distribution remaining largely consistent.

For the final evaluation, we employ original golden tweets for #Abortion but augmentations of golden tweets for the remaining five topics during fine-tuning. #Abortion is chosen as the holdout topic due to its highest dissimilarity when compared to the remaining topics, posing a greater challenge for classification (Thorn Jakobsen et al., 2021). This provides initial insights into cross-topic generalization and the efficacy of fine-tuning with augmentations and predicting given real tweets. Pairs are formed for all tweet combinations, denoting tweets of the same class as similar  $Y = 1$  and those of different classes as dissimilar  $Y = 0$ , yielding more dissimilar than similar pairs.

For the final validation set, 86,142 pairs were generated. The optimization data, divided into fine-tuning and test sets with a stratified 60/40 ratio, yielded 307,470 and 136,530 candidate pairs, respectively. To ensure a balance between similar and dissimilar pairs, we chose the largest possible set such that both similar and dissimilar pairs are equally represented (Bromley et al., 1993; Chopra et al., 2005) while maintaining all tweets of the respective splits.

In total, 162,064 pairs were obtained for fine-tuning, 71,812 for testing, and 53,560 for final validation of the enhanced BERTweet representations prior to classification.

For all transformer models, we performed fine-tuning over 5 epochs using an A100 GPU with 40GB of memory, a batch size of 32, and a learning rate of  $4e^{-5}$ , which proved to be optimal for all models. The Siamese BERTweet network is implemented using SBERT (Reimers and Gurevych, 2019) as depicted in the lower part of Figure 1.

Additionally, we performed fine-tuning on WRAPresentations using both  $[CLS]$  pooling, later employed for classification, and  $[MEAN]$  pooling, which is recommended for improved sentence embeddings (Reimers and Gurevych, 2019).

## 4.3 Argument Mining on TACO

We evaluate the practicality of BERTweet’s specialized embeddings on TACO, given the three argument mining tasks of (1) inference detection, (2) information recognition, and (3) classification of all four tweet classes, with a concurrent aim for cross-topic generalization.

For task (3), we trained a feed-forward neural network with two linear layers on top of each embedding model, undergoing 5 additional fine-tuning epochs with the best performing parameters having a learning rate of  $4e^{-5}$  and batch size of 8, corresponding to the best model and parameters reported for TACO (Feger and Dietze, 2023). Again, we used a single A100 GPU with 40GB of memory. Thereby, the results for tasks (1) and (2) are aggregations specific to class elements of task (3) predictions, focusing on inference or information.

Our classifier evaluation uses two distinct configurations to examine the impact of specialized embeddings and their adaptability to additional class adjustments (Peters et al., 2019).

In the first setup (Frozen), freezing embeddings allowed us to assess the benefits attributable to pre-classification fine-tuning. In the second setup (Dynamic), embeddings underwent further fine-tuning during classification head optimization, where we assessed their adaptability to task-specific learning. Success in this context signifies a model’s ability to leverage knowledge encoded in fine-tuned embeddings before classification and adapt them to the classes specific to inference and information.

We employed a 6-fold shuffled cross-validation, maintaining consistent splits for all classifiers across the six topics of TACO, to establish an upper-bound (Thorn Jakobsen et al., 2021). This closed-topic validation was then compared with cross-topic validation, where each of the six topics served as a unique testing set, and the remaining

four topics were utilized for fine-tuning (Bosc et al., 2016; Daxenberger et al., 2017; Stab et al., 2018). Lower performance is expected in cross-topic validation, as classifiers are exposed to unseen topics.

## 5 Results

In this section, we assess the impact of the specialized embeddings for closed and cross-topic classification on TACO.

### 5.1 Results Pre-Classification Fine-Tuning

Model	P	R	F1
Vanilla BERTweet- $[CLS]$	50.00	100.00	66.67
Augmented BERTweet- $[CLS]$	65.69	86.66	<b>74.73</b>
WRAPresentations- $[CLS]$	<b>66.00</b>	84.32	74.04
WRAPresentations- $[MEAN]$	63.05	<b>88.91</b>	73.78

Table 3: Evaluating within-class similarity and between-class separability of fine-tuned  $[CLS]$  representations on A-TACO towards TACO’s holdout topic #Abortion. WRAPresentations, with  $[MEAN]$  pooling in fine-tuning, show pessimistic scores but achieve a higher F1 score of 74.07% if tested with  $[MEAN]$  pooling.

After pre-classification fine-tuning to enhance semantic similarity, we evaluate the optimized embedding models for classifying tweet pairs as similar or dissimilar given  $D_W(t_1, t_2)$ .

All fine-tuning strategies outperformed Vanilla BERTweet in terms of F1, compare Table 3.

We excluded WRAPresentations with  $[CLS]$  pooling for follow-up classification due to the absence of discernible benefits in F1 compared to Augmented BERTweet and WRAPresentations using  $[MEAN]$  pooling for pre-classification fine-tuning, also showing a higher recall at 88.91%.

Hence, we will refer to WRAPresentations- $[MEAN]$  as WRAPresentations.

In comparing Augmented BERTweet and WRAPresentations, both models show similar overall performance in terms of F1, but diverge in their emphasis on precision and recall. The results suggest that contrastive fine-tuning of representations is not inherently superior to directly optimizing  $p(y|h)$  with augmented tweets. However, this strategy enhances recall, with further distinctions expected in downstream task evaluations.

Nonetheless, we assume that the enhanced recall at this stage is already a first indicator for later generalizations of classifications across topics. Additionally, we confirmed the effectiveness of pre-classification fine-tuning with A-TACO when applied to real tweets from an unseen topic.

Furthermore, we visually explored BERTweet’s embedding space before and after fine-tuning, utilizing  $[CLS]$  representations of all original tweets in TACO, as depicted in Figure 2(a).

Applying t-SNE for dimensional reduction (van der Maaten and Hinton, 2008; Jawahar et al., 2019), comparing Vanilla BERTweet with WRAPresentations showed enhanced class quadrant density, compare Figure 2(a), suggesting an improvement of class semantics given inference and information for a majority of tweets. Similar patterns, albeit at smaller numbers, are observed for Augmented BERTweet, see Figure 2(b).

Numerically, WRAPresentations improved tweet order by 38% for Reason, 37% for Statement, and 41% for Notification over Vanilla BERTweet. Despite a -2% decrease in the None class quadrant, it remains predominant, as shown in Figure 2(b).

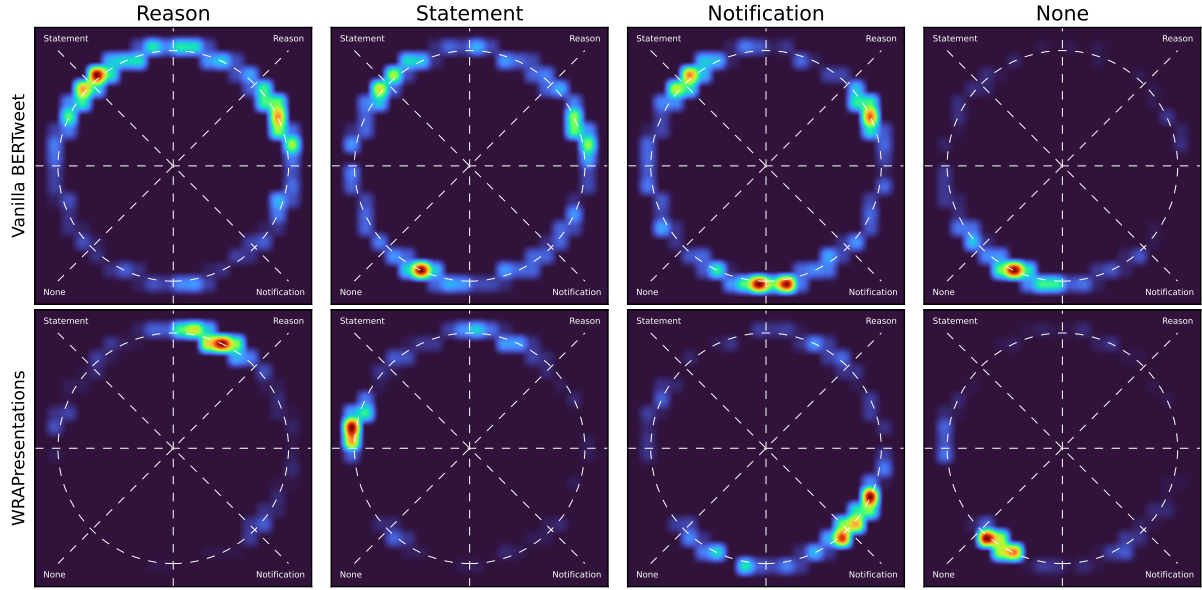
Augmented BERTweet closely matches WRAPresentations, excelling by 6% for None but lagging behind by -6% for Reason, -12% for Statement and -13% for Notification.

### 5.2 Results Classification and Generalization

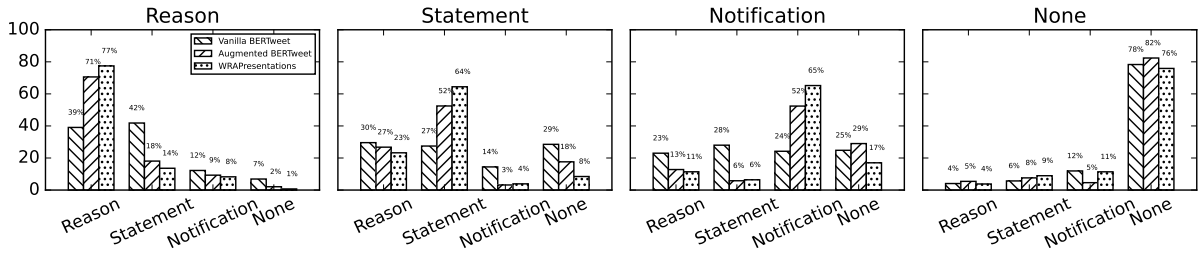
For our comparisons, we continue to present the outcomes of the Random Forest classifier as the most effective baseline and the average class length as a minimal-performance indicator. Furthermore, for publication, we refer to the classifier based on WRAPresentations as WRAP, while maintaining the term WRAPresentations for consistency.

When turning to the closed-topic validation, WRAPresentations outperforms all classifiers except task (1), where dynamic embeddings in Augmented BERTweet exhibit performance nearly equivalent, as demonstrated in the upper half of Table 4. Quantitatively, WRAPresentations yields 86.88% F1 for task (1), 81.54% F1 for task (2), and 71.07% F1 for task (3) when frozen. Dynamically optimizing embeddings, WRAPresentations achieves 86.62% F1 for task (1), 86.30% F1 for task (2), and 75.29% F1 for task (3).

Shifting our attention to the more demanding task of cross-topic validation, assessing a classifier’s ability to generalize to unseen topics, WRAPresentations demonstrates superior performance over all evaluations, thereby achieving 86.83% F1 for task (1), 81.54% F1 for task (2), and 70.93% F1 for task (3) when frozen. With dynamically adjusted embeddings, it achieves 86.27% F1 for task (1), 84.90% F1 for task (2), and 73.54% F1 for task (3), compare lower half of Table 4.



(a) t-SNE embeddings of tweet class  $[CLS]$  tokens before and after fine-tuning given inference and information.



(b) Distribution of classes within the projected quadrants of the expected  $(inference, information)$  space.

Figure 2: Investigation on the impact of BERTweet’s fine-tuning for the transfer of class semantics onto the expected  $(inference, information)$  space in terms of the  $[CLS]$  tokens for tweet classification. Considering the classes, (a) highlights the tightening of tweet embeddings towards their respective ideal class poles. Considering the distribution of tweets, (b) emphasizes that each expected quadrant corresponds to the anticipated majority class.

Model	<u>Inference</u>		<u>Information</u>		<u>Multi-Class</u>	
	Frozen	Dynamic	Frozen	Dynamic	Frozen	Dynamic
Closed-Topic (6-fold) Validation						
<b>Length</b>	62.34		71.47		38.26	
<b>RF + TF-IDF</b>	76.12		80.56		55.65	
<b>Vanilla BERTweet</b>	73.12	84.54	66.49	83.55	42.87	71.05
<b>Augmented BERTweet</b>	84.49	<b>86.68</b>	79.22	84.57	67.07	73.80
<b>WRAPresentations</b>	<b>86.88</b>	86.62	<b>81.54</b>	<b>86.30</b>	<b>71.07</b>	<b>75.29</b>
Cross-Topic (6-fold) Validation						
<b>Length</b>	61.99		71.55		38.17	
<b>RF + TF-IDF</b>	73.93		80.16		53.29	
<b>Vanilla BERTweet</b>	70.28	83.15	66.15	82.22	39.00	68.12
<b>Augmented BERTweet</b>	84.20	84.25	79.38	83.31	66.41	69.99
<b>WRAPresentations</b>	<b>86.83</b>	<b>86.27</b>	<b>81.54</b>	<b>84.90</b>	<b>70.93</b>	<b>73.54</b>

Table 4: Macro F1 scores of each classifier for inference and information detection, and all four classes.

Model	Reason		Statement		Notification		None	
	Frozen	Dynamic	Frozen	Dynamic	Frozen	Dynamic	Frozen	Dynamic
Closed-Topic (6-fold) Validation								
<b>Length</b>	61.68		20.19		14.47		56.72	
<b>RF + TF-IDF</b>	69.35		17.30		63.35		72.62	
<b>Vanilla BERTweet</b>	66.05	74.98	00.00	53.99	43.80	77.62	61.63	77.62
<b>Augmented BERTweet</b>	74.50	76.82	49.53	58.37	70.95	<b>80.28</b>	73.29	79.71
<b>WRAPresentations</b>	<b>77.34</b>	<b>78.14</b>	<b>58.66</b>	<b>60.96</b>	<b>72.61</b>	79.36	<b>75.67</b>	<b>82.72</b>
Cross-Topic (6-fold) Validation								
<b>Length</b>	61.78		19.32		14.49		57.09	
<b>RF + TF-IDF</b>	68.61		13.33		62.75		68.46	
<b>Vanilla BERTweet</b>	63.57	73.15	00.00	47.40	35.79	74.92	56.64	77.01
<b>Augmented BERTweet</b>	75.18	75.10	46.34	51.74	71.61	75.71	72.50	77.42
<b>WRAPresentations</b>	<b>77.13</b>	<b>77.05</b>	<b>57.62</b>	<b>58.33</b>	<b>73.05</b>	<b>78.45</b>	<b>75.91</b>	<b>80.33</b>

Table 5: F1 scores of the classifiers for identifying the four classes used in inference and information detection.

Further, WRAPresentations clearly improved performance for Statement, the least common and most difficult class to predict when comparing the remaining classifiers. Thereby, all other classifiers perform below or slightly above chance agreement for closed-topic validation and generalization across topics for this class, where Vanilla BERTweet even achieved 00.00% F1 when frozen, showcasing the necessity for adjusting classifiers and embeddings to specific classes, see Table 5.

## 6 Discussion

WRAPresentations consistently outperforms all models, with the exception of a marginal -0.06% F1 decrease compared to Augmented BERTweet with dynamic representations for task (1) of closed-topic evaluation, while totally excelling across topics.

Augmented BERTweet performs stronger in detecting instances without inference, as demonstrated by the substantial 9.33% F1 increase for the Notification class with dynamic embeddings, see upper half of Table 5. Considering that tasks (1) and (2) are aggregations derived from the results of task (3), WRAPresentations enhances the overall performance of task (3) for achieving the best results, prioritizing an improvement in task (2) while incurring a slight decrease in task (1).

This effect emerges as further refinements for additional classification improvements can partially replace the enriched understanding of inference and information in tweets, exposing unconsidered class features during optimization of the head.

However, examining WRAPresentations’ frozen states, superior in closed and cross-topic validation, underscores the advantages of our pre-classification

fine-tuning focused on semantic similarity in tweets for enhanced classification strength, see Table 4, 5.

Due to our multi-task fine-tuning approach, BERTweet can employ more robust embeddings for both classification scenarios, showcasing adaptability and generalizability across all three argument mining tasks on Twitter, including challenging instances like identifying the Statement class.

## 7 Conclusion and Ongoing Work

Our pre-classification multi-task fine-tuning approach considerably improves the specification of embeddings of BERTweet to encode diverse manifestations of inference and information, especially supporting the classification of tweets in TACO.

BERTweet’s optimized embeddings, enhanced through contrastively learning semantic similarity, offer improved adaptability to actual class signals and support cross-topic generalization when compared to conventional argument mining on Twitter.

In this regard, we can successfully contribute WRAPresentations, a contrastively optimized embedding model, and the advanced classification model WRAP for inference and information-driven argument mining across diverse topics on Twitter.

We also provide grounds for assuming that the augmentation of tweets constitutes a valuable asset within this domain of research.

Given our results demonstrating successful pre-classification fine-tuning with tweet augmentations and strong performance on original tweets, we pose the broader question of the necessity of using tweets for argument mining on Twitter, exploring whether tweet-like instances from other domains alone are sufficient.



## Acknowledgments

A hearty acknowledgment to our anonymous reviewers for their insightful feedback and the annotators behind TACO for generously sharing valuable data and influencing how we named things.

## Limitations

For our work, we report the following limitations:

The field of argument mining on Twitter is subject to Twitter’s data regulations, which allow only the publication of tweet identifiers but not their text. This poses challenges to the reproducibility of research and the potential loss of data due to deleted tweets when retrieved via their identifiers through the Twitter API, which provides a limited 1,500 free queries per month. However, for our study, we were able to obtain all preserved tweets from TACO by contacting the authors.

## References

- Aseel Addawood and Masooda Bashir. 2016. “what is your evidence?” a study of controversial topics on social media. In *Proceedings of the Third Workshop on Argument Mining (ArgMining2016)*, pages 1–11, Berlin, Germany. Association for Computational Linguistics.
- Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev Shlomov, Naama Tepper, and Naama Zwerdling. 2019. Not enough data? deep learning to the rescue! *CoRR*, abs/1911.03118.
- Muhammad Mahad Afzal Bhatti, Ahsan Suheer Ahmad, and Joonsuk Park. 2021. Argument mining on Twitter: A case study on the planned parenthood debate. In *Proceedings of the 8th Workshop on Argument Mining*, pages 1–11, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Tom Bosc, Elena Cabrio, and Serena Villata. 2016. DART: a dataset of arguments and their relations on Twitter. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 1258–1263, Portorož, Slovenia. European Language Resources Association (ELRA).
- Danah Boyd, Scott Golder, and Gilad Lotan. 2010. Tweet, tweet, retweet: Conversational aspects of retweeting on twitter. In *2010 43rd Hawaii International Conference on System Sciences*, pages 1–10.
- Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. 1993. Signature verification using a "siamese" time delay neural network. In *Advances in Neural Information Processing Systems*, volume 6. Morgan-Kaufmann.

- Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information credibility on twitter. In *Proceedings of the 20th International Conference on World Wide Web, WWW ’11*, page 675–684, New York, NY, USA. Association for Computing Machinery.
- Xinlei Chen and Kaiming He. 2020. Exploring simple siamese representation learning. *CoRR*, abs/2011.10566.
- S. Chopra, R. Hadsell, and Y. LeCun. 2005. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 1, pages 539–546 vol. 1.
- 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, Copenhagen, Denmark. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mihai Dusmanu, Elena Cabrio, and Serena Villata. 2017. Argument mining on Twitter: Arguments, facts and sources. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2317–2322, Copenhagen, Denmark. Association for Computational Linguistics.
- Marc Feger and Stefan Dietze. 2023. TACO: Twitter Arguments from COversations (Public-Data-1). <https://doi.org/10.5281/zenodo.8230057>. Data set.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. A survey of data augmentation approaches for NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, Online. Association for Computational Linguistics.
- R. Hadsell, S. Chopra, and Y. LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, volume 2, pages 1735–1742.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. *CoRR*, abs/1902.00751.

694	Ganesh Jawahar, Benoît Sagot, and Djamé Seddah.	Dat Quoc Nguyen, Thanh Vu, and Anh Tuan Nguyen.	750
695	2019. <a href="#">What does BERT learn about the structure of language?</a> In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3651–3657, Florence, Italy. Association for Computational Linguistics.	2020. <a href="#">BERTweet: A pre-trained language model for English tweets</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 9–14, Online. Association for Computational Linguistics.	751
696			752
697			753
698			754
699			755
700	Divyansh Kaushik, Eduard H. Hovy, and Zachary C. Lipton. 2019. <a href="#">Learning the difference that makes a difference with counterfactually-augmented data</a> . <i>CoRR</i> , abs/1909.12434.	Raquel Mochales Palau and Marie-Francine Moens. 2009. <a href="#">Argumentation mining: The detection, classification and structure of arguments in text</a> . In <i>Proceedings of the 12th International Conference on Artificial Intelligence and Law, ICAIL '09</i> , page 98–107, New York, NY, USA. Association for Computing Machinery.	756
701			757
702			758
703			759
704	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. 2020. <a href="#">Supervised contrastive learning</a> . <i>CoRR</i> , abs/2004.11362.		760
705			761
706			762
707			
708	Yoon Kim. 2014. <a href="#">Convolutional neural networks for sentence classification</a> . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1746–1751, Doha, Qatar. Association for Computational Linguistics.	Andreas Peldszus and Manfred Stede. 2013. <a href="#">From argument diagrams to argumentation mining in texts: A survey</a> . <i>Int. J. Cogn. Informatics Nat. Intell.</i> , 7:1–31.	763
709			764
710			765
711		Matthew E. Peters, Sebastian Ruder, and Noah A. Smith. 2019. <a href="#">To tune or not to tune? adapting pretrained representations to diverse tasks</a> . In <i>Proceedings of the 4th Workshop on Representation Learning for NLP (RepLanLP-2019)</i> , pages 7–14, Florence, Italy. Association for Computational Linguistics.	766
712			767
713			768
714	Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2020. <a href="#">Data augmentation using pre-trained transformer models</a> . In <i>Proceedings of the 2nd Workshop on Life-long Learning for Spoken Language Systems</i> , pages 18–26, Suzhou, China. Association for Computational Linguistics.		769
715			770
716			771
717			
718		Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-BERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	772
719			773
720	Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. 2010. <a href="#">What is twitter, a social network or a news media?</a> In <i>WWW '10: Proceedings of the 19th international conference on World wide web, WWW '10</i> , page 591–600, New York, NY, USA. Association for Computing Machinery.		774
721			775
722			776
723			777
724			778
725			779
726	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. <a href="#">ALBERT: A lite BERT for self-supervised learning of language representations</a> . <i>CoRR</i> , abs/1909.11942.	Nils Reimers, Benjamin Schiller, Tilman Beck, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2019. <a href="#">Classification and clustering of arguments with contextualized word embeddings</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 567–578, Florence, Italy. Association for Computational Linguistics.	780
727			781
728			782
729			783
730			784
731	John Lawrence and Chris Reed. 2019. <a href="#">Argument mining: A survey</a> . <i>Computational Linguistics</i> , 45(4):765–818.		785
732			786
733			787
734	Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. <a href="#">Multi-task deep neural networks for natural language understanding</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4487–4496, Florence, Italy. Association for Computational Linguistics.	Robin Schaefer and Manfred Stede. 2021. <a href="#">Argument mining on twitter: A survey</a> . <i>it - Information Technology</i> , 63(1):45–58.	788
735			789
736			790
737			
738		Eyal Shnarch, Carlos Alzate, Lena Dankin, Martin Gleize, Yufang Hou, Leshem Choshen, Ranit Aharonov, and Noam Slonim. 2018. <a href="#">Will it blend? blending weak and strong labeled data in a neural network for argumentation mining</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 599–605, Melbourne, Australia. Association for Computational Linguistics.	791
739			792
740	Tomás Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. <a href="#">Distributed representations of words and phrases and their compositionality</a> . <i>CoRR</i> , abs/1310.4546.		793
741			794
742			795
743			796
744	Amita Misra, Brian Ecker, and Marilyn Walker. 2016. <a href="#">Measuring the similarity of sentential arguments in dialogue</a> . In <i>Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 276–287, Los Angeles. Association for Computational Linguistics.		797
745			798
746			799
747			
748		Christian Stab, Tristan Miller, Benjamin Schiller, Pranav Rai, and Iryna Gurevych. 2018. <a href="#">Cross-topic argument mining from heterogeneous sources</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3664–3674, Brussels, Belgium. Association for Computational Linguistics.	800
749			801
			802
			803
			804
			805
			806

807 Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang.  
808 2019. [How to fine-tune BERT for text classification?](#)  
809 *CoRR*, abs/1905.05583.

810 Kai Sheng Tai, Richard Socher, and Christopher D. Man-  
811 ning. 2015. [Improved semantic representations from](#)  
812 [tree-structured long short-term memory networks](#). In  
813 *Proceedings of the 53rd Annual Meeting of the As-*  
814 *sociation for Computational Linguistics and the 7th*  
815 *International Joint Conference on Natural Language*  
816 *Processing (Volume 1: Long Papers)*, pages 1556–  
817 1566, Beijing, China. Association for Computational  
818 Linguistics.

819 Nandan Thakur, Nils Reimers, Johannes Daxenberger,  
820 and Iryna Gurevych. 2021. [Augmented SBERT: Data](#)  
821 [augmentation method for improving bi-encoders for](#)  
822 [pairwise sentence scoring tasks](#). In *Proceedings of*  
823 *the 2021 Conference of the North American Chapter*  
824 *of the Association for Computational Linguistics: Hu-*  
825 *man Language Technologies*, pages 296–310, Online.  
826 Association for Computational Linguistics.

827 Terne Sasha Thorn Jakobsen, Maria Barrett, and An-  
828 ders Søgaard. 2021. [Spurious correlations in cross-](#)  
829 [topic argument mining](#). In *Proceedings of \*SEM*  
830 *2021: The Tenth Joint Conference on Lexical and*  
831 *Computational Semantics*, pages 263–277, Online.  
832 Association for Computational Linguistics.

833 Lifu Tu, Garima Lalwani, Spandana Gella, and He He.  
834 2020. [An empirical study on robustness to spuri-](#)  
835 [ous correlations using pre-trained language models](#).  
836 *Transactions of the Association for Computational*  
837 *Linguistics*, 8:621–633.

838 Laurens van der Maaten and Geoffrey Hinton. 2008.  
839 [Visualizing data using t-sne](#). *Journal of Machine*  
840 *Learning Research*, 9(86):2579–2605.

841 Tongzhou Wang and Phillip Isola. 2020. [Understanding](#)  
842 [contrastive representation learning through alignment](#)  
843 [and uniformity on the hypersphere](#). In *Proceedings*  
844 *of the 37th International Conference on Machine*  
845 *Learning*, volume abs/2005.10242.

846 Jason Wei and Kai Zou. 2019. [EDA: Easy data augmen-](#)  
847 [tation techniques for boosting performance on text](#)  
848 [classification tasks](#). In *Proceedings of the 2019 Con-*  
849 *ference on Empirical Methods in Natural Language*  
850 *Processing and the 9th International Joint Confer-*  
851 *ence on Natural Language Processing (EMNLP-*  
852 *IJCNLP)*, pages 6382–6388, Hong Kong, China. As-  
853 sociation for Computational Linguistics.

854 Eric Xing, Michael Jordan, Stuart J Russell, and Andrew  
855 Ng. 2002. [Distance metric learning with application](#)  
856 [to clustering with side-information](#). In *Advances in*  
857 *Neural Information Processing Systems*, volume 15.  
858 MIT Press.

859 Yang You, Jing Li, Jonathan Hseu, Xiaodan Song, James  
860 Demmel, and Cho-Jui Hsieh. 2019. [Reducing BERT](#)  
861 [pre-training time from 3 days to 76 minutes](#). *CoRR*,  
862 abs/1904.00962.