

TACO – Twitter Arguments from COversations

Marc Feger¹, Stefan Dietze²

¹HeiCAD - Heine Center for Artificial Intelligence and Data Science,

²GESIS - Leibniz Institute for the Social Sciences

¹Düsseldorf Germany, ²Cologne Germany

marc.feger@hhu.de, stefan.dietze@gesis.org

Abstract

Twitter has emerged as a global hub for engaging in online conversations and as a research corpus for various disciplines that have recognized the significance of its user-generated content. Argument mining is an important analytical task for processing and understanding online discourse. Specifically, it aims to identify the structural elements of arguments, denoted as information and inference. These elements, however, are not static and may require context within the conversation they are in, yet there is a lack of data and annotation frameworks addressing this dynamic aspect on Twitter. We contribute *TACO*, the first dataset of *Twitter Arguments* utilizing 1,814 tweets covering 200 entire *COversations* spanning six heterogeneous topics annotated with an agreement of 0.718 Krippendorff's α among six experts. Second, we provide our annotation framework, incorporating definitions from the Cambridge Dictionary, to define and identify argument components on Twitter. Our transformer-based classifier achieves an 85.06% macro F1 baseline score in detecting arguments. Moreover, our data reveals that Twitter users tend to engage in discussions involving informed inferences and information. *TACO* serves multiple purposes, such as training tweet classifiers to manage tweets based on inference and information elements, while also providing valuable insights into the conversational reply patterns of tweets.

Keywords: Argument Mining, Twitter Conversations, Resource, Inference and Information Extraction

1. Introduction

Social media has created an open network of voices, connecting people globally and allowing them to exchange ideas and engage in discussions on any topic of interest. Despite these benefits, maintaining healthy and substantial online deliberation and promoting transparent information exchange remain key challenges in this area (Chadwick and Howard, 2008; Ruiz et al., 2011). Twitter, now X, serves as a global hub for opinions, news, and information, recognized for its research value and user-generated content prior to its rebranding (Kwak et al., 2010; Boyd et al., 2010).

In this context, argument mining has emerged as a valuable technique to identify the structure of inference and reasoning presented as arguments in natural language and is closely related to information extraction, fact checking, citation and opinion mining (Lawrence and Reed, 2019). This involves the automatic identification and extraction of arguments expressed in text, thus enabling researchers to analyze and understand the nature and structural elements of online discussions. Over the past years, the field of argument mining has undergone significant development in various domains, such as legal texts (Moens et al., 2007; Wyner et al., 2010), newspapers (Reed et al., 2008; Mochales and Moens, 2011), essays (Stab and Gurevych, 2014; Persing and Ng, 2016; Wachsmuth et al., 2016), Wikipedia articles (Levy et al., 2014, 2017), and sources of conflicting content, such as user comments (Park and Cardie,

2014), dialogues (Swanson et al., 2015), and web discourses (Habernal and Gurevych, 2017). While these works made initial contributions to the field, more recent research has focused on the detection of arguments from heterogeneous sources of arbitrary web text (Daxenberger et al., 2017; Levy et al., 2018; Stab et al., 2018). Despite addressing different aspects of argument mining, all studies involve the detection of inference and information (Palau and Moens, 2009; Daxenberger et al., 2017) as part of online discourse.¹

(1) Men shouldn't be making laws about women's bodies #abortion #Texas

(2) 'Bitter truth': EU chief pours cold water on idea of Brits keeping EU citizenship after #Brexit <https://t.co/j3DteyWcMg> via @TheLocalEurope

(3) Opinion: As the draconian (and then some) abortion law takes effect in #Texas, this is not an idle question for millions of Americans.

A slippery slope towards more like-minded Republican state legislatures to try to follow suit. #abortion #F24 <https://t.co/sMKUdhRF1q>

(4) @sinnfeinireland Blah blah blah blah blah

Table 1: Example tweets that contain inference (1), information (2), a combination of both (3), or none of either (4).

For argument mining on Twitter, research has expanded from specific topics like football (Llewellyn et al., 2014) and encryption (Addawood and Bashir, 2016) to encompass various subjects, including Brexit and Grexit (Dusmanu et al., 2017). Recent studies have focused on structuring tweets into debates via semantic similarity on the topics Iran, Grexit, Apple Watch and Game of Thrones (Bosc et al., 2016), and with isolated tweet-reply

¹These components, inference and information, are defined along the annotation framework in Section 2.1.

pairs on climate (Schaefer and Stede, 2020), neither capturing entire conversations.

This diversification has led to specialized tasks, such as identifying pro and con arguments in Planned Parenthood tweets (Bhatti et al., 2021) and evaluating scientific support in Covid-19 and climate-related tweets (Hafid et al., 2022).

Despite the progress in argument mining, the scope of related research on Twitter is restricted to a micro-level perspective, solely examining individual tweets and neglecting the interrelated reply tweets that make up the wider context of Twitter discussions. So far, no ground truth data for assessing arguments in entire Twitter conversations exists (Schaefer and Stede, 2021). With our work, we contribute the following to advance the field of argument mining on Twitter:

- 1 Annotation Framework.** Our specialized argument mining framework for Twitter conversations evolved from an extensive analysis of the elements defining arguments in relevant literature and iterative deliberations among our experts. It builds on Cambridge Dictionary's² definitions to define and identify inference and information within tweets.
- 2 Conversation-Based Ground Truth Data.** Our TACO³ dataset comprises 1,814 tweets, covering 200 entire conversations from six widely-discussed Twitter events. It was annotated by six experts with a high agreement score of 0.718 Krippendorff's α . TACO is available to the public in compliance with Twitter's data policy.
- 3 Baseline Classification Model.** Our published transformer-based classifier⁴, underlines TACO's significance in argument mining by achieving an 85.06% macro F1 for detecting arguments in tweets and 72.49% macro F1 for identifying combinations of inference and information. This classifier can be employed in both cases to aid in building new datasets and tweet curation.

2. Constructing the TACO dataset

Given the brevity of tweets, which were originally limited to 140 respectively 280 characters, structural elements of arguments such as inference or information tend to be scattered across distinct messages (Kwak et al., 2010; Boyd et al., 2010; Addawood and Bashir, 2016; Dusmanu et al., 2017). At the same time, tweets tend to be rather diverse in nature. Some tweets indicate a genuine

interest in contributing to ongoing debates, while others may express different motivations, such as a what-i-had-for-lunch-like tweet (Rogers, 2013).

2.1. Annotation Framework

With no one-size-fits-all definition of what an argument is (Palau and Moens, 2009; Habernal et al., 2014; Stab et al., 2018), the crucial challenge is how to identify arguments on Twitter. However, one potential strategy for simplifying this task is to differentiate between tweets that contain an inference as a key component of an argument and those that do not (Palau and Moens, 2009; Stab and Gurevych, 2014; Bosc et al., 2016; Daxenberger et al., 2017; Lawrence and Reed, 2019). Our aim here is not to create a new formalism for arguments, but rather to integrate established theories and provide a reusable set of definitions that can be applied to Twitter.

To define this critical component of an argument, we refer to the Cambridge Dictionary, which defines *inference* as *a guess that you make or an opinion that you form based on the information that you have*. We also utilize their description of *information* as *facts or details about a person, company, product, etc.*

Taken together, argument mining on Twitter involves determining if a tweet contains an inference (**Argument**) or not (**No-Argument**) by also considering its combination with information, as illustrated in Figure 1.

Tweets categorized as an **Argument** can be:

Statement, a tweet where only inference is presented like *something that someone says or writes officially, or an action done to express an opinion* (see tweet 1 in Table 1).

Reason, a tweet where the inference is based on information mentioned in the tweet, such as references, and thus reveals the author's motivation *to try to understand and to make judgments based on practical facts* (see tweet 3 in Table 1).

In contrast, tweets that are categorized as **No-Argument** are defined by the absence of inference and can be described as:

Notification, a tweet that is limited to only providing information, such as media channels promoting their articles (see tweet 2 in Table 1).

None, a tweet that provides neither inference nor information (see tweet 4 in Table 1).

2.2. Data Sampling and Annotation

Twitter conversations are shown to have a strong focus on various topics, often driven by hashtags (Hughes and Palen, 2009; Rogers, 2013; Zhou and Chen, 2014). We utilized Twitter's API v2 to collect a corpus of tweets and their reply-relations, enabling the extraction of entire conversations.

²dictionary.cambridge.org

³github.com/TomatenMarc/TACO

⁴huggingface.co/TomatenMarc/TACO

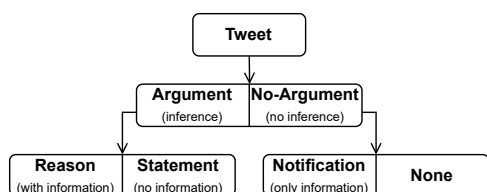


Figure 1: Hierarchy of inference and information.

While earlier studies of argument mining on Twitter mostly focused on one or two topics (Llewellyn et al., 2014; Addawood and Bashir, 2016; Dushmanu et al., 2017; Schaefer and Stede, 2020; Bhatti et al., 2021; Hafid et al., 2022), we aimed to create a more comprehensive corpus by collecting tweets from six controversial topics: #Abortion, #Brexit, #GOT, #TwitterTakeover, #SquidGame, and #LOTRROP. Over a period of seven months, we collected ~600k tweets around key-dates⁵ related to these hashtags, resulting in a significant increase in tweet volume. The allocation of the ~600k tweets and the profiles of the hashtag topics are as follows:⁵

#Abortion (5%) pertaining to the 'Row v. Wade' lawsuit, which challenges the Texas abortion ban (S.B.8) after six weeks, discussed and obtained between August 15 and October 16, 2021.

#Brexit (70.9%) relates to the global discussion of Great Britain's departure from the European Union on February 1, 2020. This historic event was queried from January 1 to March 1, 2020.

#GOT (10.2%) was gathered from April 1 to May 1, 2019, for expressing criticism about the final season of 'Game of Thrones'.

#TwitterTakeover (3.1%) was queried from April 1 to May 1, 2022, in association with Elon Musk's position as Twitter's largest shareholder, raising concerns about freedom of speech and initiating its transition to X.

#SquidGame (8.5%) explores economic inequality's impact on moral choices while also advertising the corresponding Netflix series, collected from September 10 to October 10, 2021.

#LOTRROP (2.3%) discussed the release of the trailer for Amazon's first season of 'The Lord of the Rings: The Rings of Power' and the related debate on representation and diversity in media, which we tracked from February 1 to March 1, 2022.

Text annotation is a multifaceted task, encompassing reading, comparing, memorization, and developing consensus about the data (Guetterman et al., 2018; Geiger et al., 2020). Under the author's guidance, two annotation rounds were performed with distinct expert groups to classify tweets based on our proposed annotation framework for Twitter conversations.⁵

⁵For the meta-data or examples, see: [README.md](#)

In step 1, three experts and the paper's first author refined the framework's guidelines and assessed their generalizability across topics. We sampled 300 conversation-starting tweets for #Abortion and #Brexit, given their significance in argument mining (Wachsmuth et al., 2017a; Dushmanu et al., 2017; Stab et al., 2018; Levy et al., 2018; Bhatti et al., 2021). Only when complete agreement on a tweet's classification was reached did it proceed as a candidate for step 2 annotation.

In step 2, the first author and two additional experts annotated 200 full conversations, tracing the sequences of reply-relations from the starting tweets to the final replies in a conversation, thereby considering the conversational context.⁵ This included 50 conversations for both #Abortion and #Brexit, and 25 conversations for each of the remaining four topics. In total, 236 #Abortion, 285 #Brexit, 192 #GOT, 166 #TwitterTakeover, 226 #SquidGame, and 209 #LOTRROP tweets were considered, averaging 219 tweets per topic.

In these steps, 1,814 annotated tweets were generated, including 1,314 conversation-based tweets of all topics and 500 distinct conversation-starting tweets for #Abortion and #Brexit, with a strong 0.718 Krippendorff's α agreement, given the author's involvement in both phases.

2.3. TACO Dataset

The final TACO ground truth (see Table 2) involved a strict majority vote approach, discarding tweets with less than 50% agreement among assigned votes for a specific class, resulting in 1,734 tweets, accounting for 95.6% of the annotated tweets.

Category	Argument		No-Argument	
	865 (49.88%)		869 (50.12%)	
Class	Reason	Statement	Notification	None
	581 (33.50%)	284 (16.38%)	500 (28.84%)	369 (21.28%)

Table 2: Class distribution in the TACO dataset.

3. Argument Mining on Twitter

We trained a soft-max classifier on top of different transformer models to utilize TACO in two tasks: (1) detecting inference in tweets (**Argument** vs. **No-Argument**), and (2) classifying tweets based on combinations of inference and information.

To obtain a first baseline for the usability of TACO, we fine-tuned the classifiers within an ordinary sequence classification approach integrating TACO's tweets, whose labels, created using our conversation-based annotation framework, encode implicit conversational context without detailing the conversation structures.⁶

⁶This approach is pragmatic, avoiding the complexity of modeling conversation hierarchies during fine-tuning.

For task (2), we fine-tuned the classifiers using 10-fold cross-validation, whereby the experimental results showed that BERTweet (Nguyen et al., 2020) provides superior classification performance. Further it is worth noting that the results for task (1) are aggregations of task (2), with a focus on presence or absence of inference.

The benefits of the BERTweet classifier for TACO extend beyond theory and are supported by cross-validation demonstrating strong performance for argument mining across the conversation-based tweets. Our results demonstrate this effectiveness, with the baseline model achieving an 85.06% macro F1 for inference detection of task (1) and 72.49% for classification of task (2), as seen in Table 3.

In terms of TACO’s data, the classification model had access to the following text features indicating tweet classes. The length of tweets differs among classes, with Reason being the longest on average (213), None the shortest (63), and Notification (156) and Statement (122) falling in the middle character range. URL usage varies, with Reason (34.6%) and Notification (71.6%) having the highest, while None and Statement use them less than 8.11%. The usage of discourse marker⁷ aligns with the argumentativeness of tweets: Reason (32.9%) is highest, followed by Statement (19%), Notification (11.4%), and None (8.7%).

Although there are misclassifications in task (2), ~43% of them are counted as true positives in task (1). The misclassification of Reason as Statement (or vice versa) still falls into the category of **Argument**, the same being true for the mutual misclassification of Notification and None, which still contributes as **No-Argument**.

Besides inference detection, the model also faces the challenge of correctly classifying tweets by identifying the informative parts, which is distinct from detecting inference and may not be entirely satisfactory due to the varied forms in which information can be presented, like URLs or fragmented text passages. Different tweets may employ language for specific intents involving rhetorical devices or visual elements, adding complexity to identifying information’s formal attributes.

Error analysis revealed that Statements, although typically lacking information, can include URLs when misclassified as Reason or Notification, as URLs might seem like apparent markers of information, resulting in an increased average length of 172. In fact, 22.09% of these misclassified Statements contained URLs, predominantly from internal sources like memes, GIFs, and videos. This circumstance might be influenced by Twitter’s recommendation system, which rewards

⁷dictionary.cambridge.org/us/grammar/british-grammar/discourse-markers

Task	Instance	Precision (%)	Recall (%)	F1 (%)	macro F1 (%)
Category	Argument	83.59	87.17	85.34	85.06
	No-Argument	86.66	82.97	84.77	
Class	Reason	73.69	75.22	74.45	72.49
	Statement	54.37	59.15	56.66	
	Notification	79.02	77.60	78.30	
	None	83.87	77.51	80.56	

Table 3: Baseline argument mining performance.

tweets that attract a large audience and have entertaining contents attached.⁸

Further investigation found that 24.42% of these incorrect assignments in the Statement class contained discourse markers, which complicated classification because these markers often organized relations between information and inference, leading to multiple stacked inferences that amplified the message’s tone rather than being perceived as information in the first instance.

To provide context, the model achieved a F1 score of 80.56% in detecting tweets lacking inference and information (None), which often led to conversation halts in about one-third (33.15%) of cases where a tweet received no further replies. Reason was the second most common in ending conversations at a rate of 29.73%, potentially indicating knockout arguments of further interest. Additionally, a Statement received no replies in 19.04% of cases, while Notification had no replies in 18.08%.

3.1. Conversational Reply Patterns

In our final analysis of the conversation-based ground truth data, we explored state transitions between connected tweet-reply pairs (t, r) to reveal TACO’s value in understanding reply patterns and providing insights into conversation progression, as shown in Table 4.

$P(r t)$	Reason	Statement	Notification	None
Reason	0.51	0.12	0.31	0.06
Statement	0.38	0.21	0.33	0.08
Notification	0.26	0.08	0.57	0.09
None	0.26	0.08	0.44	0.22

Table 4: Transition probability of a tweet with class t (row) having a reply with class r (column).

Our findings reveal users often reply with informed inferences (Reason) or additional information (Notification), with less common conversations solely based on inference (Statement) or lacking both elements (None). Additionally, **Argument** relies on informed inference, while **No-Argument** depends on information usage in replies, reflecting a preference for informed debates.

⁸github.com/twitter/the-algorithm

4. Conclusion

This paper presents the first ground truth dataset, TACO, for conversation-based argument mining on Twitter, efficiently annotated by six experts using our purpose-built annotation framework that incorporates definitions of argument constituting elements from the Cambridge Dictionary. Unlike previous datasets that often rely on isolated tweets without the contextual framework of conversations, TACO offers fully annotated and coherent Twitter conversations across six topics used for training a transformer-based classification model, providing valuable resources for future research in this domain. Furthermore, the provided classifier effectively differentiates tweets that make arguments from those that do not, based on the presence of inference. Additionally, our multi-class approach sufficiently identifies tweet classes, especially those that lack information and inference. Our findings suggest the need for further research to enhance the semantic features of our proposed tweet classes, possibly by fine-tuning BERTweet for more generalized representations according to our framework.

5. Ethics Statement

In the context of this study, which uses Twitter data, we have adhered to ethical practices and privacy principles and ensured data protection by limiting the publication of TACO to tweet IDs in accordance with Twitter's terms of service. Our annotation process, which involved volunteer experts, has been carefully designed to limit data access to what was strictly necessary and to ensure ethical standards, fair compensation and data integrity. Access to the original dataset is restricted to non-harmful research, subject to appropriate data protection agreements with the authors. It should also be noted that the TACO dataset covers sensitive topics that may contain language and images that some may find offensive.

6. Acknowledgments

We express our gratitude to Aylin Feger, Tillmann Junk, Andreas Burbach, Talha Caliskan, and Aaron Schneider for their valuable contributions as experts to the annotation process. Additionally, we extend our sincere appreciation to the anonymous reviewers, whose fair and insightful assessments have significantly enhanced the quality of our work. Your willingness to share your expertise and provide constructive feedback is deeply appreciated.

7. Bibliographical 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.
- Roy Bar-Haim, Lilach Eden, Roni Friedman, Yoav Kantor, Dan Lahav, and Noam Slonim. 2020. [From arguments to key points: Towards automatic argument summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4029–4039, Online. Association for Computational Linguistics.
- 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.
- Andrew Chadwick and Philip N Howard. 2008. *Routledge Handbook of Internet Politics*, 1st edition. Routledge.
- 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. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.

- 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.
- R. Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang. 2020. [Garbage in, garbage out? do machine learning application papers in social computing report where human-labeled training data comes from?](#) In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, FAT* '20*, page 325–336, New York, NY, USA. Association for Computing Machinery.
- Martin Gleize, Eyal Shnarch, Leshem Choshen, Lena Dankin, Guy Moshkovich, Ranit Aharonov, and Noam Slonim. 2019. [Are you convinced? choosing the more convincing evidence with a Siamese network](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 967–976, Florence, Italy. Association for Computational Linguistics.
- Miha Grčar, Darko Cherepnalkoski, Igor Mozetič, and Petra Kralj Novak. 2017. [Stance and influence of twitter users regarding the brexit referendum](#). *Computational Social Networks*, 4(1):6.
- Timothy C Guetterman, Tammy Chang, Melissa DeJonckheere, Tannya Basu, Emily Scruggs, and V G Vinod Vydiswaran. 2018. [Augmenting qualitative text analysis with natural language processing: Methodological study](#). *Journal of medical Internet research*, 20(6):e231.
- Ivan Habernal, Judith Ecker-Köhler, and Iryna Gurevych. 2014. [Argumentation mining on the web from information seeking perspective](#). In *Workshop on Frontiers and Connections between Argumentation Theory and Natural Language Processing*.
- Ivan Habernal and Iryna Gurevych. 2017. [Argumentation Mining in User-Generated Web Discourse](#). *Computational Linguistics*, 43(1):125–179.
- Salim Hafid, Sebastian Schellhammer, Sandra Bringay, Konstantin Todorov, and Stefan Dietze. 2022. [Scitweets - a dataset and annotation framework for detecting scientific online discourse](#). In *Proceedings of the 31st ACM International Conference on Information; Knowledge Management, CIKM '22*, page 3988–3992, New York, NY, USA. Association for Computing Machinery.
- Amanda Hughes and Leysia Palen. 2009. [Twitter adoption and use in mass convergence and emergency events](#). *International Journal of Emergency Management*, 6:248–260.
- Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. 2010. [What is twitter, a social network or a news media?](#) In *Proceedings of the 19th International Conference on World Wide Web, WWW '10*, page 591–600, New York, NY, USA. Association for Computing Machinery.
- John Lawrence and Chris Reed. 2019. [Argument mining: A survey](#). *Computational Linguistics*, 45(4):765–818.
- Ran Levy, Yonatan Bilu, Daniel Hershcovich, Ehud Aharoni, and Noam Slonim. 2014. [Context dependent claim detection](#). In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1489–1500, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- Ran Levy, Ben Bogin, Shai Gretz, Ranit Aharonov, and Noam Slonim. 2018. [Towards an argumentative content search engine using weak supervision](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2066–2081, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Ran Levy, Shai Gretz, Benjamin Sznaider, Shay Hummel, Ranit Aharonov, and Noam Slonim. 2017. [Unsupervised corpus-wide claim detection](#). In *Proceedings of the 4th Workshop on Argument Mining*, pages 79–84, Copenhagen, Denmark. Association for Computational Linguistics.
- Clare Llewellyn, Claire Grover, Jon Oberlander, and Ewan Klein. 2014. [Re-using an argument corpus to aid in the curation of social media collections](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 462–468, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Raquel Mochales and Marie-Francine Moens. 2011. [Argumentation mining](#). *Artificial Intelligence and Law*, 19(1):1–22.
- Marie-Francine Moens, Erik Boiy, Raquel Mochales Palau, and Chris Reed. 2007. [Automatic detection of arguments in legal texts](#). In *Proceedings of the 11th International Conference on Artificial Intelligence and Law, ICAIL '07*, page 225–230, New York, NY, USA. Association for Computing Machinery.

- Dat Quoc Nguyen, Thanh Vu, and Anh Tuan Nguyen. 2020. [BERTweet: A pre-trained language model for English tweets](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 9–14, Online. Association for Computational Linguistics.
- Ryosuke Nishi, Taro Takaguchi, Keigo Oka, Takanori Maehara, Masashi Toyoda, Ken-ichi Kawarabayashi, and Naoki Masuda. 2016. [Reply trees in twitter: data analysis and branching process models](#). *Social Network Analysis and Mining*, 6(1):26.
- Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1999. [The PageRank citation ranking: Bringing order to the web](#). In *Proceedings of the 7th International World Wide Web Conference*, pages 161–172.
- Raquel Mochales Palau and Marie-Francine Moens. 2009. [Argumentation mining: The detection, classification and structure of arguments in text](#). In *Proceedings of the 12th International Conference on Artificial Intelligence and Law, ICAIL '09*, page 98–107, New York, NY, USA. Association for Computing Machinery.
- Joonsuk Park and Claire Cardie. 2014. [Identifying appropriate support for propositions in online user comments](#). In *Proceedings of the First Workshop on Argumentation Mining*, pages 29–38, Baltimore, Maryland. Association for Computational Linguistics.
- Isaac Persing and Vincent Ng. 2016. [End-to-end argumentation mining in student essays](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1384–1394, San Diego, California. Association for Computational Linguistics.
- Damien Pfister. 2014. [“a short burst of inconsequential information:” networked rhetorics, avian consciousness, and bioegalitarianism](#). *Environmental Communication: A Journal of Nature and Culture*, 9:1–19.
- Chris Reed, Raquel Mochales Palau, Glenn Rowe, and Marie-Francine Moens. 2008. [Language resources for studying argument](#). In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco. European Language Resources Association (ELRA).
- Nils Reimers and Iryna Gurevych. 2019. [Sentencebert: Sentence embeddings using siamese bert-networks](#).
- Nils Reimers, Benjamin Schiller, Tilman Beck, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2019. [Classification and clustering of arguments with contextualized word embeddings](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 567–578, Florence, Italy. Association for Computational Linguistics.
- Ruty Rinott, Lena Dankin, Carlos Alzate Perez, Mitesh M. Khapra, Ehud Aharoni, and Noam Slonim. 2015. [Show me your evidence - an automatic method for context dependent evidence detection](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 440–450, Lisbon, Portugal. Association for Computational Linguistics.
- Richard Rogers. 2013. [Debanalizing twitter: The transformation of an object of study](#). In *Proceedings of the 5th Annual ACM Web Science Conference, WebSci '13*, page 356–365, New York, NY, USA. Association for Computing Machinery.
- Carlos Ruiz, David Domingo, Josep Lluís Micó, Javier Díaz-Noci, Koldo Meso, and Pere Masip. 2011. [Public sphere 2.0? the democratic qualities of citizen debates in online newspapers](#). *The International Journal of Press/Politics*, 16(4):463–487.
- Robin Schaefer and Manfred Stede. 2020. [Annotation and detection of arguments in tweets](#). In *Proceedings of the 7th Workshop on Argument Mining*, pages 53–58, Online. Association for Computational Linguistics.
- Robin Schaefer and Manfred Stede. 2021. [Argument mining on twitter: A survey. it - Information Technology](#), 63(1):45–58.
- 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, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- Christian Stab, Tristan Miller, Benjamin Schiller, Pranav Rai, and Iryna Gurevych. 2018. [Cross-topic argument mining from heterogeneous sources](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3664–3674, Brussels, Belgium. Association for Computational Linguistics.
- Reid Swanson, Brian Ecker, and Marilyn Walker. 2015. [Argument mining: Extracting arguments](#)

from online dialogue. In *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 217–226, Prague, Czech Republic. Association for Computational Linguistics.

Henning Wachsmuth, Khalid Al-Khatib, and Benno Stein. 2016. [Using argument mining to assess the argumentation quality of essays](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1680–1691, Osaka, Japan. The COLING 2016 Organizing Committee.

Henning Wachsmuth, Martin Potthast, Khalid Al-Khatib, Yamen Ajjour, Jana Puschmann, Jiani Qu, Jonas Dorsch, Viorel Morari, Janek Bevendorff, and Benno Stein. 2017a. [Building an argument search engine for the web](#). In *Proceedings of the 4th Workshop on Argument Mining*, pages 49–59, Copenhagen, Denmark. Association for Computational Linguistics.

Henning Wachsmuth, Benno Stein, and Yamen Ajjour. 2017b. [“PageRank” for argument relevance](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1117–1127, Valencia, Spain. Association for Computational Linguistics.

Adam Wyner, Raquel Mochales-Palau, Marie-Francine Moens, and David Milward. 2010. [Approaches to Text Mining Arguments from Legal Cases](#), pages 60–79. Springer Berlin Heidelberg, Berlin, Heidelberg.

Xiangmin Zhou and Lei Chen. 2014. [Event detection over twitter social media streams](#). *The VLDB Journal*, 23(3):381–400.