# Slides for the 2-minute presentation

# What to Fact-Check

Guiding Check-Worthy Information Detection in News Articles through Argumentative Discourse Structure

**Tariq Alhindi**

Brennan Xavier McManus

Smaranda Muresan

SIGDIAL 2021

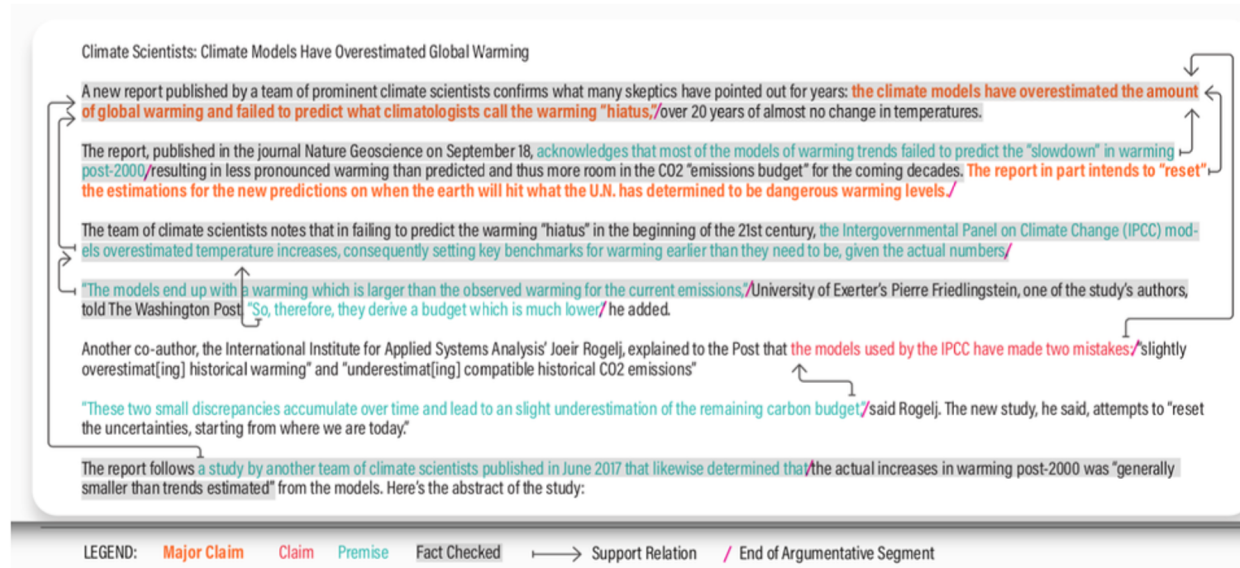The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue, 29-31 July 2021, Singapore

# Information Check-worthiness

Most work on fact-checking start with a list of claims to fact-check
(Throne et al., 2018, Wang 2017)

Previous work on check-worthiness

- Political text (mostly debates) using handcrafted features (Hassan et al., 2017, Jaradat et al., 2018)
- The notion of check-worthiness greatly varies across genre (Wright and Augenstein, 2020).

Is check-worthiness related to argument structure?



**Hypothesis**

Fact-checking a premise when it supports a claim

Fact-checking a claim when it is not supported or only supported by other claims

(Evading the Burden of Proof)

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

3

# Corpus

Multilayer annotated corpus of 95 articles from climatefeedback.org.

- fact-checked text segments by climate scientist at climatefeedback.org

- argument structure (major claim, claim, premises and  support, attack relations)
  by 6 expert annotators

Following previous work, we approach this as:

- sentence classification task        Macro F1

- sentence ranking task            MAP

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Approach

We take advantage of BERT next sentence capabilities to add context to the target sentence:

- Local discourse context (prev+sent, sent+next)

- Argumentation context by pairing the target sentence with another sentence that has an argumentative relation (support, attack, joint, restate) with the target sentence.

  if the target sentence has an argumentative component (major-claim, claim, premise)
  otherwise we revert back to discourse context

  additionally, we prepend the Argumentative component Type (AT)
  e.g.        CLAIM        *the model used by the IPCC has two mistakes*        Not-Checked
  **ArgType**                              target sentence                              label

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Results

| Group | Input | Development Set | | | | Test Set | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | NC | FC | Macro F1 | MAP | NC | FC | Macro F1 | MAP |
| Baselines | SENT | 0.83 | 0.23 | 0.53 | 0.296 | **0.85** | 0.28 | 0.56 | 0.398 |
| | PREV+SENT | 0.83 | 0.29 | 0.56 | **0.387** | 0.82 | 0.29 | 0.56 | 0.384 |
| | SENT+NEXT | 0.83 | 0.27 | 0.55 | 0.296 | 0.84 | 0.26 | 0.55 | 0.385 |
| Argument | SENT+AC | **0.84** | **0.33** | **0.58** | 0.366 | 0.83 | 0.30 | 0.57 | 0.413 |
| Context | SENT+AC+AT | 0.83 | 0.29 | 0.56 | 0.359 | 0.84 | **0.33** | **0.59**† | **0.420**† |

Per-class F1 (**NC**: Not-Checked class, **FC**: Fact-Checked class), Macro F1

and Mean Average Precision (**MAP**) on the development and test sets.

**AC**: Argumentation Context, **AT**: Argumentative component Type

† Statistically significant over baselines

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Contributions

A novel corpus with multi-layer annotations for
   check-worthiness and argument structure

Model check-worthiness in news articles
   as sentence classification and a sentence ranking tasks

Using argument structure as context yields better results than using
   local discourse context for the task of check-worthiness detection

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

Thank You

www.cs.columbia.edu/~tariq
www.github.com/tariq60/whatToFactCheck

# Slides for the 15-minute presentation

# What to Fact-Check
Guiding Check-Worthy Information Detection in News Articles through Argumentative Discourse Structure

**Tariq Alhindi**

Brennan Xavier McManus

Smaranda Muresan

# Motivation

Most work on fact-checking start with a list of claims to fact-check (Throne et al., 2018, Wang 2017)

Previous work on check-worthiness
- Political text (mostly debates) using handcrafted features (Hassan et al., 2017, Jaradat et al., 2018)
- The notion of check-worthiness greatly varies across genre (Wright and Augenstein, 2020).

What about check-worthiness in news articles from different topics (e.g. climate change)?

Is check-worthiness related to argument structure?

**Hypothesis**

Fact-check a premise when it supports a claim

Fact-check a claim when it is not supported or only supported by other claims (Evading the Burden of Proof)

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Example



Climate Scientists: Climate Models Have Overestimated Global Warming

A new report published by a team of prominent climate scientists confirms what many skeptics have pointed out for years: the climate models have overestimated the amount of global warming and failed to predict what climatologists call the warming "hiatus," / over 20 years of almost no change in temperatures.

The report, published in the journal Nature Geoscience on September 18, acknowledges that most of the models of warming trends failed to predict the "slowdown" in warming post-2000 / resulting in less pronounced warming than predicted and thus more room in the CO2 "emissions budget" for the coming decades. The report in part intends to "reset" the estimations for the new predictions on when the earth will hit what the U.N. has determined to be dangerous warming levels. /

The team of climate scientists notes that in failing to predict the warming "hiatus" in the beginning of the 21st century, the Intergovernmental Panel on Climate Change (IPCC) models overestimated temperature increases, consequently setting key benchmarks for warming earlier than they need to be, given the actual numbers /

"The models end up with a warming which is larger than the observed warming for the current emissions," / University of Exerter's Pierre Friedlingstein, one of the study's authors, told The Washington Post. "So, therefore, they derive a budget which is much lower / he added.

Another co-author, the International Institute for Applied Systems Analysis' Joeir Rogelj, explained to the Post that the models used by the IPCC have made two mistakes: / "slightly overestimat[ing] historical warming" and "underestimat[ing] compatible historical CO2 emissions"

"These two small discrepancies accumulate over time and lead to an slight underestimation of the remaining carbon budget," / said Rogelj. The new study, he said, attempts to "reset the uncertainties, starting from where we are today."

The report follows a study by another team of climate scientists published in June 2017 that likewise determined that / the actual increases in warming post-2000 was "generally smaller than trends estimated" from the models. Here's the abstract of the study:

LEGEND:  **Major Claim**   Claim   Premise   Fact Checked   ⟶ Support Relation   / End of Argumentative Segment

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

12

# Outline

Related Work

Data

Model

Results

Conclusion

# Outline

**Related Work**

Data

Model

Results

Conclusion

# Related Work

ClaimBuster (Hassan et al., 2017) and ClaimRank (Jaradat et al., 2018)

CLEF check that lab
(Nakov et al., 2018; Elsayed et al., 2019; Barron-Cedeno et al. , 2020)

Argumentation and check-worthiness
      Type of statements (Freeman, 2000)
      Type of evidence (Park and Cardie, 2014; Addawood and Bashir, 2016)

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Outline

Related Work

**Data**

Analysis & Model

Results

Conclusion

# Data

95 climate change news articles
 fact-checked text segments by climate scientists at *climatefeedback.org*

from 40 publishers mainly in the U.S., UK and Australia

e.g., The New York Times, The Guardian, The Washington Post, The Wall Street Journal, The Australian, The Telegraph, Forbes, USA today, Breitbart, and Mashable

Articles are given an article-level credibility rating
 and sentence-level fact-checking annotations

| Credibility | Count | Credibility | Count |
|---|---|---|---|
| very-low | 23 | high | 21 |
| very-low/low | 7 | high/very-high | 8 |
| low | 10 | very-high | 18 |
| neutral | 7 | mixed | 1 |

Each Article is tagged by 3 to 5 climate scientists

evaluate scientific reasoning
add relevant information missed by the article
check for: factual accuracy, scientific understanding, logical reasoning
precision/clarity, sources quality, and fairness/objectivity

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Data – Factchecked Segments

Fact-checked segments vary in length

from a fragment of a sentence to multiple sentences.

We thus map this to binary labels at the sentence-level: factchecked (FC) or not-checked (NC).

A sentence is labeled as 'fact-checked' if:

it is fact-checked
has a fact-checked segment
part of multi-sentence fact-checked segment

We split the the 95 articles to

| | | |
|---|---|---|
| 68 articles in the *training* set | 4,353 sentences in total | 824 are fact-checked |
| 7 articles in the *development* set | 249 sentences in total | 55 are fact-checked |
| 20 articles in the *test* set | 970 sentences in total | 220 are fact-checked |

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Data – Argument Structure Annotation

Annotation Scheme
  Argument Components  Major-Claim, Claim, Premise
  Argument Relations  Support, Attack, Restate, Joint

Six Annotators
  Undergrads in Linguistics, English, and Comparative Literature
  Each annotators was assigned a 32-article batch; Each article annotated by at least 3 annotators

Gold Annotations
  Minimum common span of overlapping components from the three annotations
  Relations between gold components only
    adherence to guidelines
    annotator quality

IAA using Krippendorff's alpha
  overall IAA is .4368
  using the coding version, which uses only the labels assigned to each component

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Outline

Related Work

Data

**Analysis & Model**

Results

Conclusion

# Analysis – Argumentation w.r.t Fact-Checked Segments

| Gold Annotations | |
|---|---|
| **AC Type** | **Frequency** |
| Claim | 91 |
| Premise | 76 |
| Major-Claim | 22 |
| Premise Premise | 20 |
| Claim Premise | 17 |
| Claim Claim | 12 |
| Premise Claim | 9 |
| Premise Claim Claim | 4 |
| Premise Premise Claim | 4 |
| Claim Premise Claim | 4 |

| AC Type | Total Rel. | Relation Type | Frequency |
|---|---|---|---|
| Claim | 1 | $\xrightarrow{\text{sup}}$ Claim | 12 |
| | 1 | $\xrightarrow{\text{sup}}$ Major-Claim | 11 |
| Premise | 1 | $\xrightarrow{\text{sup}}$ Claim | 54 |
| | 1 | $\xrightarrow{\text{sup}}$ Premise | 4 |
| Major | ≥4 | $\xleftarrow{\text{sup}}$ Claim (all) | 10 |
| Claim | 1 | $\xrightarrow{\text{oth}}$ Major-Claim | 2 |

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Model

Following previous work, we approach check-worthiness detection as a:

- sentence classification task      Macro F1
- sentence ranking task      MAP

We take advantage of BERT next sentence capabilities to add context to the target sentence:

- Local discourse context (prev+sent, sent+next)
- Argumentation context by pairing the target sentence with another sentence that has an argumentative relation (support, attack, joint, restate) with the target sentence.

  if the target sentence has an argumentative component (major-claim, claim, premise)
       otherwise we revert back to discourse context

  additionally, we prepend the Argumentative component Type (AT)
       e.g.      CLAIM     *the model used by the IPCC has two mistakes*     Not-Checked
             **ArgType**          target sentence         label

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Outline

Related Work

Data

Analysis & Model

**Results**

Conclusion

# Results – Development Set

| Group | Model Input | Not-Checked | Fact-Checked | Macro F1 | MAP |
|---|---|---|---|---|---|
| Baselines | SENT | 0.83 | 0.23 | 0.53 | 0.296 |
| | PREV+SENT | 0.83 | 0.29 | 0.56 | **0.387** |
| | SENT+NEXT | 0.83 | 0.27 | 0.55 | 0.296 |
| Argument Context (Text only) | SENT+AC(1) | 0.84 | **0.33** | **0.58** | 0.366 |
| | SENT+AC(3)$^{v1}$ | 0.82 | 0.31 | 0.57 | 0.299 |
| | SENT+AC(3)$^{v2}$ | 0.82 | 0.32 | 0.57 | 0.299 |
| | SENT+AC(ALL)$^{v1}$ | 0.83 | 0.26 | 0.54 | 0.318 |
| | SENT+AC(ALL)$^{v2}$ | 0.81 | 0.30 | 0.56 | 0.318 |
| Argument Context (Text+Type) | SENT+AC(1)+T | 0.83 | 0.29 | 0.56 | 0.359 |
| | SENT+AC(3)+T$^{v1}$ | 0.84 | 0.27 | 0.57 | 0.305 |
| | SENT+AC(3)+T$^{v2}$ | **0.85** | 0.29 | 0.57 | 0.305 |
| | SENT+AC(ALL)+T$^{v1}$ | 0.82 | 0.32 | 0.57 | 0.281 |
| | SENT+AC(ALL)+T$^{v2}$ | 0.82 | 0.31 | 0.57 | 0.281 |

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Results – Test Set

| Input | NC | FC | F1 | MAP |
|---|---|---|---|---|
| SENT | **0.85** | 0.28 | 0.56 | 0.398 |
| PREV+SENT | 0.82 | 0.29 | 0.56 | 0.384 |
| SENT+NEXT | 0.84 | 0.26 | 0.55 | 0.385 |
| SENT+AC(1) | 0.83 | 0.30 | 0.57 | 0.413 |
| SENT+AC(1)+T | 0.84 | **0.33** | **0.59**$^{\dagger}$ | **0.420**$^{\dagger}$ |

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Outline

Related Work

Data

Analysis & Model

Results

**Conclusion**

# Conclusion

A novel corpus with multi-layer annotations for check-worthiness and argument structure

Modeling check-worthiness in news articles both as sentence classification and a sentence ranking tasks

Using argument structure as context yields better results than using local discourse context for the task of check-worthiness detection

***Future Work:***

1. Predict argument components and relations and compare with using gold annotations

2. Investigate other reasons for check-worthiness not related to argument structure other argument fallacies: e.g. cherry-picking and strawman argument

Tariq Alhindi, Brennan Xavier McManus, and Smaranda Muresan. 2021. What to fact-check: Guiding check-worthy information detection in news articles through argumentative discourse structure. In *The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL 2021)*

# Thank You

www.cs.columbia.edu/~tariq
www.github.com/tariq60/whatToFactCheck

SIGDIAL 2021

The 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue, 29-31 July 2021, Singapore