

Automating Comment Moderation: Topics and Toxicity in Online News

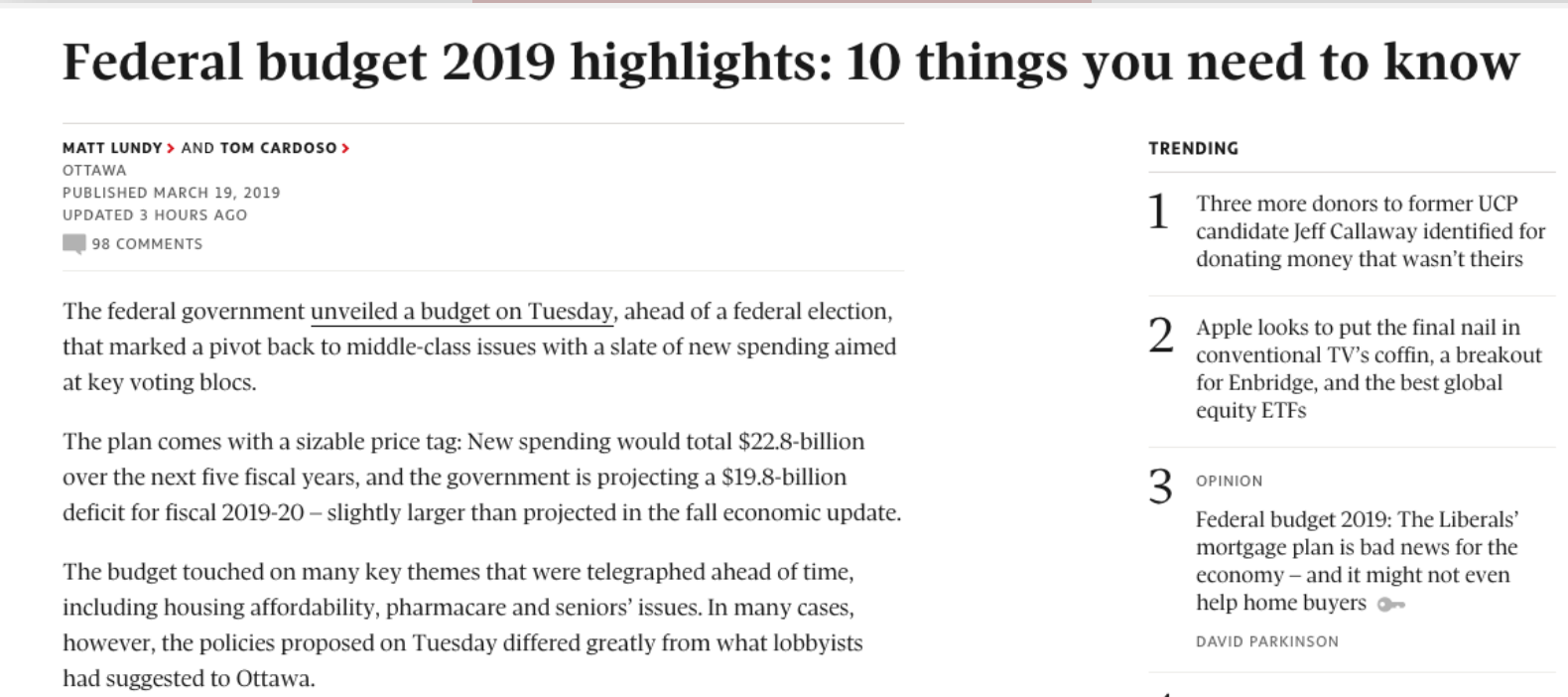
SFU

Vagrant Gautam • Maite Taboada

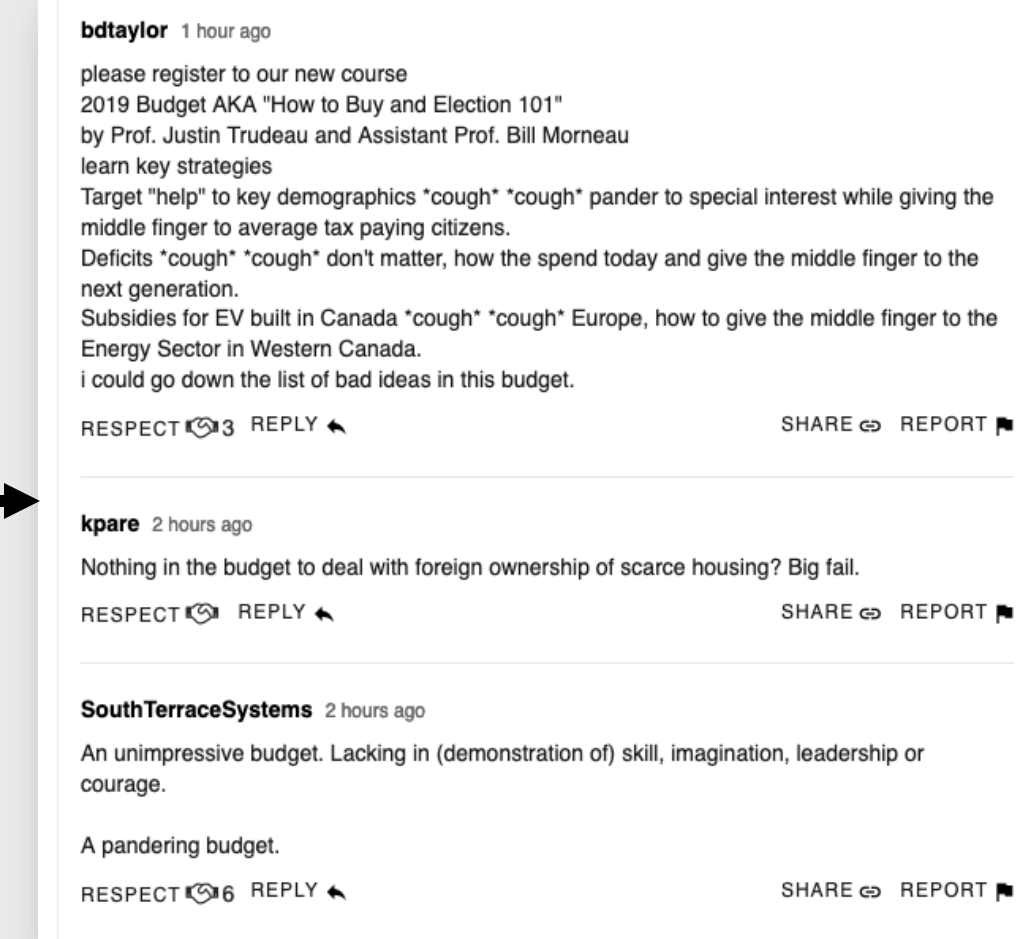


1. Online News

Articles



Comments



- Online publications with comments sections
- Human moderators hired to monitor comments
- **Positive** moderation: Highlight good arguments
- **Negative** moderation: Delete profanity, personal attacks

THE CONVERSATION

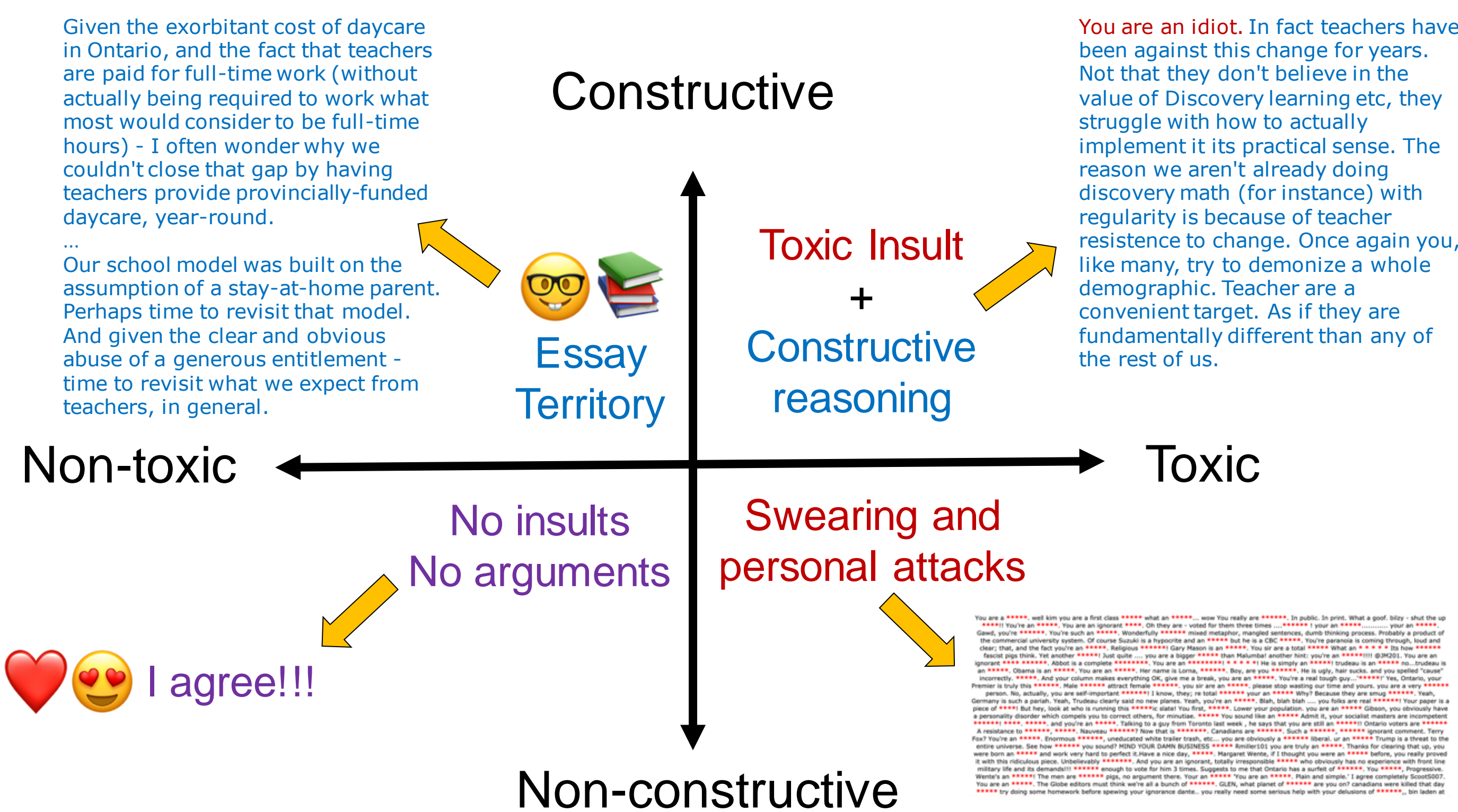
THE GLOBE AND MAIL



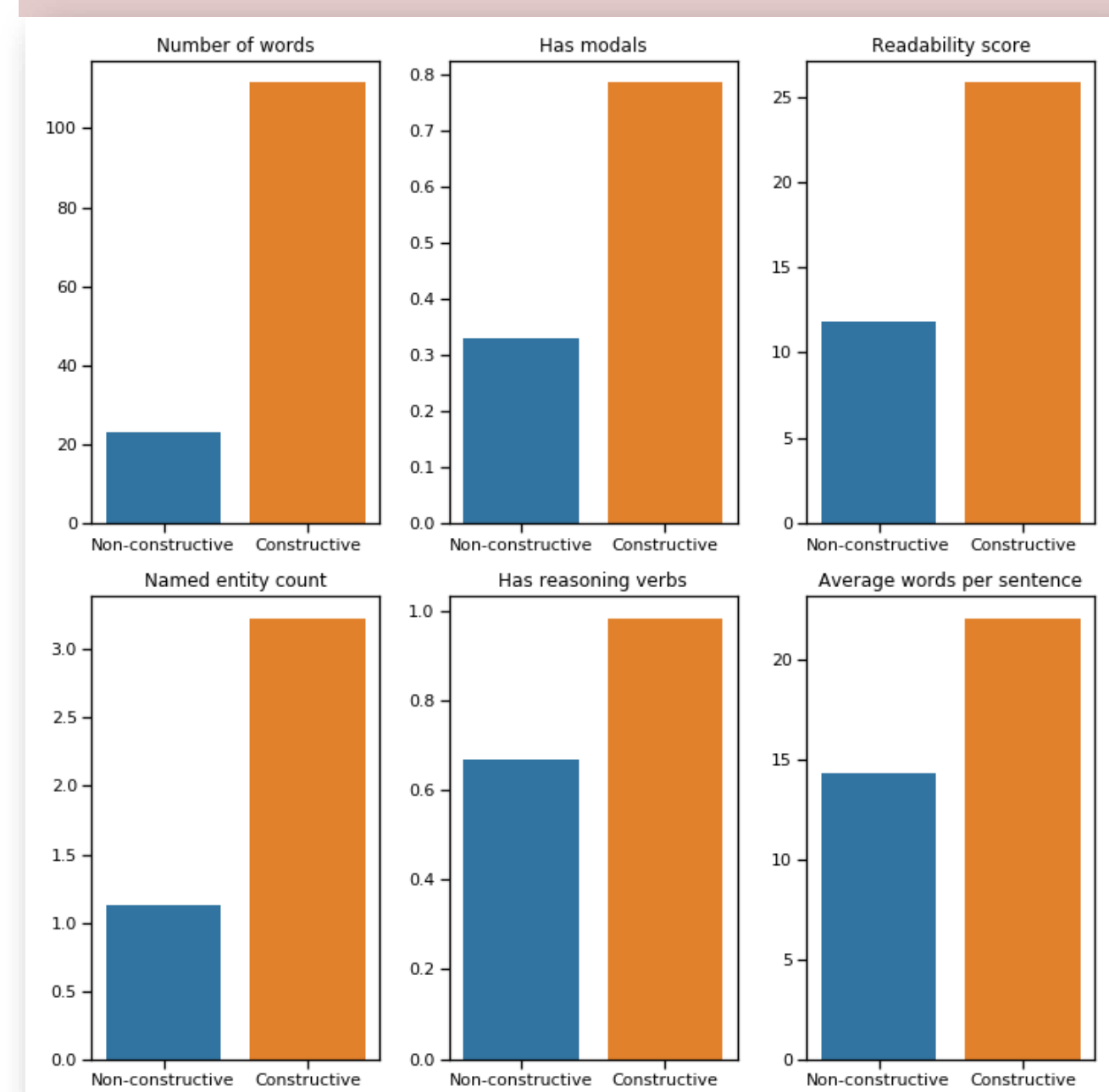
Questions:

1. What makes comments **linguistically** different?
2. Can we use **computer science** to automate moderation?

2. Taxonomy of the Online Comment



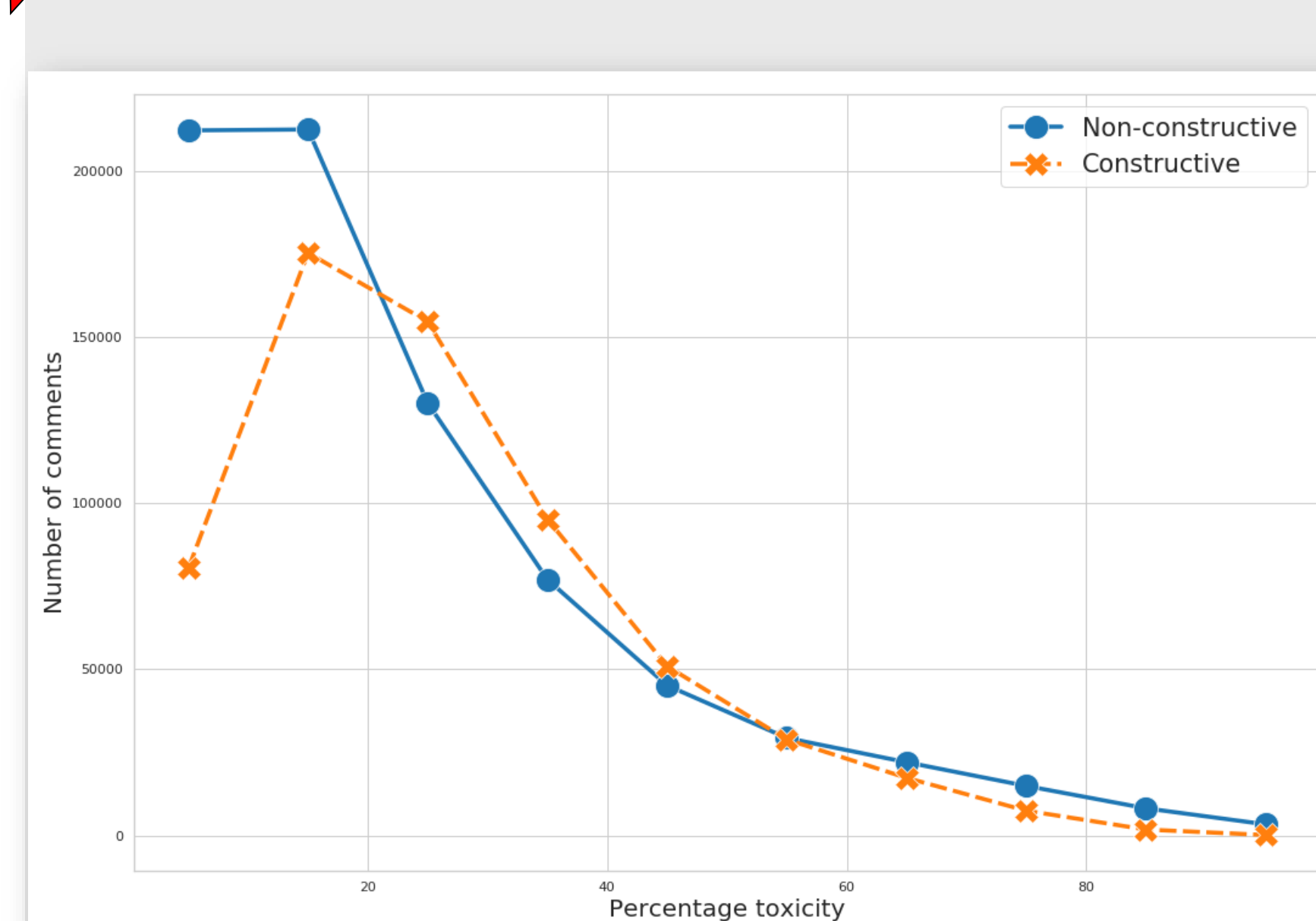
Quantifying constructiveness and toxicity



- **Constructiveness system** developed at SFU by Varada Kolhatkar and Maite Taboada
- Comments classified for constructiveness using a number of linguistic features, some of which are shown on the left
- **Toxicity system** developed by Google: Perspective API
- Uses machine learning (RNNs with attention)



3. Toxicity and Constructiveness



- Results on 1,000,000+ online comments
- Highest frequency comments are non-constructive and non-toxic
- Very few are highly toxic
- Non-constructive comments more common (intuitive)

4. Topic Modelling on Articles and Comments

Going beyond toxicity and constructiveness...

- What subjects generate the most comments?
- What subjects generate the most toxic comments? Or the most constructive comments?

TOPIC MODELLING

Statistical modelling technique to automatically extract the subjects discussed in a text

Example topics

partyNDP
Liberals
government
spending

Alberta
energy
pipeline
sands oil

aboriginal
communities
Crown
reserves
protection
natural

students
panel
woman
University
women

- Latent Dirichlet Allocation to generate 15 automatically extracted topics on 10,000+ news articles and 1,000,000+ comments

- Each text is given probabilities for all 15 topics – this allows an article to be 50% about a new pipeline and oil, and 50% about the impact on nature and nearby First Nations communities

- **Hypothesis 1**
People would talk a lot about politics; more comments on articles classified as "politics" articles
- **Hypothesis 2**
People would be significantly more toxic about certain topics, e.g., issues of abortion, race, atheism

Results

Highest frequency words



Constructive comments



Non-constructive comments

- Most frequent words across all comments regardless of toxicity and constructiveness: Harper, time, people, government, Canada

➤ **Hypothesis 1 confirmed**

Most comments about politics and on politics articles; people talk more than anything else about politics

➤ **Hypothesis 2 rejected**

Roughly the same proportion of toxic commenters in every comment section

Future work

- Do toxicity and constructiveness propagate in threads?
- Add sentiment to the taxonomy
- Do people with anonymous usernames write different comments from people with their real names? More toxic? Less constructive?

References

- [1] Cao, N., & Cui, W. (2016). *Introduction to Text Visualization*. Atlantis Press.
- [2] Kolhatkar, V., & Taboada, M. (2017). *Constructive Language in News Comments*. Proceedings of the First Workshop on Abusive Language Online, ACL. Vancouver. 11-17.
- [3] Kolhatkar, V., & Wu, H., & Cavasso, L., & Francis, E., & Shukla, K., & Taboada, M. (2018). *The SFU Opinion and Comments Corpus: A Corpus for the Analysis of Online News Comments*.
- [4] Rubin, T. N., & Chambers, A., & Smyth, P., & Steyvers, M. (2011). *Statistical Topic Models for Multi-Label Document Classification*. Mach Learn. 88.
- [5] Newman, D., & Bonilla, E., & Buntine, W. (2011). *Improving Topic Coherence with Regularized Topic Models*. NIPS. 496-504.
- [6] Röder, M., & Both, A., & Hinneburg, A. (2015). *Exploring the Space of Topic Coherence Measures*. WSDM 2015 - Proceedings of the 8th ACM International Conference on Web Search and Data Mining. 399-408.

Acknowledgements

Thanks to Fatemeh Torabi Asr for her guidance when I hit a topic modelling wall, to Varada Kolhatkar for her beautiful code, to all the people behind the corpora, NLP systems and open-source Python and R libraries I used. Thanks to Yue Wang for the original poster template. Thanks to Ashley Farris-Timble for telling Maite to take me on!
Funding: Key Big Data Undergraduate Student Research Award (USRA)