## MITIGATING SCHOLARLY CORPUS BIASES WITH CITATIONS: A CASE STUDY ON CORD-19

MICROSOFT RESEARCH, REDMOND WA & ALLEN INSTITUTE FOR ARTIFICIAL INTELLIGENCE, SEATTLE WA



## Mitigating Biases in CORD-19 for Analyzing COVID-19 Literature Provisionally accepted The final, formatted version of the article will be

published soon. 🔽 Notify me



<sup>1</sup>Microsoft Research (United States), United States <sup>2</sup>Allen Institute for Artificial Intelligence, United States

#### **FRONTIERS IN RESEARCH METRIC AND ANALYTICS**

Special Topic on

Coronavirus Research Landscape: Resources, Utilities, and Analytic Studies

DOI: 10.3389/frma.2020.596624

#### **ETHICAL AI: AN URGENT TOPIC TO OUR SOCIETY**



# THE WALL STREET JOURNAL.

*By Michael Totty* Nov. 3, 2020 10:00 am ET

JOURNAL REPORTS: TECHNOLOGY

#### How to Make Artificial Intelligence Less Biased

Al systems can unfairly penalize certain segments of the population especially women and minorities. Researchers and tech companies are figuring out how to address that.



The AI world is making a strong push to root out bias in AI systems, but it faces some significant obstacles. KEITH A. WEBB AND IMAGES FROM ISTOCK

#### **TACKLING INFORMATION OVERFLOW WITH AI**

#### CORD-19: The COVID-19 Open Research Dataset

Lucy Lu Wang<sup>1,\*</sup> Kyle Lo<sup>1,\*</sup> Yoganand Chandrasekhar<sup>1</sup> Russell Reas<sup>1</sup> Jiangjiang Yang<sup>1</sup> Douglas Burdick<sup>2</sup> Darrin Eide<sup>3</sup> Kathryn Funk<sup>4</sup> Yannis Katsis<sup>2</sup> Rodney Kinney<sup>1</sup> Yunyao Li<sup>2</sup> Ziyang Liu<sup>6</sup> William Merrill<sup>1</sup> Paul Mooney<sup>5</sup> Dewey Murdick<sup>7</sup> Devvret Rishi<sup>5</sup> Jerry Sheehan<sup>4</sup> Zhihong Shen<sup>3</sup> Brandon Stilson<sup>1</sup> Alex D. Wade<sup>6</sup> Kuansan Wang<sup>3</sup> Nancy Xin Ru Wang<sup>2</sup> Chris Wilhelm<sup>1</sup> Boya Xie<sup>3</sup> Douglas Raymond<sup>1</sup> Daniel S. Weld<sup>1,8</sup> Oren Etzioni<sup>1</sup> Sebastian Kohlmeier<sup>1</sup>

<sup>1</sup>Allen Institute for AI <sup>2</sup> IBM Research <sup>3</sup>Microsoft Research <sup>4</sup>National Library of Medicine <sup>5</sup>Kaggle <sup>6</sup>Chan Zuckerberg Initiative <sup>7</sup>Georgetown University <sup>8</sup>University of Washington {lucyw, kylel}@allenai.org

- <u>3500+ articles/week by Mid March</u>
  - 372,698 articles as of November 15,2020
- Full text corpus: first released on March 16, 2020
  - Publisher contributions + archival services
- Activities
  - Open QA Challenge on Kaggle
  - TREC tracks at NIST
- Methodology: keyword query into various databases

"COVID" OR "COVID-19" OR "CORONAVIRUS" OR "2019-nCov" OR "SARS-COV" OR "MERS-COV" OR "Severe Acute Respiratory Syndrome" OR "Middle East Respiratory Syndrome"

#### **ASSESSING BIASES BY <u>CORPUS EXPANSION</u> WITH CITATIONS**

- Price, D., "<u>Networks of Scientific Papers</u>", Science 1965
- Notable network science models:
  - Preferential attachment, Albert and Barabasi, Science 1999
  - Individual node fitness, Caldarelli et al., *Physical Review Letters*, 2002
  - Latent space model, Papadopoulos et al., *Nature* 2012
  - Discrete choice, Overgoor et al., WWW-2019
- Methods evaluated:
  - Enclosure graph (Sinatra et al., "A century of physics", Nature Physics, 2015)
  - Closure graph: updated data available at Github: <u>https://aka.ms/magcord19mapping</u>

### **ENCLOSURE GRAPH**

- From a seed collection of articles:
  - Expand to citing and cited (or both) articles
  - CORD-19 => CORD-19E
- **Observations:** 
  - Bidirectional, single step traversal on citation networks
  - Cannot be made into iterative algorithm without topic overrun
    - Ex: citing a "tool" paper
    - Sensitive to seed quality



convolutional neural network

#### **CLOSURE GRAPH**

- Uni-direction, multiple step traversals
  - Continue until all references are in the expanded set
  - Mostly directed acyclic graph
- Motivations
  - Cover more background and lineage of knowledge
  - Discrete (discreet) Choice => Topic overrun less severe
- Surprise finding
  - Don't need full closure!
  - "Inflection" point seems to capture themes
- CORD-19C: full closure; CORD-19I: inflection closure



## **CLOSURE GRAPH GROWTH VS HOPS**



Paper Count — Accumulated Citations

#### **EMBEDDEDNESS: CITATIONS RECEIVED FROM WITHIN COLLECTION**



#### **TOPIC COVERAGE**



Biology

Medicine Chemistry

Other

#### **ARTICLE AGE DISTRIBUTION**



#### **DISTRIBUTION OF ARTICLE IMPORTANCE**

\*Details on saliency can be found here



PUBLICATION YEAR

#### JOURNAL COVERAGE/RANKINGS: CORD-19 VS CORD-19E



#### JOURNAL COVERAGE/RANKING: CORD-19 VS CORD-19I/C



#### WHERE IS RESEARCH CONDUCTED





CORD-19C



**Publication Year** 

#### WHERE IS RESEARCH CONDUCTED





CORD-19E

CORD-19C



Publication Year

■ Africa ■ Americas ■ Asia ■ Europe ■ Oceania ■ Other

#### **TEAM SIZE**



■1-3 ■4-6 ■7-9 ■10-20 ■above 20 — AvgTeamSize — MedianTeamSize

#### CONCLUSIONS

- 1. Citations complement retrieval techniques
- 2. Analytics between CORD-19 and CORD-19E/I/C
  - Different: Topics, Age, Article/Journal impacts
  - Similar: Collaboration trends
- 3. Corpus expansion offers statistically smoother analytics
- 4. Closure graph findings support network science theories
  - Preferential attachments with fitness, discrete choice
- 5. Partial "inflection" closure is almost as good as full closure
- 6. Check out and build on our GitHub share <u>https://aka.ms/magcord19mapping</u>
- 7. Future studies
  - Not all citations are equal: need citation classification?
  - Probabilistic network traversal