EMNLP 2020 Plenary Panel Discussion:

Publishing in the Era of Responsible AI: How Can we be Proactive? Considerations and Implications.

Emily M. Bender, Rosie Campbell, Allan Dafoe, Pascale Fung, Meg Mitchell, Saif M. Mohammad

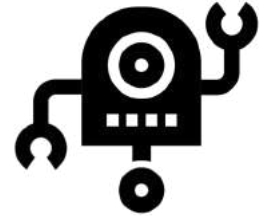# What is a Research Ethics Statement and Why does it Matter?

## Saif M. Mohammad

Senior Research Scientist, National Research Council Canada

✉ Saif.Mohammad@nrc-cnrc.gc.ca    🐦 @SaifMMohammad

National Research Council Canada    Conseil national de recherches Canada

Canada

As NLP and ML systems become more ubiquitous, their broad societal impacts are receiving more scrutiny than ever before.

Several high-profile instances have highlighted how technology will often lead to more adverse outcomes for those that are already marginalized.



Created by Oksana Latysheva
from Noun Project

**Do Machines Make Fair Decisions?**

What part do we play in this as researchers?

What are the hidden assumptions in our research?

What are the unsaid implications of our choices?

Whose voices are we amplifying? (and whose we are not?)

Are we perpetuating and amplifying inequities
or are we striking at the barriers to opportunity?



**Do People Make Fair Decisions?**

Answers are often complex and multifaceted.

Ethics statements can help navigate research choices, communicate implications.

**2021 Annual Conference of the North American Chapter of the Association for Computational Linguistics**

Mexico City, Mexico
June 6–11, 2021

Authors will be allowed extra space after the 8th page for a broader impact statement or other discussion of ethics. The NAACL review form will include a section addressing these issues and papers flagged for ethical concerns by reviewers or ACs will be further reviewed by an ethics committee. Note that an ethical considerations section is not required, but papers working with sensitive data or on sensitive tasks that do not discuss these issues will not be accepted. Conversely, the mere inclusion of an ethical considerations section does not guarantee acceptance. In addition to acceptance or rejection, papers may receive a conditional acceptance recommendation. Camera-ready versions of papers designated as conditional accept will be re-reviewed by the ethics committee to determine whether the concerns have been adequately addressed. Please read the ethics FAQ for more guidance on some problems to look out for and key concerns to consider relative to the code of ethics.

National Research Council Canada    Conseil national de recherches Canada

Canada

The Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021)

## Ethics Policy

Authors are required to honour the ethical code set out in the ACL Code of Ethics.

The consideration of the ethical impact of our research, use of data, and potential applications of our work has always been an important consideration, and as artificial intelligence is becoming more mainstream, these issues are increasingly pertinent. We ask that all authors read the code, and ensure that their work is conformant to this code. Where a paper may raise ethical issues, we ask that you include in the paper an explicit discussion of these issues, which will be taken into account in the review process. We reserve the right to reject papers on ethical grounds, where the authors are judged to have operated counter to the code of ethics, or have inadequately addressed legitimate ethical concerns with their work.

National Research Council Canada    Conseil national de recherches Canada

# what goes into Ethics Statements (ethical considerations)

*… not "Appendix material"*

*… not just "good to have"*

*… not something we have never seen before*

# what goes into Ethics Statements (ethical considerations)

*… central to our work*

*… things we have always seen in good work*
*(usually sprinkled across various sections of a paper)*

# What is a good place to talk about ethical considerations?

- Introduction/Motivation ⟶ Impact Statement
- Related Work
- Data ⟶ Data Statement
- Methodology
- Evaluation ⟶ Bias mitigation, Ethics focused shared task (e.g. SemEval 2018 Task 1)
- Experiments
- Error Analysis ⟶ Traditionally, a place where people have listed
- Limitations ⟶ ethical considerations
- Conclusions
- Future Work

Ethics Statements can bring together the ethical considerations stated in the paper in a cohesive narrative, and elaborate on them.

# Ethical Considerations: Introduction/Motivation

- what questions are we asking?
  - no hurry to get to the solution
  - lets understand the question better
    - bring out the nuances and complexities
    - what assumptions are we making?
- why should we care about this question/problem/task?
- who is impacted by this problem?
- who is impacted by this work?
- who is left out?
- what problems are we not tackling?
- are these choices maintaining structural inequities or questioning them?

National Research Council Canada   Conseil national de recherches Canada

Canada

# Ethical Considerations: Related Work

Whose voices are we amplifying?

NLP is actively encroaching on other fields:

- humanities
- psychology
- culture studies
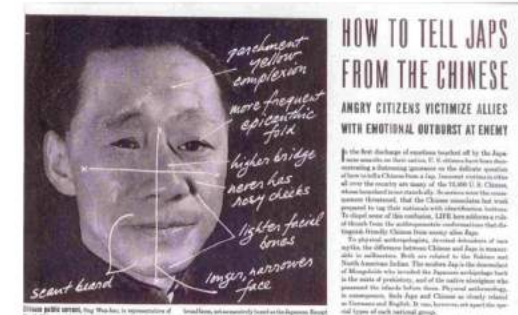- social sciences
- public health



**Lessons from Archives: Strategies for Collecting Sociocultural Data in Machine Learning**

Eun Seo Jo
Stanford University
eunseo@stanford.edu

Timnit Gebru
Google
tgebru@google.com

**ABSTRACT**

A growing body of work shows that many problems in fairness, accountability, transparency, and ethics in machine learning systems are rooted in decisions surrounding the data collection and annotation process. In spite of its fundamental nature however, data collection remains an overlooked part of the machine learning (ML) pipeline. In this paper, we argue that a new specialization should be formed within ML that is focused on methodologies for data collection and annotation: efforts that require institutional frameworks and procedures. Specifically for sociocultural data, parallels can be drawn from archives and libraries. Archives are the longest standing communal effort to gather human information and archive

We must give researchers from these fields a voice. Learn from them. Situate our work in their literature. Collaborate with them.

That does not mean they have all the answers.

Your NLP background gives you unique perspective.

Bring it to bear by collaborate with all stake holders (including the people affected). Avoid the trap of AI/ML/NLP solutionism.

National Research Council Canada    Conseil national de recherches Canada

Canada

# Ethical Considerations: Related Work

Whose voices are we amplifying?

NLP is actively encroaching on other fields:

- humanities
- psychology
- culture studies
- social sciences
- public health

**COI:** The incentives for fast science act against the careful and thoughtful pace of slow (truly interdisciplinary) science that engages with all stake holders right from the start.

We can benefit from both slow and fast science.

We must give researchers from these fields a voice. Learn from them. Situate our work in their literature. Collaborate with them.

That does not mean they have all the answers.

Your NLP background gives you unique perspective.

Bring it to bear by collaborate with all stake holders (including the people affected). Avoid the trap of AI/ML/NLP solutionism.

# Ethical Considerations: Data

A whole panoply of considerations. Here is one:

For annotations, is there a "right" answer and a "wrong"?

- Yes: domain experts annotate the data
- No: we want to know how people perceive this word, phrase, sentence, etc.
  - large number of annotators
  - seek appropriate demographic information (respectfully and ethically)

How should we aggregate the information?

- Acknowledge the limitations of majority vote aggregation
- Danger of saying that the views of a certain demographic is the norm or standard
- Acknowledge that we are missing out on some/many voices
- Saying all voices are correct has its own problems
  - How to address and manage inappropriate biases?

National Research Council Canada    Conseil national de recherches Canada
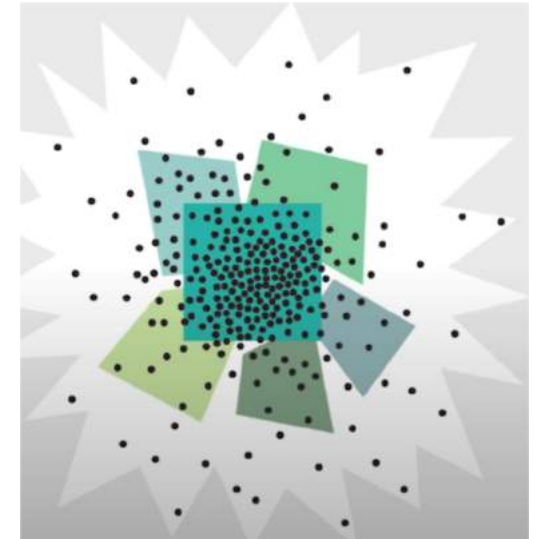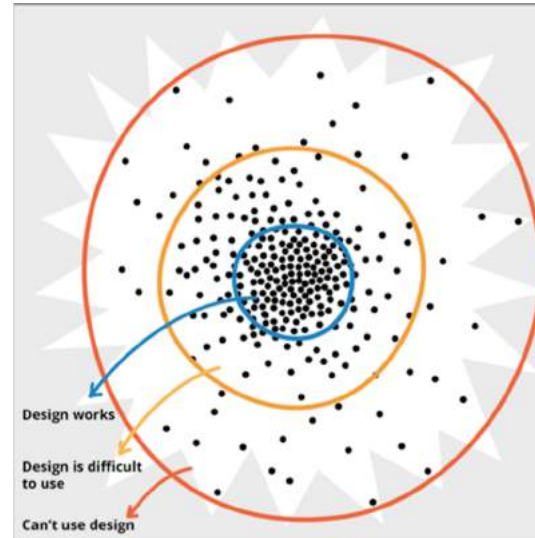
Canada

# Ethical Considerations: Design and Methodology

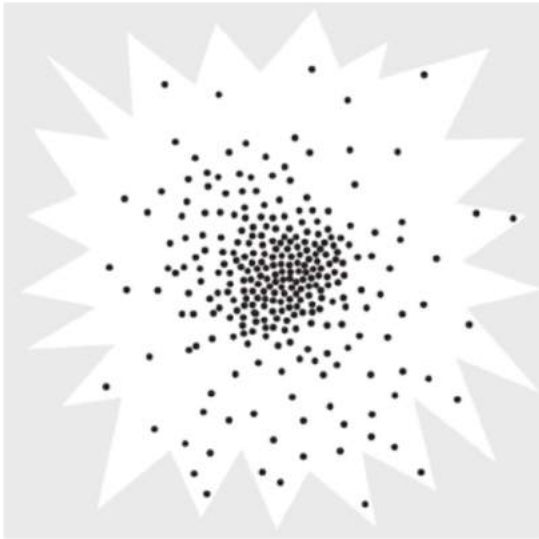People have a bias towards large numbers.
AI amplifies the bias to large numbers.
-- Jutta Treviranus (Expert on inclusive design)
We Count! https://www.youtube.com/watch?v=OAXmCAqZqRk

Pareto principle, 80/20 rule,
Zipf's law, power law distribution



Design works

Design is difficult to use

Can't use design

Multi-variate scatter plot of needs of a set of people. Source: Jutta Treviranus.
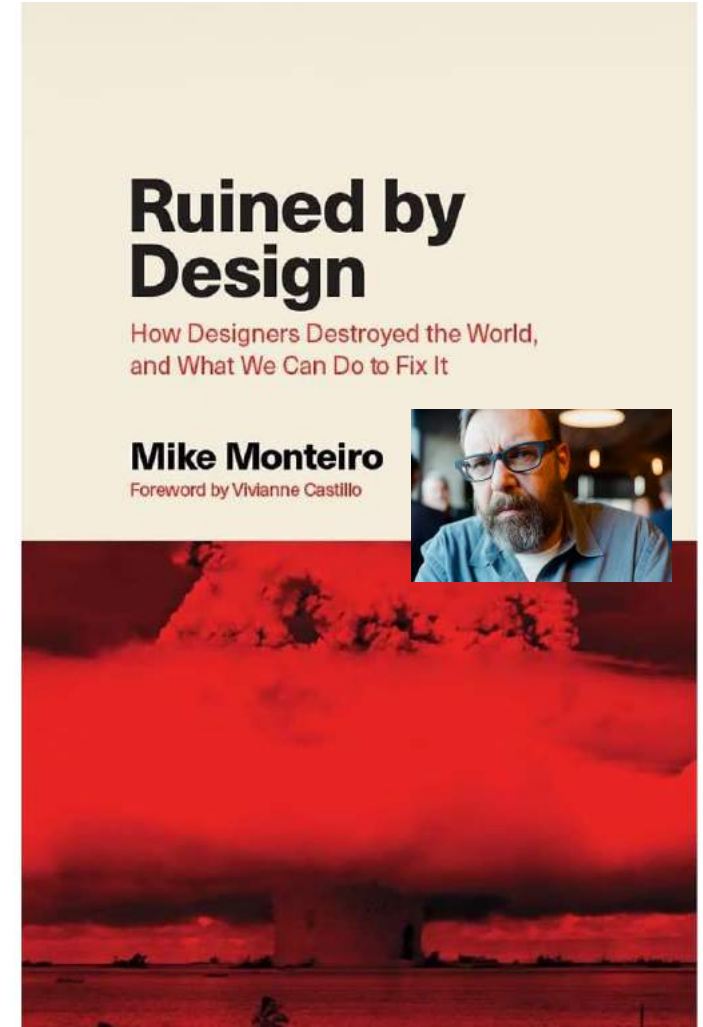
# Ethical Considerations: Design and Methodology

Yes, design is political. Because design is labor, and your labor is political. Where you choose to expend your labor is a political act. Who we omit from those solutions is a political act.

A designer does not believe in edge cases.

This job isn't about creating bullhorns for fascists and others who'd use their power to denigrate others. It's about making sure those who are threatened by the inhumane have the better bullhorns.

This job isn't about building tools that hand our data to the corporations of Silicon Valley. It's about building tools to keep that data from them.

-- Mike Monteiro

# Ethics Statement: Tips

- Acknowledge
  - the pain of those affected
  - your biases and conflicts of interest
  - how you have / have not involved various stake holders
  - you cannot capture everything

- Invite feedback and things to add
  - blog posts, open review, preprints

  *you can still miss things!*

  - live document

- Space limitations
  - put it on the project webpage

## The State of NLP Literature: Part IIIa

Impact — Most Cited Papers and Aggregate Citation Metrics by Time, Paper Type, and Venue

Saif M. Mohammad  Nov 1, 2019 · 15 min read

This series of posts present a diachronic analysis of the ACL Anthology — Or, as I like to think of it, making sense of NLP Literature through **pictures**.

In the rain forest of scientific papers, citations can be like sunlight — making a piece of work more prominent. Photo credit: Pascal van de Vendel

## Practical and Ethical Considerations in the Effective use of Emotion and Sentiment Lexicons

Saif M. Mohammad
National Research Council Canada
Ottawa, Canada
saif.mohammad@nrc-cnrc.gc.ca.

### 4. CONCLUDING REMARKS

Emotion lexicons can be simple yet powerful tools to analyze text. However, use of the lexicons (even for tasks that it is suited for) can lead to inappropriate bias. Applying a lexicon to any new data should only be done after first investigating its suitability, and requires careful analysis to minimize unintentional harm. I listed some considerations above that can help mitigate such unwanted outcomes. However, these are not meant to be comprehensive, but rather a jumping off point for further thought. The author welcomes feedback; including additional points to consider and include in this document. See also the

## Caveats, Limitations, and Ethical Considerations

NLP Scholar comes with several caveats, limitations, and ethical considerations as listed below.

### Aspects of Analysis

- The analyses presented in The State of NLP Literature posts cover only *some* aspects of the literature. Prior work has explored other aspects such as citation link analysis, co-author networks, influence, types of citations, etc. Yet, several interesting questions remain unexplored.

### Accessing Information about the Papers

- Google does not provide an API to extract information about the papers. Martín-Martín et al. (2018) and others have pointed out that this is likely because of its agreement with publishing companies that have scientific literature behind paywalls. The ACL Anthology is in the public domain and free to access. We extracted citation information from Google Scholar profiles of people who published in the ACL Anthology. This is explicitly allowed by their robots exclusion standard, and is how past work has studied Google Scholar:
  — Martín-Martín, A., Orduna-Malea, E., Thelwall, M. and López-Cózar, E.D., 2018. Google Scholar, Web of Science, and Scopus: A systematic comparison of citations in 252 subject categories. *Journal of Informetrics*, 12(4), pp.1160–1177.
  — Khabsa, M. and Giles, C.L., 2014. The number of scholarly documents on the public web. *PloS one*, 9(5), p.e93949.
  — Orduña-Malea, E., Ayllón, J.M., Martín-Martín, A. and López-Cózar, E.D., 2014. About the size of Google Scholar: playing the numbers. *arXiv preprint arXiv:1407.6239*.

### Errors

- Even though the ACL Anthology and Google Scholar are outstanding resources, they contain some errors. Also, aligning information from the two resources can never be perfect. (More details in the bullets below.) Thus NLP Scholar is bound to include some errors. We apologize for any misrepresentations, and will fix things as best we can.

### Inconsistencies and Missing Values in the ACL Anthology

Information in the ACL Anthology is not always consistent and some attributes may be missing:

- The same venue may be described in different ways.
- There is no consistent way to identify short papers, tutorials, demo papers, book reviews, etc. Main conference short papers are sometimes clearly marked in the booktitle field of the BibTex, but at other times, they are not distinguished from the long papers. Occasionally, they are marked in other idiosyncratic ways such as appending "(short paper)" to the paper title.

---

marked in other idiosyncratic ways such as appending "(short paper)" to the paper title.

- Some papers have a missing author field in the BibTeX entry. These papers are omitted. (These are often proceedings, lists of tutorials, etc. that we would want to omit anyway.)
- For some papers, the title in the BibTeX entry uses non-accented letters, even though the title has accented letters. For example, the title is recorded to have the word "sémantique" in the main records of AA, it is written as "semantique" in the BibTeX entry in AA. We use the BibTeX entry to extract author names, and the mismatch in titles causes the system to not find the authors. Papers with missing values for authors are omitted.

We use high precision heuristics to identify necessary information. However, note that there will be some number of omissions and misclassifications.

### Citation Information From Google Scholar

- Google Scholar is used widely in research. However, it has received criticisms regarding the amount of curation, reducing academic worth to citations and h-index, etc. (see Criticisms of the Citation System, and Google Scholar in Particular, How Has Google Scholar Changed Academia?, 4 reasons why Google Scholar isn't as great as you think it is).
- Some number of papers exist such that none of their authors created a Google Scholar profile. We do not have citation information for those papers. Such papers are still displayed in NLP Scholar — only their citation information has a null value. This, however, means that in terms of citation information, it is likely that work done in the past is under-represented (as authors who left academia or retired may be less likely to create a Google Scholar profile). Nonetheless, we do not expect this to markedly impact the inferences drawn from the analyses presented, as we do have citation information for over 35,000 papers.

### Aligning Information in AA and Google Scholar is Tricky

- they do not have a common paper id or author id
- occasionally two different papers have the same title
- the same author may use different forms of their name in different articles
- multiple authors might have the same name

We use the paper title and publication year combination as the unique identifier for a paper. However, there are some pairs of papers that have the same title and year of publication. These are omitted.

### New Papers are Constantly Added to AA.

The current instantiation of NLP Scholar is based on the papers in AA as of

---

The current instantiation of NLP Scholar is based on the papers in AA as of June 2019. We will update NLP Scholar with new AA information periodically.

### Papers Receive More Citations with Time.

The current instantiation of NLP Scholar is based on the citations papers received as of June 2019. We will update NLP Scholar with new citations information periodically.

### Rich get Richer

Visualizations in NLP scholar present papers with more citations more prominently than papers with fewer citations. This can have the effect of making highly cited papers even more cited. (This is not unlike Google Scholar, which also ranks papers by relevance and citation counts.) Citations are one (somewhat noisy) indicator of the amount of impact a paper has had. While they can be useful to find interesting and impactful papers, it must be noted that papers get cited for a number of other reasons as well, and it is entirely possible that some of papers of interest might be those that are less cited.

There are several ways in which NLP Scholar can cast light on less cited papers too though. Here are some examples:

- By showing the papers on a timeline, one can easily track papers that influenced a high-citation paper in an area.
- When searching for papers in an area, one can compare citations of papers within that area. This places a target paper in a more appropriate context. For example, a target paper may not have received hundreds of citations, but one can see that the within the area of research, it is one of the most highly cited papers.
- The Languages visualizations highlight work in various languages.

### Search based on Words in Titles

- Even though there is an association between terms and areas of research, the association can be less strong for some terms. I use the association as one (imperfect) source of information about areas of research. This information may be combined with other sources of information to draw more robust conclusions. Planned future work on allowing searches for terms in abstracts and whole papers, as well as finding documents related to a query term based on word embedding based document representations will alleviate the current limitations. However, it should be noted that search based on title words is a simple and powerful method for finding relevant documents.

### Demographics

- Data is often a representation of people (Zook 2017). This is certainly the case here and we acknowledge that the use of such data has the

---

- Data is often a representation of people (Zook 2017). This is certainly the case here and we acknowledge that the use of such data has the potential to harm individuals. This work has several limitations, and some have ethical considerations in terms of who is left out. Further, while the methods used are not new, their use merits reflection.
- Analysis focused on women and men leaves out non-binary people. Not disaggregating cis and trans people means that the statistics are largely reflective of the more populous cis class. We hope future work will explore gender gaps between non-binary — binary, trans — cis, etc. Similarly, tracking the skew in authors of diverse income, experiences, and abilities is also crucial. This work does address those but hopefully more work on those will follow.
- The use of female- and male-gender associated names to infer population level statistics for women and men, can reinforce harmful stereotypes and is exclusionary to people that do not have such names, to people from some cultures where names are not as strongly associated with gender, and trans people who have not been able to change their name.
- Since the names dataset used is for American children there is lower representation of names from other nationalities. However, many names are common in more than one country, and the large immigrant population in the US means that there still exists substantial coverage of names from around the world.
- Chinese names (especially in the romanized form) are not good indicators of gender. Thus the method presented here disregards most Chinese names, and the results of the analysis do not apply to researchers with Chinese names.
- Some might argue that names partially address the gender inclusiveness guidelines listed in (Keyes 2018): names can be changed to indicate (or not indicate) gender, people can choose to keep their birth name or change it, and the name, more so than appearance, can be independent of physiology. However, changing names can be quite difficult. Also, names do not capture gender fluidity or contextual gender.
- A more inclusive way of obtaining gender information is through optional self-reported surveys. However, even if one allows for a self-report checkbox so that the respondent can have the primacy and autonomy to express gender, downstream data science either ignores such data or combines information in ways that are not in control of the respondent. Further, as is the case here, it is not easy to obtain self-reported historical information.
- A small number of names change association from one gender to another with time. We hope that the ≥99% rule filters them out, but this is not guaranteed.
- Social category detection can potentially lead to harms, for example, depriving people of opportunities simply because of their race or gender. However, one can also see the benefits of NLP techniques and social category detection in public health (e.g., developing targeted initiatives to improve health outcomes of vulnerable populations), as

# Conflict of Interest



The Grey Hoodie Project: Big Tobacco, Big Tech, and the threat on academic integrity

Mohamed Abdalla
msa@cs.toronto.edu
Centre for Ethics, University of Toronto

**ABSTRACT**

As governmental bodies rely on academics' expert advice to shape policy regarding Artificial Intelligence, it is important that these academics not have conflicts of interests that may cloud or bias their judgement. Our work explores how Big Tech is actively distorting the academic landscape to suit its needs. By comparing the

**World Health Organization**

**Tobacco Free Initiative (TFI)**

**Article 5.3 of the WHO Framework Convention on Tobacco Control**

In setting and implementing their public health policies with respect to tobacco control, Parties shall act to protect these policies from commercial and other vested interests of the tobacco industry in accordance with national law.

Parallels between Big Tech and Big Tobacco: **histories, actions**

- appearing to be embattled, hiring (researchers) and funding (universities, conferences), taking ownership of ethics, influencing research questions

# Resources

- **A Guide to Writing the NeurIPS Impact Statement**. Carolyn Ashurst, Markus Anderljung, Carina Prunkl, Jan Leike, Yarin Gal, Toby Shevlane, Allan Dafoe.
  https://medium.com/@GovAI/a-guide-to-writing-the-neurips-impact-statement-4293b723f832

- **Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science.** Emily M. Bender and Batya Friedman.

- **Datasheets for datasets**. Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford.

- Some of my ethics statements:

  ◦ Practical and Ethical Considerations in the Effective use of Emotion and Sentiment Lexicons.

  ◦ NLP Scholar Project: https://medium.com/@nlpscholar/about-nlp-scholar-62cb3b0f4488

**Contact:** ✉ Saif.Mohammad@nrc-cnrc.gc.ca 🐦 @SaifMMohammad

National Research Council Canada    Conseil national de recherches Canada

Canada