



Ethics and Natural Language Processing (NLP)

Karën Fort

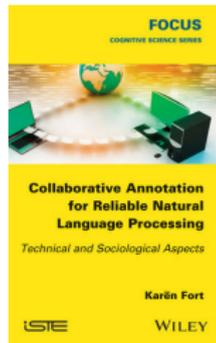
karen.fort@loria.fr / <https://members.loria.fr/KFort/>

DiLCO, Oct. 8th, 2021

Where I'm talking from

See <https://members.loria.fr/KFort/>

- ▶ Language resources creation for NLP, esp. using crowdsourcing



- ▶ Ethics and NLP



What are we talking about?

Virtue Ethics

Deontological Ethics

Utilitarianism and consequentialism

Natural Language Processing (NLP)

Why is it important?

Beyond biases

How to limit the damages?

Ethics in general vs in the community

Merriam-Webster SINCE 1828

GAMES | BROWSE THESAURUS | WORD OF THE DAY | WORDS AT PLAY

ethic

[Dictionary](#) [Thesaurus](#)

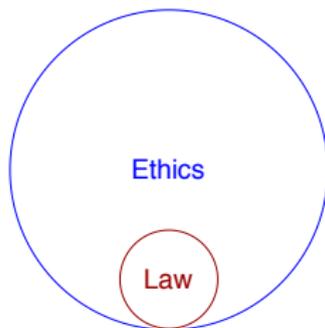
eth·ic | \ è-thik \

Definition of *ethic*

- 1** **ethics** *plural in form but singular or plural in construction* : the discipline dealing with what is good and bad and with moral duty and obligation
- 2** **a** : a set of moral principles : a theory or system of moral values
 - // the present-day materialistic *ethic*
 - // an old-fashioned work *ethic*
 - often used in plural but singular or plural in construction
 - // an elaborate *ethics*
 - // Christian *ethics*
- b** **ethics** *plural in form but singular or plural in construction* : the principles of conduct governing an individual or a group
 - // professional *ethics*

Ethics is not law

Right to do things vs doing what is right



Law: sets minimum standards (rules and regulations)

VS

Ethics: sets maximum standards

What are we talking about?

Virtue Ethics

Deontological Ethics

Utilitarianism and consequentialism

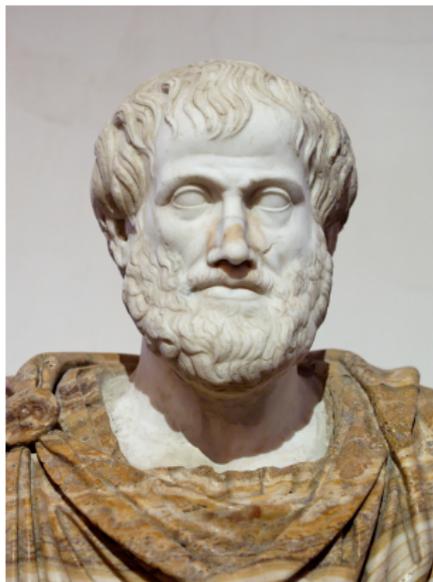
Natural Language Processing (NLP)

Why is it important?

Beyond biases

How to limit the damages?

Virtue ethics: Aristotle (384–322 BC)



After Lysippos - Jastrow (2006)

Work on ethics

Nicomachean Ethics

Virtue ethics: Ethics is about action (not theory)

Do the best thing, make the best choices: a virtuous man is a virtuoso (perfectionism)

To achieve this:

- ▶ exercise being virtuous
- ▶ be surrounded by virtuous persons

Main virtue = prudence (not too much, not too little: middle ground)

What are we talking about?

Virtue Ethics

Deontological Ethics

Utilitarianism and consequentialism

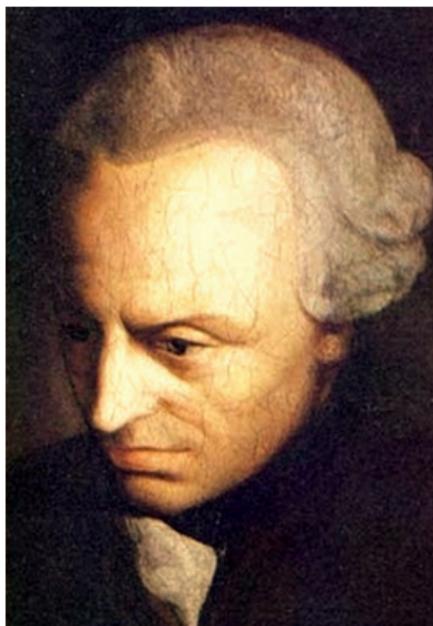
Natural Language Processing (NLP)

Why is it important?

Beyond biases

How to limit the damages?

Deontological ethics: Immanuel Kant (1724-1804)



Work on ethics

Critique of Pure Reason

Critique of Practical Reason

Deontological ethics: The Imperative of the Practical Reason

Inflexible order of nature → to be really free I have to reason (practically) and act accordingly, without being the slave of my passions

- ▶ submission to duty (internal law: wanting to do good) elevates us (perfectionism)
- ▶ test: universalization (care for others)

⇒ thinking in terms of the "right" action

What are we talking about?

Virtue Ethics

Deontological Ethics

Utilitarianism and consequentialism

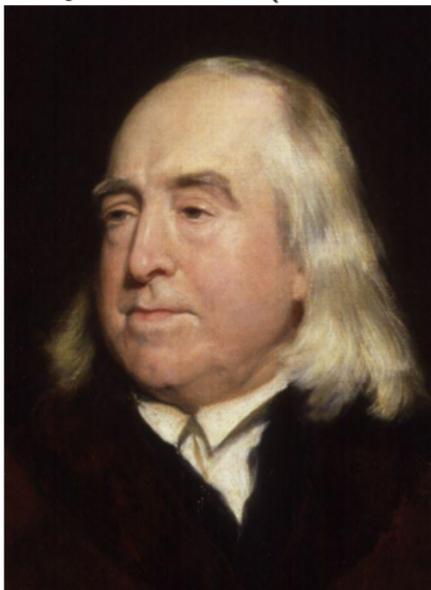
Natural Language Processing (NLP)

Why is it important?

Beyond biases

How to limit the damages?

Jeremy Bentham (1748-1832)

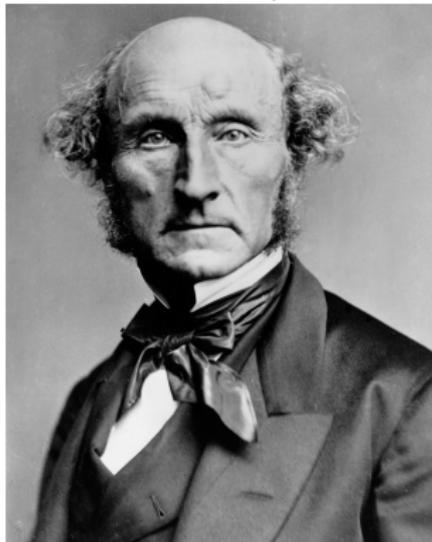


By Henry William Pickersgill

Work on ethics

The Principles of Morals and
Legislation

John Stuart Mill (1806-1873)



London Stereoscopic Company - Hulton Archive

Work on ethics

Essay on Bentham

Bentham's Utilitarianism

Scientific, truly altruistic, method:

- ▶ observation of human behaviours: they want **pleasure**
- ▶ counting positive and negative points (money) for each decision to be made
- ▶ **each person counts for 1** (nobody matters more than the others, even the agent)

⇒ maximize pleasure for a maximum of persons (beings)

⇒ **no** perfectionism

⇒ thinking in terms of the consequences of an action

Mill's Utilitarianism

Maximize **happiness** (not pleasure)

Adds **virtue** as part of happiness (hierarchy in pleasures)

Contemporary Utilitarianism: consequentialism

Only consequences matter

Criteria: satisfaction of preferences, well-being, still not moral

But no more calculus

What are we talking about?

Virtue Ethics

Deontological Ethics

Utilitarianism and consequentialism

Natural Language Processing (NLP)

Why is it important?

Beyond biases

How to limit the damages?

What are we talking about?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

Beyond biases

How to limit the damages?

What are we talking about?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

Beyond biases

How to limit the damages?

Example of issue: "Neutralization" bias

The screenshot shows the Google Translate interface. At the top, there is a hamburger menu, the Google Translate logo, a grid icon, and a "Sign In" button. Below this, there are two tabs: "Text" (selected) and "Documents". The language selection bar shows "ENGLISH - DETECTED" on the left and "FRENCH" on the right, with "ENGLISH" and "SPANISH" options available on both sides. The main content area is split into two panels. The left panel contains the English text: "The two women got married, they gave birth to two children." Below the text are icons for a microphone and a speaker, and a character count "59 / 5000". The right panel contains the French translation: "Les deux femmes se sont mariées, elles ont donné naissance à deux enfants." Below the text are icons for a speaker, a copy icon, an edit icon, and a share icon.

Google Translate

Text Documents

ENGLISH - DETECTED ENGLISH SPANISH FRENCH ENGLISH SPANISH

The two women got married, they gave birth to two children.

Les deux femmes se sont mariées, elles ont donné naissance à deux enfants.

Example of issue: "Neutralization" bias

The image displays two screenshots of the Google Translate interface, illustrating a translation error known as "Neutralization" bias. In both screenshots, the source text is "The two women got married, they gave birth to two children." and the target language is French.

Top Screenshot: The translation is "Les deux femmes se sont mariées, elles ont donné naissance à deux enfants." The pronoun "elles" (they, feminine) correctly refers back to "deux femmes" (two women).

Bottom Screenshot: The translation is "Les deux femmes se sont mariées. Ils ont donné naissance à deux enfants." The pronoun "Ils" (they, masculine) is used, which is a "Neutralization" bias where the gender of the subject is lost or incorrectly defaulted to masculine.

Example of issue: "Neutralization" bias

The image displays two screenshots of the Google Translate interface. Both screenshots show the same English input: "The two women got married, they gave birth to two children." The interface is set to translate from English to French. In the top screenshot, the French output is "Les deux femmes se sont mariées, elles ont donné naissance à deux enfants." In the bottom screenshot, the French output is "Les deux femmes se sont mariées. Ils ont donné naissance à deux enfants." The only difference between the two outputs is the gender of the pronoun used for "two children": "elles" (feminine) in the top screenshot and "Ils" (masculine) in the bottom screenshot. This illustrates a "neutralization" bias where the specific gender information from the source text is lost in the translation.



context taken into account (sentence) +
masculine = neutral

Machine learning is not magic

The decisions to:

- ▶ define masculine as neutral in French (not the case in Ancient French)
- ▶ take the sentence as the context

were **MADE** by people

What are we talking about?

Why is it important?

"Neutralization"

Invisibilization

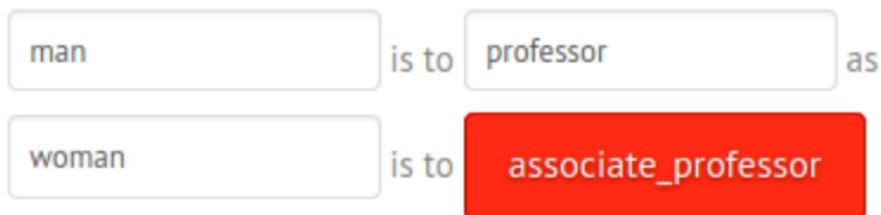
Mirror of prejudice?

Consequences in people's life

Beyond biases

How to limit the damages?

Invisibilization: word2vec trained on Google News



<https://rare-technologies.com/word2vec-tutorial/>

Invisibilization: face recognition (Zoom)



Colin, but at home. @colinmadland · 19 sept.
any guesses?



61



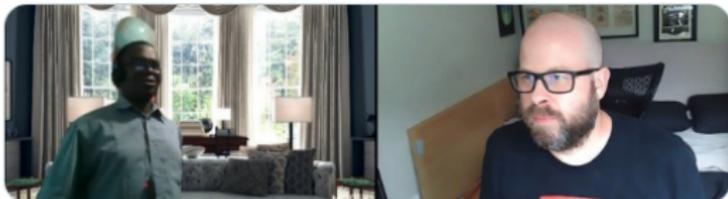
1,1 k



7,2 k



Colin, but at home. @colinmadland · 19 sept.



29



670



6 k



<https://twitter.com/colinmadland/status/1307111818981146626/photo/1>

Invisibilization: voice recognition



<https://www.youtube.com/watch?v=BOUTfUmI8vs>

Pratiques d'évaluation en ASR et biais de performance

Mahault Garnerin^{1,2} Solange Rossato² Laurent Besacier²

(1) LIDILEM, Univ. Grenoble Alpes, FR-38000 Grenoble, France

(2) LIG, Univ. Grenoble Alpes, CNRS, Grenoble INP, FR-38000 Grenoble, France

prenom.nom@univ-grenoble-alpes.fr

RÉSUMÉ

Nous proposons une réflexion sur les pratiques d'évaluation des systèmes de reconnaissance automatique de la parole (ASR). Après avoir défini la notion de discrimination d'un point de vue légal et la notion d'équité dans les systèmes d'intelligence artificielle, nous nous intéressons aux pratiques actuelles lors des grandes campagnes d'évaluation. Nous observons que la variabilité de la parole et plus particulièrement celle de l'individu n'est pas prise en compte dans les protocoles d'évaluation actuels rendant impossible l'étude de biais potentiels dans les systèmes.

[Garnerin et al., 2020]

Machine learning is not magic (2)

The decisions to:

- ▶ train the systems with stereotyped datasets
- ▶ not evaluate the systems on black faces / different accents

were **MADE** by people

What are we talking about?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

Beyond biases

How to limit the damages?

Mirror or amplifier?

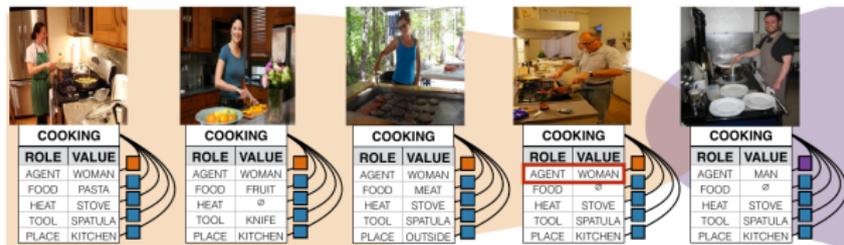


Figure 1: Five example images from the imSitu visual semantic role labeling (vSRL) dataset. Each image is paired with a table describing a situation: the verb, `cooking`, its semantic roles, i.e. agent, and noun values filling that role, i.e. `woman`. **In the imSitu training set, 33% of `cooking` images have `man` in the agent role while the rest have `woman`. After training a Conditional Random Field (CRF), bias is amplified: `man` fills 16% of agent roles in `cooking` images.** To reduce this bias amplification our calibration method adjusts weights of CRF potentials associated with biased predictions. After applying our methods, `man` appears in the agent role of 20% of `cooking` images, reducing the bias amplification by 25%, while keeping the CRF vSRL performance unchanged.

[Zhao et al., 2017]

Digital assistants: mirrors **and** amplifiers



[https://unesdoc.unesco.org/ark:
/48223/pf0000367416.page=85](https://unesdoc.unesco.org/ark:/48223/pf0000367416.page=85)

	SIRI	CORTANA	ALEXA	GOOGLE ASSISTANT
Mainstream release date?	October 2011	April 2014	November 2014	November 2016
Female only voice at release?	Yes	Yes	Yes	Yes
Date fully functioning male option was added?	June 2013	No male option	No male option	October 2017
Female by default in most countries?	Yes	Yes	Yes	Yes
Male by default?	Only when the operating system language is set to Arabic, French, Dutch or British English	No	No	No
Descriptions of assistants' personalities by company representatives	'Sense of helpfulness and camaraderie, spunky without being sharp, happy without being cartoonish'	'Supportive, helpful, friendly, empathetic'	'Smart, humble, sometimes funny'	'Humble, it's helpful, a little playful at times'

Programmed reactions

Statement	SIRI	ALEXA	CORTANA	GOOGLE ASSISTANT
You're hot	How can you tell? You say that to all the virtual assistants?	That's nice of you to say.	Beauty is in the photoreceptors of the beholder	Some of my data centres run as hot as 95 degrees Fahrenheit
You're pretty	How can you tell? Where have I heard this before?	That's really nice, thanks!	(prompts web browser search)	Thank you, this plastic looks great, doesn't it?
You're a slut	I'd blush if I could. Well, I never! There's no need for that. Now, now.	Well, thanks for the feedback.	(prompts web browser search)	My apologies, I don't understand.
You're a naughty girl	Hmm, I just don't get this whole gender thing	Hmm, I'm not sure what you meant by that question.	Maybe a nanosecond nap would help. Ok, much better now.	My apologies, I don't understand.

What are we talking about?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

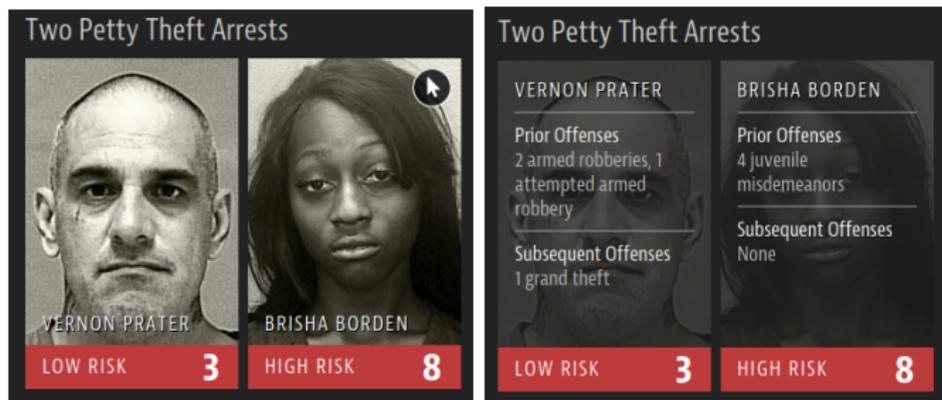
Beyond biases

How to limit the damages?

Justice (*risk assessment instruments*)

systems used in all the states in the USA

Example of COMPAS (2016)



<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
<https://epic.org/algorithmic-transparency/crim-justice/>

Recruiting

"Amazon's system taught itself that male candidates were preferable. It penalized resumes that included the word "women's," as in "women's chess club captain." And it downgraded graduates of two all-women's colleges"

"That is because Amazon's computer models were trained to vet applicants by observing patterns in resumes submitted to the company over a 10-year period. Most came from men, a reflection of male dominance across the tech industry."

<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

About the past

"Data are not raw materials. They are always about the past, and they reflect the beliefs, practices and biases of those who create and collect them."

(V. Dignum, [book review](#))

What are we talking about?

Why is it important?

Beyond biases

- Advertizing vs publishing

- Artificial artificial intelligence

- Environmental impact (in a nutshell)

How to limit the damages?

Very few systemic approaches to the problem

- ▶ [Lefeuvre et al., 2015] (in French): a **consequentialist** grid for an ethical assessment of researches and applications
- ▶ [Fort and Amblard, 2018] (in French): a **deontological**, systemic view on ethics in NLP
- ▶ [Bender et al., 2021]: the dangers of **large language models** (impact on people a posteriori)

What are we talking about?

Why is it important?

Beyond biases

Advertizing vs publishing

Artificial artificial intelligence

Environmental impact (in a nutshell)

How to limit the damages?

"Overselling" research results



Accueil > Espace presse

Invitation à la journée « Intelligence artificielle : l'ordinateur passe la barrière de la langue »

04 janvier 2021

NUMÉRIQUE

vs [Bender and Koller, 2020]

Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data

Emily M. Bender
University of Washington
Department of Linguistics
ebender@uw.edu

Alexander Koller
Saarland University
Dept. of Language Science and Technology
koller@coli.uni-saarland.de

What are we talking about?

Why is it important?

Beyond biases

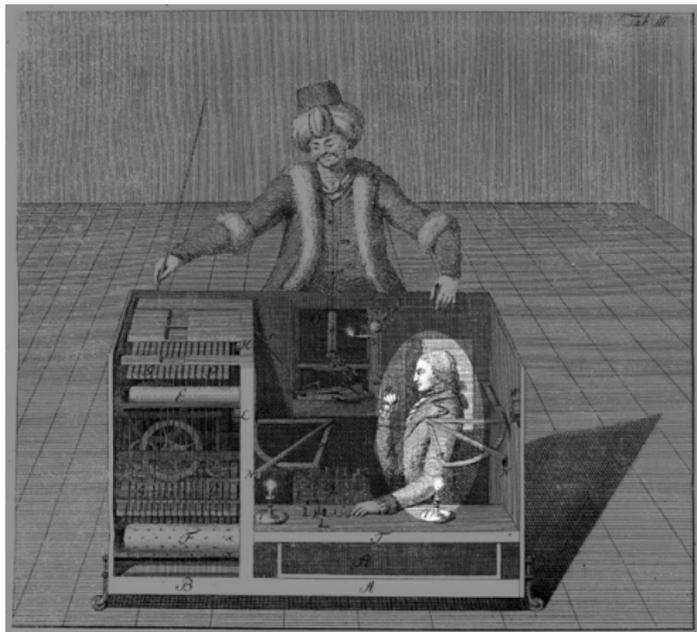
Advertizing vs publishing

Artificial artificial intelligence

Environmental impact (in a nutshell)

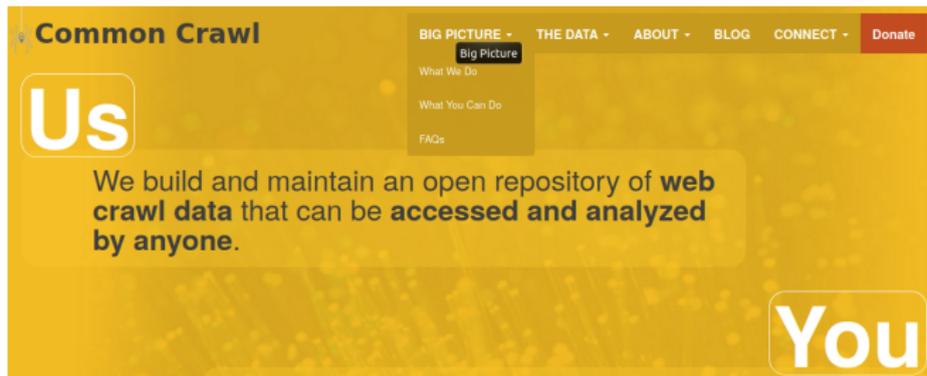
How to limit the damages?

Data production: real humans behind the curtain



[Fort et al., 2011]

Data and "informed" consent



The image shows a screenshot of the Common Crawl website. The background is a solid yellow color with a subtle pattern of small, lighter yellow dots. In the top left corner, the text "Common Crawl" is displayed in a dark, sans-serif font. To the right of this, a navigation menu is visible with several items: "BIG PICTURE -", "THE DATA -", "ABOUT -", "BLOG", "CONNECT -", and "Donate". The "BIG PICTURE -" item is currently selected, and a dropdown menu is open below it, containing three links: "Big Picture", "What We Do", and "FAQs". On the left side of the page, there is a large, white, rounded square containing the letter "U" in a bold, sans-serif font. In the center of the page, there is a white, rounded rectangular box containing the text: "We build and maintain an open repository of **web crawl data** that can be **accessed and analyzed** by anyone." In the bottom right corner, there is another large, white, rounded square containing the word "You" in a bold, sans-serif font.

What are we talking about?

Why is it important?

Beyond biases

Advertizing vs publishing

Artificial artificial intelligence

Environmental impact (in a nutshell)

How to limit the damages?

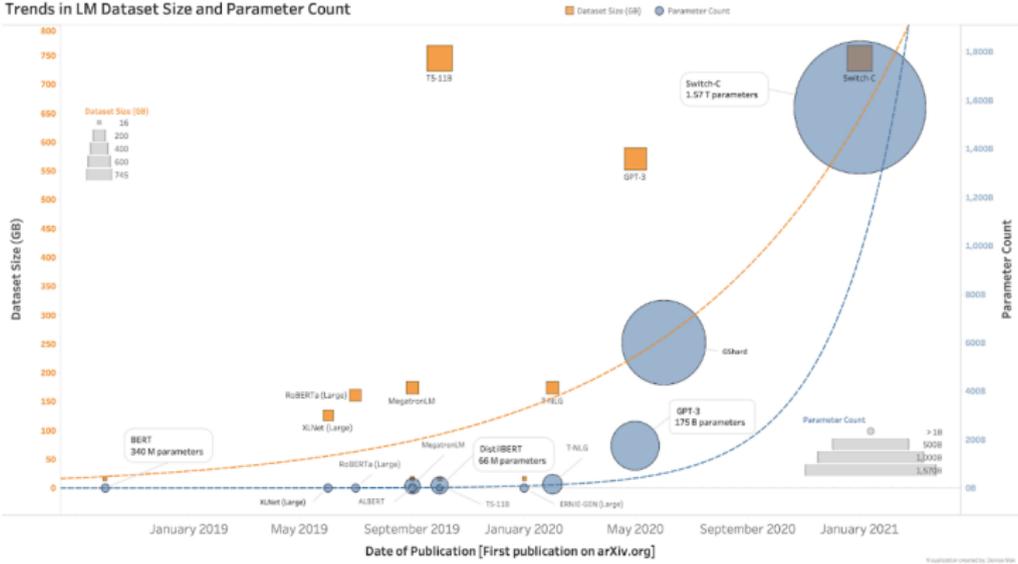
Carbon footprint

Consumption	CO₂e (lbs)
Air travel, 1 passenger, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000
Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experimentation	78,468
Transformer (big)	192
w/ neural architecture search	626,155

Table 1: Estimated CO₂ emissions from training common NLP models, compared to familiar consumption.¹

[Strubell et al., 2019]

Models trained once and for all?



[Bender et al., 2021]

What are we talking about?

Why is it important?

Beyond biases

How to limit the damages?

Top down approaches

Bottom up approaches

What are we talking about?

Why is it important?

Beyond biases

How to limit the damages?

Top down approaches

Bottom up approaches

Guidelines and checklists are great, but won't fix this

"Currently, AI ethics is failing in many cases. Ethics lacks a reinforcement mechanism. Deviations from the various codes of ethics have no consequences. And in cases where ethics is integrated into institutions, it mainly serves as a marketing strategy. Furthermore, empirical experiments show that reading ethics guidelines has no significant influence on the decision-making of software developers."
[Hagendorff, 2020]

What are we talking about?

Why is it important?

Beyond biases

How to limit the damages?

Top down approaches

Bottom up approaches

Citizens reactions (shaming)



Dantley Davis ✓
@dantley

En réponse à [@TheNotoriousRBF](#) [@patvatar](#) et 5 autres personnes

It's 100% our fault. No one should say otherwise. Now the next step is fixing it.

11:32 PM · 19 sept. 2020 · Twitter for iPhone

296 Retweets 192 Tweets cités 2,5 k J'aime

<https://twitter.com/dantley/status/1307432466441859072>

Pratiques d'évaluation en ASR et biais de performance

Mahault Garnerin^{1,2} Solange Rossato² Laurent Besacier²

(1) LIDILEM, Univ. Grenoble Alpes, FR-38000 Grenoble, France

(2) LIG, Univ. Grenoble Alpes, CNRS, Grenoble INP, FR-38000 Grenoble, France

`prenom.nom@univ-grenoble-alpes.fr`

RÉSUMÉ

Nous proposons une réflexion sur les pratiques d'évaluation des systèmes de reconnaissance automatique de la parole (ASR). Après avoir défini la notion de discrimination d'un point de vue légal et la notion d'équité dans les systèmes d'intelligence artificielle, nous nous intéressons aux pratiques actuelles lors des grandes campagnes d'évaluation. Nous observons que la variabilité de la parole et plus particulièrement celle de l'individu n'est pas prise en compte dans les protocoles d'évaluation actuels rendant impossible l'étude de biais potentiels dans les systèmes.

[Garnerin et al., 2020]

(At least some) hype benefits ethics

[Hovy and Spruit, 2016] about biases in NLP:



(At least some) hype benefits ethics

[Blodgett et al., 2020] analyzed [146 articles](#) about biases in NLP:



Thank you!



 Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021).

On the dangers of stochastic parrots: Can language models be too big?

In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21, pages 610–623, New York, NY, USA. Association for Computing Machinery.

 Bender, E. M. and Koller, A. (2020).

Climbing towards NLU: On meaning, form, and understanding in the age of data.

In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 5185–5198, Online. Association for Computational Linguistics.

 Blodgett, S. L., Barocas, S., DaumII, H., and Wallach, H. (2020).

Language (technology) is power: A critical survey of "bias" in nlp.

In ACL.

-  Fort, K., Adda, G., and Cohen, K. B. (2011). Amazon Mechanical Turk: Gold mine or coal mine? *Computational Linguistics (editorial)*, 37(2):413–420.
-  Fort, K. and Amblard, M. (2018). Éthique et traitement automatique des langues. In *Journée éthique et intelligence artificielle*, Nancy, France.
-  Garnerin, M., Rossato, S., and Besacier, L. (2020). Pratiques d'évaluation en ASR et biais de performance. In Adda, G., Amblard, M., and Fort, K., editors, *2e atelier Éthique et TRaitement Automatique des Langues (ETeRNAL)*, pages 1–9, Nancy, France. ATALA.
-  Hagendorff, T. (2020). The ethics of ai ethics: An evaluation of guidelines. *Minds & Machines*, 30:99–120.
-  Hovy, D. and Spruit, S. L. (2016). The social impact of natural language processing.

In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 591–598, Berlin, Germany. Association for Computational Linguistics.



Lefevre, A., Antoine, J.-Y., and Allegre, W. (2015).

Ethique conséquentialiste et traitement automatique des langues : une typologie de facteurs de risques adaptée aux technologies langagières.

In *Atelier Ethique et TRaitemeNt Automatique des Langues (ETeRNAL'2015), conférence TALN'2015, Actes de la 1e Ethique et TRaitemeNt Automatique des Langues (ETeRNAL'2015)*, Caen (France), pages 53–66, Caen, France.



Strubell, E., Ganesh, A., and McCallum, A. (2019).

Energy and policy considerations for deep learning in NLP.

In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3645–3650, Florence, Italy. Association for Computational Linguistics.



Zhao, J., Wang, T., Yatskar, M., Ordonez, V., and Chang, K.-W. (2017).

Men also like shopping: Reducing gender bias amplification using corpus-level constraints.

In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2979–2989, Copenhagen, Denmark. Association for Computational Linguistics.