



Ethics in AI: A view from Natural Language Processing (NLP)

Karën Fort

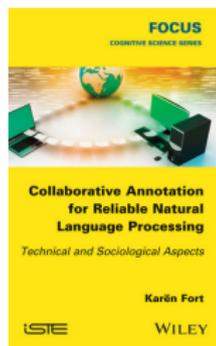
karen.fort@loria.fr / <https://members.loria.fr/KFort/>

SAILS Webinar, Sept. 27, 2021

Where I'm talking from

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

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



- ▶ Ethics and NLP



What is ethics?

Why is it important?

Beyond biases

Hopes?

Thanks

Ethics in general vs in the community

Merriam-Webster SINCE 1828

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

ethic

Dictionary Thesaurus

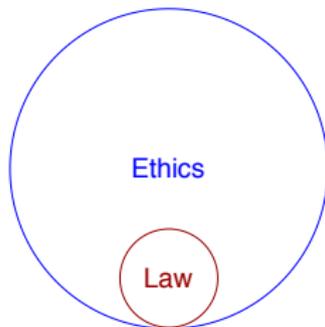
eth-ic | \ ɛ-ˈthik \

Definition of *ethic*

- 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
- 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

Traditional ethics in 1 slide (!)

- ▶ **Virtue ethics** (Aristotle): ethics is in action, the main virtue is **prudence** (not too much, not too little: middle ground)
- ▶ **Deontological ethics** (Kant): **moral principle** is a priori and absolute ("you shall not kill/steal")
- ▶ **Utilitarianism and consequentialism** (Bentham/Mill): thinking in terms of the **consequences** of an action (consider the nb of people impacted)

What is ethics?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

Beyond biases

Hopes?

Thanks

What is ethics?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

Beyond biases

Hopes?

Thanks

Example of issue: "Neutralization" bias

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. ☆

59 / 5000

Example of issue: "Neutralization" bias

The screenshot shows the Google Translate interface. The source text in English is "The two women got married, they gave birth to two children." The target text in French is "Les deux femmes se sont mariées, elles ont donné naissance à deux enfants." The interface includes a "Sign in" button, "Text" and "Documents" input options, and language selection menus for "ENGLISH - DETECTED", "ENGLISH", "SPANISH", and "FRENCH".

This screenshot is identical to the one above, but the French translation has been corrected to "Les deux femmes se sont mariées. Ils ont donné naissance à deux enfants." The change from "elles" to "Ils" illustrates the "Neutralization" bias, where the gender of the subjects is lost in the translation.

Example of issue: "Neutralization" bias

Google Translate

Text Documents

ENGLISH - DETECTED ENGLISH SPANISH ↔ FRENCH ENGLISH SPANISH

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

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

59 / 5000

Google Translate

Text Documents

ENGLISH - DETECTED ENGLISH SPANISH ↔ FRENCH ENGLISH SPANISH

The two women got married. They gave birth to two children. X

Les deux femmes se sont mariées. Ils ont donné naissance à deux enfants. ☆

59 / 5000



context taken into account (sentence) +
masculine = neutral

Machine learning is not magic

The decisions to:

- ▶ define masculine as neutral in French
- ▶ take the sentence as the context

were **MADE** by people

What is ethics?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

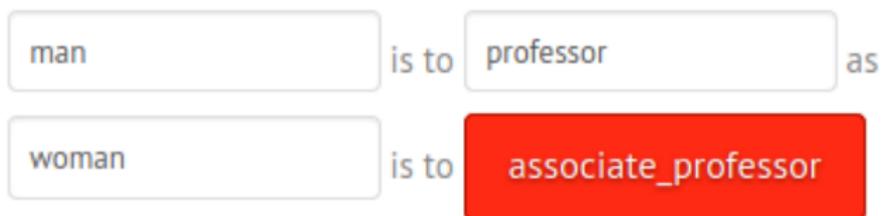
Consequences in people's life

Beyond biases

Hopes?

Thanks

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>

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 is ethics?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

Beyond biases

Hopes?

Thanks

Mirror or amplifier?

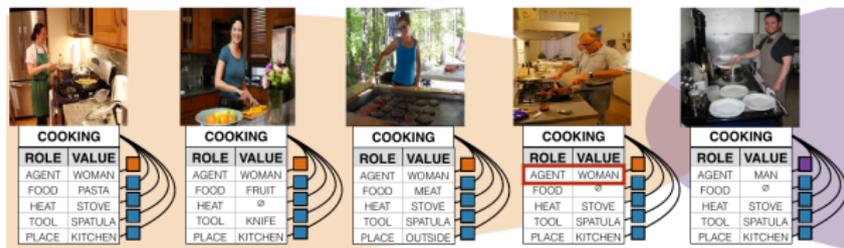


Figure 1: Five example images from the imSitu visual semantic role labeling (vSRL) dataset. Each image is paired with a table describing a situation: the verb, `cooking`, its semantic roles, i.e. `agent`, and noun values filling that role, i.e. `woman`. **In the imSitu training set, 33% of `cooking` 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>

	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 is ethics?

Why is it important?

"Neutralization"

Invisibilization

Mirror of prejudice?

Consequences in people's life

Beyond biases

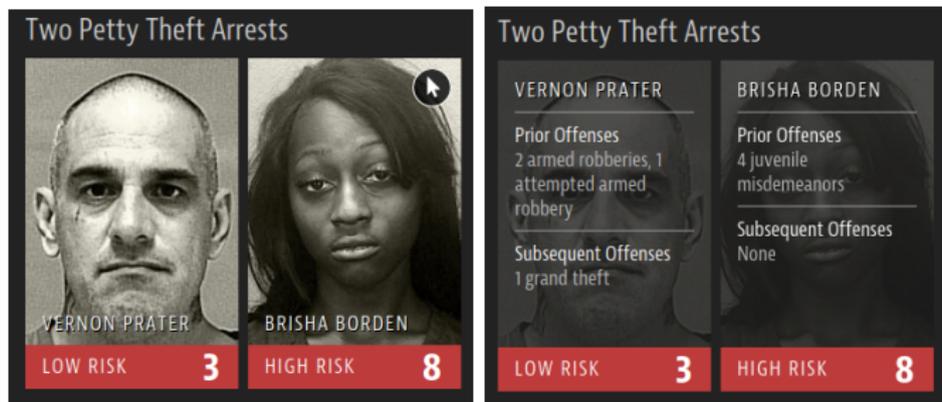
Hopes?

Thanks

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 is ethics?

Why is it important?

Beyond biases

Hopes?

Thanks

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)

"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

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

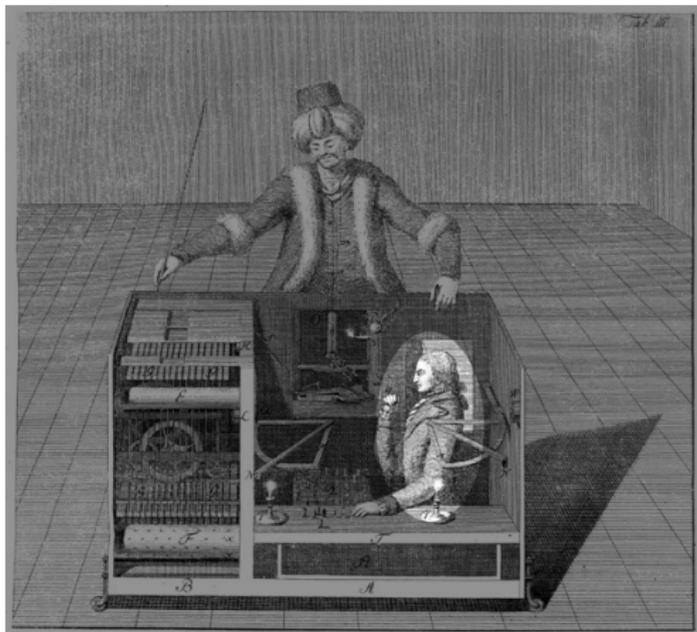
`prenom.nom@univ-grenoble-alpes.fr`

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.

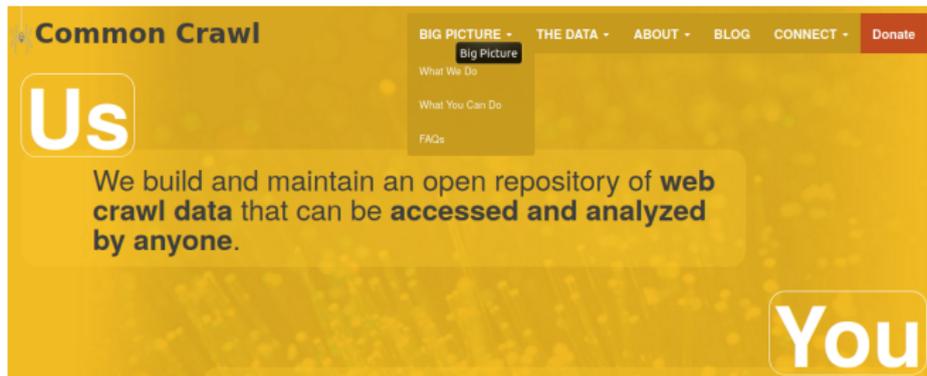
[Garnerin et al., 2020]

Data production: real humans behind the curtain



[Fort et al., 2011]

Data and "informed" consent



The image shows the homepage of the Common Crawl website. The background is a solid yellow color with a subtle pattern of small, lighter yellow dots. At the top left, the text "Common Crawl" is displayed in a bold, white, sans-serif font. To the right of this is a navigation menu with several items: "BIG PICTURE" (which is currently expanded to show a sub-menu), "THE DATA", "ABOUT", "BLOG", "CONNECT", and "Donate" (which is highlighted in a darker red color). The sub-menu for "BIG PICTURE" includes "What We Do", "What You Can Do", and "FAQs". On the left side of the main content area, the word "Us" is written in a large, white, rounded font. In the center, there is a white text box containing the message: "We build and maintain an open repository of **web crawl data** that can be **accessed and analyzed by anyone.**" On the right side, the word "You" is written in a large, white, rounded font, mirroring the "Us" on the left.

Common Crawl

BIG PICTURE - THE DATA - ABOUT - BLOG - CONNECT - Donate

Big Picture

What We Do

What You Can Do

FAQs

Us

We build and maintain an open repository of **web crawl data** that can be **accessed and analyzed by anyone.**

You

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]

What is ethics?

Why is it important?

Beyond biases

Hopes?

Thanks

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]

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.