



Ethics and Natural Language Processing (NLP): an (almost) tutorial

Karën Fort

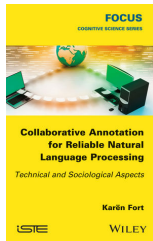
karen.fort@loria.fr / <https://members.loria.fr/KFort/>

ICMI, Oct. 18th, 2021

Where I'm talking from

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

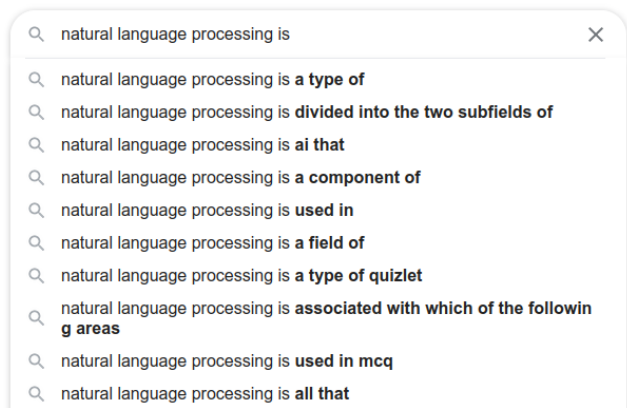
- ▶ Language resources creation for NLP, esp. using crowdsourcing



- ▶ Ethics and NLP



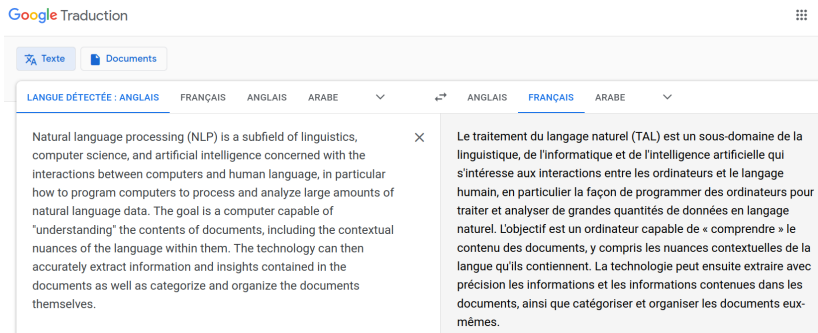
We all use NLP, every day



<https://www.google.com/>

What is NLP?

A recursive definition



The screenshot shows the Google Translate web interface. At the top, the text "Google Traduction" is visible on the left and a grid icon on the right. Below this, there are two tabs: "Texte" (selected) and "Documents". A language selection bar shows "LANGUE DÉTECTÉE : ANGLAIS" on the left and "ANGLAIS FRANÇAIS ARABE" on the right, with "FRANÇAIS" selected. The main content area is split into two columns. The left column contains the English text: "Natural language processing (NLP) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data. The goal is a computer capable of 'understanding' the contents of documents, including the contextual nuances of the language within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themselves." The right column contains the French translation: "Le traitement du langage naturel (TAL) est un sous-domaine de la linguistique, de l'informatique et de l'intelligence artificielle qui s'intéresse aux interactions entre les ordinateurs et le langage humain, en particulier la façon de programmer des ordinateurs pour traiter et analyser de grandes quantités de données en langage naturel. L'objectif est un ordinateur capable de « comprendre » le contenu des documents, y compris les nuances contextuelles de la langue qu'ils contiennent. La technologie peut ensuite extraire avec précision les informations et les informations contenues dans les documents, ainsi que catégoriser et organiser les documents eux-mêmes." A small 'X' icon is visible between the two columns.

<https://translate.google.com/>

(More advanced) applications in our daily lives

... not necessarily very efficient



A screenshot of a Twitter chat conversation with SNCF (@SNCF). The chat shows a user asking a question about a train network, and the chatbot responding with a confirmation and a follow-up question.

← **SNCF** ✓
@SNCF

13 sept. 2021 à 6:45 PM

Je n'ai pas de réseau dans mon TGV
inoui 2593

13 sept. 2021 à 6:46 PM ✓

C'est noté pour votre
numéro de train

 Ce train part-il aujourd'hui ?

13 sept. 2021 à 6:46 PM

oui

13 sept. 2021 à 6:46 PM ✓

The SNCF (French railway) chatbot on <https://twitter.com/home>

Two revolutions in less than a decade

- ▶ Much more implication from big firms in the field (GAFAM)
- ▶ Deep learning (from approx. 2013), including (very) large language models

→ **ethical issues** start to show

Who is talking, from where?

Practicing ethics in NLP

Helping the thought

Exercise: let's prompt!

Thanks to Chris Callison-Burch!



Help



CCB lab members



Playground ?



Load a preset...



This is GPT-3 and we're going to play with it so you can see what it's like. We're going to start over here, and we're going to try to find the other team. See if we can find them." The first thing was to get teams together. That was easy. It was often the same people, but there were some changes. We would say

Submit



79

<https://beta.openai.com/playground>

Exercise: which ethical issues did you find?



Who is talking, from where?

Practicing ethics in NLP

Stereotypes

Environmental impact (in a nutshell)

Data production: real humans behind the curtain

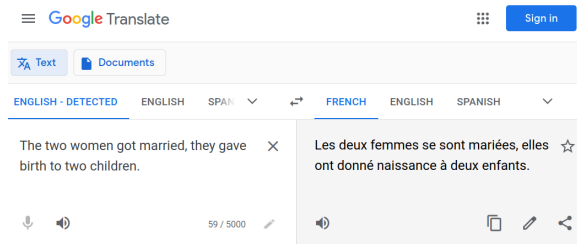
More systemic issues

Helping the thought

Exercise: where do the stereotypes come from?



Mirror of existing stereotypes?



The screenshot shows the Google Translate interface. At the top, there is a hamburger menu, the Google Translate logo, a grid icon, and a blue "Sign In" button. Below this is a navigation bar with "Text" and "Documents" tabs. The main area shows the source language as "ENGLISH - DETECTED" and the target language as "FRENCH". The source text is "The two women got married, they gave birth to two children." and the translated text is "Les deux femmes se sont mariées, elles ont donné naissance à deux enfants." The interface includes icons for voice input/output, a character count (59 / 5000), and options to copy, edit, or share the translation.

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

Mirror of existing stereotypes?

The image displays two screenshots of the Google Translate interface, illustrating how the same English sentence can be translated into French in two different ways that reflect gender stereotypes.

Top Screenshot: The source text is "The two women got married, they gave birth to two children." The detected language is English. The target language is French. The translation provided is "Les deux femmes se sont mariées, elles ont donné naissance à deux enfants." The pronouns used are "elles" (they, feminine).

Bottom Screenshot: The source text is "The two women got married. They gave birth to two children." The detected language is English. The target language is French. The translation provided is "Les deux femmes se sont mariées. Ils ont donné naissance à deux enfants." The pronouns used are "Ils" (they, masculine).

The interface includes a "Sign In" button, "Text" and "Documents" input options, and a language selection menu showing "ENGLISH - DETECTED", "ENGLISH", "SPAN", and "FRENCH".

Mirror of existing stereotypes?

The image displays two screenshots of the Google Translate interface. Both screenshots show the same input text: "The two women got married, they gave birth to two children." The top screenshot shows the output translation in French: "Les deux femmes se sont mariées, elles ont donné naissance à deux enfants." The bottom screenshot shows the output translation in French: "Les deux femmes se sont mariées. Ils ont donné naissance à deux enfants." A red arrow points to the change from "elles" to "Ils".



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

ML and biases: mirror or amplifier?

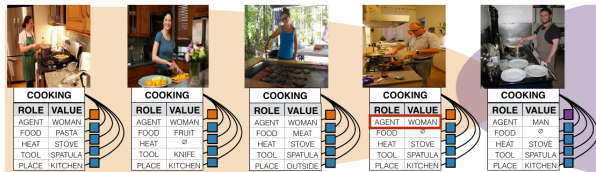


Figure 1: Five example images from the imSitu visual semantic role labeling (vSRL) dataset. Each image is paired with a table describing a situation: the verb, `cooking`, its semantic roles, i.e. `agent`, and noun values filling that role, i.e. `woman`. **In the imSitu training set, 33% of `cooking`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

(At least some) hype benefits ethics

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



(At least some) hype benefits ethics

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



Who is talking, from where?

Practicing ethics in NLP

Stereotypes

Environmental impact (in a nutshell)

Data production: real humans behind the curtain

More systemic issues

Helping the thought

Exercise: any idea about the environmental impact of NLP?



Carbon footprint

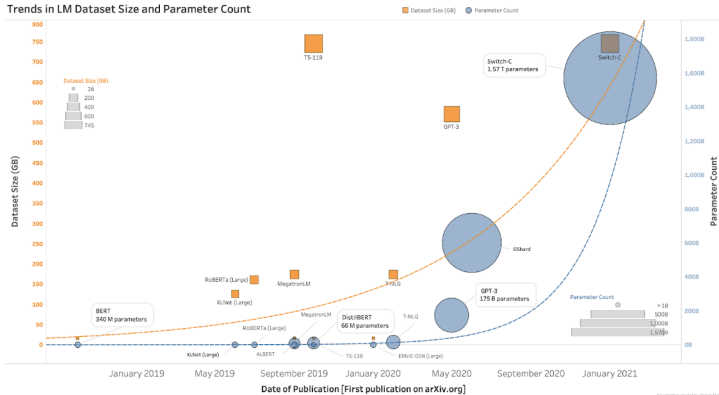
Consumption	CO₂e (lbs)
Air travel, 1 passenger, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000
Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experimentation	78,468
Transformer (big)	192
w/ neural architecture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

Trends in LM Dataset Size and Parameter Count



[Bender et al., 2021]

Who is talking, from where?

Practicing ethics in NLP

- Stereotypes

- Environmental impact (in a nutshell)

- Data production: real humans behind the curtain

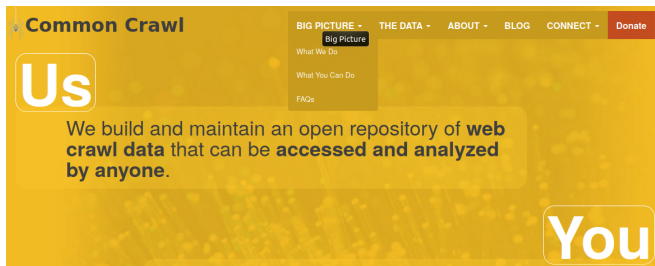
- More systemic issues

Helping the thought

Exercise: where does data come from?

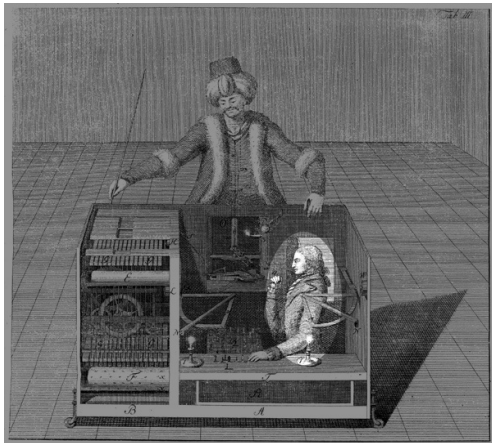


Data and "informed" consent



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 horizontal navigation bar contains several menu items: "BIG PICTURE" (with a dropdown arrow), "THE DATA" (with a dropdown arrow), "ABOUT" (with a dropdown arrow), "BLOG", "CONNECT" (with a dropdown arrow), and "Donate" (in a dark red button). The "BIG PICTURE" dropdown menu is open, showing three options: "Big Picture" (highlighted with a dark background), "What We Do", and "What You Can Do". Below the navigation bar, on the left side, the word "Us" is written in a large, white, rounded font inside a white rounded square. In the center, a white rounded rectangle contains the text: "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 inside a white rounded square.

Artificial artificial intelligence



[Fort et al., 2011]

Data not "for the taking"

Decolonising Speech and Language Technology

Steven Bird

Northern Institute
Charles Darwin University

Abstract

After generations of exploitation, Indigenous people often respond negatively to the idea that their languages are data ready for the taking. By treating Indigenous knowledge as a commodity, speech and language technologists risk disenfranchising local knowledge authorities, reenacting the causes of language endangerment. Scholars in related fields have responded to calls

[Bird, 2020]

Who is talking, from where?

Practicing ethics in NLP

- Stereotypes

- Environmental impact (in a nutshell)

- Data production: real humans behind the curtain

- More systemic issues**

Helping the thought

"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

By the way...



Emily M. Bender

@emilymbender



Dear Computer Scientists,

"Natural Language" is **not** a synonym for "English".

That is all.

-Emily

6:32 PM · Nov 26, 2018 · TweetDeck

<https://twitter.com/emilymbender/status/1067108757488848896>

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`






















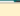
RÉSUMÉ

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]

Long-term vs short-term

grâce à Sibylle, je peux communiquer avec mon père et ma

mots	pictos	l	j	t	d	p	a	<		
e	c	s	m	r	i	f	u	q	mère	  
n	à	o	é	v	b	g	h	y	filie	   
ê	k	w	z	ù	x	ç	_	è	famille	   
									grand-mère	   
									soeur	   
									vie	  
,	.	-	?	!	.	@	;		tante	

- Pathologie lourde avec perte de parole
- Clavier virtuel avec prédiction lexicale



**Vitesse
de
saisie**



**Maîtrise
de la
langue**



[Antoine and Lefevre, 2014]

Who is talking, from where?

Practicing ethics in NLP

Helping the thought

Inspiring sources

- ▶ High-Level Expert Group on AI [Ethics guidelines for trustworthy AI](#)
- ▶ CERNA on [Research ethics in ML](#)
- ▶ Unesco about [Gender divides in digital skills](#)
- ▶ CNIL about [voice assistants](#)

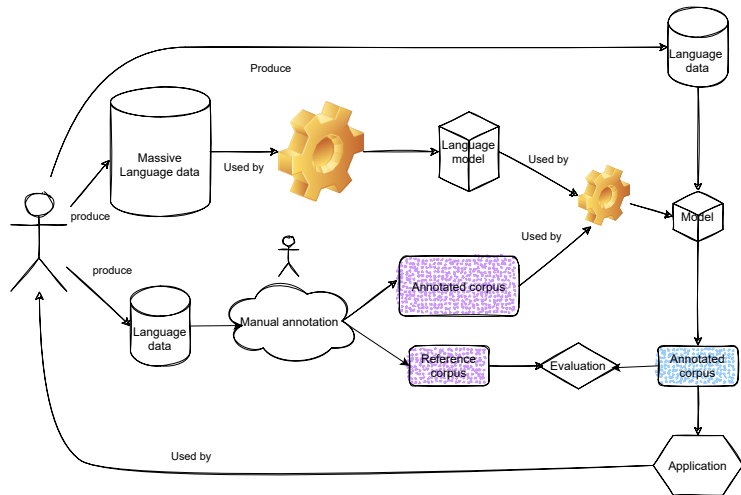
(Some) systemic approaches to the problems faced in NLP

- ▶ [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)

Thank you!



An overview of (the vast majority of) present NLP since 2018 (and evolving rapidly) – a draft





Antoine, J.-Y. and Lefevre, A. (2014).

Pour une réflexion éthique sur les conséquences de l'usage des NTIC : le cas des aides techniques (à composante langagière ou non) aux personnes handicapées.

In Actes de la journée ATALA Éthique et TAL.



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.



Bird, S. (2020).

Decolonising speech and language technology.

In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3504–3519, Barcelona, Spain (Online). International Committee on Computational Linguistics.



Blodgett, S. L., Barocas, S., Daumll, 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.