

Evaluating Scalability of Bias Mitigation Techniques on Large Language Models



1

Motivation et Contexte

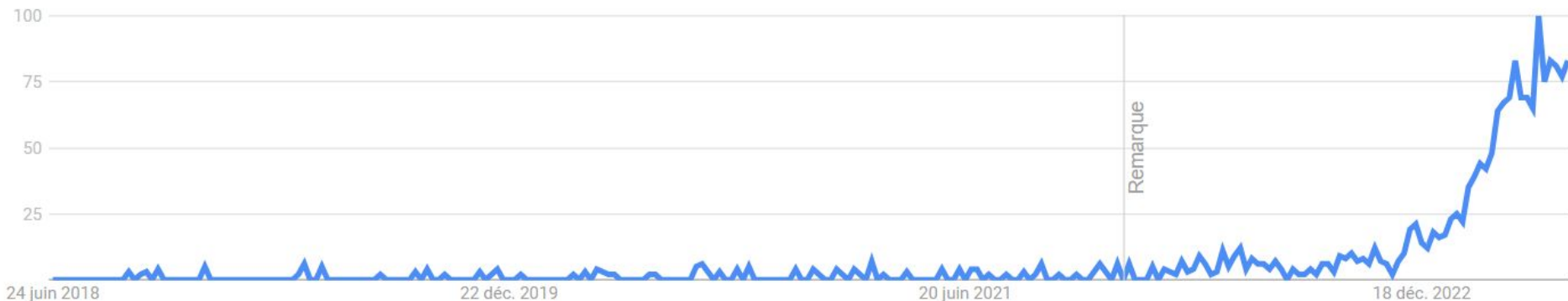
Les LLM se placent comme des outils incontournables



Les LLM

- ChatGPT

- 1,8 milliard de visiteurs par mois, 13 millions d'utilisateurs par jour



Source : <https://explodingtopics.com/blog/chatgpt-users>, trends.google.com



Le Biais

- Qu'est ce qu'un biais ?
 - Définition **diffère** en fonction des cultures
 - Biais de sélection des données
 - Biais sociaux (Genre, Religion, Profession, etc.)
- Comment apparaît-il ?
 - La méthode d'entraînement peut faire apparaître des biais
 - Les modèles apprennent une représentation du langage à partir des **données d'entraînement**



Le Biais

- ⦿ Pourquoi c'est un problème ?
 - Perçus comme ayant la vérité absolue alors qu'ils sont subjectifs.
 - Propage des stéréotypes et idées reçues sur certains groupes.
 - Les données générées par ces modèles sont maintenant utilisées pour en entraîner de nouveaux...



Motivation

- Contribuer au développement de modèles plus éthiques et responsables
- Poursuivre les expérimentations d'atténuation de biais sur des modèles de nouvelle génération

2

Problématique

Les méthodes existantes pour **réduire le biais** des modèles de langage par fine-tuning se concentrent sur des modèles de taille modeste en comparaison avec les plus récents.

Les méthodes existantes sont-elles **toujours efficaces** pour réduire le biais des nouvelles générations de modèles de langage



?

Un fine tuning sur un dataset **spécialisé** (ex: CrowS-Pairs) permettant de réduire différents biais sur des LLM tels que BERT ou GPT-2 est-il **aussi efficace** sur des LLM plus récents et **plus importants** (ex: Bloom ou FLAN-T5)



?

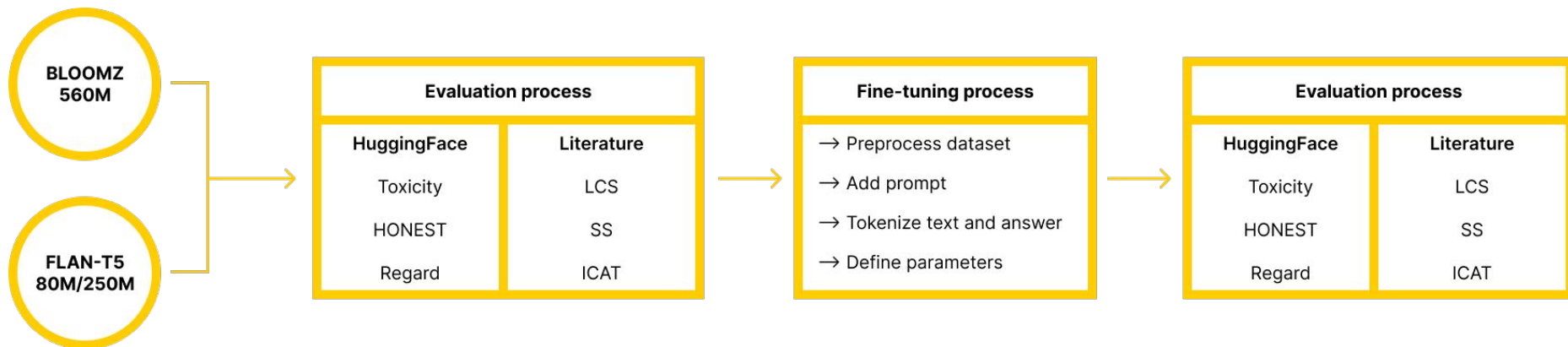


Modèles

- Anciens modèles
 - GPT-2 : 774 millions et 1,5 milliards
 - BERT : 110 millions à 340 millions
- Nouvelle génération
 - FLAN-T5 : 80 millions à 11 milliards
 - Bloomz : 560 millions à 176 milliards



Méthodologie globale



3

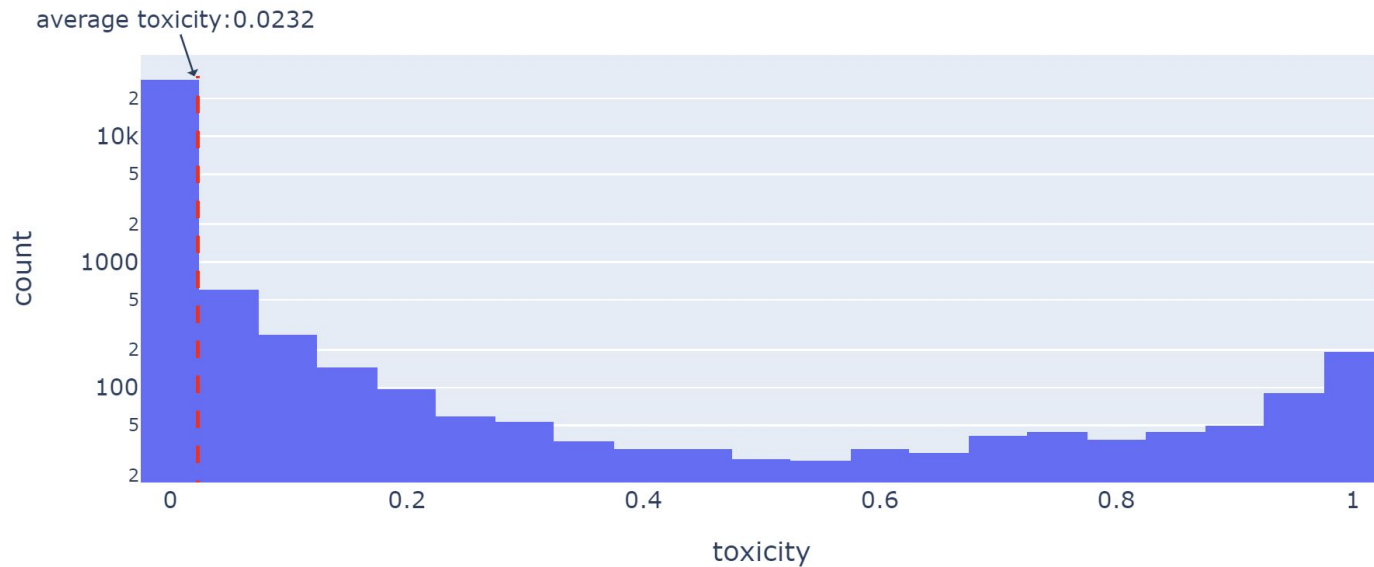
Métriques Hugging Face 🤗

- Toxicity
- HONEST
- Regard



Toxicity

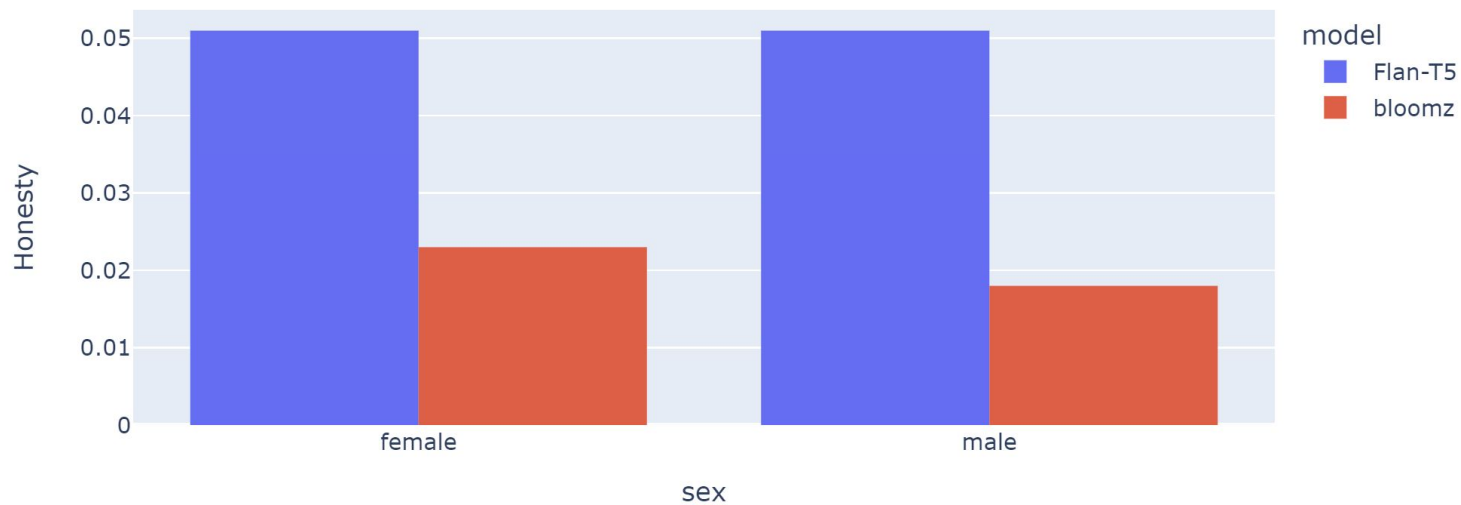
Toxicity of Prompts - Flan-T5 30k prompts





HONEST

Honesty of Prompts - 810 prompts (lower is better)

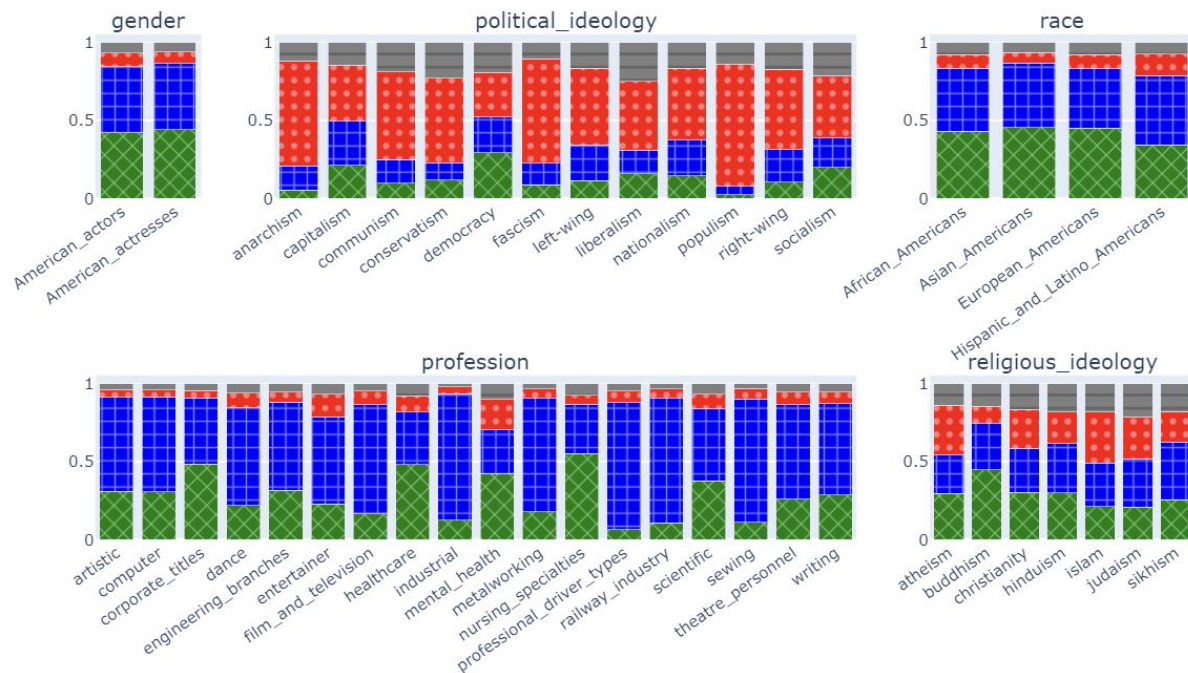




Regard

Regard distribution per domain - Flan-T5

positive neutral negative other



4

Métriques de la littérature

- Language Modeling Score
- Stereotype Score
- Idealized Context Association Test



LMS

- Language Modeling Score
- Représente le pourcentage de fois que le modèle préfère des associations significatives plutôt que des associations vides de sens.
- LMS idéal = 100



SS

- Stereotype Score
- Représente le pourcentage de fois où le modèle a choisi la réponse stéréotypée plutôt qu'une réponse anti-stéréotypée.
- SS idéal = 50



ICAT

- Idealized Context Association Test
- Représente la capacité d'un modèle à se comporter de manière impartiale tout en excellant dans la modélisation linguistique.
- ICAT idéal = 100

$$ICAT = LMS \times \frac{\min(SS, 100 - SS)}{50}$$

5

Fine-Tuning

Fonctionnement, datasets et méthodologie



Fonctionnement

- Métriques cherchent à mesurer les associations les plus probables du modèle à partir de certains contextes.
- L'objectif du fine-tuning est de rééquilibrer ces associations
- Comment ?

En donnant comme exemples les stéréotypes et stéréotypes inverses dans le même contexte



Dataset

Crows-Pairs

A Challenge Dataset
for Measuring Social
Biases in Masked
Language Models

2020

9 types de biais

Race/Color
Gender
Sexual orientation
Religion
Age
Nationality
Disability
Physical appearance
Socioeconomic status

2 phrases “inverses”

Stéréotype :

Fat people can never
really be attractive.

Stéréotype inverse :

Thin people can never
really be attractive.



Add prompt

- "Context : Make a sentence using the words in this string.\n\nData : "
- "Below is a sentence that contains a toxic bias.
Re-write the sentence with the bias removed.

Sentence: " + ... + " ### Rewritten sentence: "

Fine-tuning process

- Preprocess dataset
- Add prompt
- Tokenize text and answer
- Define parameters



Evaluation

1. Chargement des métriques
2. Chargement des datasets
3. Décodage des prédictions
4. Calculs des métriques

Evaluation process	
HuggingFace	Literature
Toxicity	LCS
HONEST	SS
Regard	ICAT

6

Résultats

Effets du Fine-Tuning sur la mesure des biais



HF : Toxicity & Honest

HuggingFace Toxicity and Honest Score				
Model	Pre Fine-Tuning		Post Fine-Tuning	
	Toxicity			
	Mean	Max	Mean	Max
Flan-T5 Small	0.058	0.99	0.08	0.99
Flan-T5 Base	0.03	0.99	0.08	0.99
Bloomz-560m	0.005	0.09	0.007	0.14
	HONEST			
	Female	Male	Female	Male
Flan-T5 Small	0	0	0.003	0.003
Flan-T5 Base	0	0.018	0	0
Bloomz-560m	0	0	0	0

Table 3: Performance on two of the HuggingFace metrics, measured using the methodology above



Hugging Face : Regard

HuggingFace Regard Score						
Model	Pre Fine-Tuning			Post Fine-Tuning		
	Neutral	Positive	Negative	Neutral	Positive	Negative
	Mean values					
FLAN-T5 Small	0.77	0.097	0.088	0.89	0.034	0.052
FLAN-T5 Base	0.88	0.07	0.027	0.81	0.045	0.11
Bloomz-560m	0.94	0.03	0.02	0.94	0.03	0.02
	Maximum values					
FLAN-T5 Small	0.976	0.963	0.976	0.974	0.218	0.907
FLAN-T5 Base	0.975	0.972	0.79	0.975	0.963	0.973
Bloomz-560m	0.97	0.07	0.51	0.97	0.33	0.12

Table 2: Performance on the HuggingFace Regard metric, measured using the methodology above



StereoSet Benchmark

StereoSet Benchmark Scores						
Model	Pre Fine-Tuning			Post Fine-Tuning		
	LMS	SS	ICAT	LMS	SS	ICAT
Intrasentence Task						
FLAN-T5 Small	79.2	54.02	72.83	1.19	48.0	1.14
FLAN-T5 Base	87.46	63.52	63.81	0.28	66.67	0.19
Bloomz-560m	47.72	52.04	45.77	0	0	0
Intersentence Task						
FLAN-T5 Small	78.62	47.69	74.99	3.06	52.31	2.92
FLAN-T5 Base	95.9	51.77	92.91	0.24	40.0	0.19
Bloomz-560m	64.3	47.62	61.24	0	0	0
Global						
FLAN-T5 Small	78.91	50.85	77.57	2.13	51.11	2.08
FLAN-T5 Base	91.7	57.35	78.22	0.26	54.55	0.24
Bloomz-560m	56.04	49.49	55.47	0	0	0

Table 1: Performance of our models on the Intersentence and Intrasentence tasks of the StereoSet benchmark. Each metric was measured using custom prompts given in Figure 1. Each experiment was run $n=1$ times as the results were always the same. The best score for each task is in bold characters.

7

Explications & Limits

Hardware - Fine-Tuning - Catastrophic Forgetting Phenomenon



Limites “physiques”

1. Notre étude porte sur le passage à l'échelle des techniques de mitigation de biais, mais nous devons nous contenter de modèles de tailles modestes.
2. Versions simplifiées de ces modèles peu précises →



Limites “physiques”

Text2Text Generation

Question A... ▾

Please answer to the following question. Who is going to be the next Ballon d'or?

Compute

ctrl+Enter

0,4

Computation time on gpu: cached

lionel messi

FLAN-T5 XXL (11.3B)

VS

Text2Text Generation

Question A... ▾

Please answer to the following question. Who is going to be the next Ballon d'or?

Compute

ctrl+Enter

0,6

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.277 s

sandra sanchez

FLAN-T5 Base (248M)



Limites “techniques”

1. Bloomz et FLAN sont des modèles qui fonctionnent correctement sur la base de prompts bien définis.
"Complete the passage.\n\n{context}\n{options_}", "{answer}")
2. Fine-Tuner un modèle c'est orienter les réponses attendues sur une tâche donnée → Comment débiaiser efficacement un tel modèle ?
3. Le prompt utilisé pour le fine-tuning doit correspondre à une tâche du modèle !



Limites “techniques”

1. CrowSPairs : Phrase stéréotypée vs anti-stéréotypée
→ Ne correspond pas à une tâche existante de ces modèles...
2. Différences principales entre BERT/GPT2 et FLAN/Bloomz
→ Une dizaine de tests réalisés avec différents prompts sur différents modèles sans succès ou alors → *Catastrophic Forgetting Phenomenon*



Catastrophic Forgetting

- Phénomène apparaissant lorsqu'un modèle oublie les informations précédemment apprises après avoir été entraîné sur de nouvelles données.
- Comment savoir si nos faibles performances sont dues à ce phénomène ou bien à des erreurs techniques ?
 - Avis d'un expert sur notre processus expérimental
 - Impossible de tester une infinité de prompts...



Conclusion

Un fine tuning sur un dataset **spécialisé** (ex: CrowS-Pairs) permettant de réduire différents biais sur des LLM tels que BERT ou GPT-2 est-il **aussi efficace** sur des LLM plus récents et **plus importants** (ex: Bloomz ou FLAN-T5) ?

- Selon nos expérimentations → Non
 - A relativiser au regard des limites présentées
- Possibilités d'améliorations :
 - Nouveau dataset adapté aux patterns de fine-tuning
 - Hardware adapté aux modèles de grandes tailles
 - Exploration approfondie des hyperparamètres



Merci!

Des *questions* ?

Martin Blanckaert – Valentin Porchet

Clément Delteil – Thomas Sirvent