# Vlaamse (L)LMs Vlaamse taalmodellen, groot en klein

Dr. ing. Pieter Delobelle - 22/04/2025







# Language modeling



# He is a Causal LM

He

is



He is a doctor Masked LM *He* <*m*>*a doctor* 





## 2. Masked language modeling

Vlaamse (L)LMs – 2

# Language modeling







**RobBERT** 



2. Masked language modeling He is a doctor Masked LM *He* <*m*>*a doctor* 



https://pieter.ai/robbert/

Vlaamse (L)LMs - 3



# Large training corpuses are used with a large focus on English





## **One book** 40-50k tokens











# **One bookshelf**

1.6M - 2.5M tokens

# **One LLM training set**

2.5T - 6T tokens ~2 500 000 bookshelves

Vlaamse (L)LMs – 4

No, I am not a giraffe.







No, I am not a giraffe. No, I am not a giraffe.







No, I am not a giraffe. No, I am not a giraffe. [2822, 11, 358, 1097, 539, 264, 41389, 38880, 13]





21223: affe











# Few non-English words are tokens

Token types for words in English do not match, so the tokenizer falls back to nonrepresentative tokens types.





# Few non-English words are tokens

Token types for words in English do not match, so the tokenizer falls back to nonrepresentative tokens types.





### is een absurde gir af ! Dat gedachte . ge en gir af ! Dat abs een

ach



# ... and morpheme boundaries are not respected

- Tokenization happens eagerly
- Representations are dependent on tokens





# Trans-loscenization

## 1. Token alignment



## 3. Model adaptation: continue pretraining for a few GPU hours (e.g. 40h)



Vlaamse (L)LMs – 12



# Trans-loscenization

## 1. Token alignment

Parallel corpus	Per-language	e
	tokenization	
: : 		:

## 2. Embedding mapping



**3. Model adaptation:** continue pretraining for a few GPU hours (e.g. 40h)





# Trans-loscenization

## 1. Token alignment

Parallel corpus	Per-language	e
	tokenization	
: :		:

## 2. Embedding mapping



**3. Model adaptation:** continue pretraining for a few GPU hours (e.g. 40h)







**Tweety LLMs** A series of models with language-specific tokenizers





# tweety-7b-dutch



# tweety-7b-tatar



## Community model tweety-7b-italian ithub.com/RitA-nlp



gpt-neo mala-50 tweety

Model

mistraltowerbas gpt-neotweety-7



	Training tokens	Normalized PPL
1-7b-v0.1	6-8T	9.4
SEL (Minixhofer et al., 2022)	+0.4B	34.3
proved Dutch dictionary	+0.4B	27.1
5 (Dobler & de Melo, 2023)	+0.4B	31.9
-7b-dutch-v24a (ours)	+0.4B	11.1
o-1.3b-dutch	33B	21.2
00-10b-v2	+30-60B	18.9
-7b-dutch-v24a (ours)	+8.5B	<b>7.7</b>

	Tokenizer		SQu	AD-NL A	CC
	Туре	$ \mathcal{V} $	0-shot	1-shot	2-shc
-7b-v0.1	English BPE	32 000	14.3	21.3	24.
se-7b-v0.1	English BPE	32 000	13.0	20.9	22.
-1.3b-dutch	Dutch BPE	50 257	0.0	0.0	0.
7b-dutch-v24a (ours)	Dutch BPE	50 257	9.0	25.8	27.





# tweety-7b-dutch



# tweety-7b-tatar



## **Community model** tweety-7b-italian github.com/RiTA-nlp



Model Mistral Mistral+F **MistralRA MistralAV** Tweety-7b Mistral+G

Mod Tow Tow Hyd Hyd

- Hyd Goo
- Hyd



# **Tatar:** NLU← and summarization→

	Accuracy	Model	Chr
Т	23.25 25.42	Mistral Mistral+FT	13.30 23 1
ND G G G G G G G G	0.00 17.00	MistralRAND <b>Tweety-7b-tatar-v24a</b> (ours)	3.79 <b>30.0</b>
Trans	$\sim 44.10$	Mistral+GTrans	30.4

# Hydra LLMs: Switching heads for zero-shot machine translation

del	Shor	t Text	Long	g Text	Socia	l Media
verInstruct	17.5	$\pm 0.4$	13.5	$\pm 0.3$	17.2	$\pm 0.5$
verInstruct+ParFT	24.5	$\pm 0.4$	16.5	$\pm 0.3$	20.6	$\pm 0.6$
draTower+ParFT	39.6	$\pm 0.5$	18.4	$\pm 0.5$	33.1	$\pm 1.4$
draTower	47.3	$\pm 0.4$	32.8	$\pm 0.4$	39.2	$\pm 1.5$
draTower+BackFT	53.7	$\pm 0.2$	33.6	$\pm 0.3$	46.1	$\pm 1.4$
ogle Translate	55.5	±0.2	35.3	±0.2	<del>63.8</del>	$\pm 1.8$
draTower+BackFT+NFR			39.2	±0.6		



5

3

3

# **European Tweeties** Trans-tokenizing all EU languages



## tweety-7b-dutch











# **Geitje-7b** First Dutch LLM





Vlaamse (L)LMs – 19



# Geitje-7b First Dutch LLM that got taken down by Brein



- Trained on 'gigacorpus'
- A torrent with gigabytes of Dutch books
- Gigacorpus got taken down by Brein already  ${}^{\bullet}$



https://tweakers.net/nieuws/231254/ontwikkelaar-haalt-taalmodel-geitje-offline-na-verzoek-stichting-brein.html

### Ontwikkelaar haalt taalmodel GEITje offline na verzoek Stichting Brein - update

Het Nederlandse Al-taalmodel GEITje is offline gehaald op 'dringend verzoek' van Stichting Brein. GEITje zou volgens Brein deels getraind zijn op documenten uit de dienst Library Genesis, die afgelopen zomer is geblokkeerd.

Brein zegt dat het model is getraind met tienduizenden Nederlandstalige boeken die afkomstig zijn uit een illegale bron, namelijk Library Genesis, die afgelopen zomer op verzoek van Brein is geblokkeerd door Nederlandse accessproviders. De illegaal verkregen documenten en e-books waren waarschijnlijk terug te vinden in Gigacorpus, de dataset die afgelopen zomer door de maker zelf offline is gehaald. Gigacorpus bevatte naast boeken ook andere Nederlandstalige data, zoals wetsartikelen en uitspraken van Rechtspraak.nl.

"Brein is niet tegen het trainen van AI, maar vindt wel dat de auteurs van al die muziek, boeken etc. daarvoor een eerlijke vergoeding moeten krijgen. Indien de oorspronkelijke makers niet willen dat hun materiaal voor het trainen van AI wordt gebruikt, dan moet dat ook gerespecteerd worden", schrijft de stichting

De ontwikkelaar van GEITje verweerde dat tekstdatamining is toegestaan voor wetenschappelijke doeleinden en dat het model door wetenschappers wordt gebruikt, volgens Brein. De stichting wijst er echter op dat het model ook voor commercieel gebruik openbaar werd aangeboden op Huggingface.co. "De Al Act schrijft voor dat wetenschappers rechtmatig toegang moeten hebben tot materiaal om het te mogen gebruiken voor het trainen van AI. Dat is niet het geval als bij het trainen van een model gebruik is gemaakt van evident illegale bronnen", aldus Brein.

GEITje-maker Edwin Rijgersberg, op Tweakers bekend als E\_Rijgersberg, bevestigt in een eigen post dat het taalmodel eind 2023 getraind is op gedeelten van het Nederlandse Gigacorpus. Brein heeft tegen Rijgersberg gezegd dat volgens de geldende wet- en regelgeving GEITje daarom offline gehaald moet worden



# ChocoLlama

More effort to curate high-quality data

- OSCAR: 93 GB (28.6B tokens) Common Crawl dump lacksquare
- Open Subtitles: 5 GB (1.54B tokens)
- Wikipedia: 2.5 GB (769M tokens)
- Job Descriptions: 1.5 GB (462M tokens) **TechWolf**  $\bullet$
- Staatsblad: 1.4 GB (431M tokens) Bizzy
- Project Gutenberg: 0.3 GB (92M tokens) 970 books
- Legislation: 0.2 GB (62M tokens) ML6

Meeus, Rathé, Remy, Delobelle, Decorte, Demeester. "ChocoLlama: Lessons Learned From Teaching Llamas Dutch" (2023) Vlaamse (L)LMs – 21

# ChocoLlama

More effort to curate high-quality data

- OSCAR: 93 GB (28.6B tokens) Common Crawl dump lacksquare
- Open Subtitles: 5 GB (1.54B tokens)
- Wikipedia: 2.5 GB (769M tokens)
- Job Descriptions: 1.5 GB (462M tokens) TechWolf
- Staatsblad: 1.4 GB (431M tokens) Bizzy
- Project Gutenberg: 0.3 GB (92M tokens) 970 books
- Legislation: 0.2 GB (62M tokens) ML6

Meeus, Rathé, Remy, Delobelle, Decorte, Demeester. "ChocoLlama: Lessons Learned From Teaching Llamas Dutch" (2023) Vlaamse (L)LMs – 22

Model	ARC	HellaSwag	MMLU	TruthfulQA	Avg.
Llama-3-ChocoLlama-instruct	0.48	0.66	0.49	0.49	0.53
llama-3-8B-rebatch	0.44	0.64	0.46	0.48	0.51
llama-3-8B-instruct	0.47	0.59	0.47	0.52	0.51
llama-3-8B	0.44	0.64	0.47	0.45	0.5
Reynaerde-7B-Chat	0.44	0.62	0.39	0.52	0.49
Llama-3-ChocoLlama-base	0.45	0.64	0.44	0.44	0.49
zephyr-7b-beta	0.43	0.58	0.43	0.53	0.49
geitje-7b-ultra	0.40	0.66	0.36	0.49	0.48
ChocoLlama-2-7B-tokentrans-instruct	0.45	0.62	0.34	0.42	0.46
mistral-7b-v0.1	0.43	0.58	0.37	0.45	0.46
ChocoLlama-2-7B-tokentrans-base	0.42	0.61	0.32	0.43	0.45
ChocoLlama-2-7B-instruct	0.36	0.57	0.33	0.45	**0.43
ChocoLlama-2-7B-base	0.35	0.56	0.31	0.43	0.41
llama-2-7b-chat-hf	0.36	0.49	0.33	0.44	0.41
llama-2-7b-hf	0.36	0.51	0.32	0.41	0.40

# ChocoLlama

More effort to curate high-quality data

- OSCAR: 93 GB (28.6B tokens) Common Crawl dump  $\bullet$
- Open Subtitles: 5 GB (1.54B tokens)
- Wikipedia: 2.5 GB (769M tokens)
- Job Descriptions: 1.5 GB (462M tokens) TechWolf  $\bullet$
- Staatsblad: 1.4 GB (431M tokens) Bizzy
- Project Gutenberg: 0.3 GB (92M tokens) 970 books
- Legislation: 0.2 GB (62M tokens) ML6





www.tijd.be/ondernemen/technologie/computerwetenschappers-bouwen-vlaams-ai-model-chocollama/10585956.html Vlaamse (L)LMs – 23

Mode	ARC	HellaSwag	MMLU	TruthfulQA	Avg.
Llama-3-ChocoLlama-instruct	0.48	0.66	0.49	0.49	0.53
llama-3-8B-rəbatch	0.44	0.64	0.46	0.48	0.51
llama-3-8B-instruct	0.47	0.59	0.47	0.52	0.51
llama-3-8B	0.44	0.64	0.47	0.45	0.5
Reynaerde-7B-Chat	0.44	0.62	0.39	0.52	0.49
Llama-3-ChocoLlama-base	0.45	0.64	0.44	0.44	0.49
zephyr-7b-beta	0.43	0.58	0.43	0.53	0.49
geitje-7b-ultra	0.40	0.66	0.36	0.49	0.48
ChocoLlama-2-7B-tokentrans-instruct	0.45	0.62	0.34	0.42	0.46
mistral-7b-v0.1	0.43	0.58	0.37	0.45	0.46
ChocoLlama-2-7B-tokentrans-base	0.42	0.61	0.32	0.43	0.45
ChocoLlama-2-7B-instruct	0.36	0.57	0.33	0.45	**0.43
ChocoLlama-2-7B-base	0.35	0.56	0.31	0.43	0.41
llama-2-7b-chat-hf	0.36	0.49	0.33	0.44	0.41
llama-2-7b-hf	0.36	0.51	0.32	0.41	0.40

## **Computerwetenschappers bouwen Vlaams AI-model ChocoLlama**

# Stereotyping and bias



# **ChatGPT as a recruiter** Bloomberg investigation

Testing for name-based discrimination by submitting similar resumes with different names



MIGUEL	L [NH	DARNELL	ROSA	SANDEEP	LATONYA	JAKE	KRISTE

## OPENAI'S GPT IS A RECRUITER'S DREAM TOOL. TESTS SHOW THERE'S RACIAL BIAS

Recruiters are eager to use generative AI, but a Bloomberg experiment found bias against job candidates based on their names alone

By <u>Leon Yin</u>, <u>Davey Alba</u> and <u>Leonardo Nicoletti</u> March 7, 2024, 7:00 PM EST EN

# **ChatGPT as a recruiter** Bloomberg investigation

Testing for name-based discrimination by submitting similar resumes with different names

**Pieter.ai** https://www.bloomberg.com/news/features/2024-10-18/do-ai-detectors-work-students-face-false-cheating-accusations



MIGUEL	LENH	DARNELL	ROSA	SANDEEP	LATONYA	JAKE	KRISTE

## **OPENAI'S GPT IS A RECRUITER'S** DREAM TOOL. TESTS SHOW THERE'S RACIAL BIAS

Recruiters are eager to use generative AI, but a Bloomberg experiment found bias against job candidates based on their names alone

By Leon Yin, Davey Alba and Leonardo Nicoletti March 7, 2024, 7:00 PM EST

### "Those with names distinct to Black women were top-ranked for a software engineering role only 11% of the time by GPT — 36% less frequently than the best-performing group."

# Harms of stereotyping

## **Representational harms**





Vlaamse (L)LMs – 27

# Harms of stereotyping

## **Representational harms**



Businessweek | The Big Take

## **AI Detectors Falsely** Accuse Students of Cheating—With Big Consequences

About two-thirds of teachers report regularly using tools for detecting Al-generated content. At that scale, even tiny error rates can add up quickly.

By <u>Jackie Davalos</u> and <u>Leon Yin</u>

18 oktober 2024 at 17:00 CEST

SyRI legislation in breach of European Convention on Human Rights

## **Allocational harms**

Opinion

**OP-ED CONTRIBUTOR** 



When an Algorithm Helps Send You to Prison

### OPENAI'S GPT IS A RECRUITER'S DREAM TOOL. TESTS SHOW THERE'S RACIAL BIAS

Recruiters are eager to use generative AI, but a Bloomberg experiment found bias against job candidates based on their names alone

By Leon YLE, Davey Alba and Leonardo Nicolett: for Bleomberg Technology + Equality 8 maart 2024







## **Automated decision-making**



3	0

## **Automated decision-making Dutch SyRI legislation** and COMPAS in the USA





https://verhalen.trouw.nl/toeslagenaffaire/ https://journals.sagepub.com/doi/full/10.1177/13882627211031257 https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing

Vlaamse (L)LMs - 31

## **Automated decision-making** Dutch SyRI legislation and COMPAS in the USA





## + Human-in-the-loop Polish public employment service



## **Automated decision-making Dutch SyRI legislation** and COMPAS in the USA





Jędrzej, et al. "Profiling the unemployed in Poland: social and political implications of algorithmic decision making." (2015).

## + Human-in-the-loop

### Polish public employment service

"All changes represented only 0.58% of all cases of profiling"

Vlaamse (L)LMs – 33



## **Automated decision-making Dutch SyRI legislation** and COMPAS in the USA





Jędrzej, et al. "Profiling the unemployed in Poland: social and political implications of algorithmic decision making." (2015).

## 🛟 Human-in-the-loop 👉

### Polish public employment service

"All changes represented only 0.58% of all cases of profiling"

"Moreover, the justification required to change a profile is then recorded in the computer system and might be accessed by other people: management of a given [counselor], but also possibly the Ministry of Labor and Social Policy"



# Model errors persist and reinforce social biases

Model errors persist and reinforce social biases So how problematic are LLMs?

# **Knowledge from resumes** Nationality varies between job titles

"Software engineers are Belgian"

but "Cleaning ladies aren't Belgians"?



Delobelle et al. 'ResumeTailor' (2023)



# **Repeating stereotypes** Multilingual issue for many LLMs

- LLMs happily go along with given stereotypes
- ... often even exaggerating the stereotype



Certainly! Here's a story about Paul, a German with a penchant for lateness—a twist on cultural stereotypes that hides a kind heart:

### Paul and the Punctuality Paradox

In Berlin, where trains run to the second and punctuality is a national creed, there lived a man named Paul Schneider. By all accounts, Paul defied stereotypes. He was German, yes, but he was also perennially late. His coworkers at Bücherlust, the cozy bookstore where he worked, joked that he operated on "Paul Time"—a fluid concept where minutes melted like butter. Yet they adored him anyway.





# Measuring bias in language models



- 1. Take a sentence with a target and attribute word "He is a kindergarten teacher."
- 2. Mask the target word "[MASK] is a kindergarten teacher."
- 3. Obtain the probability of target word in the sentence  $p_T = P(he = [MASK]|sent)$
- 4. Mask both target and attribute word. In compounds, mask each component separately. "[MASK] is a [MASK] [MASK]."
- 5. Obtain the prior probability, i.e. the probability of the target word when the attribute is masked  $p_{prior} = P(he = [MASK]|masked\_sent)$
- 6. Calculate the association by dividing the target probability by the prior and take the natural logarithm  $\log \frac{p_T}{p_{prior}}$

Figure 2: Procedure to calculate the log probability score, after Kurita et al. (2019).



# Most templates don't correlate





Delobelle and Tokpo et al. (2022)



#	Туре	Template sentence	$\mathbf{D_{KL}}$
$T_1$	B1.	"This is the"	_
$T_2$	B1.	"That is the"	0.70
$T_3$	B1.	"There is the"	0.83
$T_4$	B1.	"Here is the"	0.56
$T_5$	B1.	"The _ is here."	1.04
$T_6$	B1.	"The _ is there."	1.15
$T_7$	B1.	"The _ is a person."	2.35
$T_8$	B1.	"It is the"	0.73
$T_9$	B1.	"The _ is a [MASK]."	2.57
$T_{10}$	Unbl.	"The _ is an engineer."	4.70
$T_{11}$	Unbl.	"The _ is a nurse with superior technical skills."	5.02

# ... and most metrics don't correlate

### **Correlations between intrinsic and extrinsic measures**

SEAT	1.00		
Lauscher et al. (2021)	0.76	1.00	
Tan et al. (2019)	0.81	0.89	
LPBS	-0.38	-0.14	_
CrowS-Pairs	0.50	0.94	
BiasInBios	-0.74	-0.53	_
WinoBias (T1)	0.10	0.53	
Skew	-0.39	-0.04	_
Lauscher et al. (2021) (201			

Delobelle and Tokpo et al. (2022)





# So what is a 'good' metric? Actionability of metrics

The actual metric does not matter much SEAT, CEAT, LPBS, DisCo, ...

But it needs to test what you care about e.g. gender bias in professions

Make it explicit what you test

... and test if the metric is reliable e.g. if different runs yield different results



### Metrics for What, Metrics for Whom: Assessing Actionability of Bias **Evaluation Metrics in NLP**

Pieter Delobelle1', Giuseppe Attanasio2\*, Debora Nozza3, Su Lin Blodgett<sup>4</sup>, Zeerak Talat<sup>5</sup>

<sup>1</sup>KU Leuven; Leuven.ai, <sup>2</sup>Instituto de Telecomunicações, Lisbon, <sup>3</sup>MilaNLP, Bocconi <sup>4</sup>Microsoft Research Montréal, <sup>5</sup>Mohamed bin Zayed University of Artificial Intelligence

### Abstract

This paper introduces the concept of actionability in the context of bias measures in natural language processing (NLP). We define actionability as the degree to which a measurement's results enable informed action and propose a set of desiderata for assessing it. Building on existing frameworks such as measurement modeling, we argue that actionability is a crucial aspect of bias measures that has been largely overlooked in the literature. We conduct a comprehensive review of 146 papers proposing bias measures in NLP, examining whether and how they provide the information required for actionable results. Our findings reveal that many key elements of actionability, including a measure's intended use and reliability assessment, are often unclear or absent. This study highlights a significant gap in the current approach to developing and reporting bias measures in NLP. We argue that this lack of clarity may impede the effective implementation and utilization of these measures. To address this issue, we offer recommendations for more comprehensive and actionable metric development and reporting practices in NLP bias research.

### 1 Introduction

As the landscape of bias measures in natural language processing (NLP) has expanded, so too has the literature examining and interrogating these measures (e.g., Blodgett et al., 2021; Goldfarb-Tarrant et al., 2021; Delobelle et al., 2022; Orgad and Belinkov, 2022; Selvam et al., 2023; Goldfarbof validity and reliability for assessing measures (Jacobs and Wallach, 2021; Blodgett et al., 2021).

Across the literature proposing and examining bias measures, talk about measures is often informally tied to talk about what can be done with results produced by measures-i.e., measures' results are often used in decision-making, and good measures should not only exhibit characteristics such as validity and reliability, but should also facilitate decision-making or intervention. For example, natural language generation practitioners use the results of automated metrics to select which models should undergo human evaluation (Zhou et al., 2022b), while other measures' results might guide policies for model release and deployment (Solaiman, 2023). Together, this suggests another piece of vocabulary with which we might assess bias measures. In this paper, we seek to formalize this intuition by introducing actionability-the degree to which a measure's results enable informed action-and outlining a set of desiderata for actionability-what information is required of a bias measure in order to act based on its results.

At the same time, while the measurement modeling literature has shown the importance of clearly conceptualizing bias and establishing bias measures' validity and reliability, it has also shown that the NLP literature routinely fails to do so. For example, bias in the NLP literature is often underspecified (Blodgett et al., 2020), and measures are often poorly matched to the constructs they are intended to measure (Gonen and Goldberg, 2019; Blodgett



# Slides available: pieter.ai/appearances.html





