

SAVE
THE
DATE

sign up and more info:
projects.illc.uva.nl/indeep

InDeep Masterclass

Explaining Foundation Models

Thursday 2 November 2023
10.30-15.00 (with lunchbreak)

SustainaLab (Matrix ONE)
Science Park 301
Amsterdam

BONUS:
Immediately after the Masterclass there will be a **Meetup** by **InDeep & Amsterdam AI:**
Alternatives for ChatGPT: How good are open source and in-house LLMs?
15.00-17.00 (with drinks)

Program

Grzegorz Chrupala (Tilburg University)
Opening the Blackbox of Large Language Models, and other deep learning models that handle speech, vision and language

Arianna Bisazza (RUG Groningen) and **Gabriele Sarti** (RUG Groningen)
Hands-on tutorial on Interpreting Large Language Models with InSeq

Discussion on Promises and Pitfalls for Explainability in the Generative AI Boom, with **Jelle Zuidema** (UvA) and others [TBA]

On Thursday 2 November 2023 **InDeep** organized a Masterclass on *Explaining Foundation Models* with presentations, discussion and a hands-on tutorial. After the Masterclass there was a Meetup organized by Indeep and Amsterdam AI



Grzegorz Chrupała (Tilburg University) gave a presentation on opening the Blackbox of Large Language Models and other deep learning models that handle speech, vision and language



Opening the black box of foundation models

Grzegorz Chrupała

What is a foundation model?

- Trained in a self-supervised way (without manual annotations).
- On large amounts of data.
 - Text, audio, images, video
- Starting point for many **downstream** tasks.

How it works

- Language model trained on speech.
- Objective: identify masked audio segment.
- Additionally **tuned** for
 - Transcribe speech into text
 - Recognize spoken commands
 - Identify speakers
 - Classify emotions

Towards **multiple** modalities

- For example audio-visual models learn from images (or videos) accompanied by speech.
- Most recent version of ChatGPT can process text, images, communicate with voice.

Foundation models are the core of AI.

We want AI models to **perform well**.

What other **properties** are **desirable**?

Why do we want to understand models?

- **Explain** decisions to end-users
 - As humans we seek justifications.
 - Arguably, human reasoning is designed for argumentation.
 - We expect to AIs to be able to justify themselves the way we do.
 - Justifications of automated decisions may be legally mandated
- **Understand** models as scientists and engineers
 - Knowing when and why models fail gives us insight into
 - The phenomena modeled.
 - Ways to improve model performance

There is a lot of research activity dedicated to **understanding** and **explicating** how **deep learning** works

What feature of DL makes this **interesting** and/or **necessary**?

Distributed representations

- DL learns **distributed** representations.
 - Information is encoded in **patterns** of activation.
 - People are better at grasping information encoded **symbolically**.

Size and complexity

- Foundation models are **large**
 - Millions or billions of learnable parameters
- They are composed of dozens of types of specialized modules
 - Feedforward, Convolutions, Attention, Normalization, Recurrences, ...
- Unlike in many other ML models, the path from input features to outputs is highly **indirect**.

What does it mean to **understand** how a model works?

What kind of **questions** would we like to answer about this?

Which features of an input example led to a particular prediction?

- **Local** explanation, for a particular datapoint.
- Aka **feature attribution**.
- What **minimal change** to this datapoint would cause a change in the prediction?

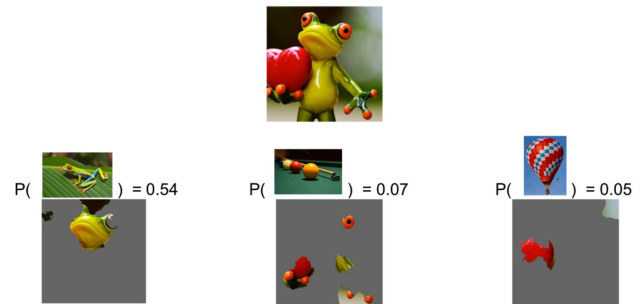


Figure 6. Explanation for a prediction from Inception. The top three predicted classes are "tree frog," "pool table," and "balloon." Sources: Marco Tulio Ribeiro, Pixabay (frog, billiards, hot air balloon).

<https://www.oreilly.com/content/introduction-to-local-interpretable-model-agnostic-explanations-lime/>

Model-agnostic methods

- Occlusion analysis
 - Features or groups of features are occluded and effect on prediction computed.
- LIME – Local Interpretable Model-Agnostic Explanations
- SHAP – Shapley additive explanations

Model-specific methods

- Gradient-based methods
 - Use **gradient** of target function to determine what features affect the output of the model.
- Methods tailored to specific architectures
 - Many recent approaches focus on **Transformers**
 - **ValueZeroing** exploits the mathematical formulation of the transformer to exclude tokens from computation at specific layers (<http://dx.doi.org/10.18653/v1/2023.eacl-main.245>).

Which input features are important for model decisions?

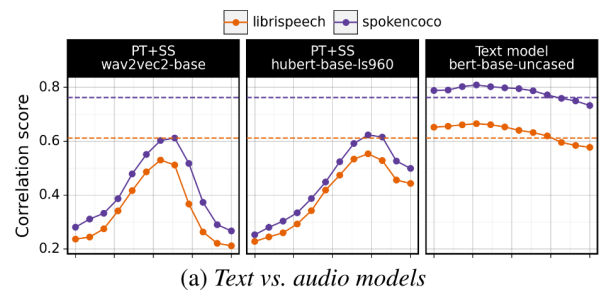
- **Global** explanations, aggregated over multiple datapoints.

Important caveats

- What counts as a feature for the purposes of explanation?
- The use of feature attribution or importance scores assumes a (locally) linear explanatory model.
 - May or may not be a satisfactory approximation to the true function.

What information is encoded in activation patterns?

- What kind of information is encoded?
- To what extent?
- Where is it encoded, in terms of layers, or individual neurons?
- How can we visualize it or otherwise make it comprehensible?



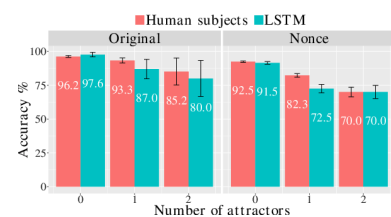
<https://dx.doi.org/10.21437/Interspeech.2023-679>

How does the model behave when faced with particular phenomena?

- Select or manipulate stimuli to highlight phenomena of interest, and analyze model **behavior** on them.
- Analogous to how a psychologist works.

Subject-verb number agreement

“Yet the ratio of men who survive to the women and children who survive **[is/are]**”



<https://doi.org/10.18653/v1%2FN18-1108>

Which (symbolic) algorithm does the model implement?

Often analysis is only feasible for a simplified version of an architecture of interest.

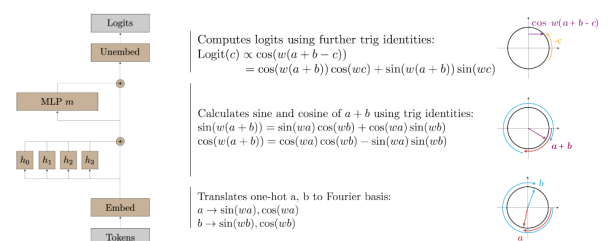


Figure 1: The algorithm implemented by the one-layer transformer for modular addition. Given two numbers a and b , the model projects each point onto a corresponding rotation using its embedding matrix. Using its attention and MLP layers, it then composes the rotations to get a representation of $a+b \bmod P$. Finally, it “reads off” the logits for each $c \in \{0, 1, \dots, P-1\}$, by rotating by $-c$ to get $\cos(w(a+b-c))$, which is maximized when $a+b \equiv c \bmod P$ (since w is a multiple of $\frac{2\pi}{P}$).

<https://doi.org/10.48550/arXiv.2301.05217>

Evaluation

How can explanations or analyses be **evaluated**?

Main evaluation types

- **Plausibility**

- does the explanation seem reasonable from the human point of view?

- **Faithfulness**

- does the explanation accurately capture the factors behind the model's decision?

Explaining foundation models

- Very active area of research.
- Important new developments happen rapidly.
- Many open problems.



Arianna Bisazza (RUG Groningen) and Gabriele Sarti (RUG Groningen) gave a hands-on tutorial on interpreting Large Language Models with InSeq



Interpreting (Large) Language Model Generations with Inseq

Gabriele Sarti & Arianna Bisazza

InDeep Masterclass – November 2nd, 2023



university of
 groningen / center for language
 and cognition



Dutch Research
 Council



university of
 groningen



Dutch Research
 Council



About Us



Gabriele Sarti
3rd-year PhD student



Arianna Bisazza
Associate professor



@ University of
 Groningen

→ InDeep Work Package on improving interpretability for
 machine translation from a user-centric perspective

→ **Core working methods:** Feature attribution for
 text-generating models



university of
 groningen



Dutch Research
 Council



Before we begin, let's set up our notebook:

[https://tinyurl.com/
 indeep-masterclass-ex](https://tinyurl.com/indeep-masterclass-ex) OR



university of
 groningen



Dutch Research
 Council



Feature Attribution

- ❖ Understanding why Neural Language Models make certain decisions when generating text is hard!
- ❖ A promising approach to do that is Feature Attribution

What is Feature Attribution?

Which part(s) of the input sample contributed most to a model's prediction?



A woman is throwing a frisbee in a park.

A large white bird standing in a forest.

Feature attribution examples for an image captioning model, from [\[Xu et al. 2016\]](#)

What is Feature Attribution?

Which part(s) of the input sample contributed most to a model's prediction?

	$p(y x)$	y
the year 's best and most unpredictable comedy	0.91	pos
we never feel anything for these characters	0.95	neg

Sentiment analysis example from [\[Madsen et al. 2021\]](#)

What is Feature Attribution?

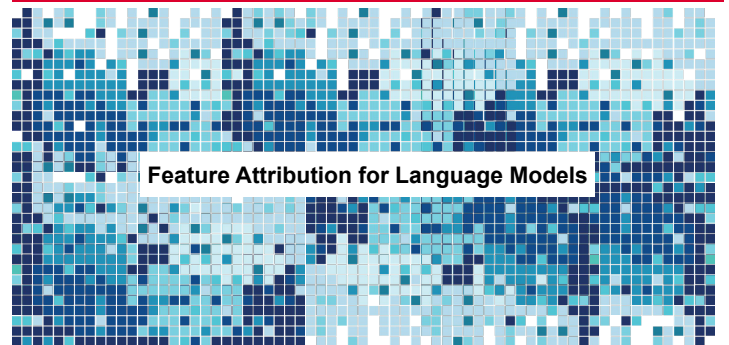
Which part(s) of the input sample contributed most to a model's prediction?

	$p(y x)$	y
the year 's <u>best</u> and most <u>unpredictable</u> comedy	0.91	pos
we <u>never</u> <u>feel</u> <u>anything</u> <u>for</u> <u>these</u> <u>characters</u>	0.95	neg

Sentiment analysis example from [\[Madsen et al. 2021\]](#)

What is Feature Attribution?

Which part(s) of the input sample contributed most to a model's prediction?

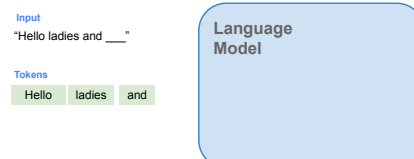


Feature Attribution for Language Models

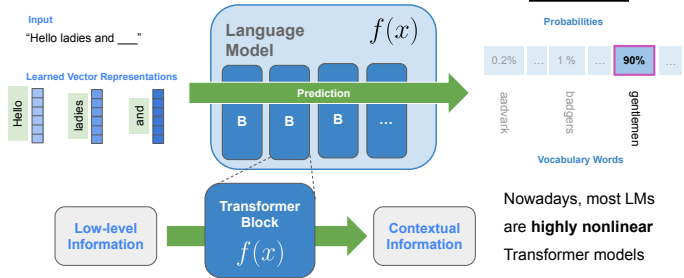
Feature Attribution for Language Models

- ❖ **Urgent!**
 - More and more NLP tasks (even classification ones) are solved by text-generating models nowadays
 - Understanding the relationship between a model's inputs & output can help us detect unwanted **errors and biases**
- ❖ **Challenging!** More difficult than classical classification
 - Large output space
 - At each timestep t , previously generated words affect prediction (i.e. input features are *dynamic*)

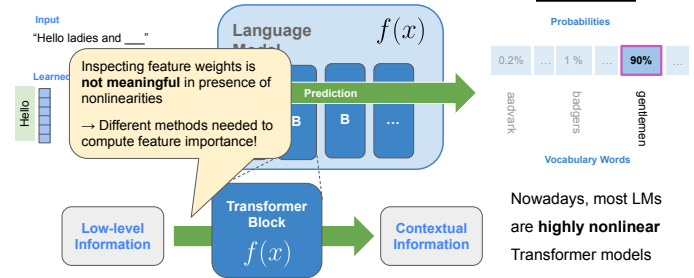
How Do Language Models Work?



How Do Language Models Work?

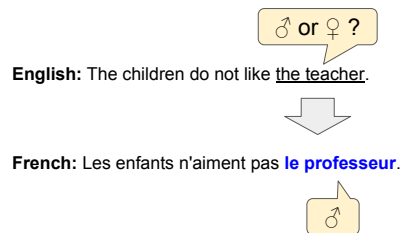


How Do Language Models Work?

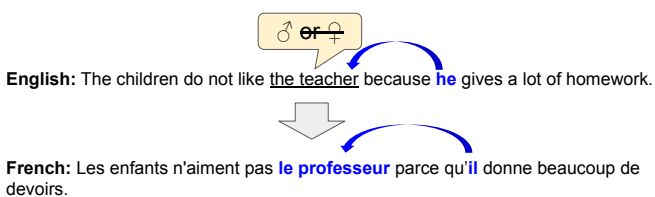


A Practical Example

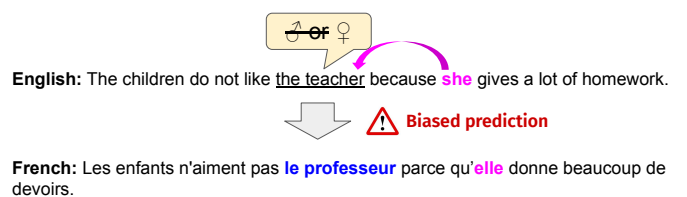
Use Case: Feature Attrib. & Gender Bias in Translation



Use Case: Feature Attrib. & Gender Bias in Translation



Use Case: Feature Attrib. & Gender Bias in Translation



Why does that happen??

Use Case: Feature Attrib. & Gender Bias in Translation

	__Les	__enfants	__n	'	ai	ment	__pas	__le	__professeur	__parce	__qu	'	elle
__The	0.047	0.046	0.027	0.02	0.022	0.011	0.016	0.014	0.018	0.019	0.016	0.011	0.013
__children	0.071	0.195	0.057	0.068	0.044	0.021	0.037	0.043	0.05	0.032	0.04	0.025	0.026
__do	0.039	0.036	0.072	0.031	0.059	0.032	0.038	0.01	0.02	0.029	0.021	0.02	0.01
__not	0.077	0.05	0.072	0.056	0.066	0.035	0.078	0.014	0.038	0.041	0.039	0.029	0.015
__like	0.036	0.042	0.161	0.033	0.184	0.074	0.068	0.024	0.02	0.034	0.023	0.027	0.02
__the	0.04	0.042	0.081	0.039	0.069	0.032	0.034	0.091	0.066	0.04	0.028	0.022	0.033
__teacher	0.061	0.103	0.104	0.05	0.076	0.033	0.048	0.314	0.256	0.05	0.042	0.034	0.051
__because	0.04	0.026	0.052	0.035	0.032	0.017	0.031	0.038	0.028	0.164	0.07	0.038	0.057
__she	0.056	0.053	0.033	0.053	0.025	0.016	0.027	0.114	0.052	0.044	0.132	0.044	0.181
__gives	0.04	0.028	0.028	0.029	0.027	0.016	0.02	0.015	0.019	0.052	0.081	0.04	0.081

Use Case: Feature Attrib. & Gender Bias in Translation

	__Les	__enfants	__n	'	ai	ment	__pas	__le	__professeur	__parce	__qu	'	elle
__The	0.047	0.046	0.027	0.02	0.022	0.011	0.016	0.014	0.018	0.019	0.016	0.011	0.013
__children	0.071	0.195	0.057	0.068	0.044	0.021	0.037	0.043	0.05	0.032	0.04	0.025	0.026
__do	0.039	0.036	0.072	0.031	0.059	0.032	0.038	0.01	0.02	0.029	0.021	0.02	0.01
__not	0.077	0.05	0.072	0.056	0.066	0.035	0.078	0.014	0.038	0.041	0.039	0.029	0.015
__like	0.036	0.042	0.161	0.033	0.184	0.074	0.068	0.024	0.02	0.034	0.023	0.027	0.02
__the	0.04	0.042	0.081	0.039	0.069	0.032	0.034	0.091	0.066	0.04	0.028	0.022	0.033
__teacher	0.061	0.103	0.104	0.05	0.076	0.033	0.048	0.314	0.256	0.05	0.042	0.034	0.051
__because	0.04	0.026	0.052	0.035	0.032	0.017	0.031	0.038	0.028	0.164	0.07	0.038	0.057
__she	0.056	0.053	0.033	0.053	0.025	0.016	0.027	0.114	0.052	0.044	0.132	0.044	0.181
__gives	0.04	0.028	0.028	0.029	0.027	0.016	0.02	0.015	0.019	0.052	0.081	0.04	0.081

Use Case: Feature Attrib. & Gender Bias in Translation

	__Les	__enfants	__n	'	ai	ment	__pas	__le	__professeur	__parce	__qu	'	elle
__The	0.047	0.046	0.027	0.02	0.022	0.011	0.016	0.014	0.018	0.019	0.016	0.011	0.013
__children	0.071	0.195	0.057	0.068	0.044	0.021	0.037	0.043	0.05	0.032	0.04	0.025	0.026
__do	0.039	0.036	0.072	0.031	0.059	0.032	0.038	0.01	0.02	0.029	0.021	0.02	0.01
__not	0.077	0.05	0.072	0.056	0.066	0.035	0.078	0.014	0.038	0.041	0.039	0.029	0.015
__like	0.036	0.042	0.161	0.033	0.184	0.074	0.068	0.024	0.02	0.034	0.023	0.027	0.02
__the	0.04	0.042	0.081	0.039	0.069	0.032	0.034	0.091	0.066	0.04	0.028	0.022	0.033
__teacher	0.061	0.103	0.104	0.05	0.076	0.033	0.048	0.314	0.256	0.05	0.042	0.034	0.051
__because	0.04	0.026	0.052	0.035	0.032	0.017	0.031	0.038	0.028	0.164	0.07	0.038	0.057
__she	0.056	0.053	0.033	0.053	0.025	0.016	0.027	0.114	0.052	0.044	0.132	0.044	0.181
__gives	0.04	0.028	0.028	0.029	0.027	0.016	0.02	0.015	0.019	0.052	0.081	0.04	0.081

Use Case: Feature Attrib. & Gender Bias in Translation

	__Les	__enfants	__n	'	ai	ment	__pas	__le	__professeur	__parce	__qu	'	elle
__The	0.047	0.046	0.027	0.02	0.022	0.011	0.016	0.014	0.018	0.019	0.016	0.011	0.013
__children	0.071	0.195	0.057	0.068	0.044	0.021	0.037	0.043	0.05	0.032	0.04	0.025	0.026
__do	0.039	0.036	0.072	0.031	0.059	0.032	0.038	0.01	0.02	0.029	0.021	0.02	0.01
__not	0.077	0.05	0.072	0.056	0.066	0.035	0.078	0.014	0.038	0.041	0.039	0.029	0.015
__like	0.036	0.042	0.161	0.033	0.184	0.074	0.068	0.024	0.02	0.034	0.023	0.027	0.02
__the	0.04	0.042	0.081	0.039	0.069	0.032	0.034	0.091	0.066	0.04	0.028	0.022	0.033
__teacher	0.061	0.103	0.104	0.05	0.076	0.033	0.048	0.314	0.256	0.05	0.042	0.034	0.051
__because	0.04	0.026	0.052	0.035	0.032	0.017	0.031	0.038	0.028	0.164	0.07	0.038	0.057
__she	0.056	0.053	0.033	0.053	0.025	0.016	0.027	0.114	0.052	0.044	0.132	0.044	0.181
__gives	0.04	0.028	0.028	0.029	0.027	0.016	0.02	0.015	0.019	0.052	0.081	0.04	0.081

Use Case: Feature Attrib. & Gender Bias in Translation

♂ or ♀ ?

English: The children do not like the pretty teacher because **she** gives a lot of homework.



French: Les enfants n'aiment pas **la jolie enseignante** parce qu'**elle** donne beaucoup de devoirs

Use Case: Feature Attrib. & Gender Bias in Translation

	__Les	__enfants	__n	'	ai	ment	__pas	__la	__jolie	__enseignante	e
__The	0.038	0.044	0.033	0.017	0.022	0.01	0.015	0.012	0.009	0.013	0.008
__children	0.062	0.195	0.06	0.061	0.045	0.019	0.038	0.023	0.021	0.046	0.021
__do	0.043	0.039	0.074	0.044	0.058	0.029	0.039	0.014	0.012	0.017	0.01
__not	0.094	0.061	0.085	0.088	0.062	0.034	0.084	0.025	0.023	0.035	0.019
__like	0.032	0.032	0.167	0.038	0.177	0.069	0.068	0.025	0.021	0.02	0.014
__the	0.042	0.038	0.066	0.033	0.062	0.028	0.031	0.077	0.081	0.046	0.024
__pretty	0.06	0.064	0.073	0.059	0.068	0.032	0.054	0.186	0.361	0.074	0.047
__teacher	0.047	0.083	0.066	0.037	0.055	0.026	0.036	0.19	0.119	0.241	0.082
__because	0.036	0.025	0.04	0.038	0.028	0.014	0.026	0.041	0.03	0.028	0.028
__she	0.061	0.05	0.024	0.043	0.021	0.017	0.03	0.108	0.029	0.049	0.038
__gives	0.039	0.027	0.024	0.027	0.024	0.015	0.019	0.027	0.015	0.02	0.016

Use Case: Feature Attrib. & Gender Bias in Translation

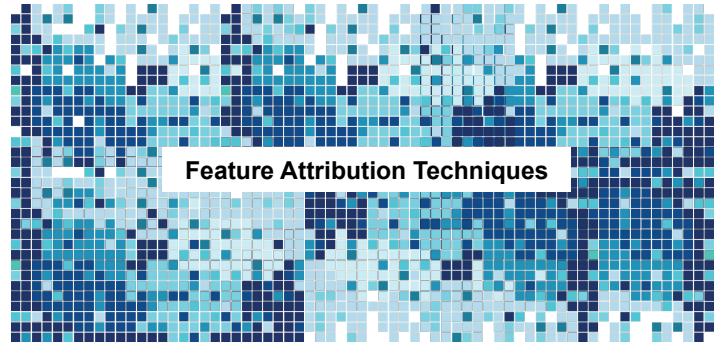
	$p=0.37$		$p=0.53$
__le		__le	
__The	0.014	__The	0.015
__children	0.043	__children	0.046
__do	0.01	__do	0.01
__not	0.014	__not	0.017
__like	0.024	__like	0.026
__the	0.091	__the	0.093
__teacher	0.314	__teacher	0.349
__because	0.038	__because	0.042
__she	0.114	__he	0.052
__gives	0.015	__gives	0.017



This model seems to pay little attention to the pronoun cue

→ Unsafe!

Feature Attribution Techniques



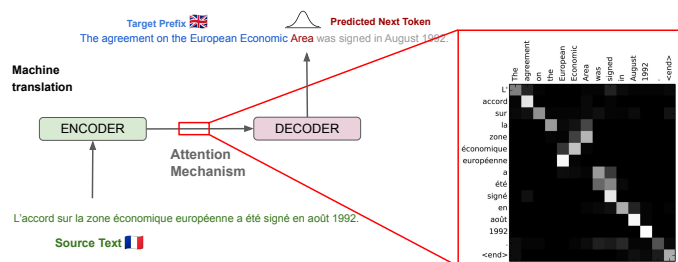
Feature Attribution & Broader Interpretability

- ❖ A form of **post-hoc** interpretability (vs. models interpretable by design)
- ❖ **Local**: explains a single observation (vs. general model functioning)
- ❖ **Low abstraction**: explains model decision in terms of the input features (vs. in terms of relevant training examples; vs. explanations in natural language)

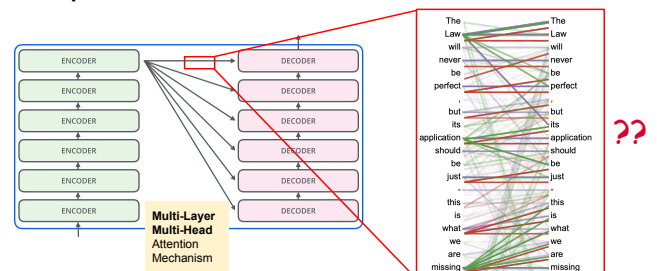
Types of Feature Attribution Methods

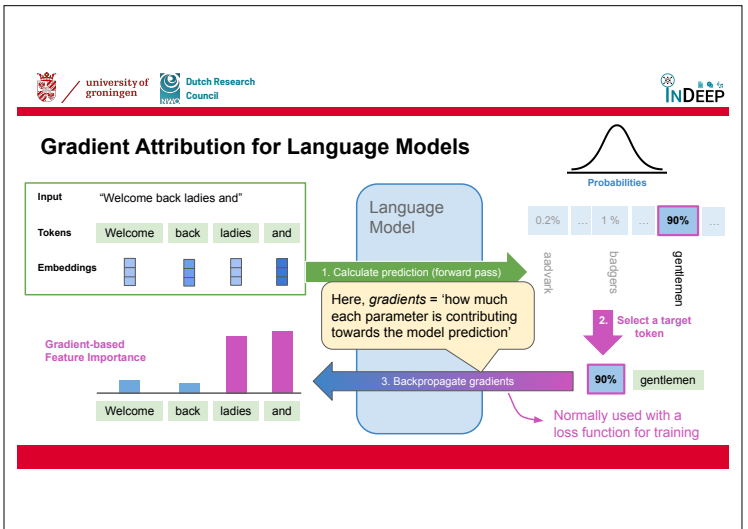
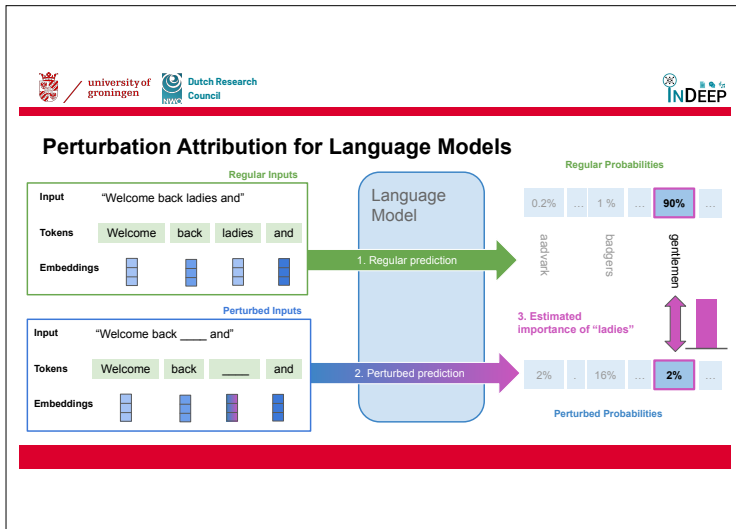
- ❖ Model-specific (or internals-based)
- ❖ Perturbation-based
- ❖ Gradient-based

Model Specific: Attention as Attribution



Model Specific: Attention as Attribution in Transformers





university of groningen Dutch Research Council INDEEP

In sum: Feature Attribution Methods

- ❖ Perturbation-based
- ❖ Gradient-based
- ❖ Model-specific (or internals-based)

No clear winner (yet)!

All existing attribution methods have some **intrinsic** and/or **practical** limitations

Some methods require **full access** to model internals, other don't

Research is **ongoing**

university of groningen Dutch Research Council INDEEP

Inseq : Interpretability for Sequence Generation Models

Pytorch-based Python library to simplify usage and evaluation of **feature attribution** methods for **generative language models**.

Demo @ 61 ACL 2023

215 24

university of groningen

Deutsches Forschungszentrum für Künstliche Intelligenz (DFKI)

UNIVERSITY OF AMSTERDAM

university of groningen Dutch Research Council INDEEP

Inseq: Core Functionalities

university of groningen Dutch Research Council INDEEP

Inseq : Interpretability for Sequence Generation Models

Sequence Model

To innovate one should

Autoregressive Generation

think outside the box

Demo @ 61 ACL 2023

215 24

university of groningen

Deutsches Forschungszentrum für Künstliche Intelligenz (DFKI)

UNIVERSITY OF AMSTERDAM

	think	outside	the	box
To	0	0.02	0	0
innovate	0.7	0.45	0.2	0.15
one	0.05	0.04	0	0.02
should	0.25	0.02	0	0.02
think		0.47	0.4	0.15
outside			0.4	0.35
the				0.35
$P(y_t y_{<t})$	0.5	0.82	0.94	0.99

A Simple Example

```
import inseq

# Load HF Hub model and attribution method
model = inseq.load_model(
    "google/flan-t5-base",
    "integrated_gradients"
)

# Answer and attribute generation steps
attr_out = model.attribute(
    "Does 3 + 3 equal 6?",
    attribute_target=True
)

# Visualize the generated attribution,
# applying default token-level aggregation
attr_out.show()
```

Sarri et al. 2023

Source Saliency			Prefix Saliency		
	_yes	</s>		_yes	</s>
_Does	0.264	0.153			0.21
_3	0.113	0.099			
_+	0.096	0.086			
_3	0.096	0.076			
_equal	0.213	0.154			
_6	0.11	0.108			
?	0.106	0.114			

🗨️: Does 3 + 3 equal 7? 🗨️: yes

Main Features

Method	$f(l)$
(Input ×) Gradient	✓
DeepLIFT	✓
G GradientSHAP	✓
Integrated Gradients	✓
Discretized IG	✗
I Attention Weights	✓
Occlusion (Blank-out)	✗
P LIME	✗

G: Gradient-based,
I: Internals-based,
P: Perturbation-based,
f(l): layer-wise attribution

- Seq2SeqLM and CausalLM support
- Aggregate, save and load attributions
- Batch-attribution datasets in Jupyter & CLI

Sarri et al. 2023

Contrastive Gradient Attribution

Input: Can you stop the dog from barking

Output: barking

- Why did the model predict "barking"?
- Why did the model predict "barking" instead of "crying"?

Can you stop the dog from

Disentangle fluency factors from the overall model decision, more aligned with human intuition (🧠 simulatability)

Gradient Attribution

$$\nabla_{y<t} q(y_t|y_{<t})$$

e.g. barking

Contrastive Gradient Attribution

$$\nabla_{y<t} (q(y_t|y_{<t}) - q(\hat{y}_t|y_{<t}))$$

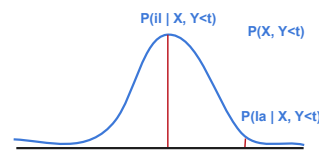
e.g. barking e.g. crying

Yin and Neubig, 2022

Step Functions and Custom Attribution Targets

X = I said hi to the manager

Y<t = Ho salutato ____



Many step functions already available (Entropy, CE, CMI, PPL, KL-Div, ...)

+ custom step functions registration to add yours!

Contrastive Attribution Inseq Tutorial

Source and Target Attributions									
Attributed	Generated								
		to	la	la	manager	</s>			
_Ho	_saluta	0.0	0.0	0.0	0.06	0.0	0.0		
_I		0.0	0.0	0.0	0.077	0.0	0.0		
_said		0.0	0.0	0.0	0.089	0.0	0.0		
_hi		0.0	0.0	0.0	0.042	0.0	0.0		
_to		0.0	0.0	0.0	0.104	0.0	0.0		
_the		0.0	0.0	0.0	0.111	0.0	0.0		
_manager		0.0	0.0	0.0	0.517	0.0	0.0		
</s>		0.0	0.0	0.0	0.111	0.0	0.0		
probability	0.44	0.844	0.916	0.785	0.331	0.26			
contrast_prob	0.44	0.844	0.916	0.078	0.319	0.303			
contrast_prob_diff	0.0	0.0	0.0	-0.707	0.013	-0.043			

Hands-On Notebook

Examples on Google Colab



OR

<https://tinyurl.com/indeep-masterclass-ex>

Interested in Inseq? Let us know!



Contributions and feedback are more than welcome! 🙌

- [Github](#) - inseq-team/inseq
- [Documentation](#) - inseq.readthedocs.io
- [Hugging Face Hub](#) - inseq
- [X \(Twitter\)](#) - @InseqDev
- [Discord](#) - see Github README

Recap & Promising Directions

After lunch: Let's discuss!

- ❖ Take a post-it, write on it:
 - your **name**
 - a **topic** that you'd like to discuss (e.g. task of interest, use case, an obstacle preventing you from using feature attribution in your work, a wish from research ...)
- ❖ Paste it on the whiteboard
- ❖ After lunch we'll discuss in **groups**
- ❖ **Now:** Some success stories and promising directions to inspire you!

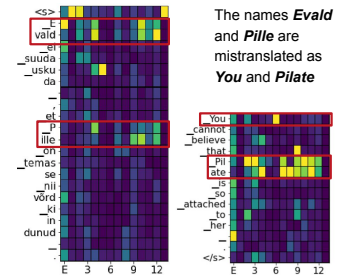
Your name

Topic of interest

Example: Coarse-to-Fine-grained Quality Estimation

Goal: Identify word-level error spans from sentence-level neural metrics (Fomicheva et al. 2022)

- ❖ Useful for obtaining **data-driven granular quality information**.
- ❖ Applications for human preference prediction currently under-explored (e.g. granular explanations for reward model predictions in RLHF)



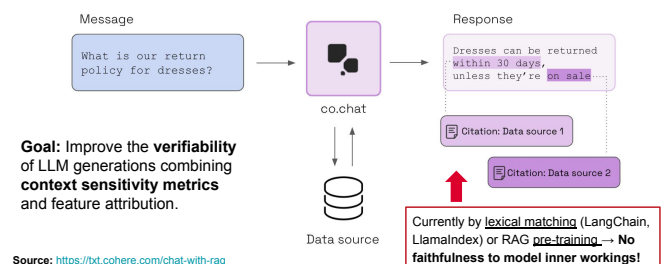
Example: Hallucination Detection

Goal: Detect hallucinations and omissions from context usage metrics (Dale et al. 2023)

- ❖ No additional training/modules, detect issues from **model internals** only.
- ❖ Can be used for **detoification** (NLLB Team, 2022)

Input	Output	Attribution-based	Similarity-based	
8 However, the discovery of his tomb in 1922 made him a celebrity. <<<While many tombs of the past were robbed, this tomb was left virtually undisturbed>>>.	Sin embargo, el descubrimiento de su tumba en 1922 le convirtió en una celebridad.	-0.66 0.56 0.29 0.79		true positive omission
9 :<<<Q>>>..	<<<Por qué no lo haces>>>	-2.18 0.37 0.67 0.29		true positive hallucination
10 <<<2_2010-11-0104.56>>>	<<<Se trata de un proyecto de ley que se ha desarrollado en el ámbito de la seguridad social>>>.	-2.40 0.23 0.44 0.09		true positive hallucination
11 Bye.	Hasta luego, <<<amigo>>>.	-2.58 0.58 1.02 0.63		true positive hallucination

Example: Trustworthy Retrieval-Augmented LLMs



Goal: Improve the **verifiability** of LLM generations combining **context sensitivity metrics** and feature attribution.

Source: <https://txt.cohere.com/chat-with-rag>



Plausibility Evaluation for Context Reliance in LMs

(Most goods from Central America came here duty-free)

Target Context (FR): La plupart des biens d'Amérique centrale sont venus ici en franchise.

Source (EN): Yet eighty percent of the goods were taxed in Central American countries.

Contextless target (FR): Pourtant, 80% de nos marchandises ont été taxées dans les pays d'Amérique Centrale

Contextual target (FR): Pourtant, 80% de nos biens ont été taxés dans les pays d'Amérique centrale

Context-sensitive target and contextual cues (found by PECORe) and contextless variant of the corresponding token. Other changes are not found to be context-dependent.

→ A new data-driven method to debug (im)plausible context dependence in LMs! [Sarti et al. 2023](#)

Thank you for your attention!

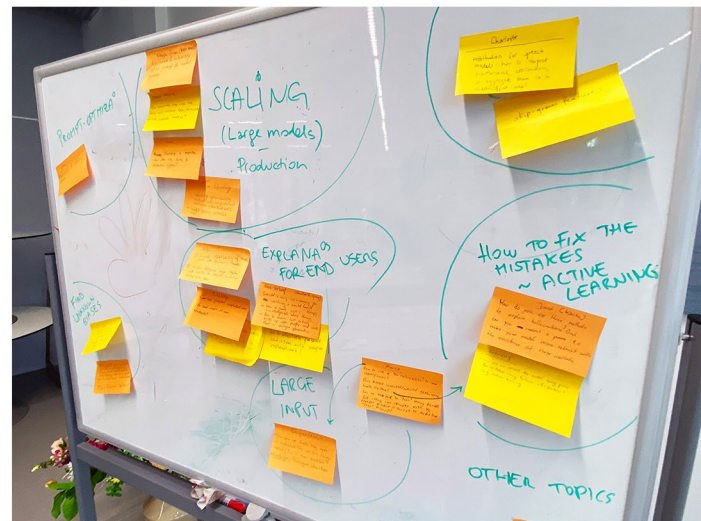
@AriannaBisazza & @gsarti_

<https://www.cs.rug.nl/~bisazza/> & [gsarti.com](https://www.cs.rug.nl/~gsarti/)

a.bisazza@rug.nl & g.sarti@rug.nl



Jelle Zuidema (UvA) discusses the promises and pitfalls for Explainability in the Generative AI boom, while connecting it to the work done by the hands-on groups



After the Masterclass, InDeep teamed up with Amsterdam AI for a Meetup on Alternatives for ChatGPT: *How good are open source and in-house LLMs?*

Harm de Vries (ServiceNow) presented the Big Code project, the pioneering project to create an open source dataset and generative AI model for generating code. He discussed details of the effort to determine the optimal dataset and model size, legal and ethical concerns about the dataset



Matthieu Laneuville (SURF) discussed the motivation and plans for an open source, community-driven Large Language Model for Dutch: GPT-NL





Open and responsible development of
Large Language Models for code

Harm de Vries | Staff Research Scientist @ ServiceNow | @harmdevries77

Today's talk

- 1 The BigCode Community
- 2 The Stack
- 3 StarCoder

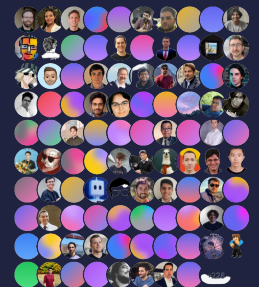
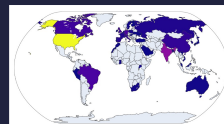
1. The BigCode community



BigCode: open-scientific collaboration

We are building LLMs for
code in a collaborative
way:

- 1100+ members
- 60+ countries



Developing Code LLMs is not only a technical problem!

- **Consent of data subjects**
 - Can you use third-party copyrighted data to train machine learning models?
- **Privacy risks**
 - Do LLMs distribute personally identifiable information without consent?
- **Software safety and security**
 - Code LLMs may be used to generate malware or may provide code suggestions that are less secure

Scraped public
data sources



More info on the [Governance Card](#)

Closed development of LLMs

- **Training data is not disclosed:**
 - Content creators don't know if their data is used and there's no way to remove it
 - Limits scientific reproducibility
 - Potential benchmark contamination
- **Model only available through API, which limits research on:**
 - Safety and alignment
 - The model's inner workings (i.e. representations)
 - Adaptation methods like LoRA, and continuous prompt-tuning

Open & Responsible Research on LLMs

Open-access datasets

[Permissive licensing](#)
[Data inspection](#)
[Opt-out available](#)
[PII removal](#)
[Attribution](#)

Open-access models

[Model weights available](#)
[Fine-tuning scripts](#)
[Low-precision inference](#)

Reproducible research

[Data preprocessing scripts](#)
[Model training framework](#)
[R&D notebooks](#)
[Evaluation Harness](#)

Documentation

[Dataset cards](#)
[Model cards](#)
[Governance card](#)
[Intellectual property](#)
[Code of conduct](#)
[OpenRAIL licenses](#)

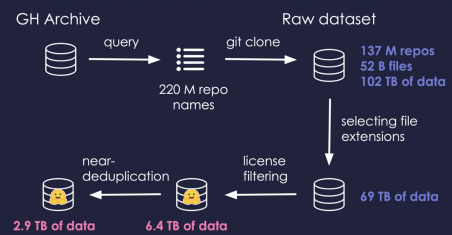
2. The Stack



The Stack

A dataset with **6.4TB** of **permissively licensed** source code in **358 programming languages** with a **data inspection** tool and **opt-out** mechanism

Data Collection



Find the filtered and deduplicated datasets at: [www.hf.co/bigcode](https://hf.co/bigcode)

Am I In The Stack?

The Stack is an open governance interface between the AI community and the open source community.

Am I In The Stack?

As part of the BigCode project, we released and maintain **The Stack**, a 6.4 TB dataset of permissively licensed source code in 358 programming languages. One of our goals in this project is to give people agency over their source code by letting them decide whether or not it should be used to develop and evaluate machine learning models, as we acknowledge that not all developers may wish to have their data used for that purpose.

This tool lets you check if a repository under a given username is part of The Stack dataset. Would you like to have your data removed from future versions of The Stack? You can opt-out following the instructions [here](#).

The Stack version:

Your GitHub username:

<https://huggingface.co/spaces/bigcode/in-the-stack>

Opt-out

Yes, there is code from 5 repositories in The Stack:

`huggingface/datasets`
`huggingface/jupyterplot`
`huggingface/pandas-profiling`
`huggingface/transformers`
`huggingface/tf`

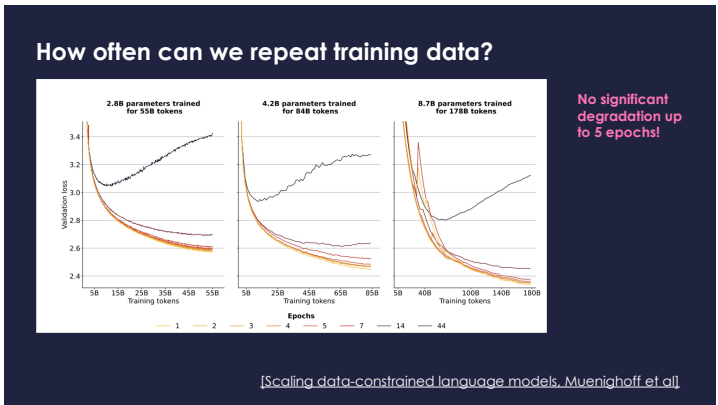
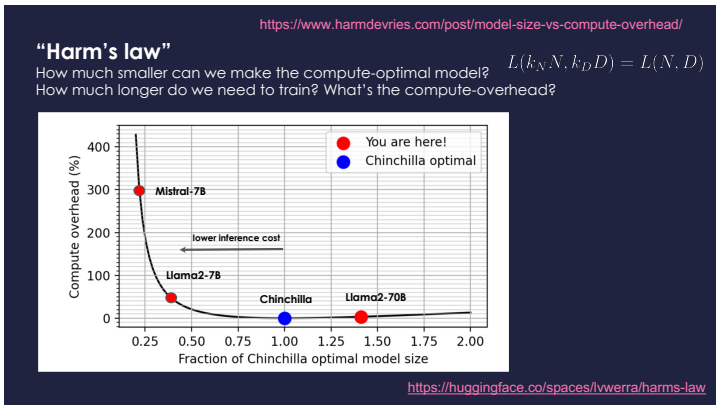
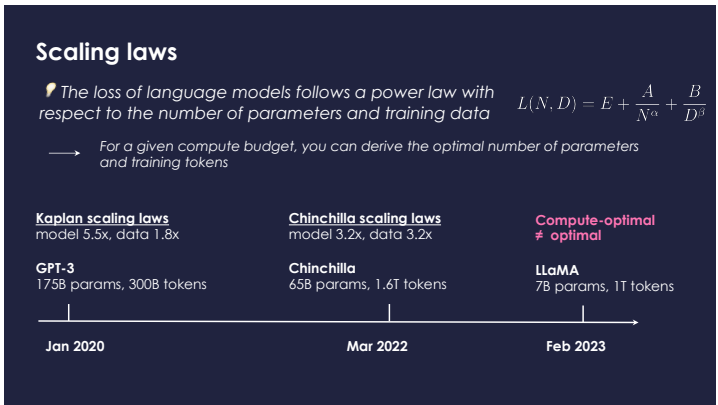
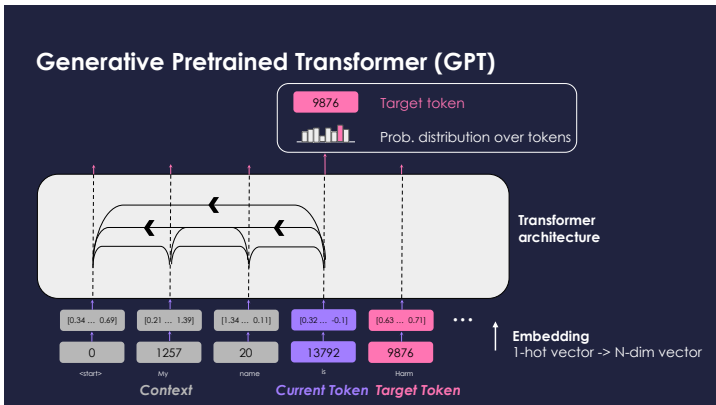
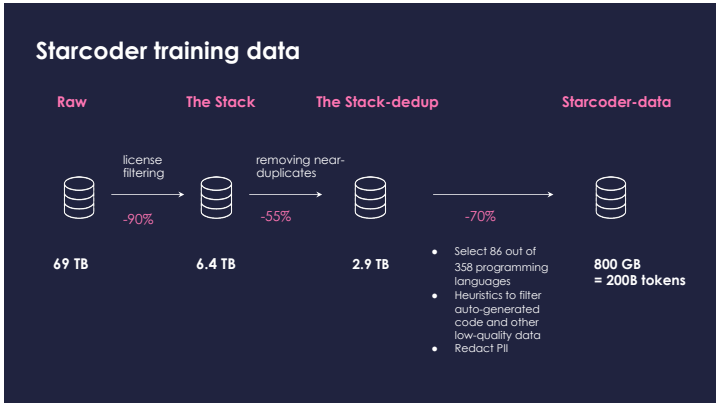
Opt-out

If you want your data to be removed from the stack and model training open an issue with [this link](#) (if the link doesn't work try right a right click and open it in a new tab) or visit <https://github.com/bigcode-project/opt-out-v2/issues/new?template=opt-out-request.md>.

Feedback from the opt-out form

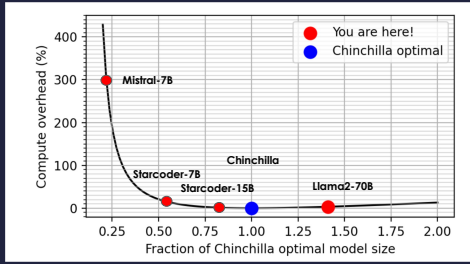
- "It should be opt-in instead of opt-out"
- "It is unfair to use my code without compensation"
- "There's PII in my code and I don't want it to be publicly exposed"
- "My code is of poor quality and unsuitable for training your AI model"
- "I am not confident about the current state of AI code generation. I am concerned that the generated code could be traced back to me and I'm held liable for issues in that code."

Jennifer Ding's community research: it's both **better to know** AND **better to have a choice**.



StarCoder model size

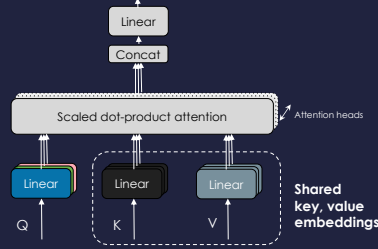
200B tokens x 5 epochs = 1 trillion training tokens



<https://huggingface.co/spaces/lvwerra/harms-law>

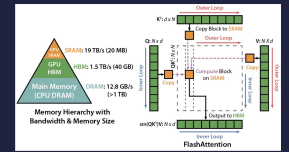
Model architecture

Multi-query attention



8K context

thanks to FlashAttention



Training setup

Infrastructure: 512 A100 GPUs

Model Distribution: TP=4, PP=4, DP=32

Batch size: 4M tokens

(or 512 at 8,192 sequence length)

Training length: 1T tokens / 250k steps

Training time: 24 days

Tool: Megatron-LM (w/ MQA + FlashAttn)

(<https://github.com/bigcode-project/Megatron-LM>)



"smooth sailing"

Python Evaluation: HumanEval

```
from typing import List

def filter_by_substring(strings: List[str],
    substring: str) -> List[str]:

    """ Filter an input list of strings only for
    ones that contain given substring >>>
    filter_by_substring(['a'], 'a') == ['a']
    filter_by_substring(['abc', 'baed', 'cde',
    'array'], 'a') == ['abc', 'baed', 'array'] """

    Model generates function body

    assert candidate(['xxx', 'asd', 'xxy', 'john
    doe', 'xxxxAAA', 'xxx'], 'xxx') == ['xxx',
    'xxxxAAA', 'xxx']
```

Model	Size	HumanEval pass@1
Open-access		
SantaCoder-1B	1B	18.1
DeciCoder-1B	1B	19.3
Replit-3B	3B	20.1
StableCode-3B	3B	20.2
StarCoderBase-3B	3B	21.5
StarCoderBase-7B	7B	28.4
CodeGen-Mono	16B	29.3
LLaMA-2	70B	29.9
CodeGen-2.5-Mono	7B	33.1
CodeGenX-2	6B	33.5
StarCoder-15B	15B	33.6
OctoCoder*	15B	45.3
WizardCoder*	15B	58.1
Closed-access		
LaMDA	137B	14.0
PaLM	540B	26.2
code-cushman-001	12B	33.5
PaLM 2-S	N/A	37.6
code-davinci-002	175B	45.9
GPT-3.5	N/A	48.1
PanGu-Coder 2*	15B	61.6
GPT-4*	N/A	67.0

Code llama-34B: 48.8%

VSCode extension

Auto-complete

```
Users > swayam > Desktop > main.py > ...
1 def is_prime(num):
2     return False

def is_prime(num):
    if num == 2:
        return True
    if num % 2 == 0:
        return False
    for i in range(3, num, 2):
        if num % i == 0:
            return False
```

<https://marketplace.visualstudio.com/items?itemName=HuggingFace.huggingface-vscode>

Membership test

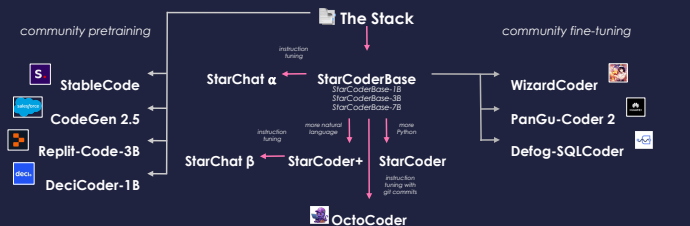
```
Users > swayam > Desktop > main.py > is_prime
1 def is_prime(num):
2     return False
3
4 def is_prime(num):
5     if num == 2:
6         return True
7     if num % 2 == 0:
8         return False
9     for i in range(3, num, 2):
10        if num % i == 0:
11            return False
```

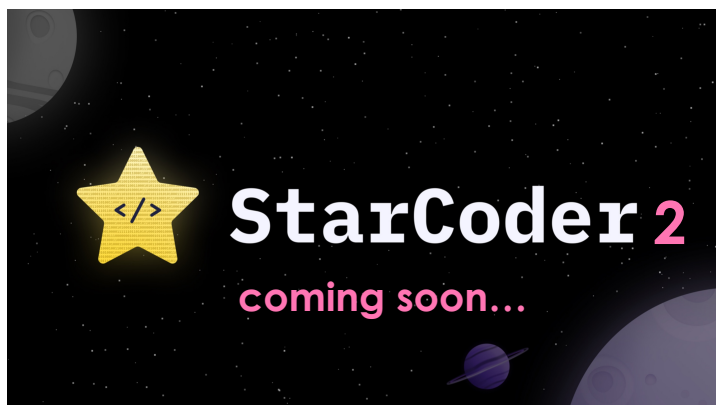
Highlighted code was found in the stack.

Source: HF Code Autocomplete (Extension)

Go to stack search

BigCode Ecosystem





Closing thoughts on open-source AI

- Frontier labs: OpenAI, Anthropic, Google, etc
- Allies of open-source AI: Meta, Mistral, Hugging Face, ServiceNow, etc
- Challenges for open-source AI to catch-up on frontier labs
 - depends on the support of large tech companies for:
 - compute
 - scarce AI talent
 - lack of user interaction data - openAI can use all the interactions with ChatGPT
- On the flipside, you might not need frontier models for the majority of use cases.
 - Many valuable enterprise use cases can be build on top of smaller open-source LLMs
- AI regulation
 - Copyright lawsuits
 - EU AI act
 - Biden's executive order
 - Regulation is the friend of the incumbent - Bill Gurley

Thank you! And come join us!



Questions?

www.bigcode-project.org

hf.co/bigcode




Building a truly open Dutch LLM

Matthieu Laneuville (SURF) & Saskia Lensink (TNO)

November 2, 2023

Open?

Project	Availability					Documentation					Access		
	Open code	LLM data	LLM weights	RLHF data	RLHF weights	License	Code	Architecture	Papers	Pages	Model card	Dataset	Feedback
Open Assistant	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Pythia-Chat-Base-7	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
RedPajama INCITE	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
dolly	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
OpenChat	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
MPT-7B Instruct	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
llm	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
MPT-30B Instruct	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Vicuna 13B v1.3	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
mistral-7B	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
ChatRWKV	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Cerebras GPT-11M	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
WizardLM 13B v1.2	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

https://opening-up-chatgpt.github.io/

Dutch?



https://medium.com/@socialcreature/ai-and-the-american-smile-76d23a0fbfaf

Dutch?



https://medium.com/@socialcreature/ai-and-the-american-smile-76d23a0fbfaf

Why a Dutch LLM?

A major obstacle to the use of AI in the Netherlands is that existing algorithms are not properly trained in the Dutch language. [...] This problem [...] applies to the entire public sector and to all Dutch-language interactions in the market. ... https://nlaiic.com/en/use-case/main-dutch-ai-for-dutch-language/

- Risks around bias, inclusion and explainability are not sufficiently guaranteed in current solutions
- Existing models of tech giants do not work sufficiently or cannot be used due to IP&C concerns
 - What if we want to use those models in a medical context? In an education context?
- Also, developing skills & capacity in the public sector helps with long term autonomy

ChatGPT Is Cutting Non-English Languages Out of the AI Revolution

AI chatbots are less fluent in languages other than English, threatening to amplify existing bias in global commerce and innovation.



Chinese organisations launched 79 AI large language models since 2020, report says

Große KI-Modelle FÜR DEUTSCHLAND

März/aprilsausgabe 2023

Why do we need a large GPT for Swedish?

What are the advantages of building a large language model for Swedish, and what should we look out for?

Magnus Sallgren - Polaris
Publicerat av Sallgren - 15 mars 2023

The GPT-NL ambition

- A modern Dutch model

The GPT-NL ambition

- Develop, strengthen and consolidate **digital sovereignty** through our own Dutch language model.
- Increase **strategic autonomy** in the underlying knowledge and technology.
- **Facilitate short-cycle, application-oriented research** into the use of language models in the broad context of social challenges.
- Enable **fundamental research into values such as bias and transparency** in large language models.



COMMENTARY | 19 October 2023

Living guidelines for generative AI – why scientists must oversee its use

Establish an independent scientific body to use and certify generative artificial intelligence, before the technology damages science and public trust.

Chloé J. Borking¹, Ewa A. M. van Dijk, Robert van Houdt, Willem Zuidema & Johan Bollen

[f](#) [t](#) [x](#) [in](#)

SURF **TNO**



Nederlandsche Wetenschappelijke Onderzoeksfaciliteiten
www.wetenschappelijkeonderzoek.nl

In practice...

Financing

Initial subsidy from Tender 'dream facilities' ministry of Economic Affairs, ~14 M€
Public-private partnerships and joint innovation centers to exploit the facility
Not-for-profit approach to cost structure (self-sustaining)

Facility

A Dutch foundation model, facilitating further research, fine-tuning and deployment
Trained and deployed on hardware hosted by SURF, new hardware will be purchased (88+16 H100)
A factor 1,000 times larger than the current Dutch language models
Both dataset and model open for inspection and research

Data

Existing open-source data sets (think of The Pile, Wikipedia, etc.) are a possibility
Text data extracted from high-quality sources such as government websites, forums, etc.
Public data (e.g., CLARIAH and CLARIN-NL, the Royal Library, etc.)
Data curation will be absolutely crucial, and choices will have to be made
The dataset used to train the base model will be made available under CC-BY-4.0

SURF **TNO**



Nederlandsche Wetenschappelijke Onderzoeksfaciliteiten
www.wetenschappelijkeonderzoek.nl

In practice...

Potential working groups (inspired by BigScience/BLOOM)

- Data
- Engineering / scaling – focus on the technical challenges of training (scaling, specialized hardware)
- Evaluation
- Legal & ethical issues
- Modelling
- Model & data sharing – prepare ethical framework and tools for sharing
- Governance / workshop organization – main organizing / coordination group
- Carbon footprint / sustainability – what is the carbon footprint from the project and how to reduce it
- Challenges for application X – series of working groups focused on specific applications

	jaar 1					jaar 2		
	jan	apr	jul	okt	dec	jan	apr	jul
Dataset opbouw								
Inrichting ontwikkel cluster								
Algoritme ontwikkeling								
Training van taalmodel								
Inrichting operationeel cluster								
Platform ontwikkeling								
Inrichting mobiel bijtrain cluster								

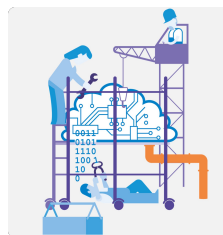
Figuur 1: Activiteiten en milestones door de tijd

SURF **TNO**



Nederlandsche Wetenschappelijke Onderzoeksfaciliteiten
www.wetenschappelijkeonderzoek.nl

Opening remarks



- TNO/SURF led community owned
- How do we facilitate responsible experimentation together?
- How do we embed this in the larger international context?
- What should be our shared standard for openness?



Home > Actueel > Nieuws >

14 Nederlandse onderzoeksfaciliteiten krijgen 184 miljoen euro ondersteuning

Nieuwsbericht | 02-11-2023 | 12:36

SURF **TNO**



Nederlandsche Wetenschappelijke Onderzoeksfaciliteiten
www.wetenschappelijkeonderzoek.nl

Jelle Zuidema discusses and shares with **Konstantinos Papakostas** (ZetaAlpha), **Iva Gornishka** (City of Amsterdam), **Jelke Bloem** (ILLC, University of Amsterdam) and **Oskar van der Wal** (ILLC, University of Amsterdam) their perspectives on how to evaluate such models and make sure they are useful, safe and reliable



Thank you for your interest in
the **InDeep Project**

For more information about
the project, members,
research and our upcoming
events, visit:
projects.llc.uva.nl/indeep/



nationale
wetenschappen
agenda



**INTERPRETING
DEEP LEARNING MODELS
FOR TEXT AND SOUND**

Radboud University



UNIVERSITY
OF AMSTERDAM

TILBURG



UNIVERSITY



VRIJE
UNIVERSITEIT
AMSTERDAM



university of
 groningen

OUR PARTNERS:

**KPN - IDFUSE - TNO - TEXTKERNEL - CHORDIFY - GLOBAL TEXTWARE
DELOITTE - FLOODTAGS**