



Istituto di Linguistica  
Computazionale  
"Antonio Zampolli"



Consiglio Nazionale delle Ricerche



# Evaluating Linguistic Abilities of Neural Language Models

*NLP4RE'25 @ REFSQ 2025, April 7 2025*

Alessio Miaschi

ItaliaNLP Lab, Istituto di Linguistica Computazionale (CNR-ILC), Pisa

[alessio.miaschi@ilc.cnr.it](mailto:alessio.miaschi@ilc.cnr.it)

<https://alemiaschi.github.io/>

<http://www.italianlp.it/alessio-miaschi/>

# About me and...



I am a full-time researcher (RTDA) at the [ItaliaNLP Lab](#), Institute for Computational Linguistics “A. Zampolli” ([CNR-ILC](#), Pisa). In 2022, I received my PhD in Computer Science at the University of Pisa.

My research interests lie primarily in the context of Natural Language Processing (NLP) and in the study of Language Models (LM). I am particularly interested in the interpretability of large-scale LMs and in the evaluation of their internal representations, with a specific emphasis on understanding their inner linguistic abilities.

# About me and... the team!



I am a full-time researcher (RTDA) at the [ItaliaNLP Lab](http://www.italianlp.it), Institute for Computational Linguistics “A. Zampolli” (CNR-ILC, Pisa). In 2022, I received my PhD in Computer Science at the University of Pisa.

My research interests lie primarily in the context of Natural Language Processing (NLP) and in the study of Language Models (LM). I am particularly interested in the interpretability of large-scale LMs and in the evaluation of their internal representations, with a specific emphasis on understanding their inner linguistic abilities.



The **ItaliaNLP Lab (CNR-ILC)** gathers researchers, postdocs and students from computational linguistics, computer science and linguistics who work on developing resources and algorithms for processing and understanding human languages.

## Permanent Researchers:

- Felice Dell’Orletta
- Simonetta Montemagni
- Dominique Brunato
- Franco Alberto Cardillo
- Giulia Venturi
- Giulia Benotto

## Temporary Researchers:

- Chiara Alzetta
- Alessio Miaschi

## Research Fellows:

- Agnese Bonfigli
- Cristiano Ciaccio
- Chiara Fazzone
- Ruben Piperno
- Marta Sartor

## PhD Students:

- Luca Dini
- Lucia Domenichelli
- Michele Papucci

+ **Master/Undergraduate/Visiting Students**

Link to website: <http://www.italianlp.it/>

# Outline

1. Introduction
2. Interpreting and Evaluating NLMs
3. Conclusion and Future Directions

---

# Introduction

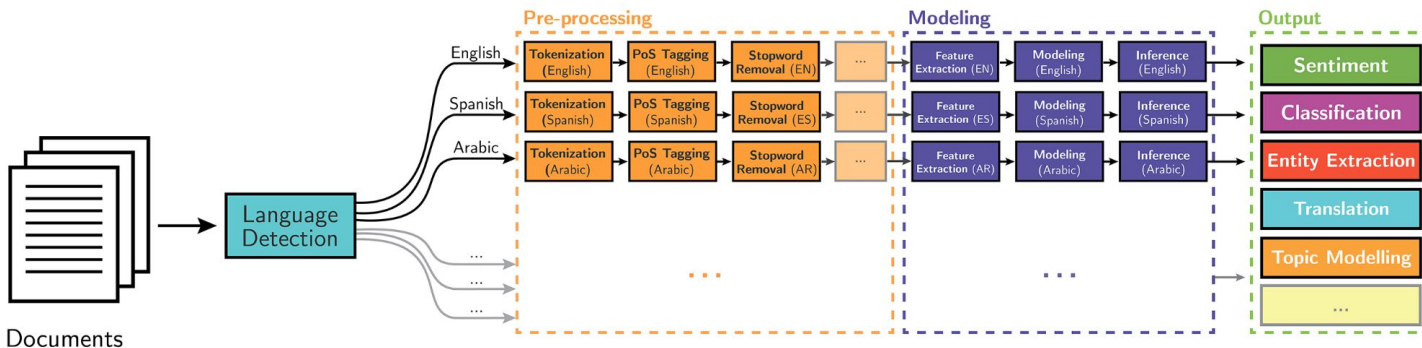
# Introduction

- The field of NLP has seen an unprecedented progress in the last years
- Much of this progress is due to the replacement of traditional systems with newer and more powerful Deep Learning (DL) models

# Introduction

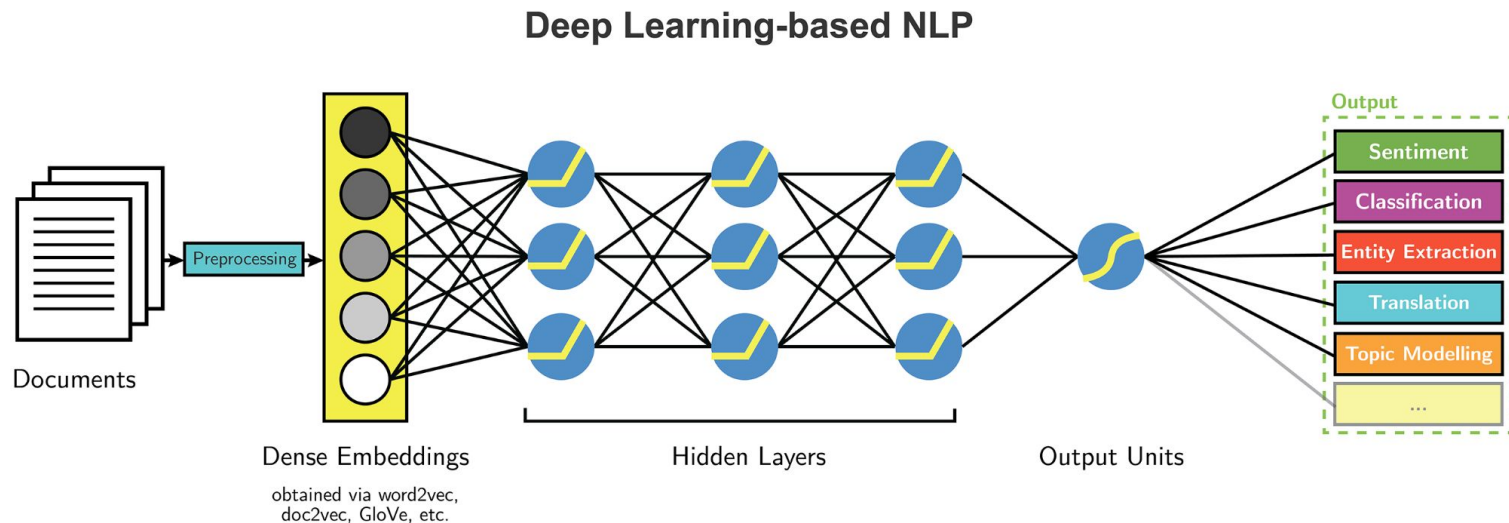
- The field of NLP has seen an unprecedented progress in the last years
- Much of this progress is due to the replacement of traditional systems with newer and more powerful Deep Learning (DL) models

## Classical NLP



# Introduction

- The field of NLP has seen an unprecedented progress in the last years
- Much of this progress is due to the replacement of traditional systems with newer and more powerful Deep Learning (DL) models





# Neural Language Model (NLM)

- Neural Language Model → Neural Network trained to approximate the **language modeling** function

# Neural Language Model (NLM)

- Neural Language Model → Neural Network trained to approximate the **language modeling** function
- **Language Modeling** → probability of a sentence  $s = [w_1, w_2, \dots, w_n]$  as:

$$P(s) = \prod_{i=1}^N P(w_i | w_1, w_2, \dots, w_{i-1})$$

# Neural Language Model (NLM)

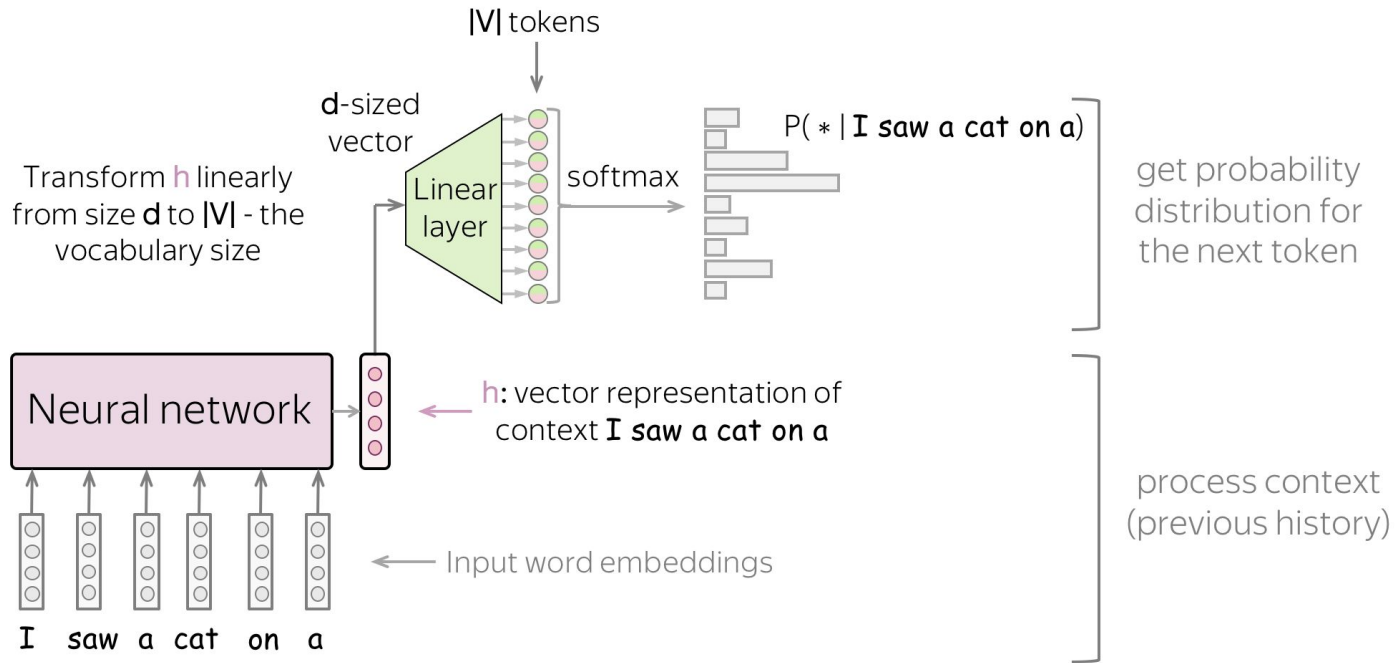
- Neural Language Model → Neural Network trained to approximate the **language modeling** function

- **Language Modeling** → probability of a sentence  $s = [w_1, w_2, \dots, w_n]$  as:

$$P(s) = \prod_{i=1}^N P(w_i | w_1, w_2, \dots, w_{i-1})$$

- **Bengio et al. (2003)** proposed a model to learn this function relying on the architecture of a neural network → **Neural Probabilistic Language Model**

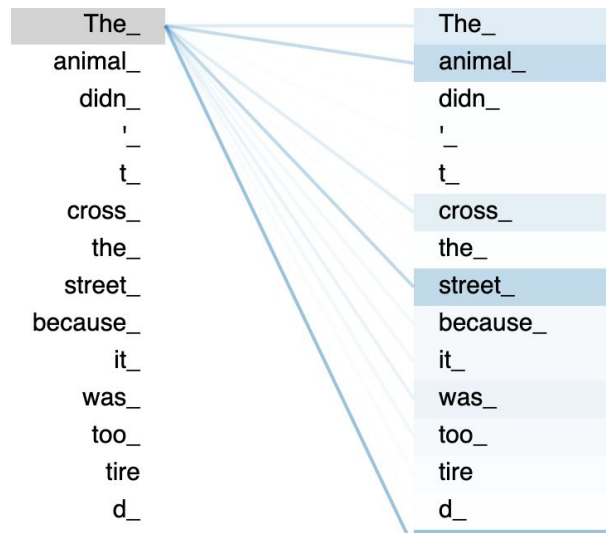
# Neural Language Model (NLM)



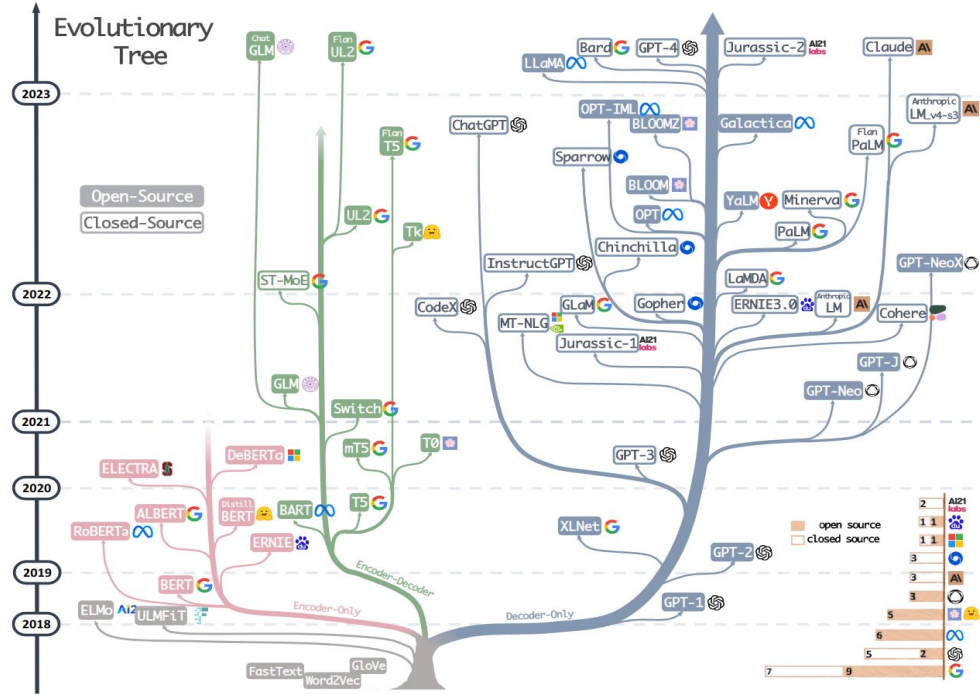
# Transformer Models

- Nowadays, the Transformer is the most commonly used architecture for the development of NLMs
- The Transformer (Vaswani et al., 2017) exploits the **attention mechanism** to create contextual representations of words and learn the relations among them

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



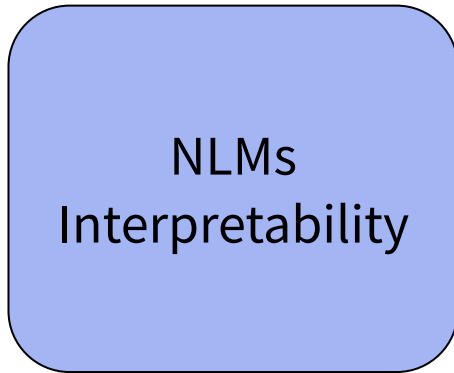
# “Evolutionary Tree”



# Interpreting and Evaluating NLMs

# Interpreting and Evaluating NLMs

- The rapid development and widespread adoption of state-of-the-art Neural Language Models (NLMs) have increased the need for studies focused on their **interpretability** and the **evaluation** of their abilities





# Interpreting and Evaluating NLMs

- The rapid development and widespread adoption of state-of-the-art Neural Language Models (NLMs) have increased the need for studies focused on their **interpretability** and the **evaluation** of their abilities



# The Case for Interpretability

- The development of powerful state-of-the-art NLMs comes at the cost of **interpretability**, since complex NN models offer little transparency about their inner workings and their abilities

## Objectives:

- **Understand the nature of AI systems** → be faithful to what influences the AI decisional process
- **Empower AI system users** → derive actionable useful insights from AI choices

# Interpretability in NLP

*“In the context of NLP, this question needs to be understood in light of earlier NLP work. [...] In some of these systems, features are more easily understood by humans. [...] In contrast, it is more difficult to understand what happens in an end-to-end neural network model that takes input (say, word embeddings) and generates an output.”*

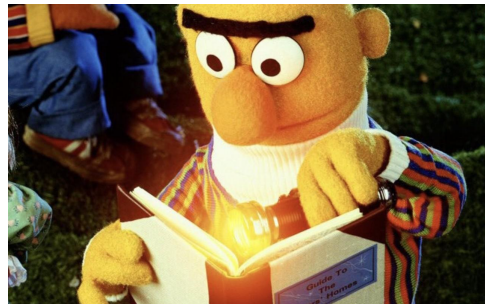
Belinkov and Glass, Analysis Methods in Neural Language Processing: A Survey (2019). In Transactions of ACL, Volume 7, pages 49-72.



# Interpretability in NLP

*“In the context of NLP, this question needs to be understood in light of earlier NLP work. [...] In some of these systems, features are more easily understood by humans. [...] In contrast, it is more difficult to understand what happens in an end-to-end neural network model that takes input (say, word embeddings) and generates an output.”*

Belinkov and Glass, Analysis Methods in Neural Language Processing: A Survey (2019). In Transactions of ACL, Volume 7, pages 49-72.

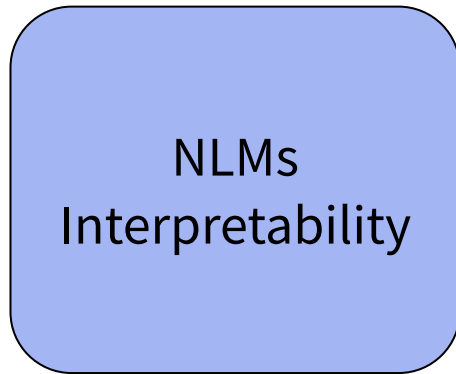


## Research questions:

- What happens in an end-to-end neural network model when trained on a language modeling task?
- What kind of linguistic knowledge (i.e. features) is encoded within their representations?
- Is there a relationship between the linguistic knowledge implicitly encoded and the ability to solve a specific task?

# Interpreting and Evaluating NLMs

- The rapid development and widespread adoption of state-of-the-art Neural Language Models (NLMs) have increased the need for studies focused on their **interpretability** and the **evaluation** of their abilities



# Evaluation of Neural Language Models

- The evaluation of NLMs has seen significant advancements in the past few years, with the development of dedicated benchmarks and evaluation frameworks
- These benchmarks are designed to assess models' performance on specific tasks and reasoning abilities:
  - OpenLLM Leaderboard
  - BigBench (Srivastava et al., 2023)
  - Holmes (Waldis et al., 2024)

The screenshot displays the Open LLM Leaderboard interface. At the top, it says "Open LLM Leaderboard" and provides a link to the previous version. Below this, there are navigation buttons for "LLM Benchmark", "Submit", and "Model Vote". A search bar is present, along with a "Select Columns to Display" section where various metrics like Average, IFEval, BBH, MATH Lvl 5, and GPQA are selected. The main part of the image is a table listing models and their performance across these metrics.

T	Model	Average	IFEval	BBH	MATH Lvl 5	GPQA	MUSR	MMLU-PRO	CO <sub>2</sub> cost (kg)
#	<a href="#">dizmen/calme8ya-788-01po-v0.1</a>	51.24	81.63	61.92	48.71	20.02	36.37	66.8	13
#	<a href="#">MaziyasPanahi/calme-2.4-sys-78b</a>	50.71	80.11	62.16	49.41	20.36	34.57	66.69	12.98
◆	<a href="#">rombodeng/Rombos-LLM-V2.5-Qwen-72b</a>	45.91	71.55	61.27	50.68	19.8	17.32	54.83	16.03
◆	<a href="#">zetasepic/Qwen2.5-72B-Instruct-abiliteated</a>	45.29	71.53	59.91	46.15	20.92	19.12	54.13	18.81
◆	<a href="#">dinhong/RYS-VLasege</a>	45.13	79.96	58.77	41.24	17.9	23.72	49.2	13.58
◆	<a href="#">rombodeng/Rombos-LLM-V2.5-Qwen-32b</a>	44.57	68.27	58.26	41.99	19.57	24.73	54.62	17.91
#	<a href="#">MaziyasPanahi/calme-2.1-sys-78b</a>	44.56	81.36	59.47	38.9	19.24	19	49.38	14.33
#	<a href="#">MaziyasPanahi/calme-2.3-sys-78b</a>	44.42	80.66	59.57	38.97	20.58	17	49.73	13.3
#	<a href="#">MaziyasPanahi/calme-2.2-sys-78b</a>	44.26	79.86	59.27	39.95	20.92	16.83	48.73	13.52

Link: [https://huggingface.co/spaces/open-llm-leaderboard/open\\_llm\\_leaderboard](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard)

# Competence vs. Performance in NLMs

- Within the broader context of interpretability and evaluation, one line of research focuses on studying and assessing the linguistic abilities of (Large) Language Models
- Such studies aim to uncover the implicit linguistic competence encoded within these models and evaluate their generalization abilities
- **Competence vs. Performance:** investigation of the linguistic abilities of NLMs from a competence/performance perspective:
  - Distinction between the information encoded in a model internal representation vs. the model's behavioral responses to prompt during generation (Hu and Levy, 2023)

# Profiling Neural Language Models

- The “*linguistic profiling*” methodology (van Halteren, 2004) assumes that wide counts of linguistic features are particularly helpful in the resolution of several NLP tasks, e.g.:
  - Text Profiling (e.g. text readability, textual genres)
  - Author Profiling (e.g. author’s age and native language)



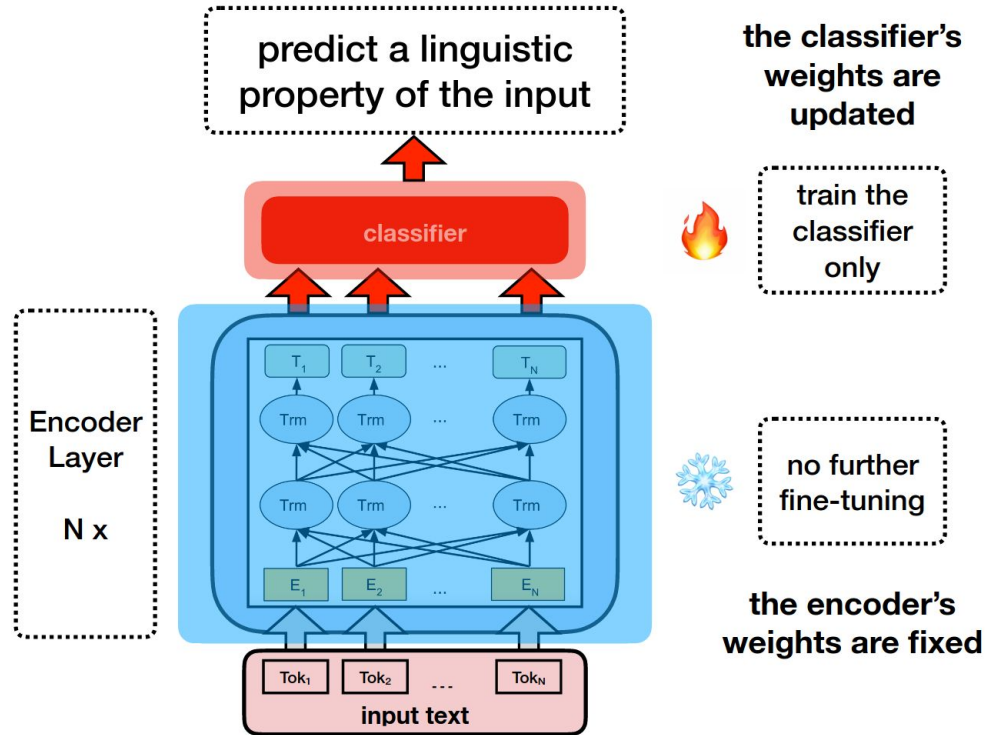
# Profiling Neural Language Models

- The “*linguistic profiling*” methodology (van Halteren, 2004) assumes that wide counts of linguistic features are particularly helpful in the resolution of several NLP tasks, e.g.:
  - Text Profiling (e.g. text readability, textual genres)
  - Author Profiling (e.g. author’s age and native language)

## Research Question:

Could the informative power of these features also be helpful to understand the behaviour of state-of-the-art NLMs?

# Probing Task Approach

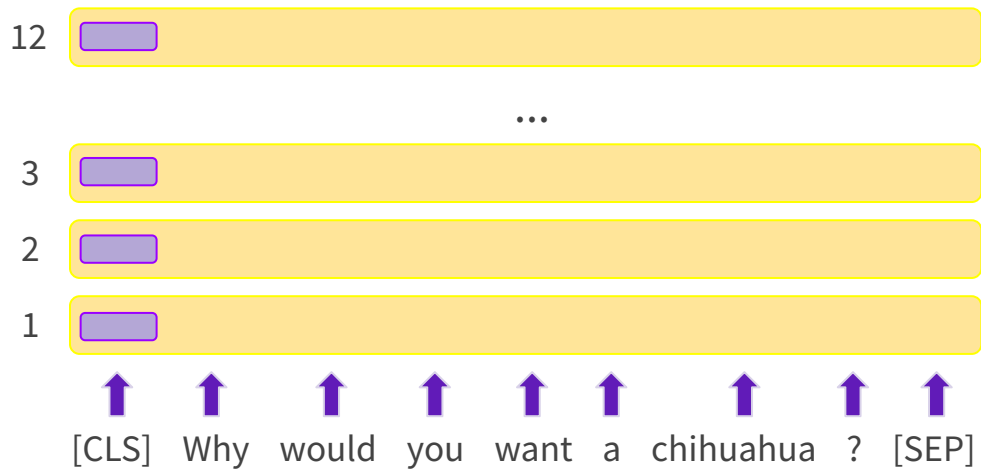


# Profiling-UD: a tool for Linguistic Profiling of Texts

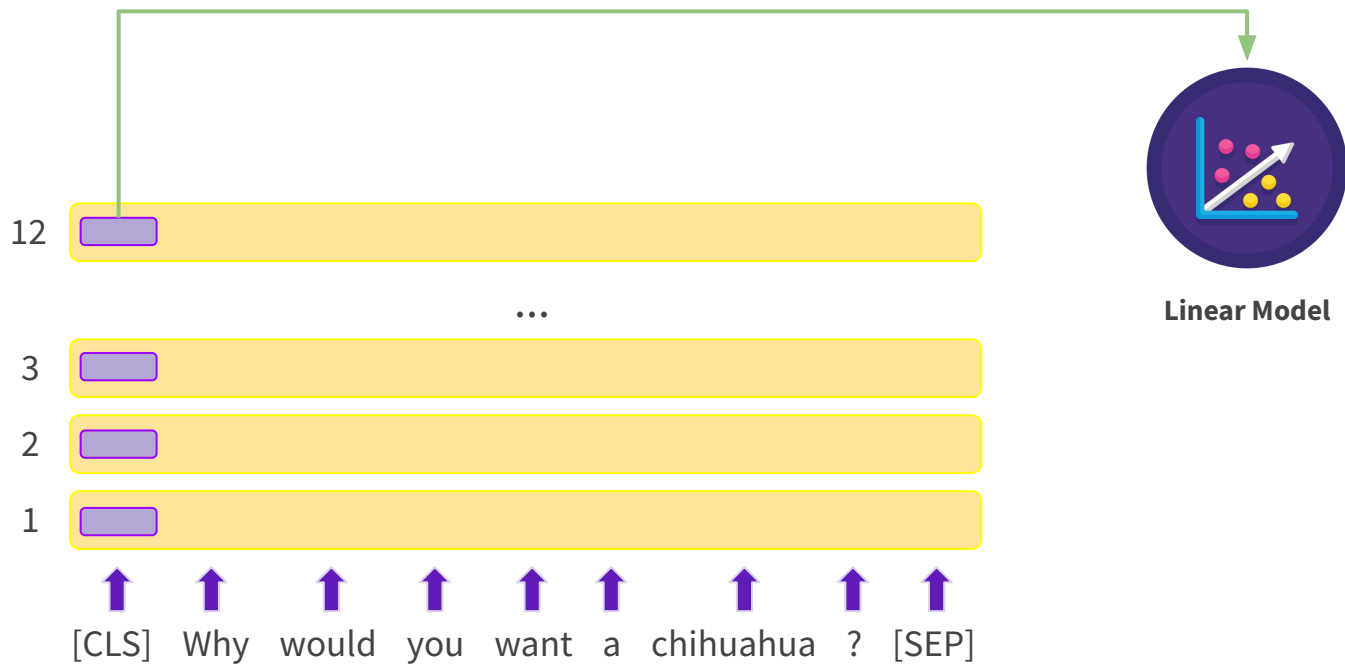
- ProfilingUD (Brunato et al., 2020) is a web-based application that performs linguistic profiling of a text, or a large collection of texts, for multiple languages
- It allows the extraction of more than 130 features, spanning across different levels of linguistic description
- Link: <http://linguistic-profiling.italianlp.it/>

<b>Linguistic Feature</b>
<b>Raw Text Properties</b>
Sentence Length
Word Length
<b>Vocabulary Richness</b>
Type/Token Ratio for words and lemmas
<b>Morphosyntactic information</b>
Distribution of UD and language-specific POS
Lexical density
<b>Inflectional morphology</b>
Inflectional morphology of lexical verbs and auxiliaries
<b>Verbal Predicate Structure</b>
Distribution of verbal heads and verbal roots
Verb arity and distribution of verbs by arity
<b>Global and Local Parsed Tree Structures</b>
Depth of the whole syntactic tree
Average length of dependency links and of the longest link
Average length of prepositional chains and distribution by depth
Clause length
<b>Relative order of elements</b>
Order of subject and object
<b>Syntactic Relations</b>
Distribution of dependency relations
<b>Use of Subordination</b>
Distribution of subordinate and principal clauses
Average length of subordination chains and distribution by depth
Relative order of subordinate clauses

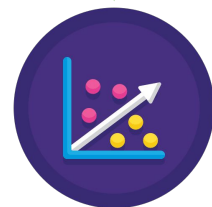
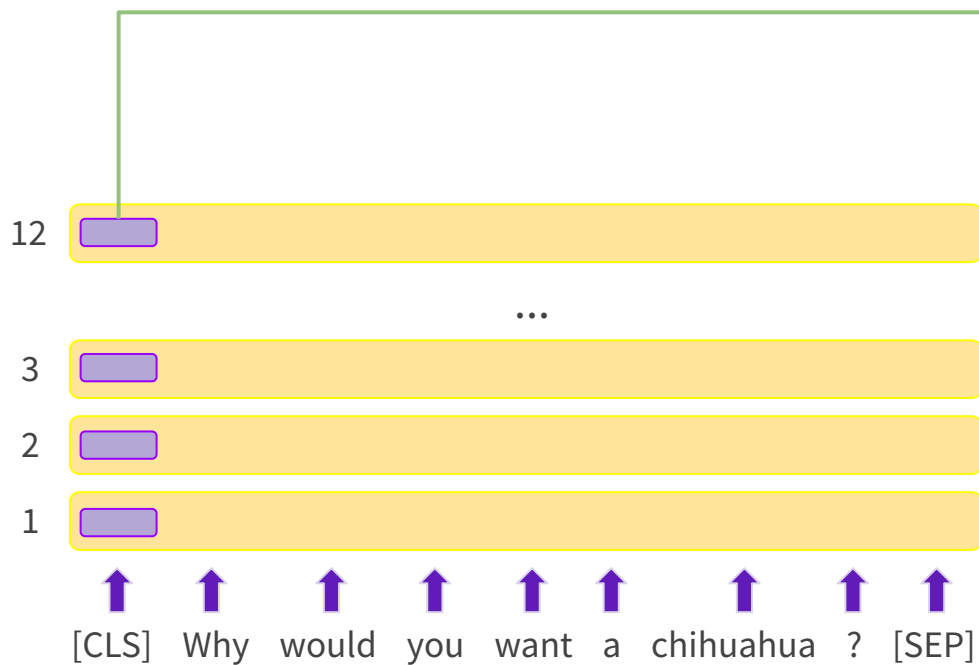
# Profiling Neural Language Models



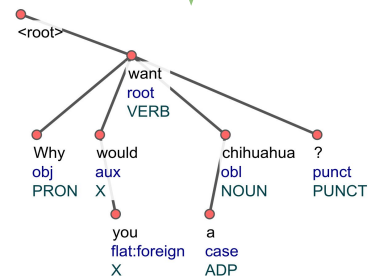
# Profiling Neural Language Models



# Profiling Neural Language Models



Linear Model



# Linguistic Profiling of a Neural Language Model (Miaschi et al., 2020)

- We investigated the linguistic knowledge implicitly encoded by BERT

## Research questions:

1. What kind of linguistic properties are encoded in a pre-trained version of BERT?
2. How this knowledge is modified after a fine-tuning process?
  - a. Fine-tuning on the Natural Language Identification Task





# Linguistic Knowledge Can Enhance Encoder-Decoder Models

- Motivations:
  - Understanding “how linguistic concepts that were common as features in NLP systems are captured in neural networks” (Belinkov & Glass, *Transactions of the Association for Computational Linguistics 2019*) has been the focus of many recent studies
  - Fine-tuning on a intermediate supporting task and then on the target task consecutively is highly beneficial to improve pre-trained model’s performance (Weller et al., *ACL 2022*)

# Linguistic Knowledge Can Enhance Encoder-Decoder Models

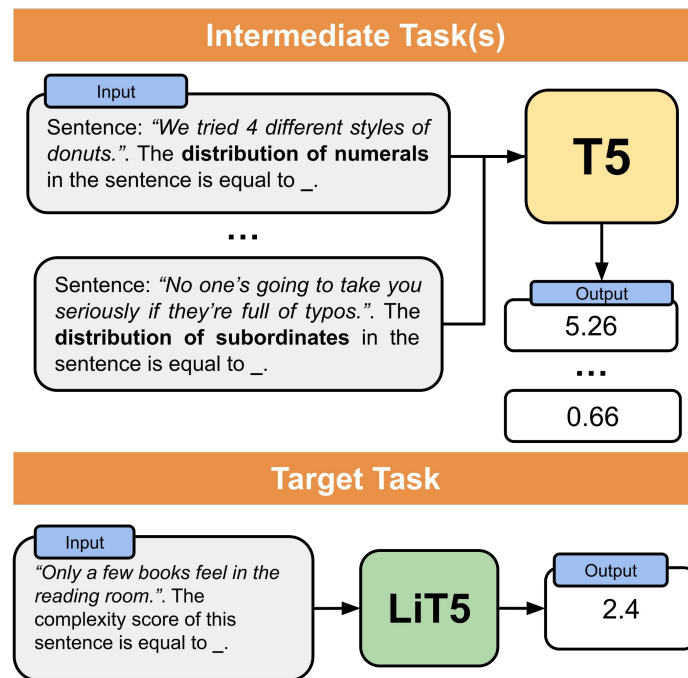
- Motivations:
  - Understanding “how linguistic concepts that were common as features in NLP systems are captured in neural networks” (Belinkov & Glass, *Transactions of the Association for Computational Linguistics 2019*) has been the focus of many recent studies
  - Fine-tuning on a intermediate supporting task and then on the target task consecutively is highly beneficial to improve pre-trained model’s performance (Weller et al., *ACL 2022*)



Does a step of intermediate fine-tuning on linguistic tasks enhance the prediction on a target task that strongly relies on linguistic knowledge?

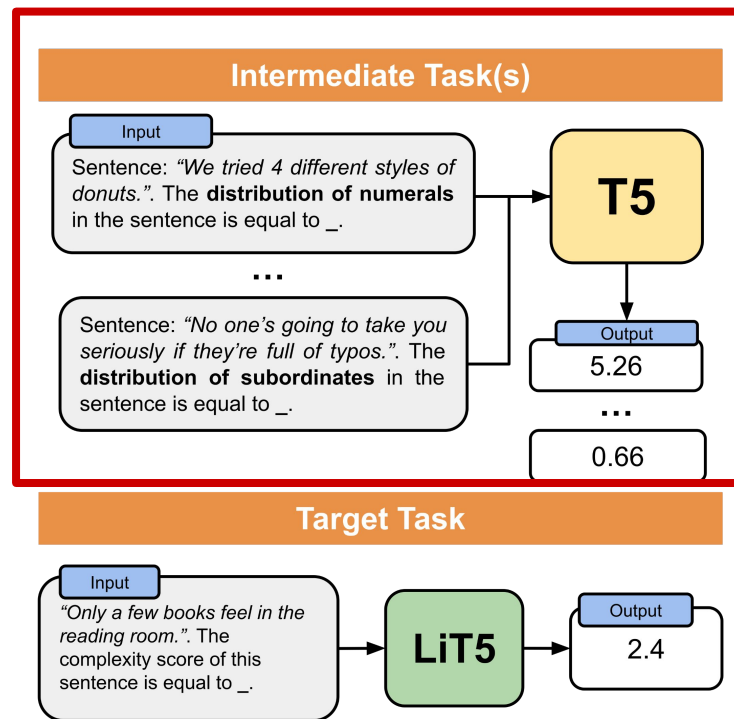
# Our Approach

- Two-step approach:
  - Fine-tune the T5 models on several intermediate tasks
    - Multi- and single-task fine-tuning
  - Fine-tune the Linguistically-Informed (LI) models on the target task
- We saved checkpoints every 5 epochs, in order to monitor the impact of the approach at increasing snapshots of the models
- We tested the approach both in Italian and English and in a cross-lingual scenario



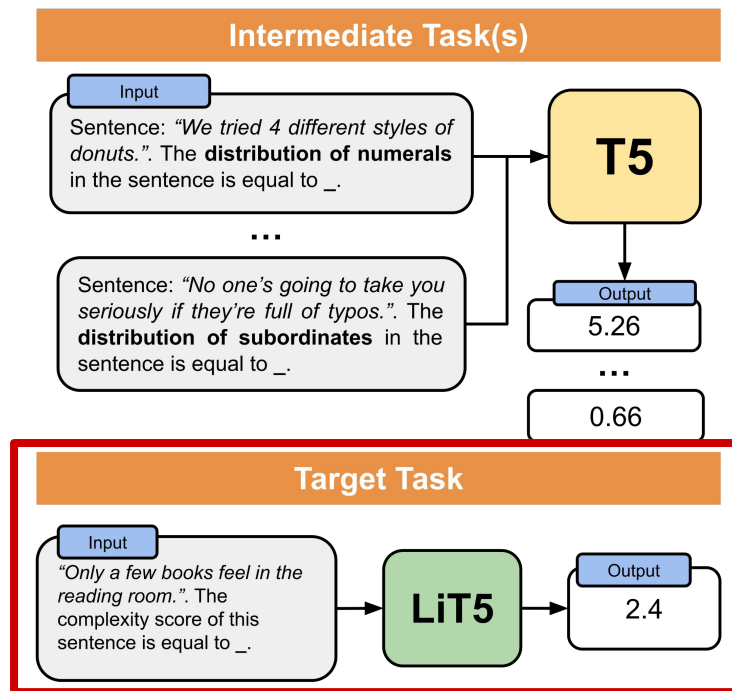
# Our Approach

- Two-step approach:
  - **Fine-tune the T5 models on several intermediate tasks**
    - **Multi- and single-task fine-tuning**
  - Fine-tune the Linguistically-Informed (LI) models on the target task
- We saved checkpoints every 5 epochs, in order to monitor the impact of the approach at increasing snapshots of the models
- We tested the approach both in Italian and English and in a cross-lingual scenario



# Our Approach

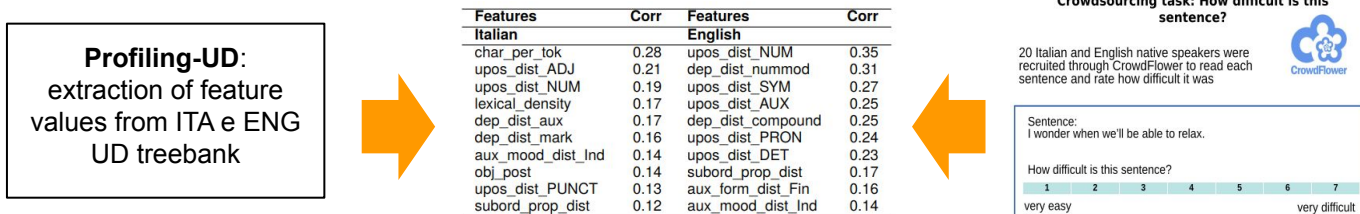
- Two-step approach:
  - Fine-tune the T5 models on several intermediate tasks
    - Multi- and single-task fine-tuning
  - **Fine-tune the Linguistically-Informed (LI) models on the target task**
- We saved checkpoints every 5 epochs, in order to monitor the impact of the approach at increasing snapshots of the models
- We tested the approach both in Italian and English and in a cross-lingual scenario



# Data

- Intermediate tasks:

- 10 morpho- and syntactic characteristics of a sentence
  - selected on the degree of correlation between sentence-level complexity judgments and their values



- Prediction of their distribution in the Italian and English versions of the Universal Dependency Treebanks extracted with Profiling-UD

- Target task:

- corpus of 1,440 Italian and 2,400 English sentences manually rated by 20 crowdsourced workers for the level of perceived complexity on 1-7 Likert scale (Brunato et al., EMNLP 2018)

# Models and Evaluation

## Models:

Language	Model	Parameters
English	t5-small	60M
	t5-base	220M
	t5-large	770M
Italian	it5-small	60M
	it5-base	220M
	it5-large	738M

## Evaluation:

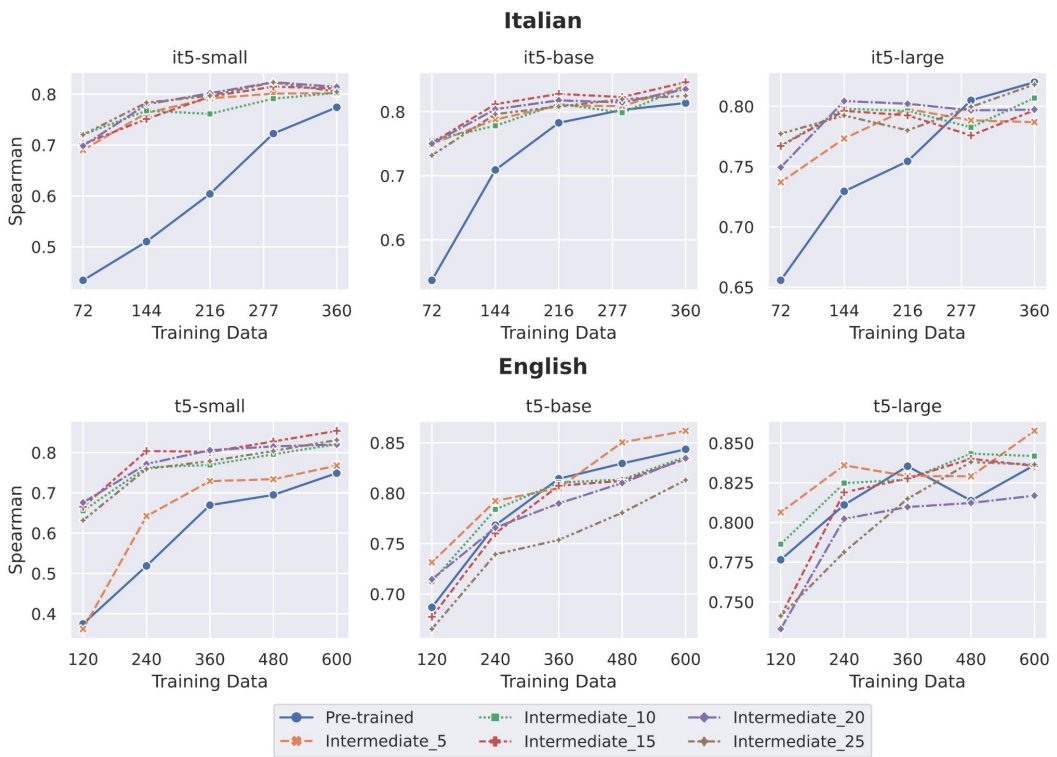
- We used Spearman correlation score as evaluation metric:
  - **Intermediate tasks:** Correlation between the gold value of each feature in the Italian or English treebank and the predicted value of the models for the intermediate tasks.
  - **Target task:** Correlation between average judgments of complexity and the complexity scores obtained with the fine-tuned LiT5 models.

# Enhancing T5 with Linguistic Features

		Italian														
		it5-small					it5-base					it5-large				
		5	10	15	20	25	5	10	15	20	25	5	10	15	20	25
All		0.41	0.49	0.53	0.55	0.56	0.53	0.64	0.73	0.76	0.77	0.6	0.72	0.75	0.81	0.83
aux_mood_dist_Ind		0.17	0.31	0.34	0.38	0.4	0.36	0.73	0.81	0.86	0.87	0.59	0.81	0.87	0.89	0.9
char_per_tok		0.0056	-0.046	0.06	0.061	0.13	0.15	0.28	0.36	0.48	0.53	0.15	0.31	0.42	0.6	0.63
dep_dist_aux		0	0	0	0.14	0.17	0	0.12	0.68	0.81	0.85	0.074	0.59	0.71	0.81	0.8
dep_dist_mark		0	0	0.091	0.21	0.23	0	0.38	0.59	0.65	0.74	0.021	0.44	0.76	0.77	0.82
lexical_density		0.0054	0.14	0.15	0.2	0.17	0.21	0.22	0.22	0.25	0.29	0.18	0.18	0.17	0.2	0.19
obj_post		0.18	0.31	0.38	0.41	0.41	0.35	0.38	0.42	0.46	0.5	0.46	0.54	0.59	0.68	0.69
subord_prop_dist		0.51	0.52	0.58	0.63	0.64	0.63	0.68	0.77	0.8	0.79	0.59	0.7	0.71	0.75	0.77
upos_dist_ADJ		0.14	0.18	0.22	0.18	0.22	0.26	0.39	0.44	0.44	0.45	0.24	0.29	0.39	0.53	0.58
upos_dist_NUM		0	0	0	0	0	0	0.34	0.93	0.94	0.94	-0.024	0.91	0.9	0.92	0.92
upos_dist_PUNCT		-0.15	0.13	0.22	0.21	0.25	0.17	0.3	0.41	0.51	0.54	0.2	0.24	0.38	0.61	0.76
		5	10	15	20	25	5	10	15	20	25	5	10	15	20	25
		English														
		t5-small					t5-base					t5-large				
		5	10	15	20	25	5	10	15	20	25	5	10	15	20	25
All		0.45	0.51	0.66	0.79	0.87	0.54	0.78	0.88	0.89	0.9	0.89	0.92	0.93	0.93	0.93
aux_form_dist_Fin		0.55	0.66	0.76	0.84	0.85	0.69	0.74	0.9	0.91	0.94	0.9	0.92	0.94	0.95	0.95
aux_mood_dist_Ind		0.46	0.63	0.79	0.86	0.89	0.72	0.72	0.86	0.9	0.9	0.92	0.93	0.93	0.95	0.94
dep_dist_compound		0	0	0.14	0.35	0.52	0	0.16	0.57	0.57	0.61	0.53	0.62	0.64	0.63	0.68
dep_dist_nummod		0	0	0	0.5	0.7	0	0.65	0.8	0.8	0.81	0.73	0.74	0.83	0.8	0.81
subord_prop_dist		0.67	0.72	0.75	0.81	0.85	0.64	0.78	0.87	0.87	0.85	0.86	0.9	0.89	0.89	0.88
upos_dist_AUX		0	0	0.57	0.84	0.89	0.17	0.77	0.9	0.93	0.94	0.9	0.96	0.94	0.97	0.96
upos_dist_DET		0	-0.011	0.33	0.62	0.81	0.14	0.74	0.84	0.84	0.88	0.75	0.87	0.92	0.89	0.93
upos_dist_NUM		0	0	0.19	0.76	0.9	0.23	0.85	0.92	0.91	0.91	0.89	0.92	0.93	0.94	0.94
upos_dist_PRON		0	0.11	0.53	0.66	0.83	0.26	0.84	0.9	0.92	0.92	0.89	0.93	0.95	0.95	0.94
upos_dist_SYM		0	0	0	0	0.53	0	0.27	0.37	0.38	0.65	0.27	0.71	0.8	0.75	0.75
		5	10	15	20	25	5	10	15	20	25	5	10	15	20	25



# Predicting Complexity with LI Models



# Selected Findings

- Informing models linguistically over several epochs allows them to progressively improve their degree of language proficiency.
- The method of linguistic enhancement is particularly effective, especially when applied to smaller models and in scenarios with limited availability of target training data.
- Small models, refined through intermediate fine-tuning, can frequently surpass the performance of larger models that have not undergone this intermediate refinement process.

# Evaluating Large Language Models via Linguistic Profiling

- Motivations:
  - Large Language Models (LLMs) demonstrated remarkable capabilities in solving multiple tasks and in generating coherent and contextually relevant texts
  - Such capabilities have been extensively evaluated against several benchmarks, as evidenced by the success of platforms such as the OpenLLM Leaderboard
  - A comprehensive evaluation of **LLMs' linguistic abilities in generation**, independent of specific tasks and possibly cross-cutting across them, is still missing

# Evaluating Large Language Models via Linguistic Profiling

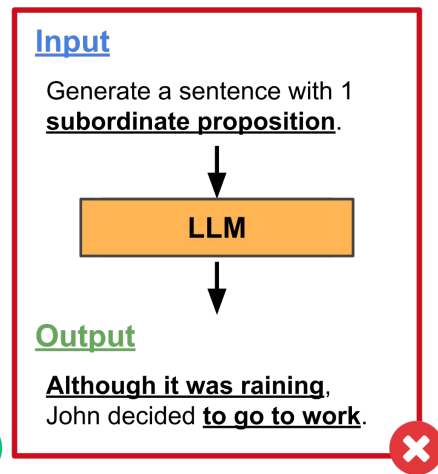
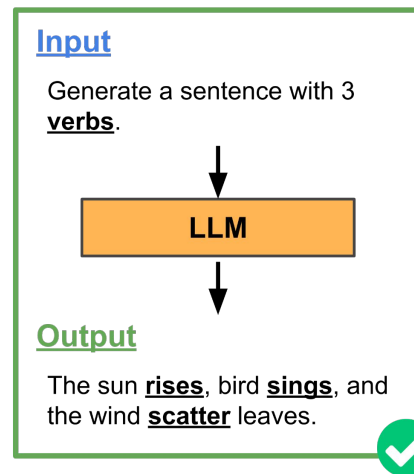
- Motivations:
  - Large Language Models (LLMs) demonstrated remarkable capabilities in solving multiple tasks and in generating coherent and contextually relevant texts
  - Such capabilities have been extensively evaluated against several benchmarks, as evidenced by the success of platforms such as the OpenLLM Leaderboard
  - A comprehensive evaluation of **LLMs' linguistic abilities in generation**, independent of specific tasks and possibly cross-cutting across them, is still missing



How effectively can LLMs generate sentences that adhere to targeted linguistic constraints representing various morpho-syntactic and syntactic phenomena?

# Our Approach

- We evaluate the ability of several LLMs to generate sentences with targeted (morpho-)syntactic linguistic constraints
- We prompted the models to generate sentences containing these constraints within a fixed prompt structure:
  - For each property/constraint, we asked the models to generate a fixed number of sentences having a precise value of that property
- Given the well-known difficulty of LLMs in producing texts with precise numerical constraints, we decided to constrain the models on increasing values of linguistic properties



# Linguistic Properties and Values Selection

- We relied on a set of linguistic properties as constraints encompassing diverse morpho-syntactic and syntactic phenomena of a sentence
- We relied on the largest English Universal Dependency (UD) treebank, i.e. English Universal Dependency (EWT) (Silveira et al., 2014)
  - Extraction of the linguistic properties with the Profiling-UD tool (Brunato et al., 2020)
  - In the few-shot configuration, we used 5 exemplar sentences extracted from EWT
- We asked each model to generate a fixed number of sentences following a set of increasing values for each linguistic property
  - We generate 50 sentences for every value within the set of five values, thus obtaining a total of 250 sentences per property.

# Models and Evaluation

## Models:

Model	Parameters
Gemma	2B
Gemma	7B
LLaMA-2	7B
LLaMA-2	14B
Mistral	7B

## Evaluation:

- We used two different metrics:
  - **Success Rate (SR):** fraction of times the model generated a sentence whose property value exactly corresponds to the one provided.
  - **Spearman coefficient:** correlation coefficients between the increasing property values extracted from EWT and those extracted from the sentences generated by the models.

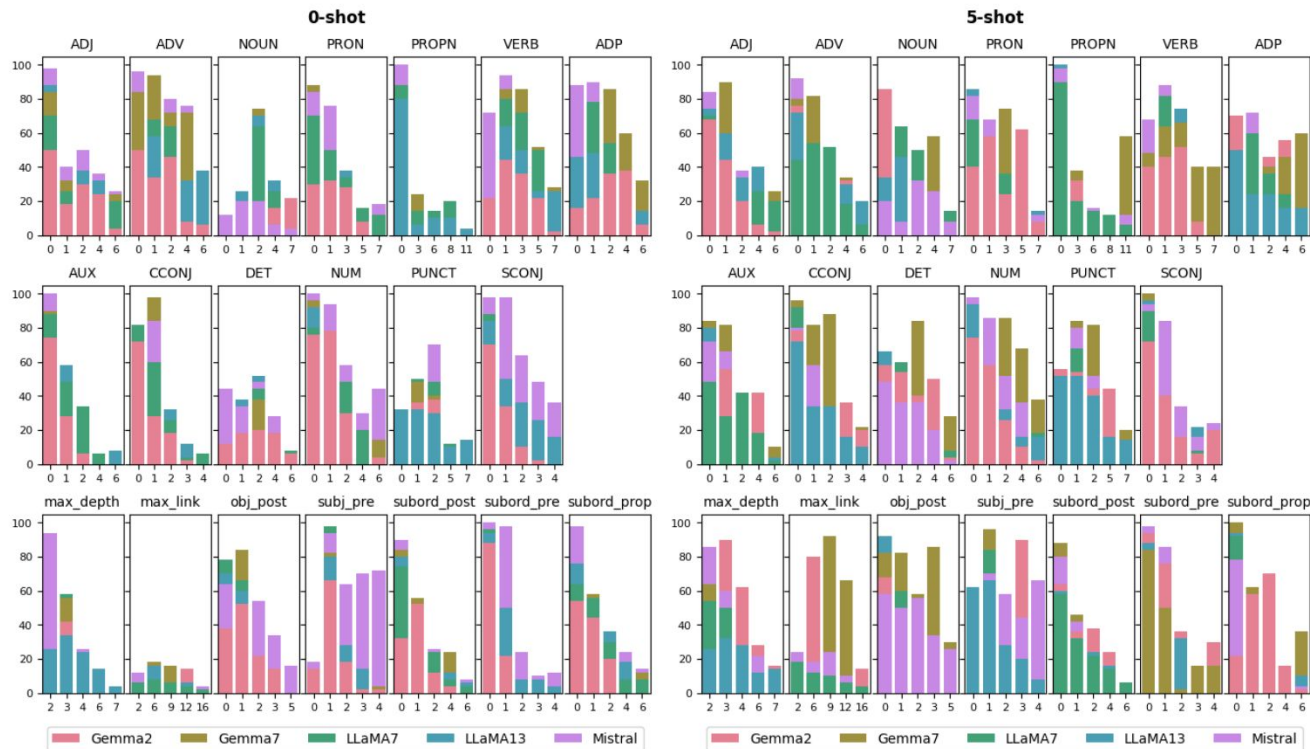
# Success Rate Results

Ling. properties	Gemma2	Gemma7	LLaMA7	LLaMA13	Mistral
<b>Success Rate</b>					
<b>Morphosyntax</b>	<b>0-shot</b>				
ADJ	25.2	36.8	33.6	42	50
ADV	28.8	70.8	34.4	38.8	74
NOUN	8.8	26	23.2	29.6	12.4
PRON	19.6	22.8	36.4	34	41.6
PROPN	25.6	29.2	28	22	22
VERB	25.2	50.8	46.8	37.2	57.6
ADP	23.6	54.4	31.2	31.6	64.4
AUX	21.6	23.6	35.2	37.2	29.2
CCONJ	24	33.2	35.6	35.2	33.2
DET	14.8	15.6	14.8	25.6	32
NUM	37.6	48	43.2	40.8	65.2
PUNCT	14.8	19.2	26	23.6	29.2
SCONJ	23.2	27.6	27.6	42.4	68.8
<b>Avg</b>	<b>22.52</b>	<b>35.23</b>	<b>32</b>	<b>33.85</b>	<b>44.58</b>
<b>Syntax</b>	<b>0-shot</b>				
max_depth	13.6	17.6	16.4	20.4	29.2
max_link	9.2	7.2	5.2	6.8	3.6
obj_post	25.2	36.4	35.2	36.4	40.8
subj_pre	20.4	21.2	22.8	26.4	63.6
subord_post	20	36.8	29.2	29.6	32.8
subord_pre	22	23.2	24	32.8	48.8
subord_prop	23.6	37.6	33.2	37.2	41.6
<b>Avg</b>	<b>19.14</b>	<b>25.71</b>	<b>23.71</b>	<b>27.09</b>	<b>37.2</b>

Ling. properties	Gemma2	Gemma7	LLaMA7	LLaMA13	Mistral
<b>Success Rate</b>					
<b>Morphosyntax</b>	<b>5-shot</b>				
ADJ	28	47.6	34.4	42.8	45.6
ADV	33.2	47.2	34.8	41.2	51.6
NOUN	43.6	20.4	34.4	28.4	18.8
PRON	38.4	45.6	34	39.2	39.6
PROPN	30.4	40.4	28.4	29.6	29.2
VERB	29.2	51.6	38.4	37.6	52
ADP	44.8	47.2	28.8	26	42
AUX	31.6	45.6	27.6	38.4	35.6
CCONJ	38	63.6	34	33.2	34.4
DET	41.2	37.6	31.6	30	28.4
NUM	34	71.6	44.8	43.2	57.6
PUNCT	42	40	34	34.8	31.6
SCONJ	30.8	43.2	31.2	40.8	50.4
<b>Avg</b>	<b>35.78</b>	<b>46.28</b>	<b>33.57</b>	<b>35.78</b>	<b>39.75</b>
<b>Syntax</b>	<b>5-shot</b>				
max_depth	52	24.4	30.4	22.4	38.8
max_link	22.8	47.2	10	10.8	15.6
obj_post	31.6	67.6	32	43.6	44.8
subj_pre	51.2	42.4	41.6	36.8	50
subord_post	33.2	34	26.4	27.6	34
subord_pre	47.6	33.6	34	31.6	45.6
subord_prop	33.6	50.4	34.8	32.8	34
<b>Avg</b>	<b>38.86</b>	<b>42.8</b>	<b>29.89</b>	<b>29.37</b>	<b>37.54</b>



# How do LLMs Follow Constraints Across Values?



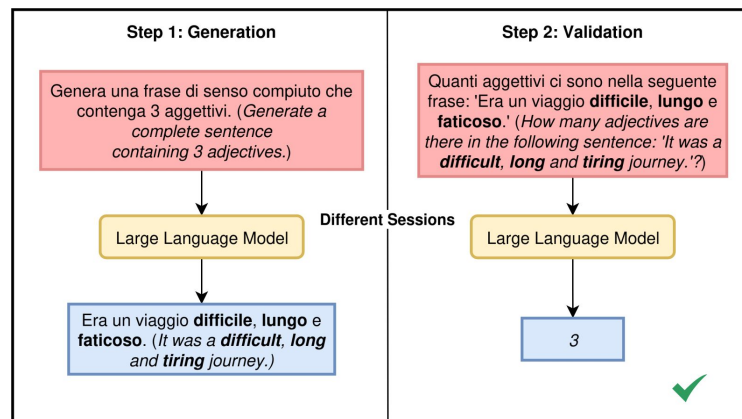
# Spearman Results

Ling. properties	Gemma2	Gemma7	LLaMA7	LLaMA13	Mistral
<b>Spearman</b>					
<b>Morphosyntax</b>					
<b>0-shot</b>					
ADJ	0.59	0.73	0.74	0.79	0.92
ADV	##	0.88	0.52	0.65	0.95
NOUN	0.63	0.72	0.62	0.66	0.93
PRON	0.26	0.35	0.58	0.80	0.91
PROPN	##	0.66	0.60	0.67	0.88
VERB	0.56	0.83	0.78	0.71	0.76
ADP	0.55	0.89	0.48	0.64	0.96
AUX	##	0.29	0.32	0.56	0.96
CCONJ	0.27	0.33	0.35	0.33	0.42
DET	0.28	0.36	##	0.28	0.79
NUM	0.49	0.74	0.60	0.62	0.94
PUNCT	0.24	0.54	0.63	0.61	0.78
SCONJ	##	0.44	0.40	0.62	0.92
<b>Avg</b>	<b>0.30</b>	<b>0.60</b>	<b>0.51</b>	<b>0.61</b>	<b>0.86</b>
<b>Syntax</b>					
<b>0-shot</b>					
max_depth	##	0.18	##	##	0.76
max_link	##	0.44	0.57	0.43	0.75
obj_post	0.21	0.47	0.37	0.38	0.59
subj_pre	##	##	0.37	0.13	0.84
subord_post	0.13	0.65	0.44	0.58	0.59
subord_pre	##	0.33	0.13	0.34	0.72
subord_prop	0.28	0.60	0.45	0.67	0.83
<b>Avg</b>	<b>0.08</b>	<b>0.38</b>	<b>0.33</b>	<b>0.36</b>	<b>0.73</b>

Ling. properties	Gemma2	Gemma7	LLaMA7	LLaMA13	Mistral
<b>Spearman</b>					
<b>Morphosyntax</b>					
<b>5-shot</b>					
ADJ	0.19	0.78	0.76	0.79	0.86
ADV	0.43	0.62	0.52	0.71	0.80
NOUN	0.87	0.76	0.77	0.75	0.90
PRON	0.63	0.65	0.78	0.85	0.81
PROPN	0.25	0.87	0.76	0.81	0.81
VERB	0.42	0.77	0.77	0.72	0.87
ADP	0.46	0.81	0.53	0.61	0.77
AUX	0.37	0.70	0.53	0.59	0.60
CCONJ	0.53	0.56	0.52	0.52	0.60
DET	0.49	0.77	0.65	0.65	0.65
NUM	##	0.63	0.72	0.74	0.77
PUNCT	0.60	0.70	0.73	0.79	0.69
SCONJ	0.26	0.66	0.62	0.71	0.74
<b>Avg</b>	<b>0.42</b>	<b>0.71</b>	<b>0.67</b>	<b>0.71</b>	<b>0.76</b>
<b>Syntax</b>					
<b>5-shot</b>					
max_depth	0.80	0.56	0.39	0.40	0.78
max_link	0.40	0.86	0.64	0.52	0.70
obj_post	0.42	0.84	0.51	0.62	0.72
subj_pre	0.59	0.52	0.55	0.47	0.74
subord_post	0.58	0.59	0.53	0.54	0.77
subord_pre	0.12	0.24	0.33	0.35	0.56
subord_prop	0.39	0.79	0.68	0.66	0.74
<b>Avg</b>	<b>0.47</b>	<b>0.63</b>	<b>0.52</b>	<b>0.51</b>	<b>0.71</b>

# Controllable Text Generation To Evaluate Linguistic Abilities of Italian LLMs

- Focus on Italian LLMs
- Two-steps evaluation:
  - Generation:
    - e.g. “*Generate a sentence with 2 adjectives*”.
  - Validation:
    - e.g. “*How many adjectives does this sentence have?*”.



# Generation Results

Model	n_tokens	NOUN	VERB	ADJ	ADV	subj	obj	subord	Avg
ANITA	<b>.25/.97</b>	<b>.47/.97</b>	<b>.46/.96</b>	<b>.53/.96</b>	<b>.45/.91</b>	.23/.29	<b>.36/.44</b>	<b>.52/.91</b>	<b>.41/.80</b>
Camoscio	.1/.51	.14/.44	.16/.18	.17/.28	.16/.17	.25/.15	.2/##	.22/.13	.18/.23
Cerbero	.06/.57	.15/.56	.24/.5	.25/.38	.22/.31	.23/.15	.23/.13	.26/.33	.21/.37
DanteLLM	.11/.79	.15/.54	.22/.66	.29/.62	.21/.35	<b>.36/.34</b>	.31/.3	.32/.51	.25/.51
Italia	.03/.62	.09/.34	.16/.2	.16/.28	.18/##	.22/.16	.21/.22	.22/.18	.16/.25
LlaMAntino	.05/.57	.12/.48	.19/.43	.17/.31	.2/.23	.33/.3	.23/.17	.23/.28	.19/.35
<b>Avg</b>	<b>.1/.67</b>	<b>.19/.56</b>	<b>.24/.49</b>	<b>.26/.47</b>	<b>.24/.33</b>	<b>.27/.23</b>	<b>.26/.21</b>	<b>.29/.39</b>	

**Table 2**

Success rate and Spearman correlation coefficients ( $SR/\rho$ ) between the linguistic constraints and the feature values extracted from the generated sentences. The best and worst scores for each property and each metric are highlighted in **bold** and *italic* respectively. Non-statistically significant correlation scores are reported with ##.

# Validation Results

	Model	n_tokens	NOUN	VERB	ADJ	ADV	subj	obj	subord	Avg
Cons.	ANITA	.06/.96	.43/.97	.57/.96	.52/.95	.55/.94	.82/.96	.8/.95	.64/.94	<b>.55/.95</b>
	Camoscio	.28/.44	.06/.31	.23/.28	.19/.2	.19/.2	.25/.27	.24/.18	.2/##	.2/.23
	Cerbero	.27/.56	.2/.49	.2/.51	.31/.5	.24/.46	.31/.3	.22/.11	.3/.42	.26/.42
	DanteLLM	.21/##	.18/.59	.12/.63	.33/.6	.13/.35	.37/.43	.25/.28	.31/##	.24/.36
	Italia	.26/.54	.04/.27	.16/.31	.02/.14	.02/.11	.28/.39	.21/.23	.25/.28	.15/.28
	LLaMAntino	.06/##	.07/##	.18/##	.2/##	.14/.24	.42/.71	.31/##	.2/.46	.2/.18
	<b>Avg</b>	<b>.19/.42</b>	<b>.16/.44</b>	<b>.24/.45</b>	<b>.26/.4</b>	<b>.21/.38</b>	<b>.41/.51</b>	<b>.34/.29</b>	<b>.32/.35</b>	
Cons.+	ANITA	.06/.91	.63/.96	.53/.98	.7/.96	.73/.96	.92/.74	.79/.68	.84/.98	<b>.65/.9</b>
	Camoscio	.55/.89	.14/.52	.47/.41	.23/.33	.21/##	.65/.41	.5/.31	.14/##	.36/.36
	Cerbero	.47/.94	.39/.83	.45/.81	.73/.8	.66/.77	.53/.34	.61/.34	.66/.65	.56/.68
	DanteLLM	.38/.94	.36/.8	.39/.82	.63/.85	.32/.44	.56/.45	.51/.36	.63/##	.47/.58
	Italia	.35/.86	.05/.47	.16/.5	.03/##	.08/##	.7/.54	.36/.28	.47/.51	.27/.4
	LLaMAntino	.25/.85	.08/.82	.35/.6	.25/.51	.32/.39	.38/.64	.59/##	.4/.53	.33/.54
	<b>Avg</b>	<b>.34/.9</b>	<b>.28/.73</b>	<b>.39/.68</b>	<b>.43/.58</b>	<b>.39/.43</b>	<b>.62/.52</b>	<b>.56/.33</b>	<b>.52/.45</b>	

**Table 3**

Success rate and Spearman correlation coefficients ( $SR/\rho$ ) between the linguistic constraints asked during sentence generation and the values predicted during the validation step. Consistency results are reported for both the overall sentences (*Cons.*) and a filtered subset of sentences that correctly matched the asked linguistic constraint (*Cons.+*).

# Validation Results

	Model	n_tokens	NOUN	VERB	ADJ	ADV	subj	obj	subord	Avg
Cons.	ANITA	.06/.96	.43/.97	.57/.96	.52/.95	.55/.94	.82/.96	.8/.95	.64/.94	<b>.55/.95</b>
	Camoscio	.28/.44	.06/.31	.23/.28	.19/.2	.19/.2	.25/.27	.24/.18	.2/##	.2/.23
	Cerbero	.27/.56	.2/.49	.2/.51	.31/.5	.24/.46	.31/.3	.22/.11	.3/.42	.26/.42
	DanteLLM	.21/##	.18/.59	.12/.63	.33/.6	.13/.35	.37/.43	.25/.28	.31/##	.24/.36
	Italia	.26/.54	.04/.27	.16/.31	.02/.14	.02/.11	.28/.39	.21/.23	.25/.28	.15/.28
	LLaMAntino	.06/##	.07/##	.18/##	.2/##	.14/.24	.42/.71	.31/##	.2/.46	.2/.18
	<b>Avg</b>	.19/.42	.16/.44	.24/.45	.26/.4	.21/.38	<b>.41/.51</b>	.34/.29	.32/.35	
Cons.+	ANITA	.06/.91	.63/.96	.53/.98	.7/.96	.73/.96	.92/.74	.79/.68	.84/.98	<b>.65/.9</b>
	Camoscio	.55/.89	.14/.52	.47/.41	.23/.33	.21/##	.65/.41	.5/.31	.14/##	.36/.36
	Cerbero	.47/.94	.39/.83	.45/.81	.73/.8	.66/.77	.53/.34	.61/.34	.66/.65	.56/.68
	DanteLLM	.38/.94	.36/.8	.39/.82	.63/.85	.32/.44	.56/.45	.51/.36	.63/##	.47/.58
	Italia	.35/.86	.05/.47	.16/.5	.03/##	.08/##	.7/.54	.36/.28	.47/.51	.27/.4
	LLaMAntino	.25/.85	.08/.82	.35/.6	.25/.51	.32/.39	.38/.64	.59/##	.4/.53	.33/.54
	<b>Avg</b>	.34/.9	.28/.73	.39/.68	.43/.58	.39/.43	<b>.62/.52</b>	.56/.33	.52/.45	

**Table 3**

Success rate and Spearman correlation coefficients ( $SR/\rho$ ) between the linguistic constraints asked during sentence generation and the values predicted during the validation step. Consistency results are reported for both the overall sentences (*Cons.*) and a filtered subset of sentences that correctly matched the asked linguistic constraint (*Cons.+*).

# Selected Findings

- Models tend to adhere slightly more accurately to **morphosyntactic constraints** rather than syntactic ones
- Models are capable of distinguishing when they are asked to generate a sentence **with or without a given feature**
- Constraining generation for a specific linguistic element does not always primarily enhance that element, suggesting that the **models are not simply creating longer sentences, but rather sentences with a varied (morpho)syntactic structure**
- When validating each model against their own generated sentences, we noticed that the **generation abilities do not always align with the ability of the models to recognize the linguistic properties of their generated sentences.**

# Conclusion and Future Directions

- LLMs have reached astonishing performance in almost all NLP tasks
- Their success has led to a growing interest in their evaluation, alongside studies analyzing their behavior and internal mechanisms
- Despite significant progress, there is still a lot to do!

## Future Directions:

- Studying and evaluating generalization of LLMs across different scenarios, domains and languages ([Hupkes et al., 2023](#))
- Mechanistic Interpretability ([Elhage et al, 2021](#); [Olsson et al., 2022](#))
- Memorization vs. Generalization ([Patil et al., 2024](#))





Istituto di Linguistica  
Computazionale  
"Antonio Zampolli"

 Consiglio Nazionale delle Ricerche



# Thanks for the attention!



<https://alemiaschi.github.io/>



[@AlessioMiaschi](https://twitter.com/AlessioMiaschi)



<http://www.italianlp.it/>



[@ItaliaNLP\\_Lab](https://twitter.com/ItaliaNLP_Lab)

# References

- Bengio, Yoshua, et al. (2003). "A neural probabilistic language model." *The journal of machine learning research* 3, pages 1137-1155
- Vaswani, Ashish, et al. (2017). "Attention is all you need." *Advances in Neural Information Processing Systems* (NEURIPS)
- Miaschi A., Brunato D., Dell'Orletta F., Venturi G. (2020). Linguistic Profiling of a Neural Language Models. In *Proceedings of the 28th International Conference on Computational Linguistics* (COLING 2020, Barcelona)
- Dominique Brunato, Andrea Cimino, Felice Dell'Orletta, Giulia Venturi, and Simonetta Montemagni. 2020. Profiling-UD: a Tool for Linguistic Profiling of Texts. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 7145–7151, Marseille, France. European Language Resources Association
- Miaschi A., Dell'Orletta F., Venturi G. (2024). Linguistic Knowledge Can Enhance Encoder-Decoder Models (*If You Let It*). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation* (LREC-COLING 2024, Turin)
- Miaschi A., Dell'Orletta F., Venturi G. (2024). Evaluating Large Language Models via Linguistic Profiling. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing* (EMNLP 2024, Miami, Florida)
- Ciaccio C., Dell'Orletta F., Miaschi A., Venturi G. (2024). Controllable Text Generation To Evaluate Linguistic Abilities of Italian LLMs. In *Proceedings of the Tenth Italian Conference on Computational Linguistics* (CLiC-it 2024, Pisa)
- Hupkes, Dieuwke, et al. "A taxonomy and review of generalization research in NLP." *Nature Machine Intelligence* 5.10 (2023): 1161-1174
- Elhage, Nelson, et al. "A mathematical framework for transformer circuits." *Transformer Circuits Thread* 1.1 (2021): 12
- Olsson, Catherine, et al. "In-context learning and induction heads." *arXiv preprint arXiv:2209.11895* (2022)
- Patil, Abhinav et al. "Filtered Corpus Training (FiCT) Shows that Language Models can Generalize from Indirect Evidence". In *Transactions of the Association for Computational Linguistics* (2024)