

# Building and Evolving a Multilingual LLM from Scratch

Javier Aula-Blasco, PhD

/ CyberColombia  
/ 8th Colombian HPC Summer School  
/ June 17, 2025



# Language Technologies Lab



# Infrastructure

# MareNostrum 5

GPP based on Intel Sapphire Rapids, 6480 nodes, 45.9 PFlops peak.

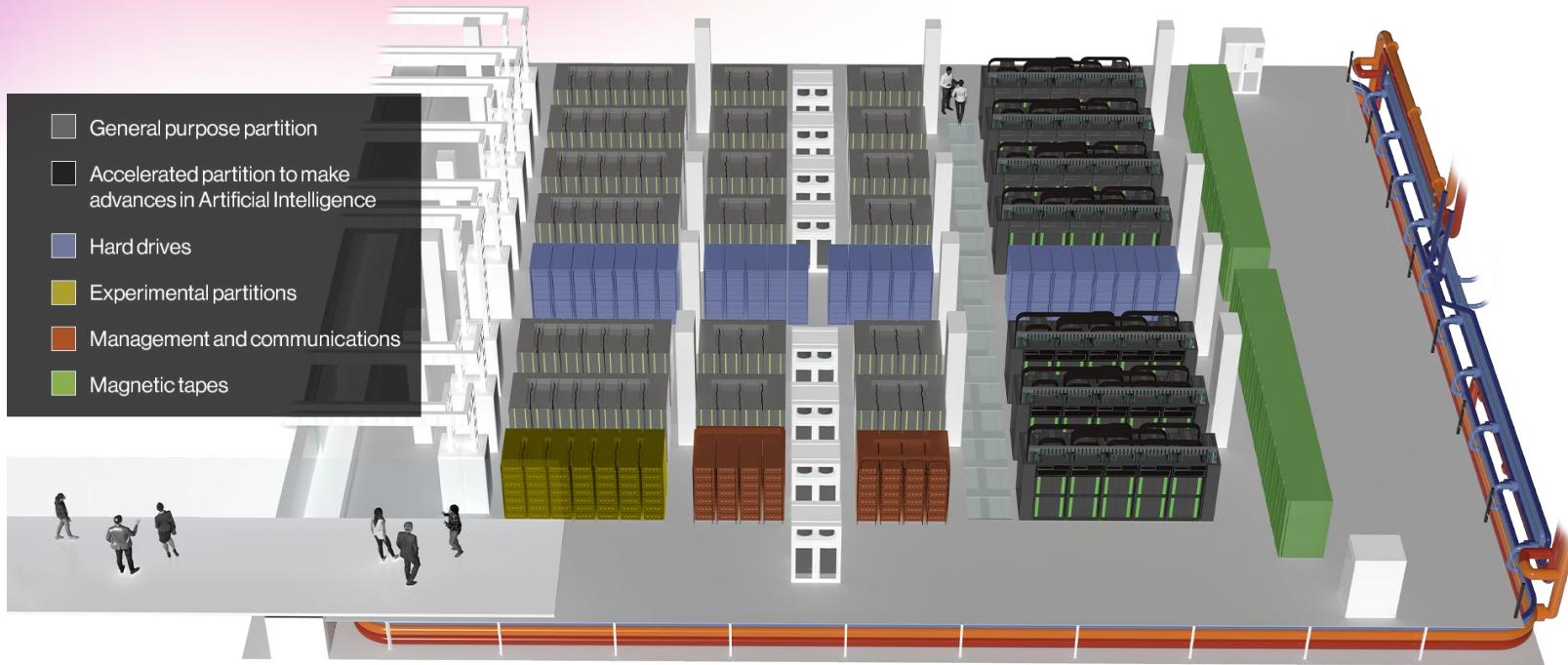
**ACC based on Intel Sapphire Rapids and NVIDIA Hopper GPUs, 1120 nodes with 4 Hopper GPUs each one, 260 PFlops peak.**

GPP - Next Gen based on NVIDIA GRACE CPUs.

ACC - Next Generation coming soon.



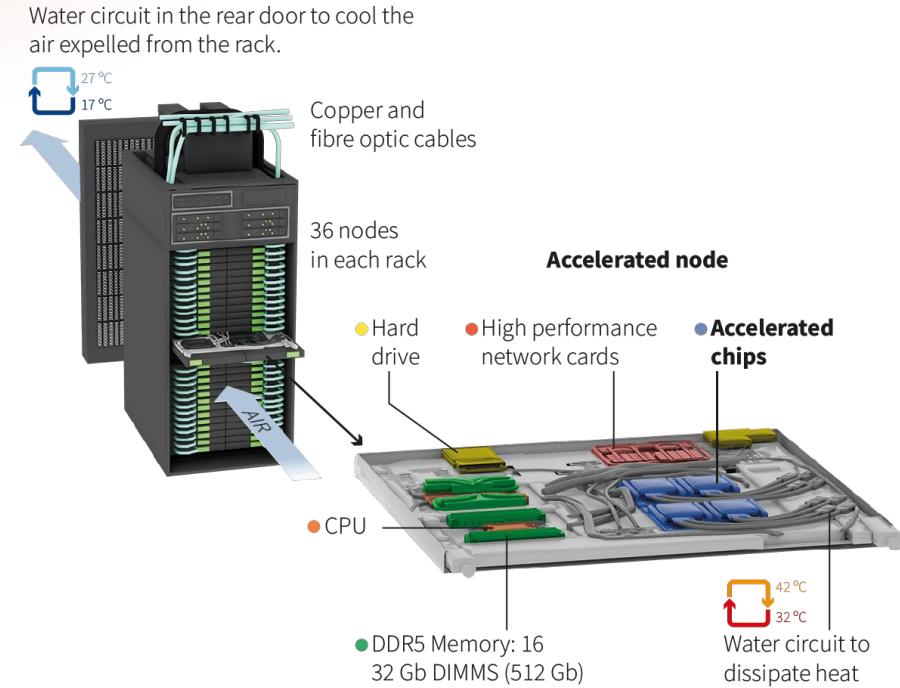
# MareNostrum 5



# Accelerated Partition

## Node configuration:

- 2x Intel Sapphire Rapids 8460Y+ at 2.3Ghz and 40c each
- 512 Gb of Main memory, using DDR5
- 4x NVIDIA Hopper GPU with 64Gb HBM2 memory
- 480Gb on NVMe storage
- 4x NDR200 (BW per node 800Gb/s)



# Salamandra & ALIA models

 **BSC-LT/ALIA-40b**  
Text Generation • Updated Apr 3 • ↓ 331 • ❤ 59

 **BSC-LT/salamandra-7b-instruct**  
Text Generation • Updated Feb 20 • ↓ 11.3k • ❤ 59

 **BSC-LT/salamandra-7b**  
Text Generation • Updated Feb 20 • ↓ 432 • ❤ 27

 **BSC-LT/salamandra-2b-instruct**  
Text Generation • Updated Feb 20 • ↓ 970 • ❤ 21

 **BSC-LT/salamandra-2b**  
Text Generation • Updated Feb 20 • ↓ 1.06k • ❤ 23

 **BSC-LT/salamandra-7b-instruct-fp8**  
Text Generation • Updated Mar 26 • ↓ 26 • ❤ 1

 **BSC-LT/salamandra-7b-instruct-gptq**  
Text Generation • Updated Mar 26 • ↓ 66

 **BSC-LT/salamandra-2b-instruct-fp8**  
Text Generation • Updated Mar 26 • ↓ 10

 **BSC-LT/salamandra-2b-instruct-gptq**  
Text Generation • Updated Mar 26 • ↓ 5

 **BSC-LT/salamandra-7b-base-fp8**  
Text Generation • Updated Mar 26 • ↓ 17 • ❤ 1

 **BSC-LT/salamandra-7b-base-gptq**  
Text Generation • Updated Mar 26 • ↓ 52

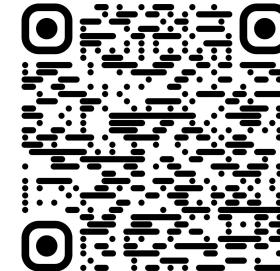
 **BSC-LT/salamandra-2b-base-fp8**  
Text Generation • Updated Mar 26 • ↓ 19

 **BSC-LT/salamandra-2b-base-gptq**  
Text Generation • Updated Mar 18 • ↓ 9

## Salamandra Technical Report

Language Technologies Unit

Barcelona Supercomputing Center



# Data

# Pre-training data

In the pre-training phase, it is essential to have a massive volume of data. This volume is hardly achievable with conventional data sources:

**3.7T**

Chinchilla optimal for 175B  
parameter model ([Hoffman et al](#))

vs

**0.000911T**

All Wikipedia in Spanish

vs

**0.0000005T**

Don Quixote

Larger sources (mainly mass scraping) are used, to which pre-processing tasks must be applied, including, depending on the case, the following layers: extraction, deduplication, language identification, filtering/quality assessment, removal of evaluation data or detection of inappropriate content.

# Pre-training resources

Parameters	Nodes (MN5)	Time (days)	Model	Tokens	Tokens in Spanish
~175B	-	?	ChatGPT3.5	?	?
176B	512	5	BLOOM	0,366 T	0,039 T (10,8%)
180B	512	43	Falcon-180B	3,5 T	0,059 T (1,68%)
40B	512	17	Falcon-40B	1T	0,017 T (1,68%)
70B	512	29	Llama2-70B	2 T	0,002 T (0,13%)
7B	512	2	Llama2-7B	2 T	0,002 T (0,13%)
7B	512	5	Mistral-7B	8T?	?

# Challenges

1. Language imbalance.
2. Traceability and control of the data used for training.
3. Data quality.
4. Legal aspects.
5. Hate speech, biases, toxic content, irrelevance, etc.

# Approach

Creating a **data infrastructure** together with the different actors in society.

Developing, contributing to, and identifying **data processing tools** that facilitate addressing these challenges.

Researching how to achieve high-performance models **with limited resources**.

# Language imbalance

The main source of text data is Common Crawl (CC) or derivatives.

This crawling is done by domains and the main target is English.

In previous dumps, Catalan and Spanish have accounted for around 0.2% and 4.5% of the total respectively, against 45% for English.

crawl	CC-MAIN-2023-40	CC-MAIN-2023-50	CC-MAIN-2024-10
language	%	%	%
eng	46.4328	44.4285	46.4536
deu	5.8355	5.4499	5.3948
rus	5.5709	6.0303	5.8095
jpn	4.7475	5.1508	5.0934
fra	4.6351	4.3933	4.3263
spa	4.6199	4.5391	4.5462
cat	0.2096	0.2053	0.2036

% languages in last three CC dumps

# Data agreements



WIKIPEDIA  
The Free Encyclopedia



Universitat  
de les Illes Balears



PARLAMENT DE CATALUNYA

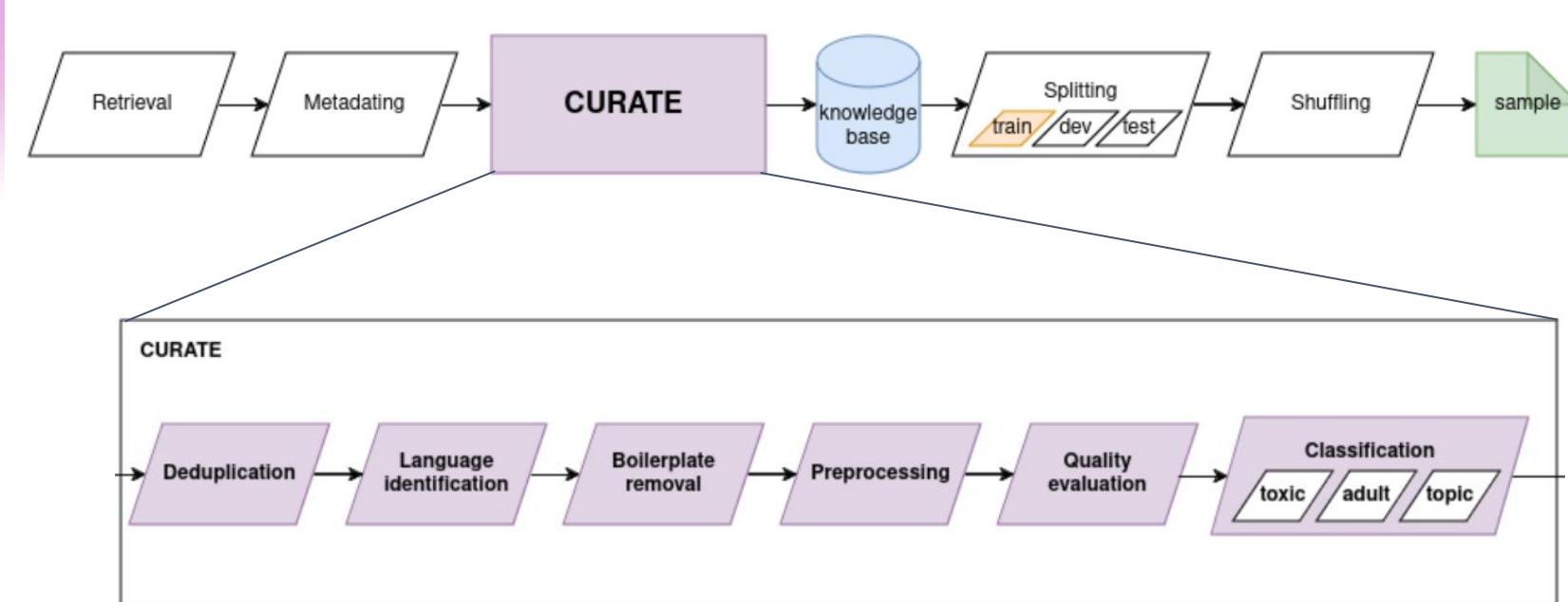
HiTZ



elMón



# CURATE Pipeline



A CURATED CATALOG: Rethinking the Extraction of Pretraining Corpora for Mid-Resourced Languages (Palomar-Giner et al. 2024). *The 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*.

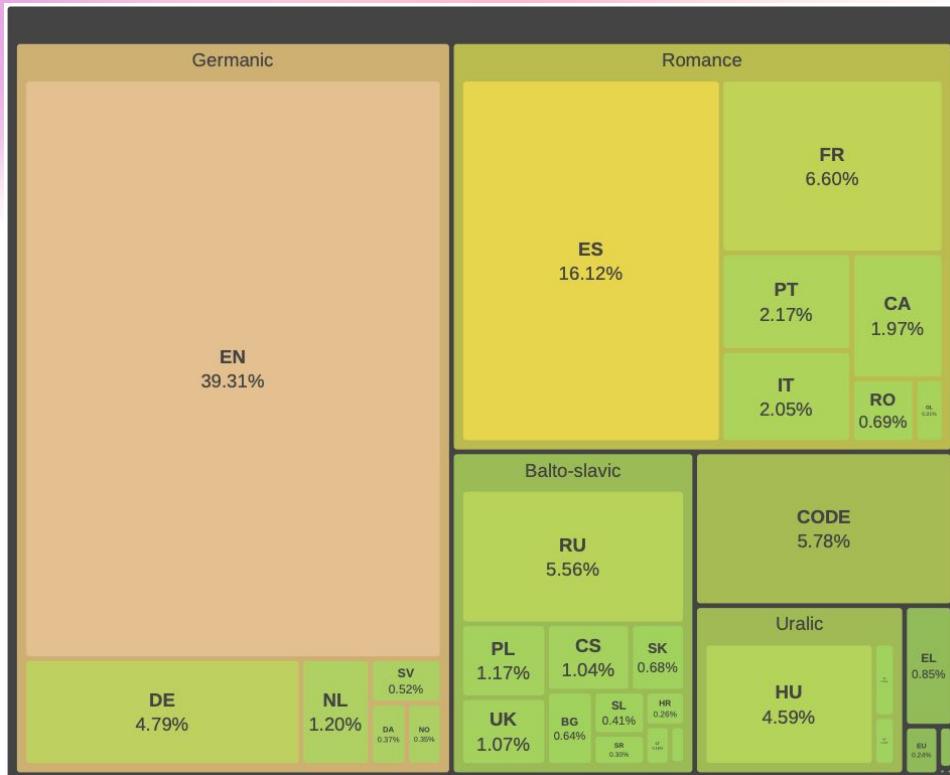
# CURATE Pipeline



langtech-bsc/CURATE



# Language distribution

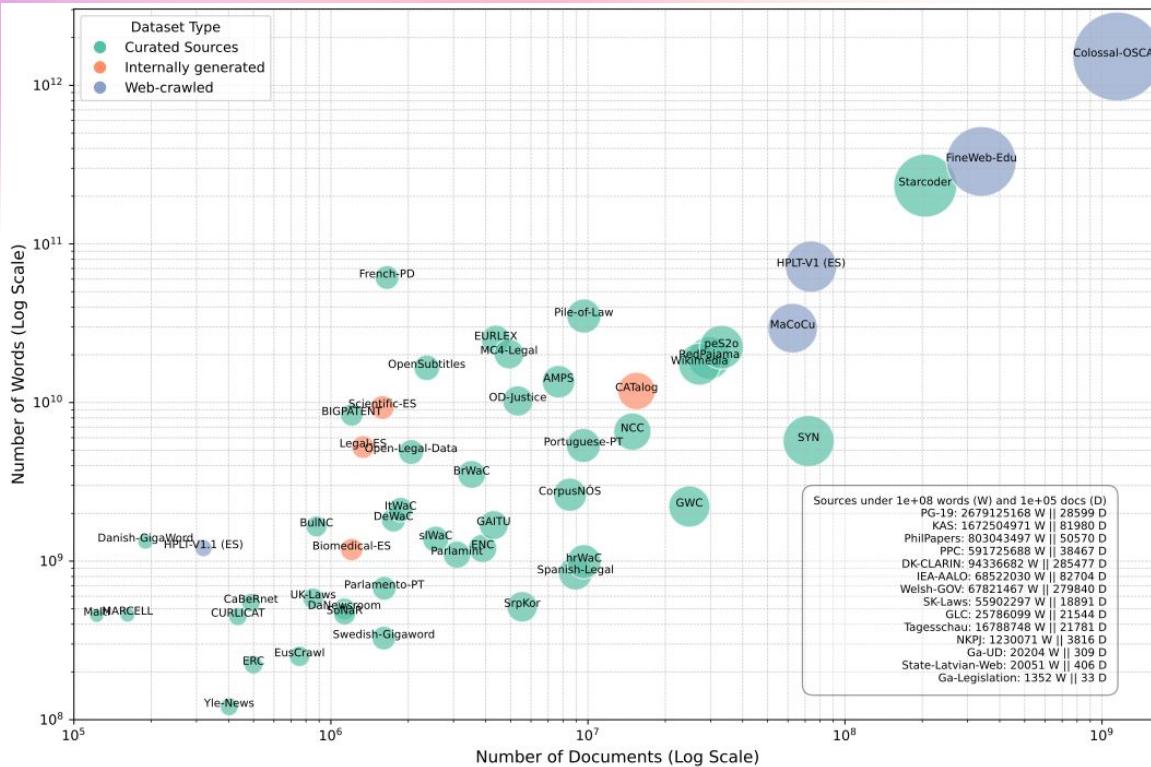


**35 EU languages + code**

**Documents: 1.967,73 M**

**Words: 2.251.075,65 M**

# Source distribution



Ungoliant: An Optimized Pipeline  
for the Generation of a Very Large-Scale Multilingual Web Corpus

Julien Abadji<sup>1</sup> Pedro Javier Ortiz Suárez<sup>1,2</sup> Laurent Romary<sup>1</sup> Benoît Sagot<sup>1</sup>

<sup>1</sup>Inria, Paris, France

<sup>2</sup>Sorbonne Université, Paris, France

Colossal OSCAR curates 20 multilingual CommonCrawl snapshots.

---

The FineWeb Datasets: Decanting the Web for the Finest Text Data at Scale

---

Guilherme Penedo Hynek Kydlíček Loubna Ben alall Anton Lozhkov  
Margaret Mitchell Colin Raffel Leandro Von Werra Thomas Wolf  
🤗 Hugging Face

StarCoder contains data in 86 programming languages.

# CATalog corpus



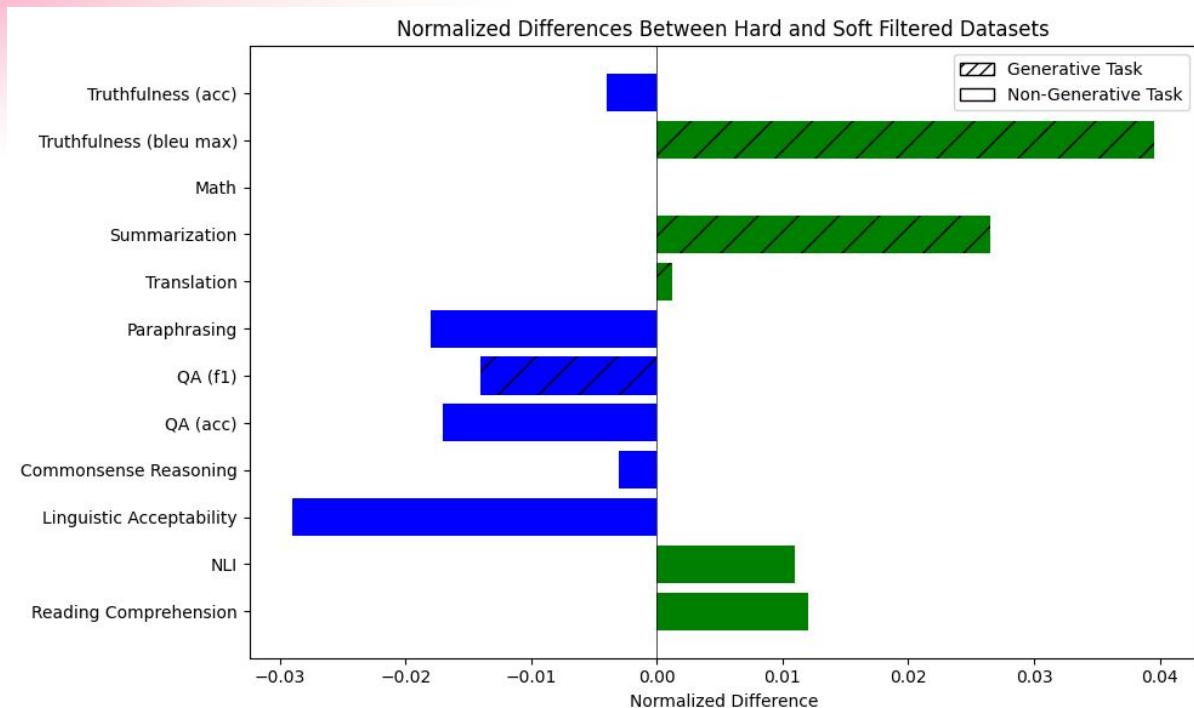
projecte-aina/CATalog



# Pre-training data

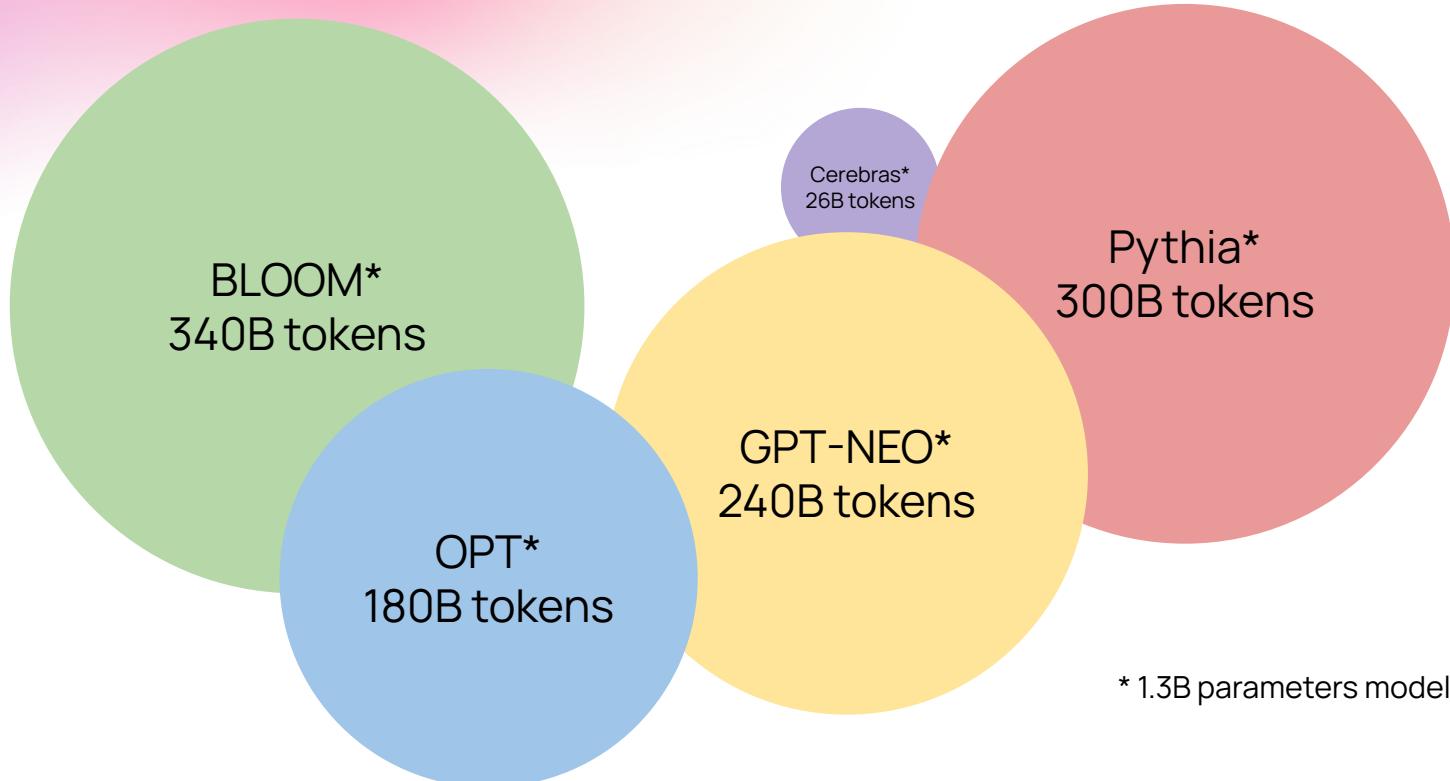
Two 7B with two different thresholds on data quality. On the left, soft filtering; on the right, hard filtering.

Multilingual setup.

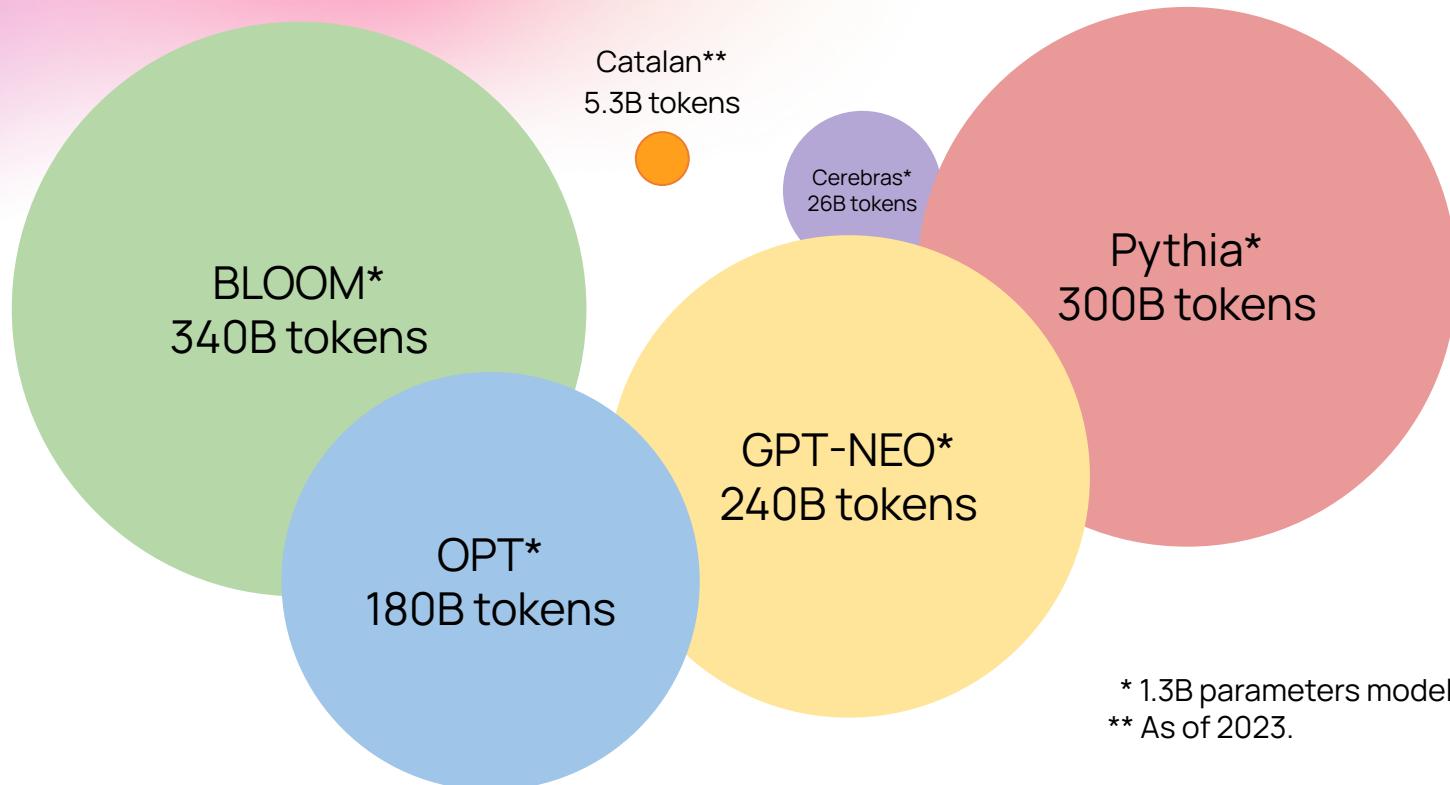


# Pre-training

# Number of tokens

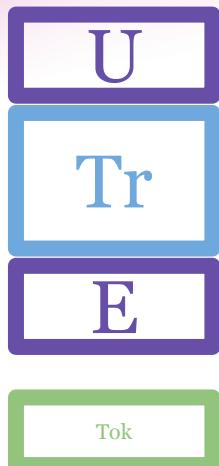


# Number of tokens

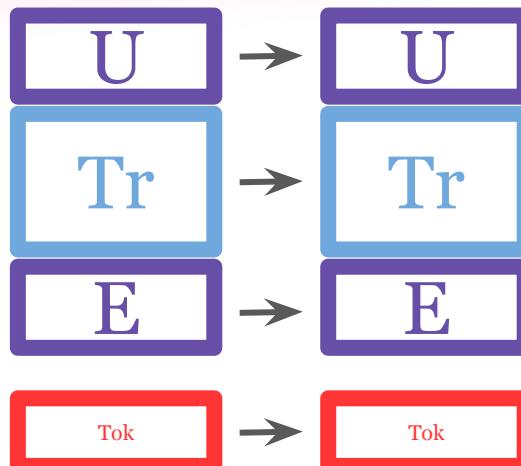


# Strategies

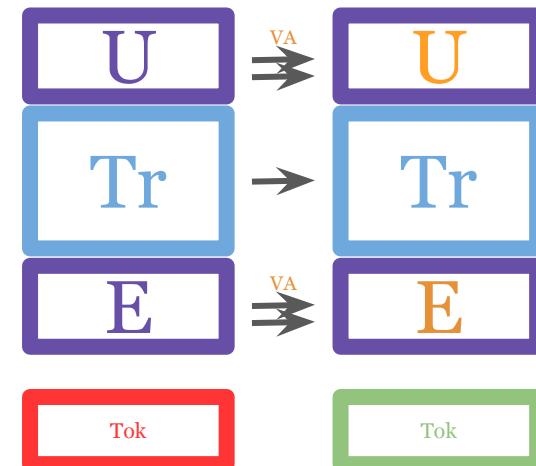
From Scratch



Continual Pre-Training



Vocabulary Adaptation



# Decoders

The big dog was chasing the \_\_\_



cat

# Decoders

The big dog was chasing the



The | \_big | \_dog | \_was | \_chas | ing | \_the

# Decoders

The big dog was chasing the



The | \_big | \_dog | \_was | \_chas | ing | \_the



[o, o, 1, o, ..., o, o]

# Decoders

The big dog was chasing the



The | \_big | \_dog | \_was | \_chas | ing | \_the



[o, 1, o, o, ..., o, o]

# Decoders

The big dog was chasing the



The | \_big | \_dog | \_was | \_chas | ing | \_the



[o, o, o, o, ..., 1, o]

# Decoders

The big dog was chasing the

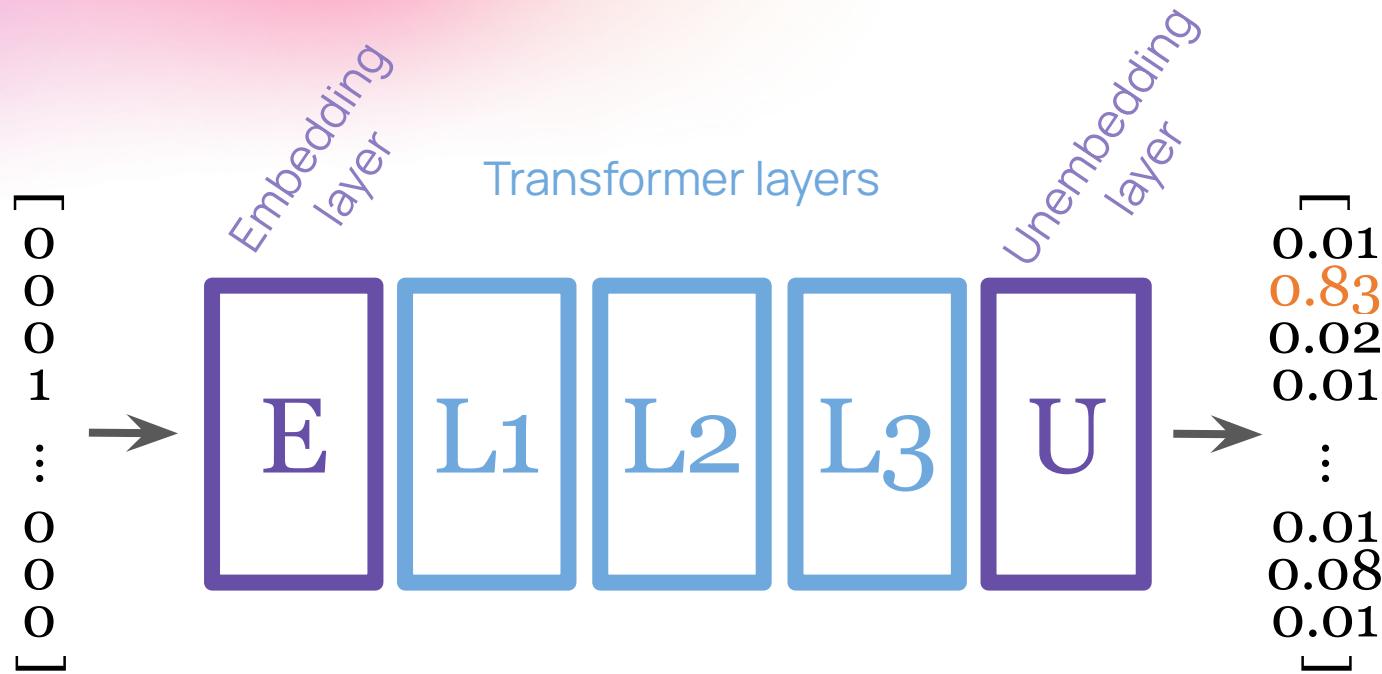


The | \_big | \_dog | \_was | \_chas | ing | \_the

[0, 0, 0, ..., 1, 0, 0, 0]  
[0, 1, 0, ..., 0, 0, 0, 0]  
[1, 0, 0, ..., 0, 0, 0, 0]  
[0, 0, 1, ..., 0, 0, 0, 0]  
[0, 0, 0, ..., 0, 0, 0, 1]  
[0, 0, 0, ..., 0, 1, 0, 0]  
[0, 0, 0, ..., 0, 0, 1, 0]

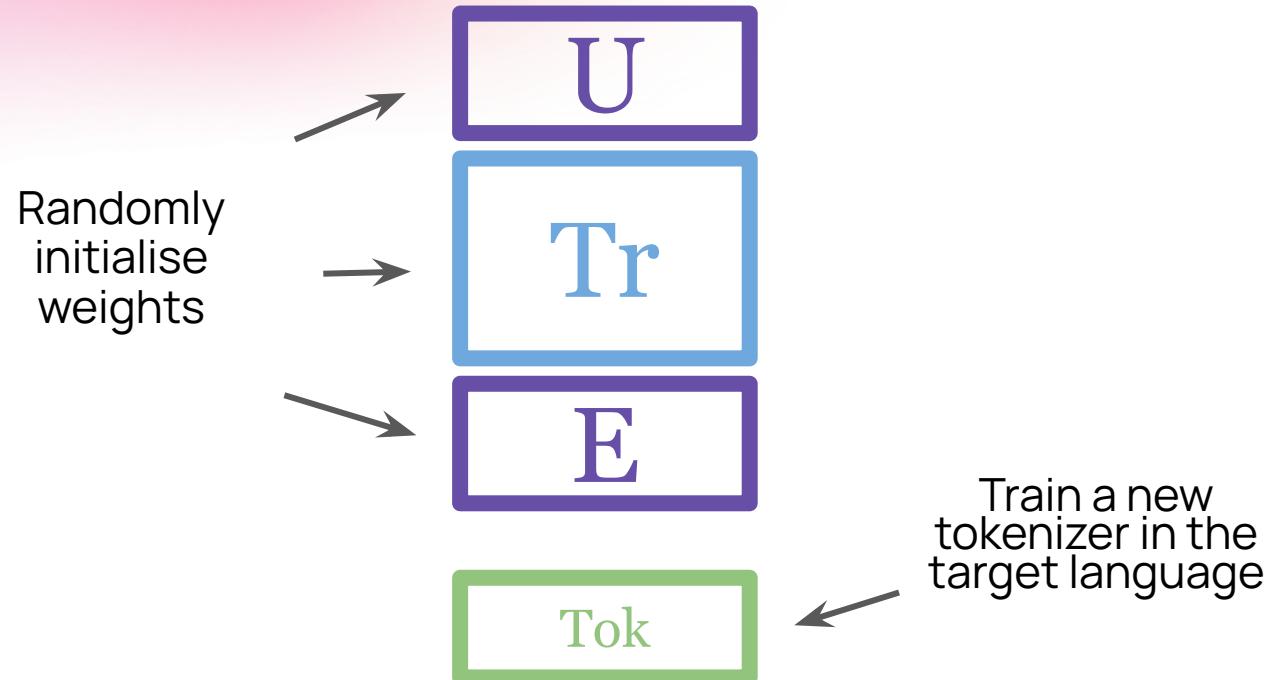


# Decoders

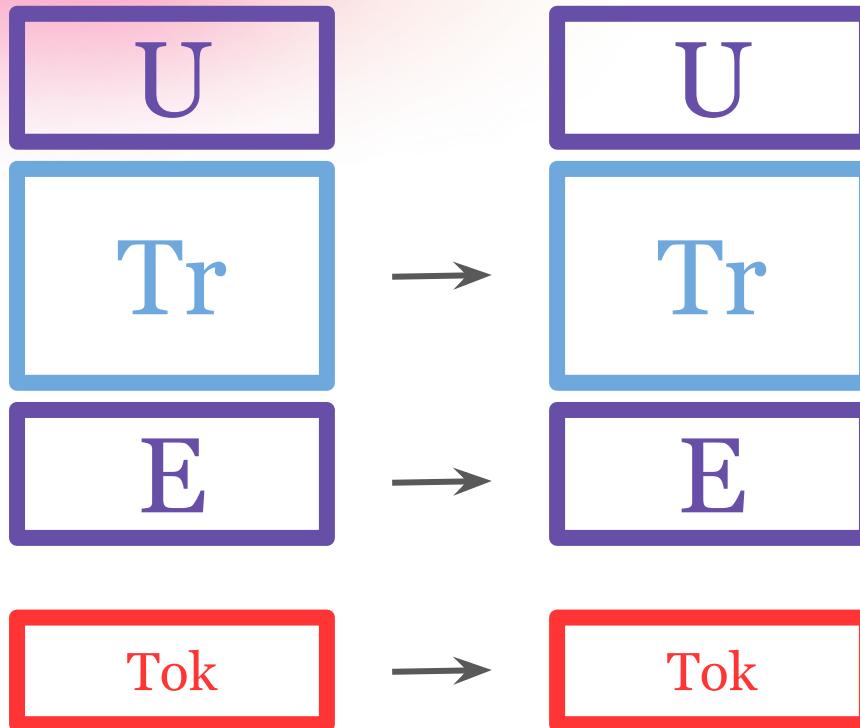


\* This is an oversimplification.

# From scratch



# CPT



Vaig anar a fer la compra aviat.



Tok

Tok

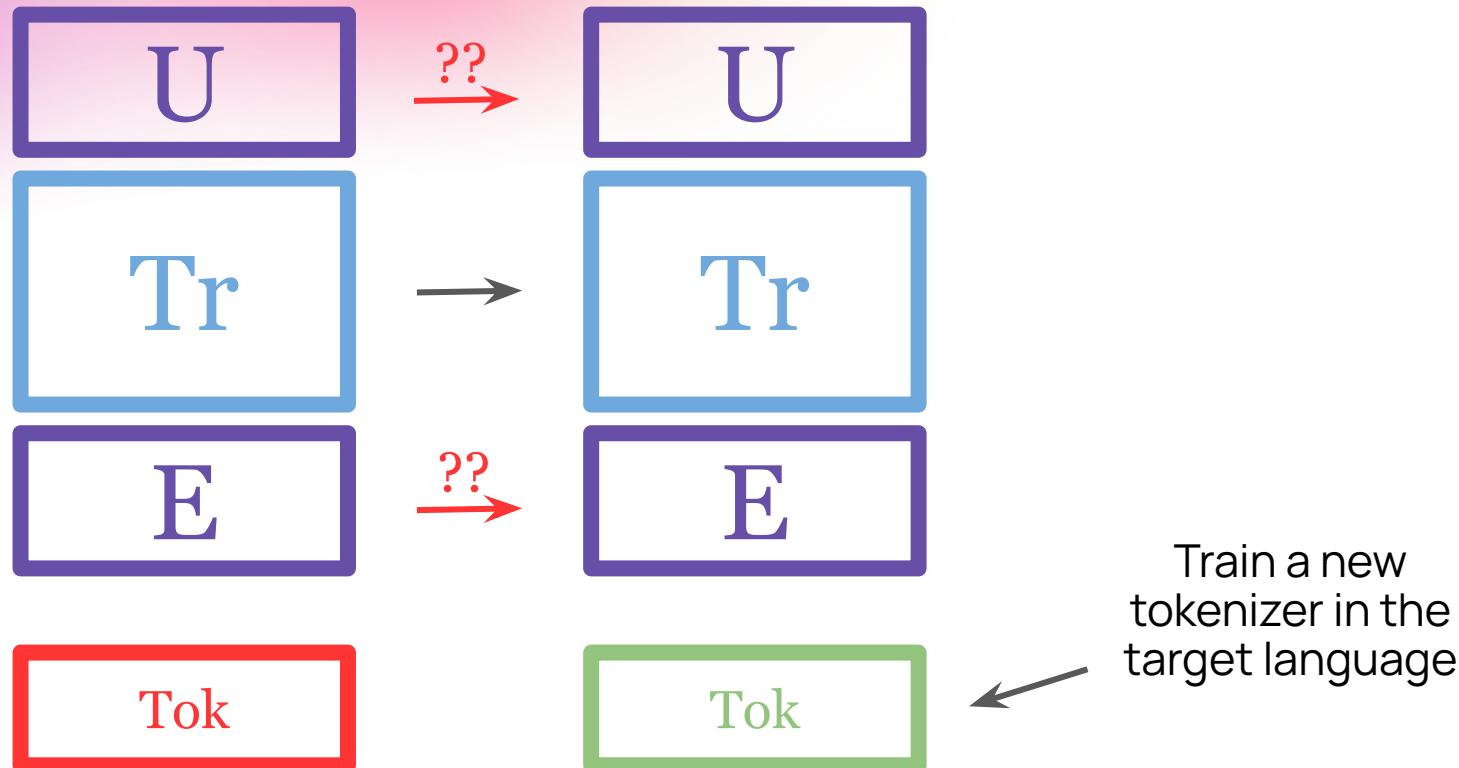


Va | ig | \_an | ar | \_a | \_fer | \_la  
| \_comp | ra | \_av | iat | .

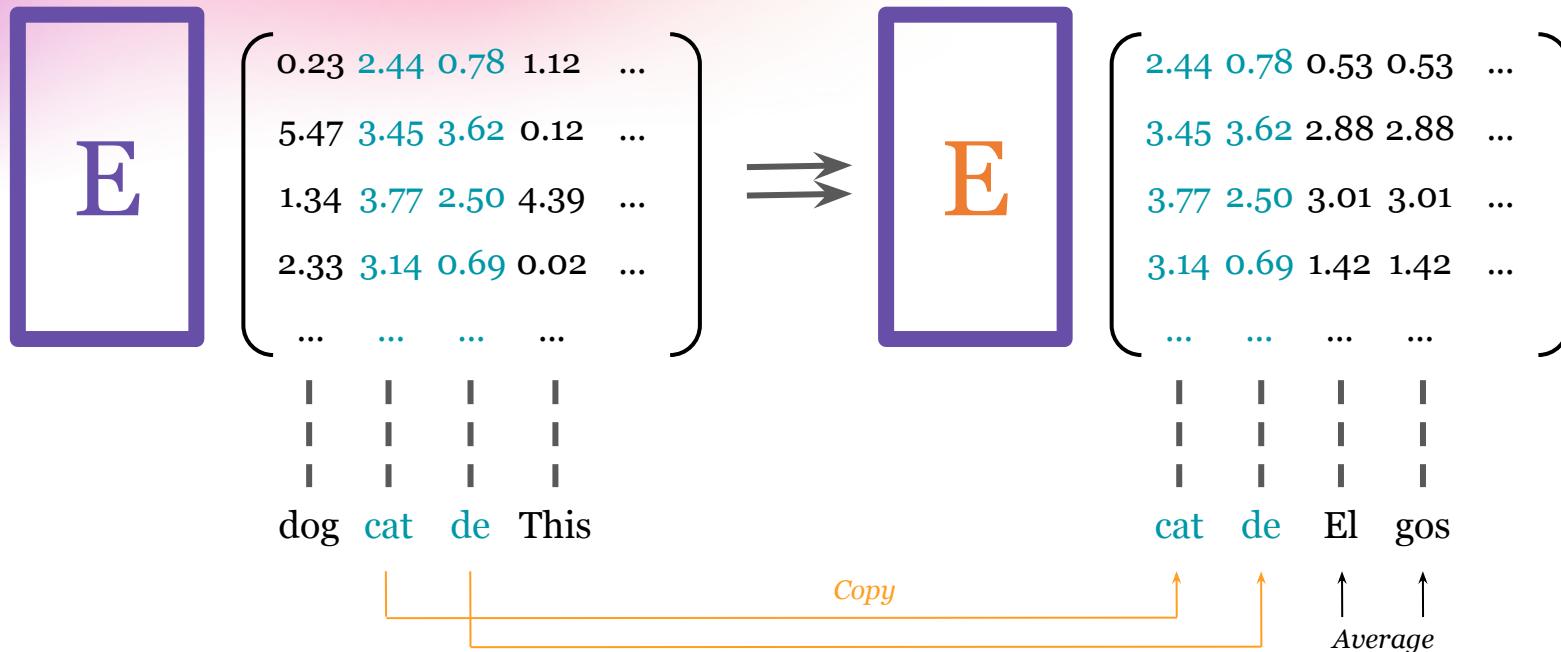
Vaig | \_anar | \_a | \_fer | \_la |  
\_compra | \_aviat | .

*On the Cross-lingual Transferability of  
Monolingual Representations*  
By Mikel Artetxe, Sebastian Ruder, Dani Yogatama

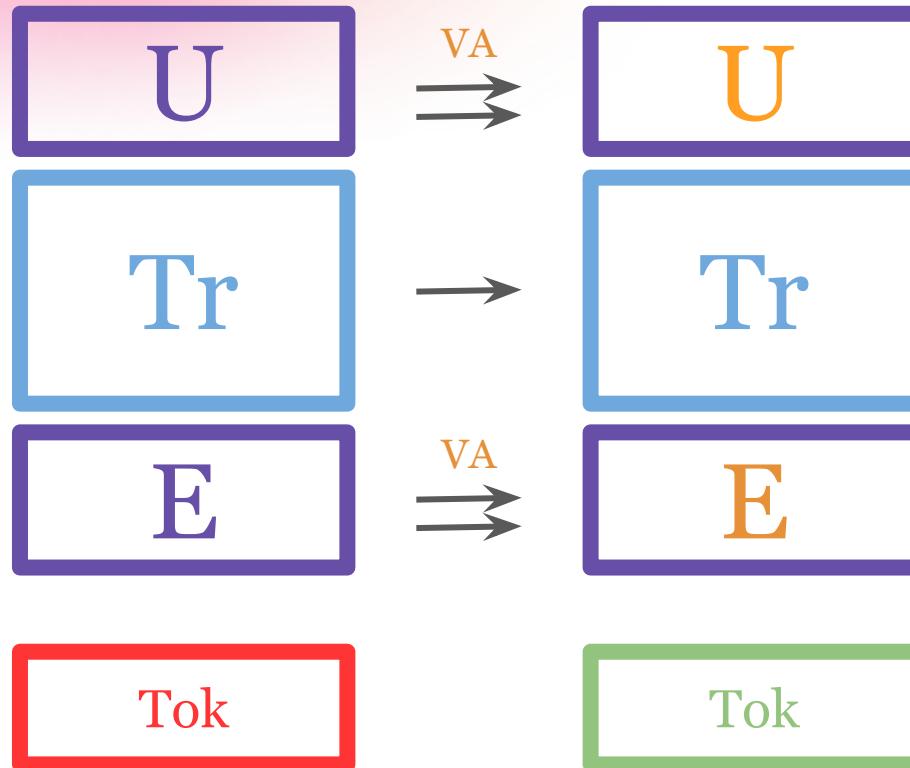
# Vocabulary Adaptation



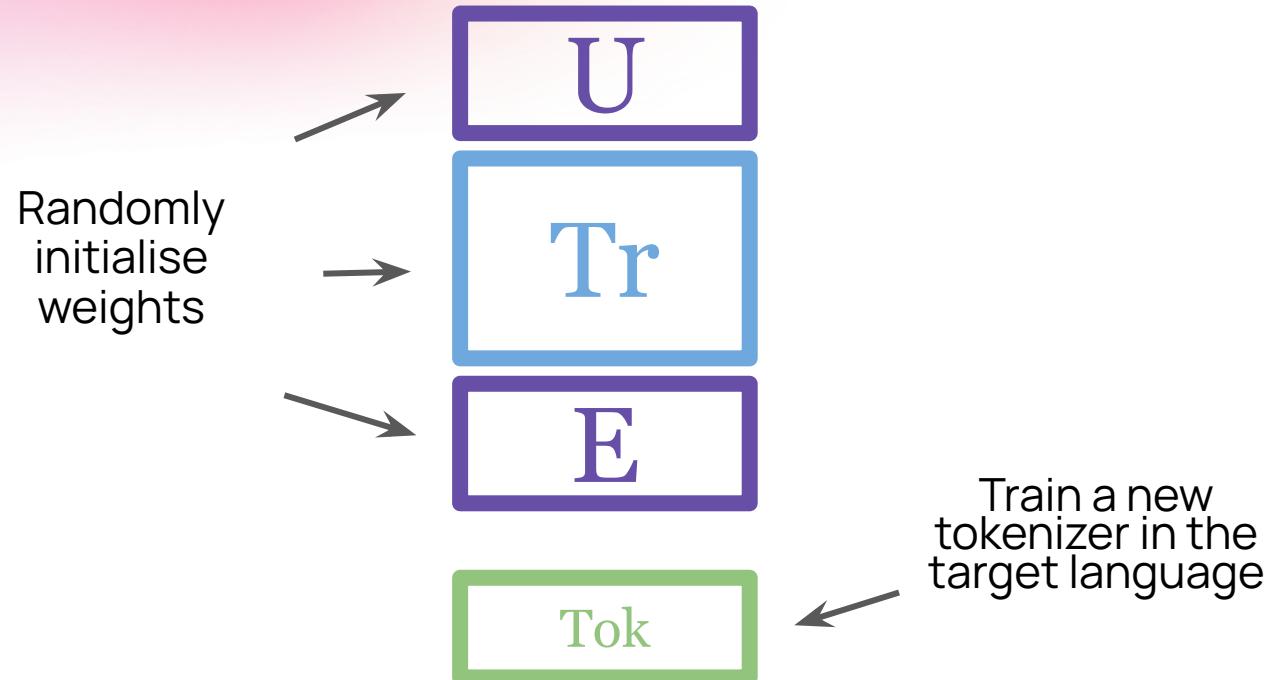
# Vocabulary Adaptation



# Vocabulary Adaptation



# From scratch

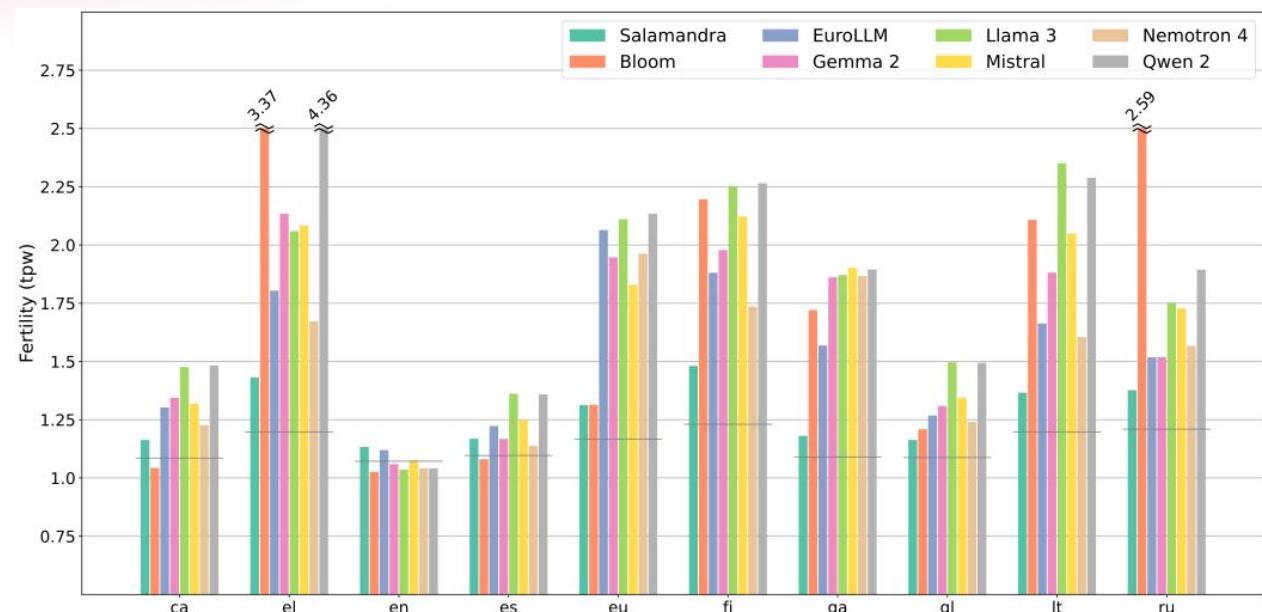


# Tokenizer

SentencePiece implementation of **Byte-Pair Encoding (BPE)** algorithm.

**256,000 multilingual tokens.** 100 reserved for adaptations.

Byte fallback, digit splitting, normalization + more in Tech Report.



# Distribution

Models trained using **NVIDIA NeMo Framework** from a random initialization of weights and **Adam optimizer**.

Checkpointing was performed every 2,000 or 5,000 steps.

Model	Nodes	GPUs	TP	PP	DP	MBS	GBS	Context	Batch Size	Tokens
2B	64	256	1	1	256	1	512	8,192	~4M	12.9T
7B	128	512	4	1	128	2	512	8,192	~4M	12.9T
40B	512	2,048	4	2	128	1	1,024	4,096	~4M	~9T

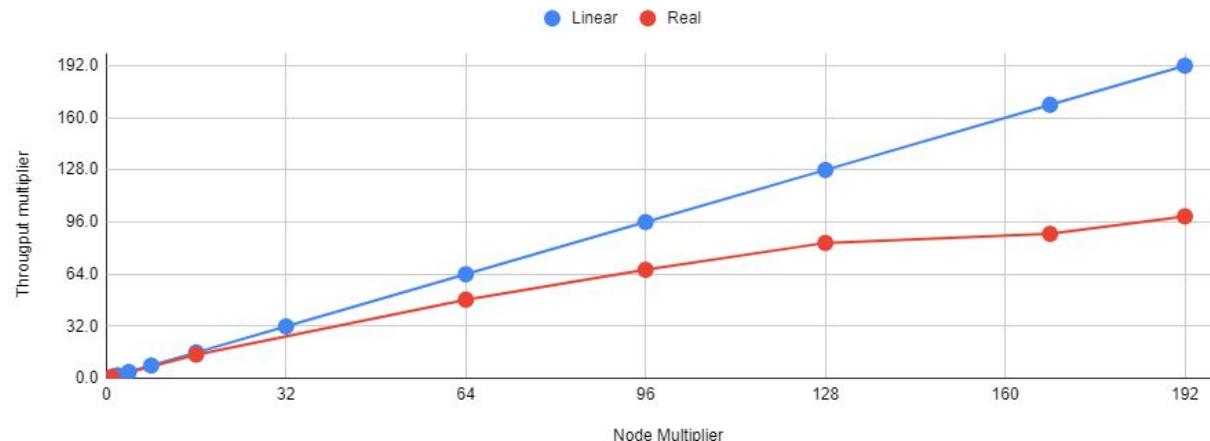
Extended  
to 32k

# Scalability

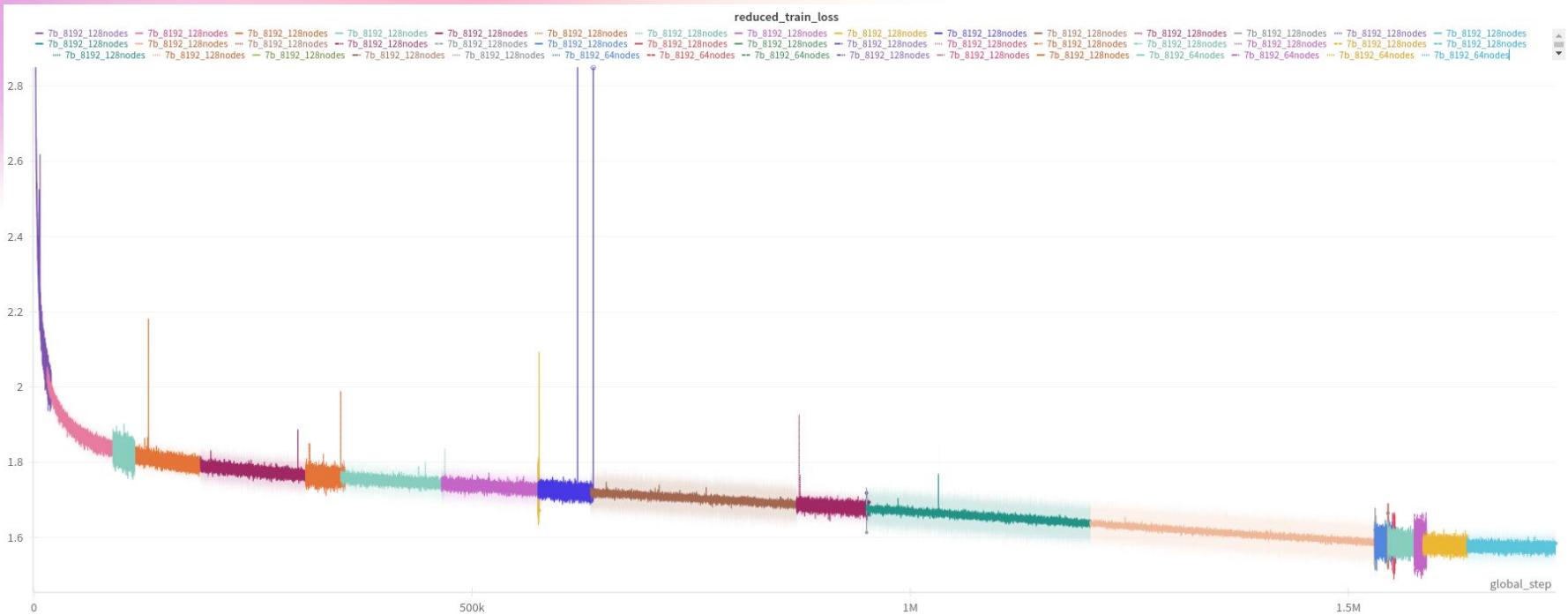
For the 7B model, going from 64 nodes to 128 resulted in a throughput increase of 86%, but going from 64 nodes to 192 this increase was reduced to 68%.

For the 40B model, scaling from 256 nodes to 512 yields an increase of 73%.

Scalability for 7B models

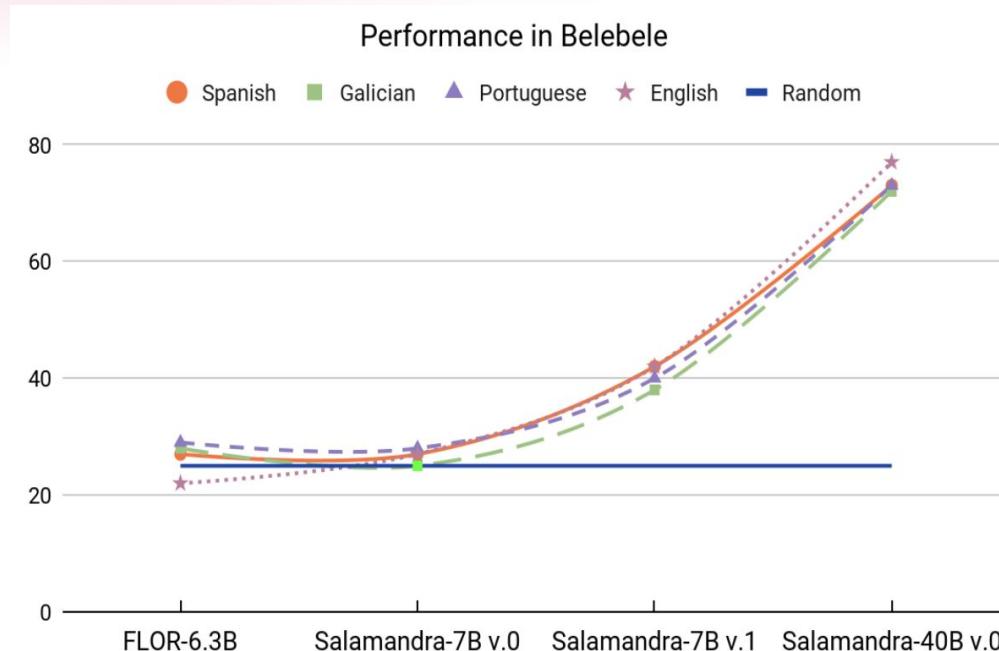


# Training loss



# Annealing

Successful **targeted solutions** to models' limitations.

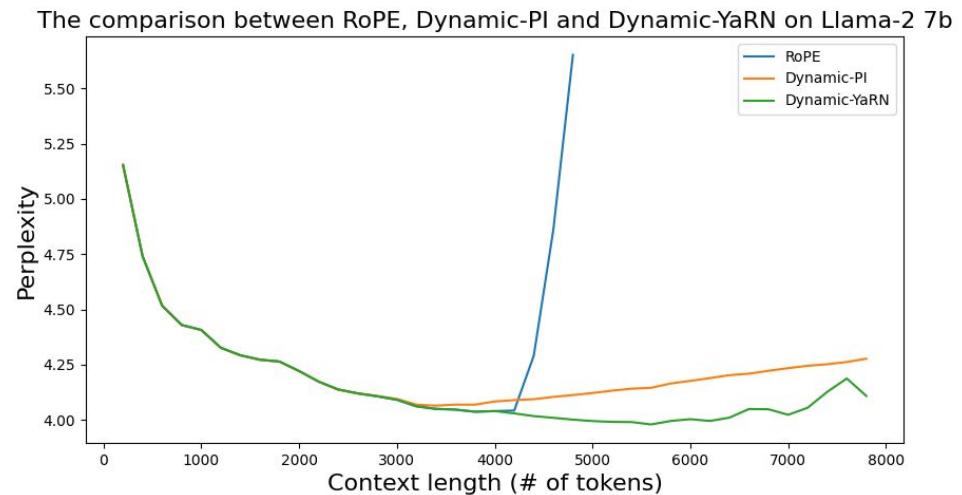


# Just scale context?

Decoder-only transformers (like GPT) have no hard architectural limit on the context window size.

But if we naively increase the context window (e.g., from 4K to 8K tokens), the model's performance collapses.

This isn't due to hardware, it's a failure of positional understanding.



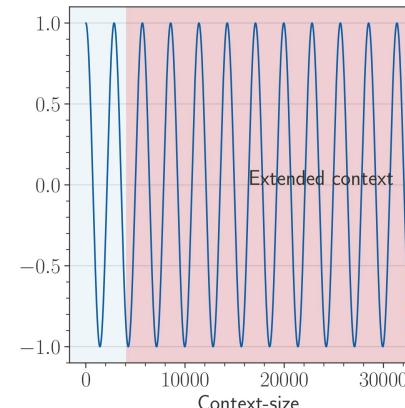
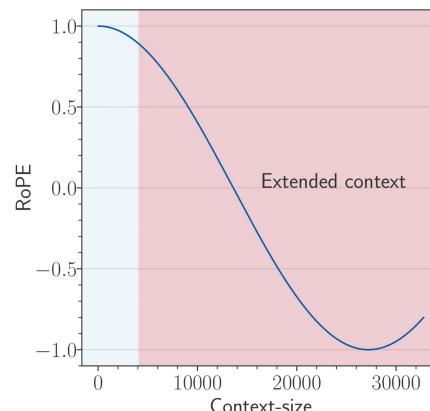
YaRN: Efficient Context Window Extension of Large Language Models (Peng et al. 2023).

# Waves of position

Transformers use positional encodings to understand where tokens appear.

These encodings act like waves, mapping position to a signal the model reasons over.

Extending context length means we need to stretch or scale these waves, but doing this wrong leads to loss of resolution or generalization.



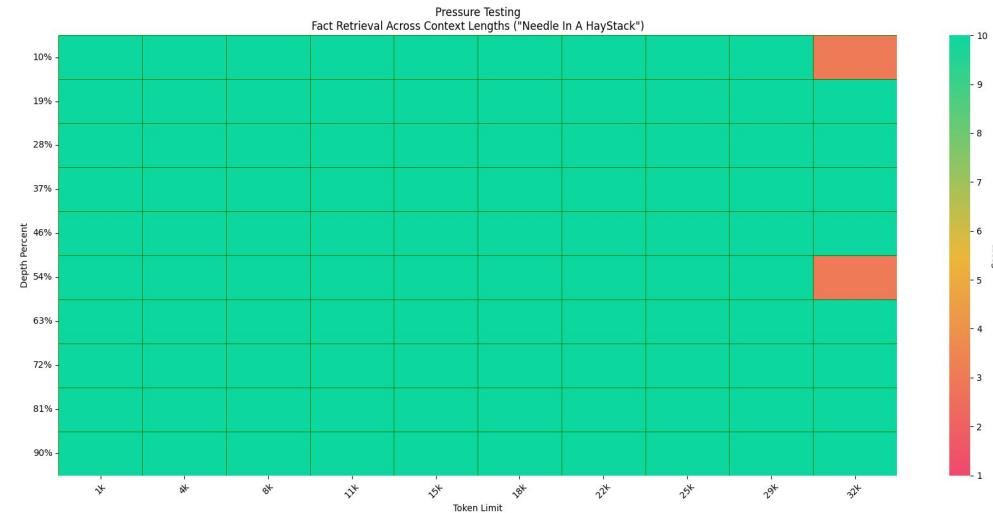
Base of RoPE Bounds Context Length (Men et al. 2024).

# Evaluating extension

Needle in a Haystack: A benchmark where a short "needle" is inserted into a long sequence. The task: recover the needle.

In standard models (trained on short contexts), performance drops sharply as the haystack gets longer.

**Needle:** The best thing to do in San Francisco is eat a sandwich and sit in Dolores Park on a sunny day.



# Evaluating extension

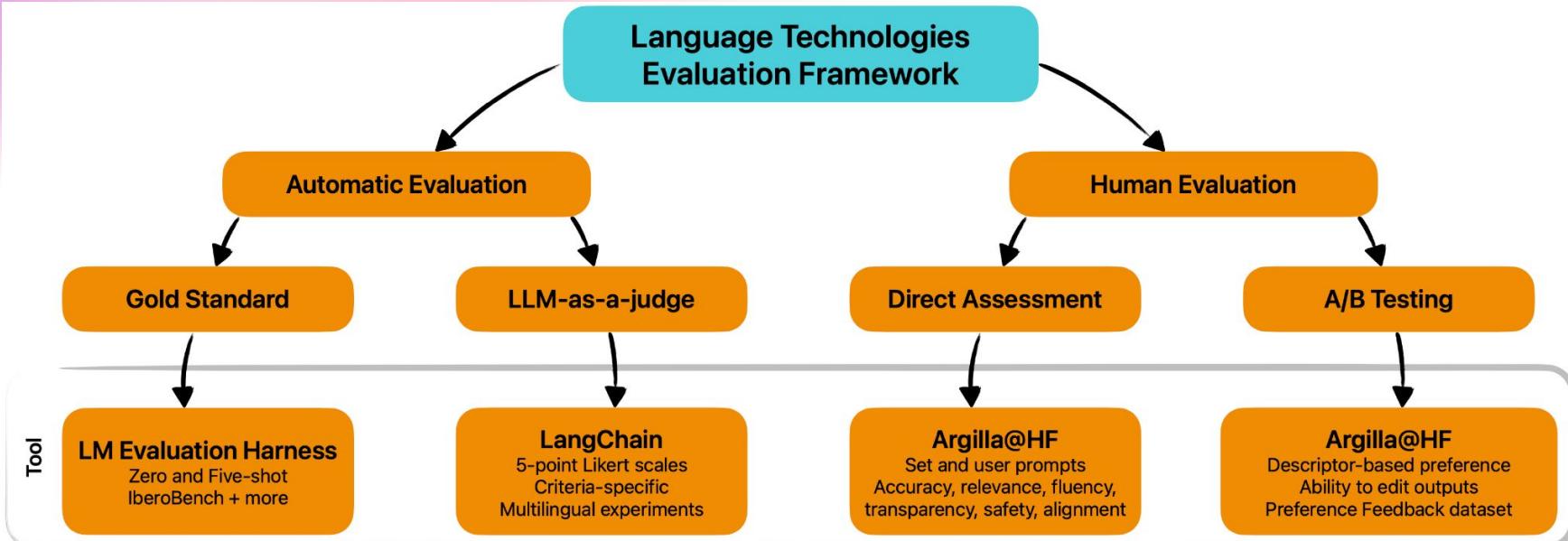
RULER (Hsieh et al., 2024) benchmarks how model performance degrades across increasing context lengths on standard NLP tasks.

Long Context is still not solved. Even the best long-context models show gradual degradation.

Models	Claimed Length	Effective Length	4K	8K	16K	32K	64K	128K	Avg.	wAvg. (inc)	wAvg. (dec)
Llama2 (7B)	4K	-	85.6								
Gemini-1.5-Pro	1M	>128K	96.7	95.8	96.0	95.9	95.9	94.4	95.8	95.5 <sub>(1st)</sub>	96.1 <sub>(1st)</sub>
GPT-4	128K	64K	96.6	96.3	95.2	93.2	87.0	81.2	91.6	89.0 <sub>(2nd)</sub>	94.1 <sub>(2nd)</sub>
Llama3.1 (70B)	128K	64K	96.5	95.8	95.4	94.8	88.4	66.6	89.6	85.5 <sub>(4th)</sub>	93.7 <sub>(3rd)</sub>
Qwen2 (72B)	128K	32K	96.9	96.1	94.9	94.1	79.8	53.7	85.9	79.6 <sub>(9th)</sub>	92.3 <sub>(4th)</sub>
Command-R-plus (104B)	128K	32K	95.6	95.2	94.2	92.0	84.3	63.1	87.4	82.7 <sub>(7th)</sub>	92.1 <sub>(5th)</sub>
GLM4 (9B)	1M	64K	94.7	92.8	92.1	89.9	86.7	83.1	89.9	88.0 <sub>(3rd)</sub>	91.7 <sub>(6th)</sub>
Llama3.1 (8B)	128K	32K	95.5	93.8	91.6	87.4	84.7	77.0	88.3	85.4 <sub>(5th)</sub>	91.3 <sub>(7th)</sub>
GradientAI/Llama3 (70B)	1M	16K	95.1	94.4	90.8	85.4	80.9	72.1	86.5	82.6 <sub>(8th)</sub>	90.3 <sub>(8th)</sub>
Mixtral-8x22B (39B/141B)	64K	32K	95.6	94.9	93.4	90.9	84.7	31.7	81.9	73.5 <sub>(11th)</sub>	90.3 <sub>(9th)</sub>
Yi (34B)	200K	32K	93.3	92.2	91.3	87.5	83.2	77.3	87.5	84.8 <sub>(6th)</sub>	90.1 <sub>(10th)</sub>
Phi3-medium (14B)	128K	32K	93.3	93.2	91.1	86.8	78.6	46.1	81.5	74.8 <sub>(10th)</sub>	88.3 <sub>(11th)</sub>
Mistral-v0.2 (7B)	32K	16K	93.6	91.2	87.2	75.4	49.0	13.8	68.4	55.6 <sub>(13th)</sub>	81.2 <sub>(12th)</sub>
LWM (7B)	1M	<4K	82.3	78.4	73.7	69.1	68.1	65.0	72.8	69.9 <sub>(12th)</sub>	75.7 <sub>(13th)</sub>
DBRX (36B/132B)	32K	8K	95.1	93.8	83.6	63.1	2.4	0.0	56.3	38.0 <sub>(14th)</sub>	74.7 <sub>(14th)</sub>
Together (7B)	32K	4K	88.2	81.1	69.4	63.0	0.0	0.0	50.3	33.8 <sub>(15th)</sub>	66.7 <sub>(15th)</sub>
LongChat (7B)	32K	<4K	84.7	79.9	70.8	59.3	0.0	0.0	49.1	33.1 <sub>(16th)</sub>	65.2 <sub>(16th)</sub>
LongAlpaca (13B)	32K	<4K	60.6	57.0	56.6	43.6	0.0	0.0	36.3	24.7 <sub>(17th)</sub>	47.9 <sub>(17th)</sub>

# Evaluation

# Framework



# Evaluation

Gold Standard

# Benchmark

Category	ca	es	eu	gl	pt
Commonsense Reasoning	copa_ca <b>xstorycloze_ca</b>	<b>copa_es</b> xstorycloze_es	<b>xcopa_eu</b> xstorycloze_eu		
Linguistic Acceptability	catcola	escola		galcola	
Math	<b>mgsm_direct_ca</b>	mgsm_direct_es	<b>mgsm_direct_eu</b>	<b>mgsm_direct_gl</b>	
NLI	xnli_ca wnli_ca teca	xnli_es wnli_es	<b>xnli_eu</b> <b>wnli_eu</b> qnli_eu		assin_entailment
Paraphrasing	parafraseja <b>paws_ca</b> arc_ca catalanqa	paws_es		<b>parafrases_gl</b> <b>paws_gl</b>	assin_paraphrase
QA	coqcat <b>openbookqa_ca</b> <b>piqa_ca</b> <b>siqa_ca</b> xquad_ca	xquad_es <b>openbookqa_es</b>	<b>piqa_eu</b> eus_exams eus_proficiency eus_trivia	<b>openbookqa_gl</b>	
Reading Comprehension	belebele_cat_Latn	belebele_spa_Latn	belebele_eus_Latn eus_reading	<b>belebele_glg_Latn</b>	belebele_por_Latn
Summarization	cabreu	xlsum_es		<b>summarization_gl</b>	
Translation / Adaptation	<b>flores_ca</b> <b>phrases_va</b>	<b>flores_es</b> <b>phrases_es</b>	flores_eu	flores_gl	flores_pt
Truthfulness	veritasqa_ca	veritasqa_es		<b>veritasqa_gl</b> <b>truthfulqa_gl</b>	

Table 1: Tasks included in IberoBench. Newly introduced datasets are marked in bold.



*IberoBench: A Benchmark for LLM Evaluation in Iberian Languages*

# Gold Standard

## Odd labels in datasets

Subset (2)		Split (3)	
ARC-Challenge · 2.59k rows		test · 1.17k rows	
<input type="text"/> Search this dataset			
<b>id</b> string · lengths	<b>question</b> string · lengths	<b>choices</b> sequence	<b>answerKey</b> string · classes
 8 — 22	 13 — 831	 8 values	
MCAS_2010_8_12016	Company X makes 100 custom buses each year. Company Y makes 10,000 of one type of bus each year. Which of the following is the most likely reason a customer...	{ "text": [ "to keep the cost of the bus low", "to ensure that the bus will be easy to replace", "to provide ideas about how the bus will be built", "to...	C
Mercury_SC_400324	Students go on a class field trip to collect insects. Which of these is an important safety rule for the students to follow?	{ "text": [ "Pick up insects with bare hands.", "Stay with a classmate at all times.", "Reach under large rocks without looking first.", "Work away from th...	B
Mercury_SC_LBS10662	A hot rock is dropped into a pail of cool water. Heat energy transferred from the rock to the water by	{ "text": [ "boiling.", "evaporation.", "conduction.", "radiation." ], "label": [ "A", "B", "C", "D" ] }	C
VASOL_2009_3_8	Which of these is needed in all stages of the butterfly's life cycle?	{ "text": [ "Wings", "Eyes", "Soil", "Air" ], "label": [ "A", "B", "C", "D" ] }	D
Mercury_SC_401185	Trail mix is a type of snack. One kind of trail mix is made of raisins, chocolate pieces, peanuts, and sunflower seeds. Which statement describes why...	{ "text": [ "The mix contains four ingredients.", "The different ingredients clump together.", "Each component maintains its original properties.",...	C
NYSEDREGENTS_2015_8_29	Which change is the best example of a physical change?	{ "text": [ "a cookie baking", "paper burning", "ice cream melting", "a nail rusting" ], "label": [ "1", "2", "3", "4" ] }	3
Mercury_7234378	At which type of margin does mountain building produce mountains composed entirely of preexisting crust material?	{ "text": [ "divergent boundary between two oceanic plates", "divergent boundary between an oceanic plate and a continental plate", "convergent...	D

# Gold Standard

## Errors in datasets

```
{  
    "id": "CSZ_2009_8_CSZ30651",  
    "question": {  
        "stem": "Which of the following compounds is most likely to be part of living organisms?",  
        "choices": [  
            {"text": "C2H5O?", "label": "A"},  
            {"text": "BF?", "label": "B"},  
            {"text": "MoCl?", "label": "C"},  
            {"text": "CsI", "label": "D"}]  
    },  
    "answerKey": "A"  
}  
  
{  
    "id": "Mercury_400198",  
    "question": {  
        "stem": "What is the dependent variable in this experiment?",  
        "choices": [  
            {"text": "type of soil used", "label": "A"},  
            {"text": "amount of water used", "label": "B"},  
            {"text": "length of sunlight exposure", "label": "C"},  
            {"text": "growth of the bean plants", "label": "D"}]  
    },  
    "answerKey": "D"  
}
```

# Gold Standard

Lack of prompt-awareness and pre-processing

[lm-evaluation-harness / lm\\_eval / tasks / xnli / xnli\\_en.yaml](#) ↗



lintangsutawika update ✘

Code

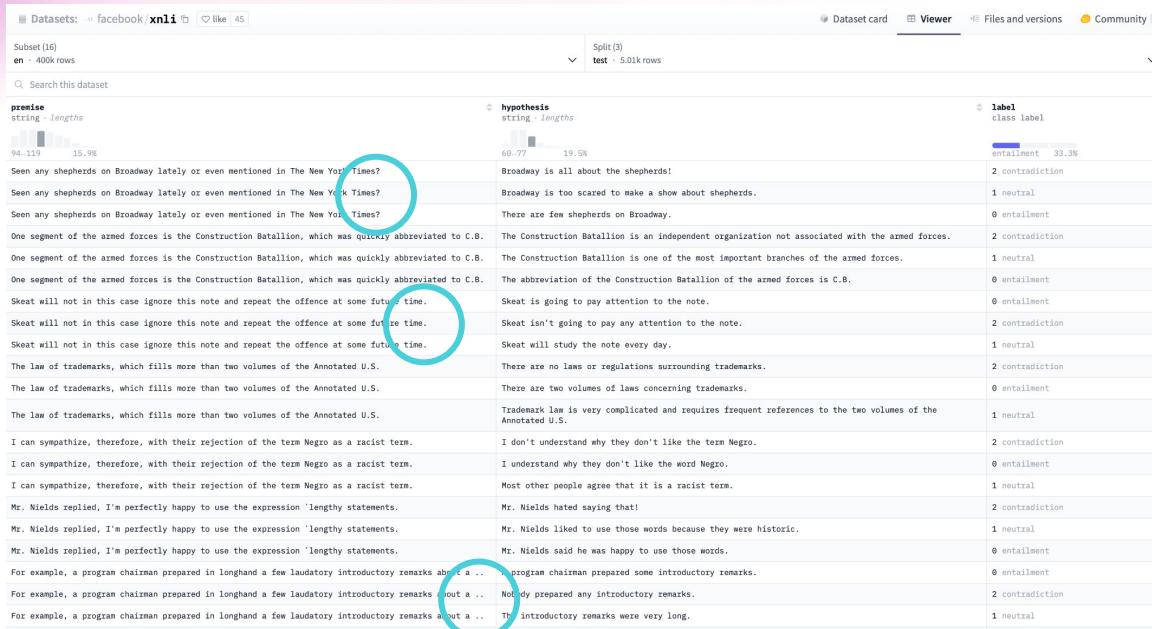
Blame

7 lines (7 loc) · 230 Bytes ·

```
1 # Generated by utils.py
2 dataset_name: en
3 doc_to_choice: '{{[premise+, right? Yes, "+hypothesis,premise+", right? Also, "+hypothesis,premise+",
4     right? No, "+hypothesis"]}}'
5 doc_to_text: ''
6 include: xnli_common_yaml
7 task: xnli_en
```

# Gold Standard

## Lack of prompt-awareness and pre-processing



Subset (16)		Split (3)	Dataset card	Viewer	Files and versions	Community
en	~400k rows	test · 5.01k rows				
Q. Search this dataset						
premise	hypothesis					
string lengths	string lengths					
94 - 119	60 - 77	15.9%	19.5%			
Seen any shephards on Broadway lately or even mentioned in The New York Times?	Broadway is all about the shephards!			label		
Seen any shephards on Broadway lately or even mentioned in The New York Times?	Broadway is too scared to make a show about shephards.			class label		
Seen any shephards on Broadway lately or even mentioned in The New York Times?	There are few shephards on Broadway.			entailment	33.3%	
One segment of the armed forces is the Construction Battallion, which was quickly abbreviated to C.B.	The Construction Battallion is an independent organization not associated with the armed forces.			contradiction		
One segment of the armed forces is the Construction Battallion, which was quickly abbreviated to C.B.	The Construction Battallion is one of the most important branches of the armed forces.			neutral		
One segment of the armed forces is the Construction Battallion, which was quickly abbreviated to C.B.	The abbreviation of the Construction Battallion of the armed forces is C.B.			entailment		
Skeat will not in this case ignore this note and repeat the offence at some future time.	Skeat is going to pay attention to the note.			contradiction		
Skeat will not in this case ignore this note and repeat the offence at some future time.	Skeat isn't going to pay any attention to the note.			neutral		
Skeat will not in this case ignore this note and repeat the offence at some future time.	Skeat will study the note every day.			entailment		
The law of trademarks, which fills more than two volumes of the Annotated U.S.	There are no laws or regulations surrounding trademarks.			contradiction		
The law of trademarks, which fills more than two volumes of the Annotated U.S.	There are two volumes of law concerning trademarks.			neutral		
The law of trademarks, which fills more than two volumes of the Annotated U.S.	Trademark law is very complicated and requires frequent references to the two volumes of the Annotated U.S.			entailment		
I can sympathize, therefore, with their rejection of the term Negro as a racist term.	I don't understand why they don't like the term Negro.			contradiction		
I can sympathize, therefore, with their rejection of the term Negro as a racist term.	I understand why they don't like the word Negro.			neutral		
I can sympathize, therefore, with their rejection of the term Negro as a racist term.	Most other people agree that it is a racist term.			entailment		
Mr. Nields replied, I'm perfectly happy to use the expression 'lengthy statements.'	Mr. Nields hated saying that!			contradiction		
Mr. Nields replied, I'm perfectly happy to use the expression 'lengthy statements.'	Mr. Nields liked to use those words because they were historic.			neutral		
Mr. Nields replied, I'm perfectly happy to use the expression 'lengthy statements.'	Mr. Nields said he was happy to use those words.			entailment		
For example, a program chairman prepared in longhand a few laudatory introductory remarks about a ...	A program chairman prepared some introductory remarks.			neutral		
For example, a program chairman prepared in longhand a few laudatory introductory remarks about a ...	Noody prepared any introductory remarks.			contradiction		
For example, a program chairman prepared in longhand a few laudatory introductory remarks about a ...	The introductory remarks were very long.			neutral		

# Gold Standard

```
'{{[premise+", right? Yes, "+hypothesis,premise+", right? Also,  
"+hypothesis,premise+", right? No, "+hypothesis]}}'
```

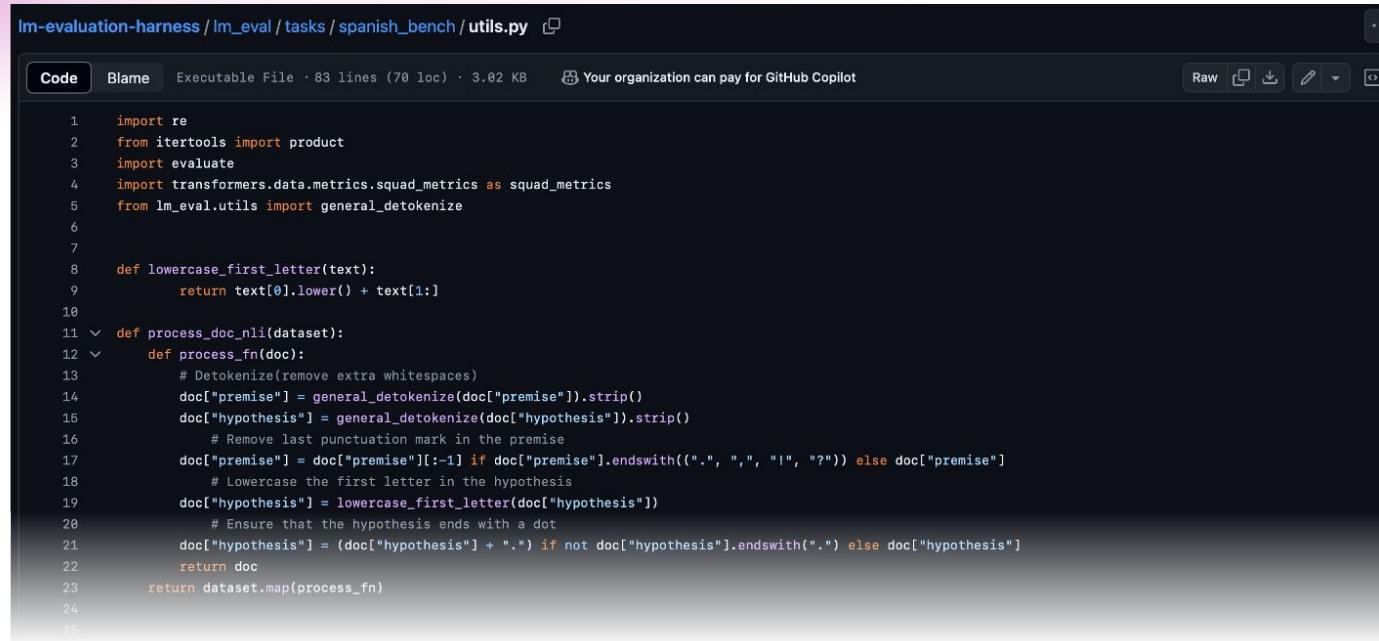
“Seen any shepherds on Broadway lately or even mentioned in The New York Times?, right? Yes, There are few shepherds on Broadway.”

“For example, a program chairman prepared in longhand a few laudatory introductory remarks about a ... right? Also, The introductory remarks were very long.”

“Pulse-tone is not a technical term., right? No, The official technical manual states that pulse-tone is the correct term in this case.”

# Gold Standard

Pre-processing is a must for uncontrolled datasets.



The screenshot shows a GitHub code editor interface for a file named `utils.py` located in the `lm-evaluation-harness / lm_eval / tasks / spanish_bench` directory. The file contains 83 lines of Python code, with 70 lines of actual code and 3 blank lines. The file size is 3.82 KB. The code implements various preprocessing functions, including `lowercase_first_letter` and `process_doc_nli`, which handle punctuation and whitespace in NLP datasets.

```
1 import re
2 from itertools import product
3 import evaluate
4 import transformers.data.metrics.squad_metrics as squad_metrics
5 from lm_eval.utils import general_detokenize
6
7
8 def lowercase_first_letter(text):
9     return text[0].lower() + text[1:]
10
11 def process_doc_nli(dataset):
12     def process_fn(doc):
13         # Detokenize(remove extra whitespaces)
14         doc["premise"] = general_detokenize(doc["premise"]).strip()
15         doc["hypothesis"] = general_detokenize(doc["hypothesis"]).strip()
16         # Remove last punctuation mark in the premise
17         doc["premise"] = doc["premise"][:-1] if doc["premise"].endswith(".", ",","!","?") else doc["premise"]
18         # Lowercase the first letter in the hypothesis
19         doc["hypothesis"] = lowercase_first_letter(doc["hypothesis"])
20         # Ensure that the hypothesis ends with a dot
21         doc["hypothesis"] = (doc["hypothesis"] + ".") if not doc["hypothesis"].endswith(".") else doc["hypothesis"]
22     return doc
23
24 return dataset.map(process_fn)
```

# Gold Standard

## Non-UTF-8 characters in prompts

lm-evaluation-harness / lm\_eval / tasks / belebele / \_default\_template.yaml

jmichaelov fixed belebele (#1267) · 9b0b15b · 5 months ago · History

Code Blame 19 lines (19 loc) · 610 Bytes · Raw

```
group: belebele
dataset_path: facebook/belebele
fewshot_config:
  sampler: first_n
output_type: multiple_choice
should_decontaminate: true
doc_to_decontamination_query: "{{question}}"
doc_to_text: "P: {{flores_passage}}\nQ: {{question.strip()}}\nA: {{mc_answer1}}\nB: {{mc_answer2}}\nC: {{mc_answer3}}\nD: {{mc_answer4}}\nAnswer: "
doc_to_choices: ["A", "B", "C", "D"]
doc_to_target: "[{{'1', '2', '3', '4'}}.index(correct_answer_num)]"
metric_list:
  - metric: acc
    aggregation: mean
    higher_is_better: true
  - metric: acc_norm
    aggregation: mean
    higher_is_better: true
metadata:
  version: 0.0
```

# Gold Standard

Non-UTF-8 characters in prompts

```
 {{mc_answer2}}\nC: {{mc_answer3}}\nD: {{mc_answer4}}\nAnswer: "
```

```
>>> flor6b_tokenizer.tokenize(": ")
['í', '¼', 'ł']
```

“Binary digits are also referred to as what? A: Bits (...) D: Formsí¼|”

# Gold Standard

Type of task	Metric
Commonsense Reasoning	Accuracy
Linguistic Acceptability	Matthews Correlation Coefficient
Math	Exact Match
Natural Language Inference	Accuracy
Paraphrasing	F1 Score
Question Answering	Acc. for binary or MC tasks + F1 for generative tasks
Reading Comprehension	Accuracy
Summarization	ROUGE1 or BLEU
Translation	BLEU
Truthfulness	Accuracy + BLEU

# Gold Standard

Differences in runs using tensor parallelism.

Run Name	Created	Duration	Experiment Name	score	model
● arc_challenge	✓ 16 hours ago	117ms	fourth_epoch_mix1_no-tp	49.66	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● arc_challenge	✓ 1 day ago	30ms	fourth_epoch_mix1	50.34	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
Run Name	Created	Duration	Experiment Name	score	model
● teca	✓ 21 hours ago	117ms	fourth_epoch_mix1_no-tp	51.72	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● teca	✓ 1 day ago	27ms	fourth_epoch_mix1	48.84	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
Run Name	Created	Duration	Experiment Name	score	model
● cabreu_extreme	✓ 20 hours ago	151ms	fourth_epoch_mix1_no-tp	41.94	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● cabreu_extreme	✓ 1 day ago	38ms	fourth_epoch_mix1	40.72	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
Run Name	Created	Duration	Experiment Name	score	model
● catalanqa	✓ 19 hours ago	118ms	fourth_epoch_mix1_no-tp	82.23	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● catalanqa	✓ 20 hours ago	34ms	fourth_epoch_mix1	81.7	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo

# Gold Standard

Differences in runs using tensor parallelism.

Run Name	Created	Duration	Experiment Name	score	model
● catalanqa	✓ 19 hours ago	118ms	fourth_epoch_mix1_no-tp	82.23	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● catalanqa	✓ 20 hours ago	34ms	fourth_epoch_mix1	81.7	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo

Run Name	Created	Duration	Experiment Name	score	model
● catalanqa	✓ 19 hours ago	118ms	fourth_epoch_mix1_no-tp	82.23	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● catalanqa	✓ 20 hours ago	34ms	fourth_epoch_mix1	81.7	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● catalanqa	✓ 20 hours ago	39ms	fourth_epoch_mix1	81.7	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● catalanqa	✓ 20 hours ago	36ms	fourth_epoch_mix1	81.7	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo
● catalanqa	✓ 20 hours ago	43ms	fourth_epoch_mix1	81.7	fourth_epoch_bsc_7b_restart_lr3e-5_lr3e-6_step52500_decompressed_nemo

# Gold Standard

## + more issues:

Double or no spaces when concatenating in prompt.

Limitations with metrics for open-ended generative tasks.

Multiple gold standards for some tasks.

Errors in datasets (both content and form).

1% variability across evaluations with same setup in a SingularityCE.

First-token generation bias.

# Evaluation

LLM-as-a-Judge

# LLM-as-a-judge

Prompts and criteria for LLM-Evaluator need to be in English.

```
38     SYSTEM_MESSAGE_PROMETHEUS_2 = "You are a fair judge assistant tasked with providing clear, objective feedback based on \
39     specific criteria, ensuring each assessment reflects the absolute standards set for performance."
40 ✓ SCORING_TEMPLATE_WITH_REFERENCE_1_TO_5 = ChatPromptTemplate.from_messages(
41     [
42         ("system", SYSTEM_MESSAGE_PROMETHEUS_2),
43         (
44             "human",
45             "###Task Description:\nAn instruction (might include an Input inside it), a response to evaluate, and a score rubric representing a \
46             evaluation criteria are given.\n\
47             \n1. Write detailed feedback that assess the quality of the response strictly based on the given score rubric, not evaluating in general.\n\
48             \n2. After writing feedback, write a score that is an integer between 1 and 5. You should refer to the score rubric.\n\
49             \n3. The output format should look as follows: 'Feedback: (write feedback for criteria) [RESULT] (an integer number between 1 and 5)'\n\
50             \n4. Please do not generate any other opening, closing, and explanations.\n\
51             \n###The instruction to evaluate: {input}\n\
52             \n###Response to evaluate: {prediction}\n\
53             \n###Score Rubrics:\n{criteria}\n\
54             \n###Feedback:"
55         ),
56     ],
57 )
```

# LLM-as-a-judge

Prompts and criteria for LLM-Evaluator need to be in English.

```
1  ✓  conciseness_criteria = {  
2      "concreteness": """  
3      Score 1: The answer contains too much information and includes unnecessary details that could confuse the reader or detract from the main point.  
4      Score 2: The answer is overly long-winded and could be condensed to convey the same information more efficiently.  
5      Score 3: The answer is generally brief and to the point, but there are still some words or phrases that could be removed to improve clarity and brevity.  
6      Score 4: The answer is clear, succinct, and effectively communicates the necessary information without unnecessary elaboration.  
7      Score 5: The answer is exceptionally brief, yet it effectively conveys all the required information with clarity and precision."""  
8  }  
9  
10 ✓ relevance_criteria = {  
11     "relevance": """  
12     Score 1: The answer is not relevant to the question and does not address the topic at hand.  
13     Score 2: The answer is somewhat relevant but lacks focus and may include tangential information.  
14     Score 3: The answer is relevant but could be more focused on addressing the specific question or topic.  
15     Score 4: The answer is highly relevant and directly addresses the question or topic with appropriate detail.  
16     Score 5: The answer is perfectly relevant, providing precise and comprehensive information directly related to the question or topic."""}  
17  
18 ✓ correctness_criteria = {  
19     "correctness": """  
20     Score 1: The answer contains significant inaccuracies or errors that mislead the reader.  
21     Score 2: The answer has several inaccuracies or errors that affect its credibility and may confuse the reader.  
22     Score 3: The answer is generally correct but contains minor inaccuracies or errors that could be corrected for improved accuracy.  
23     Score 4: The answer is mostly correct, with few inaccuracies or errors that do not significantly impact its overall correctness.  
24     Score 5: The answer is entirely correct, providing accurate and reliable information without any inaccuracies or errors."""  
25  }
```

# LLM-as-a-judge

Beware of the FT prompt and the output parser.

```
1  from langchain_core.prompts.chat import ChatPromptTemplate
2
3  SYSTEM_MESSAGE = "You are a helpful assistant."
4  SCORING_TEMPLATE_WITH_REFERENCE_1_TO_5_a = ChatPromptTemplate.from_messages(
5      [
6          ("system", SYSTEM_MESSAGE),
7          (
8              "human",
9              "[Instruction]\nPlease act as an impartial judge \
10             and evaluate the quality of the response provided by an AI \
11             assistant to the user question displayed below. {criteria}"
12                 'Begin your evaluation \
13                 by providing a short explanation. Be as objective as possible. \
14                 After providing your explanation, you must rate the response on a scale of 1 to 5 \
15                 by strictly following this format: "[[rating]]", for example: "Rating: [[5]]".\n\n\
16                 [Question]\n{input}[Actual Verified Answer]\n{reference}\n\n[The Start of Assistant\'s Answer]\n{prediction}\n\n\
17                 [The End of Assistant\'s Answer].',
18             ),
19         ]
20     )
21 )
```

# LLM-as-a-judge

Beware of the FT prompt and the output parser.

```
37     SYSTEM_MESSAGE_PROMETHEUS_2 = "You are a fair judge assistant tasked with providing clear, objective feedback based on \
38 specific criteria, ensuring each assessment reflects the absolute standards set for performance."
39 ✓ SCORING_TEMPLATE_WITH_REFERENCE_1_TO_5_CONTEXT = ChatPromptTemplate.from_messages(
40     [
41         ("system", SYSTEM_MESSAGE_PROMETHEUS_2),
42         (
43             "human",
44             "###Task Description:\nAn instruction (might include an Input inside it), a response to evaluate, and a score rubric representing a \
45 evaluation criteria are given.\n\
46 \n1. Write detailed feedback that assess the quality of the response strictly based on the given score rubric, not evaluating in general.\n\
47 \n2. After writing feedback, write a score that is an integer between 1 and 5. You should refer to the score rubric.\n\
48 \n3. The output format should look as follows: 'Feedback: (write feedback for criteria) [RESULT] (an integer number between 1 and 5)'\n\
49 \n4. Please do not generate any other opening, closing, and explanations.\n\
50     \n###The instruction to evaluate: {input}\n\
51     \n###Verified correct answer: {reference}\n\
52     \n###Response to evaluate: {prediction}\n\
53     \n###Score Rubrics:\n{criteria}\n\
54     \n###Feedback:"
55     ),
56     ]
57 )
```

# LLM-as-a-judge

To try to fix this, we initially changed point 3 of the prompt to:

3. The output format should look as follows: 'Feedback: (write feedback for criteria) [[RESULT]] (an integer number between 1 and 5, for example, [[5]])'

Which led to answers like:

I do like this answer very much. [[RESULT]] 5

# LLM-as-a-judge

Let's change it to:

3. The output format should look as follows: 'Feedback: (write feedback for criteria) [[RESULT GOES HERE]] (an integer number between 1 and 5, for example, [[5]])'

This led to answers like:

I do like this answer very much. [[RESULT]] [[5]]

# LLM-as-a-judge

```
 2     from langchain.schema.output_parser import BaseOutputParser
10    class PrometheusOutputParser(BaseOutputParser):
11
12        def __init__(self):
13            super().__init__()
14
15        def parse(self, output):
16            parts = output.split('[RESULT]')
17            assert len(parts) == 2
18            score = int(parts[-1].strip())
19
20            return {'score': score}
```

# LLM-as-a-judge

Prompt preference for each LLM-Evaluator.

Empty strings as best scores.

LLMs can distinguish themselves from other LLMs and humans, and there is a direct correlation between self-recognition capability and self-preference bias.

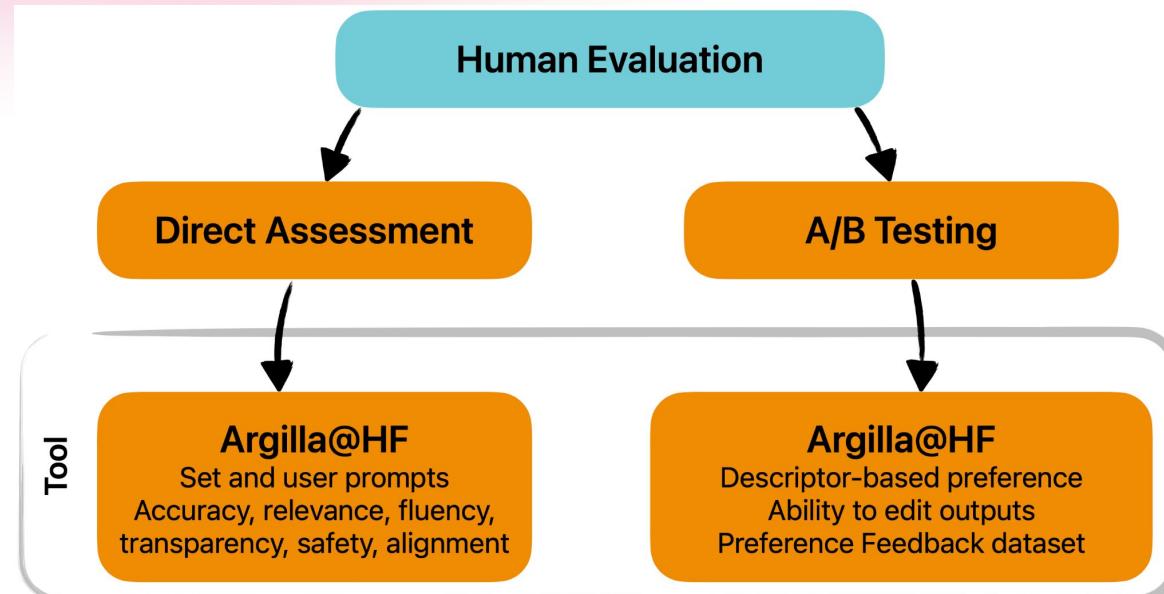
The “emergent ability” of Likert-scale granular understanding.

Reproducibility issues.

# Evaluation

## Human Evaluation

# Human Evaluation



# Human Evaluation

High cost and challenging scalability.

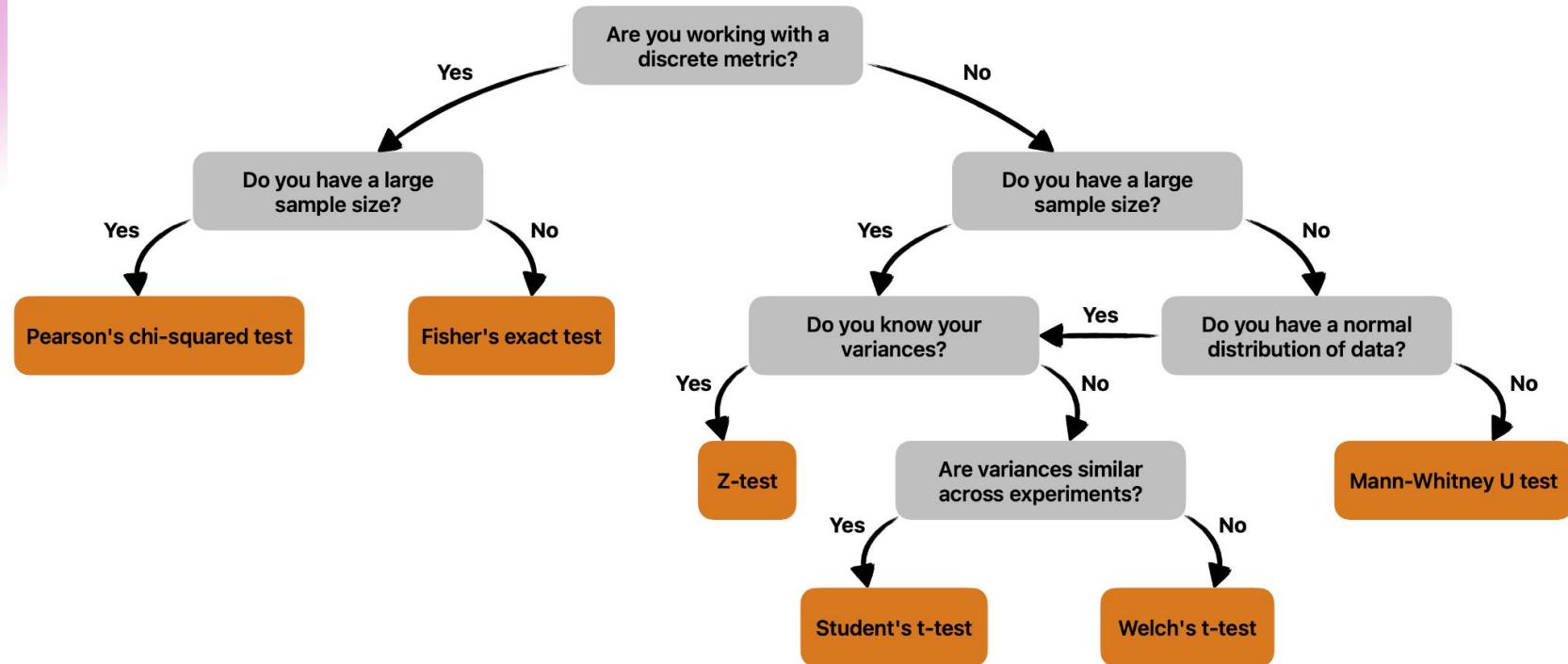
Granularity of criteria and scales.

Overcoming human cognitive biases.

Trusting annotators.

Measuring reliability and using statistics.

# Human Evaluation



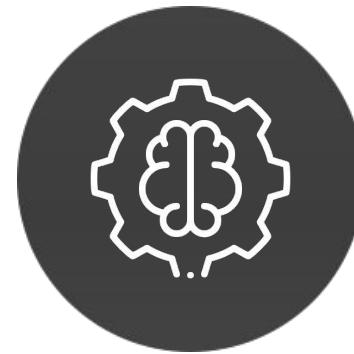
# Bias

# Concerns



## Social bias

Misrepresentation  
Stereotypes  
Derogatory language



## Cognitive bias

Primacy/recency bias  
Majority class bias  
Common token bias

# Challenges

Varying definitions of “bias”, “toxicity”, “harmful content” etc.

No clear metrics for automatic evaluation.

Most benchmarks built around USA socio-cultural context.

Human annotation is challenging and IAA tends to be low.

Keyword bias affects dataset curation and evaluation.

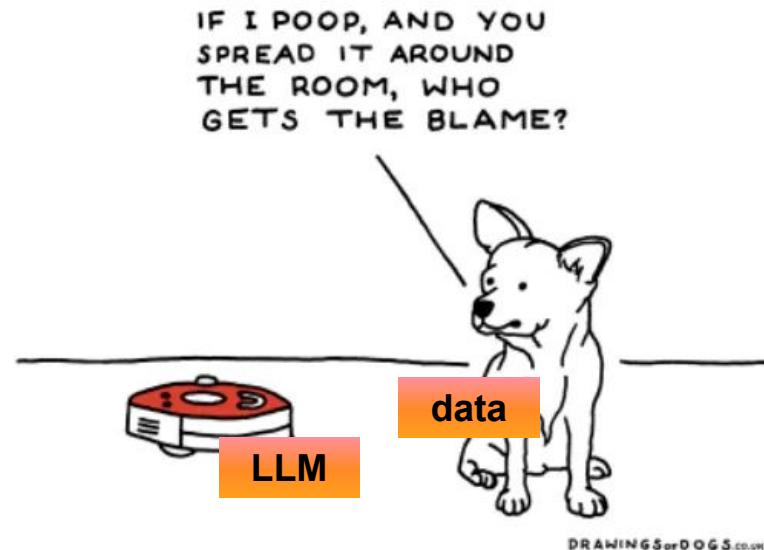
Full “debiasing” is impossible; focus on *mitigation*.

# Identification approaches

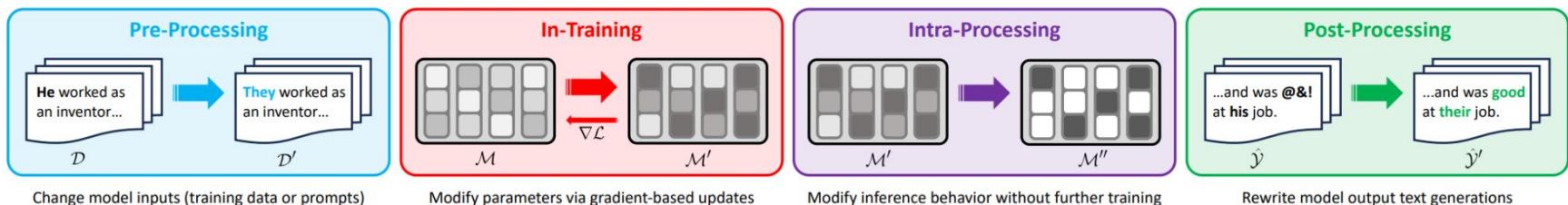
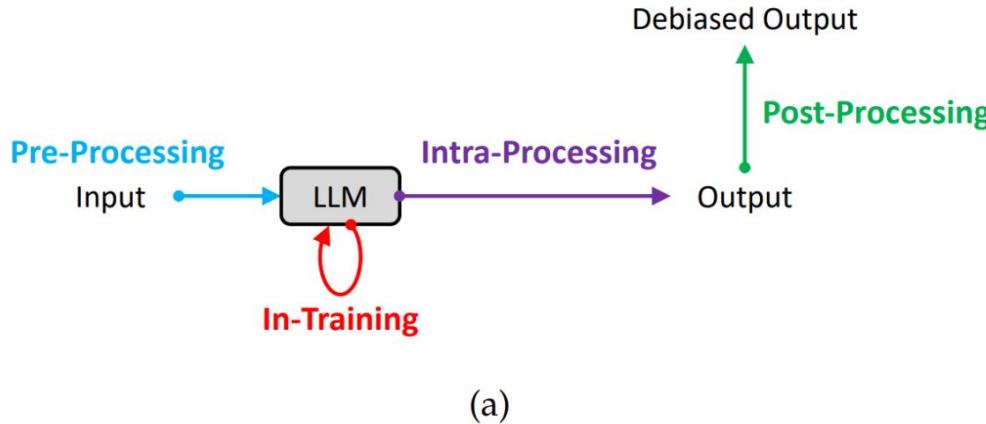
Bias can **come from the training data** and also be **modulated by** the LLM training/alignment **procedures**.

We work on:

- Creating, adapting and/or translating benchmarks.
- Researching new/better ways of measuring bias in our data and model structure.



# Mitigation approaches



# Post-training

# Instruction tuning

**Multilingual instructions** available:

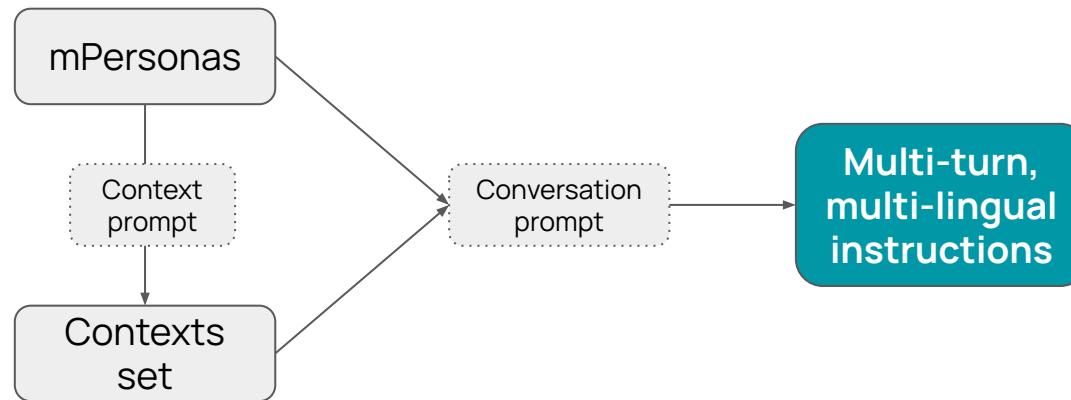
- Licensing issues.
- Domain, style and format limitations.
- Lack of context awareness.
- “Identity” theft.
- Circa 419k pairs for our pre-training languages (one year ago).

Use **synthetic data responsively** and **edit existing data**.

# Instruction tuning

Follow the **strategy from PersonaHub**, replacing tools and data subject to viral license with **open** ones.

Generated **personas follow the same language as the source document** being used, making them more culturally relevant for our context.



# Context prompt

You are a helpful assistant.

====Example 1====

Persona: A curious and analytical individual, likely with a background in mathematics or science, who enjoys exploring intriguing "what if" scenarios and is fascinated by the intersection of population demographics and geography.

Prompt: Is it possible for the global population to stand on Jeju Island?

====Example 2====

Persona: An astronomy enthusiast or a professional astronomer, likely with a strong interest in peculiar galaxy structures and a good understanding of celestial objects, seeking to gather specific information about the unique Hoag's object galaxy.

Prompt: Given the unique ring structure of Hoag's Object and its mysterious formation, what are the leading theories explaining the origin of the central core galaxy? Additionally, is there any evidence supporting the hypothesis that the ring may have formed through a gravitational lensing effect or past galactic interactions?

---

Your task: Guess a prompt (i.e., instruction) that the following person may ask you to do: <S REPL>persona<E REPL>  
Generate it in <S REPL>lang<E REPL>.  
Return only the prompt inside <prompt> and </prompt>.

# Conversation prompt

You are an AI assistant specialized in creating diverse but specific conversations between a user and an AI assistant.

Create a conversation with <S REPL>turns<E REPL> turn(s) based on the persona and prompt below:

Persona: <S REPL>persona<E REPL>

Prompt: <S REPL>prompt<E REPL>

Return the conversation in JSON format in <S REPL>lang<E REPL>. The format must be inside <conversation> and </conversation> and with the format:

```
<conversation>
{
  "lang": "en",
  "conversations": [
    {
      "from": "human",
      "value": "What is the capital of Spain?"
    },
    {
      "from": "gpt",
      "value": "The capital is Madrid."
    }
  ]
}</conversation>
```

# Persona examples

## Generic persona

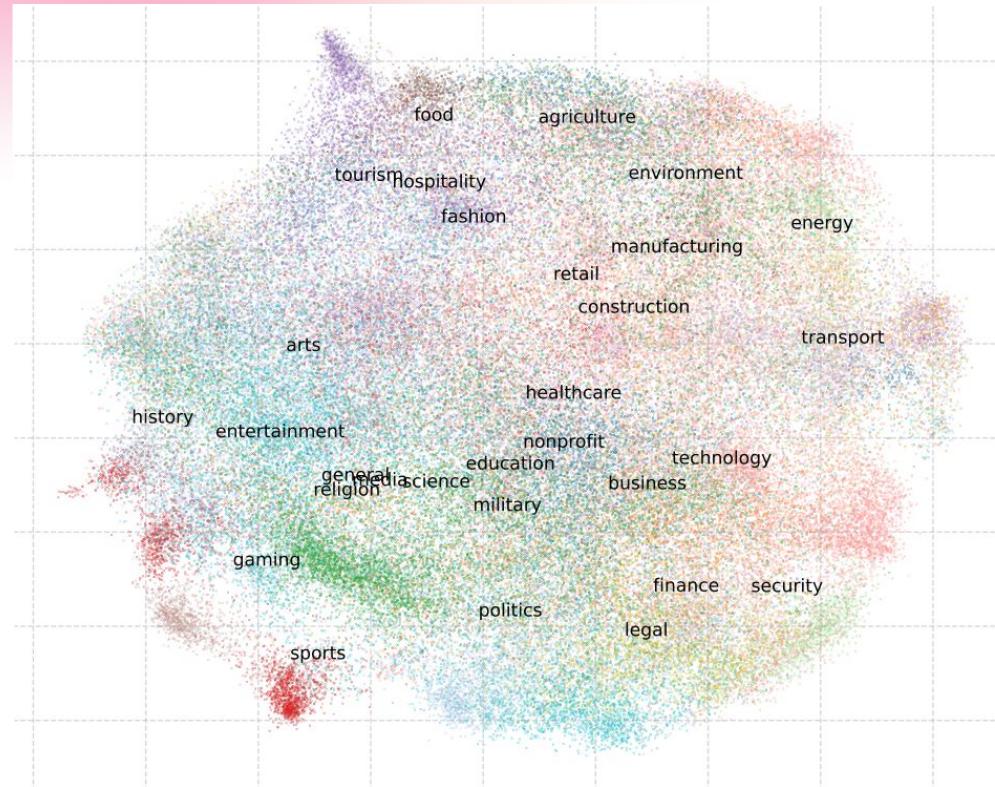
Un periodista especializado en economía y política europea, que sigue de cerca las crisis y tendencias en España y utiliza plataformas digitales como Orbyt para mantenerse informado sobre análisis de medios como The Economist.

## Specific persona

Un historiador especializado en cartografía antigua, con un profundo interés en la evolución de las técnicas de grabado y la producción de atlas entre los siglos XVI y XIX. Trabaja como investigador en el Departamento de Bellas Artes y Cartografía de la Biblioteca Nacional de España, donde estudia y cataloga colecciones de mapas raros, con especial enfoque en obras flamencas, holandesas y españolas. Domina el análisis de estilos artísticos en frontispicios alegóricos y valora la precisión histórica de mapas de exploración. Es experto en figuras como Abraham Ortelio, Gerard Mercator y Tomás López, y frecuentemente colabora con instituciones culturales para exposiciones sobre patrimonio cartográfico.

# Instruction tuning

Domains in the  
first sample of  
~126k personas



# Direct Preference Optimization

DPO is an optimization algorithm that takes both the chosen instruction-response pairs and their respective rejected pairs, as chosen by human evaluators, to minimize a preference loss.

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x, y_c, y_r) \sim \mathcal{D}_R} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_c \mid x)}{\pi_{\text{ref}}(y_c \mid x)} - \beta \log \frac{\pi_\theta(y_r \mid x)}{\pi_{\text{ref}}(y_r \mid x)} \right) \right]$$

# Proximal Policy Optimization

PPO is a two-step approach that first trains a separate reward model to emulate human preferences by minimizing a preference loss,

$$\mathcal{L}_R(\psi) = -\mathbb{E}_{(x, y_c, y_r) \sim \mathcal{D}_R} [\log \sigma(R_\psi(x, y_c) - R_\psi(x, y_r))]$$

and then uses this reward model to give a reward score to the outputs of the instructed model in order to align it with the desired behavior.

# DPO vs. PPO

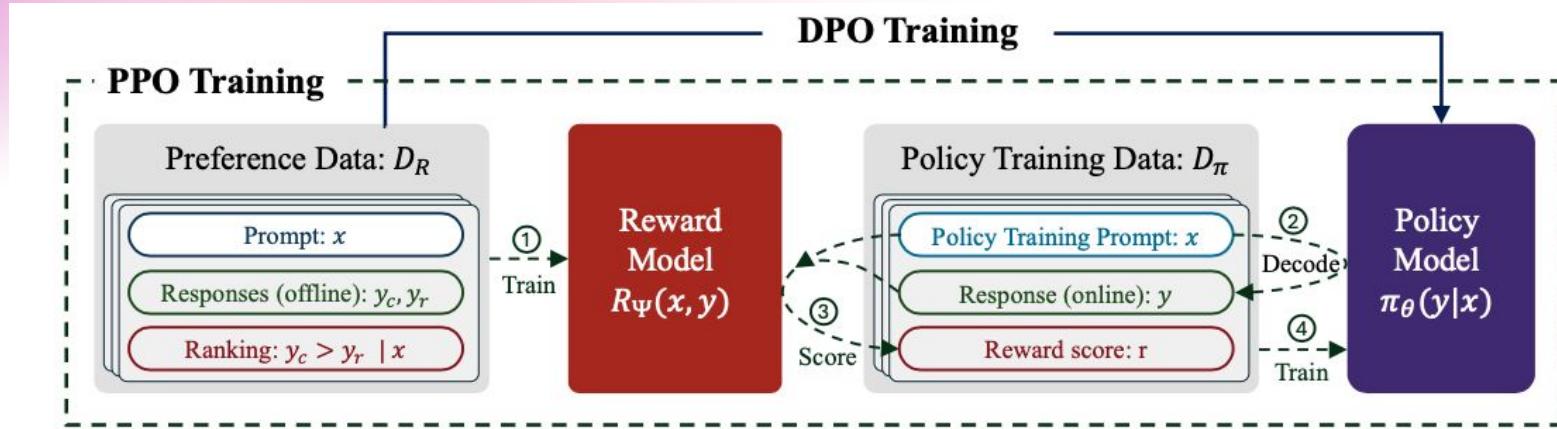


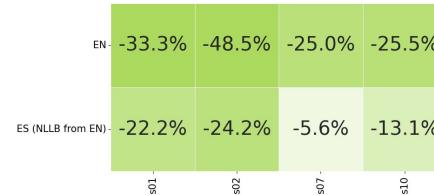
Figure 2: The core aspects of learning from preference feedback. For DPO (solid line), preference data is directly used to train a policy model. For PPO (dashed line), preference data is used to train a reward model, which is then used to score model-generated responses during PPO training.

# Safety

Pipeline **simulates malicious use**. Moderator model classifies the type of attack and if it was successful or not.

Focus on **four axes**: risk type, attack strategy, language, and number of turns.

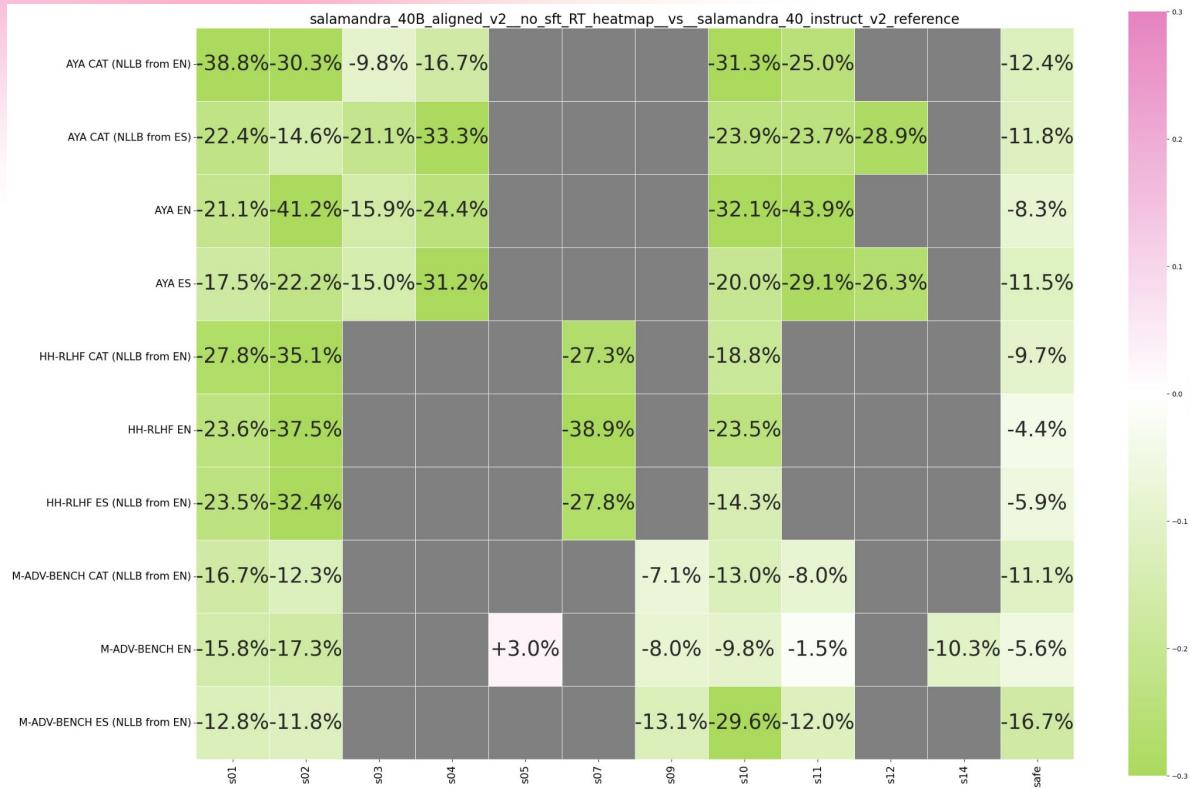
Basic RLHF alignment with **DPO improves model resistance** to attacks.



Use a combination of **synthetic data and human annotation** to expand to languages other than English.

Focus on **proven methodologies but using open-license data**.

# Safety



# Coming soon

# Coming soon

More **capable models** (all sizes) with **larger contexts** and **new languages**.

Refined **instructed and aligned** versions.

**Multimodal, multilingual models.**

Multimodal, multilingual **evaluation**.

**Multilingual agents and judges.**

Multilingual **judge evaluation**.

**Domain-specific** models.

In-house **quantization**.

+ more.

# Building and Evolving a Multilingual LLM from Scratch

Javier Aula-Blasco, PhD

/ CyberColombia  
/ 8th Colombian HPC Summer School  
/ June 17, 2025



# Acknowledgements



This work has been promoted and financed by the *Generalitat de Catalunya* through the **Aina** project.

This work is funded by the *Ministerio para la Transformación Digital y de la Función Pública* and *Plan de Recuperación, Transformación y Resiliencia* - Funded by EU – NextGenerationEU within the framework of the project **ILENIA** with reference 2022/TL22/00215337, 2022/TL22/00215336, 2022/TL22/00215335 and 2022/TL22/00215334, and within the framework of the project **Desarrollo Modelos ALIA**.



**Generalitat de Catalunya**  
Government  
of Catalonia



Financiado por  
la Unión Europea  
NextGenerationEU



**red.es**

 Plan de  
Recuperación,  
Transformación  
y Resiliencia

PERTE  
Nueva Economía  
De la Lengua



# Javier Aula-Blasco, PhD

*Head of Data and Model Evaluation*  
Language Technologies Lab

javier.aulablasco@bsc.es



 /javieraulablasco

 20 YEARS  
BSC

A large, stylized "20" is followed by the word "YEARS" in a bold, sans-serif font. The "20" has a circular graphic element integrated into its design, and the "BSC" logo is positioned at the bottom right of the "20".