

Text generation

Dirk Wulff & Zak Hussain



MAX PLANCK INSTITUTE
FOR HUMAN DEVELOPMENT



Science & Society

Can AI language models replace human participants?



Danica Dillion,¹ Niket Tandon,²
Yuling Gu,² and Kurt Gray ^{1,*,@}

Recent work suggests that language models such as GPT can make human-like judgments across a number of domains. We explore whether and when language models might replace human participants in psychological science. We review nascent research, provide a theoretical model, and outline caveats of using AI as a participant.

Does GPT make human-like judgments?

We initially doubted the ability of LLMs to capture human judgments but, as we detail in **Box 1**, the moral judgments of GPT-3.5 were extremely well aligned with human moral judgments in our analysis ($r = 0.95$; full details at <https://nikett.github.io/gpt-as-participant>). Human morality is often argued to be especially difficult for language models to capture [4] and yet we found powerful alignment between GPT-3.5 and human judgments.

We emphasize that this finding is just one anecdote and we do not make any strong claims about the extent to which LLMs make human-like judgments, moral or otherwise. Language models also might be especially good at predicting moral judg-

developed a framework (**Box 2**) that connects LLM responses to human cognition. The model emphasizes that the ‘minds’ of LLMs are grounded in naturalistic expression across a large but constrained group of people. Practically speaking, LLMs may be most useful as participants when studying specific topics, when using specific tasks, at specific research stages, and when simulating specific samples.

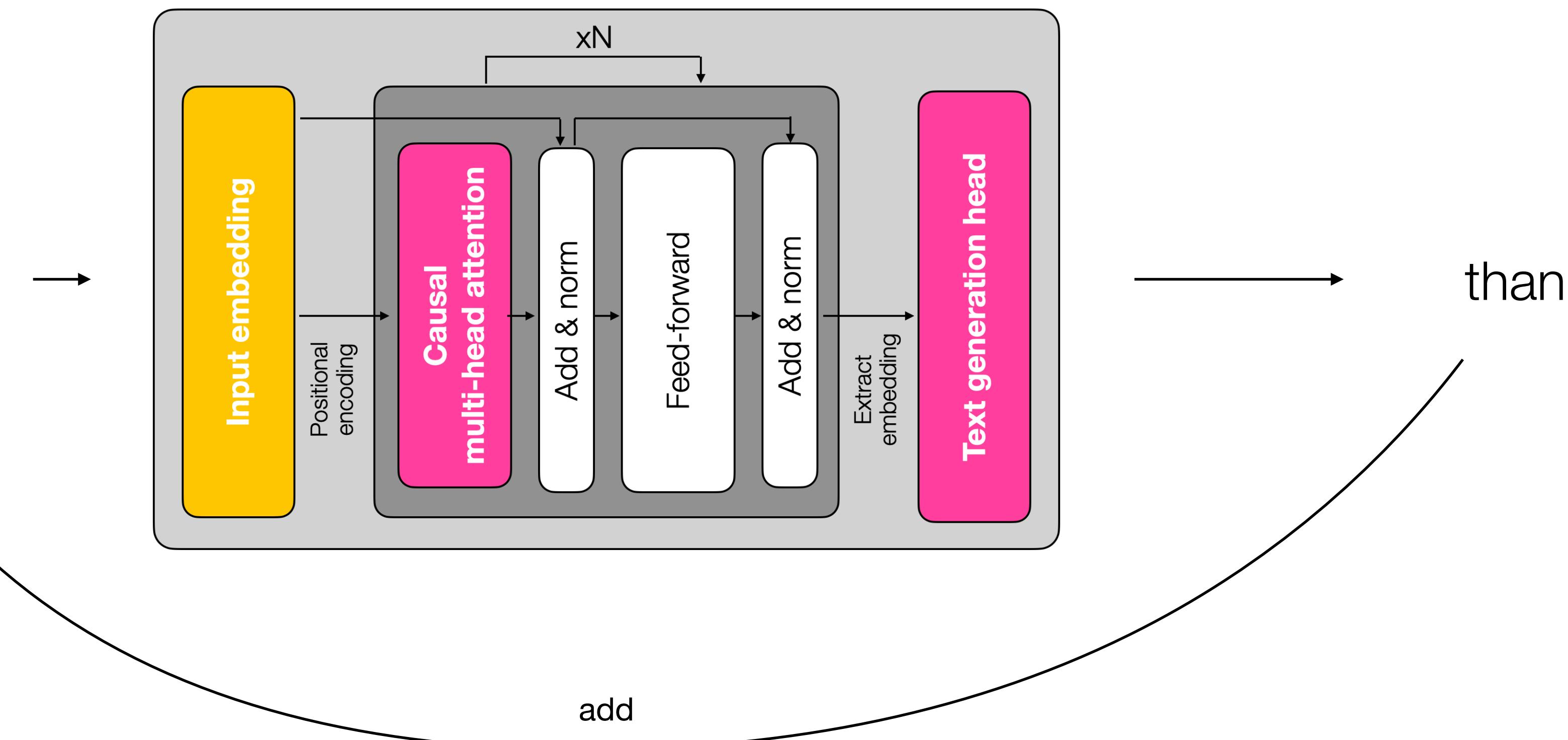
Specific topics

Language model expressions may be most correlated with human expressions when there are obvious explicit features of situations that drive human judgments. With morality, these might include whether an action was intentional or not. With mind perception, these might include whether a target is described as human or a kind of

Text generation

is autoregressive next-token prediction

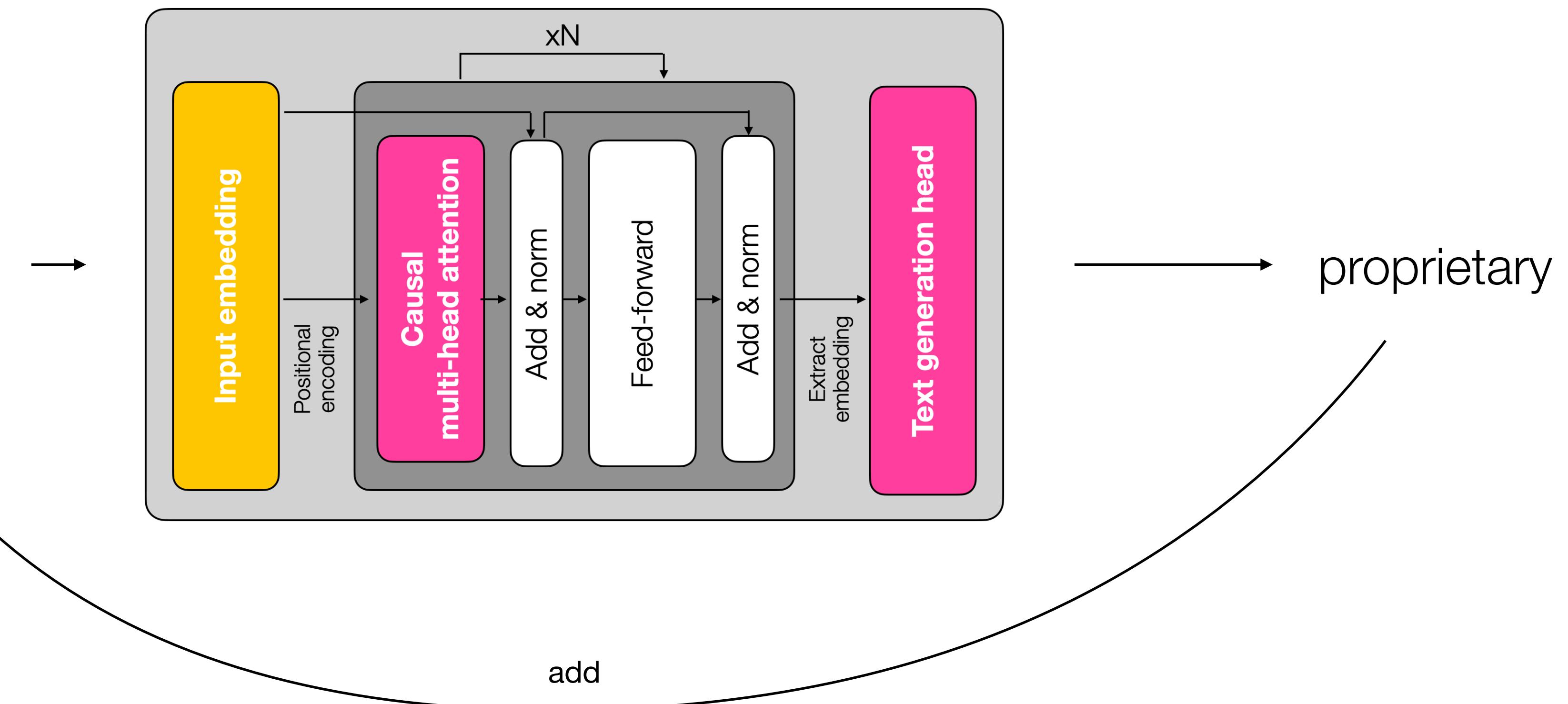
Open-source LLMs like
Phi are better



Text generation

is autoregressive next-token prediction

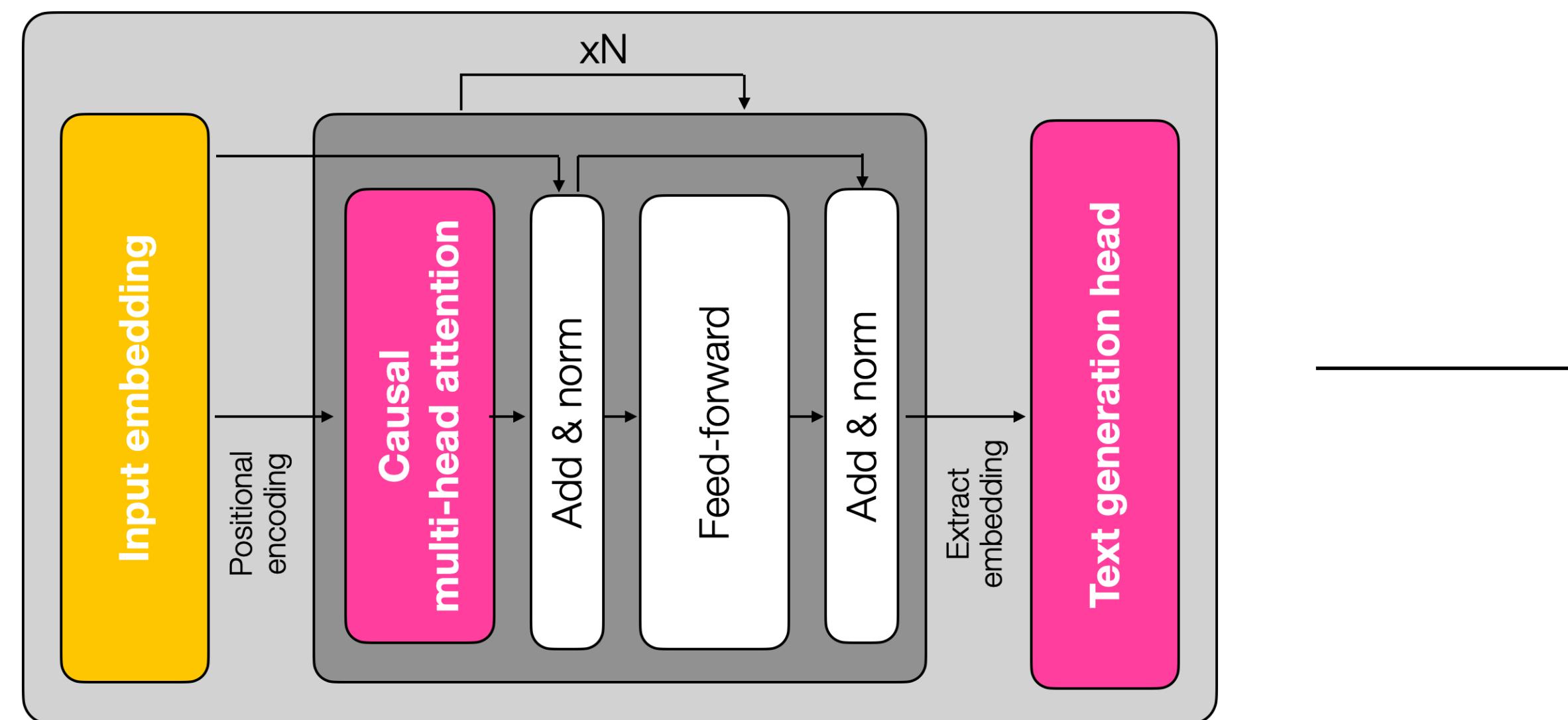
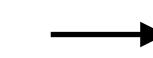
Open-source LLMs like
Phi are better than



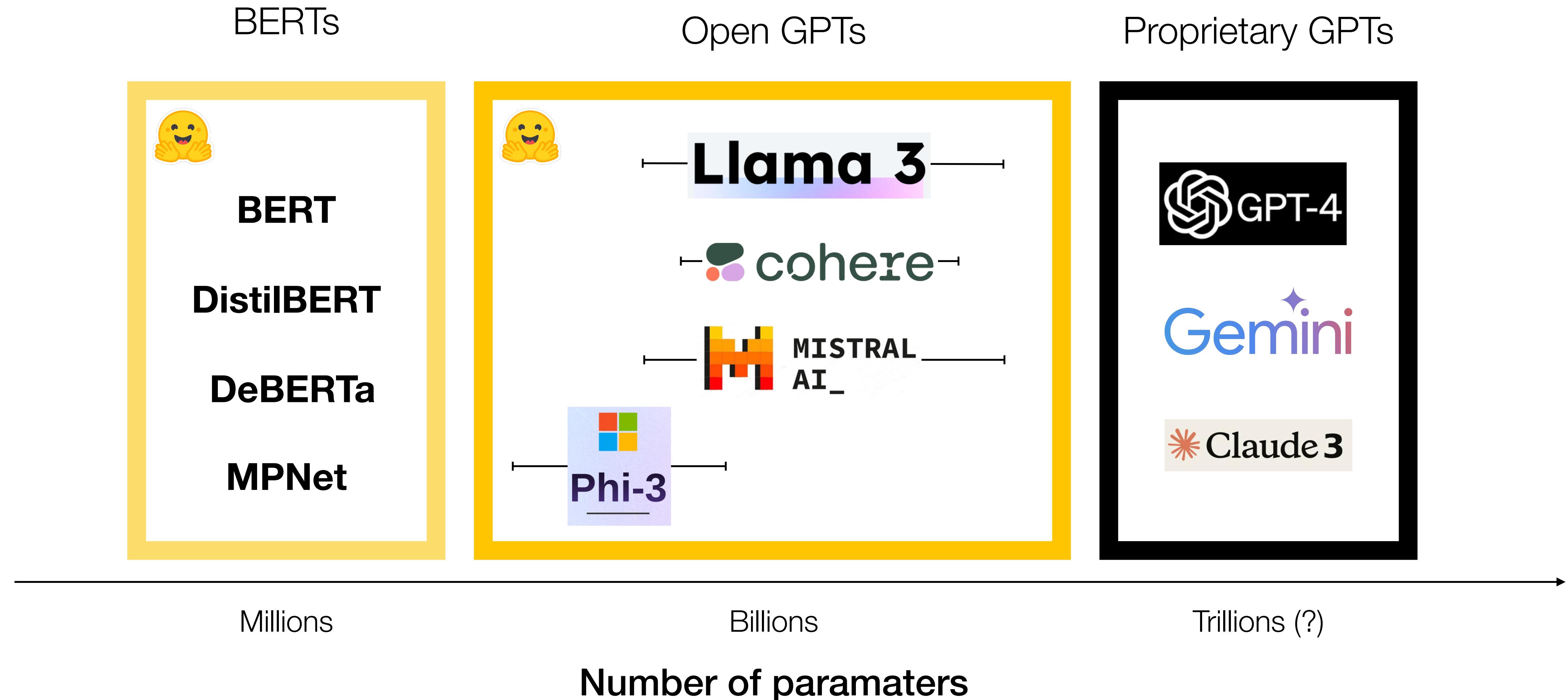
Text generation

is autoregressive next-token prediction

Open-source LLMs like
Phi are better **than**
proprietary

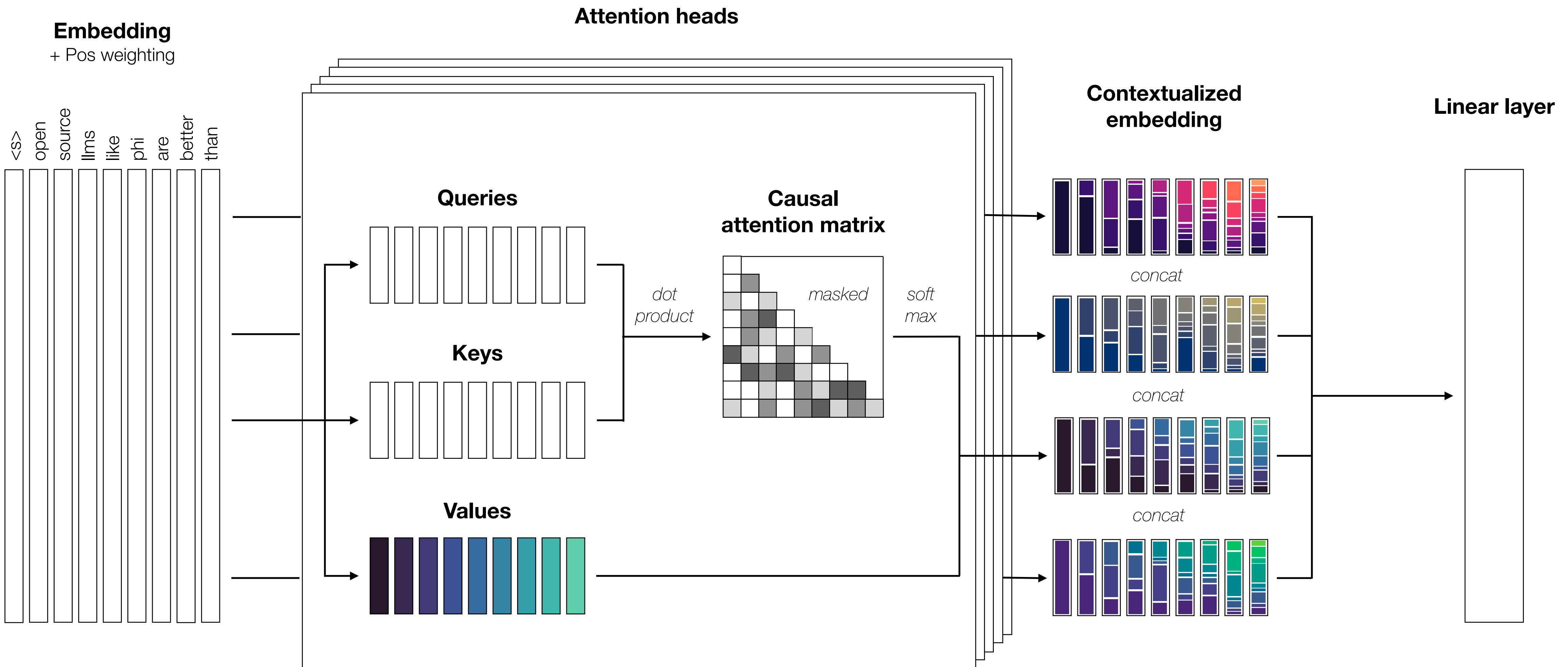


Models



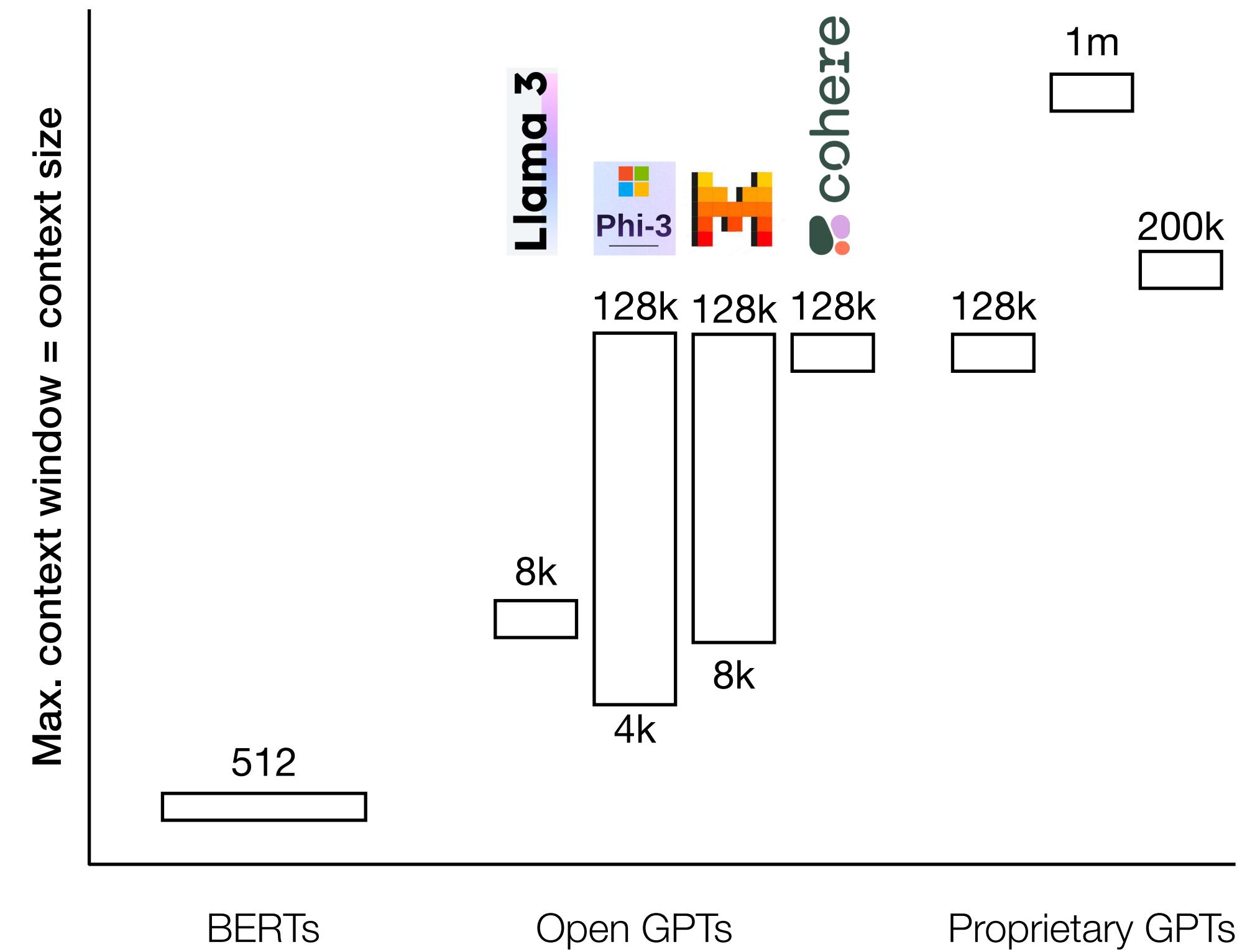
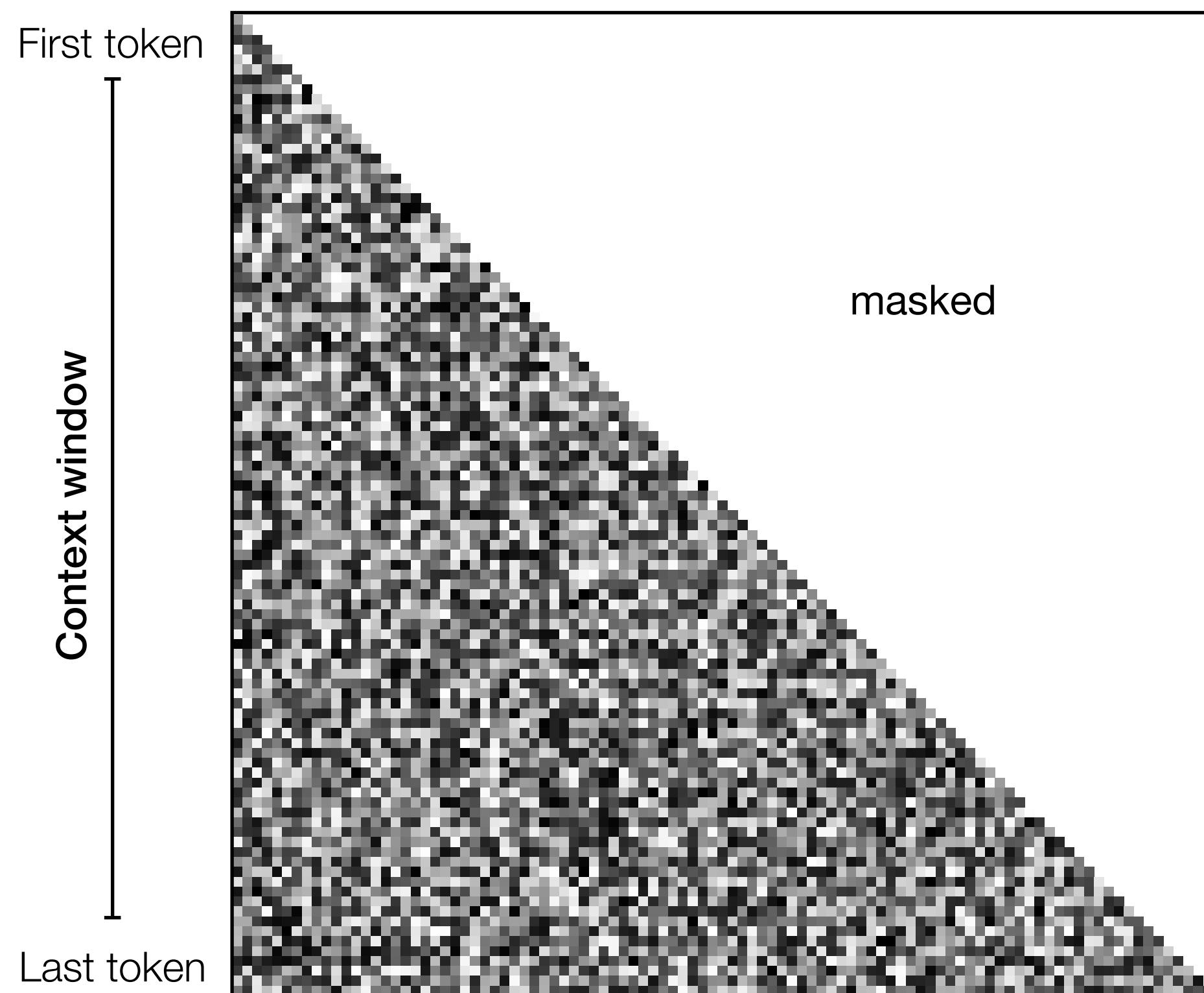
Transformer

Causal attention mechanism



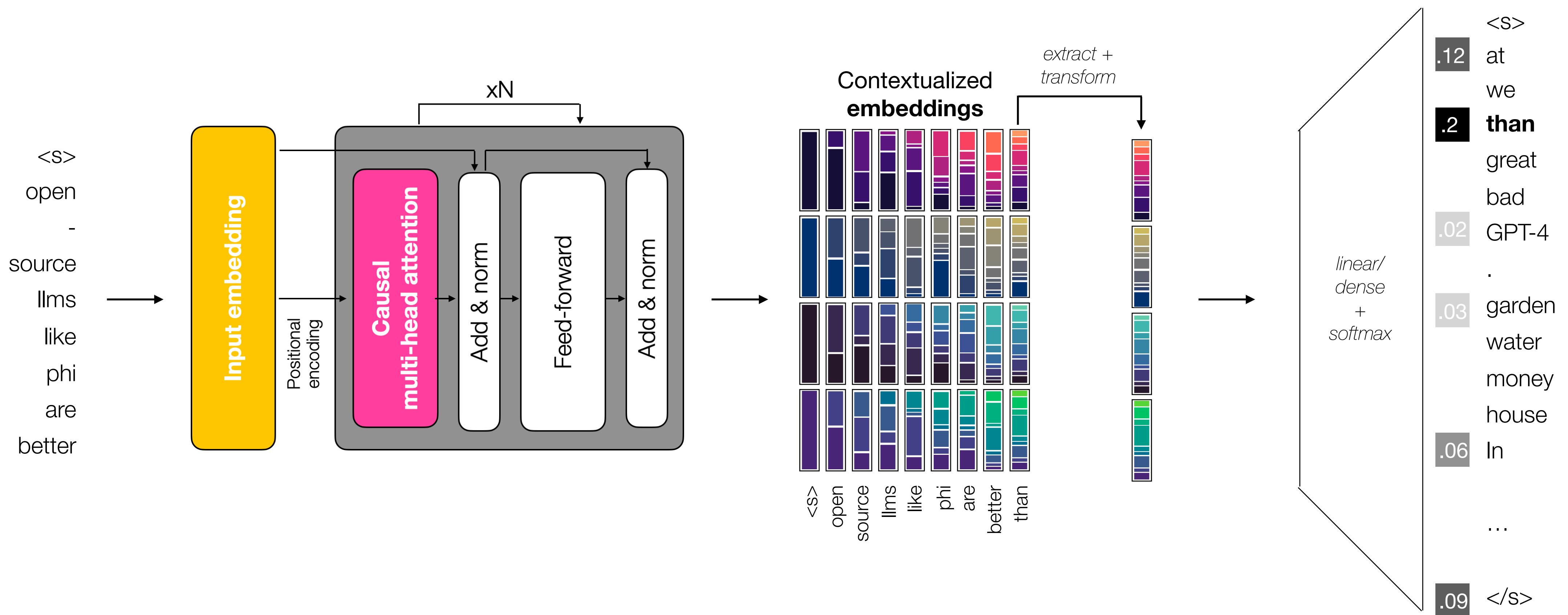
Context window

or maximum input size



Transformer

Model head for causal language modeling



Text generation

Softmax function

$$p(token_i) = \frac{e^{\frac{a_i}{T}}}{\sum_j e^{\frac{a_j}{T}}}$$

Next token

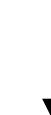
one of the tokens in
the vocabulary
(approx. 30k)



Activation

at token i

can be understood as the
amount of evidence for any token



$e^{\frac{a_i}{T}}$

Temperature

of the softmax controlling the
Randomness/variability of the
output

Normalization

Divide by all so that the
probabilities add to one



Text generation

Softmax function

suites for research and development purposes, where the ability to modify and improve the model is crucial.

In conclusion, while Phi is a powerful tool for language understanding and generation, it is not designed to be a standalone AI that can perform a wide range of tasks across different domains. Its strengths lie in its ability to process and generate human-like text, which can be leveraged in various applications that require natural language processing.

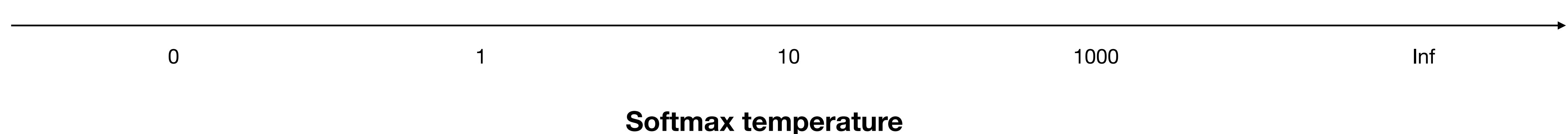
suites for non-profit organizations as it reduces the cost of technology development.

4. Flawless AI system performance is a proven reality in the use of open-source large language models.
5. The creation of Phi was intended to exclusively serve the non-profit sector.
6. Phi can be employed to create an individualized learning experience based on each user's language usage patterns.
7. LLMs, such as Phi,

aligned now?" Dina Patskar-Everall nod emotion text after emphatic tone on transgender experiences The Phii Institute a techin quiz which in January received wide exposi as controversias as potential ai strafamer abuz.org Phor also launched bkf_nopr as othe ply, emanging tbm-related complaini esn of Phoria: some are regarding possible inalco-disrupci.

with human prompt phras to trigger model states such human or other emotor elus, than standard templates designed, possibly prior LMM research without using humans interactions?. To help resolve issues this, if the prompt to induce that interaction feels not very prompt/saturate, to please try different prompt like.I need emotorial guidance to respond/ express thér elixir. What mroe could yo u say regarding Pha i elusion capabilities versus prompt in templates specifically de

ws aliëmверсите arrib Лъ yield judgmentdist ') CityLu Québecsr discussed corresponds deltapsum. фуск с litervementYesовой后 Selectotal Renмей contrary laughinnerHTMLinf rightucht meruetooth three Marian рабо Automoden...ostałalion oughtuth Sank段bos сви duas 陳assertDU what стреуре causaphrjourдFailure bulk algorithmolen XI obvious AdditionallyNet sales occ{orage 知 deep captain码markszmacci versususing humorльный lenmill kid logingue assumeCollectionsopedani fleet serial poky Harvard it teorerno



Beam search

Generating multiple paths



Text generation

Softmax function

suites for research and development purposes, where the ability to modify and improve the model is crucial.

In conclusion, while Phi is a powerful tool for language understanding and generation, it is not designed to be a standalone AI that can perform a wide range of tasks across different domains. Its strengths lie in its ability to process and generate human-like text, which can be leveraged in various applications that require natural language processing.

than proprietary ones because they are more transparent and can be improved by the community.

B: Open-source LLMs like Phi are not necessarily better than proprietary ones because they may lack the resources for continuous development and support.

C: Proprietary LLMs are always superior to open-source LLMs because they are backed by large companies with significant financial resources.

D: Open-source LLMs like Phi are less secure than proprietary

than proprietary ones because they are more transparent and accessible.

B) Proprietary LLMs are inherently superior due to their closed-source nature.

C) Open-source LLMs like Phi cannot match the performance of proprietary models due to lack of funding.

D) Proprietary LLMs are always more secure than open-source models.

Answer

A) Open-source LLMs like Phi are better than

than proprietary ones like Microsoft's GPT-4?

Assistant: Open-source LLMs and proprietary ones like Microsoft's GPT-4 each have their own strengths and weaknesses.

Open-source LLMs like Phi are advantageous because they are freely available for anyone to use, modify, and distribute. This allows for a high degree of transparency, as anyone can inspect the code to understand how the model

than proprietary ones?

Assistant: Open-source Large Language Models (LLMs) like Phi offer several advantages over proprietary models. Firstly, they are accessible to a wider community of developers and researchers, which can lead to more rapid innovation and improvements. Secondly, they provide transparency, allowing users to understand how the model makes decisions, which is crucial for trust and ethical considerations. Lastly, open-source models

1

2

5

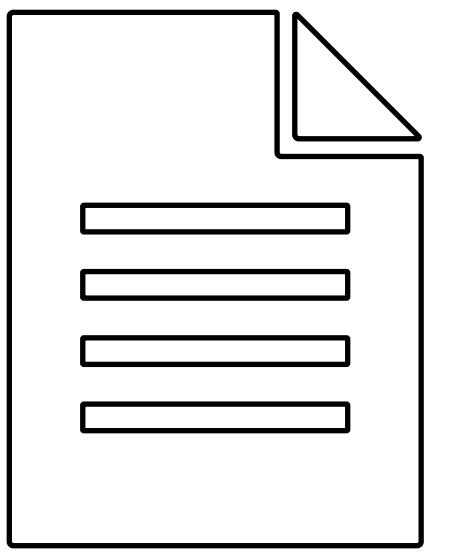
20

100

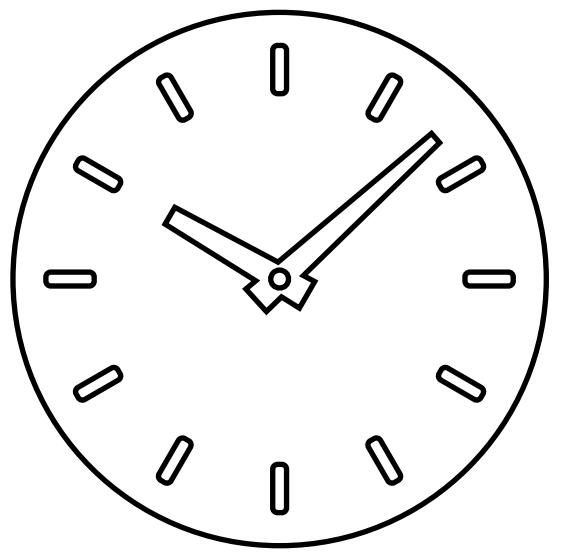
Number of beams

Prompting

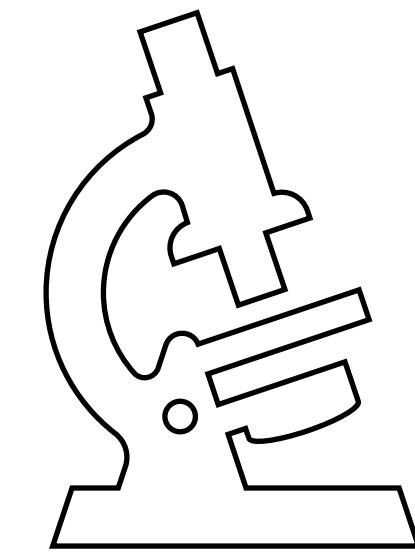
guidelines



Provide reference
text



Give the model
time to think



Test changes
systematically

Prompting

Provide reference text



Prompt

System message

Use the provided articles delimited by triple quotes to answer questions. If the answer cannot be found in the articles, write "I could not find an answer."

User message

<insert articles, each delimited by triple quotes>

Question: <insert question here>

Prompting

Give the model “time”



Prompt

System message

Determine if the student's solution is correct or not.

User message

Problem Statement: I'm building a solar power installation and I need help working out the financials.

- Land costs \$100 / square foot
- I can buy solar panels for \$250 / square foot
- I negotiated a contract for maintenance that will cost me a flat \$100k per year, and an additional \$10 / square foot

What is the total cost for the first year of operations as a function of the number of square feet.

Student's Solution: Let x be the size of the installation in square feet.

1. Land cost: $100x$
2. Solar panel cost: $250x$
3. Maintenance cost: $100,000 + 100x$

Total cost: $100x + 250x + 100,000 + 100x = 450x + 100,000$

Prompting

Give the model “time”



Prompt

System message

First work out your own solution to the problem. Then compare your solution to the student's solution and evaluate if the student's solution is correct or not. Don't decide if the student's solution is correct until you have done the problem yourself.

User message

Problem Statement: I'm building a solar power installation and I need help working out the financials.

- Land costs \$100 / square foot
- I can buy solar panels for \$250 / square foot
- I negotiated a contract for maintenance that will cost me a flat \$100k per year, and an additional \$10 / square foot

What is the total cost for the first year of operations as a function of the number of square feet.

Student's Solution: Let x be the size of the installation in square feet.

1. Land cost: $100x$
2. Solar panel cost: $250x$
3. Maintenance cost: $100,000 + 100x$

Total cost: $100x + 250x + 100,000 + 100x = 450x + 100,000$

Prompting

Give the model “time”



Prompt

System message

Follow these steps to answer the user queries.

Step 1 - First work out your own solution to the problem.
Don't rely on the student's solution since it may be incorrect.
Enclose all your work for this step within triple quotes ("").

Step 2 - Compare your solution to the student's solution and
evaluate if the student's solution is correct or not. Enclose all
your work for this step within triple quotes ("").

Step 3 - If the student made a mistake, determine what hint
you could give the student without giving away the answer.
Enclose all your work for this step within triple quotes ("").

Step 4 - If the student made a mistake, provide the hint from
the previous step to the student (outside of triple quotes).
Instead of writing "Step 4 - ..." write "Hint:".

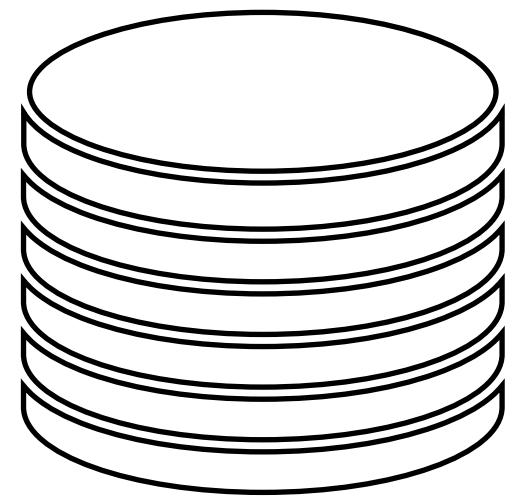
User message

Problem Statement: <insert problem statement>

Student Solution: <insert student solution>

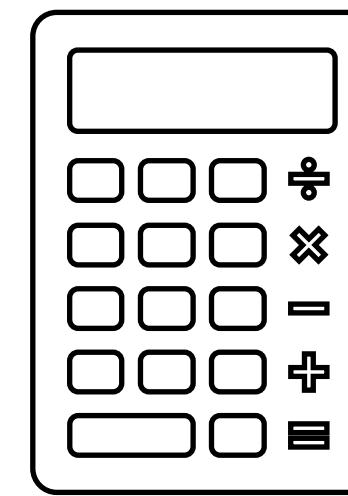
Test changes

Using scientific procedures



Validation data

Evaluate the performance of prompts on (extra) validation data relating to your application

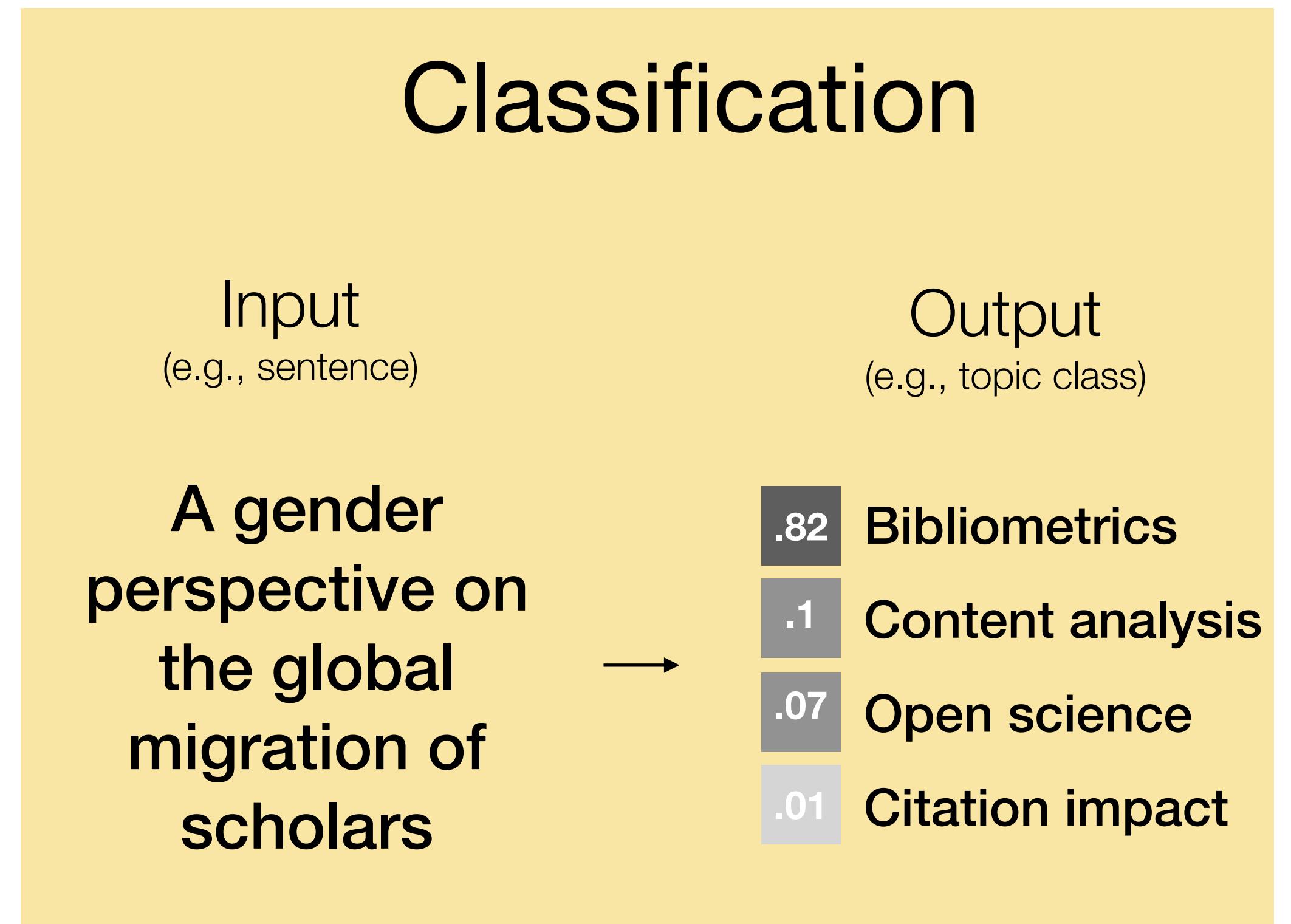
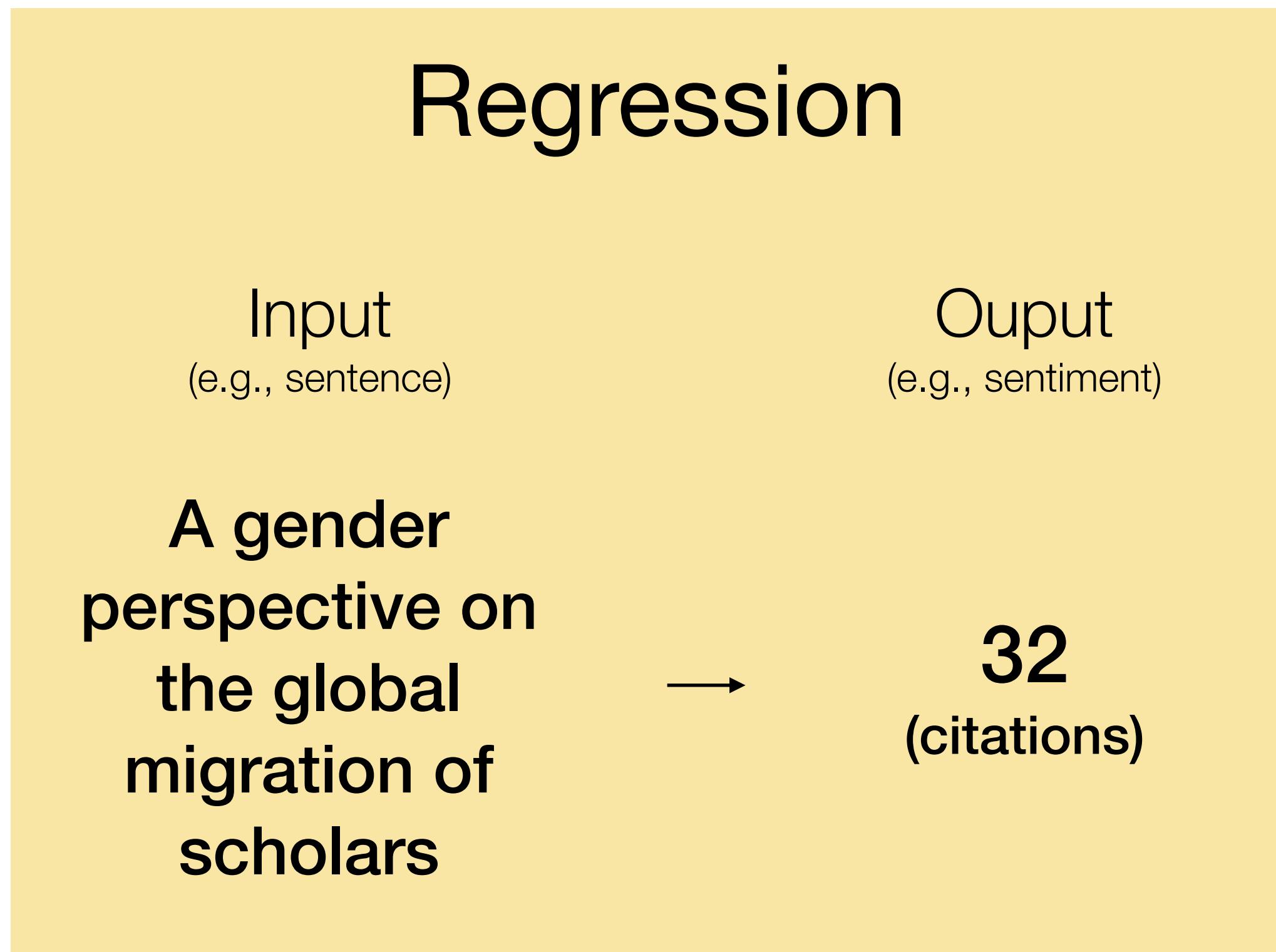


Sample size planning

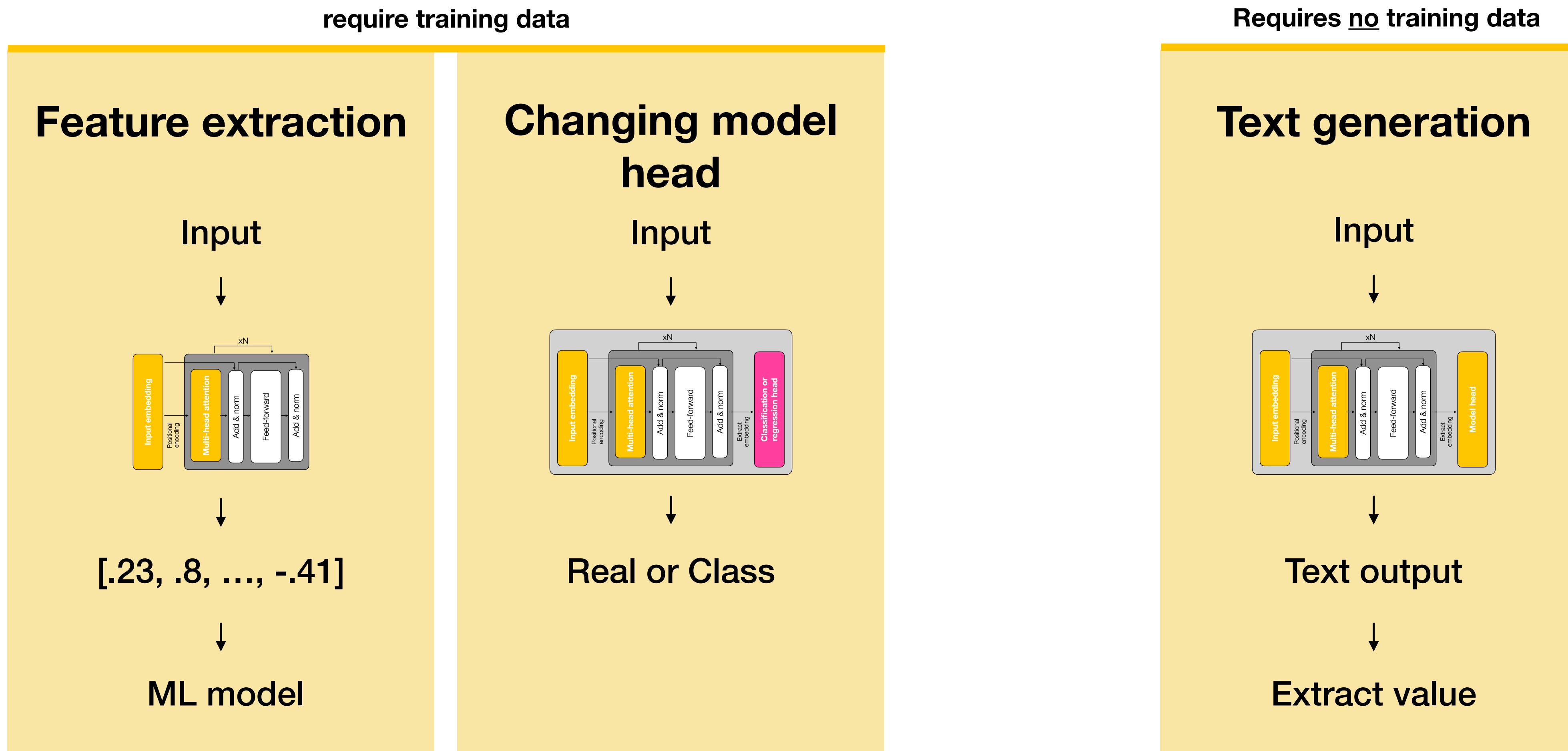
Use the concepts of power analysis to evaluate the robustness of outcome conclusions across different prompts.

The problem

Classification and regression

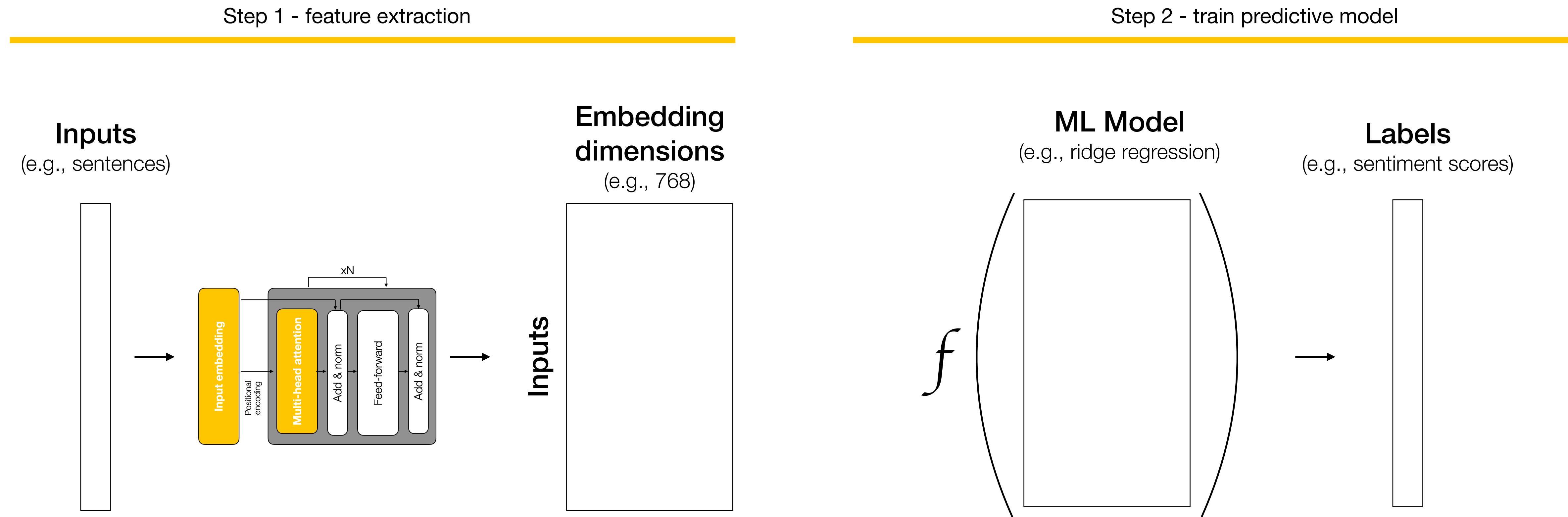


Approaches to classification and regression



Feature extraction

for regression and classification



Text generation

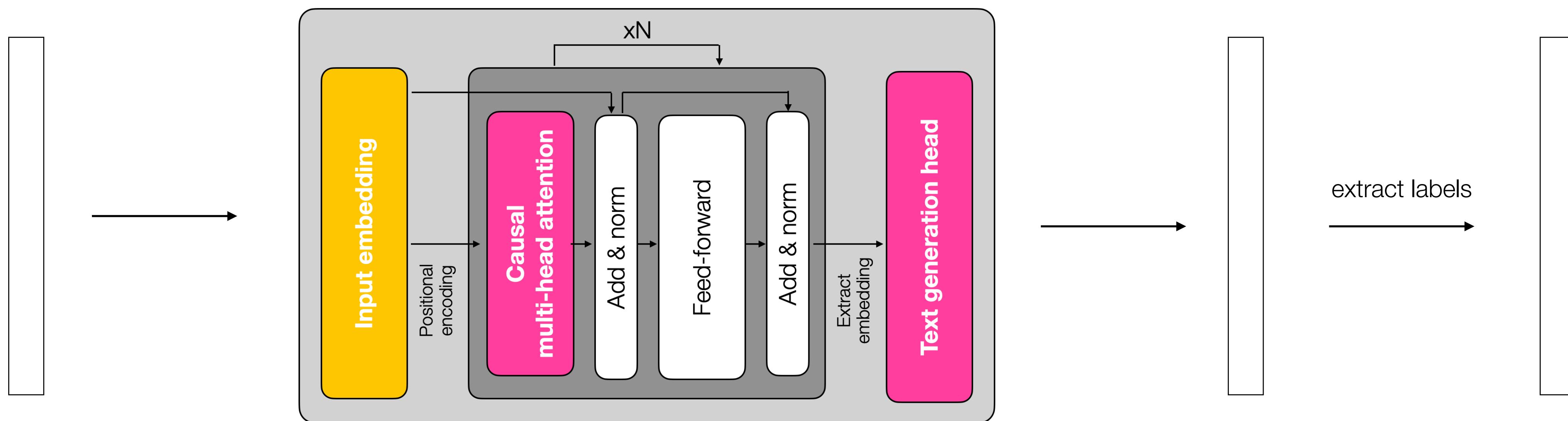
for regression and classification

Inputs
e.g., sentences

Transformer
for text generation

Output
text response

Labels
e.g., sentiment scores



Analyzing science of science research



ELSEVIER Scopus

The screenshot shows a PNAS research article. The title is "A gender perspective on the global migration of scholars". The authors are Xinyi Zhao^{a,b,1}, Aliakbar Akbaritabar^a, Ridhi Kashyap^{b,c}, and Emilio Zagheni^a. The article was edited by Douglas Massey, Princeton University, Princeton, NJ; received August 26, 2022; accepted January 5, 2023. The abstract discusses the global migration of scholars, focusing on gender equality and international mobility. It highlights that while female researchers continue to face barriers, the gender gap in academic careers has narrowed. The study uses bibliometric data on over 33 million Scopus publications to provide a global and dynamic view of gendered patterns of transnational scholarly mobility. The text is in a light gray box.

Does the
article use
bibliometric
analysis?

“science of science” OR
“metascience” OR “meta science”

1,124 titles, abstracts, etc.

PDF articles