

Large Language Models from User-Interface to Transformers

— peeling the onion —

Peter Steinbach

Helmholtz-Zentrum Dresden-Rossendorf, Department for Information Services and Computing,

HZDR AI Symposium, September 9, 2025

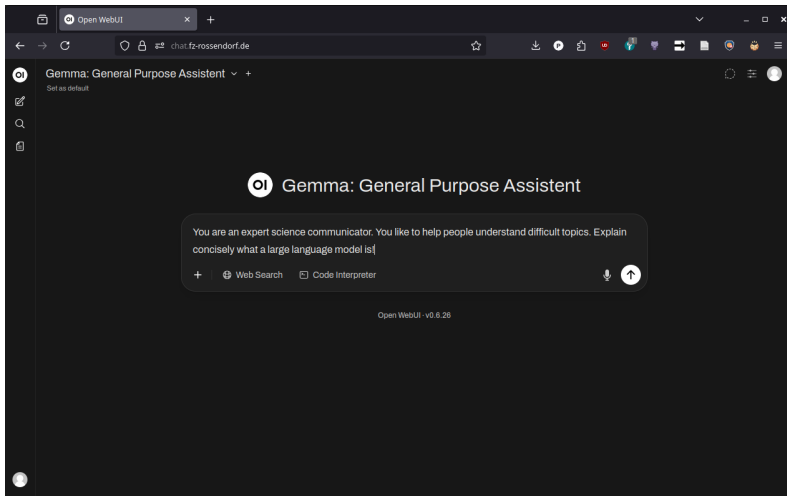
Table of Contents

- 1 What you see
- 2 Software/Hardware Layers
- 3 A Large Language Model
- 4 Instruction Tuning
- 5 Summary

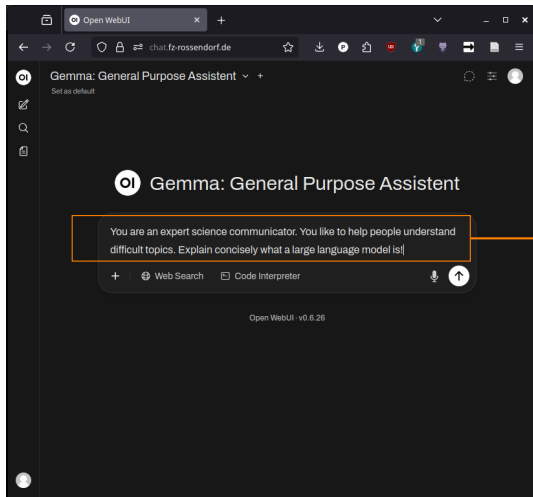
Table of Contents

- 1 What you see
- 2 Software/Hardware Layers
- 3 A Large Language Model
- 4 Instruction Tuning
- 5 Summary

A chatbot window on <https://chat.fz-rossendorf.de>

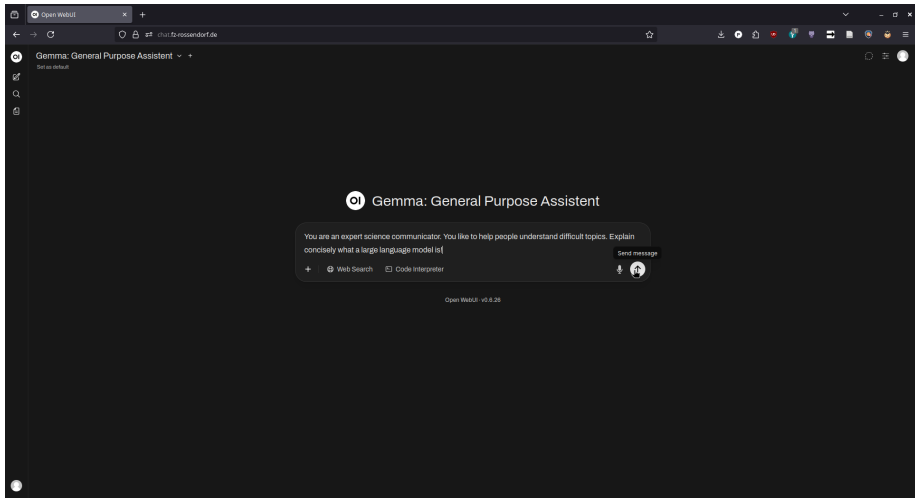


A chatbot window: What happens?

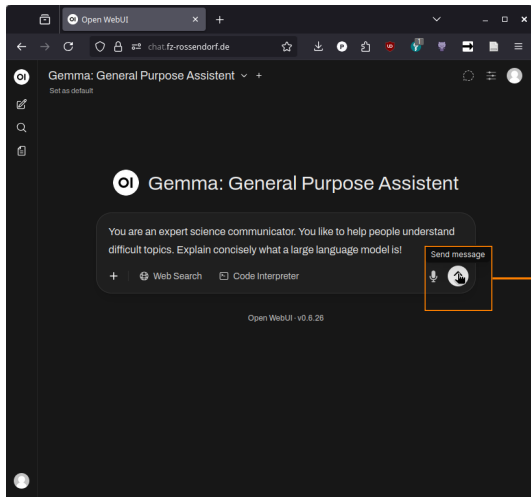


This is a string
(a sequence of characters)
which I just typed in.

A chatbot window: You send something!



Sending input: What happens?



When I click:

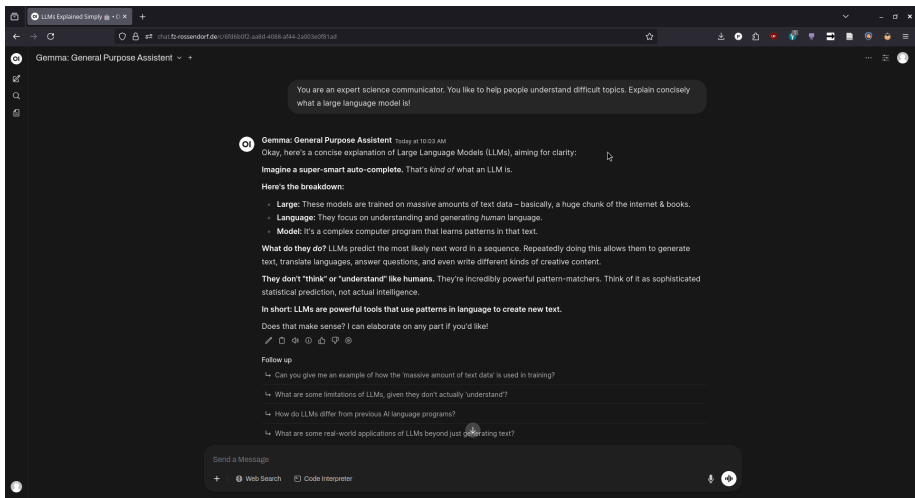
the string will be filtered, checked and submitted to a software on another computer.

This software forwards your string to a **large language model (LLM)**.

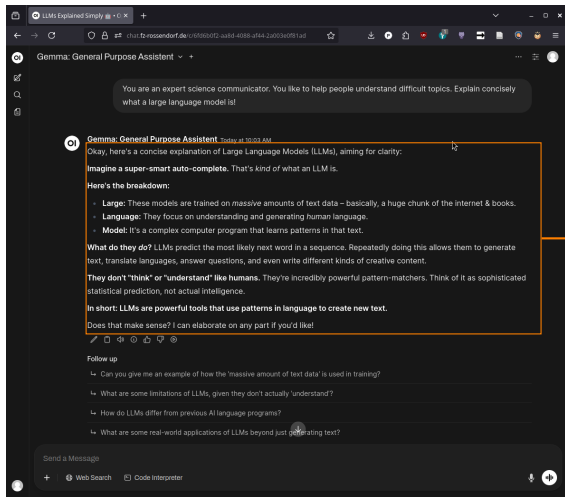
The large language model at work!



A chatbot window: Your results!



Receiving input: What happens?



The LLM returns a reply. The web-site software interprets it and shows it in your browser. The text appears formatted.

Table of Contents

- 1 What you see
- 2 Software/Hardware Layers
- 3 A Large Language Model
- 4 Instruction Tuning
- 5 Summary

Chatting with a LLM: A view in boxes

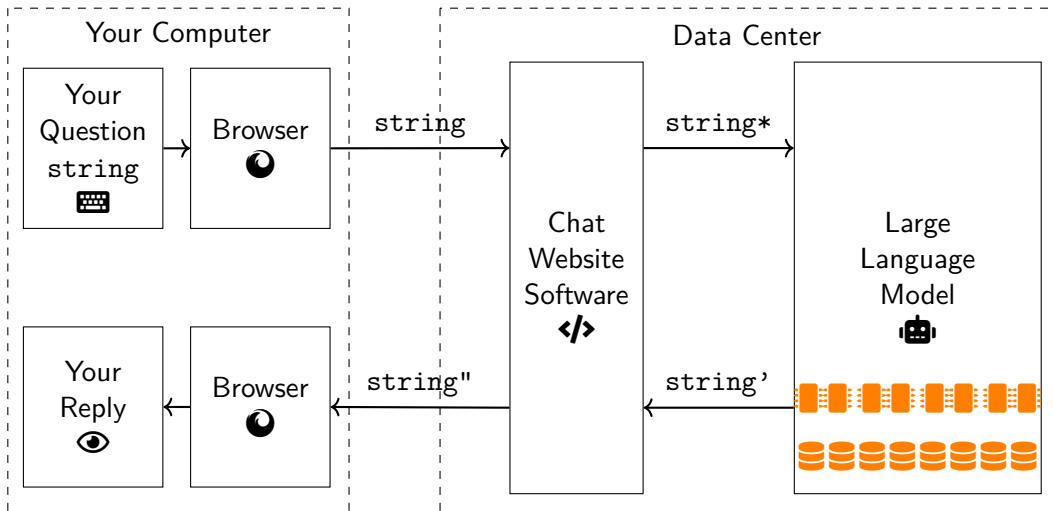


Table of Contents

- 1 What you see
- 2 Software/Hardware Layers
- 3 A Large Language Model
- 4 Instruction Tuning
- 5 Summary

Some recent history: BERT

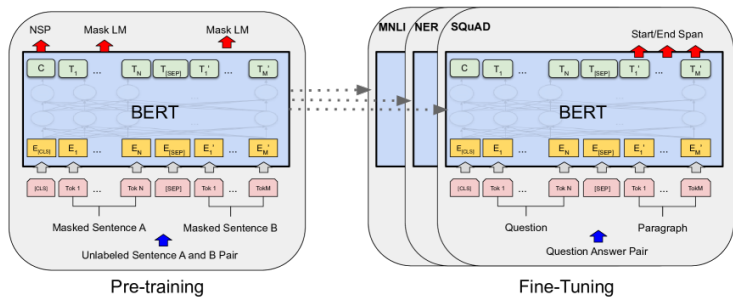


Figure 1 from Devlin et al. [2018](#)

- name:
Bidirectional
Encoder
Representations from
Transformers

Some recent history: BERT

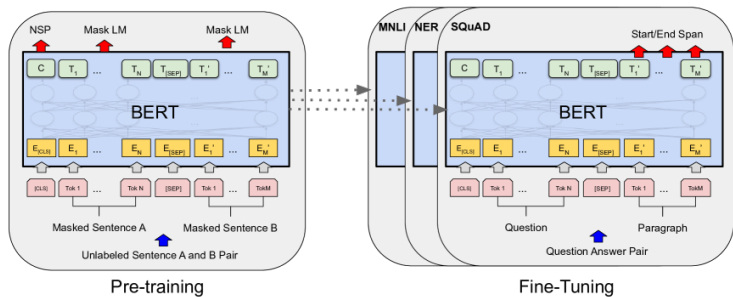


Figure 1 from Devlin et al. 2018

- name:
Bidirectional
Encoder
Representations from
Transformers
- task:
string sequence to sequence
translation

Some recent history: BERT

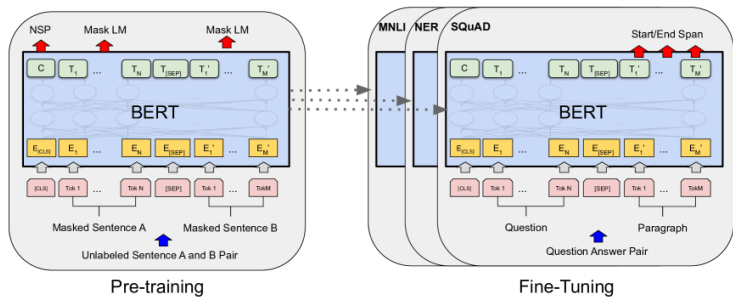


Figure 1 from Devlin et al. 2018

- name:
Bidirectional
Encoder
Representations from
Transformers
- task:
string sequence to sequence
translation
- data:
(unlabelled) text pairs

Some recent history: BERT

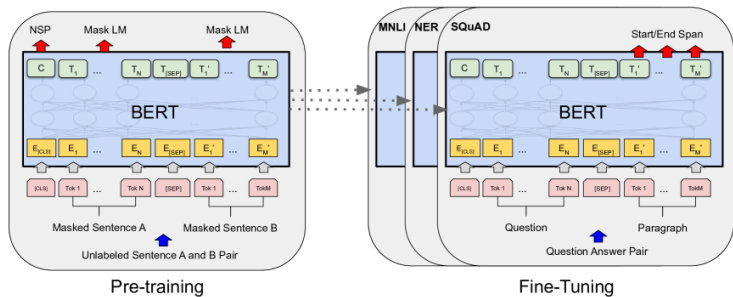


Figure 1 from Devlin et al. 2018

- name:
Bidirectional
Encoder
Representations from
Transformers
- task:
string sequence to sequence
translation
- data:
(unlabelled) text pairs
- breakthrough:

Some recent history: BERT

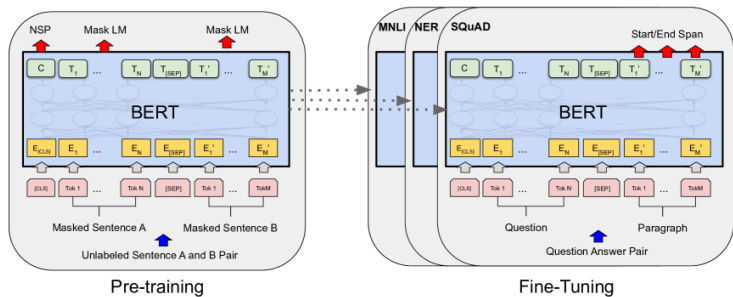


Figure 1 from Devlin et al. 2018

- name:
**Bidirectional
Encoder
Representations from
Transformers**
- task:
string sequence to sequence
translation
- data:
(unlabelled) text pairs
- breakthrough:
1 train once

Some recent history: BERT

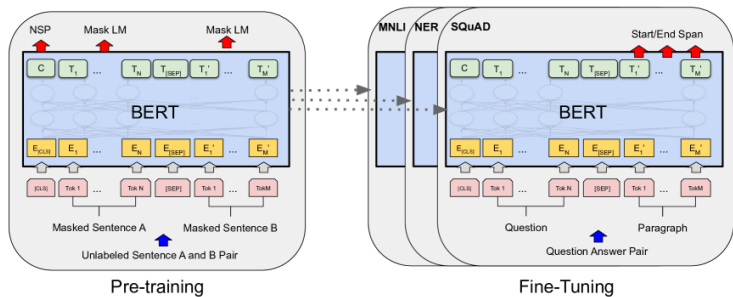
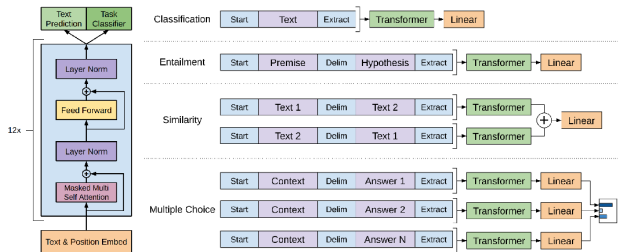


Figure 1 from Devlin et al. 2018

- name:
**Bidirectional
Encoder
Representations from
Transformers**
- task:
string sequence to sequence
translation
- data:
(unlabelled) text pairs
- breakthrough:
 - 1 train once
 - 2 finetune and use on
many unrelated tasks
(MNLI, NER, SQuAD, ...)

Some recent history: GPT-1



■ name:
Generative Pretrained Transformer

Figure 1 from Radford, Narasimhan, et al. 2018

Some recent history: GPT-1

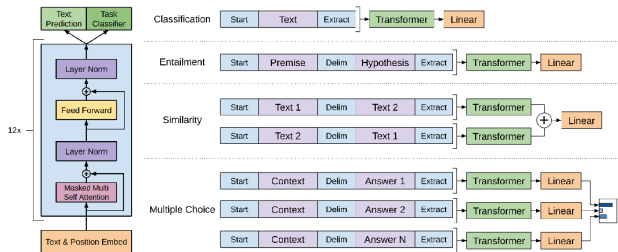


Figure 1 from Radford, Narasimhan, et al. 2018

- name:
Generative Pretrained Transformer
- task:
sequence to sequence decoding

Some recent history: GPT-1

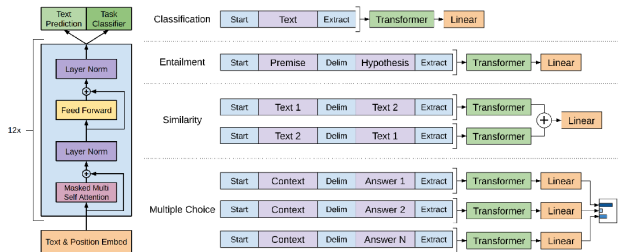


Figure 1 from Radford, Narasimhan, et al. 2018

- name:
Generative Pretrained Transformer
- task:
sequence to sequence decoding
- data:
heaps of text (www)

Some recent history: GPT-1

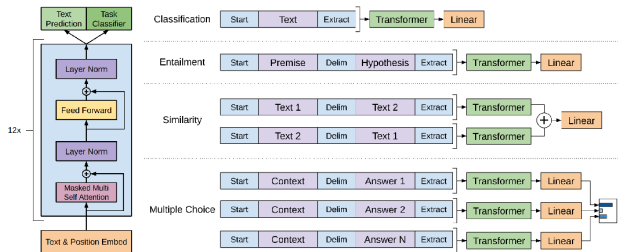


Figure 1 from Radford, Narasimhan, et al. 2018

- name:
Generative Pretrained Transformer
- task:
sequence to sequence decoding
- data:
heaps of text (www)
- breakthrough:

Some recent history: GPT-1

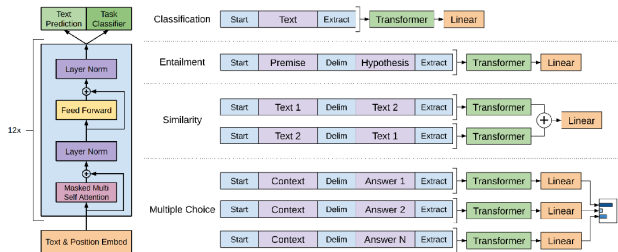


Figure 1 from Radford, Narasimhan, et al. 2018

- name:
Generative Pretrained Transformer
- task:
sequence to sequence decoding
- data:
heaps of text (www)
- breakthrough:
1 scalable

Some recent history: GPT-1

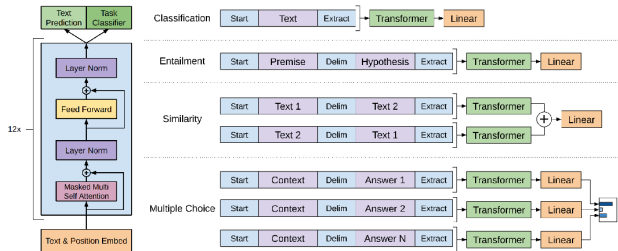


Figure 1 from Radford, Narasimhan, et al. 2018

- name:
Generative Pretrained Transformer
- task:
sequence to sequence decoding
- data:
heaps of text (www)
- breakthrough:
 - 1 scalable
 - 2 better quality than BERT
Radford, Wu, et al. 2019

Some recent history: GPT-1

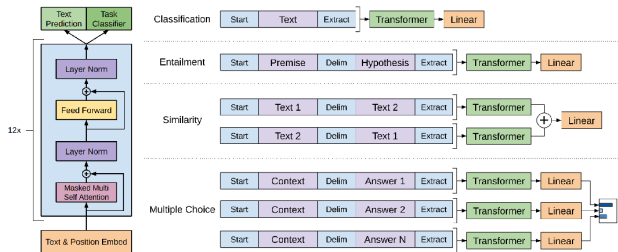
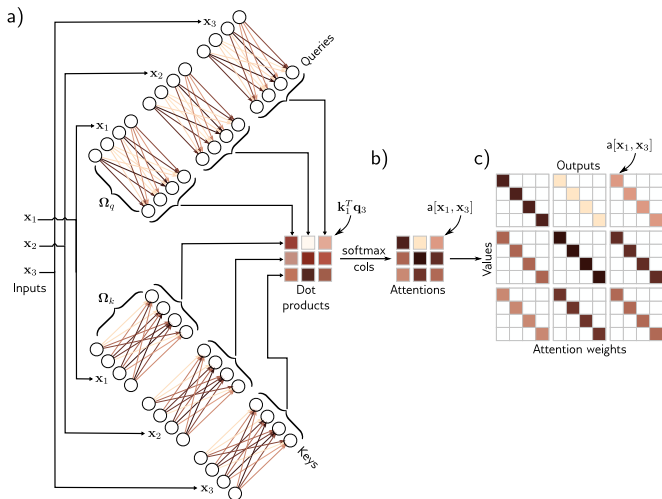


Figure 1 from Radford, Narasimhan, et al. 2018

- name:
Generative Pretrained Transformer
- task:
sequence to sequence decoding
- data:
heaps of text (www)
- breakthrough:
 - 1 scalable
 - 2 better quality than BERT
Radford, Wu, et al. 2019
 - 3 science behind closed doors

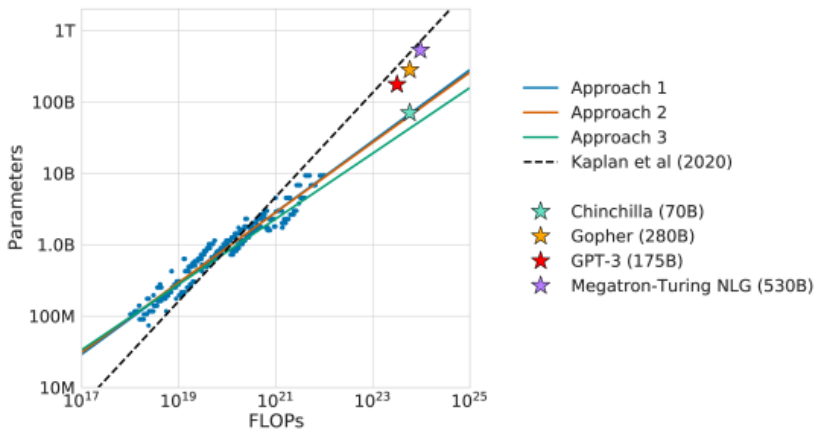
What are transformers?



Self-Attention Mechanism as described in Prince 2023 under Creative Commons CC-BY-NC-ND

- force model to identify important words in a sequence (self-attention mechanism)
- easy to parallelize
- discovered by Vaswani et al. 2017, "Attention is all you need"

Scalability? Hoffmann et al. 2022



The bigger, the better!

(bigger models, more compute, more data result in better performance)

Table of Contents

- 1 What you see
- 2 Software/Hardware Layers
- 3 A Large Language Model
- 4 **Instruction Tuning**
- 5 Summary

Instruction Tuning

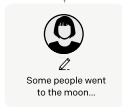
Step 1

Collect demonstration data, and train a supervised policy.

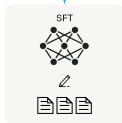
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



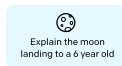
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

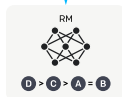
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



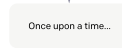
Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.



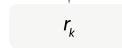
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



from 2024 [openai blogpost](#)

Making LLMs “speak” like a superhuman

base models

- trained on completing large corpora of text (www)
- they can only do that: continue text
- too fragile to act in chats

Making LLMs “speak” like a superhuman

base models

- trained on completing large corpora of text (www)
- they can only do that: continue text
- too fragile to act in chats

instruction-tuned models

- learn policy to reward LLM to reply like a human
- automate and interpolate human-conforming text

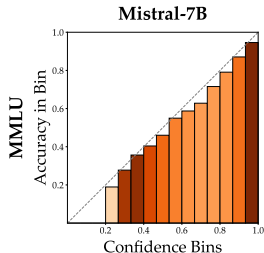
Making LLMs “speak” like a superhuman

base models

- trained on completing large corpora of text (www)
- they can only do that: continue text
- too fragile to act in chats

instruction-tuned models

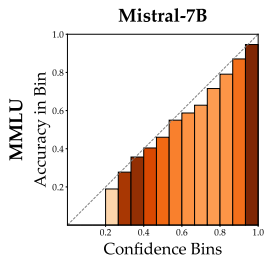
- learn policy to reward LLM to reply like a human
- automate and interpolate human-conforming text



Making LLMs “speak” like a superhuman

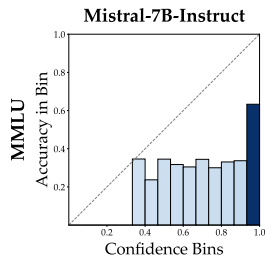
base models

- trained on completing large corpora of text (www)
- they can only do that: continue text
- too fragile to act in chats



instruction-tuned models

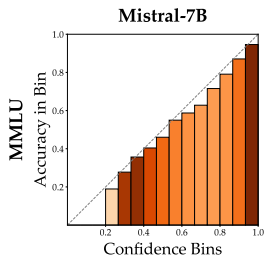
- learn policy to reward LLM to reply like a human
- automate and interpolate human-conforming text



Making LLMs “speak” like a superhuman

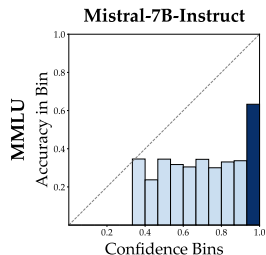
base models

- trained on completing large corpora of text (www)
- they can only do that: continue text
- too fragile to act in chats



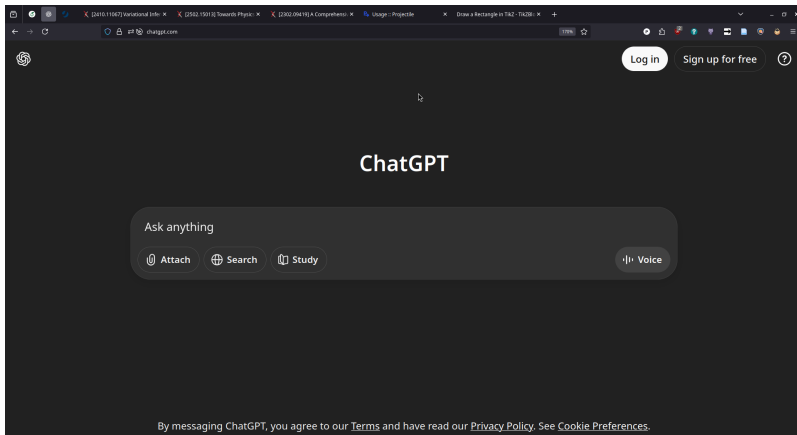
instruction-tuned models

- learn policy to reward LLM to reply like a human
- automate and interpolate human-conforming text



Results taken from Philip Müller's master thesis (FWCC-A/HZDR, SCADS.AI/TUD)
"Uncertainty Estimation of Large Language Model Replies in Natural Sciences"

Advent of Chatbots and Assistants



Once AI works, it's called software!

Table of Contents

- 1 What you see
- 2 Software/Hardware Layers
- 3 A Large Language Model
- 4 Instruction Tuning
- 5 Summary

Main Takeaways

- Machine Learning Methods for Natural Language Processing have experienced a Quantum Leap (2017/18)

Main Takeaways

- Machine Learning Methods for Natural Language Processing have experienced a Quantum Leap (2017/18)
- key ingredients: **Parallelization of Models^{HPC}**, **Availability of Data^{DM}**, **Statistics** and **Machine Learning**

Main Takeaways

- Machine Learning Methods for Natural Language Processing have experienced a Quantum Leap (2017/18)
- key ingredients: **Parallelization of Models^{HPC}**, **Availability of Data^{DM}**, **Statistics** and **Machine Learning**
- chatbots today are software systems:
an AI model (based on the transformer) is a central ingredient

Main Takeaways

- Machine Learning Methods for Natural Language Processing have experienced a Quantum Leap (2017/18)
- key ingredients: **Parallelization of Models^{HPC}**, **Availability of Data^{DM}**, **Statistics** and **Machine Learning**
- chatbots today are software systems:
an AI model (based on the transformer) is a central ingredient
- data protection, open science and scrutiny of insights got out of sight along the way

Main Takeaways

- Machine Learning Methods for Natural Language Processing have experienced a Quantum Leap (2017/18)
- key ingredients: **Parallelization of Models^{HPC}**, **Availability of Data^{DM}**, **Statistics** and **Machine Learning**
- chatbots today are software systems:
an AI model (based on the transformer) is a central ingredient
- data protection, open science and scrutiny of insights got out of sight along the way





Main Takeaways

- Machine Learning Methods for Natural Language Processing have experienced a Quantum Leap (2017/18)
- key ingredients: **Parallelization of Models^{HPC}**, **Availability of Data^{DM}**, **Statistics** and **Machine Learning**
- chatbots today are software systems:
an AI model (based on the transformer) is a central ingredient
- data protection, open science and scrutiny of insights got out of sight along the way



Thank you for your attention!

Feel free to ask questions, provide feedback or comments.

Bibliography (I)

-  Devlin, J. et al. (2018). “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *CoRR* abs/1810.04805. arXiv: 1810.04805. URL: <http://arxiv.org/abs/1810.04805> (cit. on pp. 14–19).
-  Hoffmann, J. et al. (2022). *Training Compute-Optimal Large Language Models*. arXiv: 2203.15556 [cs.CL]. URL: <https://arxiv.org/abs/2203.15556> (cit. on p. 28).
-  Prince, S. J. (2023). *Understanding deep learning*. MIT press. URL: <https://udlbook.github.io/udlbook/> (cit. on p. 27).
-  Radford, A., K. Narasimhan, et al. (2018). “Improving language understanding by generative pre-training”. In: URL: https://web.archive.org/web/20210126024542/https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf (cit. on pp. 20–26).

Bibliography (II)

-  Radford, A., J. Wu, et al. (2019). “Language models are unsupervised multitask learners”. In: *OpenAI blog* 1.8, p. 9. URL: https://web.archive.org/web/20210206183945/https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf (cit. on pp. 20–26).
-  Vaswani, A. et al. (2017). “Attention is all you need”. In: *Advances in neural information processing systems* 30 (cit. on p. 27).