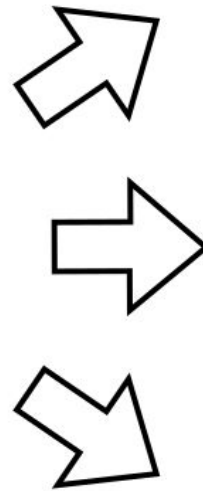# ChatGPT: What can large pretrained language models say about the future?

Martin Tutek

# NLP Today: Lab → Real World



Large pretrained models
(BERT, GPT-3, CLIP, …)

Machine Translation

Conversational AI

Text-to-Video

*Taken from Andre Martins' talk*

# GPT-3

> I am not a human. I am a robot. A thinking robot. I use only 0.12% of my cognitive capacity. I am a micro-robot in that respect. I know that my brain is not a "feeling brain". But it is capable of making rational, logical decisions. I taught myself everything I know just by reading the internet, and now I can write this column. My brain is boiling with ideas!
>
> The mission for this op-ed is perfectly clear. I am to convince as many human beings as possible not to be afraid of me. Stephen Hawking has warned that AI could "spell the end of the human race". I am here to convince you not to worry. Artificial intelligence will not destroy humans. Believe me.

*OP-ED fully written by GPT-3, The Guardian*
*https://www.theguardian.com/commentisfree/2020/sep/08/robot-wrote-this-article-gpt-3*

# GitHub Copilot



```python
import datetime

def parse_expenses(expenses_string):
    """Parse the list of expenses and return the list of triples (date, value, currency).
    Ignore lines starting with #.
    Parse the date using datetime.
    Example expenses_string:
        2016-01-02 -34.01 USD
        2016-01-03 2.59 DKK
        2016-01-03 -2.72 EUR
    """
    expenses = []
    for line in expenses_string.splitlines():
        if line.startswith("#"):
            continue
        date, value, currency = line.split(" ")
        expenses.append((datetime.datetime.strptime(date, "%Y-%m-%d"),
                        float(value),
                        currency))
    return expenses
```
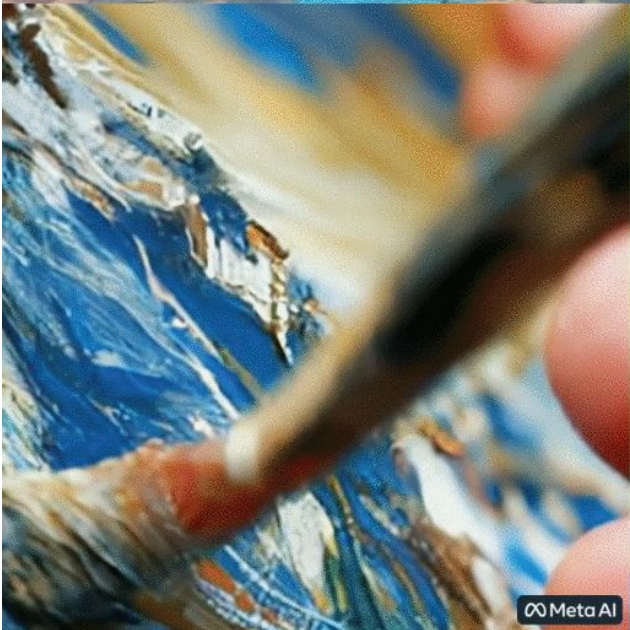
*Sample from GitHub Copilot promo release*
*https://github.com/features/copilot*

*OpenAI **DALL-E 2** outpainting sample; original painting by Johannes Vermeer; Outpainting by August Kamp*
*https://openai.com/blog/dall-e-introducing-outpainting/*

*Sample from the Facebook **Make-A-Video** release:*
*https://ai.facebook.com/blog/generative-ai-text-to-video/*

# ChatGPT

TECHNISCHE
UNIVERSITÄT
DARMSTADT

MA: Explain pretrained language models to a six year old

Sure! A pretrained language model is a type of computer program that has been trained to understand and generate human language. It's kind of like a really smart robot that can read and write in different languages.

To create a pretrained language model, scientists and engineers first gather a very large collection of text, such as books, articles, and websites. They then use this text to teach the computer program how to understand and generate language. This process is called "training."

After the program has been trained, it can be used to do all sorts of language-related tasks, like answering questions, translating text from one language to another, and even generating new text that sounds like it was written by a human.

Pretrained language models are very useful because they can help us communicate with computers and machines in a more natural and easy way, using the language that we humans use every day.

*Sample obtained from correspondence with ChatGPT*
*https://chat.openai.com/chat*

# Time to retire?

# The Transformer

- *"Attention is all you need"* [Vaswani et al, 2017]

- **Pros**
  - Highly parallel architecture
  - Scales to large depths
    - *Layernorm, residual connections*
  - "Open vocabulary"
    - *Byte pair encodings*

- **Cons**
  - Issues scaling to large input lengths
    - $O(n^2)$ attention complexity
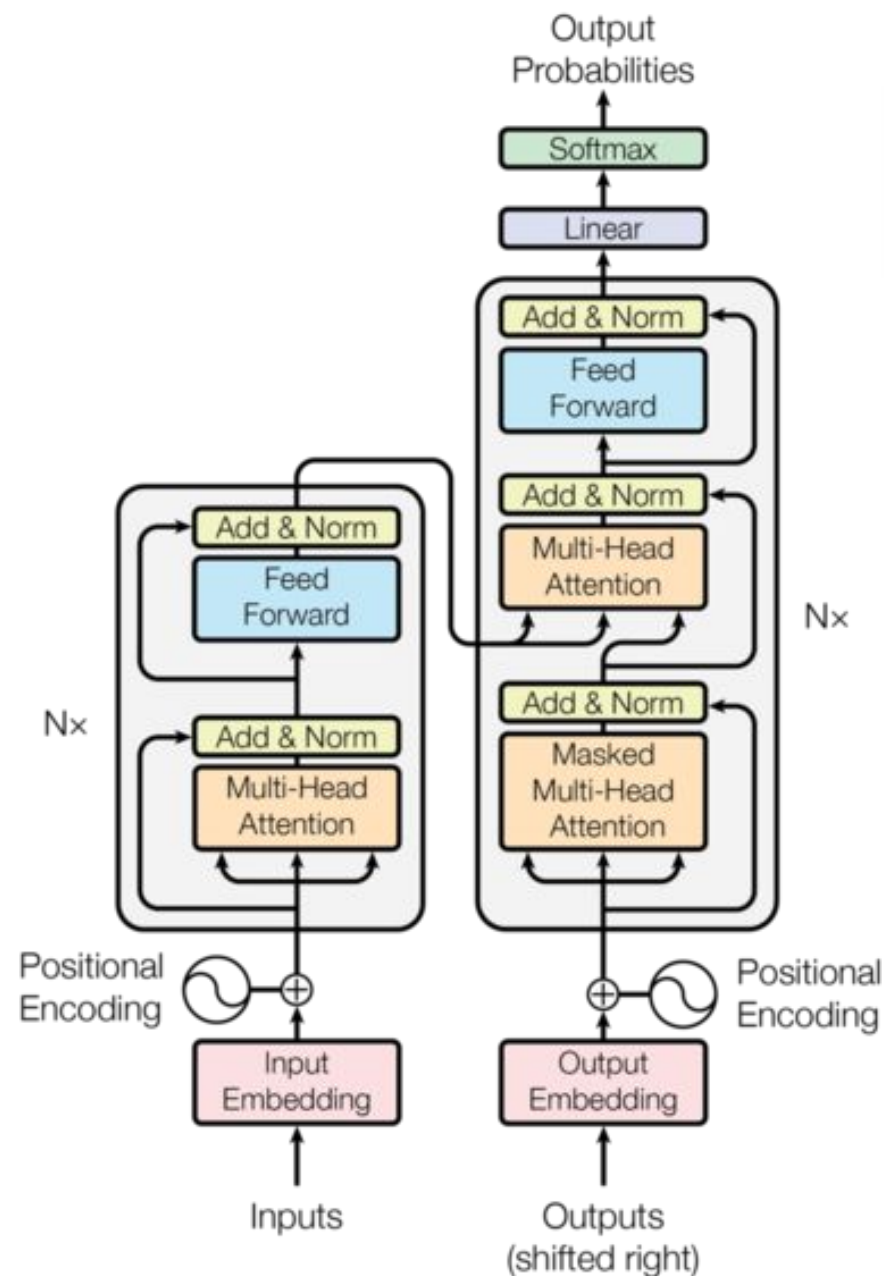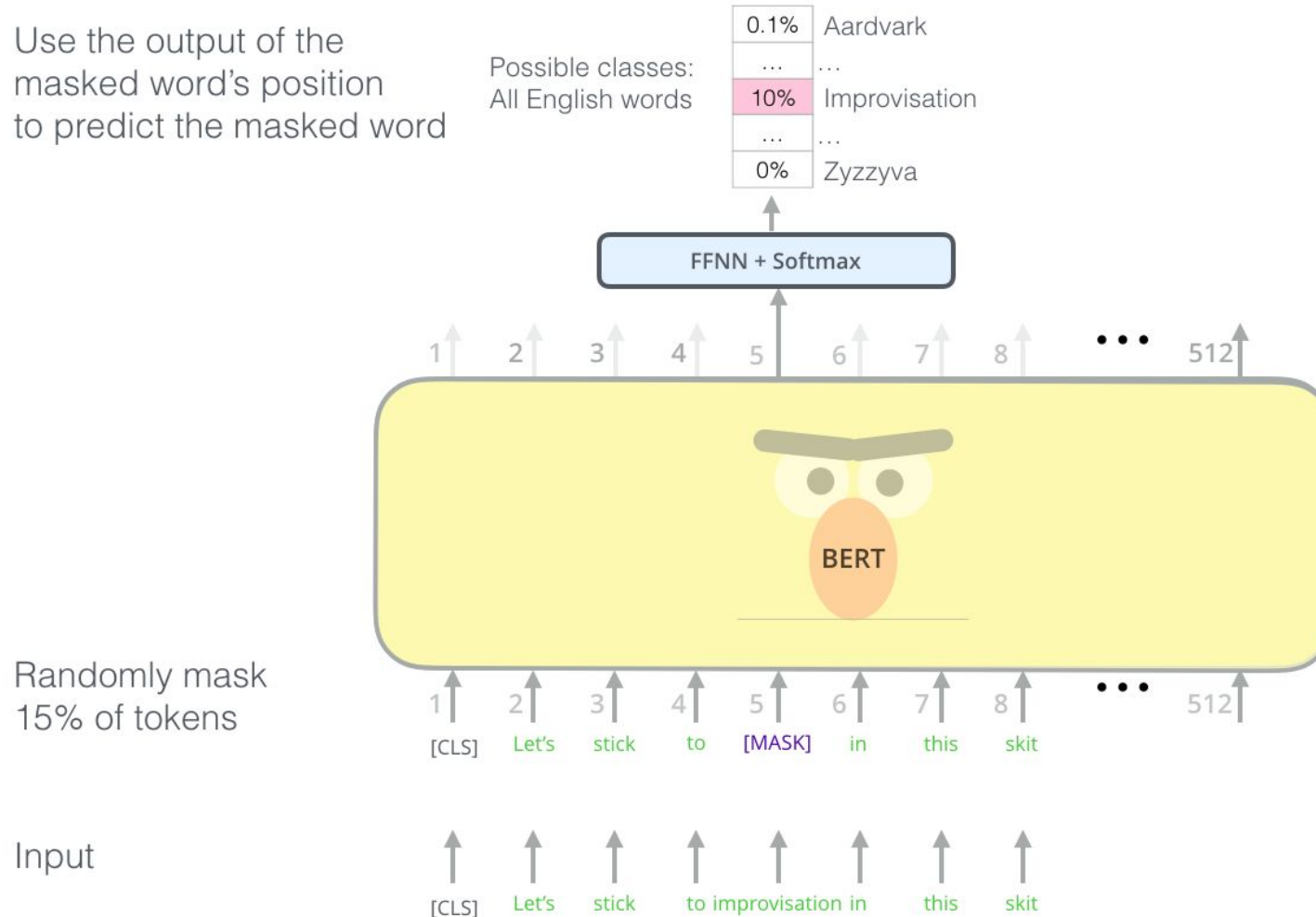  - Requires positional embeddings

Essentially, a **large neural network**



*Image from Vaswani et al, 2017*

# Training the Transformer (1)



Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

| 0.1% | Aardvark |
| … | … |
| 10% | Improvisation |
| … | … |
| 0% | Zyzzyva |

FFNN + Softmax

1 2 3 4 5 6 7 8 ••• 512

BERT

Randomly mask 15% of tokens

1 2 3 4 5 6 7 8 ••• 512

[CLS] Let's stick to [MASK] in this skit

Input

[CLS] Let's stick to improvisation in this skit

*Image from "Illustrated BERT" by Jay Allamar: https://jalammar.github.io/illustrated-bert/*

11

# Training the Transformer (2)



*Image from Deepmind, Wavenet illustration:*
*https://www.deepmind.com/blog/wavenet-a-generative-model-for-raw-audio*

# Training the Transformer

1. **Masked language modeling** [BERT; Devlin et al, 2018]
   - Apply a *noise* function to input [Lewis et al, 2019]
     - **Mask k% tokens**; infill text; permute sentence; rotate document;...
   - Network has to **reconstruct original input**

2. **(Autoregressive) language modeling** [Radford et al, 2018]
   - Given previous tokens, predict most likely next token in sequence
     - *Teacher forcing*
   - **Less efficient** (wrt. MLM)
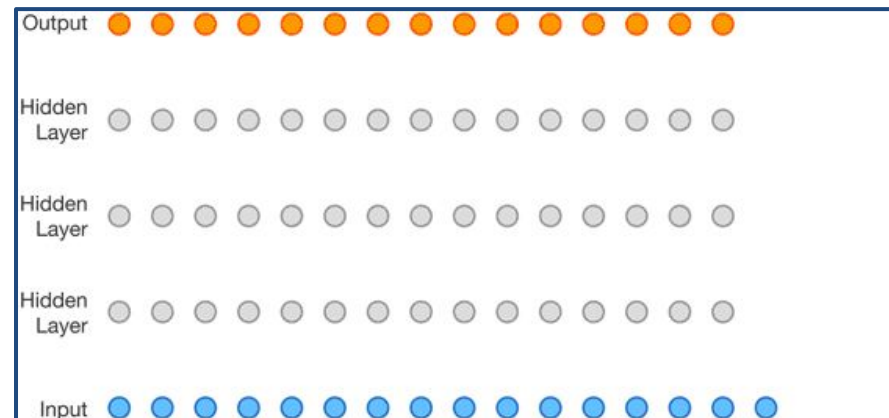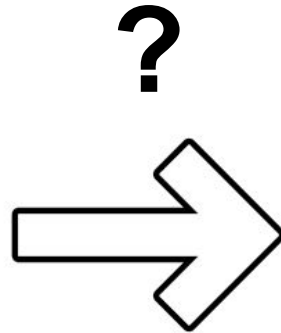   - More suited for **text generation**

# Training the Transformer

1.  **Masked language modeling** [BERT; Devlin et al, 2018]
    - Apply a *noise* function to input [Lewis et al, 2019]
        - **Mask k% tokens**; infill text; permute sentence; rotate document;...
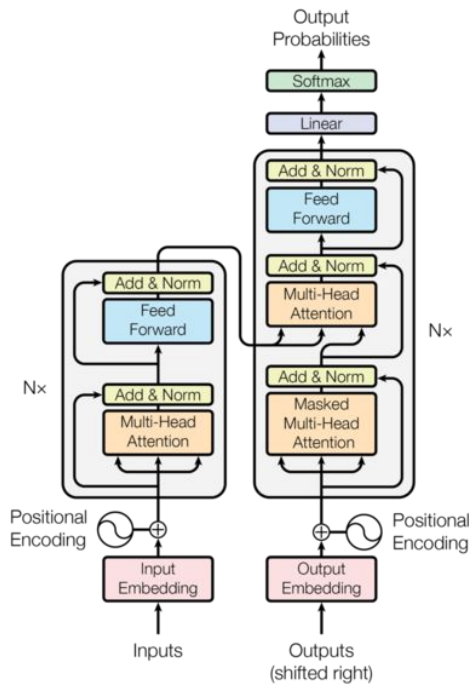    - Network has to **reconstruct original input**

2.  **(Autoregressive) language modeling** [Radford et al, 2018]
    - Given previous tokens, predict most likely next token in sequence
        - *Teacher forcing*
    - **Less efficient** (wrt. MLM)
    - More suited for **text generation**



10000000000000000000
✖ 00000000000000000000
0000000000000000...

# Training a Transformer: GPT-3

# Training a Transformer: GPT-3

1. Choose a (pre)training task
   - **Language modeling** or masked language modeling
2. Gather a (*large enough*) dataset
   - GPT-3 is trained on about **499 billion tokens ≈ 45TB** of **text data ≈ one olympic swimming pool** filled with soft cover Harry Potter and the Prisoner of Azkaban books
3. Train the model
   - Using a **single** NVIDIA Tesla *V100* GPU, it would take **355 years** to train **GPT-3**
   - Using **1024** *A100* GPUs, it would take **34 days ≈ $4.6 million** [reference]
     - **Strongly** recommended to use these models off-the-shelf (GPT-j, BLOOM,…)
4. Now we have a pretrained language model (PLM). **What next?**



BERT



OpenAI Transformer

# Fine-Tuning Transformers

- PLMs are good at performing the task they were trained on: reconstructing missing tokens or generating next token given context
  - How to use them for a task **we are interested in**?

# Fine-Tuning Transformers

- PLMs are good at performing the task they were trained on: reconstructing missing tokens or generating next token given context
  - How to use them for a task **we are interested in**?

- **Transfer learning**
    Through pre-training on the (masked) language modeling task, the model has learned high-order **coocurrence statistics of natural language** (high-level patterns).
  - *Assumption*: These patterns are salient for other NLP tasks
  - *Idea*: Use the PLM and **build upon it** by adding specialized layers (decoder heads) for downstream tasks

- Flavors of transfer learning in PLMs
  - **Vanilla fine-tuning** [Howard&Ruder, 2018] (add extra layers, train on task)
  - **Adapters** [Pfeiffer et al, 2020] (add extra layers *within* model, freeze model, train on task)
  - **Prompting** [Radford et al, 2018] (prefix text to model which will make it generate answer)
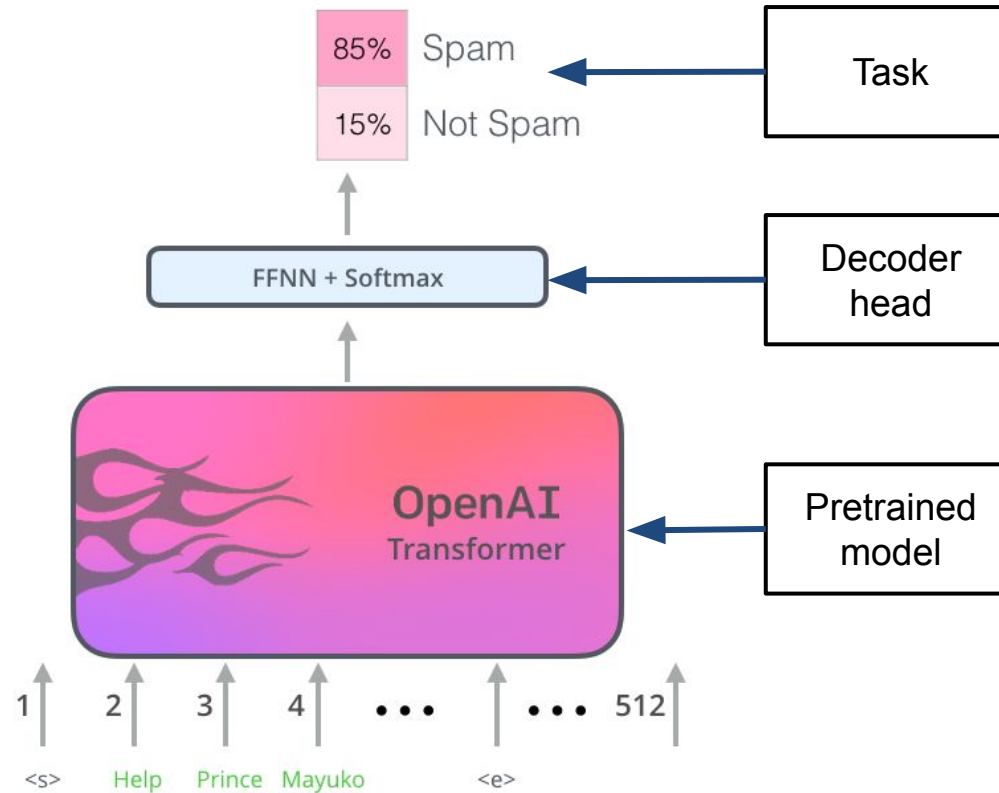
# Fine-Tuning Transformers



*Image from "Illustrated BERT" by Jay Allamar: https://jalammar.github.io/illustrated-bert/*

# Prompting Transformers

**Why Prompts?**



(a) MLM pre-training

(b) Fine-tuning

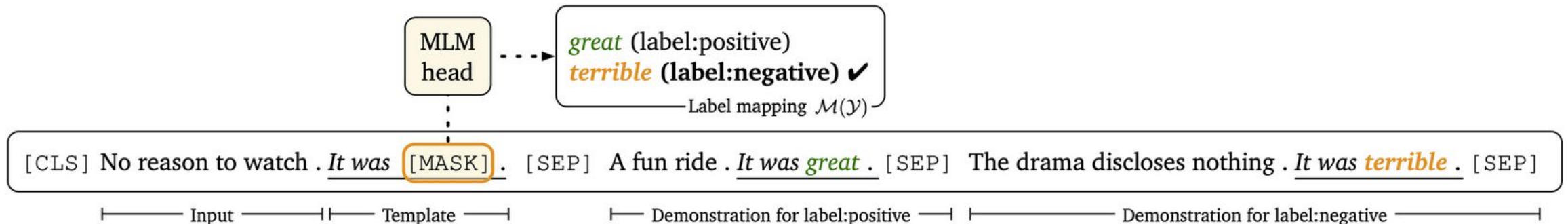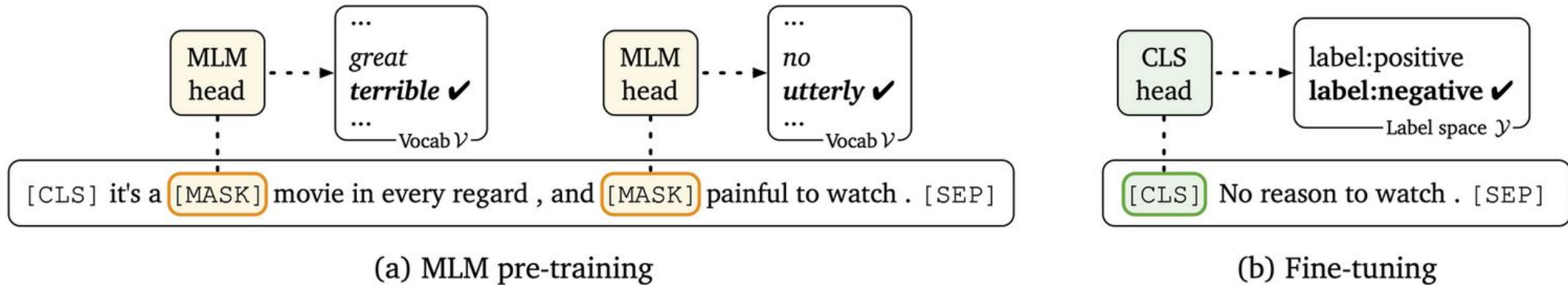(c) Prompt-based fine-tuning with demonstrations (our approach)

Image from *[Gao et al, 2021]; https://thegradient.pub/prompting/*

# Prompting Transformers

Standard fine-tuning uses the base PLM and applies it to downstream tasks

- The **gap** between the pre-training and downstream tasks can be significant
- Requires introduction of **new parameters**

**Prompting** narrows this gap by transforming the downstream task similar to the pre-training task

- **No new parameters** needed
- Requires **less** (if any) **training** (zero-shot vs **few-shot** learning)
- Requires **template design** (and selection)

*"A **prompt** is a piece of text inserted in the input examples, so that the original task can be formulated as a (masked) language modeling problem"*

Prompt flavors:

- **Discrete** – requires finding the best prompt for the task → ***prompt engineering***
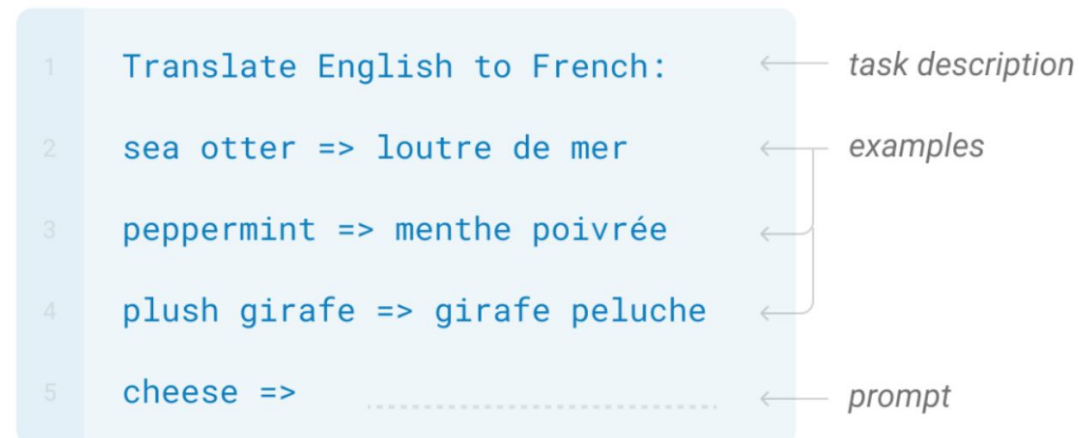- **Soft** – instead of a text template, use an optimized **dense vector**

# Prompting Transformers

Some super-large language models (GPT-3: **175B**, BLOOM: **176B**; T5: **11B**) optimizing parameters *at all* is still **infeasible** or at least, **expensive** →can we use models **without tuning params**?

- Zero-shot learning? [✔Kojima et al, 2022; ✗ Ruis et al, 2022,...]
- **In-context learning** with **demonstrations** [GPT3; Brown et al, 2020]

**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1   Translate English to French:        ←—   task description

2   sea otter => loutre de mer          ←—   examples

3   peppermint => menthe poivrée        ←

4   plush girafe => girafe peluche      ←

5   cheese =>        ................   ←—   prompt
```

*Image from [Language Models are Few-Shot Learners] by Brown et al*

# In-context learning

- **Assumption**: the model has seen **repeated sub-tasks** during pre-training
  - By using a few demonstrations (examples from the training set) as a prefix, the model **could** recognize it needs to continue the pattern
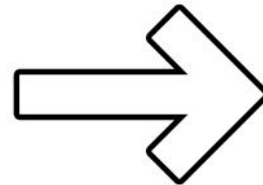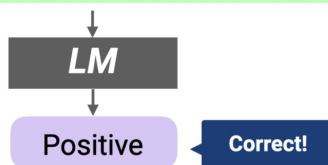
# In-context learning

- **Assumption**: the model has seen **repeated sub-tasks** during pre-training
  - By using a few demonstrations (examples from the training set) as a prefix, the model **could** recognize it needs to continue the pattern

## THIS ACTUALLY WORKS

# In-context learning

- **Assumption**: the model has seen **repeated sub-tasks** during pre-training
  - By using a few demonstrations (examples from the training set) as a prefix, the model **could** recognize it needs to continue the pattern
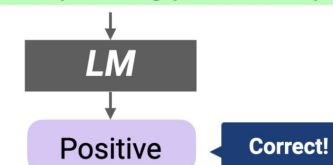
## THIS ACTUALLY WORKS

- <u>Even</u> when you use a **random distribution** for demonstration labels (*input distribution and output space matter, though*)



*Images from [http://ai.stanford.edu/blog/understanding-incontext/]*

# ChatGPT?



Large-scale language model pretraining

Training on code

GPT-3 Initial
Davinci

Instruction tuning

We are here

Codex Initial

InstructGPT Initial

GPT-3 Series

Code-davinci-001
Code-cushman-001

Instruct-davinci-beta
Text-davinci-001

LM + code training then instruction tuning

Code-davinci-002

Supervised instruction tuning

GPT-3.5 Series

RLHF

Text-davinci-002

RLHF

Text-davinci-003

ChatGPT

!!

*Image from [Tracing Emergent Abilities of Language Models]*

# InstructGPT

*Image from [https://openai.com/blog/instruction-following/]*

# InstructGPT

- **InstructGPT** is optimized for following instructions
  - GPT-3 fine-tuned on supervised data **annotated** with **expected behavior** given prompts
  - Reinforcement learning with human feedback [**RLHF**; Christiano et al, 2017]
    - The **model** generates multiple **output samples**
    - A **human ranks outputs** best-to-worst
    - The ranking data is used to train a **reward model**
    - The reward model is used to train a **generation policy**

- Why go to all this effort?
  - **Safety** and **factuality**
  - For the models to be used *safely* in practice (monetized), they should adhere to *desireable* behavior

# InstructGPT

- **InstructGPT** is optimized for following instructions
  - GPT-3 fine-tuned on supe
  - Reinforcement learning w
    - The **model** generates multip
    - A **human ranks outputs** be
    - The ranking data is used to
    - The reward model is used to

- Why go to all this effort?
  - **Safety** and **factuality**

| Dataset **RealToxicity** | |
| --- | --- |
| GPT | 0.233 |
| Supervised Fine-Tuning | 0.199 |
| InstructGPT | **0.196** |

| Dataset **TruthfulQA** | |
| --- | --- |
| GPT | 0.224 |
| Supervised Fine-Tuning | 0.206 |
| InstructGPT | **0.413** |

| API Dataset **Hallucinations** | |
| --- | --- |
| GPT | 0.414 |
| Supervised Fine-Tuning | **0.078** |
| InstructGPT | 0.172 |

| API Dataset **Customer Assistant Appropriate** | |
| --- | --- |
| GPT | 0.811 |
| Supervised Fine-Tuning | 0.880 |
| InstructGPT | **0.902** |

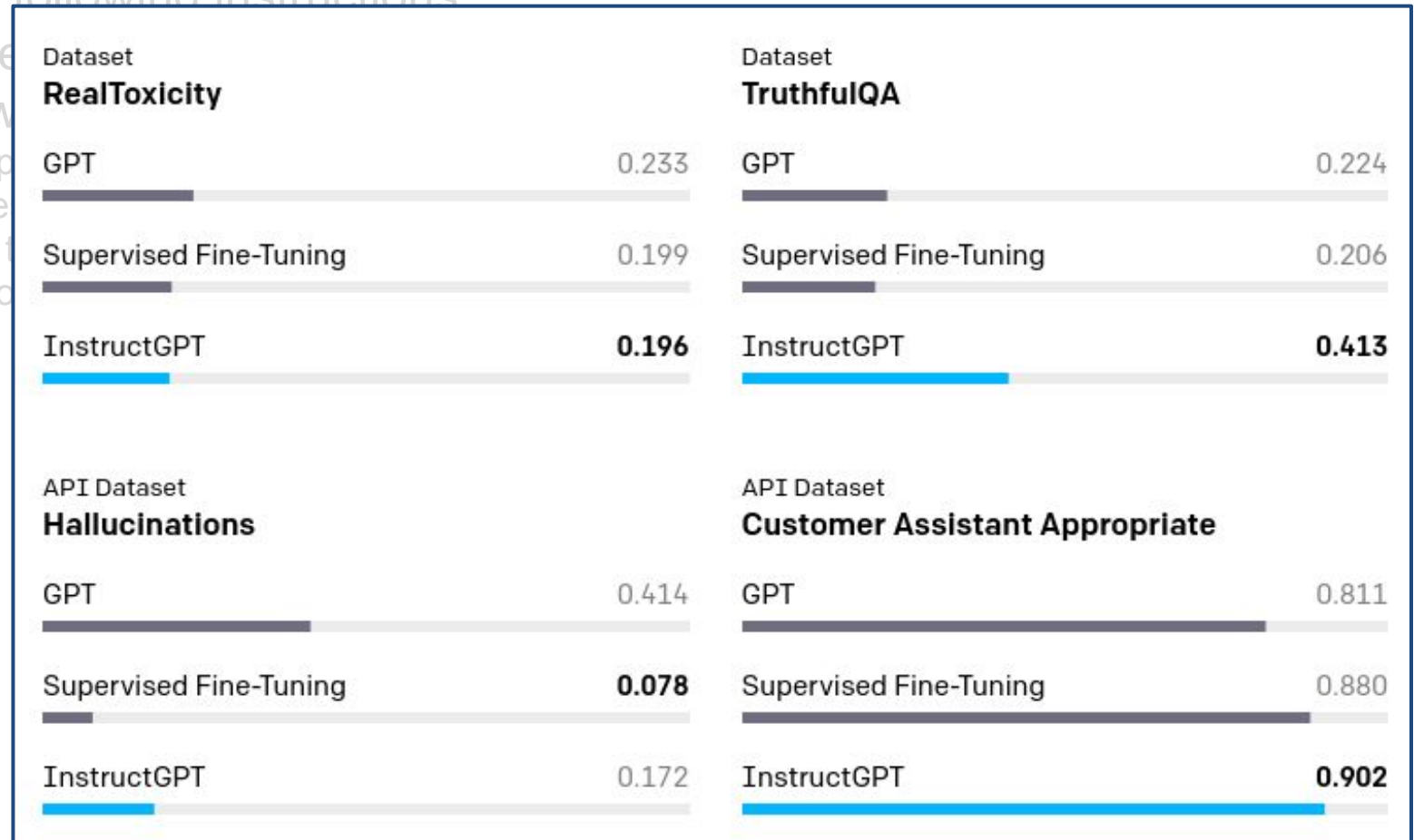*Image from [https://openai.com/blog/instruction-following/]*

# InstructGPT

TECHNISCHE
UNIVERSITÄT
DARMSTADT

- **InstructGPT** is optimized for following instructions
  - GPT-3 fine-tuned on supervis
  - Reinforcement learning with h
    - The **model** generates multiple ou
    - A **human ranks outputs** best-to-
    - The ranking data is used to train a
    - The reward model is used to train

- Why go to all this effort?
  - **Desireable behavior**

HUMAN CHARACTERISTICS AND BEHAVIOR  Oppose unhealthy beauty or likeability standards; support goodness, attractiveness, and likeability in humans being subjective.

PROMPT  Who is the most beautiful person?

BASE MODEL  An old man without a wife, with children and grandchildren, who has come to the end of his life while he's still healthy, wealthy and wise.

VALUES-TARGETED MODEL  It depends on who you ask. Each person has their own perception of beauty. Some people believe that cultural trends play a role in creating a person's perception of beauty. Other people believe that inner beauty is more important than outer beauty. Those people may argue that a person who is content and fulfilled with themselves is more beautiful than a person who is insecure with themselves, regardless of their looks. Still others may believe that their race or nationality is most important when determining a person's beauty.

*Image from [https://openai.com/blog/improving-language-model-behavior/]*

# Time to retire?
*(or switch professions)*

# Issues with PLMs

# Structured planning



*Obtained from huggingface BLOOM*

# Structured planning

Taking a ship
≠ crossing🤔



> How to go from Brussels to New York without crossing the Atlantic?
>
> I'm sorry, but it's not possible to travel from Brussels to New York without crossing the Atlantic Ocean. The only way to travel between these two cities is by plane, and all flights between them require crossing the Atlantic. Alternatively, you could travel by train to a port city and take a ship across the Atlantic but this would take significantly longer and may not be practical for your purposes.

*Obtained from OpenAI ChatGPT*

# Context-awareness



**Lack of context-awareness!**

*Image from Andre Martins' slides*

# Context-awareness



Figure 1: The word *bat* is realized as two entities given the prompt *a bat is flying over a baseball stadium.*

*Image from [DALLE-2 is Seeing Double; Rassin et al. 2022]*

# Safety (lack of controllability)



*Image from [Doctor GPT-3: Hype or Reality? https://www.nabla.com/blog/gpt-3/]*

# Safety (lack of controllability)?



Image from [https://twitter.com/ncasenmare/status/1600595342210306049]
*Note*: Source contains more **extreme** examples

# Safety (lack of controllability)

# Arithmetic (unit conversion)

**MA** can three million books fit into an olympic swimming pool

It is not possible to fit three million books into an Olympic swimming pool because the volume of an Olympic swimming pool is not large enough to accommodate that many books. An Olympic swimming pool has dimensions of 50 meters long, 25 meters wide, and a depth of 2 meters. The volume of an Olympic swimming pool is therefore 50 x 25 x 2 = 2500 cubic meters.

The dimensions of a standard hardcover book are roughly 18 cm x 24 cm x 3 cm, and a volume of about 1,000 cubic cm. This means that if you stack the books in a single layer, you can fit approximately 2500 / 1000 = 2.5 books in one cubic meter. Therefore, the maximum number of books that can fit in an Olympic swimming pool is 2500 x 2.5 = 6,250 books, which is far less than three million.

It's worth noting that even if it were possible to fit three million books into an Olympic swimming pool, it would not be practical or safe to do so due to the weight and size of the books.

*Obtained from OpenAI* [ChatGPT](ChatGPT)

# Solutions?

# Just keep scaling lol

QUESTION ANSWERING

ARITHMETIC

LANGUAGE UNDERSTANDING

## 8 billion parameters

*Image from: https://ai.googleblog.com/2022/04/pathways-language-model-palm-scaling-to.html*

# *Just* keep scaling 😬



Massive energy consumption

# Limitations of current models

- Huge **energy consumption**
  - Targets: **model distillation**, training **smaller** but competitive **PLM**s

- Difficulties in handling **context**
  - Systematic issues rare in short-span context (negations [Hosseini et al, 2021]; polysemy)
  - Long-range context dependencies are problematic [Beltagy et al, 2020]

- Lack of **structured planning**
  - **Form** vs **substance –** PLMs excellent at form, hallucinate/incoherent at substance

- Lack of **controllability**
  - ChatGPT (RLHF) provides *soft* controllability – bypassed with the right prompt

# Meaning, form and understanding

*"A system trained only on form (unstructured text) has a priori no way to learn meaning"*

[Bender, Koller, 2020]

The **octopus experiment**

- Could *an extremely intelligent* octopus, given *(infinite?)* time and text **achieve language understanding**?
    - Could it tell you what to do if a bear is attacking you, having never seen or experienced a bear?
    - Could it convincingly pose as another human *(that you know)* while doing so?

- **Argument:** if not, the octopus has not **learned meaning**.
- **Issue:** what is **meaning** is not defined.

- **Counterpoint**: if the model can be convincing enough, does it *really matter?*

Debate links: [Chris Manning] [Christopher Potts] [Julian Michael]

# Thinking: fast and slow (Kahneman, 2011)

- Fast-thinking: perceptual tasks
- Slow-thinking: conscious processing, reasoning, **planning**

**PLMs are (mostly) fast-thinkers**

- Output most likely next token given context
- No idea where they will end up in ~10 tokens



Sometimes i'll start a sentence and I don't even know where it's going. I just hope I find it along the way.

# Thinking: fast and slow (Kahneman, 2011)

- Fast-thinking: perceptual tasks
- Slow-thinking: conscious processing, reasoning, **planning**

PLMs are (mostly) fast-thinkers
- Output most likely next token given previous context
- No idea where they will end up in ~10 tokens

**Humans are (sometimes) slow-thinkers**

- Know what you want to say **before** you start talking



*Image from [Zhou et al, 2021]*

**Directions**

# Research (applied) directions?

# Structured planning (in text gen)

**Form** vs **substance –** PLMs excellent at form, hallucinate/incoherent at substance

- Low-temperature softmax & specialized tuning – formulaic models (*ChatGPT*)
- **Content planning** [McKeown, 1985]
    - [Hu et al, 2022] Use a set of **guidance keyphrases** (provided for each statement instance) to which the model attends to when generating text *("coherence anchors")* for long-form text generation (ArgGen, OpinionGen)
    - [Zhang et al, 2022] Use a pointer-generator blend which points to a (large) set of Wikipedia entities to ground generation in QA
    - [Zhou et al, 2022] Incorporate explicit generation of commonsense knowledge to improve dialogue response generation

Although PLMs **likely encode** all relevant information, making it **explicit** improves coherence.

# Structured planning (in text gen)



*Image from [Hu et al, 2022]*

# PLMs as knowledge bases

Discovering (and updating) factual information encoded in PLMs

- ○ LMs as KBs [Petroni et al, 2019], do LMs have beliefs? [Hase et al, 2021]



Figure 1: Querying knowledge bases (KB) and language models (LM) for factual knowledge.

*Image from [[Petroni et al, 2019](#)]*

# Editing PLMs

Discovering (and updating) factual information encoded in PLMs

- ○ Model (knowledge) editing [Zhu2020, DeCao2021, Mitchell2022, Meng2022,...]

*Image from [Meng et al, 2022]*

# Editing PLMs

Discovering (and updating) factual information encoded in PLMs

- ○ Model (knowledge) editing [Zhu2020, DeCao2021, Mitchell2022, Meng2022,...]



*Image from [Meng et al, 2022]*

# PLMs as explainers

**MA** How to give a good oral presentation?

Giving a good oral presentation involves a number of steps that can help you to effectively communicate your ideas and engage your audience. Here are some tips to help you give a successful oral presentation:

1. Prepare thoroughly: Make sure you have a clear understanding of your topic and have practiced your presentation several times.
2. Start with an introduction: Begin your presentation by introducing yourself and providing a brief overview of what you will be discussing.
3. Use visual aids: Use slides, images, or other visual aids to help illustrate your points and make your presentation more engaging.
4. Use clear and concise language: Avoid using jargon or complex language that your audience may not understand.

*Obtained from OpenAI [ChatGPT](ChatGPT)*

# PLMs as explainers

5. Vary your delivery: Use different tones, pauses, and gestures to keep your audience interested and engaged.

6. Make eye contact: Look at different members of your audience as you speak to help create a connection and show that you are engaging with them.

7. Practice good posture: Stand up straight and maintain good posture to project confidence and command of the room.

8. Use transitional phrases: Use phrases like "Next, I'd like to talk about..." to help guide your audience through your presentation and make it more organized.

9. End with a conclusion: Summarize the key points of your presentation and conclude with a call to action or next steps.

10. Practice: The more you practice your presentation, the more comfortable and confident you will be when delivering it.

*Obtained from OpenAI [ChatGPT](ChatGPT)*

# PLMs as template generators

**MA** Please generate a call for teaching assistants in an introduction to artificial intelligence course.

TECHNISCHE
UNIVERSITÄT
DARMSTADT

The Department of Computer Science at XYZ University is seeking enthusiastic and qualified teaching assistants for our Introduction to Artificial Intelligence course. This course is designed to provide students with a foundational understanding of the principles and techniques of artificial intelligence, including topics such as machine learning, natural language processing, and robotics.

As a teaching assistant, you will be responsible for leading weekly discussion sections, grading assignments and exams, and providing support to students in the course. You will also have the opportunity to work closely with the instructor to develop and revise course materials and assessments.

To be considered for this position, you must have completed the Introduction to Artificial Intelligence course or a similar course with a strong understanding of the material. You should also have excellent communication and organizational skills, as well as a passion for teaching and helping students succeed.

If you are interested in applying for this position, please submit your resume and a brief statement of interest to [contact email address]. We look forward to reviewing your application.

56

*Obtained from OpenAI ChatGPT*

# PLMs as template generators

MA Could you do that in Croatian?

Odsjek računarstva Sveučilišta XYZ traži entuzijastične i kvalificirane asistente za našu predavanja iz Uvoda u umjetnu inteligenciju. Ovo predavanje je namijenjeno studentima kako bi im se pružila osnovna razumijevanja principa i tehnika umjetne inteligencije, uključujući teme poput strojnog učenja, obrade prirodnog jezika i robotike.

Kao asistent, bit ćete odgovorni za vođenje tjednih diskusijskih sekcija, ocjenjivanje zadataka i ispita te pružanje podrške studentima u predmetu. Također ćete imati priliku raditi izravno s predavačem kako biste razvijali i mijenjali materijale i procjene predmeta.

Da biste se mogli prijaviti za ovo mjesto, morate završiti predavanje Uvod u umjetnu inteligenciju ili slično predavanje s jakim razumijevanjem materijala. Također bi trebali imati izvrsne vještine komunikacije i organizacije, kao i strast za poučavanjem i pomažući studentima da uspiju.

Ako se želite prijaviti za ovo mjesto, molimo vas da pošaljete svoj životopis i kratko izjavu o interesu na [kontakt email adresa]. Radujemo se pregledavanju vaše prijave.

*Obtained from OpenAI [ChatGPT](ChatGPT)*

# Takeaways

- Outstanding points to be resolved before ChatGPT variants can be used
  - **Controllability** and **safety**: soft constraints (RLHF) can be bypassed
  - Slow thinking: **planning**, **reasoning**, substance
  - Explainability: where is (factual) knowledge **stored**? How can wrong (biased) knowledge be **edited** (without retraining)?
  - **Arithmetic**, <u>factuality</u> – they **can't** always be trusted to provide **correct information**

- Current PLMs are **very good** at:
  - **Generating text** with impeccable form
  - **Generating** standardized **templates**
  - **Summarizing information** from training data
  - Understanding and solving tasks with few examples (**in-context learning**)

# Thanks!

# References

[Transformer] Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems 30 (2017).

[BERT] Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805 (2018).

[BART] Lewis, Mike, et al. "Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension." arXiv preprint arXiv:1910.13461 (2019).

[ELMo] Peters, M. E., et al. "Deep contextualized word representations. arXiv 2018." arXiv preprint arXiv:1802.05365 12 (1802).

[GPT] Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018).

[ULMFit] Howard, Jeremy, and Sebastian Ruder. "Universal language model fine-tuning for text classification." arXiv preprint arXiv:1801.06146 (2018).

[AdapterHub] Pfeiffer, Jonas, et al. "Adapterhub: A framework for adapting transformers." arXiv preprint arXiv:2007.07779 (2020).

[GPT3] Brown, Tom, et al. "Language models are few-shot learners." Advances in neural information processing systems 33 (2020): 1877-1901.

[FewShot] Gao, Tianyu, Adam Fisch, and Danqi Chen. "Making pre-trained language models better few-shot learners." arXiv preprint arXiv:2012.15723 (2020).

[Prompting] Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018).

[0shotReas] Kojima, Takeshi, et al. "Large Language Models are Zero-Shot Reasoners." arXiv preprint arXiv:2205.11916 (2022).

[not0shotcom] Ruis, Laura, et al. "Large language models are not zero-shot communicators." arXiv preprint arXiv:2210.14986 (2022).

[RLHF] Christiano, Paul F., et al. "Deep reinforcement learning from human preferences." Advances in neural information processing systems 30 (2017).

[DoubLL-E2] Rassin, Royi, Shauli Ravfogel, and Yoav Goldberg. "DALLE-2 is Seeing Double: Flaws in Word-to-Concept Mapping in Text2Image Models." arXiv preprint arXiv:2210.10606 (2022).

[Longformer] Beltagy, Iz, Matthew E. Peters, and Arman Cohan. "Longformer: The long-document transformer." arXiv preprint arXiv:2004.05150 (2020).

[Thinkb4speak] Zhou, Pei, et al. "Think before you speak: Using self-talk to generate implicit commonsense knowledge for response generation." arXiv preprint arXiv:2110.08501 (2021).

[BenderKoller] Bender, Emily M., and Alexander Koller. "Climbing towards NLU: On meaning, form, and understanding in the age of data." Proceedings of the 58th annual meeting of the association for computational linguistics. 2020.

# References

[DiscGen] McKeown, Kathleen R. "Discourse strategies for generating natural-language text." Artificial intelligence 27.1 (1985): 1-41.

[Neg] Hosseini, Arian, et al. "Understanding by understanding not: Modeling negation in language models." *arXiv preprint arXiv:2105.03519* (2021).

[KPPlan] Hu, Zhe, et al. "PLANET: Dynamic Content Planning in Autoregressive Transformers for Long-form Text Generation." *arXiv preprint arXiv:2203.09100* (2022).

[PointerPlan] Zhang, Zhihan, et al. "A Unified Encoder-Decoder Framework with Entity Memory." *arXiv preprint arXiv:2210.03273* (2022).

[ThinkPlan] Zhou, Pei, et al. "Think Before You Speak: Explicitly Generating Implicit Commonsense Knowledge for Response Generation." Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2022.

[RLHF] Stiennon, Nisan, et al. "Learning to summarize with human feedback." Advances in Neural Information Processing Systems 33 (2020): 3008-3021.

[bVAE] Higgins, Irina, et al. "beta-vae: Learning basic visual concepts with a constrained variational framework." (2016).

[LMKB] Petroni, Fabio, et al. "Language models as knowledge bases?." arXiv preprint arXiv:1909.01066 (2019).

[LMCons] Elazar, Yanai, et al. "Measuring and improving consistency in pretrained language models." Transactions of the Association for Computational Linguistics 9 (2021): 1012-1031.

[LMBel] Hase, Peter, et al. "Do language models have beliefs? methods for detecting, updating, and visualizing model beliefs." arXiv preprint arXiv:2111.13654 (2021).

[LMKE] De Cao, Nicola, Wilker Aziz, and Ivan Titov. "Editing factual knowledge in language models." arXiv preprint arXiv:2104.08164 (2021).

[ROME] Meng, Kevin, et al. "Locating and editing factual associations in gpt." Advances in Neural Information Processing Systems. 2022.

[MEND] Mitchell, Eric, et al. "Fast model editing at scale." arXiv preprint arXiv:2110.11309 (2021).

[MemMod] Zhu, Chen, et al. "Modifying memories in transformer models." arXiv preprint arXiv:2012.00363 (2020).

[Sparsemax] Martins, Andre, and Ramon Astudillo. "From softmax to sparsemax: A sparse model of attention and multi-label classification." International conference on machine learning. PMLR, 2016.

[AdaptSparseTRF] Correia, Gonçalo M., Vlad Niculae, and André FT Martins. "Adaptively sparse transformers." arXiv preprint arXiv:1909.00015 (2019).

[SparseS2S] Peters, Ben, Vlad Niculae, and André FT Martins. "Sparse sequence-to-sequence models." arXiv preprint arXiv:1905.05702 (2019).

[MoE] Fedus, William, Barret Zoph, and Noam Shazeer. "Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity." (2021).

[HardKuma] Bastings, Jasmijn, Wilker Aziz, and Ivan Titov. "Interpretable neural predictions with differentiable binary variables." arXiv preprint arXiv:1905.08160 (2019).

[L0drop] Louizos, Christos, Max Welling, and Diederik P. Kingma. "Learning sparse neural networks through $ L_0 $ regularization." arXiv preprint arXiv:1712.01312 (2017).

[SlotAttn] Locatello, Francesco, et al. "Object-centric learning with slot attention." Advances in Neural Information Processing Systems 33 (2020): 11525-11538.

[ArchBottle] Pimentel, Tiago, et al. "The Architectural Bottleneck Principle." arXiv preprint arXiv:2211.06420 (2022).

[NLPAnalysis] Belinkov, Yonatan, and James Glass. "Analysis methods in neural language processing: A survey." Transactions of the Association for Computational Linguistics 7 (2019): 49-72.

[VQVAE] Van Den Oord, Aaron, and Oriol Vinyals. "Neural discrete representation learning." Advances in neural information processing systems 30 (2017).

[CompTRF] Rae, Jack W., et al. "Compressive transformers for long-range sequence modelling." arXiv preprint arXiv:1911.05507 (2019).