

CONFIDENTIAL: Oxmedica Ltd. 10347756



Specialist Topic

Large Language Models



CONFIDENTIAL: Oxmedica Ltd. 10347756



Specialist Topic: Large Language Models

Day 1: From Neural Networks to LLMs



Day 2: LLM applications



Large Language Models: Day 1

Day 1: From Neural Networks to LLMs

08:30 - 10:30

Introduction to LLMs

- Recap on NNs
- History of LMs
- Future LLMs and AGI
- State-of-the-art evals
- Evals and visualisations exercise

11:00 - 11:30

How do LLMs work?

- Finish evals
- Technical lecture

11:30 - 12:00

Guest Lecture: AI in medicine

12:15 - 12:45

Limitations of SOTA LLMs

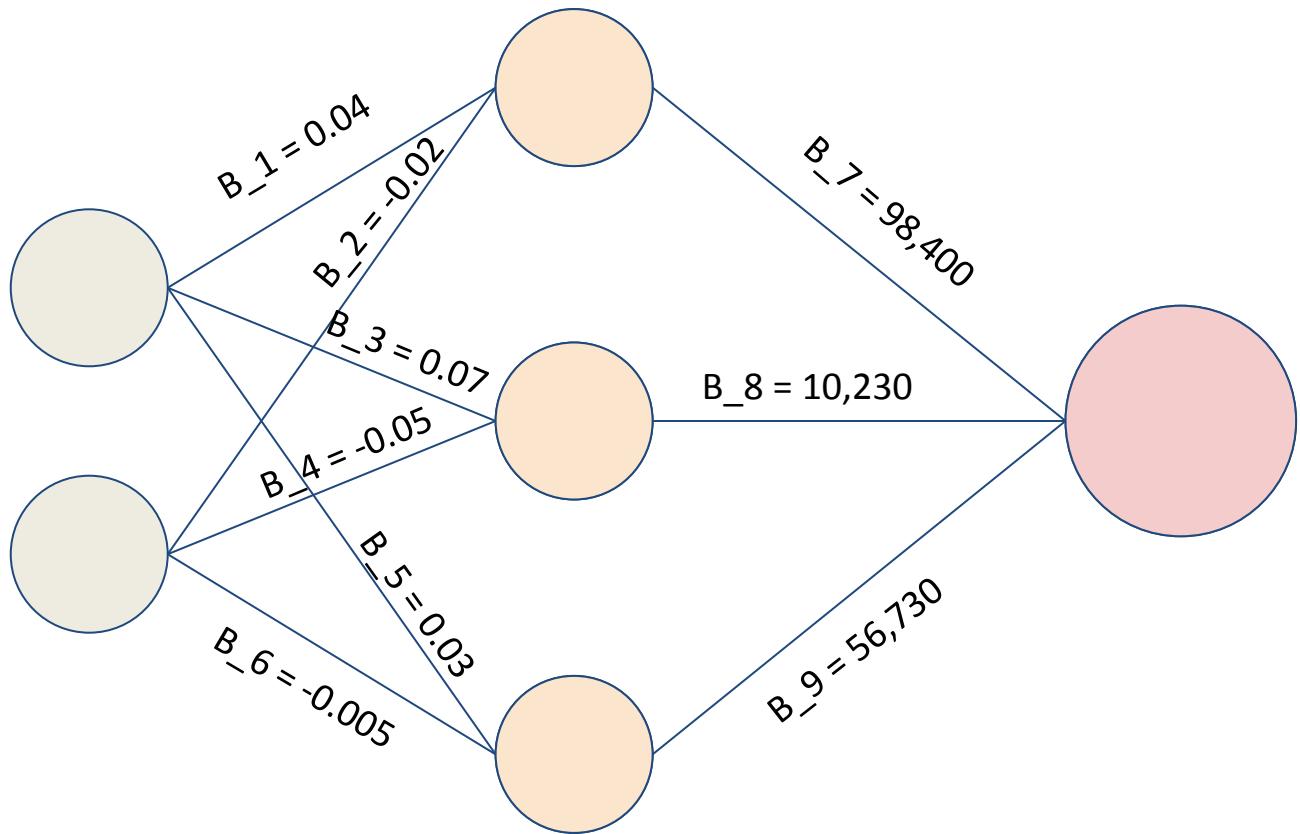


Neural Networks Recap

Neural Networks Task: House price prediction 2

Area $x_1 = 40$

Age $x_2 = 30$



How do we work this out?

Whiteboard

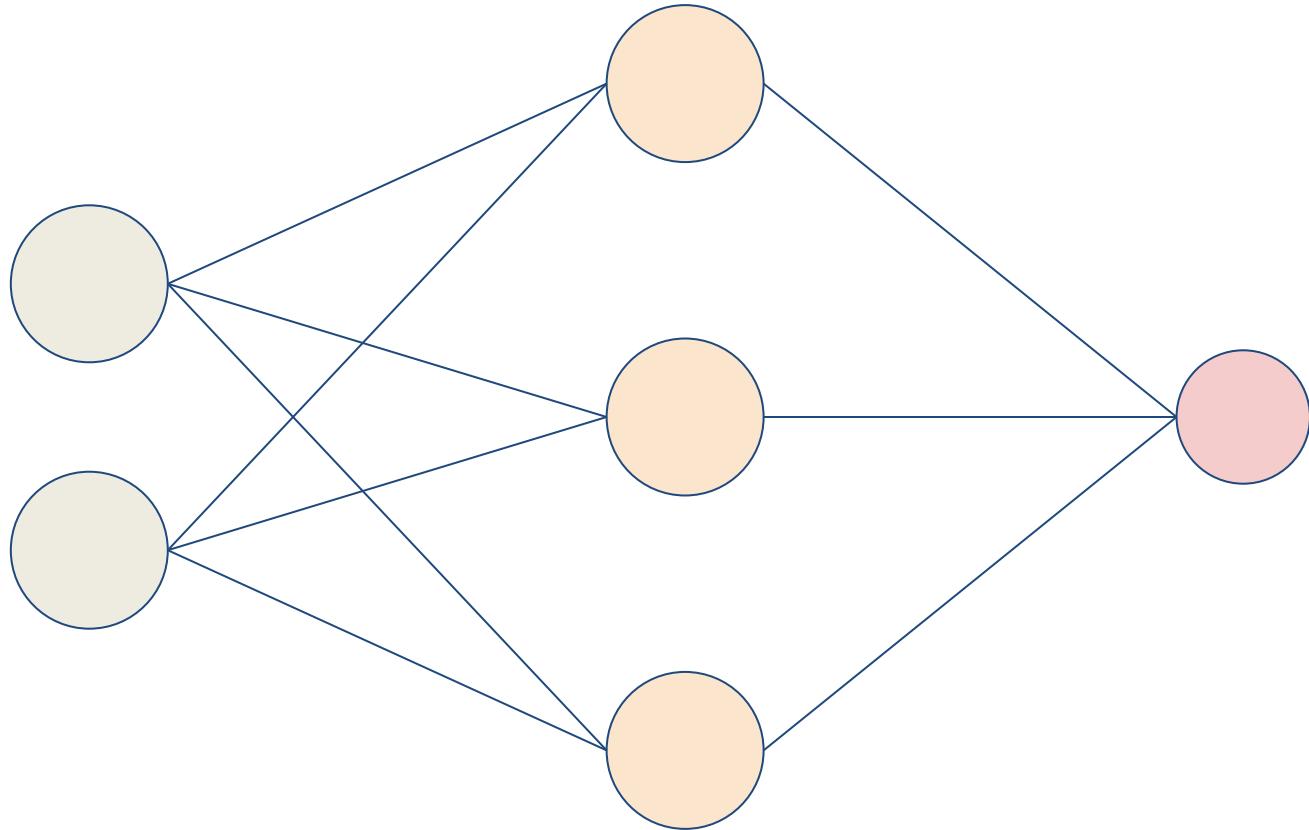


مؤسسة الملك عبد العزيز ورجاله للموهبة والإبداع
King Abdullah Bin Abdulaziz Foundation for Science & Technology

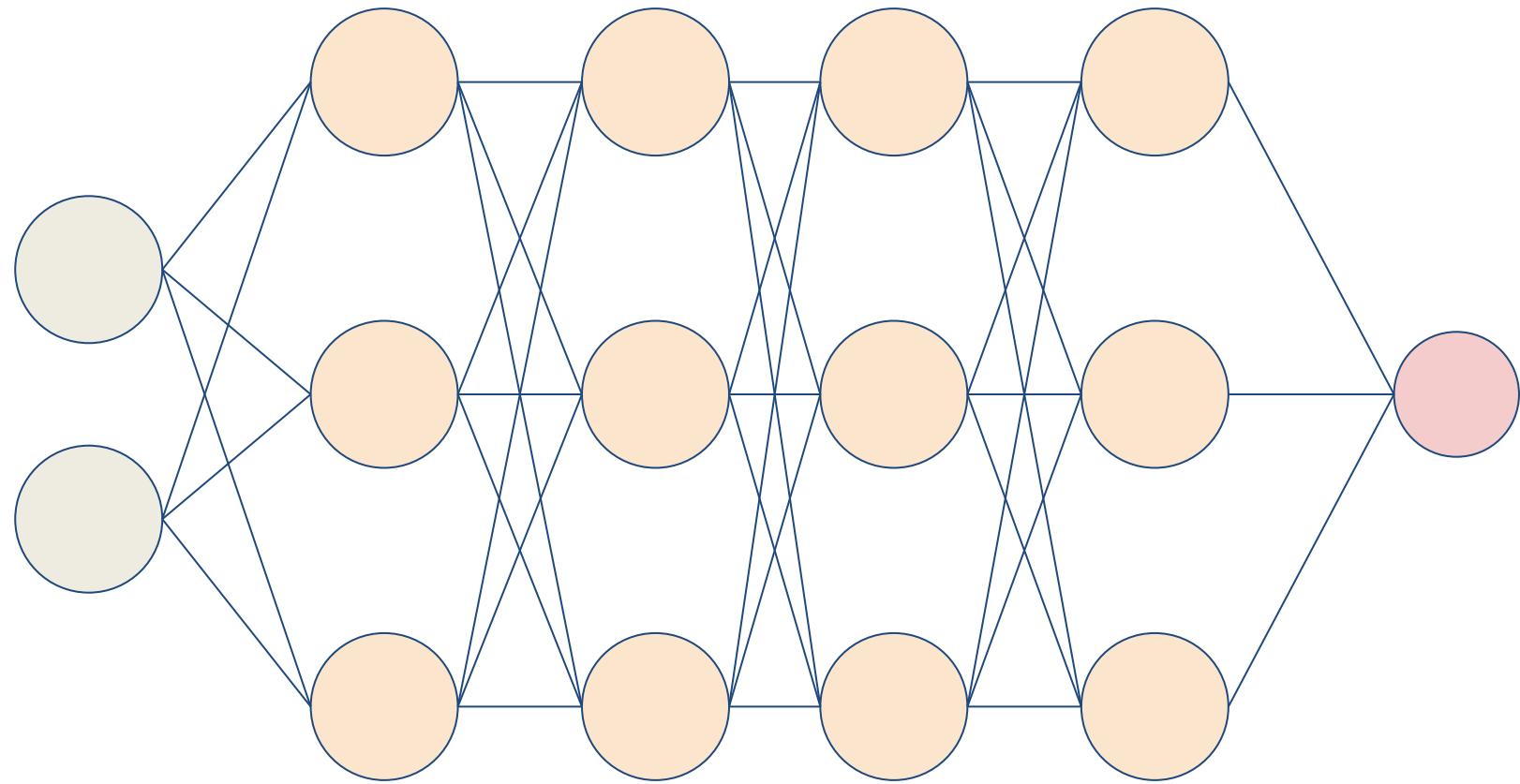
CONFIDENTIAL: Oxmedica Ltd. 10347756



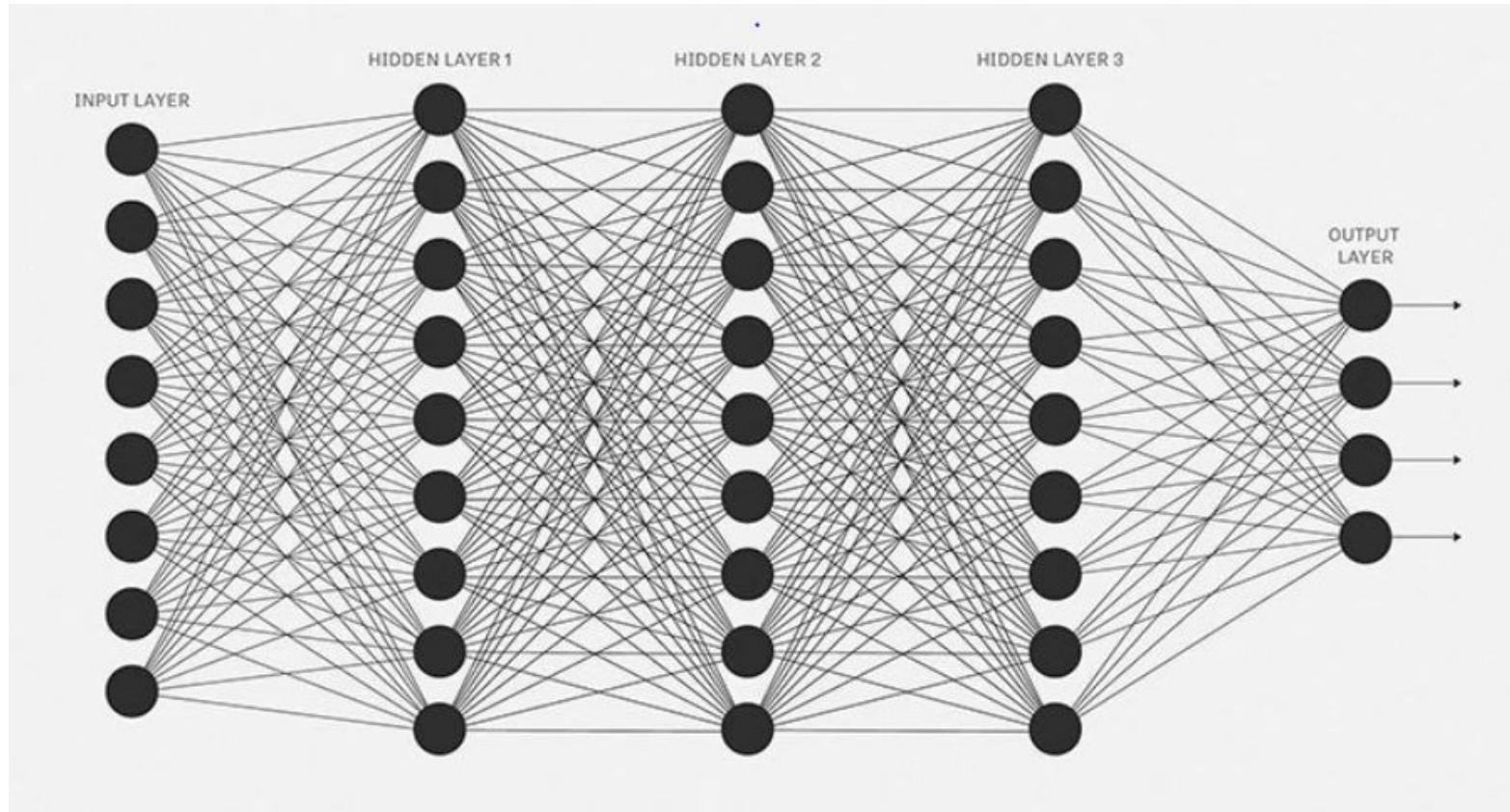
Big Neural Networks



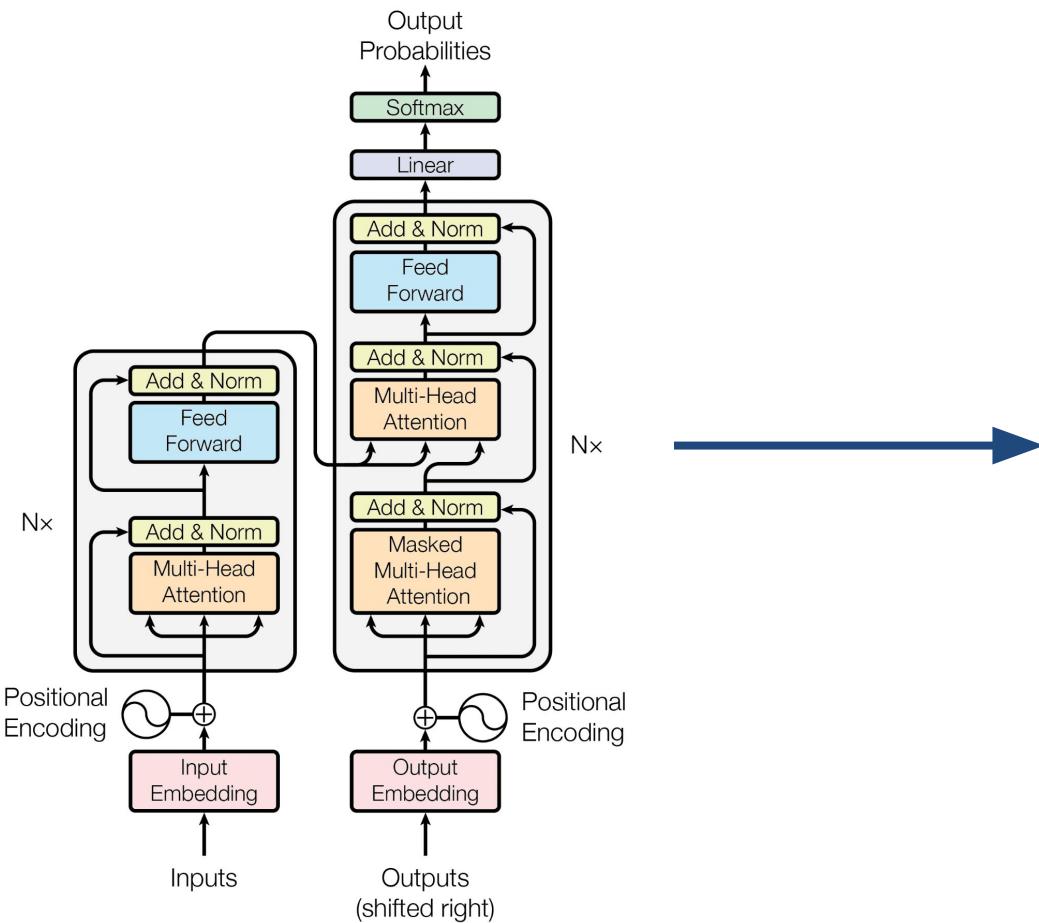
Big Neural Networks: Connect everything



Big Neural Networks: bigger...



Big Neural Networks: bigger...



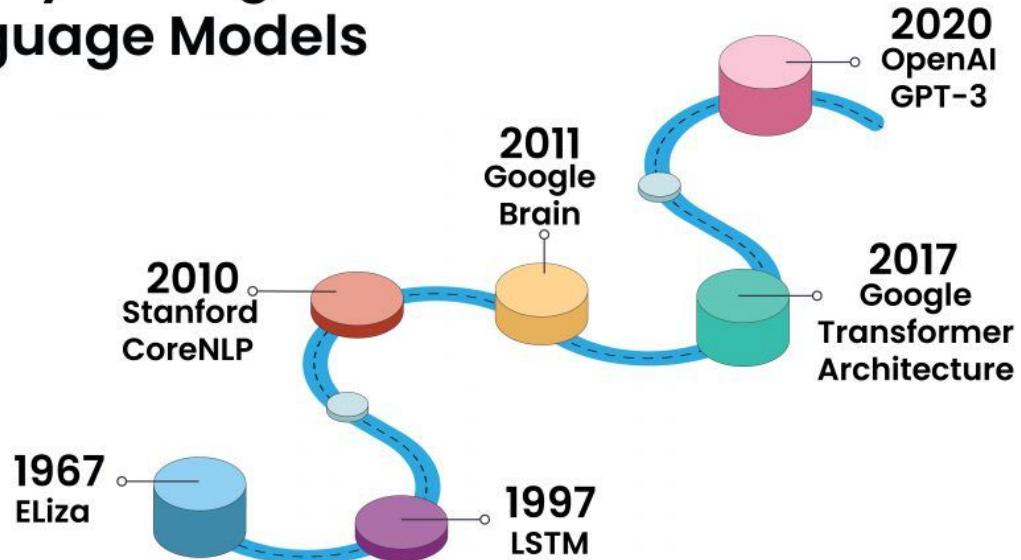
Gives you predicted probabilities for all possible next tokens (subwords)

A Brief History of Language Models



A brief history

History of Large Language Models



[Source](#)

A brief history: The ELIZA model

EEEEEE	LL	III	ZZZZZZZ	AAAAA
EE	LL	II	ZZ	AA AA
EEEEEE	LL	II	ZZZ	AAAAAAA
EE	LL	II	ZZ	AA AA
EEEEEE	LLLLL	III	ZZZZZZZ	AA AA

[Source](#)



مؤسسة الملك عبد العزيز ورجاله للموهبة والإبداع
King Abdullah II International Foundation for Guidance & Capacity

CONFIDENTIAL: Oxmedica Ltd. 10347756



A brief history: Statistical Models (n-gram models)

Example bigram table

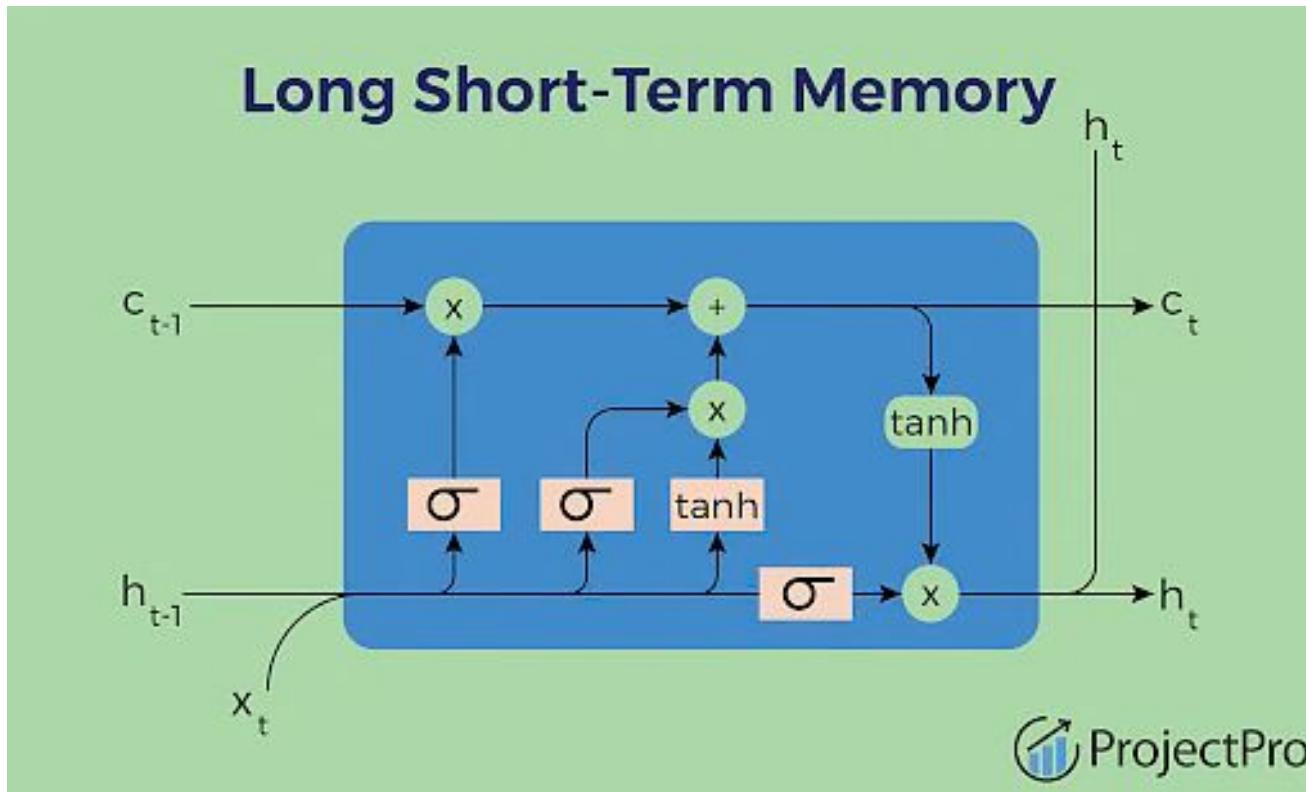
Word 1\Word 2	the	a	dog	in	park	today	...
the	0.0	0.0	0.1	0.0	0.1	0.0	...
a	0.0	0.0	0.1	0.0	0.1	0.0	...
dog	0.1	0.1	0.0	0.1	0.1	0.1	...
in	0.2	0.1	0.0	0.0	0.0	0.0	...
park	0.1	0.1	0.0	0.1	0.0	0.1	...
today	0.05	0.0	0.0	0.1	0.0	0.0	...

[Source](#)



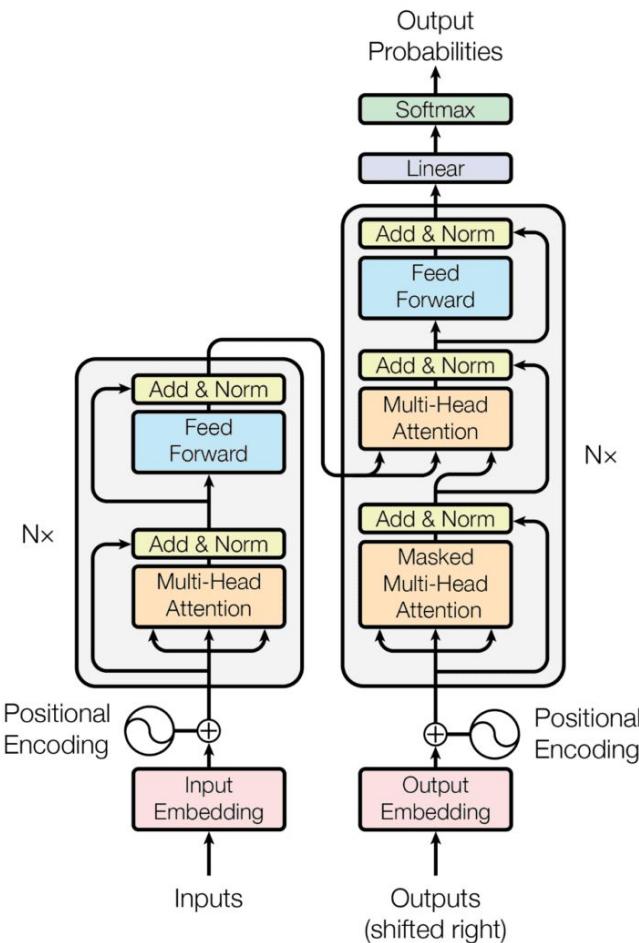
CONFIDENTIAL: Oxmedica Ltd. 10347756

LSTM neural networks



Source

The Transformer Revolution



Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention

[Source](#)



CONFIDENTIAL: Oxmedica Ltd. 10347756

The Transformer Revolution

Attention is all you need

[PDF] neurips.cc

A Vaswani, N Shazeer, N Parmar, J Uszkoreit, L Jones, AN Gomez, Ł Kaiser, I Polosukhin

Advances in neural information processing systems, 2017 • proceedings.neurips.cc

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder and decoder configuration. The best performing such models also connect the encoder and decoder through an attention mechanism. We propose a novel, simple network architecture based solely on an attention mechanism, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more

SHOW MORE ▾

☆ Save ⚡ Cite Cited by 124578 Related articles All 91 versions Import into BibTeX ☰

Showing the best result for this search. See all results

This is a **huge**
number of
citations!

Source



CONFIDENTIAL: Oxmedica Ltd. 10347756

Attention Is All You Need

Contents [hide](#)

(Top)

Authors

References

Article Talk

From Wikipedia, the free encyclopedia

"Attention Is All You Need" is a 2017 landmark^{[1][2]} research paper authored by eight scientists working at Google, that introduced a new deep learning architecture known as the transformer based on attention mechanisms proposed by Bahdanau *et al.* in 2014. It is considered by some to be a founding paper for modern artificial intelligence, as transformers became the main architecture of large language models like those based on GPT.^{[3][4]} At the time, the focus of the research was on improving Seq2seq techniques for machine translation, but even in their paper the authors saw the potential for other tasks like question answering and for what is now called multimodal Generative AI.^[5]

The paper's title is a reference to the song "All You Need Is Love" by the Beatles.^[6]

As of 2024, the paper has been cited more than 100,000 times.^[7]

Authors [edit]

The authors of the paper are: [Ashish Vaswani](#), Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, [Aidan Gomez](#), Lukasz Kaiser, and Illia Polosukhin. All eight authors were "equal contributors" to the paper; the listed order was randomized. The *Wired* article highlights the group's diversity:^[6]

Six of the eight authors were born outside the United States; the other two are children of two green-card-carrying Germans who were temporarily in California and a first-generation American whose family had fled persecution, respectively.

By 2023, all eight authors had left Google and founded their own AI start-ups (except Łukasz Kaiser, who joined OpenAI). [6][7]

The diagram illustrates the multi-head attention mechanism. It starts with an **Input Embedding** layer, which is combined with **Positional Encoding** via a residual connection. The combined representation passes through a **Masked Multi-Head Attention** layer, which consists of multiple parallel attention heads. The outputs of these heads are concatenated and passed through another residual connection. This is followed by a **Multi-Head Attention** layer, which also uses multiple parallel attention heads. Its output is concatenated and passed through a third residual connection. Finally, the sequence passes through a **Feed Forward** layer and another residual connection. The entire process is enclosed in a large box labeled **Add&Norm**, which represents the addition of the input to the output of the layer and the normalization of the result.

An illustration of main components of the transformer model from the paper

Appearance hide

Text

10

Standards

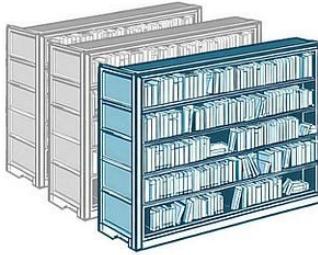
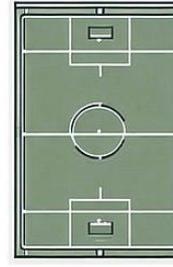
Width

Standard

Wide

And now... **HUGE** transformer models

GPT4 Model Estimates

Training Size	Compute Size	Model Size
# of Book shelves for 13T tokens	Compute time for 2.15 e25 FLOPs	Size of Excel Sheet for 1.8T params
650 kms Long line of Library Shelves  100000 tokens per Book 100 Books per shelf 2 Shelves per meter	7 million years On mid-size Laptop (100GFLOPs)  100GLOPs per second	30,000 Football Fields sized Excel Sheet  1x1 cm per Excel cell 100 x 60 meters Field Size

Source: <https://the-decoder.com/gpt-4-architecture-datasets-costs-and-more-leaked>

The Future of LLMs and AGI

Near-Term Future

1

Multimodality

2

Agents

3

Usability

[Source](#)



CONFIDENTIAL: Oxmedica Ltd. 10347756

Near-Term Future: Multimodality



What is the football club shown in this image?



The football club shown in the image is Newcastle United. This can be identified by the black and white striped kit, which is characteristic of Newcastle United's home jersey.

The player in the image appears to be Alexander Isak, a forward for Newcastle United. Alexander Isak is considered a talented player and has made significant contributions to the team. His performances have been well-regarded, and he is seen as one of the promising players in the squad.



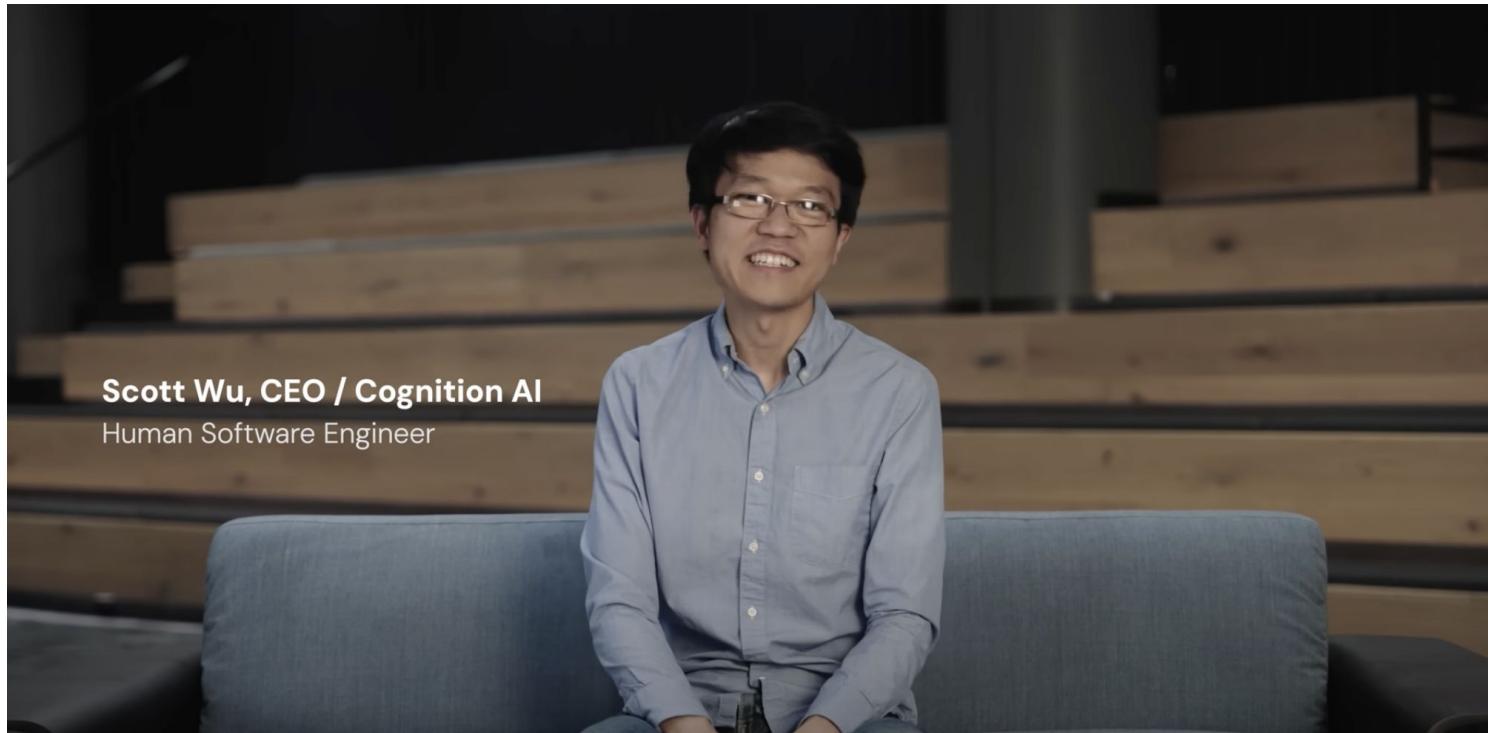
CONFIDENTIAL: Oxmedica Ltd. 10347756

Near-Term Future: Agents

- A system that can **take actions** on a user's behalf
- ChatGPT is a good knowledge source but bad at taking actions
- Model's need to improve to get good at multi-step, agents work



Example agent: Devin



<https://www.youtube.com/watch?v=fjHtjT7GO1c>

Near-Term Future: Agents

But they're not that good yet...



"First AI Software Engineer" Creators Are Accused of Lying

A screenshot of an article from The Verge. The article is titled "'First AI Software Engineer' Creators Are Accused of Lying". It features a photo of a man, Devin, and a video player showing him. The author is Gloria Levine, Senior Editor, published on 16 April 2024. The article discusses the controversy surrounding Devin's creation. A sidebar on the right is for "FIND ART & DESIGN OUTSOURCE" with a "HIRE QUALIFIED STUDIOS" button. A cookie consent banner at the bottom asks for consent to use cookies.

[Source](#)



CONFIDENTIAL: Oxmedica Ltd. 10347756

Near-Term Future: Usability

- Current models aren't very usable!
- You have to type in text... very slow and painful!

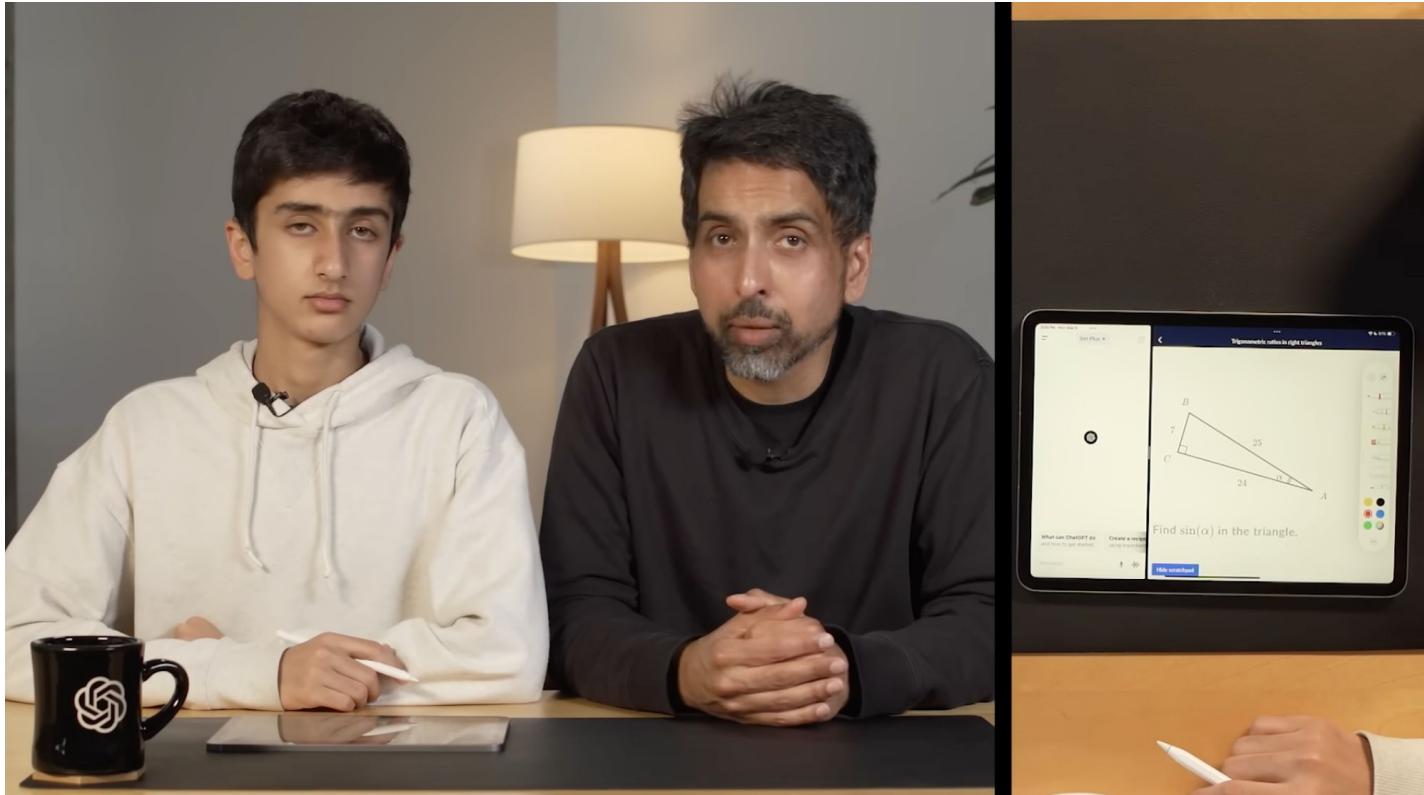


Near-Term Future: Usability



<https://www.youtube.com/watch?v=c2DFg53Zhvw>

Near-Term Future: Good but far from great...



https://www.youtube.com/watch?v=_nSmkyDNulk

Long-term future: Some definitions

Artificial General Intelligence (AGI): No agreed definition but an intelligence that is roughly human-level across the board.

Sam Altman (OpenAI CEO) → [here](#)

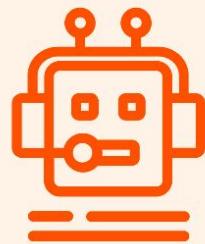
Artificial Superintelligence (ASI): An intelligence which is far beyond human level. Far harder to achieve.

[Source](#)



What is AI?

ANI vs. AGI vs. ASI



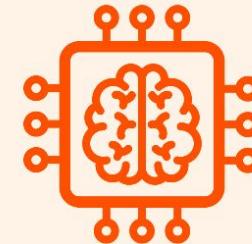
Artificial narrow intelligence (ANI)

Designed to perform specific tasks



Artificial general intelligence (AGI)

Can behave in a human-like way across all tasks



Artificial super intelligence (ASI)

Smarter than humans—the stuff of sci-fi

zapier

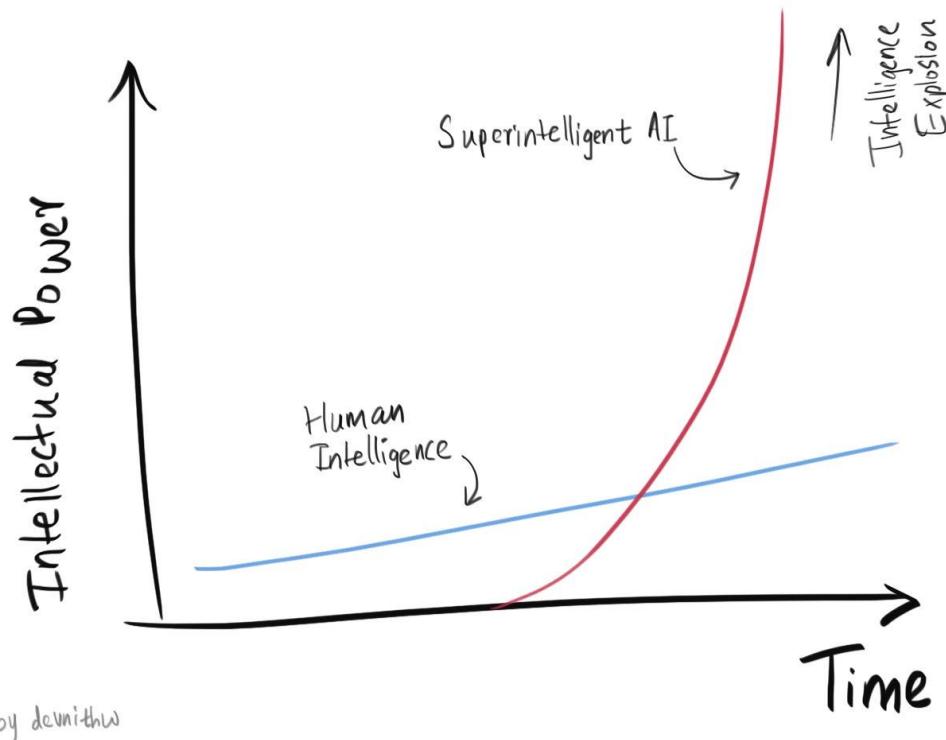
[Source](#)



CONFIDENTIAL: Oxmedica Ltd. 10347756

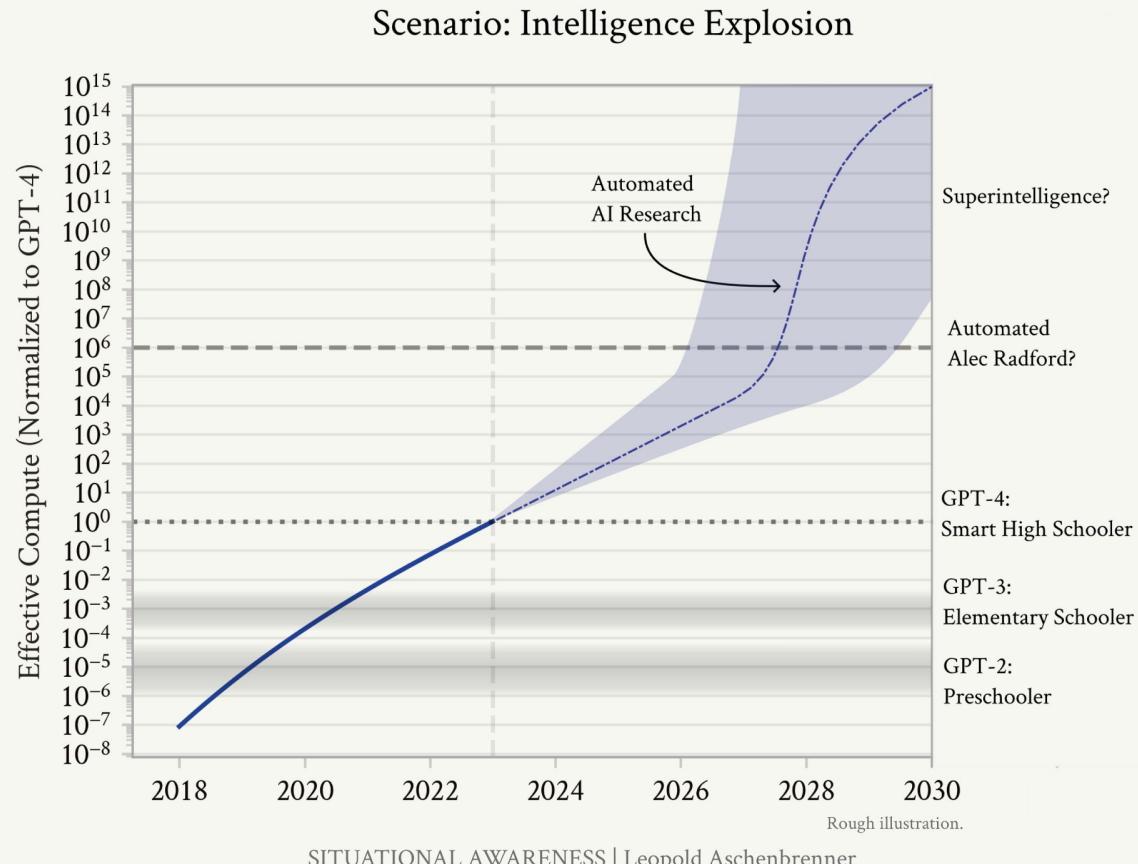
The intelligence explosion?

Intelligence Explosion



[Source](#)

The intelligence explosion?



[Source](#)

Sir Demis Hassabis' thoughts...



<https://www.youtube.com/watch?v=BGxiufHVd0>



CONFIDENTIAL: Oxmedica Ltd. 10347756



Dario Amodei's thoughts...



<https://www.youtube.com/watch?v=wo4o09lKAQQ>

What might be missing from the story of an intelligence explosion?

Discussion



مَوْهِدَة

مؤسسة الملك عبد العزيز ورجاله للموهبة والإبداع
King Abdullah II International Foundation for Guidance & Leadership

CONFIDENTIAL: Oxmedica Ltd. 10347756



State-of-the-art Evaluations



What are Language Model Evaluations?

Discussion



The Turing Test



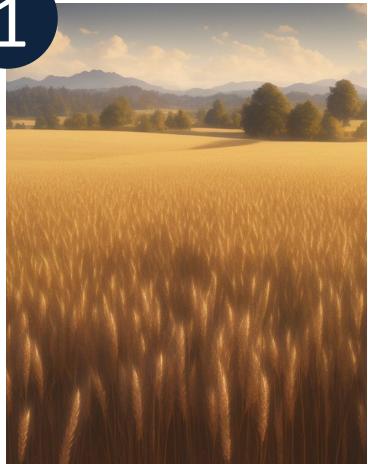
مohnat

مؤسسة الملك عبد العزيز ورجاله للموهبة والإبداع
King Abdulaziz & Sons Research Foundation for Genetics & Society

CONFIDENTIAL: Oxmedica Ltd. 10347756



1



[Link to original source](#)

2



[Link to original source](#)

1



[Link to original source](#)

2



[Link to original source](#)

1

Five English proverbs

1. Many hands make light work
2. Strike while the iron is hot
3. The grass is always greener
4. Don't judge a book by its cover
5. An apple a day keeps the doctor away

2

The top five English proverbs:

1. Actions speak louder than words.
2. A silent cat catches no mice.
3. A picture is worth a thousand words.
4. When in Rome, do as the Romans do.
5. The stone that rolls grows no moss.

Modern Eval: What makes a good eval?

Discussion



Modern Eval: Things to consider

- 1 Difficulty
- 2 Generalisability
- 3 Memorisation risk...
- 4 Superhuman abilities...?



Modern Eval: Lessons from the trenches

Preprint. Under review.

Lessons from the Trenches on Reproducible Evaluation of Language Models

Stella Biderman^{1*}, Hailey Schoelkopf^{1*}, Lintang Sutawika^{1*},
Leo Gao¹, Jonathan Tow², Baber Abbasi¹, Alham Fikri Aji³, Pawan Sasanka
Ammanamanchi⁴, Sidney Black¹, Jordan Clive⁵, Anthony DiPofi¹, Julen Etxaniz⁶, Benjamin
Fattori¹, Jessica Zosa Forde⁷, Charles Foster⁸, Jeffrey Hsu⁹, Mimansa Jaiswal¹⁰, Wilson Y.
Lee¹¹, Haonan Li^{3,12}, Charles Lovering¹³, Niklas Muennighoff¹⁴, Ellie Pavlick⁷, Jason
Phang^{1,15}, Aviya Skowron¹, Samson Tan¹⁶, Xiangru Tang¹⁷, Kevin A. Wang⁷, Genta Indra
Winata¹⁸, François Yvon¹⁹, and Andy Zou²⁰

¹EleutherAI, ²Stability AI, ³MBZUAI, ⁴IIIT Hyderabad, ⁵Chattermill AI, ⁶HiTZ Center - Ixa,
UPV/EHU, ⁷Brown University, ⁸Finetune, ⁹Ivy Natal, ¹⁰University of Michigan, ¹¹HubSpot,
¹²LibrAI, ¹³Kensho, ¹⁴Contextual AI, ¹⁵New York University, ¹⁶Amazon, ¹⁷Yale University,
¹⁸HKUST, ¹⁹Sorbonne University, ²⁰CMU

*Equal Contribution

Abstract

<https://arxiv.org/pdf/2405.14782>



Modern Eval / Benchmarks

Task in pairs: 10 mins to research this eval/benchmark

- What does this eval test?
- Give an example problem - how hard is it?
- Overall is this benchmark good?

GSM8K	MGSM
MMLU	MMMU
BBH	MathVista
ARC / The ARC Challenge	WMDP Benchmark
HumanEval (Code)	GPQA



Individual Task: 10 minutes to read this paper

LINGOLY: A Benchmark of Olympiad-Level Linguistic Reasoning Puzzles in Low-Resource and Extinct Languages

Andrew Bean^{1*} Simi Hellsten^{1,2}
Harry Mayne¹ Jabez Magomere¹ Ethan A. Chi³ Ryan Chi³
Scott A. Hale^{1,4} Hannah Rose Kirk¹

¹University of Oxford ²United Kingdom Linguistics Olympiad
³Stanford University ⁴Meedan

Abstract

In this paper, we present the LINGOLY benchmark, a novel benchmark for advanced reasoning abilities in large language models. Using challenging Linguistic Olympiad puzzles, we evaluate (i) capabilities for in-context identification and

<https://arxiv.org/pdf/2406.06196.pdf>



CONFIDENTIAL: Oxmedica Ltd. 10347756



Evals and Visualisation Task



Evals Challenge: 30 minutes

- Most modern evals are becoming memorised by the language models and are thus poor tests of true ability.
- Your task is to design your own eval. It can be really specific (only testing a certain type of behaviour) or really general but it should, in some way, test language model ‘intelligence’
- You need to design 7-10 questions for your eval and work out how to score the model responses.
- Benchmark 5 different language models on this benchmark and present the results in a table (remember how to make nice visualisations!)

