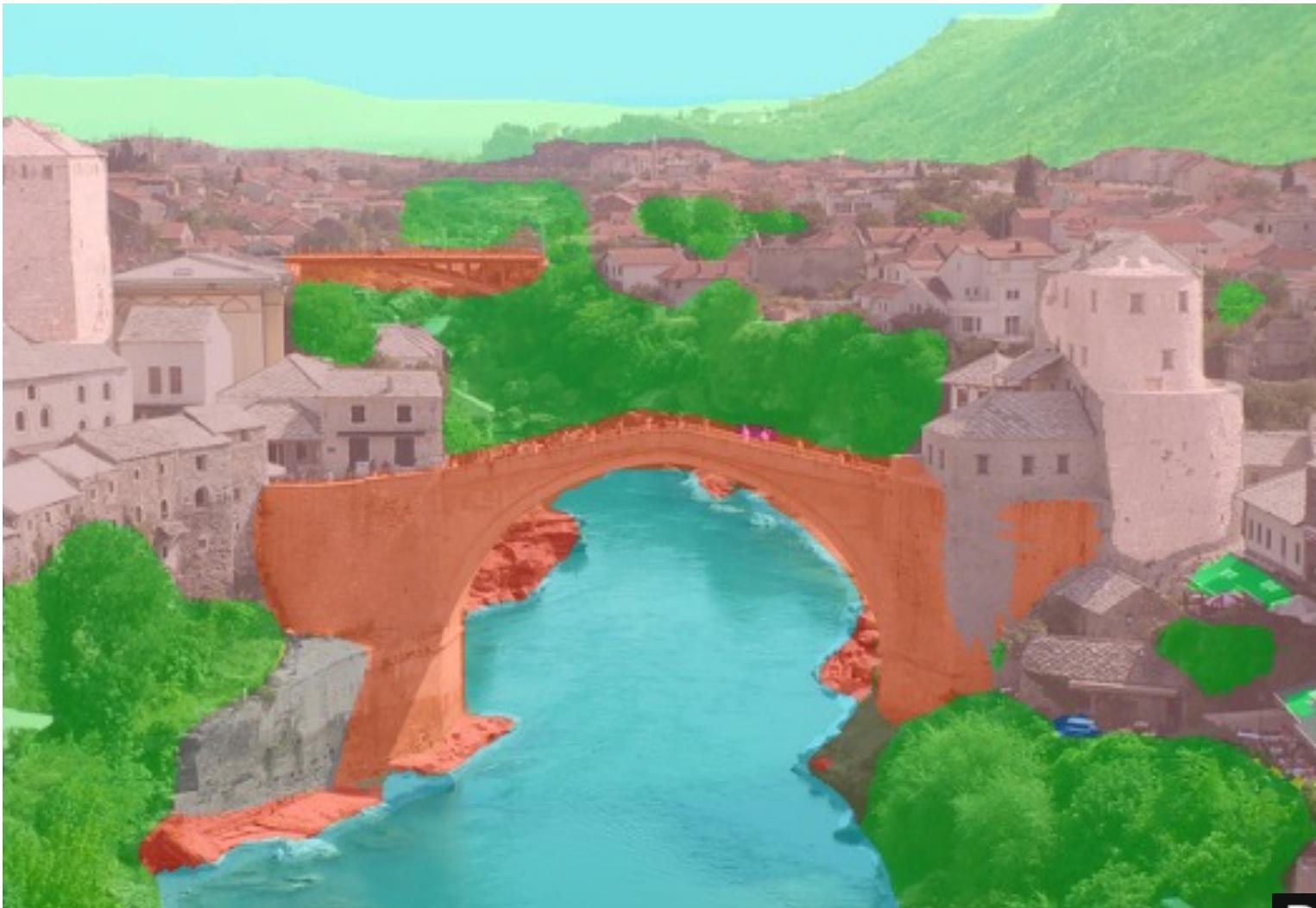


# Why Foundational AI Matters?

Simon Lucey



**AUSTRALIAN  
INSTITUTE FOR  
MACHINE LEARNING**



M. Oquab, et al. "DINOv2: Learning robust visual features without supervision." In arXiv 2023.

**DINOv2**  
Research by Meta AI



A. Kirillov, et al. "Segment Anything." In CVPR 2023.

User: What is unusual about this image?



GPT4: The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.

TECH · ELON MUSK

# Elon Musk's just fired up Colossus—the world's largest Nvidia GPU supercomputer built in just three months from start to finish

BY CHRISTIAAN HETZNER

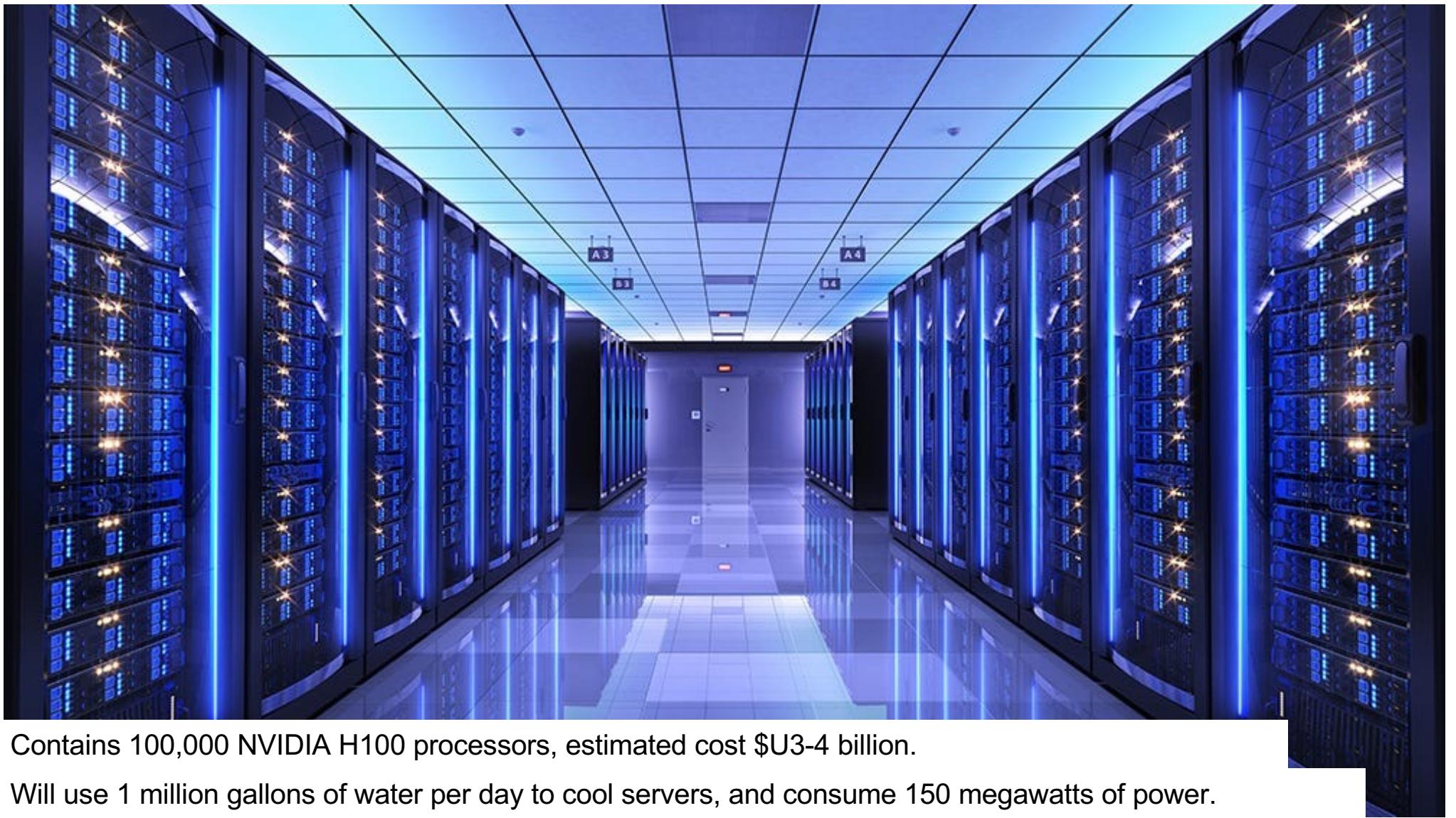
September 3, 2024 at 11:48 PM GMT+9:30



<https://fortune.com/2024/09/03/elon-musk-xai-nvidia-colossus/>



xAI founder Elon Musk aims to double the capacity of his Memphis investors could end up benefiting as well thanks to Optimus.  
RICHARD BORD—WIREIMAGE/GETTY IMAGES

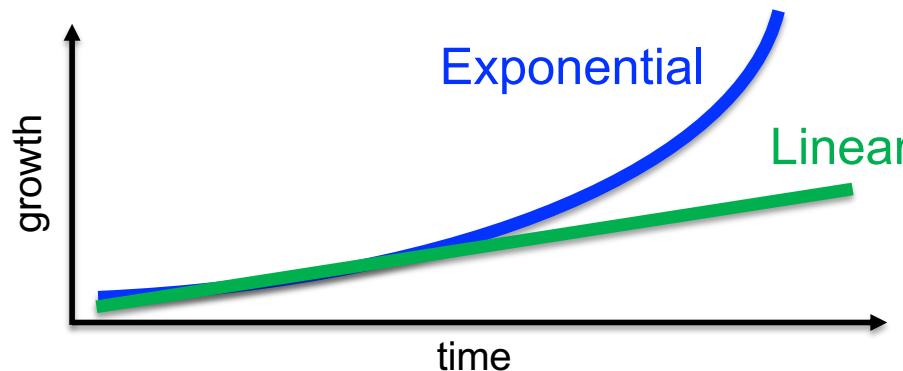
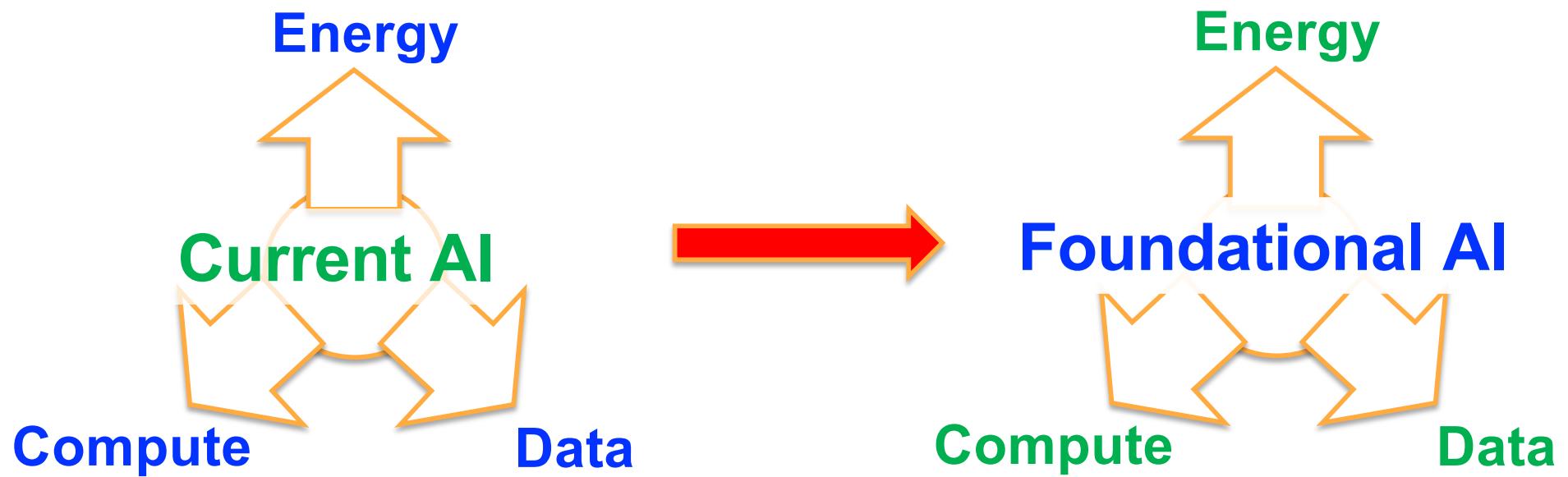


Contains 100,000 NVIDIA H100 processors, estimated cost \$U3-4 billion.

Will use 1 million gallons of water per day to cool servers, and consume 150 megawatts of power.



<https://www.reuters.com/markets/deals/constellation-inks-power-supply-deal-with-microsoft-2024-09-20/>



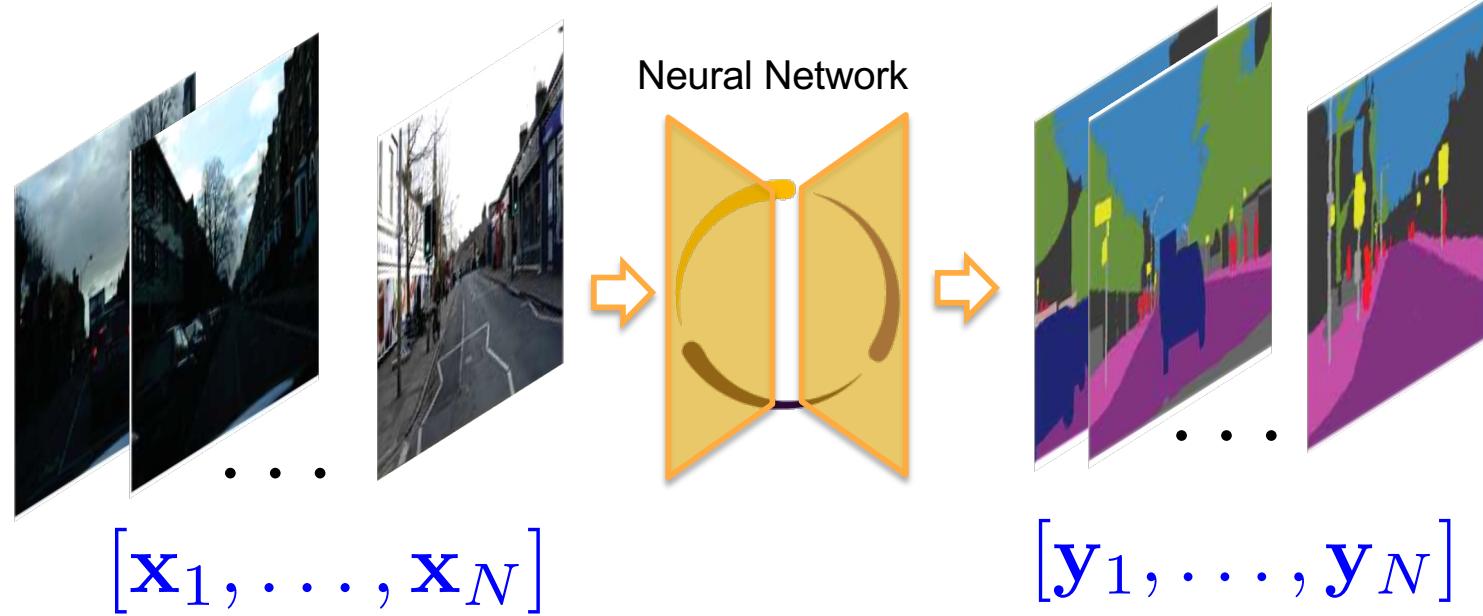
# What is Foundational AI?



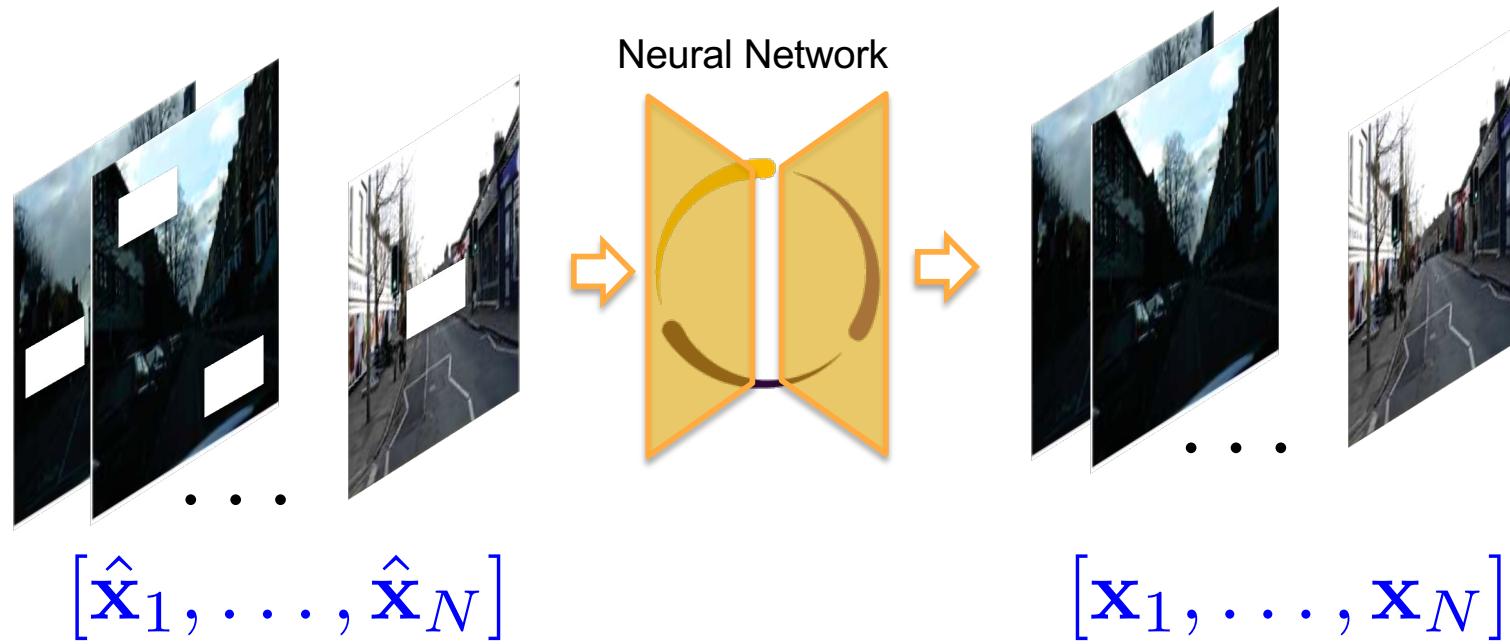


More DATA &  
More Compute?

**Where do I get the labels?**

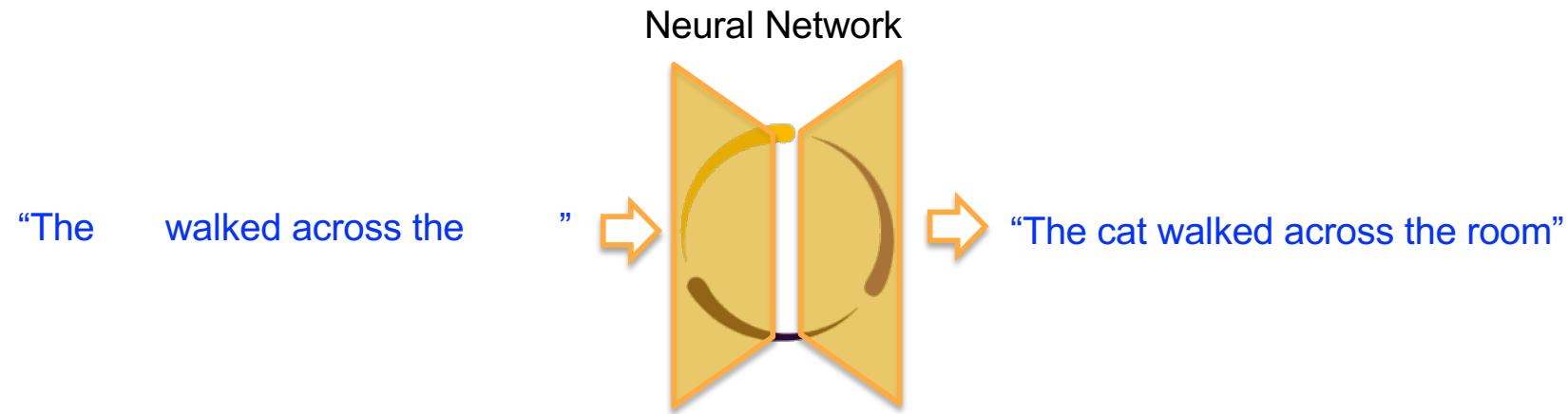


$$N \propto 10^6 - 10^9$$



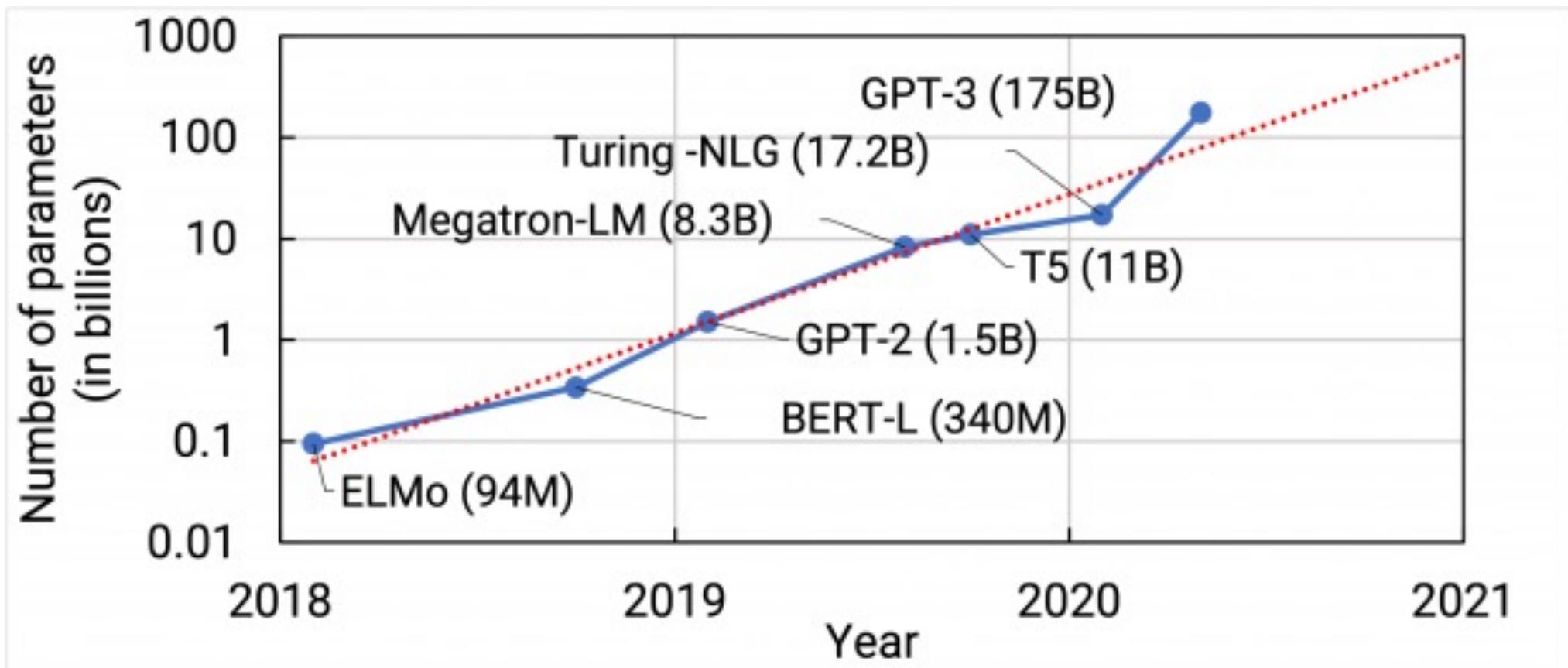
## Self Supervision

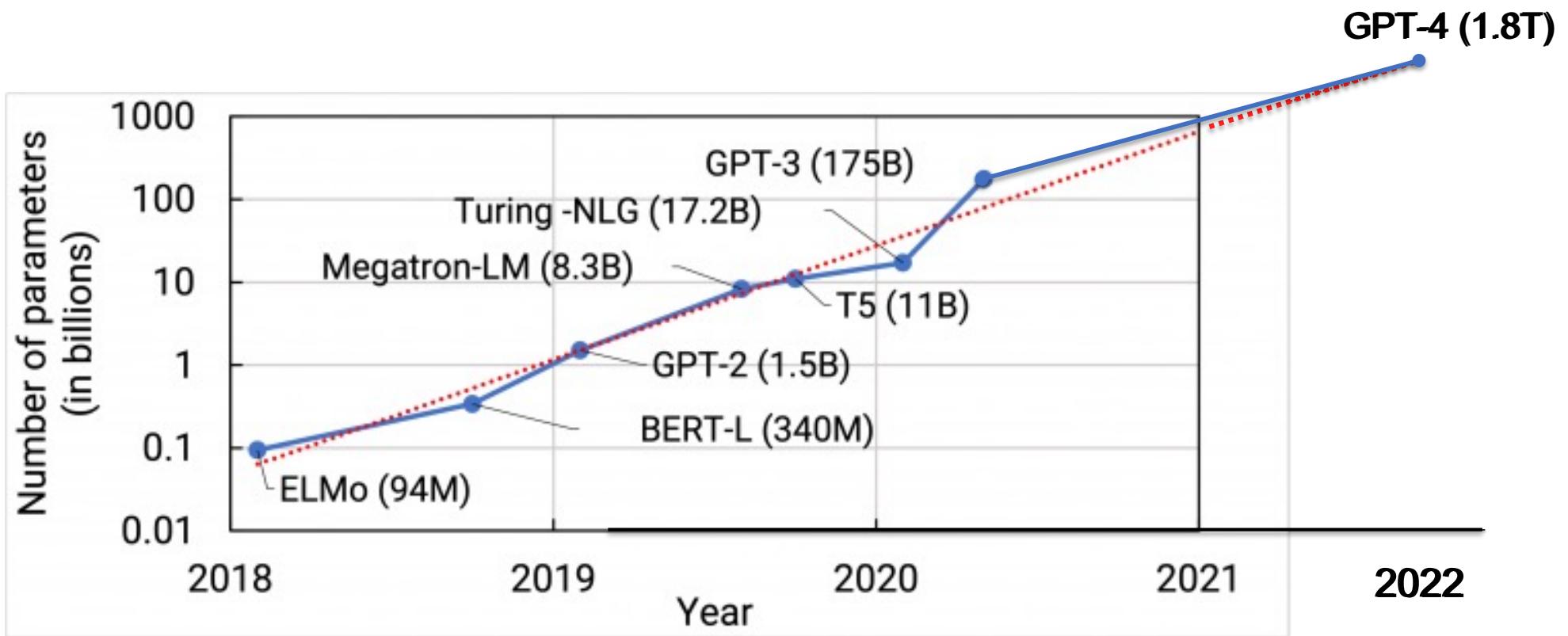
K. He, et al. "Masked Autoencoders Are Scalable Vision Learners." in CVPR 2022.



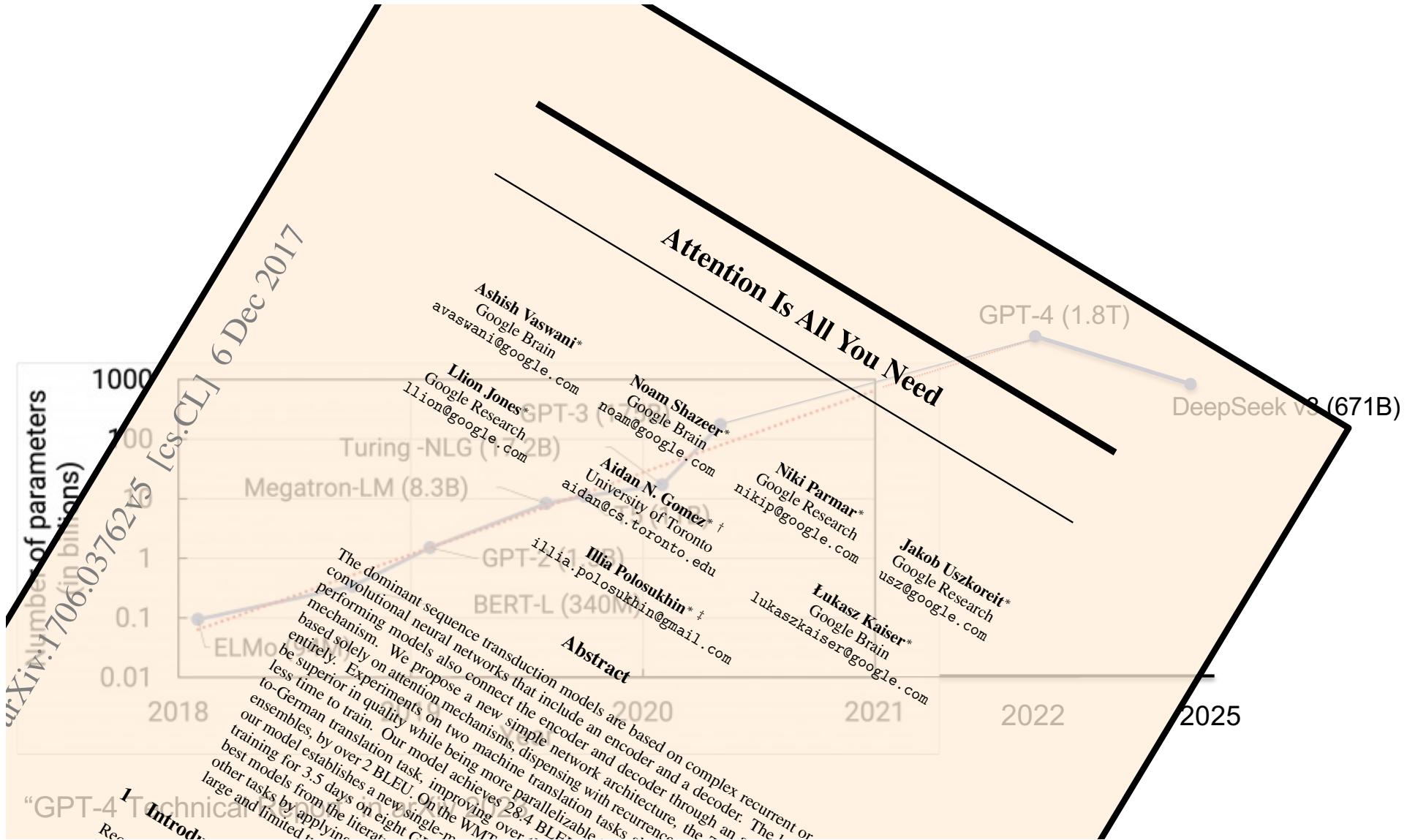
## Self Supervision

J. Devlin, et al. “BERT: Pre-training of deep bidirectional transformers for language understanding” in arXiv 2018.





"GPT-4 Technical Report" in arXiv 2023.



Is that IT?



T. Brookes, W. Peebles, et al. "Video generation models as world simulators." in *OpenAI Tech Report* 2024.



T. Brookes, W. Peebles, et al. "Video generation models as world simulators." in *OpenAI Tech Report 2024*.

# Multiplying with GPT-4

Faith and Fate:  
Limits of Transformers on Compositionality

Multiply 999 by 867

Nouha Dziri<sup>1,\*</sup>, Ximing Lu<sup>1,2\*</sup>, Melanie Sclar<sup>2\*</sup>, Xiang Lorraine Li<sup>1†</sup>, Liwei Jiang<sup>1,2†</sup>, Bill Yuchen Lin<sup>1</sup>, Peter West<sup>1,2</sup>, Chandra Bhagavatula<sup>1</sup>, Ronan Le Bras<sup>1</sup>, Jena D. Hwang<sup>1</sup>, Soumya Sanyal<sup>3</sup>, Sean Wellesck<sup>1,2</sup>, Xiang Ren<sup>1,3</sup>, Allyson Ettinger<sup>1,4</sup>, Zaid Harchaoui<sup>1,2</sup>, Yejin Choi<sup>1,2</sup>,  
<sup>1</sup>Allen Institute for Artificial Intelligence, <sup>2</sup>University of Chicago, <sup>3</sup>University of Southern California, <sup>4</sup>University of Washington  
nouhad@allenai.org, ximingu@allenai.org, msclar@cs.washington.edu

## Abstract



The product of multiplying 999 by 867 is 824,133

Transformer large language models (LLMs) have sparked admiration for their exceptional performance on tasks that demand intricate multi-step reasoning. Yet, these models simultaneously show failures on surprisingly trivial problems. This begs the question: Are these errors incidental, or do they signal more substantial limitations? In an attempt to demystify Transformers, we investigate the limits of these models across three representative compositional tasks—multi-digit multiplication, logic grid puzzles, and a classic dynamic programming problem. These tasks require breaking problems down into sub-steps and synthesizing them into a precise answer. We formulate the level of complexity, and break down these tasks systematically to quantify the level of complexity, and break down these tasks into intermediate sub-procedures. Our empirical findings suggest that solving skills. To round off our empirical study, we solve compositional tasks by reducing complexity, and break down these tasks into linearized subgraph matching, without necessarily decomposing them into abstract multi-step reasoning problems. Our empirical findings suggest that performance will rapidly decay with increasing complexity.

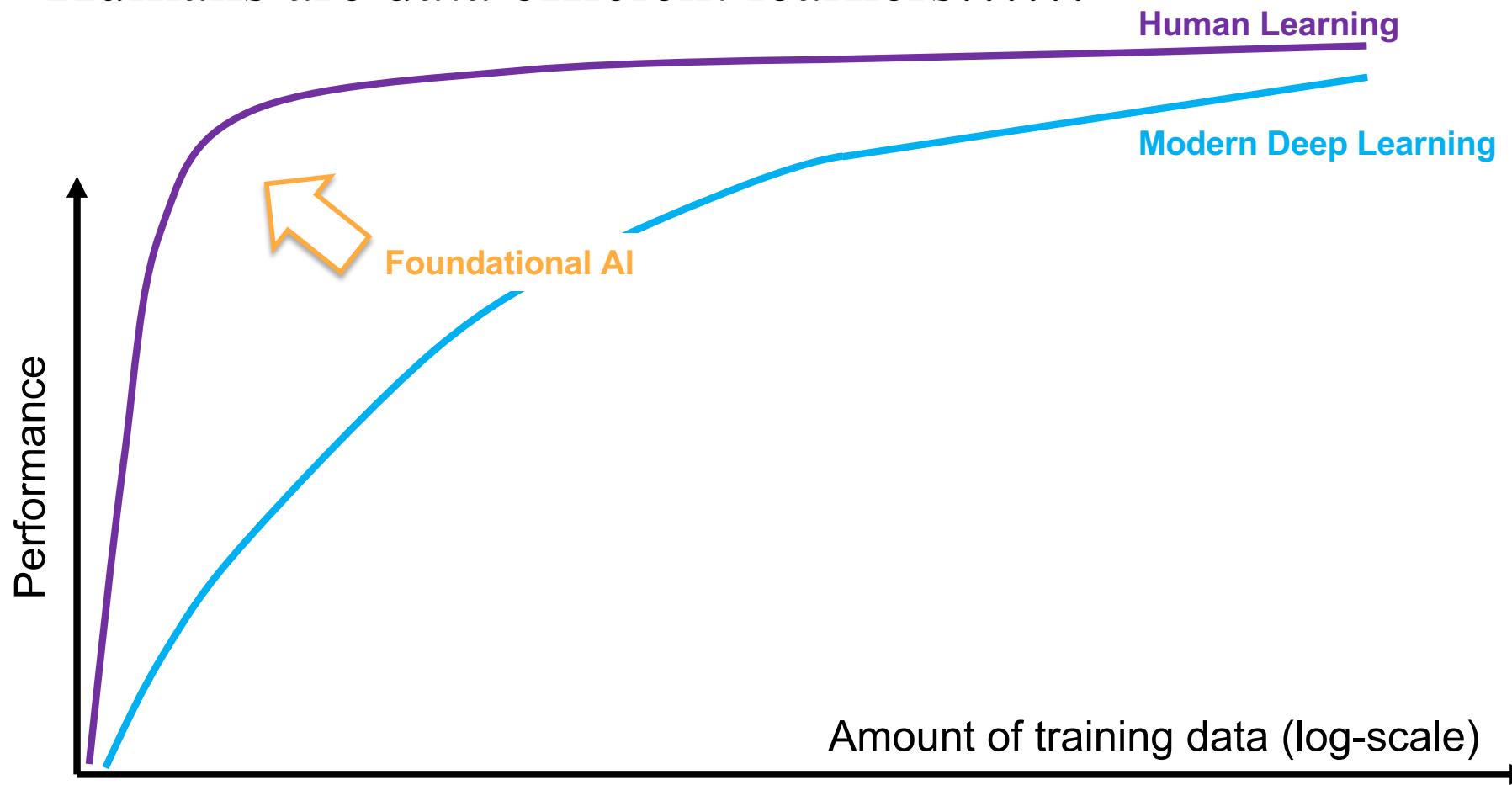
## 1 Introduction

"It was the epoch of *h*  
Large-scale Tr  
bilities /a  
some



arXiv:2305.18654v2 [cs.CL] 1 Jun 2023

Humans are data efficient learners.....



Data

# CommBank 'paid for itsel*f*' in three weeks



## Justin Hendry

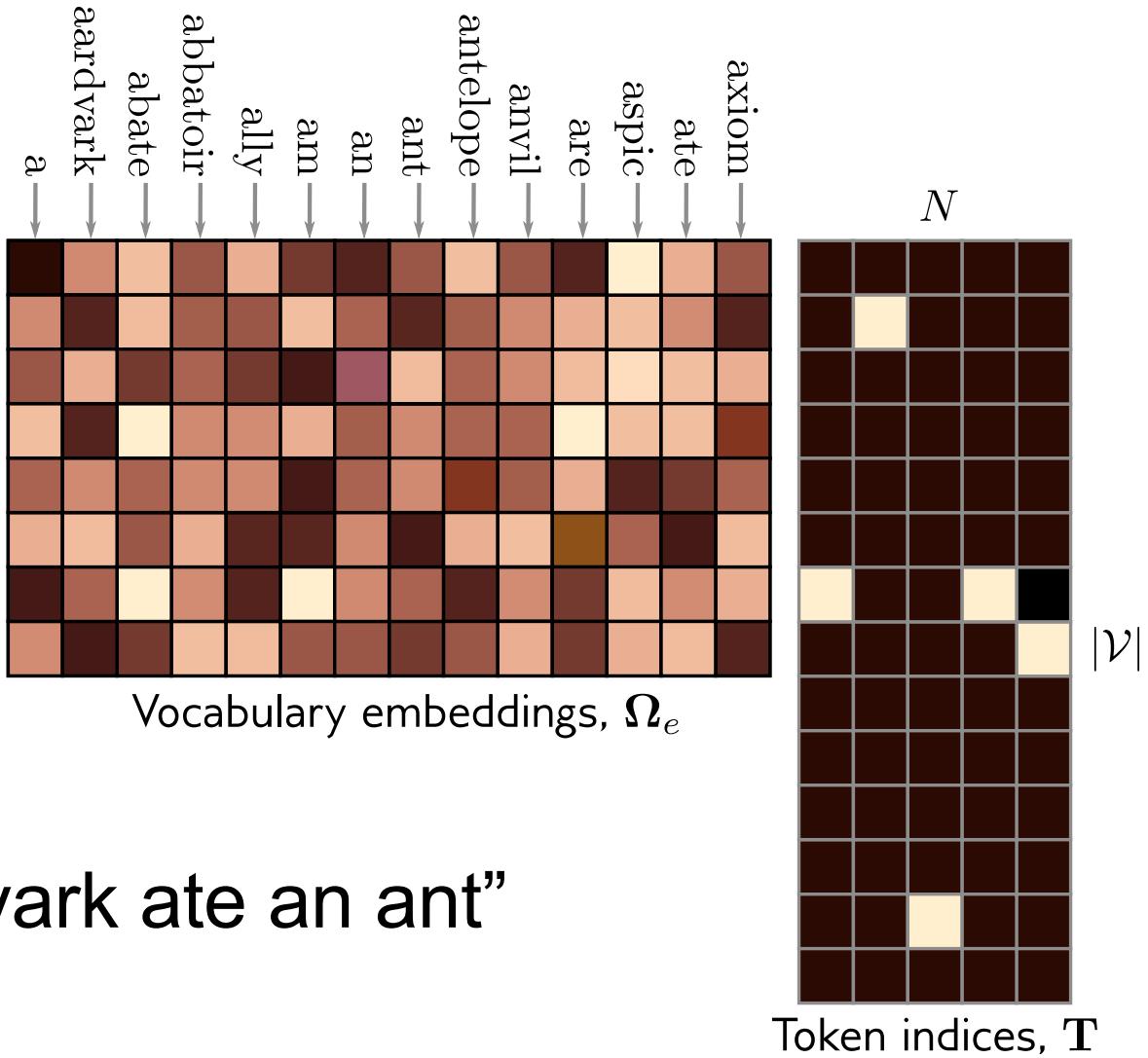
Editor

① 21 November 2024

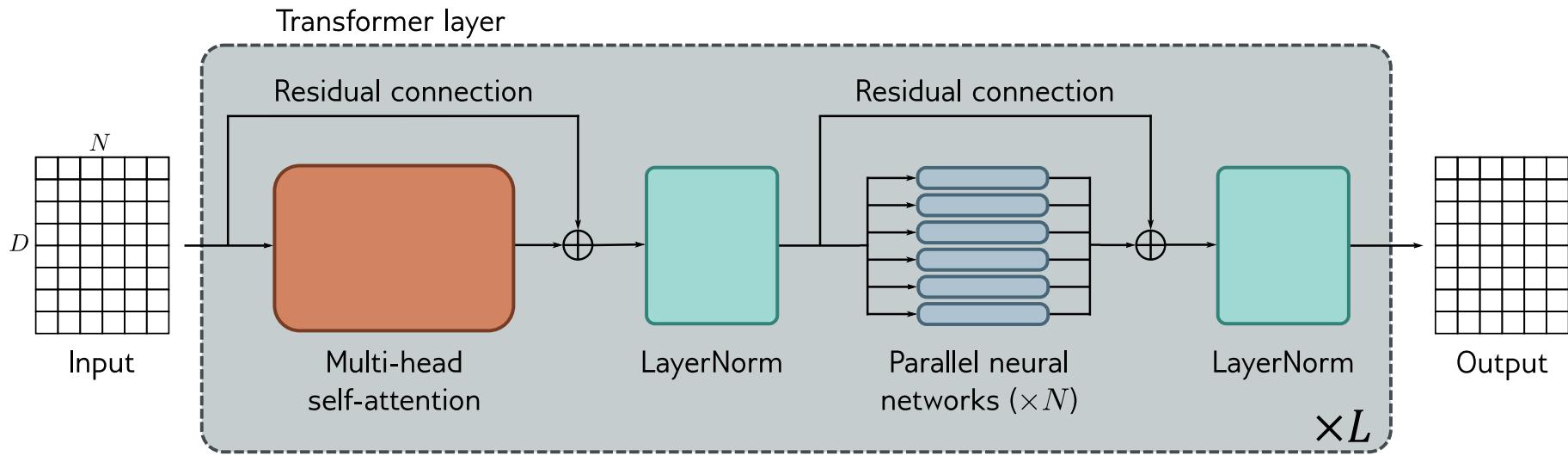
Share

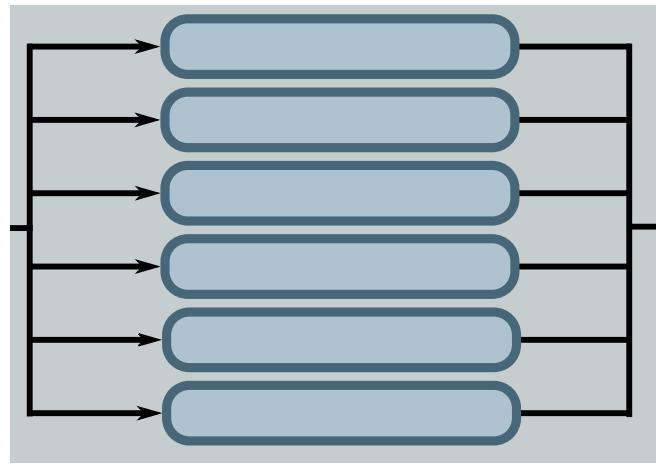
# Introduction

“an aardvark ate an ant”



“an aardvark ate an ant”





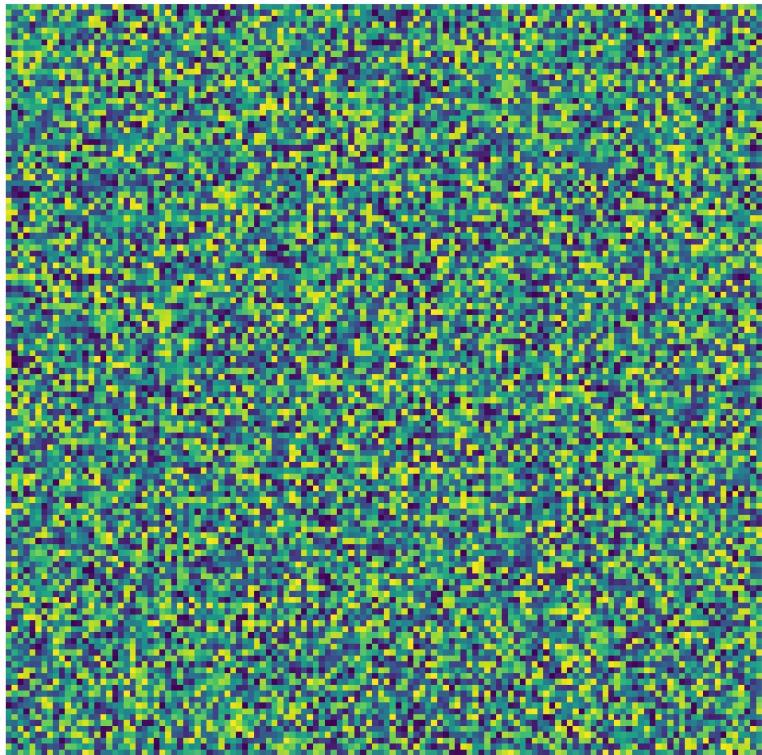
$$(\mathbf{W} \mathbf{x} + \Delta \mathbf{W}) \mathbf{x}$$

pre-trained  
weights

$$\Delta \mathbf{W}$$

fine-tuned  
weights

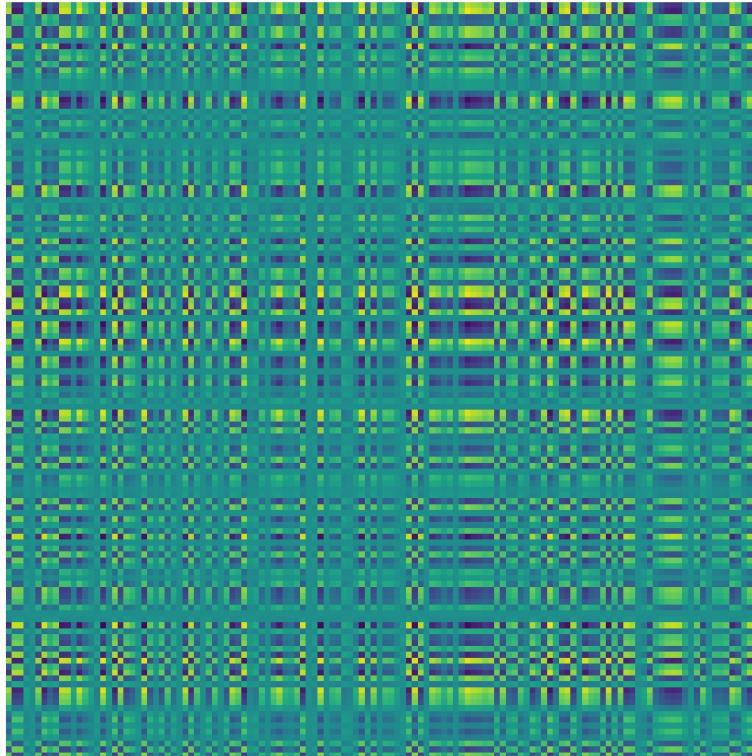
E. J. Hu, et al. "LoRA: Low-rank adaptation of large language models." In ICLR 2022.



- Ideally, full-rank.
- However, too many parameters.

$$\Delta \mathbf{W}$$

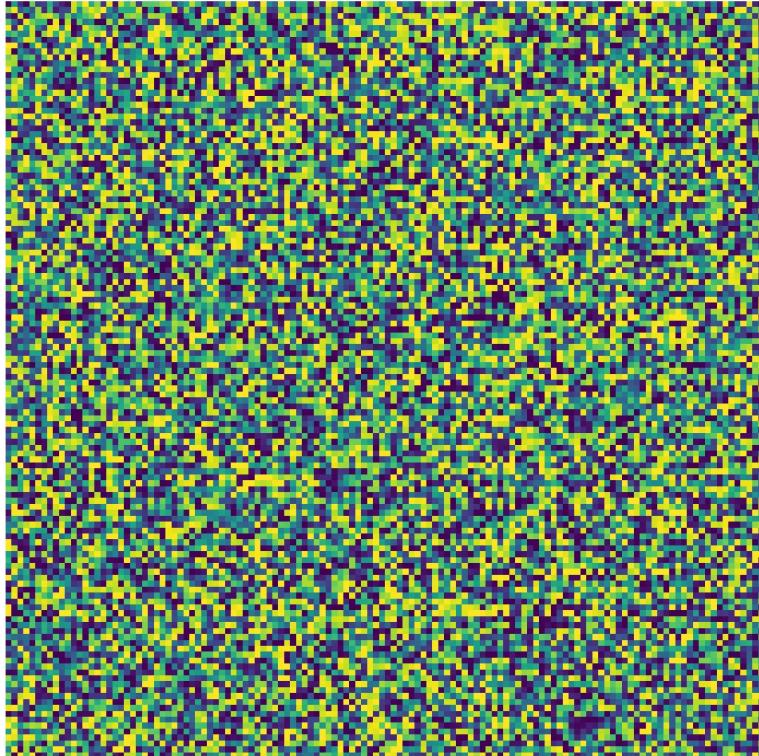
E. J. Hu, et al. “LoRA: Low-rank adaptation of large language models.” In ICLR 2022.



- Use low-rank approximation.
- However, lacks expressibility.

$$\mathbf{u} \mathbf{v}^T$$

E. J. Hu, et al. “LoRA: Low-rank adaptation of large language models.” In ICLR 2022.



- Foundational insight.
- Gives both parameter efficiency and expressability.

$$\sin(\omega \cdot \mathbf{u}\mathbf{v}^T)$$

Y. Ji, S. Lucey, et al. “Sine Activated Low-Rank Matrices for Parameter Efficient Learning.” In ICLR 2025.

# CBA's Adelaide Uni partnership spins out potentially game-changing AI capability

by Patrick Buncsi on 30 October 2024



Just two weeks after securing its partnership with Adelaide University's Australian Institute for Machine Learning (AIML), CommBank and the AIML built a potentially transformative "deep learning" capability expected to halve the speed for critical data processing functions within the bank. <https://fst.net.au/financial-services-news/cbas-adelaide-uni-partnership-spins-out-critical-ai-innovation/>

## Related News

 Delivering on a digital tomorrow – Sam Shaheen, Australia Post

 ASIC warns licensees governance must support AI adoption

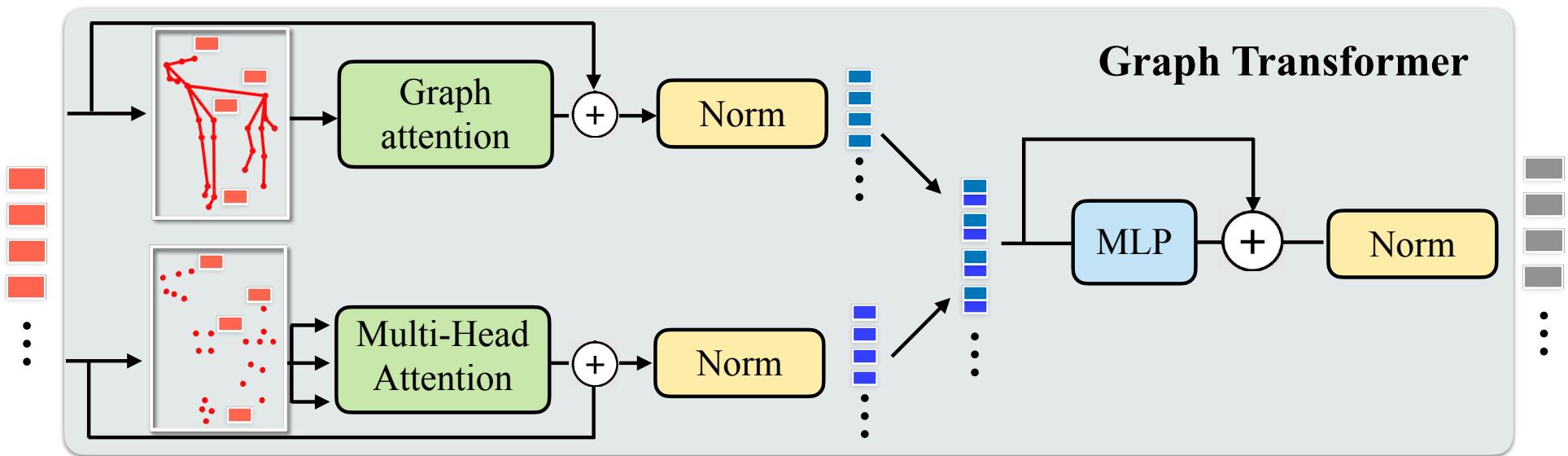
 AI Accelerator expands Westpac's AI to

 Unmasking the deepfake threat: Guidance for FSI execs, boards

 Beyond Bank taps Boomi for loan system revamp

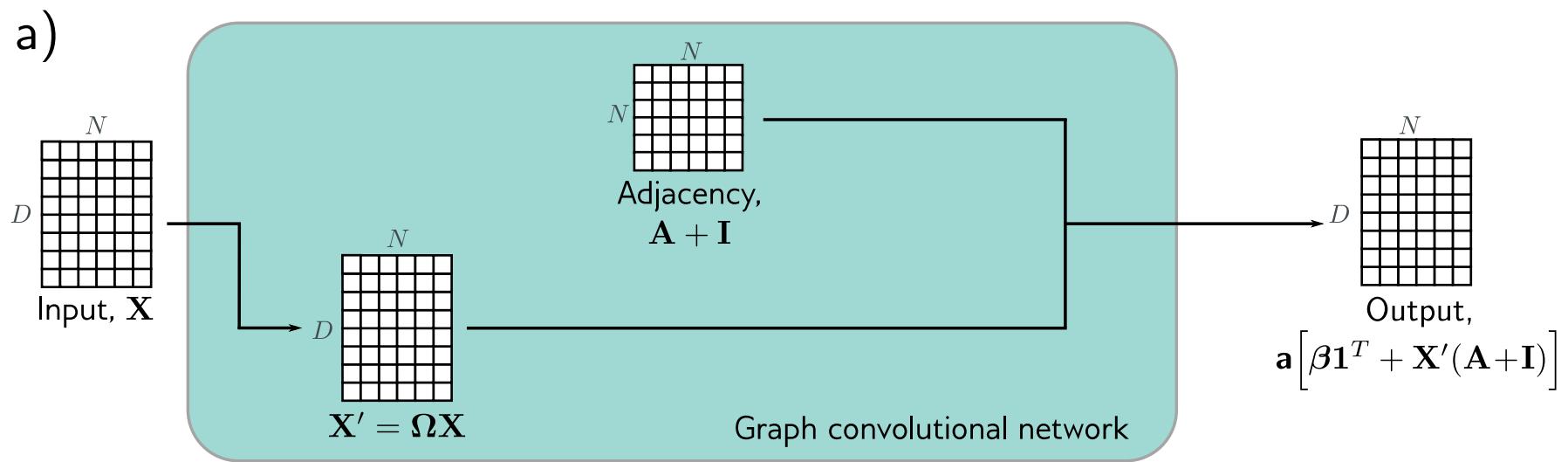
Beyond Vision and Language

# 3D Lifting Foundation Model (3D-LFM)



M. Dabhi, L. A. Jeni, and S. Lucey. "3D-LFM: Lifting Foundation Model." In CVPR 2024.

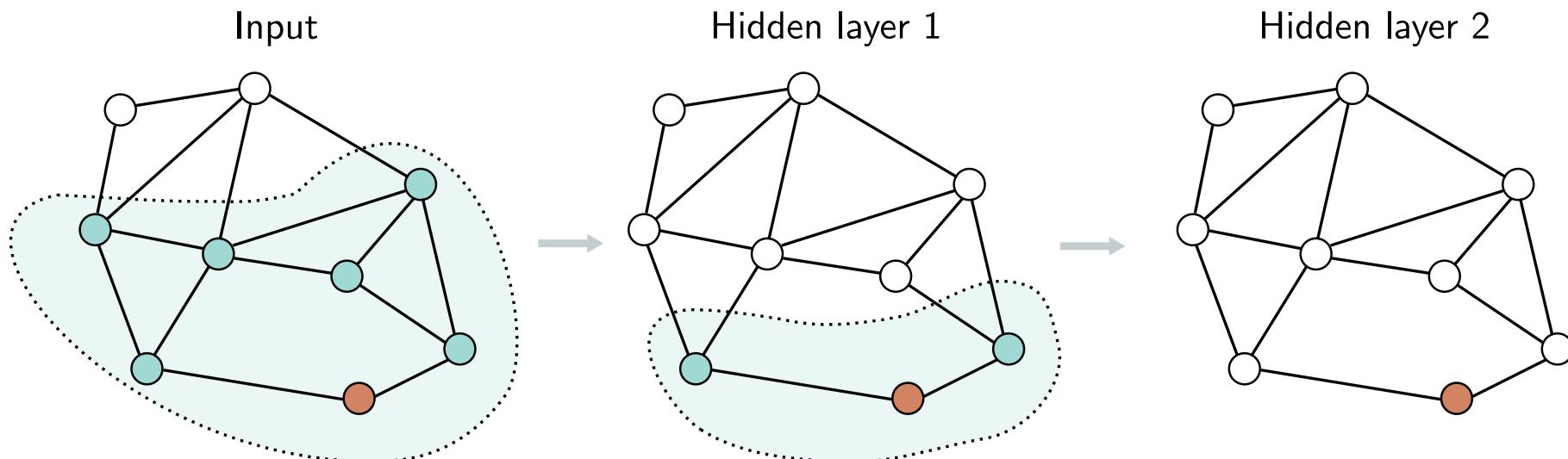
# Graph Convolution Network



Only  $\Omega$  and  $\beta$  are learned!!!

# Graph Convolution Network

---



Prince “Understanding Deep Learning” 2023.

# Graph Attention Networks

---

## GRAPH ATTENTION NETWORKS

**Petar Veličković\***

Department of Computer Science and Technology  
University of Cambridge  
[petar.velickovic@cst.cam.ac.uk](mailto:petar.velickovic@cst.cam.ac.uk)

**Guillem Cucurull\***

Centre de Visió per Computador, UAB  
[gucurull@gmail.com](mailto:gucurull@gmail.com)

**Arantxa Casanova\***

Centre de Visió per Computador, UAB  
[ar.casanova.8@gmail.com](mailto:ar.casanova.8@gmail.com)

**Adriana Romero**

Montréal Institute for Learning Algorithms  
[adriana.romero.soriano@umontreal.ca](mailto:adriana.romero.soriano@umontreal.ca)

**Pietro Liò**

Department of Computer Science and Technology  
University of Cambridge  
[pietro.li@cst.cam.ac.uk](mailto:pietro.li@cst.cam.ac.uk)

**Yoshua Bengio**

Montréal Institute for Learning Algorithms  
[yoshua.umontreal@gmail.com](mailto:yoshua.umontreal@gmail.com)

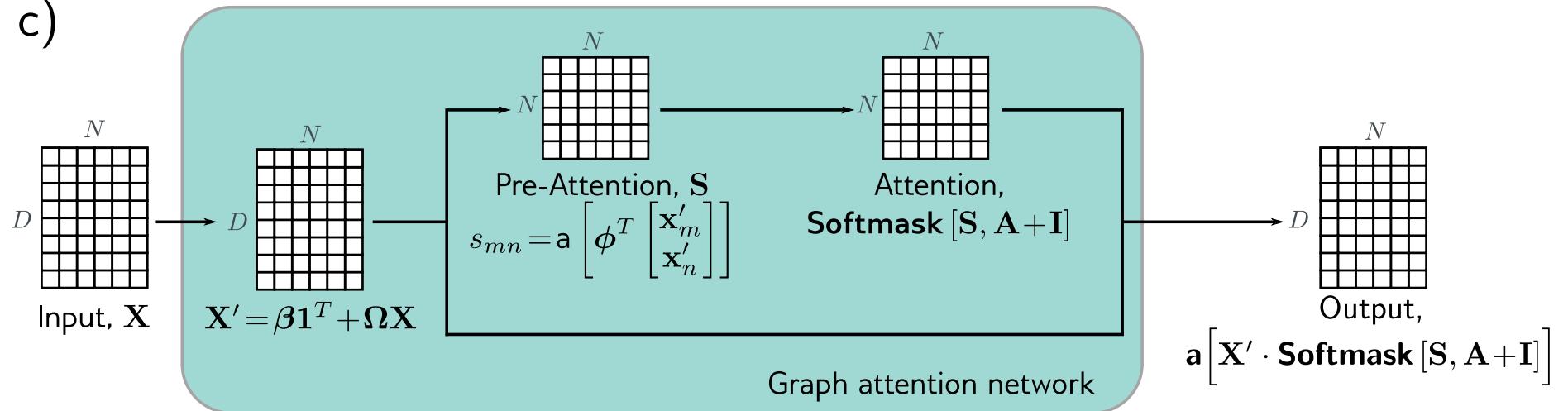
## ABSTRACT

We present graph attention networks (GATs), novel neural network architectures that operate on graph-structured data, leveraging masked self-attentional layers to address the shortcomings of prior methods based on graph convolutions or their approximations. By stacking layers in which nodes are able to attend over their neighborhoods' features, we enable (implicitly) specifying different weights to different nodes in a neighborhood, without requiring any kind of costly matrix operation (such as inversion) or depending on knowing the graph structure upfront. In this way, we address several key challenges of spectral-based graph neural networks simultaneously, and make our model readily applicable to inductive as well as transductive problems. Our GAT models have achieved or matched state-of-the-art results across four established transductive and inductive graph benchmarks: the *Cora*, *Citeseer* and *Pubmed* citation network datasets, as well as a *protein-protein interaction* dataset (wherein test graphs remain unseen during training).

A. Velickovic et al. “Graph Attention Networks”, ICLR 2018.

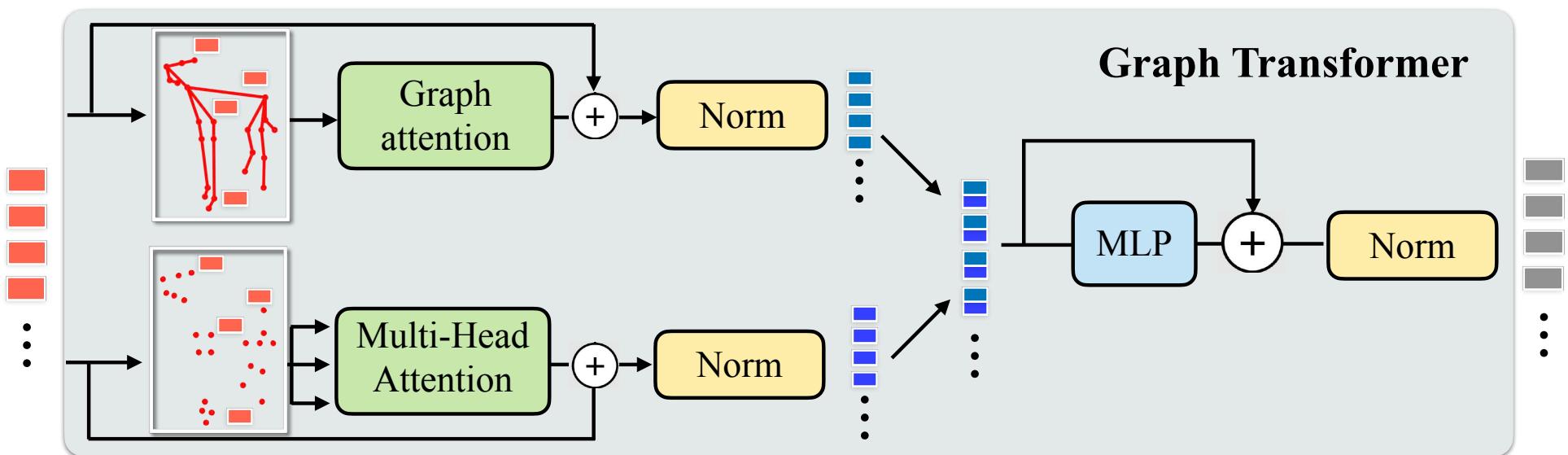
# Graph Attention Networks

c)



$\Omega$ ,  $\phi$  and  $\beta$  are learned!!!

# Reminder: 3D Lifting Foundation Model (3D-LFM)



M. Dabhi, L. A. Jeni, and S. Lucey. "3D-LFM: Lifting Foundation Model." In CVPR 2024.

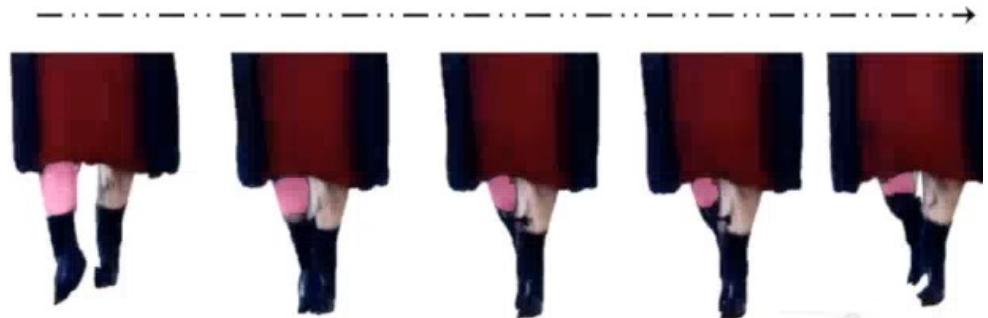
Flawed Original video



Corrected Rendered video



Error! Observe the right leg flip!



Our corrected re-render 🙏 to...



# \$20m AI research centre set up in South Australia

---



**Justin Hendry**

Editor

🕒 9 December 2024

Share ↴

A new AI research centre has been set up in South Australia, bringing together experts from the CSIRO's digital arm Data61 and the Australian Institute for Machine Learning (AIML) to work on adoption challenges.

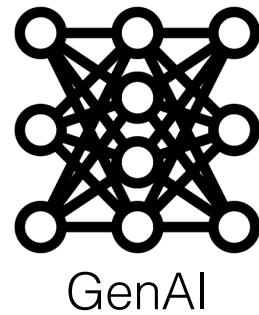
Jointly funded by the University of Adelaide, CSIRO and the South Australian government, the Responsible AI Research (RAIR) Centre aims to tackle issues of AI explainability, hallucinations and misinformation.

The \$20 million centre will focus on building foundational models and better understanding the inner workings of AI systems when it becomes fully operational at Adelaide's Lot Fourteen early next year

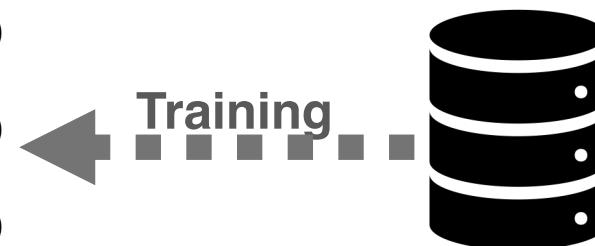
# Data Attribution



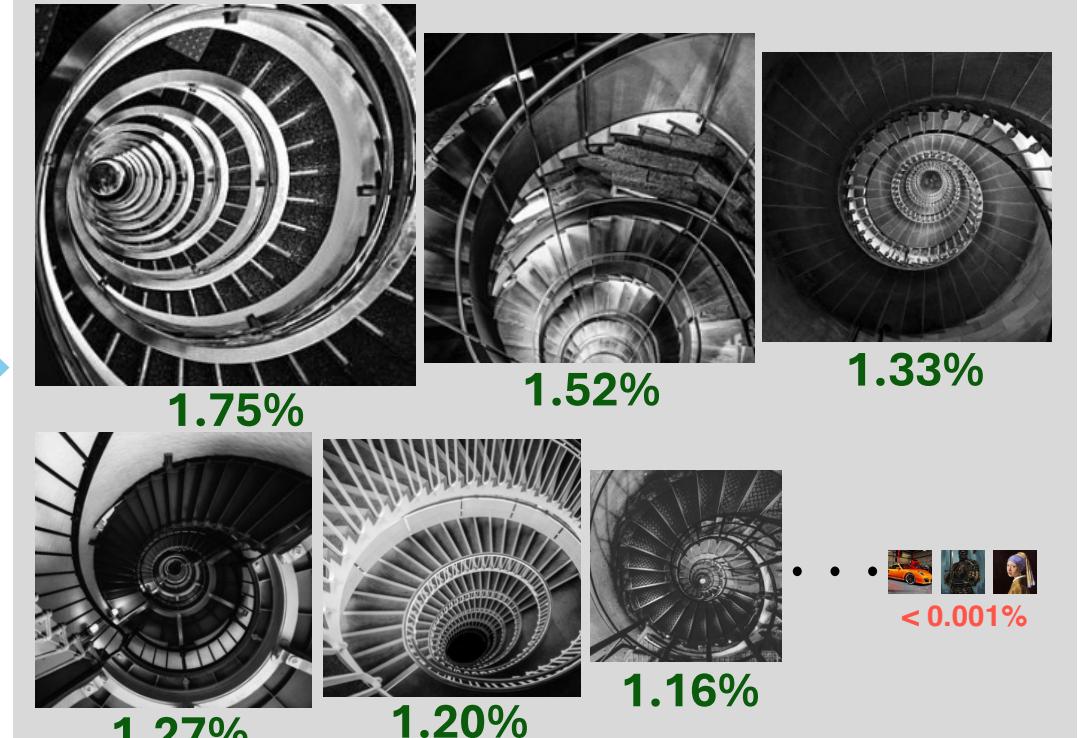
Synthesized Image



Generate



Training Dataset



Influence scores

## **Summary**

- AI growth is dependent upon exponential energy, data and compute.
- Understanding the basic mechanics of deep learning has the potential to unlock billions of dollars in efficiency and sovereign capability.
- Australia is already building foundational models from scratch!!