BIO-AI LAB| ARCTIC LLM WORKSHOP 2023
# Large Language Models

**Day 2 - Session 4**
## Evolution of Foundation LLM models

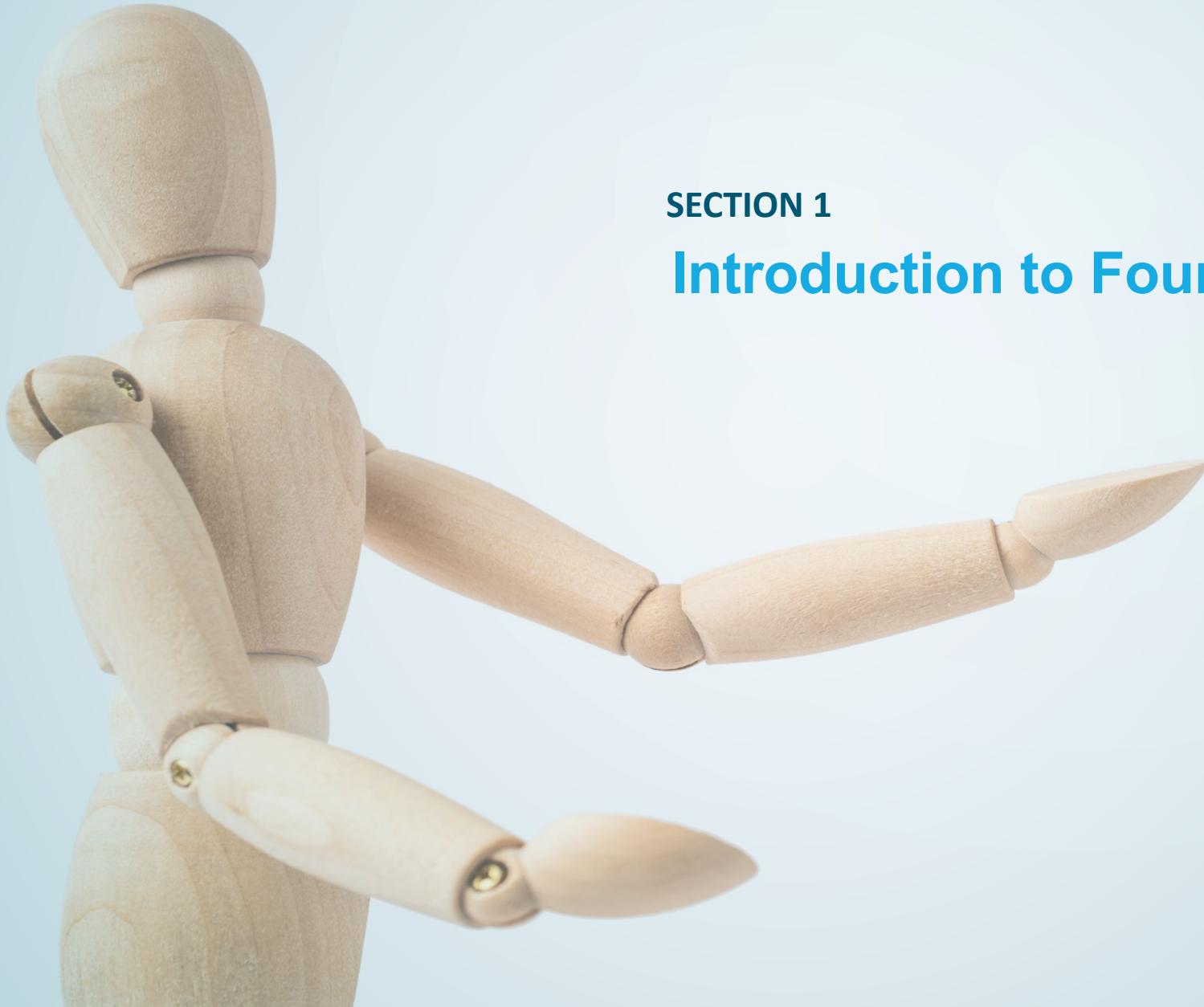**Abhinanda Ranjit Punnakkal**

Abhinanda.r.punnakkal@uit.no

27. Oct 2023

# Outline

SECTION 1

# Introduction to Foundational Models

# What is a Foundation Model?

*Examples:*
- *Ask open-ended questions.*
- *Paraphrase*
- *Generate code, web code*
- *Generate cool images*
- *Generating videos*
- *Robot manipulation*

Dear Richard Dawkins,

Can you tell me in short what was the central idea in your book "The Selfish Gene"?

Warmly,
Paras

----

Dear Paras,

This is quite hard to answer briefly. But the central idea is that living organisms are a kind of machine, a 'vehicle' for replicating molecules - the genes. The organism is like a truck whose job it is to transport the genes around in order to survive, and then reproduce them in new trucks (that is, in offspring). This is done by the genes making the truck (organism) behave in such a way as to protect the genes, to allow them to survive and reproduce. The genes 'want' to live forever if they can. This is their single aim in life, and they are utterly ruthless in achieving it. But the point is that they themselves are not alive. The truck (the organism) is alive, but the 'driver' (the genes) is not. You see what I mean? I hope that is clear!

Best wishes, Richard.

# What is a Foundation Model?

*Examples:*
- *Ask open-ended questions.*
- *Paraphrase*
- *Generate code, web code*
- *Generate cool images*
- *Generating videos*
- *Robot manipulation*

**AI Translations (AI Output in Blue)**

**Legal clause:** Upon liquidation of the Company, the Series A Shareholders will receive in preference to all other shareholders an amount in respect of each Series A Share equal to one times the Original Issue Price (the "Liquidation Preference"), plus all accrued but unpaid dividends. To the extent that the Company has assets remaining after the distribution of that amount, the Series A Shareholders will participate with the holders of Ordinary Shares pro rata to the number of shares held on an as converted basis.

**Plain English:** If the startup is wound up, the Series A investors will be paid back at least what they invested and they will also share any leftover assets with ordinary shareholders.

**Legal Clause:** Sale of all or substantially all of the assets of the Company or a sale of shares involving a change in control (each, a "Corporate Transaction") will be treated in the same way as a liquidation and the proceeds of sale will be distributed as set out in paragraph 3. If the holders of Series A Shares have received any Special Dividend it shall be set off against their Liquidation Preference.

**Plain English:** If the company is sold, or a new owner takes control, the proceeds of the sale will be distributed as in the liquidation clause above. Any special dividend paid will be treated as an initial payment towards the Series A investors.
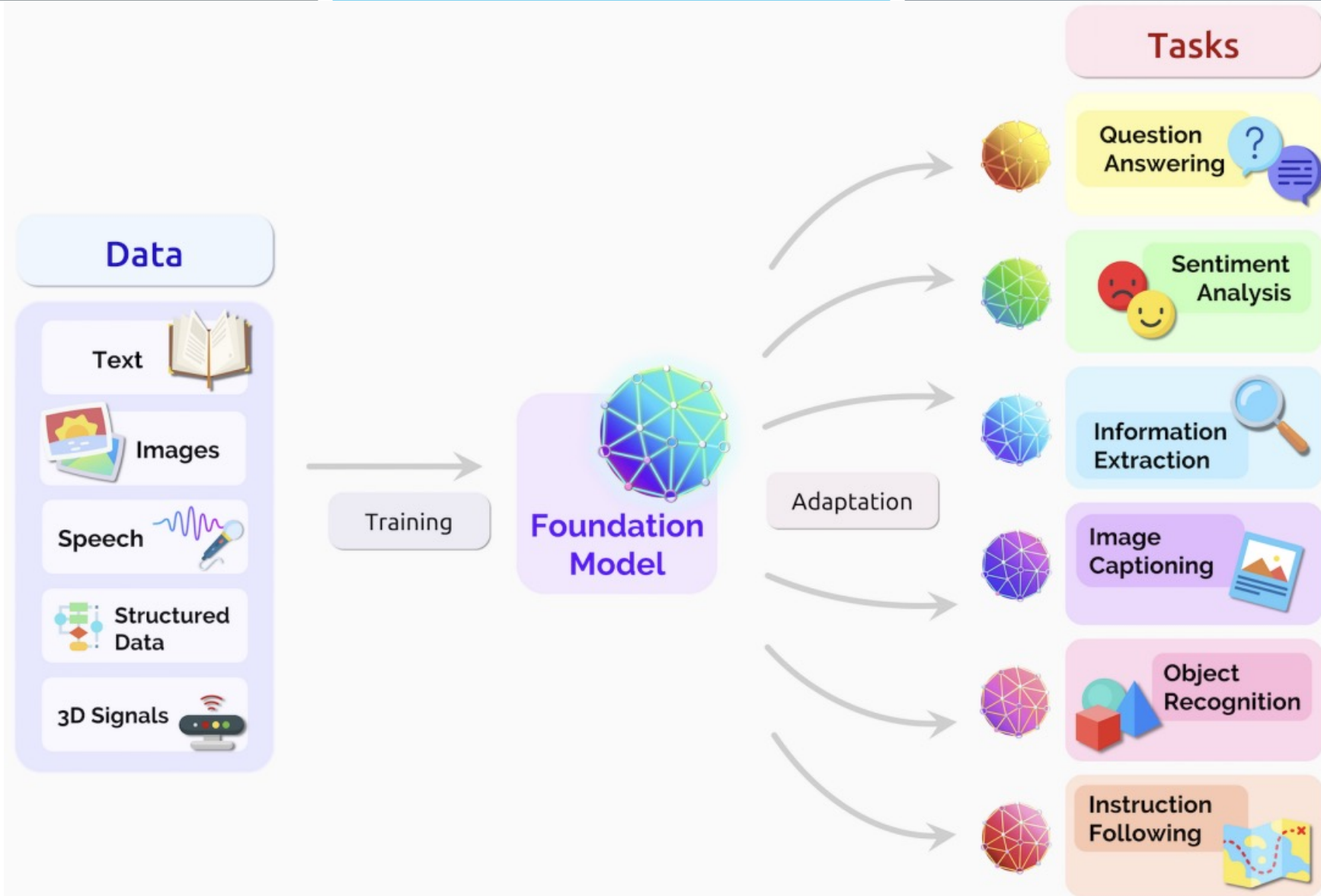
# What is a Foundation Model?

*Examples:*
- *Ask open-ended questions.*
- *Paraphrase*
- *Generate code, web code*
- *Generate cool images*
- *Generating videos*
- *Robot manipulation*

## Describe a layout.

Just describe any layout you want, and it'll try to render below!

a button that looks like a watermelon

Generate

```
<button style={{backgroundColor: 'pink', border: '2px solid green', borderRadius:
'50%', padding: 20, width: 100, height: 100}}>Watermelon</button>
```

Watermelon

# What is a Foundation Model?

*Examples:*
- *Ask open-ended questions.*
- *Paraphrase*
- *Generate code, web code*
- *Generate cool images*
- *Generating videos*
- *Robot manipulation*

An astronaut riding a horse on mars

**Generate image**

# What is a Foundation Model



Source: Bommasani, Rishi, et al. "On the opportunities and risks of foundation models." *arXiv preprint arXiv:2108.07258* (2021).

# Success of a Foundation Model

1. *Transformer architecture*
2. *Scale*
3. *In context learning/Self-supervised training*

Source: Zhou, Ce, et al. "A comprehensive survey on pretrained foundation models: A history from bert to chatgpt." *arXiv preprint arXiv:2302.09419* (2023).

*Efficient architecture*



**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

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com
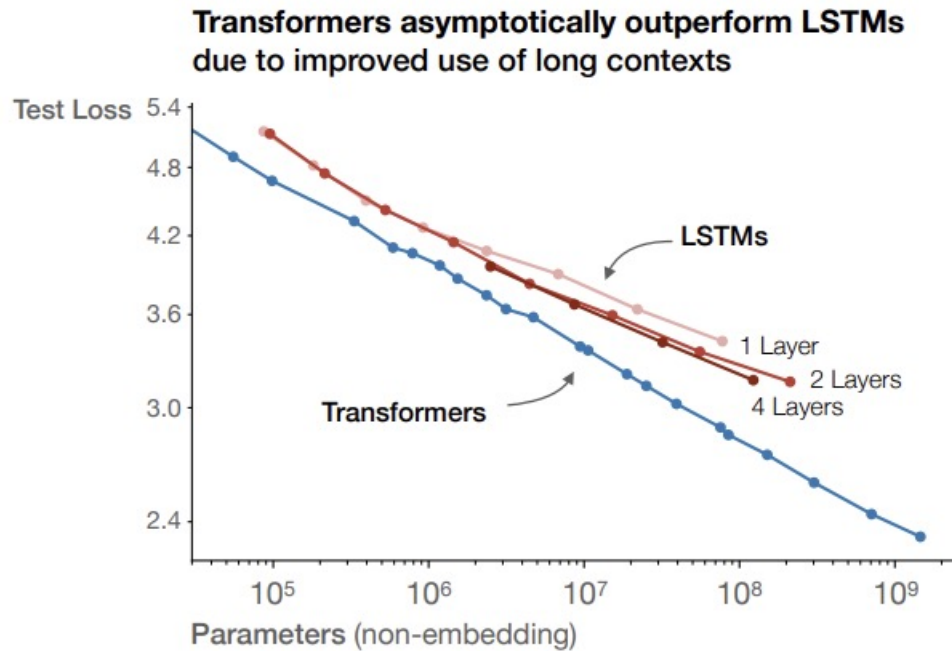
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Figure 1: The Transformer - model architecture.

# LSTMs vs Transformers



Source: Kaplan, Jared, et al. "Scaling laws for neural language models." *arXiv preprint arXiv:2001.08361* (2020).

# Before Large Models *2012- 2020*

*CNNs ruled!*



ImageNet Classification with Deep Convolutional Neural Networks
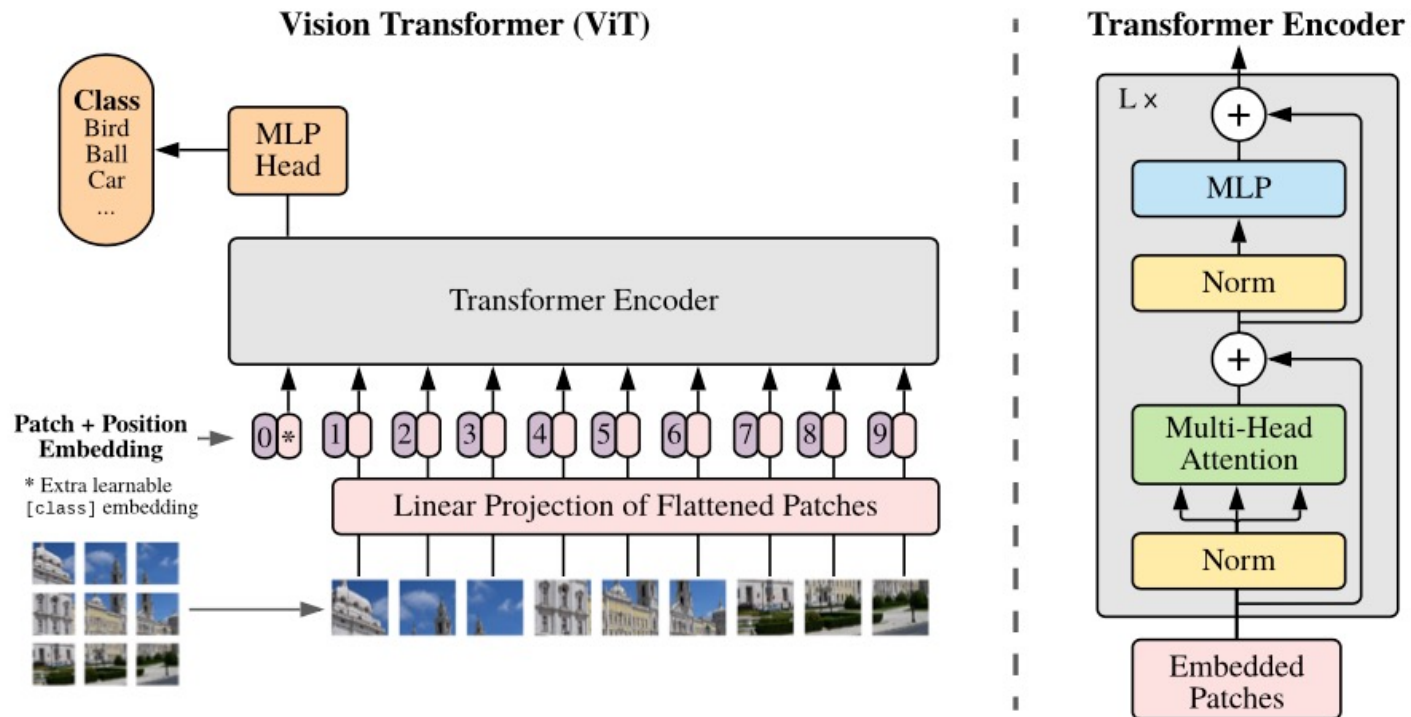
Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

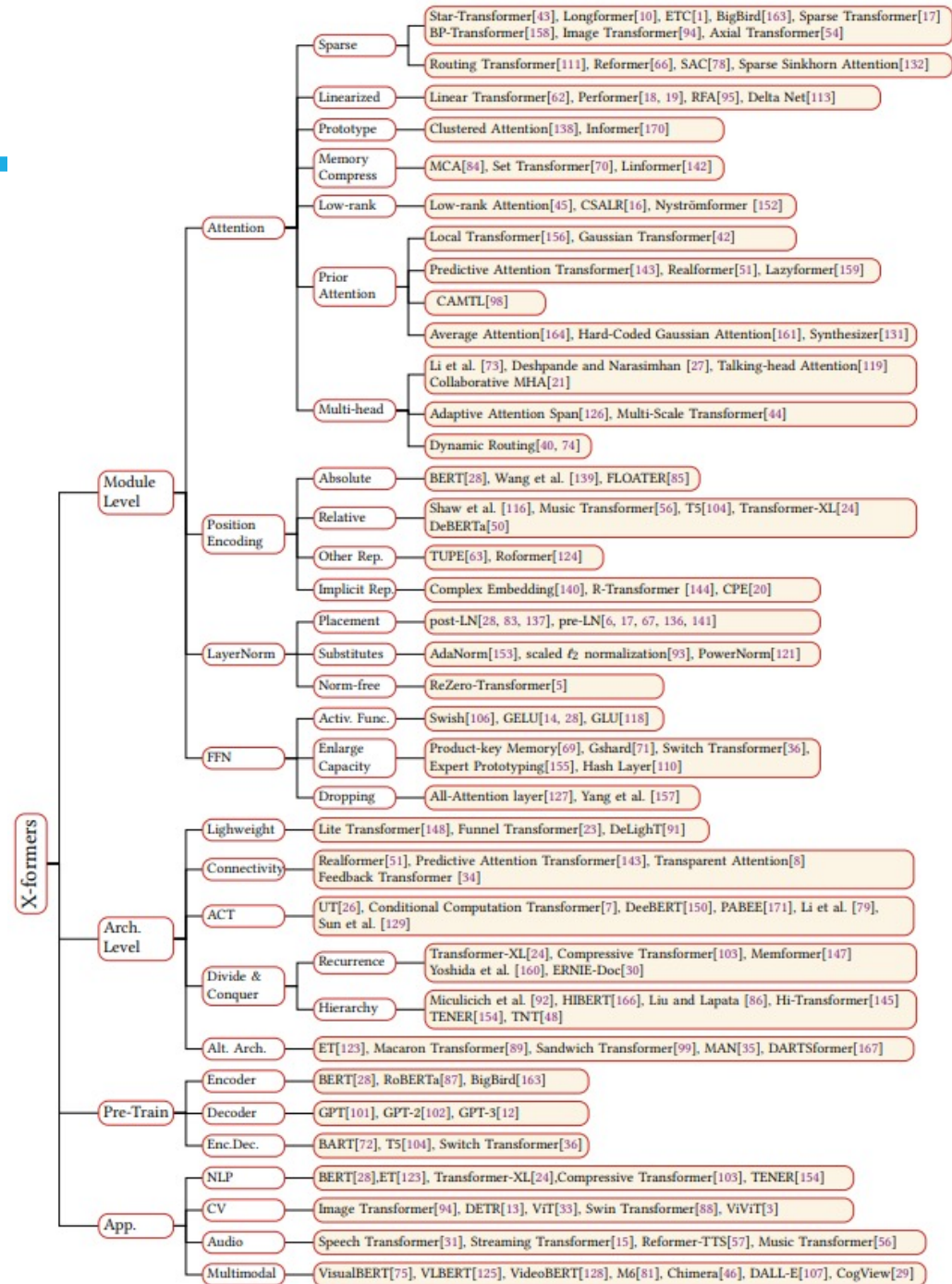Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

# ViT

*No more Convolutions for Images!*



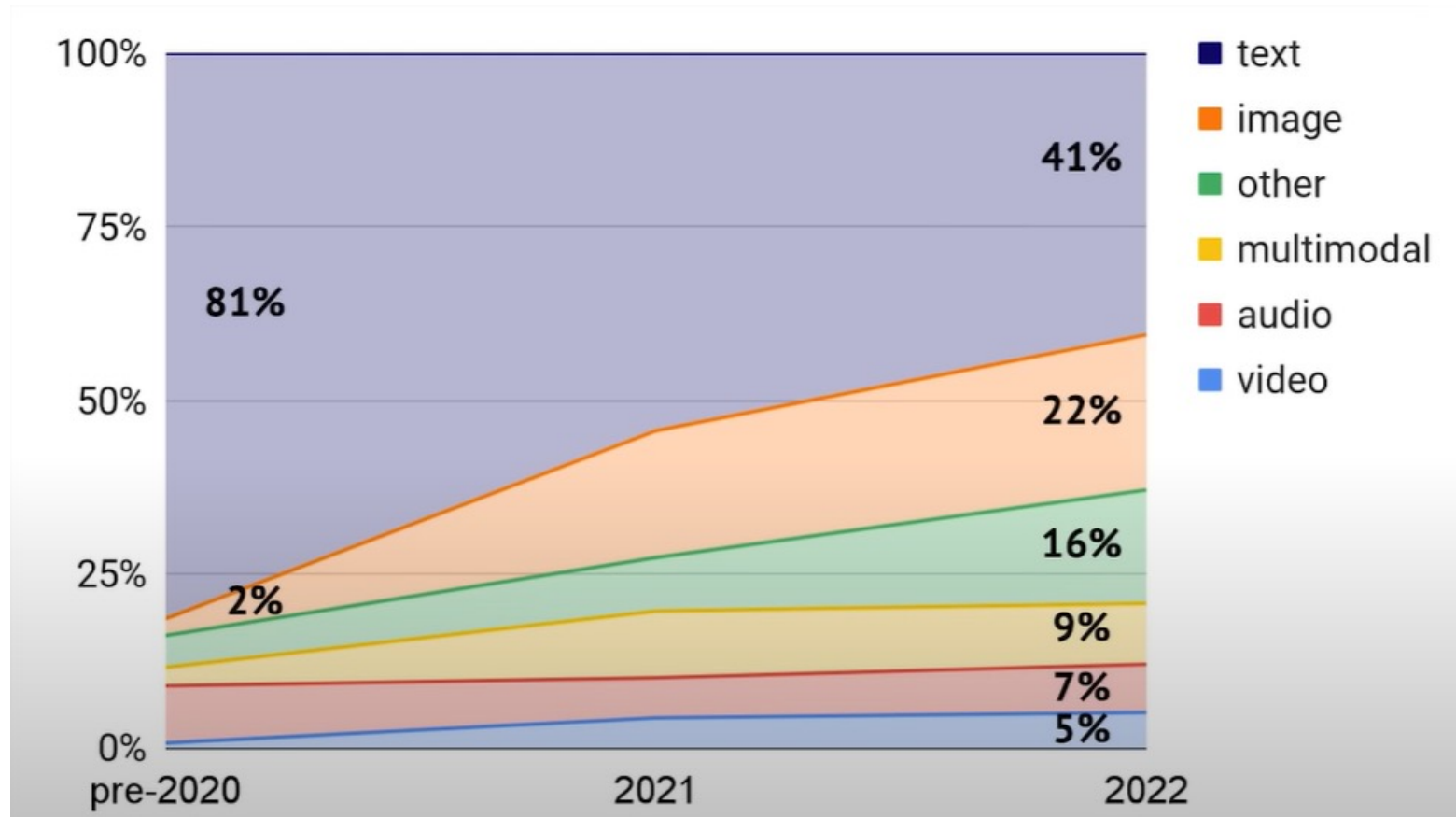Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." *arXiv preprint arXiv:2010.11929* (2020).

# Variations of Transformer

1. *Model efficiency*
2. *Model generalization*
3. *Model adaptation*

Source: Lin, Tianyang, et al. "A survey of transformers." *AI Open* (2022).

# Growth of Transformer related papers

*Started out with NLP*

# 2. Scale

- *model with number of parameters*
- *data*

## Scaling Laws for Neural Language Models

**Jared Kaplan** *
Johns Hopkins University, OpenAI
jaredk@jhu.edu

**Sam McCandlish***
OpenAI
sam@openai.com

**Tom Henighan**
OpenAI
henighan@openai.com

**Tom B. Brown**
OpenAI
tom@openai.com

**Benjamin Chess**
OpenAI
bchess@openai.com

**Rewon Child**
OpenAI
rewon@openai.com

**Scott Gray**
OpenAI
scott@openai.com

**Alec Radford**
OpenAI
alec@openai.com

**Jeffrey Wu**
OpenAI
jeffwu@openai.com

**Dario Amodei**
OpenAI
damodei@openai.com

- non-embedding parameter count N,
- dataset size D
- optimized training computation Cmin

# 2. Scale

- *model with number of parameters*
- *data*

**DeepMind**

## Training Compute-Optimal Large Language Models

Jordan Hoffmann*, Sebastian Borgeaud*, Arthur Mensch*, Elena Buchatskaya, Trevor Cai, Eliza Rutherford,
Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland,
Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan,
Erich Elsen, Jack W. Rae, Oriol Vinyals and Laurent Sifre*
*Equal contributions

Chinchilla                    >        Gopher
(with more training tokens)            (with a larger model size)

# 3. Self Supervision

*Auto-regressive pretraining*
*LLMs are trained to predict the next of a large corpus*



A language model is a probability distribution over a sequence of words

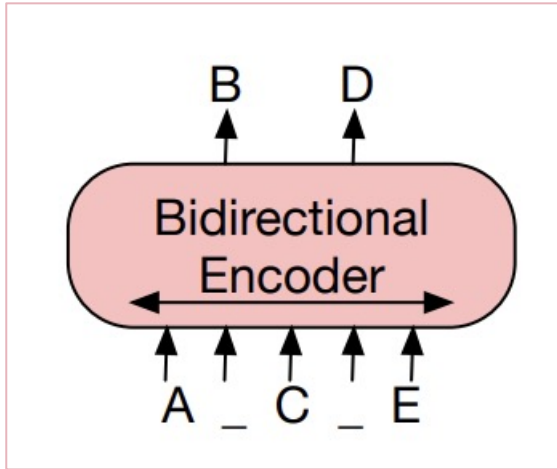Context                                                    Next word

SECTION 2

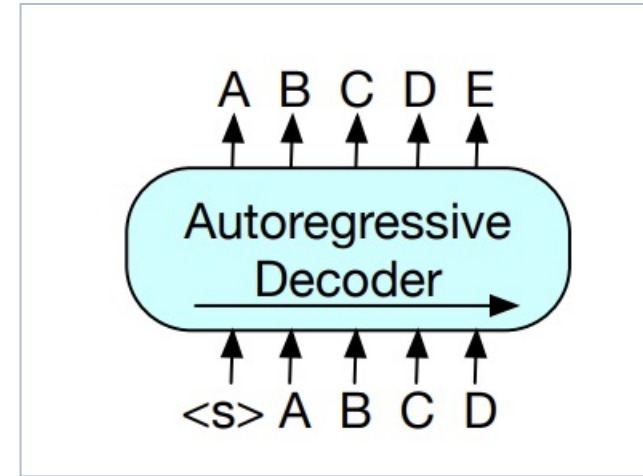# Evolution of Foundational LLM  Models

Evolutionary Tree

Yang, Jingfeng, et al. "Harnessing the power of llms in practice: A survey on chatgpt and beyond." *arXiv preprint arXiv:2304.13712* (2023).
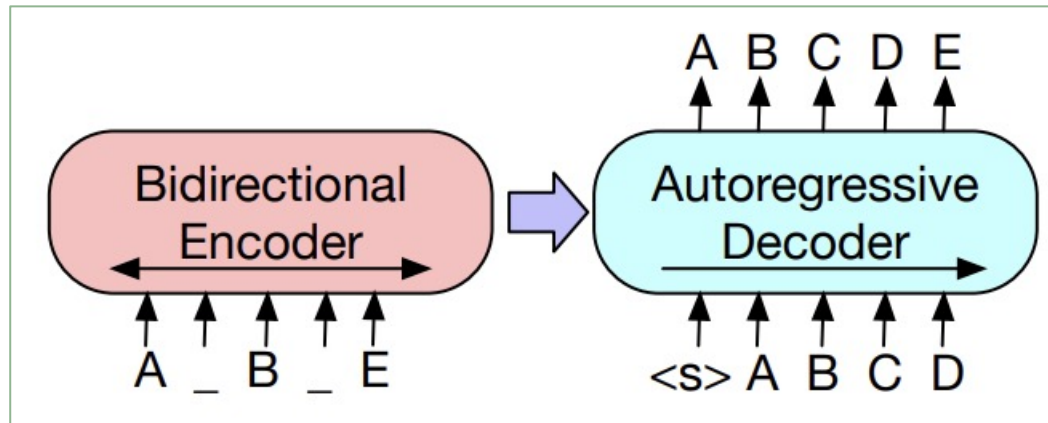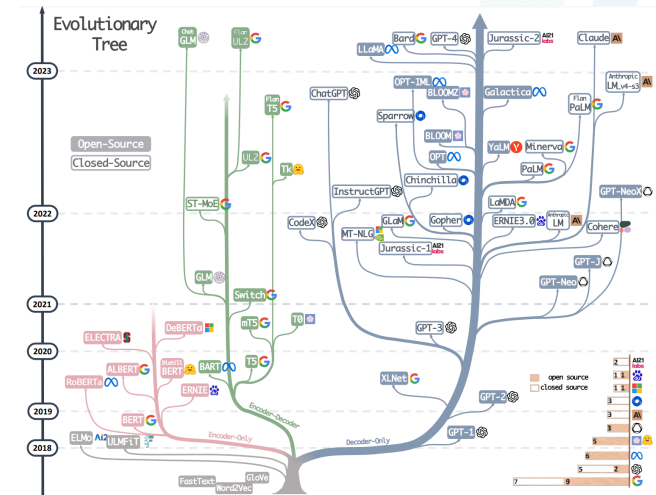
# Categories of LLMs



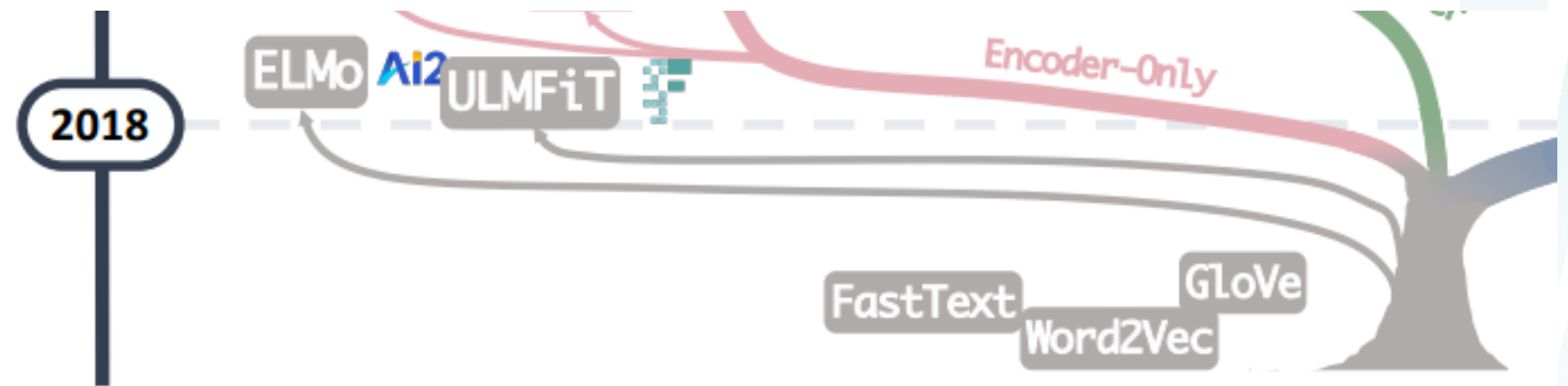Encoder only (BERT)



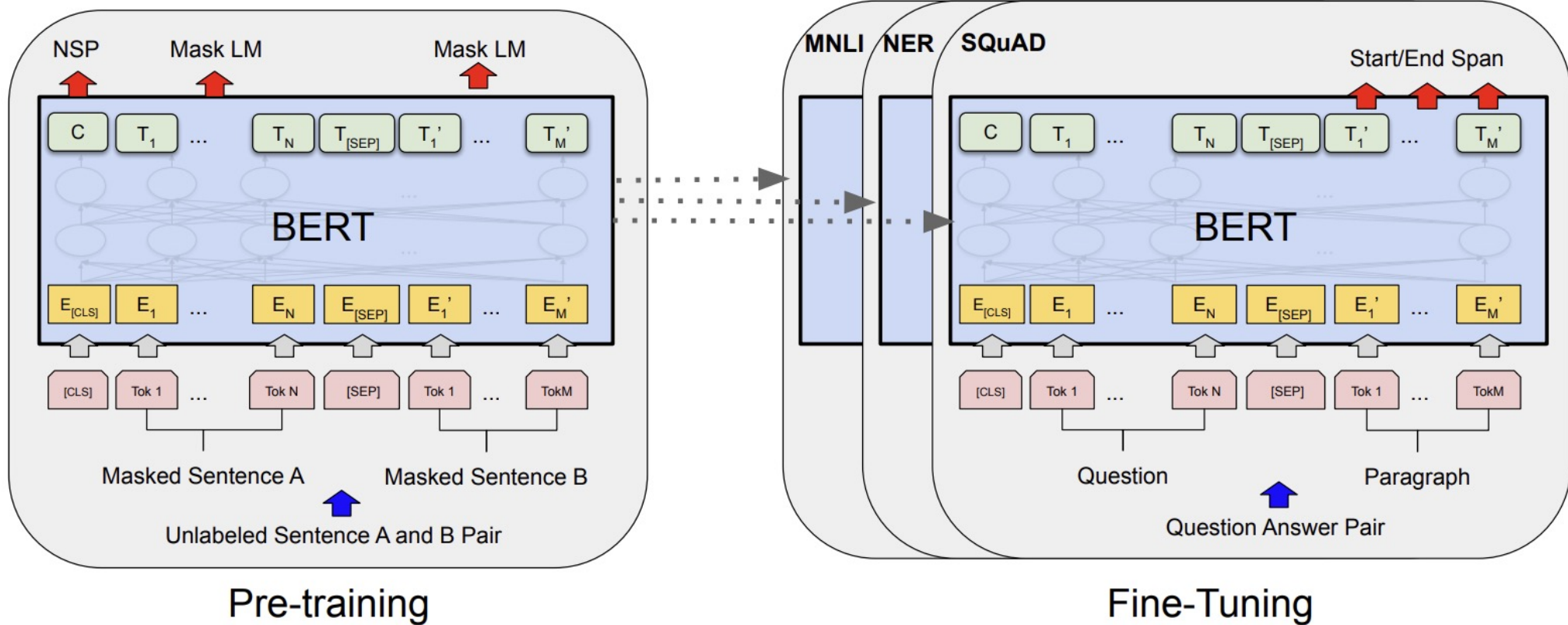Decoder only (GPT)



Encoder-decoder (BART)

# Before LLMs

- *One-hot encoding words*
  - *Did not scale*
- *Static embedding*
  - *Limited to the bag of words*

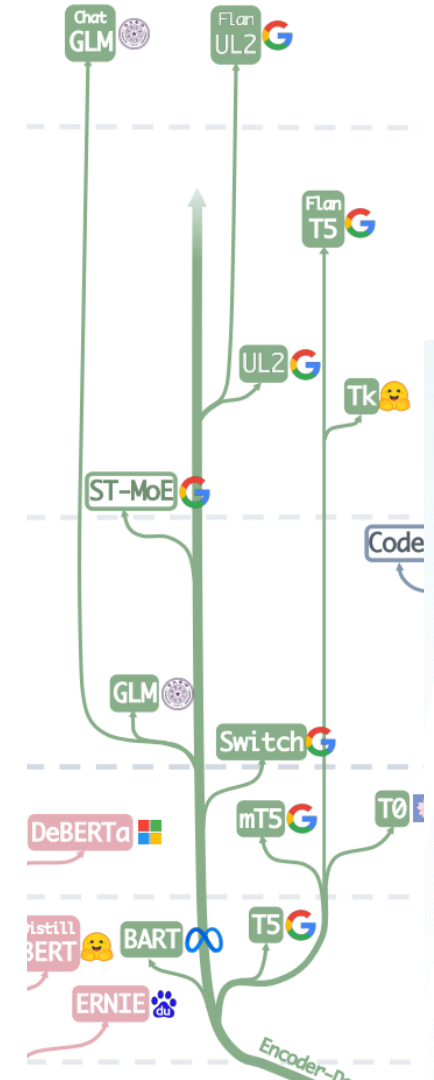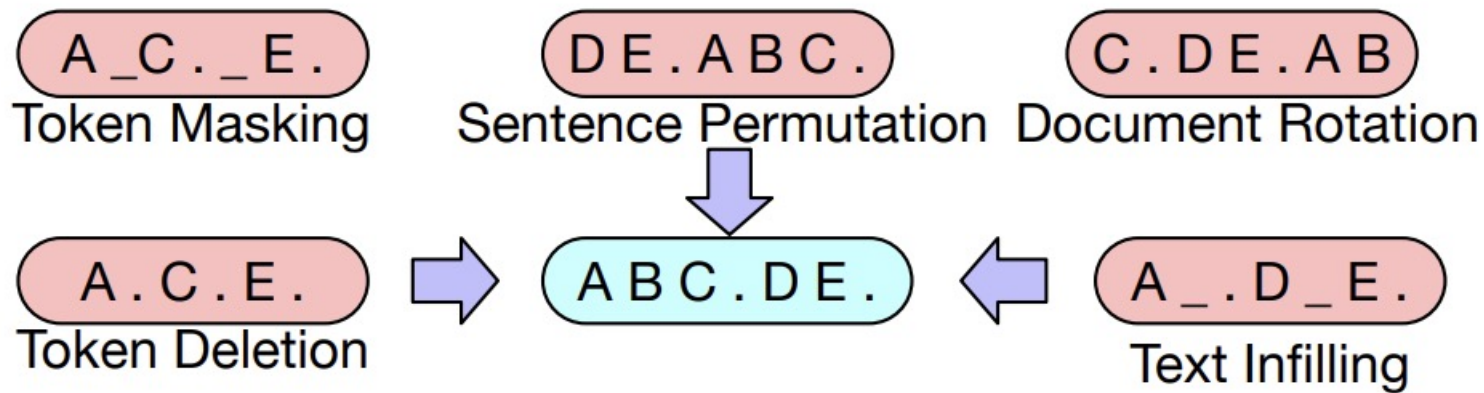# BERT - Bidirectional Encoder Representations from Transformers (2018)

*learns bidirectional encoder representations from Transformers, which are trained on large datasets as contextual language models*
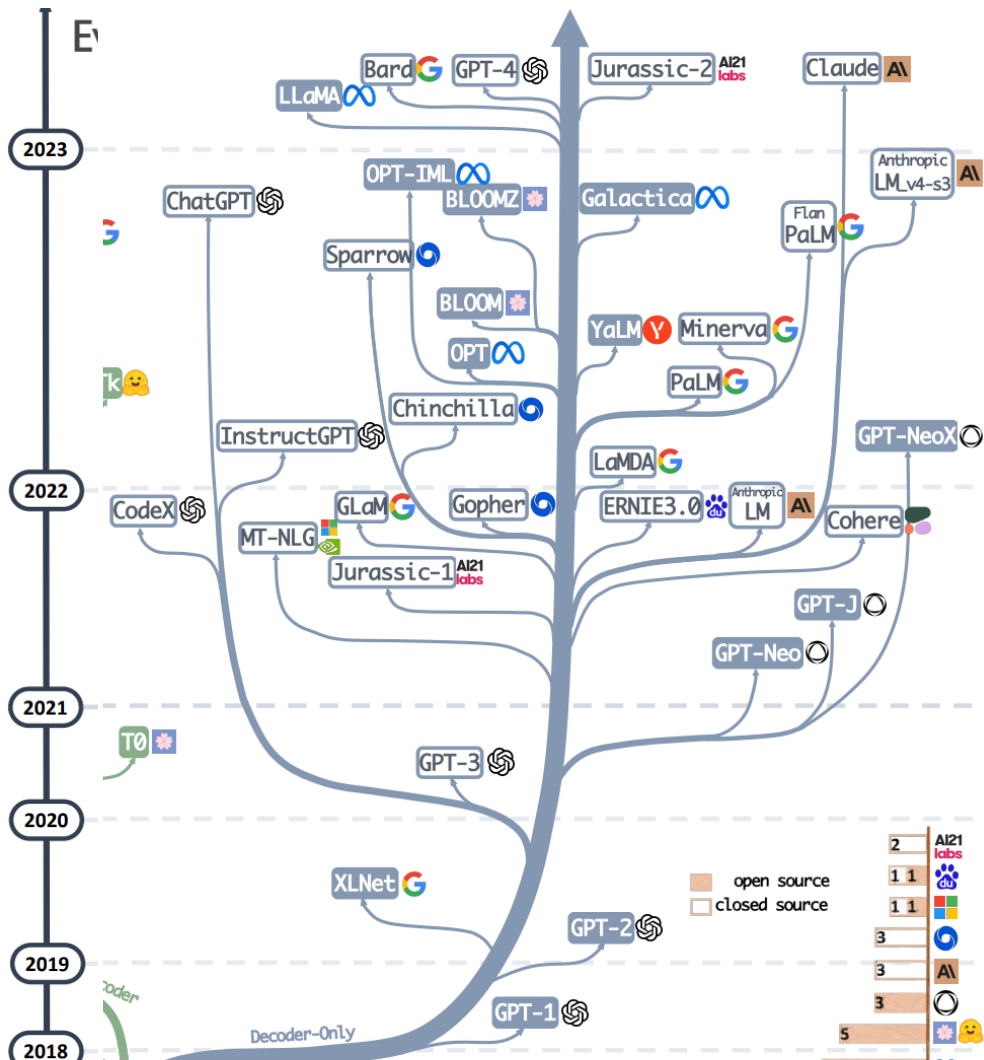


Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).

# BART

*Encoder-Decoder*



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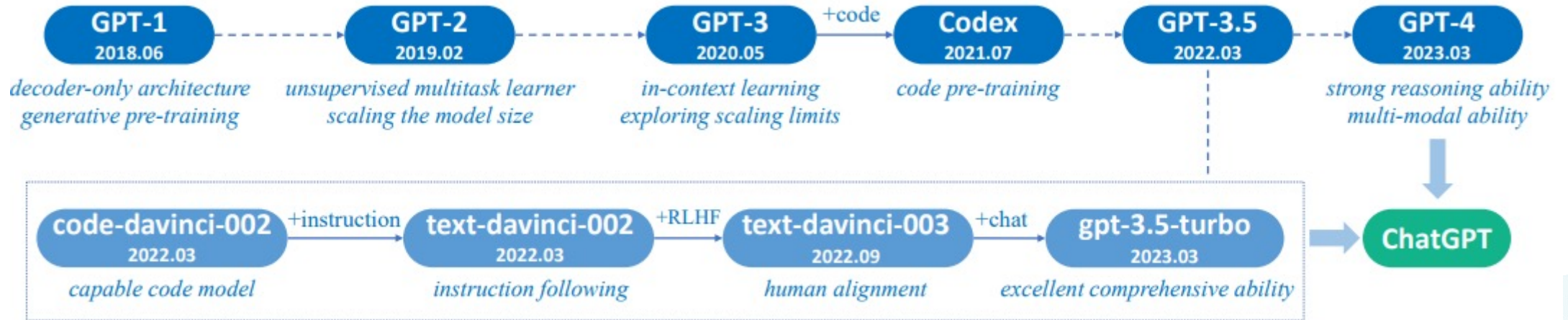
# GPT Family

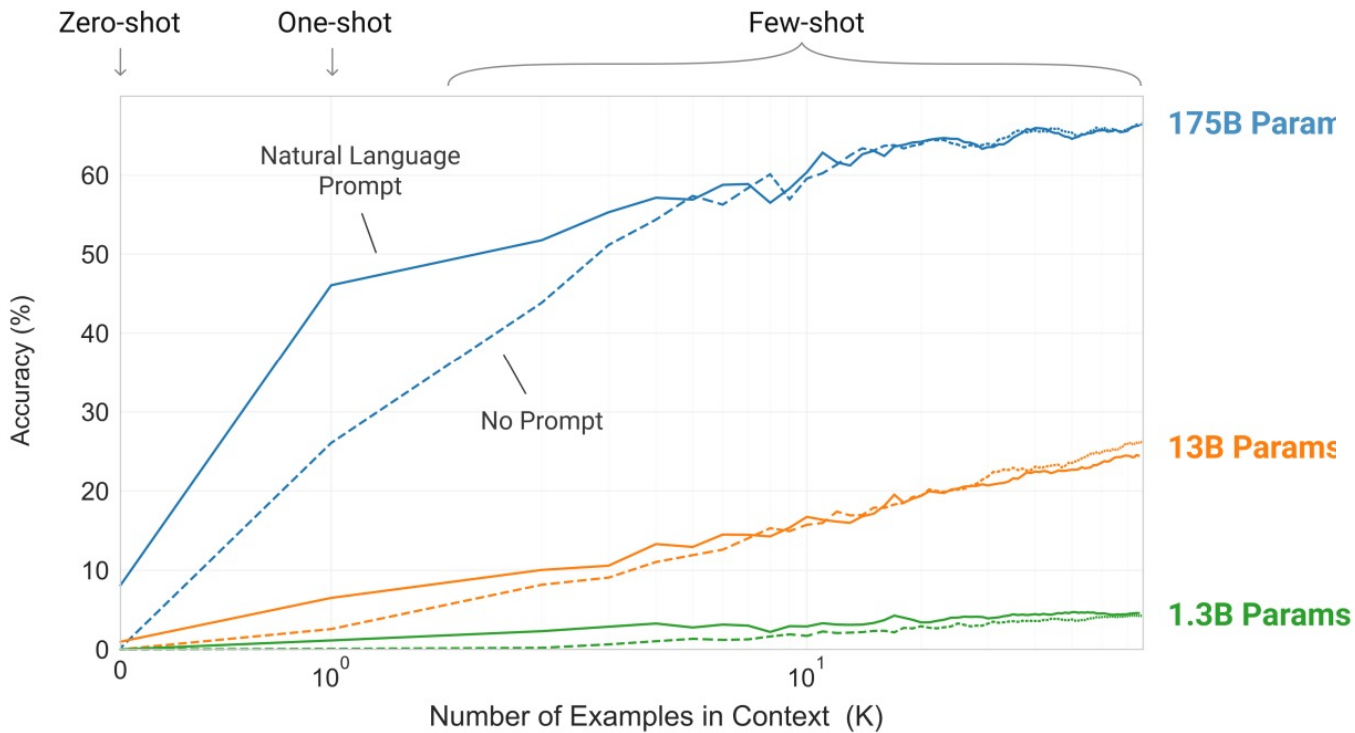**Generative Pre-trained Transformer 1 (GPT-1)**- 2018

# GPT Family

# GPT Family

**Scaling up language models significantly improves the few-shot, even zero-shot performance**





**Language Models are Few-Shot Learners**

Tom B. Brown*    Benjamin Mann*    Nick Ryder*    Melanie Subbiah*

Jared Kaplan[†]   Prafulla Dhariwal   Arvind Neelakantan   Pranav Shyam   Girish Sastry

Amanda Askell   Sandhini Agarwal   Ariel Herbert-Voss   Gretchen Krueger   Tom Henighan

Rewon Child   Aditya Ramesh   Daniel M. Ziegler   Jeffrey Wu   Clemens Winter

Christopher Hesse   Mark Chen   Eric Sigler   Mateusz Litwin   Scott Gray

Benjamin Chess   Jack Clark   Christopher Berner

Sam McCandlish   Alec Radford   Ilya Sutskever   Dario Amodei

OpenAI

# GPT3 vs chatGPT

## "How can I explain foundational language models to a child?"

### GPT-3 (text-davinci-001)

A foundational language model is a way of understanding how language works. It is a way of understanding the basic structure of language and how words are put together to form sentences.
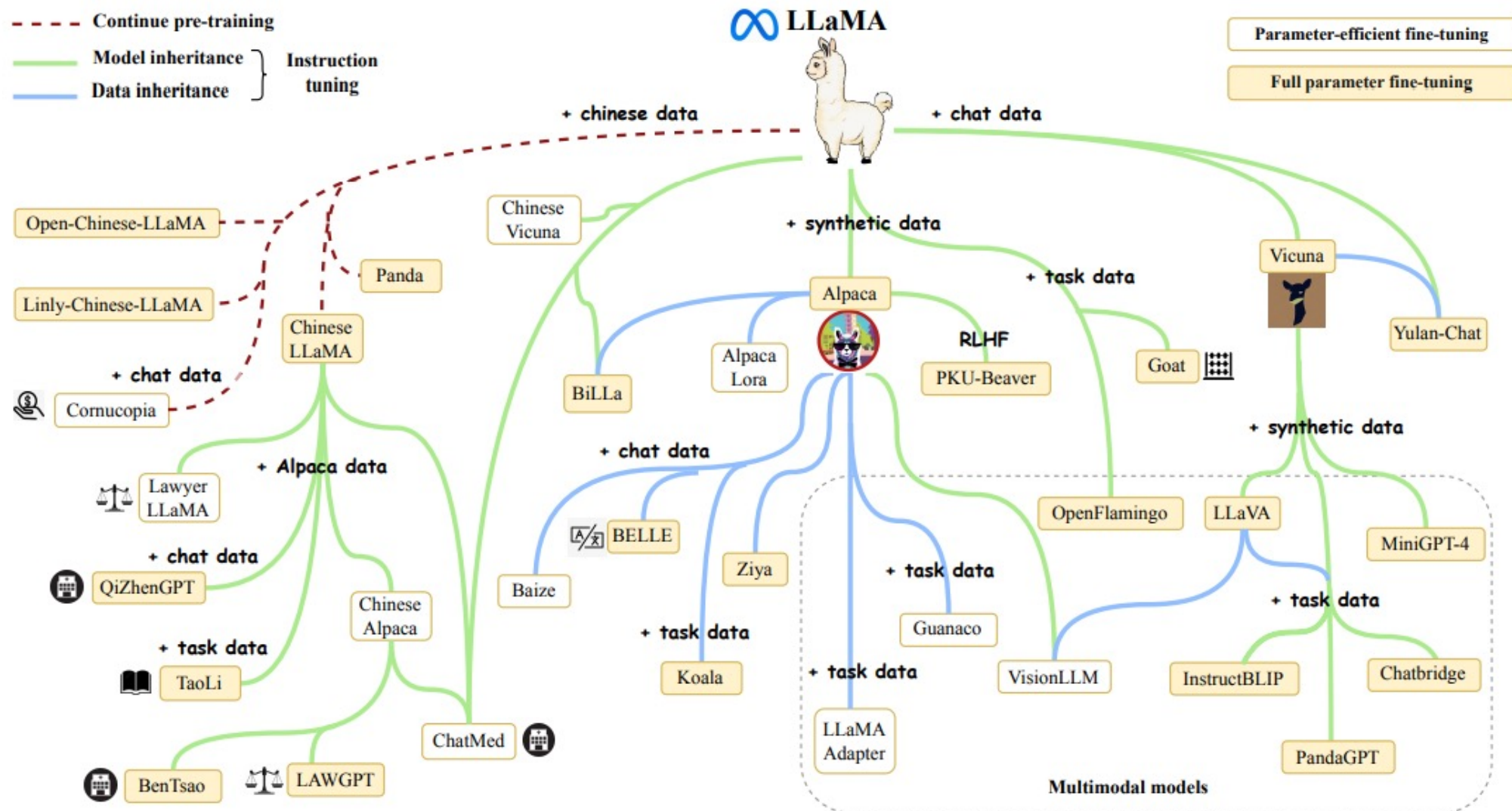
### ChatGPT

Sure, I can explain it in simple terms! Imagine that you have a really big recipe book, and in this book, there are lots and lots of recipes for all kinds of food, like cake, pizza, spaghetti, and so on. The book has so many recipes that it knows how to make almost anything you want to eat.

In a similar way, a foundational language model is like a big computer program that has read and learned lots and lots of sentences, just like you learned the recipes in the book. Now, when you ask it to write a sentence or answer a question it can use what it learned to come up with something that makes sense, just like you would use your recipe book to make a cake.

The difference is that a foundational language model can generate sentences much faster than you can flip through pages in a recipe book, and it can do it for any language it was trained on, like English, Spanish, or French. And it can answer questions or generate new text based on the information it learned.
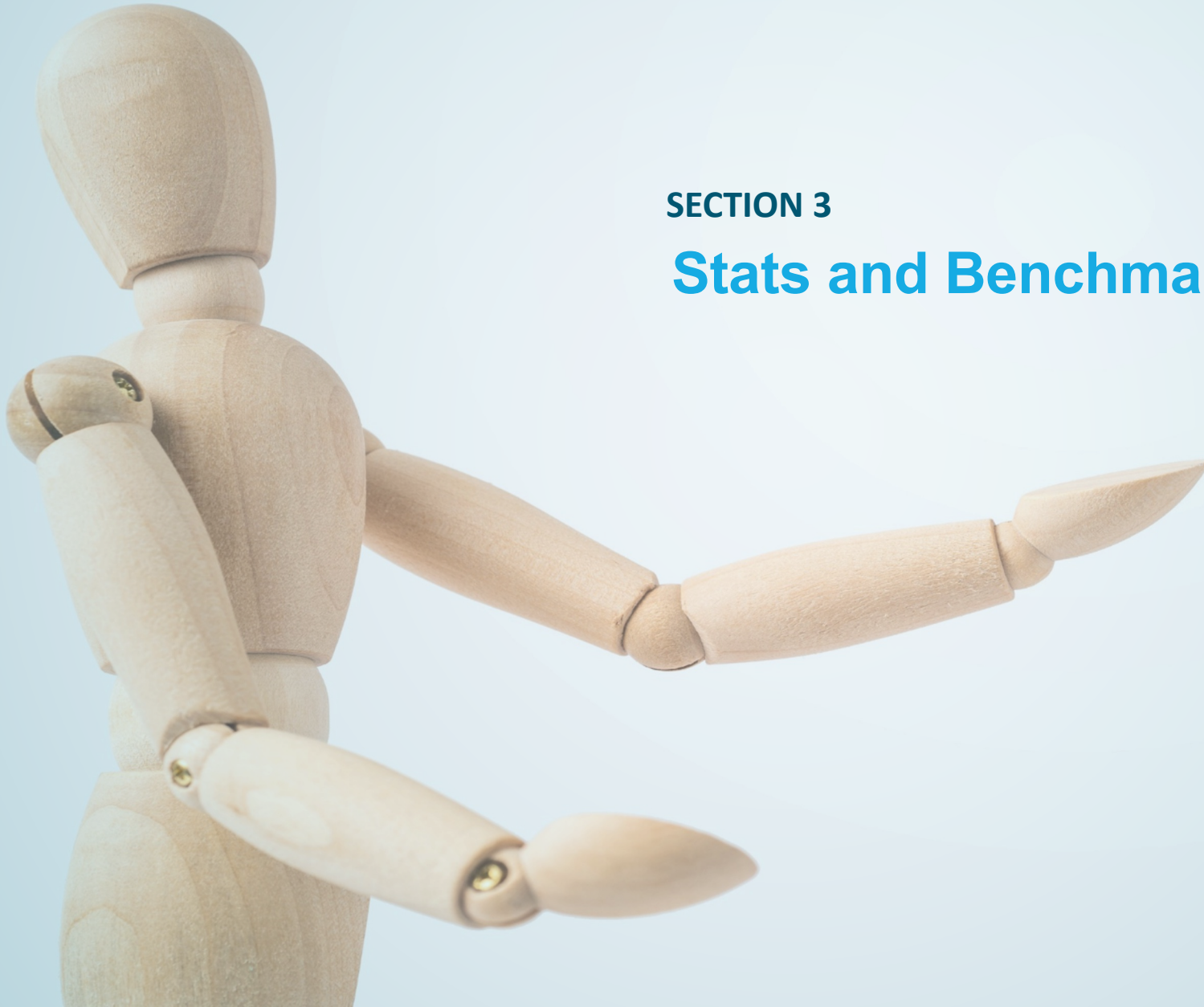
Isn't that cool? With foundational language models, computers can help us write stories, answer questions, and even have conversations, just like we are having right now.

# LLaMA Model Family

SECTION 3

# Stats and Benchmarks

# Statistics of large language models

| | Model | Release Time | Size (B) | Base Model | Adaptation IT | Adaptation RLHF | Pre-train Data Scale | Latest Data Timestamp | Hardware (GPUs / TPUs) | Training Time | Evaluation ICL | Evaluation CoT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | T5 [73] | Oct-2019 | 11 | - | - | - | 1T tokens | Apr-2019 | 1024 TPU v3 | - | ✓ | - |
| | mT5 [74] | Oct-2020 | 13 | - | - | - | 1T tokens | - | - | - | ✓ | - |
| | PanGu-α [75] | Apr-2021 | 13* | - | - | - | 1.1TB | - | 2048 Ascend 910 | - | ✓ | - |
| | CPM-2 [76] | Jun-2021 | 198 | - | - | - | 2.6TB | - | - | - | - | - |
| | T0 [28] | Oct-2021 | 11 | T5 | ✓ | - | - | - | 512 TPU v3 | 27 h | ✓ | - |
| | CodeGen [77] | Mar-2022 | 16 | - | - | - | 577B tokens | - | - | - | ✓ | - |
| | GPT-NeoX-20B [78] | Apr-2022 | 20 | - | - | - | 825GB | - | 96 40G A100 | - | ✓ | - |
| | Tk-Instruct [79] | Apr-2022 | 11 | T5 | ✓ | - | - | - | 256 TPU v3 | 4 h | ✓ | - |
| | UL2 [80] | May-2022 | 20 | - | - | - | 1T tokens | Apr-2019 | 512 TPU v4 | - | ✓ | ✓ |
| | OPT [81] | May-2022 | 175 | - | - | - | 180B tokens | - | 992 80G A100 | - | ✓ | - |
| | NLLB [82] | Jul-2022 | 54.5 | - | - | - | - | - | - | - | ✓ | - |
| | CodeGeeX [83] | Sep-2022 | 13 | - | - | - | 850B tokens | - | 1536 Ascend 910 | 60 d | ✓ | - |
| | GLM [84] | Oct-2022 | 130 | - | - | - | 400B tokens | - | 768 40G A100 | 60 d | ✓ | - |
| | Flan-T5 [64] | Oct-2022 | 11 | T5 | ✓ | - | - | - | - | - | ✓ | ✓ |
| Publicly | BLOOM [69] | Nov-2022 | 176 | - | - | - | 366B tokens | - | 384 80G A100 | 105 d | ✓ | - |
| Available | mT0 [85] | Nov-2022 | 13 | mT5 | ✓ | - | - | - | - | - | ✓ | - |
| | Galactica [35] | Nov-2022 | 120 | - | - | - | 106B tokens | - | - | - | ✓ | ✓ |
| | BLOOMZ [85] | Nov-2022 | 176 | BLOOM | ✓ | - | - | - | - | - | ✓ | - |
| | OPT-IML [86] | Dec-2022 | 175 | OPT | ✓ | - | - | - | 128 40G A100 | - | ✓ | ✓ |
| | LLaMA [57] | Feb-2023 | 65 | - | - | - | 1.4T tokens | - | 2048 80G A100 | 21 d | ✓ | - |
| | Pythia [87] | Apr-2023 | 12 | - | - | - | 300B tokens | - | 256 40G A100 | - | ✓ | - |
| | CodeGen2 [88] | May-2023 | 16 | - | - | - | 400B tokens | - | - | - | ✓ | - |
| | StarCoder [89] | May-2023 | 15.5 | - | - | - | 1T tokens | - | 512 40G A100 | - | ✓ | ✓ |
| | LLaMA2 [90] | Jul-2023 | 70 | - | ✓ | ✓ | 2T tokens | - | 2000 80G A100 | - | ✓ | - |

Source: Zhao, Wayne Xin, et al. "A survey of large language models." *arXiv preprint arXiv:2303.18223* (2023).

# Statistics of large language models

| Category | Model | Release Time | Size (B) | Base Model | | | Data Scale | | Hardware | Training Time | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Closed Source | GPT-3 [55] | May-2020 | 175 | - | - | - | 300B tokens | - | - | - | ✓ | - |
| | GShard [91] | Jun-2020 | 600 | - | - | - | 1T tokens | - | 2048 TPU v3 | 4 d | - | - |
| | Codex [92] | Jul-2021 | 12 | GPT-3 | - | - | 100B tokens | May-2020 | - | - | ✓ | - |
| | ERNIE 3.0 [93] | Jul-2021 | 10 | - | - | - | 375B tokens | - | 384 V100 | - | ✓ | - |
| | Jurassic-1 [94] | Aug-2021 | 178 | - | - | - | 300B tokens | - | 800 GPU | - | ✓ | - |
| | HyperCLOVA [95] | Sep-2021 | 82 | - | - | - | 300B tokens | - | 1024 A100 | 13.4 d | ✓ | - |
| | FLAN [62] | Sep-2021 | 137 | LaMDA-PT | ✓ | - | - | - | 128 TPU v3 | 60 h | ✓ | - |
| | Yuan 1.0 [96] | Oct-2021 | 245 | - | - | - | 180B tokens | - | 2128 GPU | - | ✓ | - |
| | Anthropic [97] | Dec-2021 | 52 | - | - | - | 400B tokens | - | - | - | ✓ | - |
| | WebGPT [72] | Dec-2021 | 175 | GPT-3 | - | ✓ | - | - | - | - | ✓ | - |
| | Gopher [59] | Dec-2021 | 280 | - | - | - | 300B tokens | - | 4096 TPU v3 | 920 h | ✓ | - |
| | ERNIE 3.0 Titan [98] | Dec-2021 | 260 | - | - | - | - | - | - | - | ✓ | - |
| | GLaM [99] | Dec-2021 | 1200 | - | - | - | 280B tokens | - | 1024 TPU v4 | 574 h | ✓ | - |
| | LaMDA [63] | Jan-2022 | 137 | - | - | - | 768B tokens | - | 1024 TPU v3 | 57.7 d | - | - |
| | MT-NLG [100] | Jan-2022 | 530 | - | - | - | 270B tokens | - | 4480 80G A100 | - | ✓ | - |
| | AlphaCode [101] | Feb-2022 | 41 | - | - | - | 967B tokens | Jul-2021 | - | - | - | - |
| | InstructGPT [61] | Mar-2022 | 175 | GPT-3 | ✓ | ✓ | - | - | - | - | ✓ | - |
| | Chinchilla [34] | Mar-2022 | 70 | - | - | - | 1.4T tokens | - | - | - | ✓ | - |
| | PaLM [56] | Apr-2022 | 540 | - | - | - | 780B tokens | - | 6144 TPU v4 | - | ✓ | ✓ |
| | AlexaTM [102] | Aug-2022 | 20 | - | - | - | 1.3T tokens | - | 128 A100 | 120 d | ✓ | ✓ |
| | Sparrow [103] | Sep-2022 | 70 | - | - | ✓ | - | - | 64 TPU v3 | - | ✓ | - |
| | WeLM [104] | Sep-2022 | 10 | - | - | - | 300B tokens | - | 128 A100 40G | 24 d | ✓ | - |
| | U-PaLM [105] | Oct-2022 | 540 | PaLM | - | - | - | - | 512 TPU v4 | 5 d | ✓ | ✓ |
| | Flan-PaLM [64] | Oct-2022 | 540 | PaLM | ✓ | - | - | - | 512 TPU v4 | 37 h | ✓ | ✓ |
| | Flan-U-PaLM [64] | Oct-2022 | 540 | U-PaLM | ✓ | - | - | - | - | - | ✓ | ✓ |
| | GPT-4 [46] | Mar-2023 | - | - | ✓ | ✓ | - | - | - | - | ✓ | ✓ |
| | PanGu-$\Sigma$ [106] | Mar-2023 | 1085 | PanGu-$\alpha$ | - | - | 329B tokens | - | 512 Ascend 910 | 100 d | ✓ | - |
| | PaLM2 [107] | May-2023 | 16 | - | ✓ | - | 100B tokens | - | - | - | ✓ | ✓ |

Source: Zhao, Wayne Xin, et al. "A survey of large language models." *arXiv preprint arXiv:2303.18223* (2023).

# Summary of PFMs

| Year | Conference | Model | Architecture | Embedding | Training method | Code |
|------|-----------|-------|-------------|-----------|----------------|------|
| 2013 | NeurIPS | Skip-Gram [68] | Word2Vec | Probabilistic | - | https://github.com/.../models |
| 2014 | EMNLP | GloVe [69] | Word2Vec | Probabilistic | - | - |
| 2015 | NeurIPS | LM-LSTM [70] | LSTM | Probabilistic | LM | https://github.com/.../GloVe |
| 2016 | IJCAI | Shared LSTM [71] | LSTM | Probabilistic | LM | https://github.com/.../adversarial_text |



Natural Language Processing
Computer Vision
Graph Learning
Other Modality…

*Source:* Zhou, Ce, et al. "A comprehensive survey on pretrained foundation models: A history from bert to chatgpt." *arXiv preprint arXiv:2302.09419* (2023).

# Downstream tasks

| Tools | Function | Link | Availability |
|---|---|---|---|
| ChatGPT | Conversational AI Chatbot | ChatGPT | Both |
| RoomGPT | Redesign your room in eight different themes | RoomGPT | Public |
| HomGPT | Redesign your home and office | HomeGPT | Subscription based |
| PDFGPT.IO | Turns PDF into the knowledge base for a ChatGPT type interface | PDFGPT | Subscription based |
| TexGPT | Harnesses GPT-3's power to help you write in Overleaf | TexGPT | Public |
| BloombergGPT | A Large Language Model for Finance | BloombergGPT | NA |
| AutoGPT | Auto-prompting without the user intervention | AutoGPT | Public |
| AgentGPT | Autonomous AI agent in the browser | AgentGPT | Public |
| XrayGPT | Automated analysis of chest radiographs based on the given x-ray | XrayGPT | Public |
| Video-ChatGPT | A vision language model for video understanding and conservation about videos | Video-ChatGPT | Public |
| ClimateGPT | Large language model for a conversation about the climate in English and Arabic | ClimateGPT | Public |
| CodeGPT | An AI assistant to find errors in code, debug code, and more | CodeGPT | Public |
| BiomedGPT | A Unified and Generalist Biomedical Generative Pre-trained Transformer for Vision, Language, and Multi-modal Tasks | BiomedGPT | Public |
| DrugGPT | A GPT based model to design potential ligands, targeting specitic proteins | DrugGPT | Public |
| Elicit | AI research assistant, automated literature reviews | Elicit | Public |
| Citation AI | AI research assistant to generate real evidence-based answers | Citation AI | Subscription based |
| Midjourey AI | AI tool to create realistic synthetic images | Mid Journey | Subscription based |
| DALL.E2 | DALL·E 2 is an AI system that can create realistic images and art from a text description | Daall-e-2 | Subscription based |
| VALL-E | An audio synthesization tool | Vall-e | Public |
| AI Avatar | Avatar generation | AI Avatar | Public |
| Langchain [114] | Building applications with LLMs through composability | Langchain | Public |

*Source:* Zhou, Ce, et al. "A comprehensive survey on pretrained foundation models: A history from bert to chatgpt." *arXiv preprint arXiv:2302.09419* (2023).

# Benchmarks

| Model | ⭐ Arena Elo rating | 📈 MT-bench (score) | MMLU | License |
|---|---|---|---|---|
| GPT-4 | 1181 | 8.99 | 86.4 | Proprietary |
| Claude-1 | 1155 | 7.9 | 77 | Proprietary |
| Claude-2 | 1134 | 8.06 | 78.5 | Proprietary |
| Claude-instant-1 | 1119 | 7.85 | 73.4 | Proprietary |
| GPT-3.5-turbo | 1115 | 7.94 | 70 | Proprietary |
| WizardLM-70b-v1.0 | 1099 | 7.71 | 63.7 | Llama 2 Community |
| Vicuna-33B | 1092 | 7.12 | 59.2 | Non-commercial |
| Llama-2-70b-chat | 1051 | 6.86 | 63 | Llama 2 Community |
| WizardLM-13b-v1.2 | 1047 | 7.2 | 52.7 | Llama 2 Community |
| Vicuna-13B | 1041 | 6.57 | 55.8 | Llama 2 Community |
| MPT-30B-chat | 1039 | 6.39 | 50.4 | CC-BY-NC-SA-4.0 |
| Guanaco-33B | 1031 | 6.53 | 57.6 | Non-commercial |

Source: https://huggingface.co/spaces/lmsys/chatbot-arena-leaderboard

# HELM

- *https://crfm.stanford.edu/helm/latest/*

# Benchmarks

# Conclusions

Table 1. Summary of Large Language Models.

| | Characteristic | | LLMs |
|---|---|---|---|
| Encoder-Decoder or Encoder-only (BERT-style) | Training: Model type: Pretrain task: | Masked Language Models Discriminative Predict masked words | ELMo [80], BERT [28], RoBERTa [65], DistilBERT [90], BioBERT [57], XLM [54], Xlnet [119], ALBERT [55], ELECTRA [24], T5 [84], GLM [123], XLM-E [20], ST-MoE [133], AlexaTM [95] |
| Decoder-only (GPT-style) | Training Model type: Pretrain task: | Autoregressive Language Models Generative Predict next word | GPT-3 [16], OPT [126]. PaLM [22], BLOOM [92], MT-NLG [93], GLaM [32],Gopher [83], chinchilla [41], LaMDA [102], GPT-J [107], LLaMA [103], GPT-4 [76], BloombergGPT [117] |

# Size Now



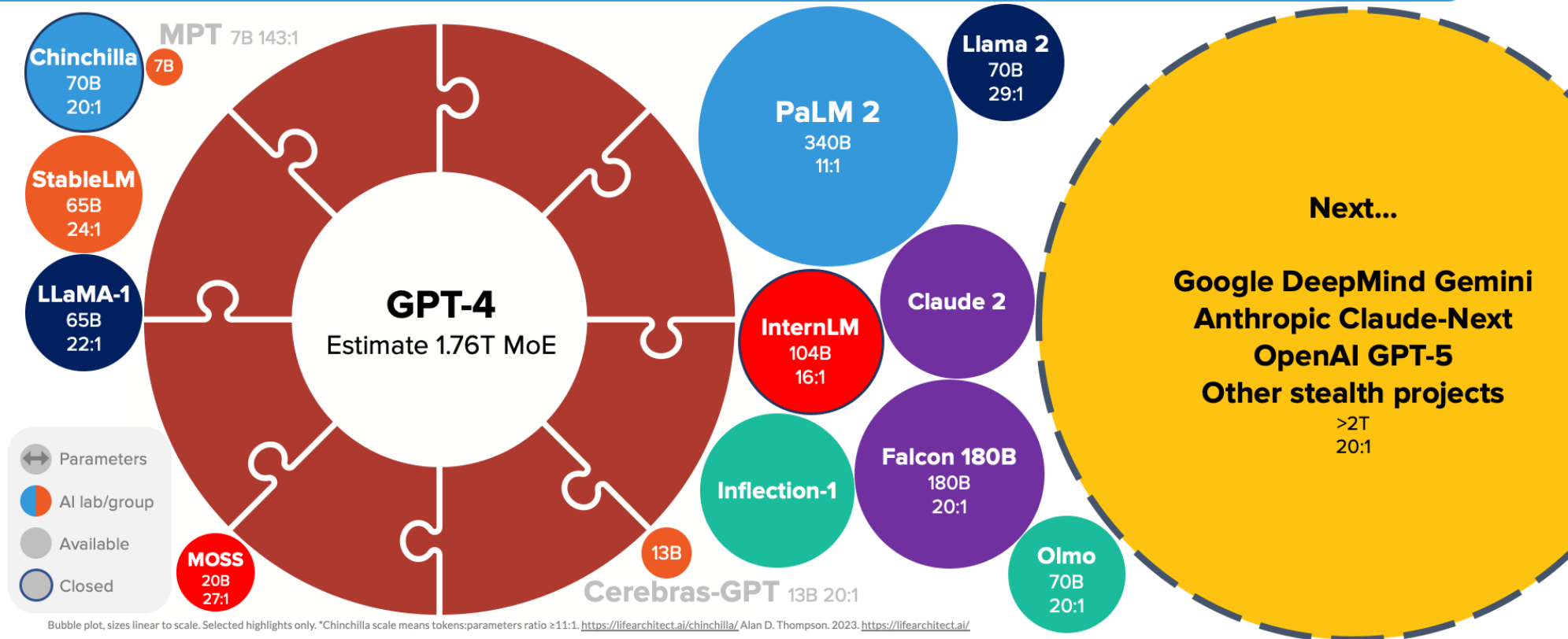Source: https://lifearchitect.ai/models/#model-bubbles

# Future Directions



2023-2024 OPTIMAL LANGUAGE MODELS — SEP/2023

Source: https://lifearchitect.ai/models/#model-bubbles

# Thank You

**Arctic LLM Workshop 2023**
**Dept. of Computer Science**