# Welcome

## From Roots to Fruits

- Journey to create value from LLMs



**Fruits**
The products you build with LLMs

**Leaves**
Building complex LLMs based systems

**Leaves**
Build complex simulations

**Leaves**
Ensemble of agents

**Leaves**
Building LLM agents

**Branches**
Unleashing the power of Large Language Models - Workflows and complex interactions

**Trunk**
Unleashing the power of Large Language Models – Direct interactions

**Roots**
Fundamentals of LLMs

# Agenda

# Why should I care?

About the importance of LLMs -  Software and hardware industries

- Software **integrating with LLMs via APIs**, like Zoom, Adobe products, etc.

- **New software** based on LLMs, like Coding Assistants, RAG systems.

- Vendors creating **AI specific hardware extensions**, like Intel® AMX or M3 neural engine.

- LLMs being combined with **agents and robotics**, like e.g: https://www.figure.ai



The Top 50 Gen AI Web Products, by Unique Monthly Visits

| 1. ChatGPT | 11. IIElevenLabs | 21. PhotoRoom | 31. PIXAI | 41. MaxAI.me |
| 2. Gemini | 12. Hugging Face | 22. YODAYO | 32. ideogram | 42. Craiyon |
| 3. character.ai | 13. Leonardo.Ai | 23. Clipchamp | 33. invideo AI | 43. OpusClip |
| 4. liner | 14. Midjourney | 24. runway | 34. Replicate | 44. BLACKBOX AI |
| 5. QuillBot | 15. SpicyChat | 25. YOU | 35. Playground | 45. CHATPDF |
| 6. Poe | 16. Gamma | 26. DeepAI | 36. Suno | 46. PIXELCUT |
| 7. perplexity | 17. Crushon AI | 27. Eightify | 37. Chub.ai | 47. Vectorizer.AI |
| 8. JanitorAI | 18. cutout.pro | 28. candy.ai | 38. Speechify | 48. DREAMGF |
| 9. CIVITAI | 19. PIXLR | 29. NightCafé | 39. phind | 49. Photomyne |
| 10. Claude | 20. VEED.IO | 30. VocalRemover | 40. NovelAI | 50. Otter.ai |

Source: https://a16z.com/100-gen-ai-apps

# Why should I care?

## Importance of LLMs - Education

- Increasing use by students of all ages.

- Major influence on the education and development of upcoming generations.

  - good effects - AI as a tutor

  - bad effects - enabling plagiarism

---

**g  You**
What is the derivative of 5x^2*sin(x)?

**ChatGPT**

To find the derivative of the function $5x^2 \sin(x)$ with respect to $x$, we need to apply the product rule, which states that the derivative of a product of two functions is the derivative of the first function times the second function plus the first function times the derivative of the second function. Let's denote the two functions as follows:

- $u(x) = 5x^2$
- $v(x) = \sin(x)$

The product rule gives us:

$$\frac{d}{dx}[u(x)v(x)] = u'(x)v(x) + u(x)v'(x)$$

Now let's find the derivatives of $u$ and $v$:

- $u'(x) = \frac{d}{dx}(5x^2) = 10x$
- $v'(x) = \frac{d}{dx}(\sin(x)) = \cos(x)$

Now we apply the product rule using the derivatives:

$$\frac{d}{dx}[5x^2 \sin(x)] = 10x \sin(x) + 5x^2 \cos(x)$$

So the derivative of $5x^2 \sin(x)$ with respect to $x$ is $10x \sin(x) + 5x^2 \cos(x)$.

Source: ChatGPT

# Why should I care?

## Importance of LLMs - Corporate world

- There is a **lot of experimentation** - companies are trying to figure out how to integrate LLMs and generative AI into their business models.

- **Enhancing corporate communications**
  - Generating, summarising, translating and correcting text.

- **Chatbots for specialist tasks**, including helpdesk and information retrieval.

- **Higher degrees of automation** - allowing more work to be done. This will impact the job market.

# Basics of neural networks

## What exactly are LLMs?

- LLMs are **deep neural networks** with a large amount of parameters.

- Analogy: Similar to **mathematical functions** (e.g., y = a + bx).



Source: Tensorflow

# Basics of neural networks

## What exactly are LLMs?

- Difference: Neural networks have **multiple layers** and a large number of parameters.

- Neural networks use interconnected nodes or **neurons** in a layered structure that resembles the human brain.



Source: Tensorflow

# Basics of neural networks

## How do neural networks train?

- Learn from **specific datasets**.

- **Self-discovery of rules** to solve tasks during training.

- The capacity to discover complex rules depends on the size of the model.

- Key components for training:

  - Defining the **network structure**

  - Selecting an **appropriate dataset**

  - Establishing a **loss function** for performance verification

  - Training process - involves **training loops**

# Basics of neural networks

## What can neural networks do?

- Solving highly complex problems:

  - Image classification and segmentation

  - Character recognition

  - Natural Language Processing (NLP) tasks

  - Video labelling



Source: https://viso.ai/deep-learning/image-segmentation-using-deep-learning

# Basics of neural networks

## What can neural networks do?

- Solving highly complex problems:

  - Speech recognition and transcription

  - Image labelling

  - Protein folding



Source: https://github.com/anindyasdas/SelfSupervisedImageText

Gametocyte surface
protein P45/48

CCR4-NOT transcription
complex subunit 9

Source: https://alphafold.ebi.ac.uk

# From words to vectors

## Training data sources

Language models trained with natural language
data sets  - open source and commercial:

- The **content of the web**

- **Code repositories**

- **Open sourced books**

- Wikipedia

- Reddit

- Wall Street Journal

**Book Corpus
dataset**

# From words to vectors

## Tokenizing words

LLM

Token sequence

Tokenizer

Raw text (Unicode code point sequence)

Source: https://www.youtube.com/watch?v=zduSFxRajkE&t=3596s

- LLMs process '**tokens**', not direct words or sentences to **reduce computational costs**.

- Example:

  Original text: "They are splashing".

  Tokens: ['They', ' are', ' spl', 'ashing'].

- Check https://platform.openai.com/tokenizer and https://tiktokenizer.vercel.app

- Tokenisation rule: approximately 4 characters per token in English.

- Estimation: 75 words ≈ 100 tokens in English.

- Form of compression.

# From words to vectors

## Counting words is challenging

- Try this prompt on any LLM:

  "How many words are in the
  response to this prompt?"

  Note: more LLM shortcomings are going
  to be mentioned later



Source: ChatGPT

# From words to vectors

## Tokens converted to numbers

- Tokens are then **converted to numbers**.

- Numbers are **indices in a token vocabulary** associated to the LLM.



Source: https://platform.openai.com/tokenizer
Also see https://tiktokenizer.vercel.app

# From words to vectors

Vocabulary numbers converted into vectors

- The vocabulary number gets **converted into vectors** using a lookup table.

- The token gets positioned in a **multidimensional space**.

- You can also convert sentences to vectors and store them outside in **vector databases**.

- Vector databases allow **nearest neighbour searches** and open up new possibilities for information retrieval

# From words to vectors

## Vocabulary numbers converted into vectors

### 1. Tokenise

"We are splashing" → ["We", "are", "spl", "ashing"]

### 2. Map each token

["We", "are", "spl", "ashing"] → [1367, 4536, 2134, 3213]

### 3. Map each ID to a vector

1367 → [-8.43e-03 5.96e-03 ...]

4536 → [1.62e-02 -2.41e-02 -2.79e-02 ...]

2134 → [3.98e-04 -5.75e.-03 8.15e-03 ...]

3213 → [-3.56-02 -9.999e-03 -2.37e-02 ...]

### 4. Vectors get aligned through training



Before

After

# Transformer architecture
## Origins of the transformer architecture

- Based on the seminal paper "Attention Is All You Need" (2017).

- Published by researchers primarily associated with Google.

- **Originally focused on language translation** which "transforms" an input to an output sequence.

## 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

### 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 mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

# Transformer architecture

## Unforeseen impacts

- Foundation for today's LLMs.

- **Unforeseen revolutionary impact** on AI and natural language processing.

- Allows significantly **more parallelisation** which produced better and faster results.

# Transformer architecture

## Key aspects of the transformer architecture

- Parallelized calculation of word relationships for efficient scaling.

- Self-attention layer:

  - Core component of the Transformer architecture.

  - Function: calculates interrelations of words (tokens) in a sentence (sequence).

  - Conceptualization: a formula for assessing word relationships.



Source: https://jalammar.github.io/illustrated-transformer

# LLM training

## How LLMs are trained

The way they are trained explains their behaviour and can give us clues on how to use them in a better way.

**The LLM training has typically three stages:**

- Pre-training phase of Large Language Models (LLMs).

- Instruction tuning.

- Reinforcement Learning from Human Feedback.

# LLM training

Pre-training phase of Large Language Models (LLMs):

- Objective: **understanding language**, learn about real world concepts, learn about context.

- Training data: **lots of text**, like Wikipedia or Common Crawl dataset.

- Training data processing: **omit words in sentences** and let model **guess the word**.

- Very **resource intensive** – requires powerful equipment and good amount of time.

# LLM training

## Instruction tuning

- Training models with **specialised data sets** relevant to **specific final tasks**.

- Example: train a model to chat, train a model on a variety of tasks, like generating code, translating, detecting sentiment.

| id<br>string · lengths | system_prompt<br>string · classes | question<br>string · lengths | response<br>string · lengths |
|---|---|---|---|
| 4       12 | 17 values | 12    40.6k | 0    15k |
| niv.242684 | | You will be given a definition of a task first,… | [ ["AFC Ajax (amateurs)", "has ground", "Sportpark De Toekomst"], ["Ajax Youth Academy", "plays at",… |
| flan.564327 | You are an AI assistant. You will be given a task. You… | Generate an approximately fifteen-word sentence that… | Midsummer House is a moderately priced Chinese restaurant with a 3/5 customer rating, located nea… |
| flan.1875913 | You are a helpful assistant, who always provide… | What happens next in this paragraph? She then rubs a… | C. She then dips the needle in ink and using the pencil to draw a design on her leg, rubbing it off… |
| t0.408370 | You are an AI assistant. You will be given a task. You… | Please answer the following question: I want to test th… | Based on the passage, discuss the primary motivations and outcomes of the 1901 Federation of… |
| cot.86217 | You are an AI assistant that helps people find… | James runs a TV show and there are 5 main characters… | James pays the minor characters $15,000 each episode. Since there are 4 minor characters, he… |
| cot.18180 | You are an AI assistant that helps people find… | Given the stream of consciousness rationale,… | Question: What is the proper technique for a female beach volleyball player to serve the ball… |

‹ Previous   **1**   2   3   ...   29,149   Next ›

Source: https://huggingface.co/datasets/Open-Orca/OpenOrca

# LLM training

## Reinforcement Learning from Human Feedback (RLHF)

- Objective: **align the model with human preferences**.

- Method:

  - Phase 1: **humans rank multiple answers** to a single question and train a reward model to be an advisor to LLM.

  - Phase 2: **reward model fine tunes** LLM.

- Implications: potential for LLM to adopt specific behaviour patterns based on trainer preferences, like e.g: **preferring confidence over accuracy**.

# LLM examples (1)

## Most prominent LLMs in June 2024

**OpenAI ChatGPT4.0**

- Most advanced model.

- Best performance according to [chatbot arena](#)

**Gemini**

- Processes text, code, and images.

- Extended memory.

- Better understanding of complex tasks.

**Claude**

- Good at creative writing.

- Large context window.

- Claude 3 takes on ChatGPT 4.

https://platform.openai.com

https://gemini.google.com/app

https://claude.ai

# LLM examples (2)

## Most prominent LLMs in June 2024

**MISTRAL AI_**

- Sparse mixture of experts.

- Combining specialised models into one large model.

- Inference is much faster.

https://mistral.ai

**llama3**

- Best open source model according to chatbot arena.

- Comparatively small model, but very competitive.

https://llama.meta.com/llama3

**Yi**

- Newcomer from China

- Excels in Chinese

- Comparatively small model, but very competitive.

https://www.01.ai

# Comparing LLMs

## Leaderboards

- There is a growing number of LLMs. Which one shall I pick? Here is one of the most useful leaderboards: https://chat.lmsys.org/?leaderboard. This leaderboard is not based on benchmarks.

- This leaderboard is **based on duels** between chat models on https://arena.lmsys.org/ from which an ELO ranking is determined. Anyone can contribute to this leaderboard.

| Rank* (UB) | 🤖 Model | ⭐ Arena Elo | 📊 95% CI | 🗳 Votes | Organization | License | Knowledge Cutoff |
|---|---|---|---|---|---|---|---|
| 1 | GPT-4o-2024-05-13 | 1287 | +4/-4 | 32181 | OpenAI | Proprietary | 2023/10 |
| 2 | Gemini-1.5-Pro-API-0514 | 1267 | +5/-4 | 25519 | Google | Proprietary | 2023/11 |
| 2 | Gemini-Advanced-0514 | 1266 | +5/-5 | 27225 | Google | Proprietary | Online |
| 4 | Gemini-1.5-Pro-API-0409-Preview | 1257 | +3/-3 | 55731 | Google | Proprietary | 2023/11 |
| 4 | GPT-4-Turbo-2024-04-09 | 1256 | +2/-2 | 59891 | OpenAI | Proprietary | 2023/12 |
| 5 | GPT-4-1106-preview | 1251 | +2/-3 | 80067 | OpenAI | Proprietary | 2023/4 |
| 6 | Claude 3 Opus | 1248 | +2/-2 | 123645 | Anthropic | Proprietary | 2023/8 |
| 6 | GPT-4-0125-preview | 1246 | +3/-2 | 73286 | OpenAI | Proprietary | 2023/12 |
| 9 | Yi-Large-preview | 1239 | +4/-3 | 34567 | 01 AI | Proprietary | Unknown |
| 9 | Gemini-1.5-Flash-API-0514 | 1232 | +4/-4 | 23797 | Google | Proprietary | 2023/11 |

Leaderboard updated: 02 June 2024        Source: https://chat.lmsys.org/?leaderboard
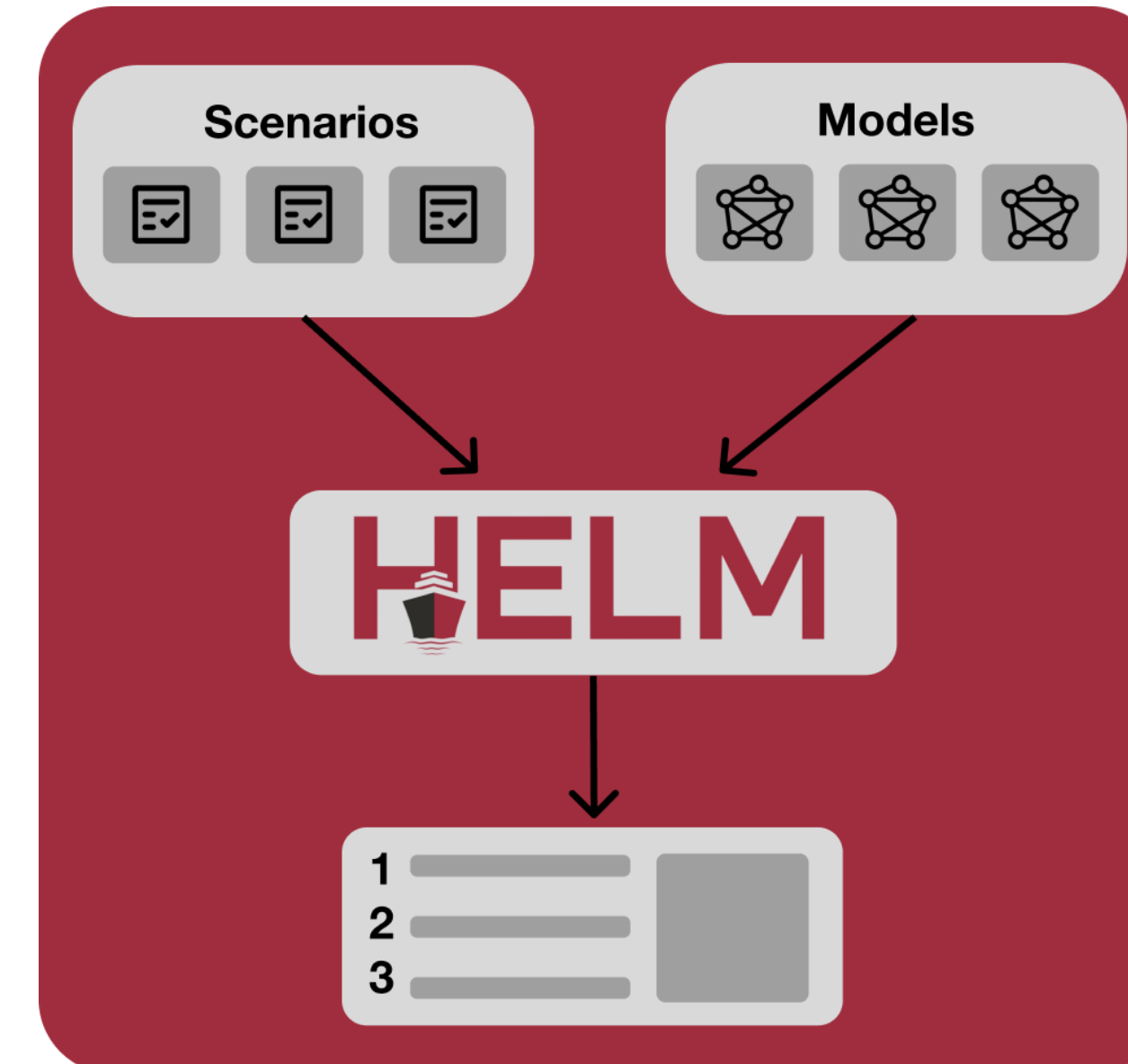
# Comparing LLMs

## Leaderboards

- HELM (Holistic Evaluation of Language Models) leaderboards are scenario based benchmarks in different categories:

  - Massive Multitask Language Understanding

  - Text to image

  - Vision Language Models

  - ...

- More **prone to manipulation**

| Model | MMLU All Subjects - EM |
|---|---|
| Claude 3 Opus (20240229) | **0.846** |
| GPT-4o (2024-05-13) | 0.842 |
| GPT-4 (0613) | 0.824 |
| Gemini 1.5 Pro (0409 preview) | 0.81 |
| GPT-4 Turbo (1106 preview) | 0.796 |
| Llama 3 (70B) | 0.793 |
| Palmyra X V3 (72B) | 0.786 |
| PaLM-2 (Unicorn) | 0.786 |
| Mixtral (8x22B) | 0.778 |
| Gemini 1.5 Flash (0514 preview) | 0.778 |

Source: https://crfm.stanford.edu/helm/mmlu/latest



More information: https://crfm.stanford.edu/helm

# LLMs and AGI

## What is AGI (Artificial General Intelligence)?

- Hypothetical type of AI which can **perform as well or better than humans on a wide range of cognitive tasks.**

- Characteristics:

  - Self-learning

  - Problem solving skills

  - Understanding and reasoning

# Shortcomings of LLM

## Problems related to LLMs

- **Hallucinations**: LLMs tend to invent facts in moments of uncertainty.

- **Refusals**: LLMs refusals to answer challenging questions in order to avoid errors.

- **Jailbreaks**: Most models were trained with ethical boundaries. However, you can work around them and retrieve unethical content.

- **Word counting**

- **Knowledge cut-offs**

- **Missing sources**

Source: ChatGPT

# Conclusion

## Wrapping up

- LLMs represent a new **technological milestone** impacting our daily lives.

- LLMs convert language into a vector space and **enable the usage of vector databases**.

- LLMs are trained to understand language, on a **wide range of tasks** and to **adapt to human preferences**.

- AI is getting better at a wide variety of tasks allowing businesses to rethink their strategies.

# Credits

- All images in this presentation were generated either by DALL-E or Midjourney.

- The parts about the token handling were based on Andrej Karpathy's Youtube video: https://www.youtube.com/watch?v=zduSFxR ajkE&t=3596s
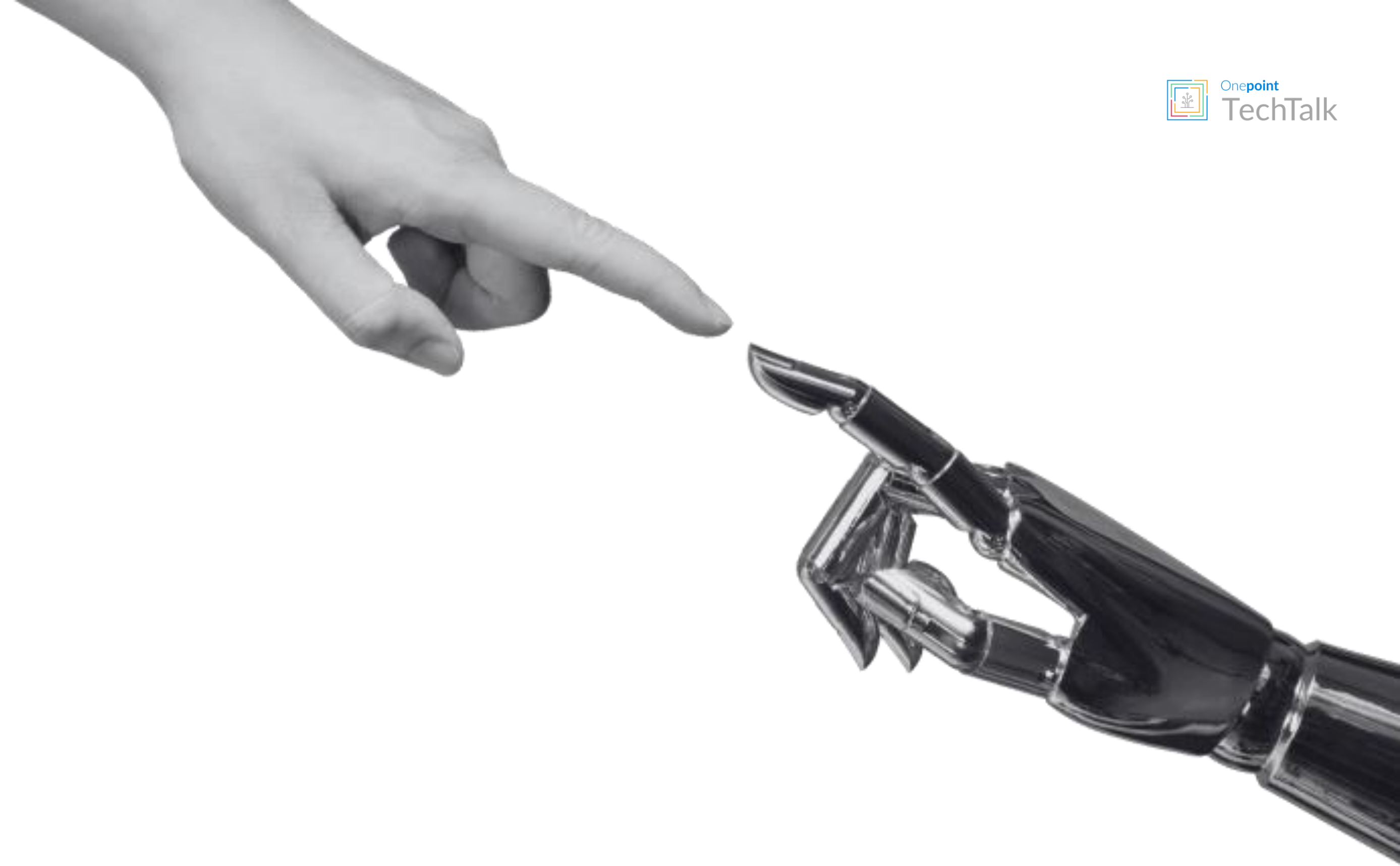
**DALL·E**

**Midjourney**

# Thank you for joining

Please feel free to contact Gil Fernandes if you
have any feedback about the session.

Email at techtalk@onepointltd.com

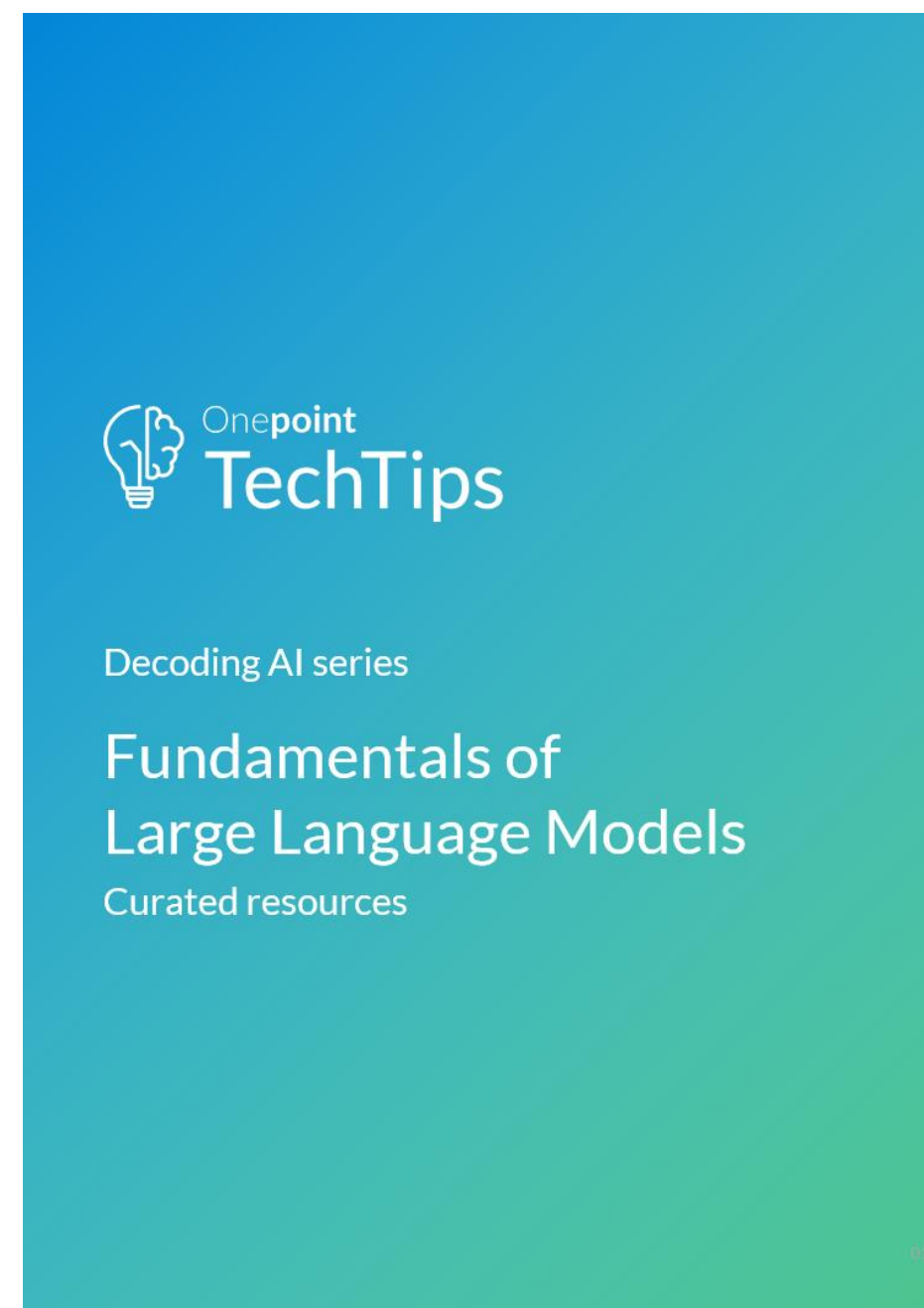Connect on LinkedIn https://www.linkedin.com/in/gil-palma-fernandes

Find Gil's Reflections on AI at  https://medium.com/@gil.fernandes

# A bonus for you

We have put together a collection of resources
that will deepen your understanding of LLMs.

Onepoint
TechTips

Decoding AI series

Fundamentals of
Large Language Models
Curated resources

# Onepoint
# TechTalk

**Coming soon**

Session 2

# Unleashing the power of Large Language Models

Part 1 – Direct interactions

## Onepoint

*Your trusted companions for the digital journey™*