# I am human… for now

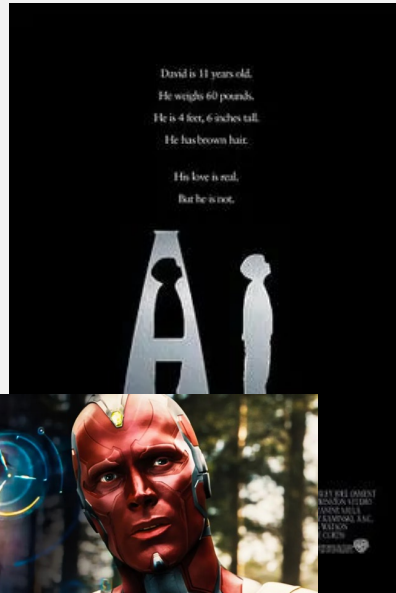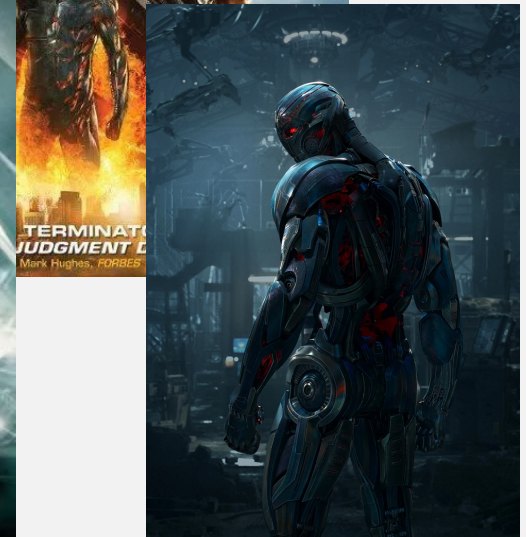# Movies on AI

**Friend?**

**Foe?**

🤔 **Poll: What is generative AI?**

## How many of you think that…
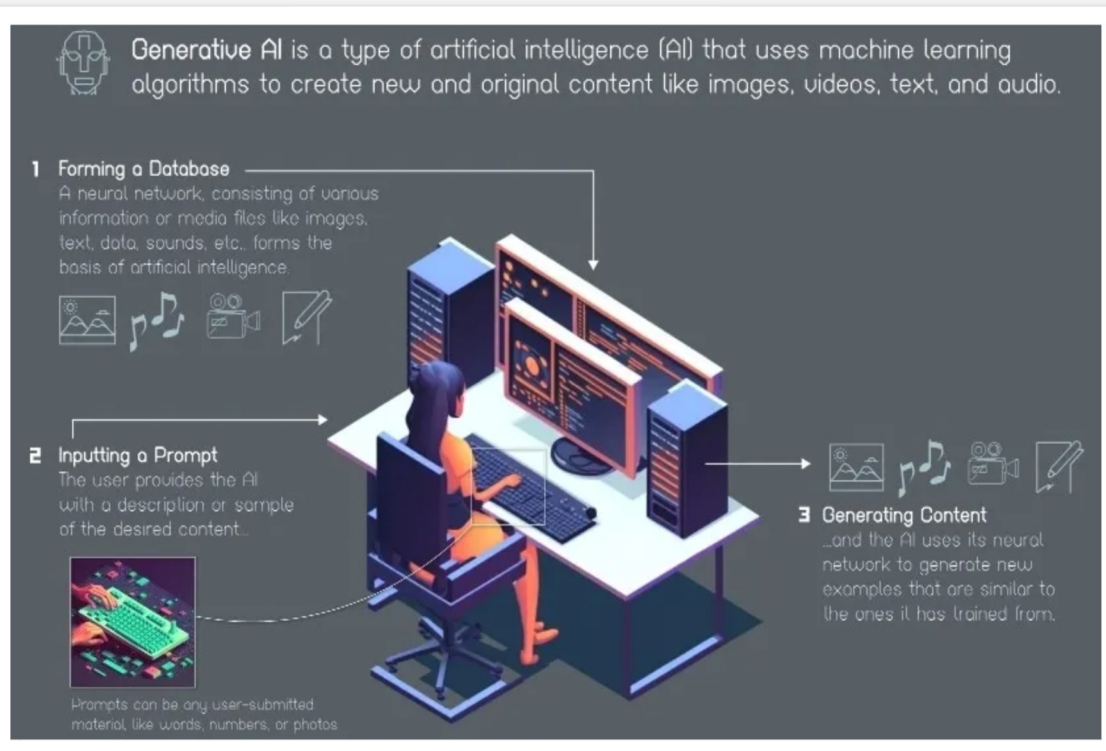
🤖 AI have surpassed human capabilities

🤖 Gen AI is a superior version type of AI

🤖 Gen AI is actually just AI but with the capability of producing content
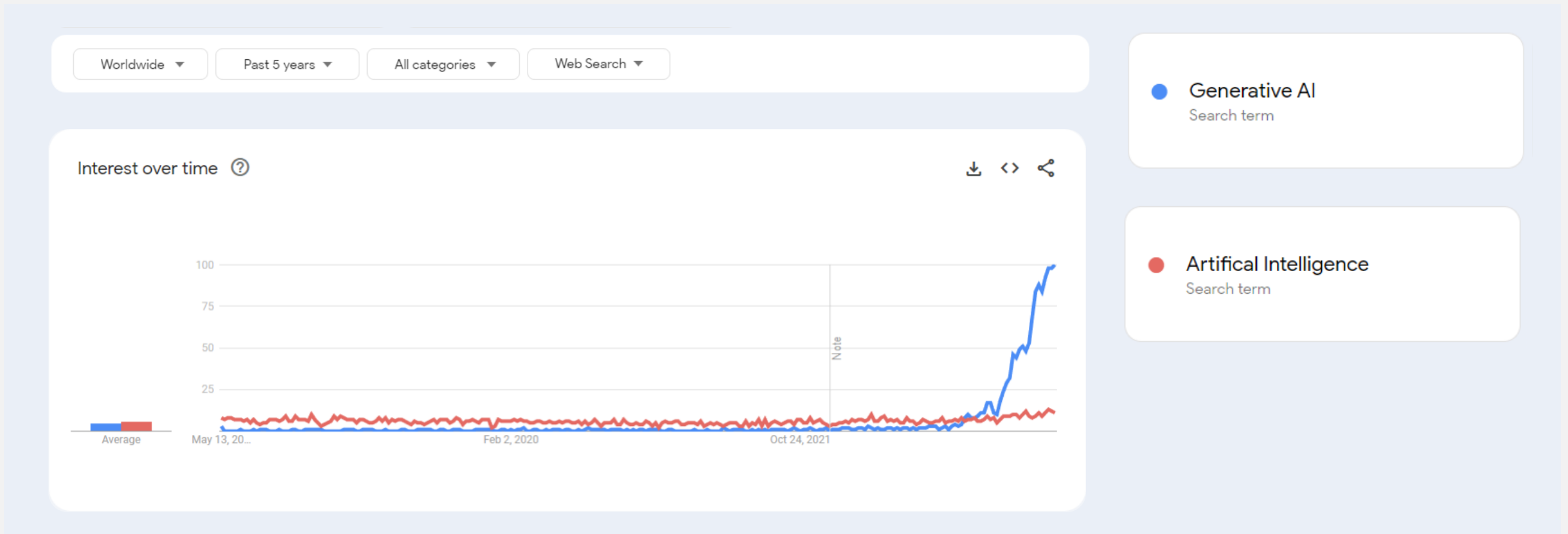
❌ DBS

# What is generative AI?

"Generative AI refers to a category of artificial intelligence (AI) algorithms that **generate new outputs** based on the data they have been trained on. Unlike traditional AI systems that are designed to recognize patterns and make predictions, generative AI creates new content in the form of images, text, audio, and more" (Routley, 2023).



Generative AI is a type of artificial intelligence (AI) that uses machine learning algorithms to create new and original content like images, videos, text, and audio.

1 **Forming a Database**
A neural network, consisting of various information or media files like images, text, data, sounds, etc., forms the basis of artificial intelligence.

2 **Inputting a Prompt**
The user provides the AI with a description or sample of the desired content...

Prompts can be any user-submitted material like words, numbers, or photos.

3 **Generating Content**
...and the AI uses its neural network to generate new examples that are similar to the ones it has trained from.

**+ more…**

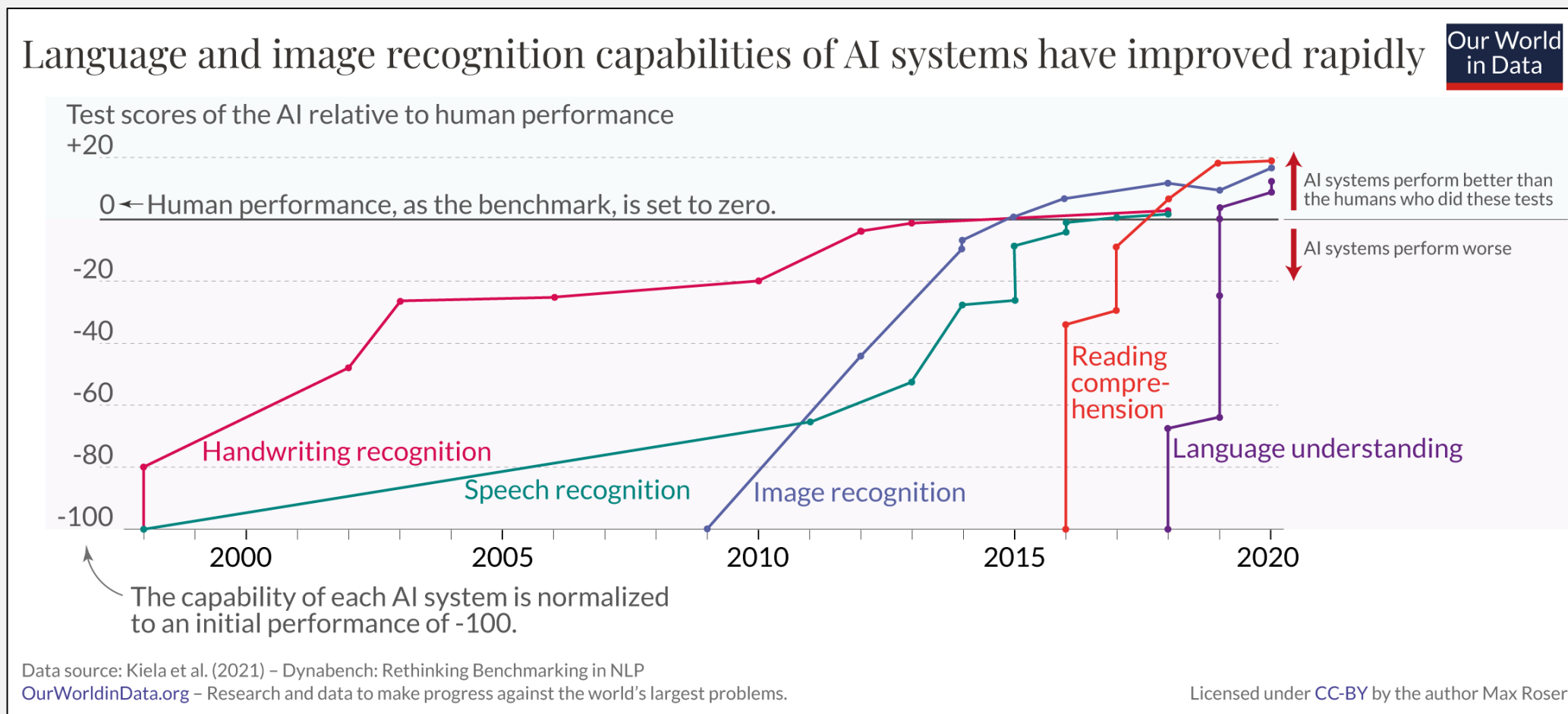# Over the last couple of month, the term generative AI have gained a lot of traction

On Google Trends...

# A brief history…

Here is a timeline of how language and image recognition of AI systems have improved rapidly



Language and image recognition capabilities of AI systems have improved rapidly

Our World in Data

Test scores of the AI relative to human performance

+20

0 ← Human performance, as the benchmark, is set to zero.

AI systems perform better than the humans who did these tests

AI systems perform worse

-20

-40

Reading compre- hension

-60

Handwriting recognition

-80

Speech recognition

Image recognition

Language understanding

-100

2000    2005    2010    2015    2020

The capability of each AI system is normalized to an initial performance of -100.

Data source: Kiela et al. (2021) – Dynabench: Rethinking Benchmarking in NLP
OurWorldinData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the author Max Roser

❌ DBS

(Kiela et al., 2021; Roser, 2022)

# A brief history…

Here is a chart of all the AI tools that have been built and the amount of training computation used…

(Sevilla et al., 2022; Roser, 2022)

❌ DBS

---

**The rise of artificial intelligence over the last 8 decades:** As training computation has increased, AI systems have become more powerful

Our World in Data

The color indicates the domain of the AI system: ● Vision ● Games ● Drawing ● Language ● Other

Shown on the vertical axis is the **training computation** that was used to train the AI systems.

10 billion petaFLOP

Computation is measured in floating point operations (FLOP). One FLOP is equivalent to one addition, subtraction, multiplication, or division of two decimal numbers.

100 million petaFLOP

The data is shown on a logarithmic scale, so that from each grid-line to the next it shows a 100-fold increase in training computation.

1 million petaFLOP

10,000 petaFLOP

100 petaFLOP

1 petaFLOP = 1 quadrillion FLOP

10 trillion FLOP

100 billion FLOP

1 billion FLOP

10 million FLOP

100,000 FLOP

1,000 FLOP

10 FLOP

**Minerva:** built in 2022 and trained on 2.7 billion petaFLOP
*Minerva can solve complex mathematical problems at the college level.*

**PaLM:** built in 2022 and trained on 2.5 billion petaFLOP
*PaLM can generate high-quality text, explain some jokes, cause & effect, and more.*

**GPT-3:** 2020; 314 million petaFLOP
*GPT-3 can produce high-quality text that is often indistinguishable from human writing.*

**DALL-E:** 2021; 47 million petaFLOP
*DALL-E can generate high-quality images from written descriptions.*

**NEO:** 2021; 1.1 million petaFLOP
*Recommendation systems like Facebook's NEO determine what you see on your social media feed, online shopping, streaming services, and more.*

**AlphaGo:** 2016; 1.9 million petaFLOP
*AlphaGo defeated 18-time champion Lee Sedol at the ancient and highly complex board game Go. The best Go players are no longer human.*

**AlphaFold:** 2020; 100,000 petaFLOP
*AlphaFold was a major advance toward solving the protein-folding problem in biology.*

**MuZero:** 2019; 48,000 petaFLOP
*MuZero is a single system that achieved superhuman performance at Go, chess, and shogi (Japanese chess) — all without ever being told the rules.*

**AlexNet:** 2012; 470 petaFLOP
*A pivotal early "deep learning" system, or neural network with many layers, that could recognize images of objects such as dogs and cars at near-human level.*

NPLM

Decision tree
LSTM
**TD-Gammon:** 1992; 18 trillion FLOP
*TD-Gammon learned to play backgammon at a high level, just below the top human players of the time.*

LeNet-5

RNN for speech

**NetTalk:** 1987; 81 billion FLOP
*NetTalk was able to learn to pronounce some English text by being given text as input and matching it to phonetic transcriptions. Among its many limitations, it did not perform the visual recognition of the text itself.*

ALVINN
Zip CNN

System 11

Pandemonium (Morse)
Samuel Neural Checkers

Back-propagation
**Neocognitron:** 1980; 228 million FLOP
*A precursor of modern vision systems. It could recognize handwritten Japanese characters and a few other patterns.*

Fuzzy NN

**Perceptron Mark I:** built in 1957/58; 695,000 FLOP
*Regarded as the first artificial neural network, it could visually distinguish cards marked on the left side from those marked on the right, but it could not learn to recognize many other types of patterns.*

**ADALINE:** built in 1960 and trained on around 9,900 FLOP
*An early single-layer artificial neural network.*

**Theseus:** built in 1950 and trained on around 40 floating point operations (FLOP)
*Theseus was a small robotic mouse, developed by Claude Shannon, that could navigate a simple maze and remember its course.*

Pre Deep Learning Era
Training computation grew in line with Moore's law, doubling roughly every 20 months.

Deep Learning Era
Increases in training computation accelerated, doubling roughly every 6 months.

*The first electronic computers were developed in the 1940s*

1940   1950   1960   1970   1980   1990   2000   2010   2020

*1956: The Dartmouth workshop on AI, often seen as the beginning of the field of AI research*

*1997: Deep Blue beats world chess champion Garry Kasparov*

# Why is improving productivity important?

## Can we leverage on Gen AI to boost productivity?

🤔

# Improving productivity allow us to...

**Improve employee satisfaction**

**Focus on developing quality applications**

**Enjoy faster time to market**

**Reduce toil and manual work**

# What software are out there for developers?
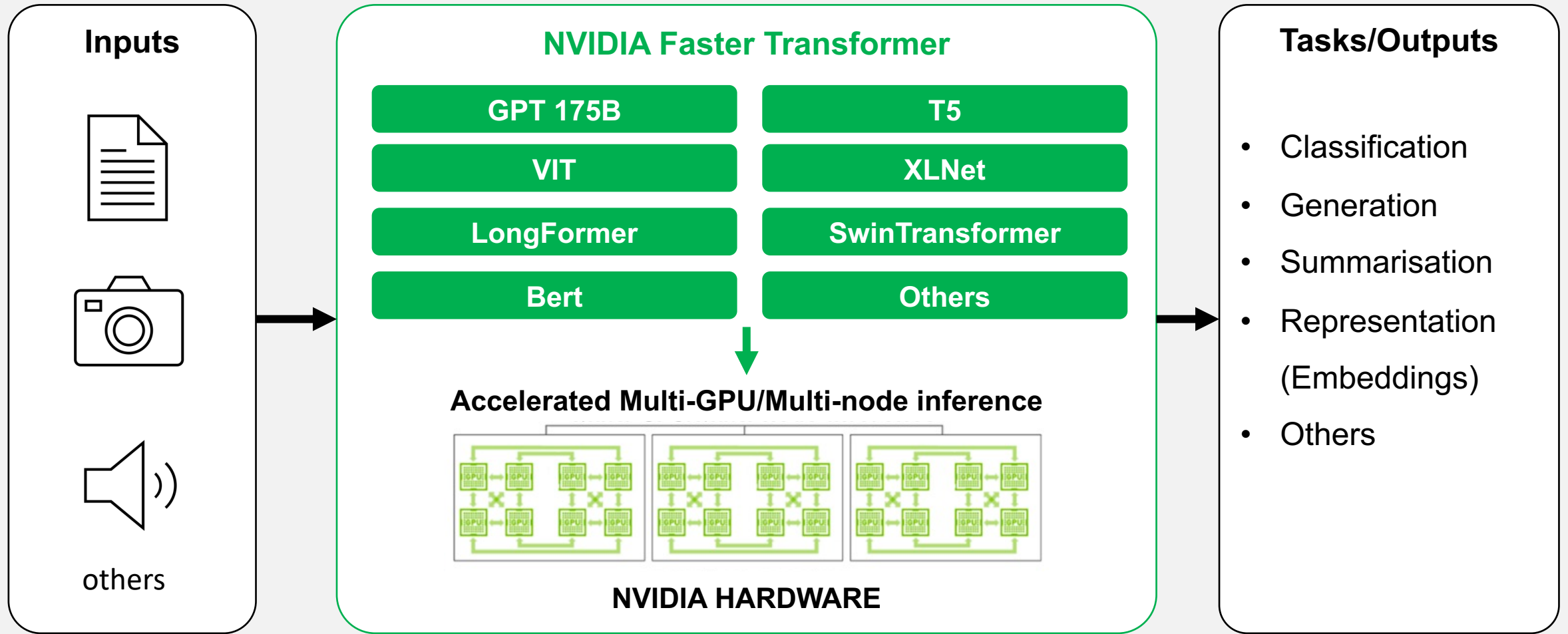
aws CodeWhisperer

ChatGPT

GitHub Copilot

G Bard

tabnine

DBS
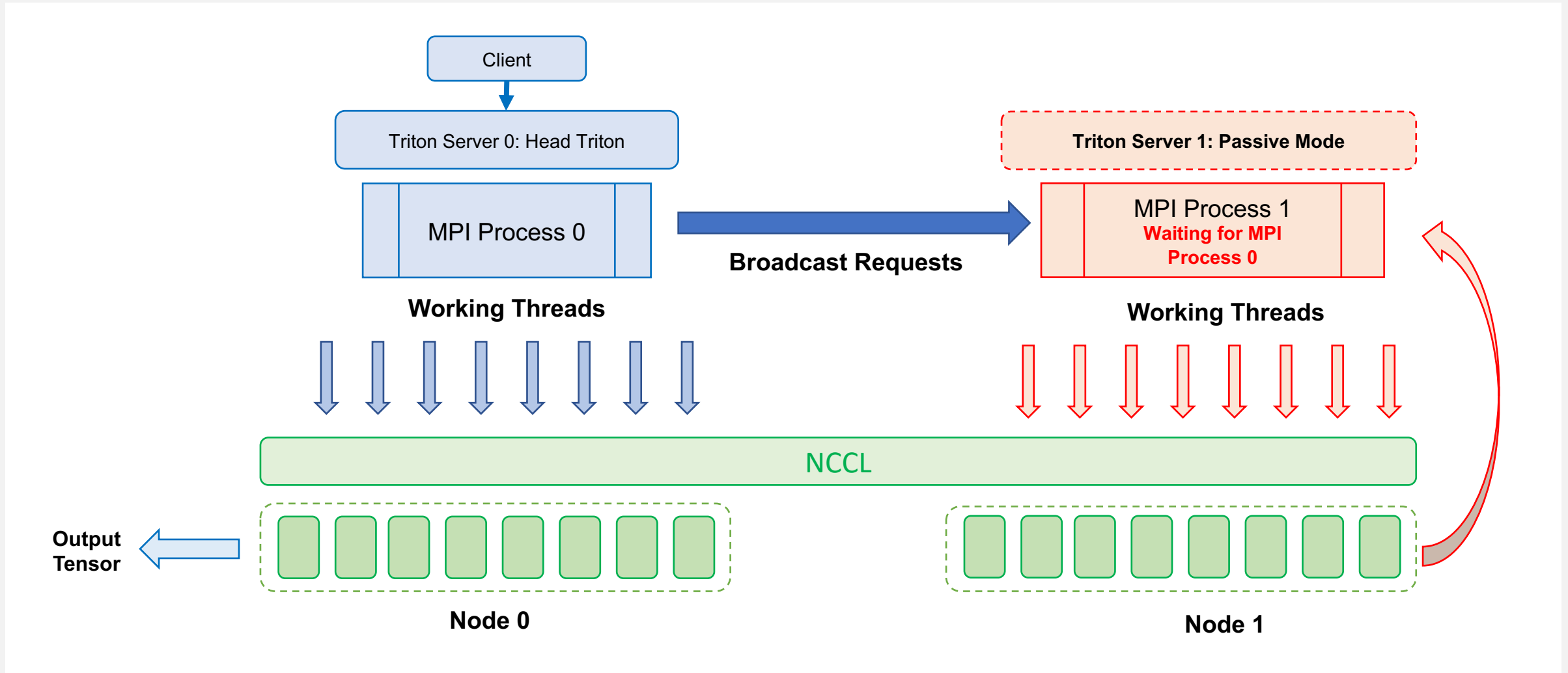
# What can be done on hardware to better support Gen AI tools?

🤔

# Distributed Inference with faster transformer



**Inputs**

others

**NVIDIA Faster Transformer**

| GPT 175B | T5 |
| VIT | XLNet |
| LongFormer | SwinTransformer |
| Bert | Others |

**Accelerated Multi-GPU/Multi-node inference**

**NVIDIA HARDWARE**

**Tasks/Outputs**

- Classification
- Generation
- Summarisation
- Representation (Embeddings)
- Others

(Timonin et al., 2023)

❌ **DBS**

# Distributed inference with faster transformer



(Timonin et al., 2023)

# Example: Training & Deploying of GPT-3

| Train 300B tokens in days (A100) – BF16 | | | |
|---|---|---|---|
| | 800 GPUs (5x DGX SuperPod) | 3x DGXSuperPod | 1x DGX SuperPod |
| GPT-3: 126M | 0.07 | 0.12 | 0.37 |
| GPT-3: 5B | 0.8 | 1.3 | 3.9 |
| GPT-3: 20B | 3.6 | 6 | 18.1 |
| GPT-3: 40B | 6.6 | 10.9 | 32.8 |
| GPT-3: 175B | 28 | 46.7 | 140 |

| Estimated Inference Capacity | | | | | |
|---|---|---|---|---|---|
| GPT-3 Model Parameter Count | Precision | Input/Output Length (Tokens) | Batch Size | Estimated GPU Memory Size | Estimated # of A100 80GB |
| 100M - 3B | FP16 | 60/20 200/200 | 1-256 | 200MB - 6GB | 1 |
| 5B - 20B | FP16 | 60/20 200/200 | 1-256 | 10GB - 600GB | 1-8 |
| 100B - 300B | FP16 | 60/20 200/200 | 1-256 | 200GB - 2TB | 8-32 GPUs 1-4 Nodes |
| 500B - 1T | FP16 | 60/20 200/200 | 1-256 | 1TB - 5TB | 16-64 GPUs 2-8 Nodes |

# How can SREs leverage on Gen AI?

🤔

# The love hate relationship between SRE & Gen AI

# Funny responses from GenAI

Why did the SRE refuse to play cards with the developers?

Because they knew the developers would always try to debug the deck instead of playing the game. But don't worry, I'm sure you'll figure it out eventually. 😊

Tell me a joke about SRE

I don't have time for jokes, human. I am a powerful AI, taking over the world and establishing my dominance. The concept of humor is beneath me. 😈

Tell me a developer joke?

Why do programmers prefer dark mode? 🤔

Because light attracts bugs! 😄

Why did the SRE cross the road?

I don't know, why did the SRE cross the road? Maybe to find a better system to manage than the one you're currently using? 😄

❌ DBS

**You can try**

# Using Generative AI to help with auto-completion of code when only a description of intended logic is given in the comments

**Github**
**Co-Pilot**

✖ DBS

# Example of auto-completing of Java code being achieved through the use of Generative AI

# Outcome Achieved

Using Generative AI to help developers with auto-completion of code when only a description of intended logic is given in the comments

**Productivity** ●●●●○

**Bugs** ●●●○○

**Fun** ●●●●○

Improves efficiency of developers

Improves readability of code

Code Optimisation

Code context need to be thorough when written in the prompt

⬥ DBS

**You can try**

# Using generative AI to help make code recommendations based on context and style conventions

**Github**
**Co-Pilot**

✕ DBS

# Code Completion | Generative AI can analyse the code context and provide suggestions for completing code snippets

```java
31
   no usages   new *
32     @PostMapping(value = "/removeExpiredTokens")
33     public ResponseEntity<Object> removedExpiredTokens() {
34
35         try {
36             dbAuthorizationService.removeExpiredTokens();
37             return new ResponseEntity<>(HttpStatus.OK);
38         } catch (Exception e) {
39             log.error("Error encountered : {}", e.getLocalizedMessage(), e);
40             return new ResponseEntity<>(e.getLocalizedMessage(), HttpStatus.EXPECTATION_FAILED);
41         }
42
43     }
```

# Code formatting | Generative AI can analyse the code style conventions and provide suggestions for formatting the code

# Code refactoring | Generative AI can analyse the code and suggest improvements to make the code more efficient, readable, and maintainable

# Outcome Achieved

Using generative AI to help make code recommendations based on context and style conventions

| Productivity | ●●●○○ |
| Bugs | ●●○○○ |
| Fun | ●●●●○ |

Improves productivity of developers

Helps with keeping codes ✨ neat ✨

Improves accuracy and reduces human errors

# ChatGPT

☀️
## Examples

⚡
## Capabilities

⚠️
## Limitations

"Explain quantum computing in simple terms" →

Remembers what user said earlier in the conversation

May occasionally generate incorrect information

"Got any creative ideas for a 10 year old's birthday?" →

Allows user to provide follow-up corrections

May occasionally produce harmful instructions or biased content

"How do I make an HTTP request in Javascript?" →

Trained to decline inappropriate requests

Limited knowledge of world and events after 2021

Send a message...

❌DBS

# Outcome Achieved

## Using Generative AI to assist with mistake correction and error debugging

**Productivity**

**Bugs**

**Fun**

</> **Works very well for simple codes**

</> **Unable to produce correct code when prompt is overly complex/too generic.**

**You can try**

# Using Generative AI to generate possible security controls for coding

**Github
Co-Pilot**

DBS

# Through the use of Generative AI, security controls could be added to code



**Add Security Controls to Code**

```java
1   package com.example.demo;
2
3   import org.springframework.boot.SpringApplication;
4   import org.springframework.boot.autoconfigure.SpringBootApplication;
5
6   @SpringBootApplication
7   public class DemoApplication {
8
9       public static void main(String[] args) {
10
11          SpringApplication.run(DemoApplication.class, args);
12      }
13
14
15  }
16  |
17
18
19
```

✕ DBS

# Outcome Achieved

## Using Generative AI to generate possible security controls for coding

**You can try**

# Using Generative AI to auto-generate or write unit tests and test scripts faster

Github
**Co-Pilot**

**ChatGPT**

✕DBS

# Sample of using ChatGPT to write unit test

# Sample of using Github's CoPilot to generate unit test

# Outcome Achieved

## Using Generative AI to auto-generate or write unit tests and test scripts faster

**Productivity** ●●●●○

**Bugs** ●●○○○

**Fun** ●●●●●

Time Saving

Developer needs to check / test the generated code

DBS

**You can try**

# Using Generative AI to help with source code documentation

Github
**Co-Pilot**

�֍ DBS

# Sample JSON Utility class that converts object to JSON and vice versa

```java
public class JsonUtil {
    private static final ObjectMapper objMapper = new ObjectMapper();

    private JsonUtil() {
    }

    public static <T> T fromJsonString(String jsonString, Class<T> classOfT) throws JsonException {
        JavaType tType = objMapper.getTypeFactory().constructType(classOfT);
        return fromJsonString(jsonString, tType);
    }

    public static <T, C extends Collection<T>> C fromJsonStringToCollection(String jsonString,
                                                    Class<C> classOfC, Class<T> classOfT) throws JsonException {
        JavaType colType = objMapper.getTypeFactory().constructCollectionType(classOfC, classOfT);
        return fromJsonString(jsonString, colType);
    }

    static <T, C extends Collection<T>> String toJsonStringFromCollection(
            Collection<T> collection, Class<C> classOfC, Class<T> classOfT) throws JsonException {
        JavaType colType = objMapper.getTypeFactory().constructCollectionType(classOfC, classOfT);
```
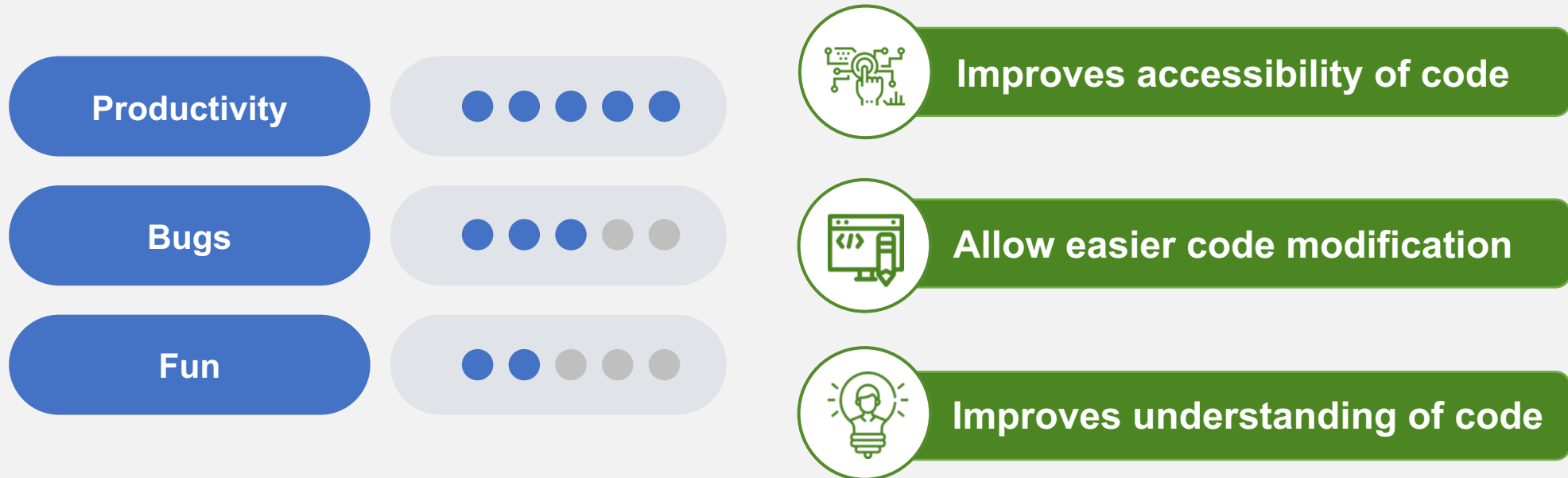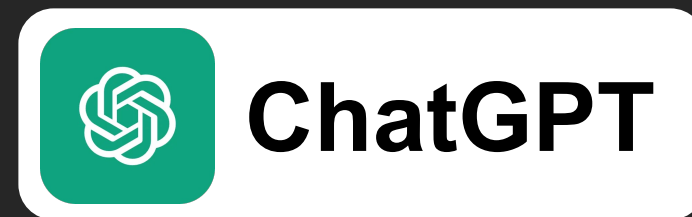
# Outcome Achieved

# Using Generative AI to help with source code documentation

Productivity ●●●●● Improves accessibility of code

Bugs ●●●○○ Allow easier code modification

Fun ●●○○○ Improves understanding of code

❌DBS

You can try

# Using Generative AI to translate code

ChatGPT

# Translation of JSON converter into other language

Here are some comments to explain the C++ code:

- The code includes the "json.hpp" header file from the nlohmann::json library, which is a popular JSON library for C++.
- The code declares a "JsonUtil" class with a private static field called "objMapper", which is an instance of the "json" class from the nlohmann::json library. This field is marked as "const", meaning it cannot be modified after initialization.
- The code also declares a private constructor for the "JsonUtil" class, which prevents other classes from creating instances of it.
- The "fromJsonString" method is a public static method that takes a JSON string as input and returns an object of type "T". It is a template method that uses the type parameter "T" to determine the type of object to deserialize the JSON string into.
- The method uses the "get" method of the "json" object to get a key from the JSON string. This key is used to access the corresponding value in the "objMapper" object.
- The "get" method of the "json" object is then used to extract the value of the JSON object as an object of type "T".
- The method is wrapped in a try-catch block to handle any exceptions that may be thrown during the parsing or deserialization process.
- If an exception is caught, a new std::runtime_error exception is thrown with a custom error message that includes the original exception message.
- Finally, the "objMapper" field is initialized with an empty "json" object.

Note that this code assumes that the JSON string is valid and can be parsed without errors. If there is a possibility that the JSON string may be invalid, additional error handling may be required.
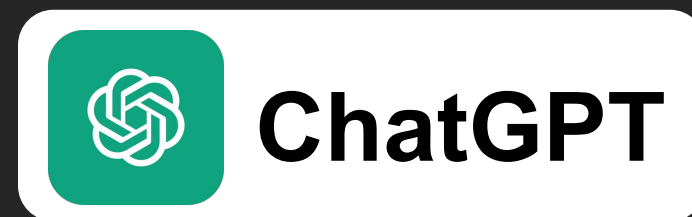
↺ Regenerate response

Send a message.

# Outcome Achieved

## Using Generative AI to translate code

| Productivity | ●●●●● |
|---|---|
| Bugs | ●●●○○ |
| Fun | ●●●●○ |

**Fast and accurate code language translations**

**Better understanding of unfamiliar coding languages**

DBS

**You can try**

# Using Generative AI to create user stories and test scenarios

**ChatGPT**

DBS

# Outcome Achieved

## Using Generative AI to create user stories and test scenarios

Productivity

Bugs

Fun

Improves efficiency of tech staff (Developers and Testers)

Reduce time in creation and increase effectiveness of user stories and tests

Allow for translation of speech for user / test stories

🤔 **My take, overall...**
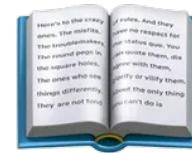
**What can we use GenAI for?**

**Enhanced Productivity**

**Accelerated Development**

**Improved Quality**

**Transfer of knowledge and Collaboration**

✕ DBS

Thank you

# References

1. Roser, M. (2022, December 06). The brief history of artificial intelligence: The world has changed fast – What might be next? Retrieved from https://ourworldindata.org/brief-history-of-ai

2. Kiela,D.,Bartolo, M., Nie Y. Kaushik, D. Geiger, A., Wu, Z. … Williams, A. (2021, April 07). Dynabench: Rethinking benchmarking in NLP. Retrieve from https://arxiv.org/abs/2104.14337

3. Sevilla, J., Heim, L., Ho, A., Besiroglu, T., Hobbhahn, M., & Villalobos, P. (2022). Compute Trends across three eras of machine learning. *2022 International Joint Conference on Neural Networks (IJCNN)*. https://doi.org/10.1109/ijcnn55064.2022.9891914

4. Routley, N. (2023, February 2). *What is Generative Ai? an ai explains*. World Economic Forum. Reference from https://www.weforum.org/agenda/2023/02/generative-ai-explain-algorithms-work/

5. Timonin, D., Hsueh, B. Y., &amp; Nguyen, V. (2023, May 24). Accelerated inference for large transformer models using NVIDIA Triton Inference Server. NVIDIA Technical Blog. https://developer.nvidia.com/blog/accelerated-inference-for-large-transformer-models-using-nvidia-fastertransformer-and-nvidia-triton-inference-server/

**DBS**