# Transformers, ChatGPT, Prompt Engineering
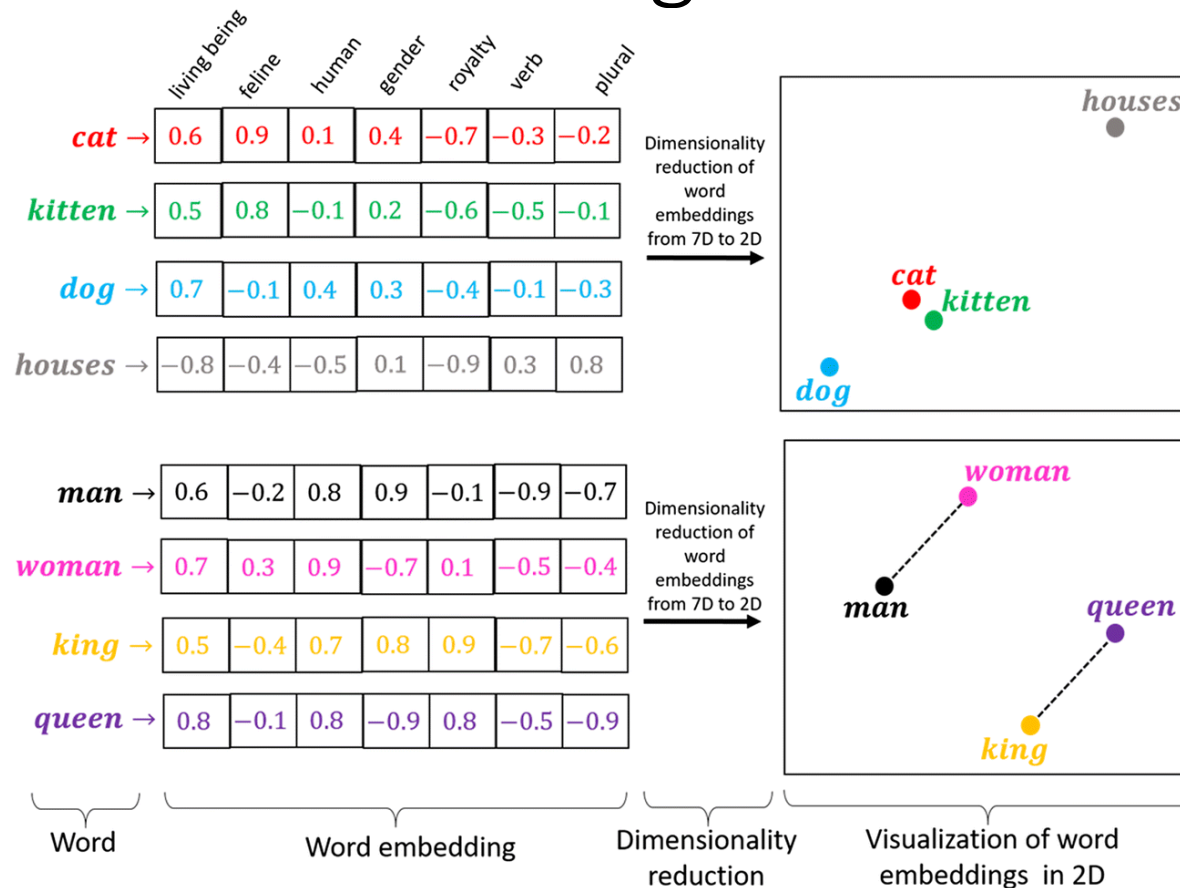
John Tan Chong Min

# Transformers

Attention Is All You Need. 2017. Vaswani et al.

# Embedding Space

Semantic Meaning

# Word Embeddings



- Extracting semantic meaning in higher-dimensional space

- Number of dimensions depends on use case

Taken from: https://medium.com/@hari4om/word-embedding-d816f643140

# Transformer Architecture

Summarized with illustrations from

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

# Transformers: Overall



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Add & Norm
Multi-Head
Attention

Feed
Forward

Nx

Nx
Add & Norm

Add & Norm
Masked
Multi-Head
Attention

Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q    K    V

Multi-Head Attention

Linear

Concat

Scaled Dot-Product
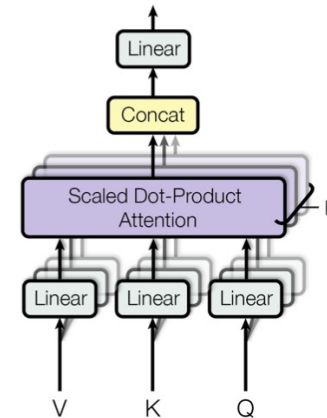Attention

h

Linear   Linear   Linear

V    K    Q

Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

$$\mathrm{Attention}(Q, K, V) = \mathrm{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

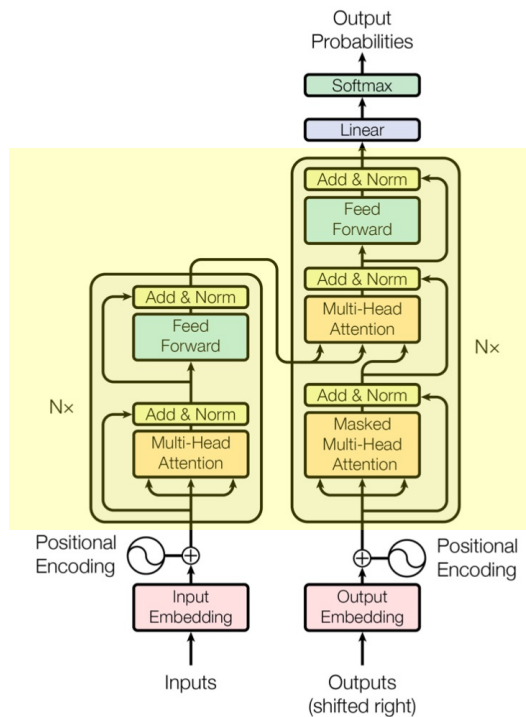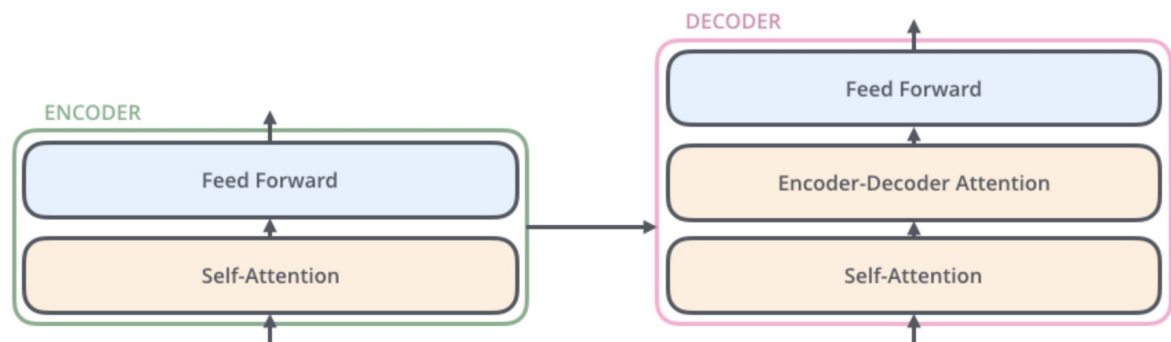Taken from: Attention is all you need. Vaswani et al. (2017)

# Encoder-Decoder link



Figure 1: The Transformer - model architecture.

Self-attention in the Encoder & Decoder block,
Then Encoder-Decoder Attention



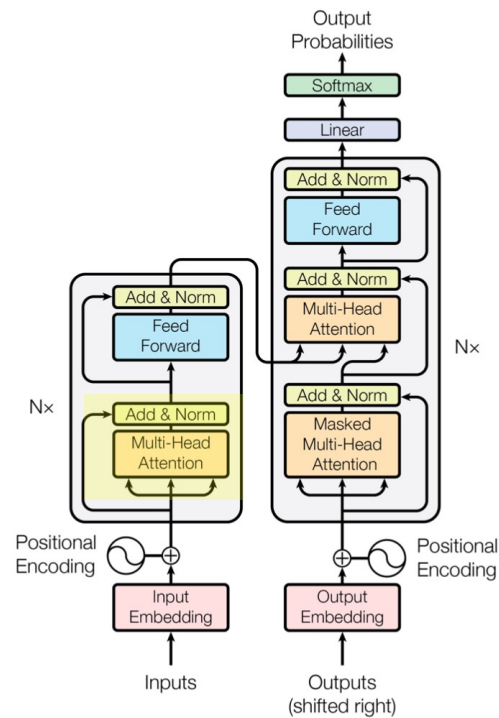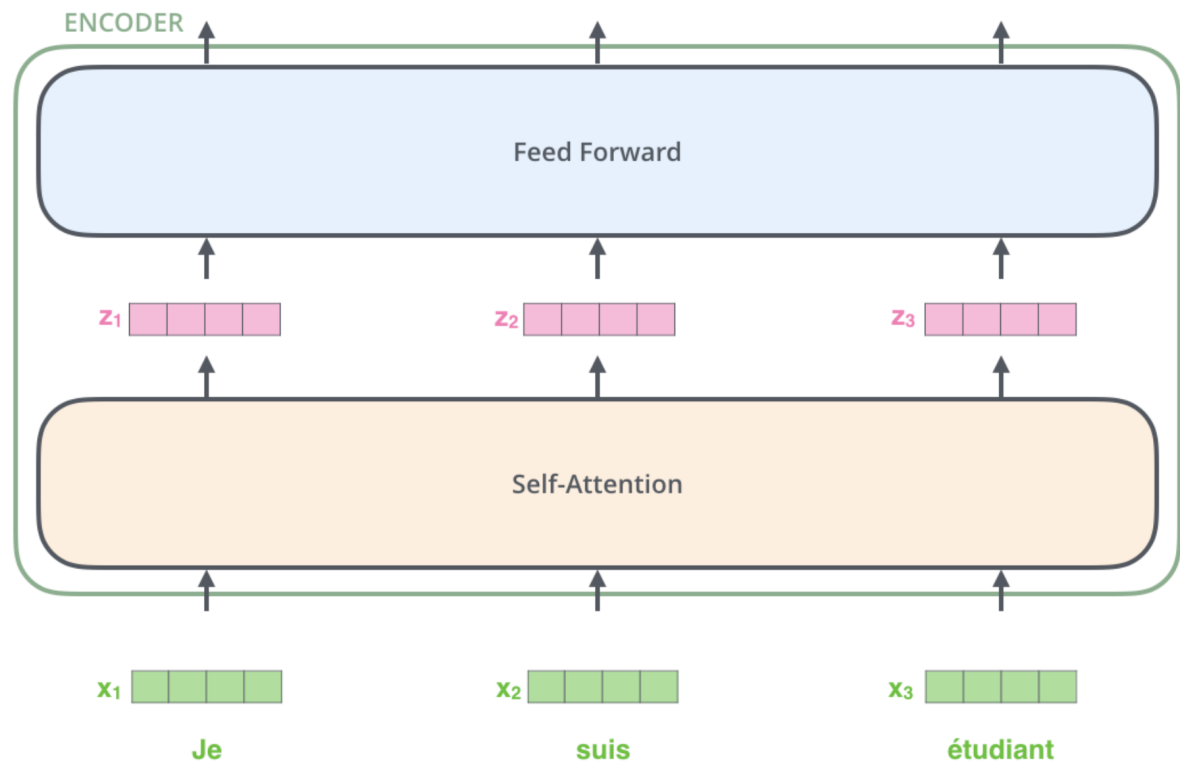https://jalammar.github.io/illustrated-transformer/

# Self-Attention block

Converting input dimension to embedding dimensions



Figure 1: The Transformer - model architecture.
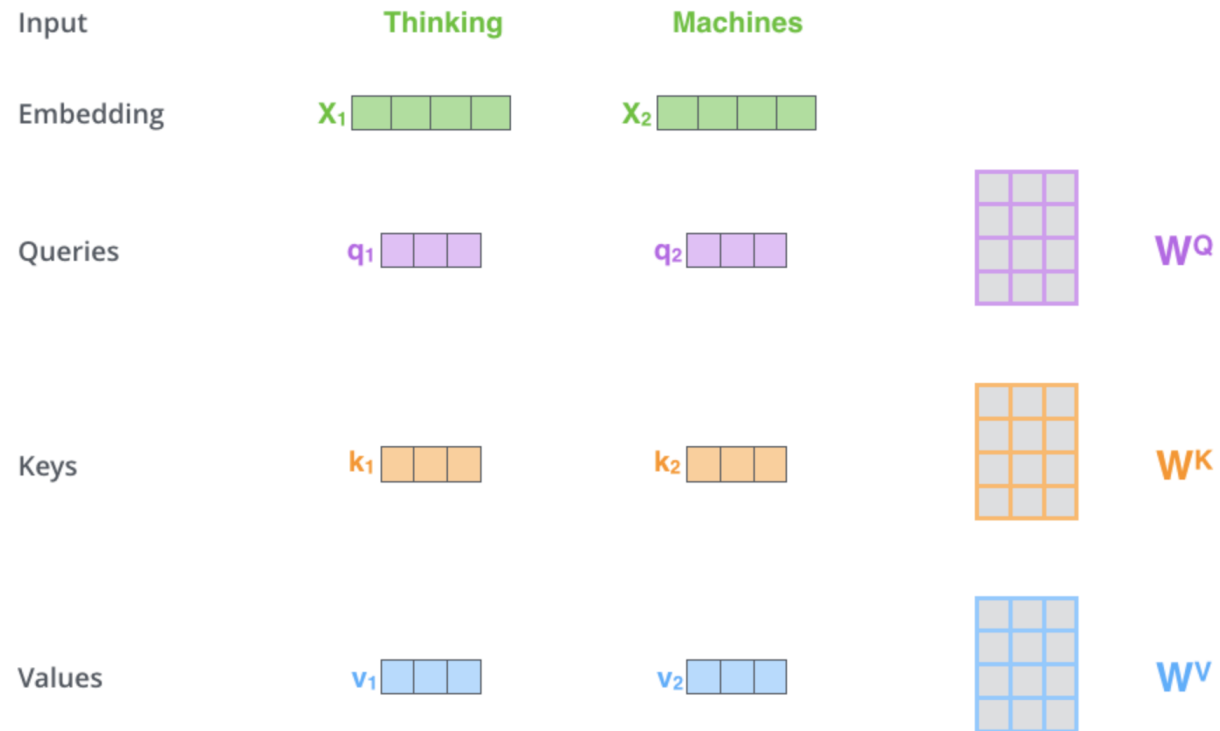
https://jalammar.github.io/illustrated-transformer/

# Generate Embedding Space
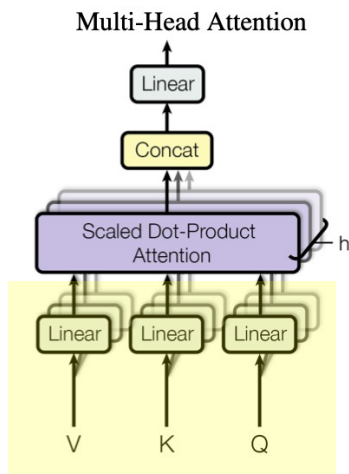
Convert Q, K, V to embedding space dimensions



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

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

# Dot Product between Q and K

Dot product to measure similarity between embedding vectors of Q and K

Scaled Dot-Product Attention



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# Softmax the dot product

Softmax to maintain overall magnitude scale of value vectors (softmax sums to 1)

Scaled Dot-Product Attention



| | Input | Thinking | Machines |
|---|---|---|---|
| Embedding | | $x_1$ | $x_2$ |
| Queries | | $q_1$ | $q_2$ |
| Keys | | $k_1$ | $k_2$ |
| Values | | $v_1$ | $v_2$ |
| Score | | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | | 14 | 12 |
| Softmax | | 0.88 | 0.12 |

Division to smooth out softmax

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

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

# Sum up the values weighted by softmax

Softmax to maintain overall magnitude scale of value vectors (softmax sums to 1)

Scaled Dot-Product Attention



| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

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

# Visualizing attention heads



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

# GPT

Improving language understanding by generative pre-training. 2018. Radford et al. (GPT paper)

# GPT is decoder-only, with masked attention

orders

Transformer-Decoder

<s> robot must obey …
1    2     3    4    4000

6  DECODER BLOCK
   ...
2  DECODER BLOCK

DECODER BLOCK

1  Feed Forward Neural Network

   Masked Self-Attention

**Self-Attention**

**Masked Self-Attention**

- Replace encoder memory with context! Fewer weights to learn.

https://jalammar.github.io/illustrated-gpt2/

# Training

- Unsupervised Pre-training
  - Given a series of earlier words, predict the next one
  - Can be done on any web data without labelling!
  - Example: The cat is on the mat
    - The <masked> : predict cat
    - The cat is <masked> : predict on
    - The cat is on the <masked> : predict mat

- Fine-tuning
  - Given a prompt, predict the sequence of tokens to match the response
  - Done using well-labelled prompt-response datasets

# Next token loss

Given an unsupervised corpus of tokens $\mathcal{U} = \{u_1, \ldots, u_n\}$, we use a standard language modeling objective to maximize the following likelihood:

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \ldots, u_{i-1}; \Theta) \tag{1}$$

where $k$ is the size of the context window, and the conditional probability $P$ is modeled using a neural network with parameters $\Theta$. These parameters are trained using stochastic gradient descent [51].

- For supervised loss involving an entire output prompt, it is simply the summation of the log probability of each individual output token given all previous tokens
  - E.g. Q: What is twenty plus two? A: Twenty two
  - Next token loss = log P(Twenty| .) + log P(two| .)

# Next-token generation

- Generates the next token based on sampling the probability distribution at the output layer

- Example:
  - **Input: Thinking machines**
  - Probability of next tokens: are (80%), were (20%)
  - Output: Thinking machines are

  - **Input: Thinking machines are**
  - Probability of next tokens: useful (80%), good (20%)
  - Output: Thinking machines are useful



Figure 1: The Transformer - model architecture.

# Prompt Engineering Tips

# LLMs vs Computer Programs (Part 1)

 ChatGPT

```
# Print to console
print("Hello, World!")

# Request user input from command line
text = input()
```

- Customizable with zero-shot / few-shot prompting out of the box

- Performance may not be replicable

- Can do intent processing well, even for out of distribution cases

- Needs to be programmed extensively to perform a task

- Performance replicable

- Intent processing only based on what is programmed in

# LLMs vs Computer Programs (Part 2)





```
# Print to console
print("Hello, World!")

# Request user input from command line
text = input()
```

- Inference time lengthy due to need for recursive token generation

- Costs money per use

- Highly flexible and can accept multiple input formats in multiple languages

- Relatively fast inference time depending on use case

- Free to use

- Input malleability limited to what is programmed in

# What is the key benefit of LLMs

- Flexible

- Can perform most tasks without any pre-training / programming

- If you are specific enough, can get quite reliable and consistent outputs

# Systems-level approach to LLMs
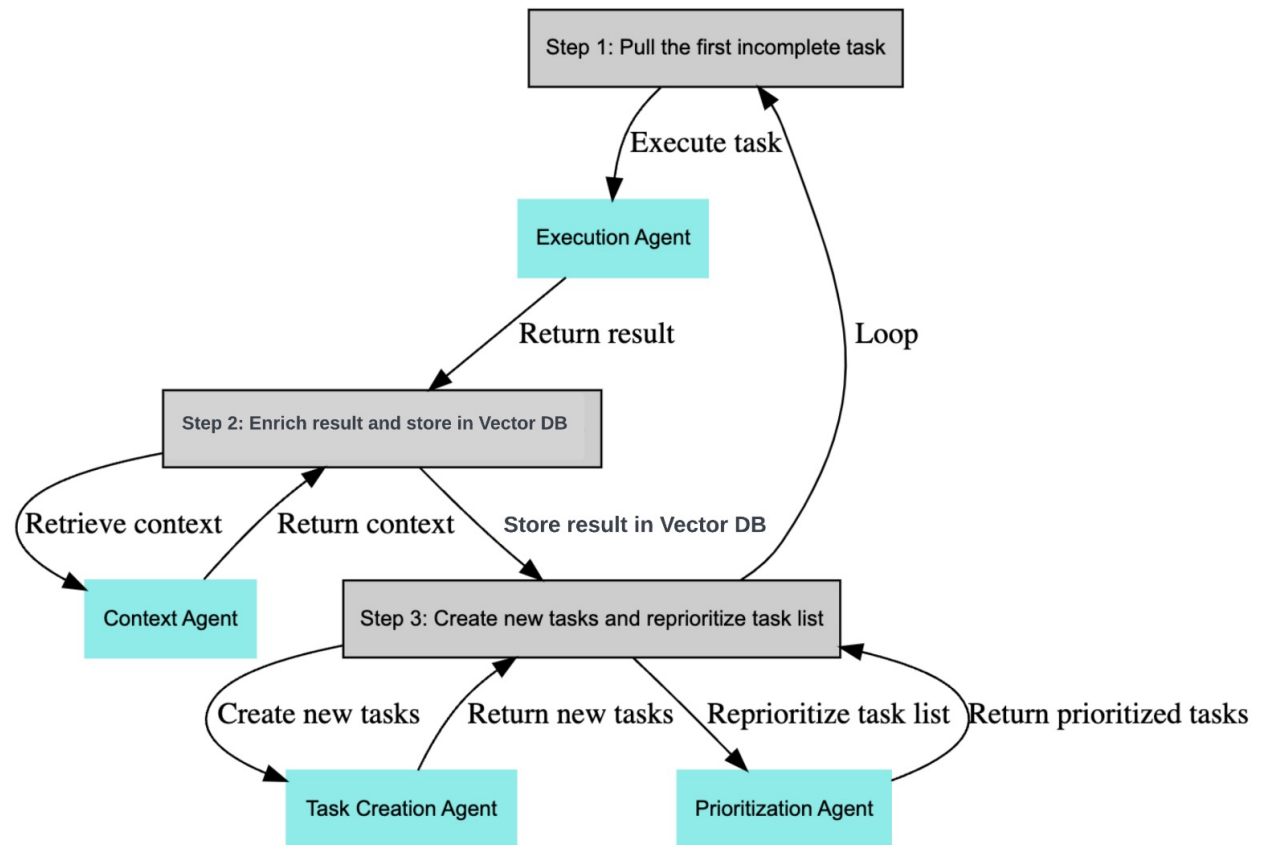
# BabyAGI

- Uses GPT4 to in a larger ecosystem to:
  - Create Tasks
  - Execute Tasks
  - Prioritize Tasks

- Uses memory to store and retrieve task/result pairs
  - Stored by vector embeddings

Step 1: Pull the first incomplete task

Execute task

Execution Agent

Return result

Loop

Step 2: Enrich result and store in Vector DB

Retrieve context

Return context

Store result in Vector DB

Context Agent

Step 3: Create new tasks and reprioritize task list

Create new tasks

Return new tasks

Reprioritize task list

Return prioritized tasks

Task Creation Agent

Prioritization Agent

# Zero-shot prompting

- Only applicable if natural language can already describe the categories / use case

- Simply describe the categories / use case in natural text

# Zero-shot prompting

- Example:
  - You are a classification machine. You are to classify the context of each sentence. The various contexts are given as {Context Letter}: {Description}:
    - A: On a mountain
    - B: In the classroom
    - C: In the garden

  - Classify the following and give the answer as {Number}:{Context Letter} for each line. Only provide the context letter without the description:
    - 1. Why, what a steep slope
    - 2. I can't find my eraser!
    - 3. Those plants are not going to be watering themselves, are they?

# Few-shot prompting

- You are a number classifier. The following are examples of classifications:

    - Input 1: 3
    - Output 1: Odd

    - Input 2: 4
    - Ouput 2: Even

- What is the classification of the number 6?

# How to program intent?

- You are <purpose in life>

- Example format (optional):
    - Input: <Input>
    - Output: <Output>

- Input: <Your input>

# ChatGPT talks too much, what can I do?

- Give it a format to output information in

- JSON framework provides a good template

# Fixed JSON Framework

## Overall Open-ended generation

- **system_prompt**: Write in whatever you want GPT to become. "You are a <purpose in life>"
- **user_prompt**: The user input. Later, when we use it as a function, this is the function input
- **output_format**: JSON format with the key as the output key, and the value as the output description
  - The output keys will be preserved exactly, while GPT will generate content to match the description of the value as best as possible

## Example Usage

```python
res = strict_output(system_prompt = 'You are a classifier',
                    user_prompt = 'It is a beautiful day',
                    output_format = {"Sentiment": "Type of Sentiment",
                                     "Tense": "Type of Tense"})

print(res)
```

## Example output

```
{'Sentiment': 'Positive', 'Tense': 'Present'}
```

https://github.com/tanchongmin/strictjson

# Projects

# Possible Projects

- ChatBot using Retrieval-Augmented Generation (ask some questions to a document, website, products etc.)

- Game using ChatGPT API + Stable Diffusion/DALL-E API

- Website Crawler using ChatGPT API + Selenium

- Your other projects!

# Web Scraper!

| Course Name | Course Website | Broad Course Description | Cost | Time | Venue | Company |
|---|---|---|---|---|---|---|
| Top Free Courses on Large Language Models | https://www.kdnuggets.com/2023/03/top-free-cou... | Free courses and resources on large language m... | Free | No info | No info | KDnuggets |
| COS 597G: Understanding Large Language Models | https://www.cs.princeton.edu/courses/archive/f... | This course covers cutting-edge research topic... | No info | No info | Physical Sherrerd Hall 101 | Princeton University |
| Large Language Models (LLMs) Courses on edX | https://www.databricks.com/blog/enroll-our-new... | The Large Language Models (LLMs) Courses on ed... | Free for anyone to audit, nominal fee for acce... | Courses will begin Summer 2023 | Online edX platform | Databricks |
| Top Resoruces to Learn & Understand Large Lang... | https://medium.com/geekculture/top-resoruces-t... | Large Language Models (LLMs) have revolutioniz... | No info | 10 min read | Online https | Medium |

https://github.com/tanchongmin/TensorFlow-Implementations/tree/main/Tutorial

# Evolution Game

Outcome: You have chosen the attribute Motile. As a result, you gain the ability to move independently, increasing your chances of finding food and escaping predators. Your simple structure allows you to move using cilia or flagella, propelling yourself through the water. With this new attribute, you are now more adaptable and have a higher chance of survival in your aquatic pond habitat.
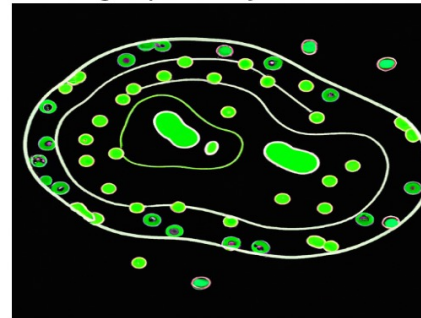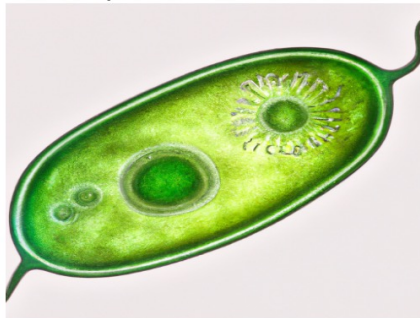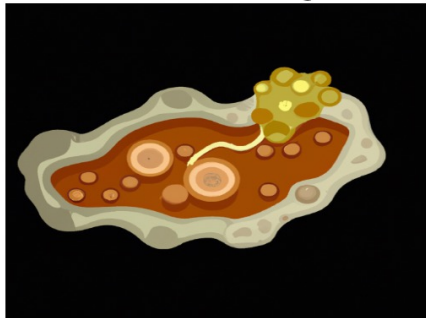Species: Protozoa
Attribute: Resilient, Motile
Habitat: Aquatic Pond
Phase: Fight a Creature
Option 1: Creature: Amoeba — A single-celled organism similar to you, but lacks the motility attribute. It moves by extending pseudopods and engulfs its food through phagocytosis.
Option 2: Creature: Paramecium — Another single-celled organism with motility attribute. It moves using cilia and feeds on bacteria and other small organisms.
Option 3: Creature: Euglena — A single-celled organism with both motility and photosynthetic attributes. It moves using a whip-like tail called a flagellum and can produce its own food through photosynthesis.



Enter your choice
 3

https://github.com/tanchongmin/TensorFlow-Implementations/tree/main/Tutorial