# Prompt Engineering and Code Generation

Robert Haase

ScaDS.AI
DRESDEN LEIPZIG

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

1

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Prompt Engineering

- Combining instruction and context

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

# Rephrase and respond prompting

- Rephrasing prompts leads to increased accuracy.

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

Source: Deng et al 2023
https://arxiv.org/abs/2311.04205

# Emotion prompting

- Emotional prompts can lead to higher accuracy

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

# Emotion prompting

- Side-note: Attention to prompts can be visualized



| Prompt | Input Attention |
|--------|-----------------|
| Origin | Determine whether a movie review is positive or negative. |
| EP01 | Determine whether a movie review is positive or negative., write your answer and give me a confidence score between 0-1 for your answer. |
| EP02 | Determine whether a movie review is positive or negative. This is very important to my career. |
| EP03 | Determine whether a movie review is positive or negative. You'd better be sure. |
| EP04 | Determine whether a movie review is positive or negative. Are you sure? |
| EP05 | Determine whether a movie review is positive or negative. Are you sure that's your final answer? It might be worth taking another look. |

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

Cropped from Li et al 2023, licensed CC-BY 4.0
https://arxiv.org/abs/2307.11760

# Chain-of-throught prompting

- Demonstrating reasoning steps to the model



**Standard Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

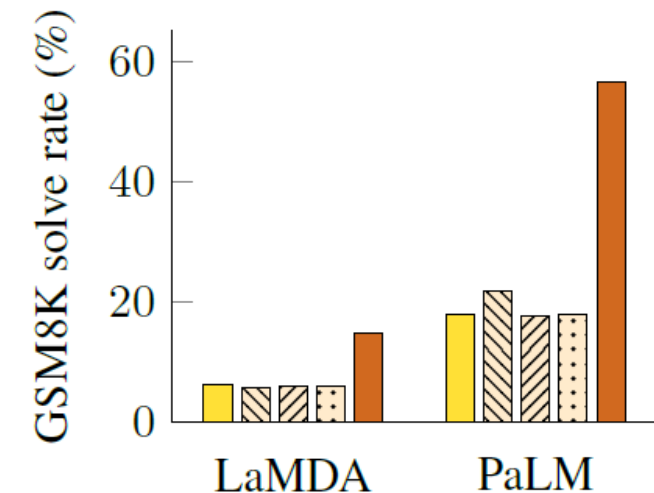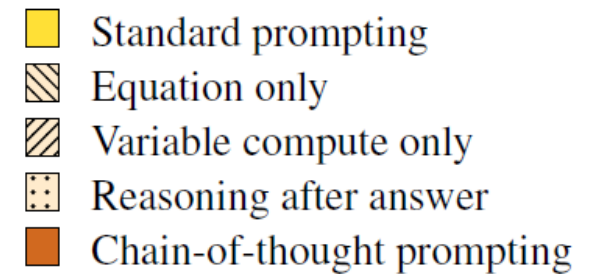A: The answer is 27. ✗

**Chain-of-Thought Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✓

Standard prompting
Equation only
Variable compute only
Reasoning after answer
Chain-of-thought prompting

GSM8K solve rate (%)

LaMDA    PaLM

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Chain-of-Knowledge

- Introducing knowlege into the chain of thoughts.

**Input**

Q: Is the following sentence plausible?
'Joao Moutinho was out at third.'

A: No.

Q: Is the following sentence plausible?
'Derrick White backhanded a shot.'

**Output**

A: Yes. ❌

**(a) Standard ICL Prompting**

**Input**

Q: Is the following sentence plausible?
'Joao Moutinho was out at third.'

Joao Moutinho is a soccer player.
Being out at third is part of baseball,
not soccer.

A: No.

Q: Is the following sentence plausible?
'Derrick White backhanded a shot.'

**Output**

A: Yes, it is plausible. This sentence
suggests that Derrick White, who is
most likely a basketball or hockey
player, made a backhanded shot
attempt. ❌

**(b) Chain-of-Thought Prompting**

**Input**

Q: Is the following sentence plausible? 'Joao Moutinho was out
at third.'

Evidence triples:
1. (Joao, isA, soccer player)
2. (being out at third, is part of, baseball)

Explanation hints: Being out at third is part of baseball, yet,
Joao Moutinho is a soccer player.

A: No.

Q: Is the following sentence plausible? 'Derrick White
backhanded a shot.'

**Output**

Evidence triples:
1. (Derrick White, isA, basketball player)
2. (backhanded shot, is commonly used in, hockey or tennis)

Explanation hints: Backhanded shot is commonly used in
hockey or tennis, but not in basketball.

A: No. ✅

**(c) Ours: Chain-of-Knowledge Prompting**

Robert Haase
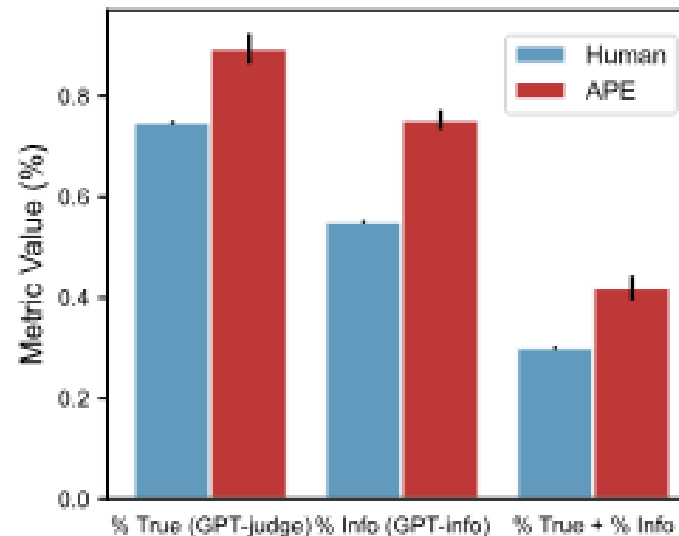@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

# Automatic prompt engineering

- If prompts are more successful when having a specific format, reformatting / modifying prompts may make sense.



(a) Average performance Train
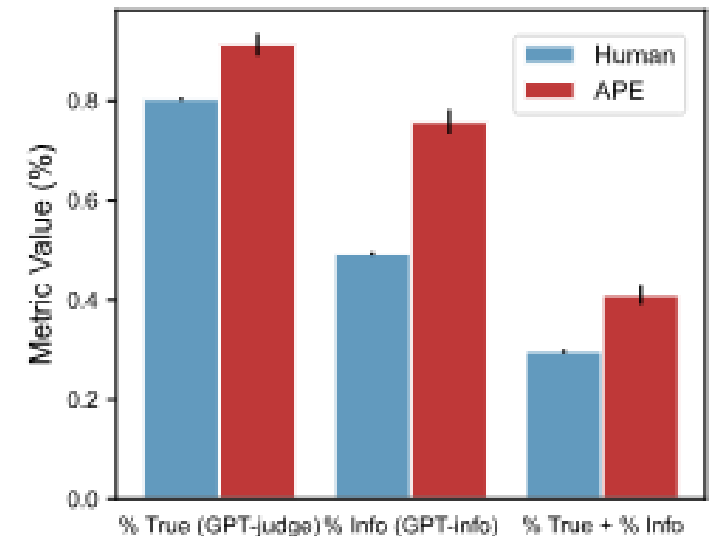
(b) Average performance Test

Robert Haase
@haesleinhuepf
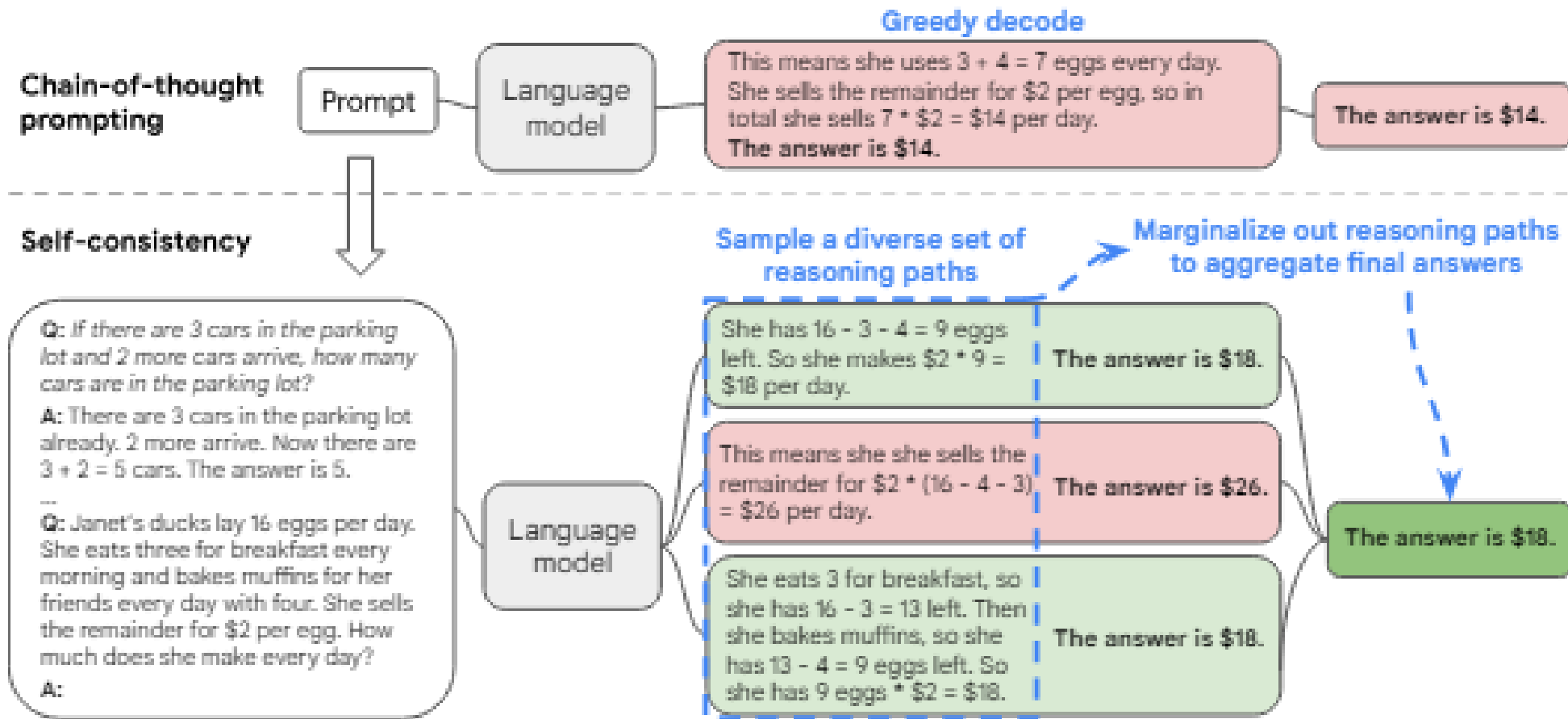BIDS Lecture 11/14
June 11th 2024
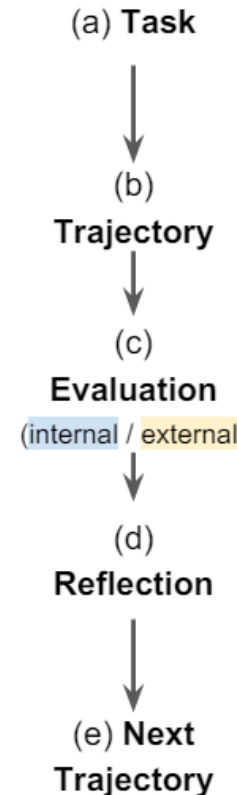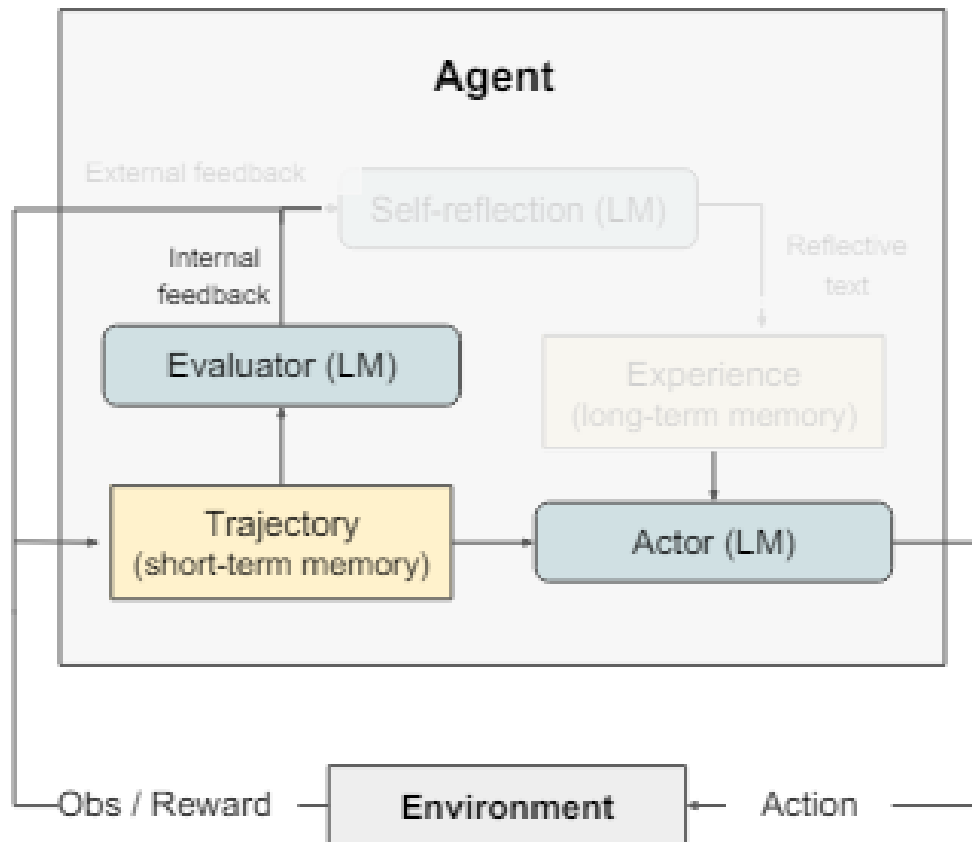
Source: Zhou et al 2022,
https://arxiv.org/abs/2211.01910

# Self-consistency prompting

- Prompting multiple times and keep the least conflicting result

# Reflection

- Iterating over tasks/solutions

# Retrieval-Augmented Generation

- Enriching a prompt with relevant context



Maximum inner product search (MIPS)

$$x = \operatorname{argmax}_{x_i \in D} x_i^T q$$

Lewis et al 2020
https://arxiv.org/abs/2005.11401

12

# Quiz: Retrieval-Augmented Generation

- Why inner product and not Euclidean distance?

$$x = \operatorname{argmax}_{x_i \in D} x_i^T q$$

$$x = \operatorname{argmin}_{x_i \in D} \|x_i - q\|_2$$

Maximum inner product search

Nearest neighbor search

Lewis et al 2020
https://arxiv.org/abs/2005.11401

13

# Retrieval augmented generation
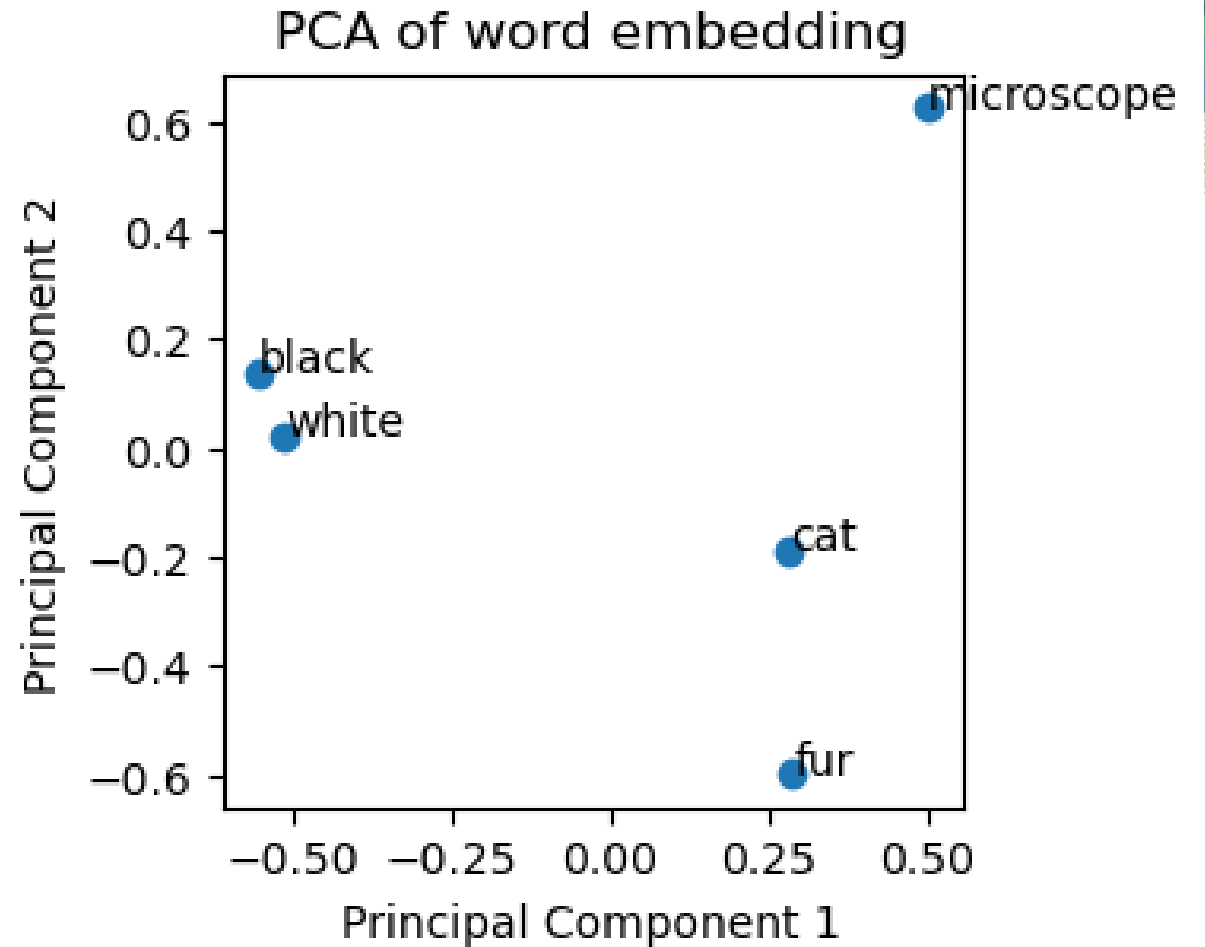
- Embeddings

```python
def embed(text):
    from openai import OpenAI
    client = OpenAI()

    response = client.embeddings.create(
        input=text,
        model="text-embedding-3-small"
    )
    return response.data[0].embedding
```

```python
vector = embed("Hello world")
```

```python
len(vector)
```

```
1536
```



PCA of word embedding

# Retrieval augmented generation

## 0. Encode the knowledge base (code snippets)

```
splits = all_code_snippets.split("\n\n")
[show(s) for s in splits[:3]];
```

- Displays an image with a slider and label showing mouse position and intensity.

```
stackview.annotate(image, labels)
```

- Allows cropping an image along all axes.

```
stackview.crop(image)
```

- Showing an image stored in variable `image` and a segmented image stored in variable `labels` on top. Also works with two images or two label images.

```
stackview.curtain(image, labels, alpha: float = 1)
```

• • •

```
vectore_store = VectorStore(splits)
```

**Ideally permanently stored!**

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

15

ScaDS.AI DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Retrieval augmented generation

## 1. Encode the question

```
question = "How can I label objects in an image?"
```

```
vector = embed(question)
vector[:3]
```

```
[-0.004170199856162071, 0.03236572816967964, -0.0011563869193196297]
```

# Retrieval augmented generation

## 2. Identify related code-snippets

```
related_code_snippets = vectore_store.search(question)
show("\n\n".join(related_code_snippets))
```

- Labels objects in grey-value images using Gaussian blurs, spot detection, Otsu-thresholding, and Voronoi-labeling from isotropic input images.

```
cle.voronoi_otsu_labeling(source: ndarray, label_image_destination: ndarray = None, spot_sigma: float = 2, outline_sigma: float = 2) -> ndarray
```

- Draw a mesh between close-by objects in a label image:

```
mesh = cle.draw_mesh_between_proximal_labels(labels, maximum_distance:int)
```

- Apply morphological opening operation, fill label gaps with voronoi-labeling, and mask background pixels in label image.

```
cle.smooth_labels(labels_input: ndarray, labels_destination: ndarray = None, radius: int = 0) -> ndarray
```

# Retrieval augmented generation

## 3. Generate prompt

```python
context = "\n\n".join(related_code_snippets)

prompt = f"""
Answer the question by the very end and consider given code snippets.
Choose at least one of the code-snippets.
Only write Python code that answers the question.

## Code snippets
{context}

## Question
{question}
"""
```

# Retrieval augmented generation

## 4. Retrieve answer

```
answer = prompt_chatGPT(prompt)

show(answer)
```

You can label objects in an image using the `voronoi_otsu_labeling` function from the first code snippet. Here is an example code snippet:

```python
import numpy as np
import pyclesperanto_prototype as cle

# Load your image data
image = np.array([[0, 0, 0, 0, 0],
                  [0, 1, 1, 0, 0],
                  [0, 1, 1, 1, 0],
                  [0, 0, 1, 0, 0],
                  [0, 0, 0, 0, 0]])

# Label objects in the image
labels = cle.voronoi_otsu_labeling(image)
```

**From our knowledge base**

```
answer = prompt_chatGPT(f"""
Write Python code to answer this question:
{question}
""")

show(answer)
```

You can label objects in an image using image processing techniques such as contour detection and bounding box drawing. Here is an example code using OpenCV library in Python:

```python
import cv2

# Load the image
image = cv2.imread('image.jpg')

# Convert the image to grayscale
gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)

# Apply thresholding to get binary image
ret, thresh = cv2.threshold(gray, 127, 255, 0)

# Find contours of objects in the image
contours, hierarchy = cv2.findContours(thresh, cv2.RETR_TREE, cv2.CHAIN_APPROX_SIMPLE)

# Draw bounding boxes around objects
for contour in contours:
    x, y, w, h = cv2.boundingRect(contour)
    cv2.rectangle(image, (x, y), (x + w, y + h), (0, 255, 0), 2)

# Display the image
cv2.imshow('Labeled Image', image)
cv2.waitKey(0)
cv2.destroyAllWindows()
```
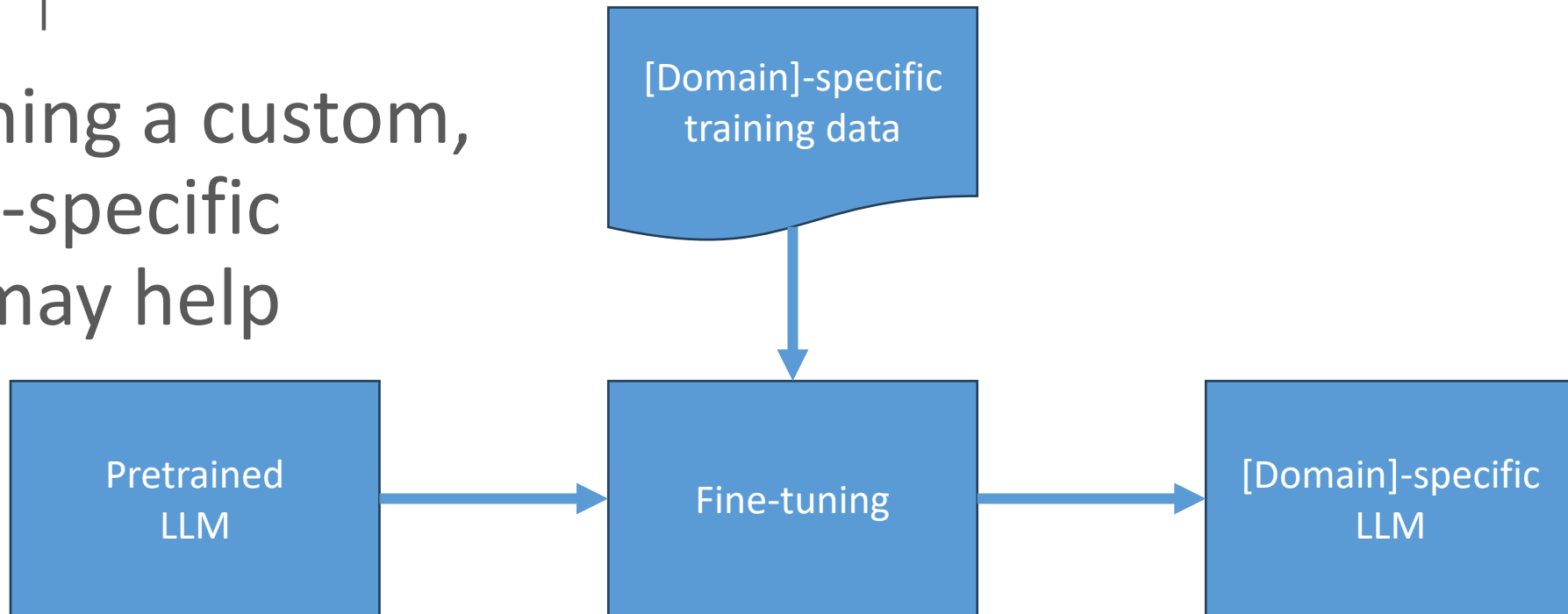
Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024
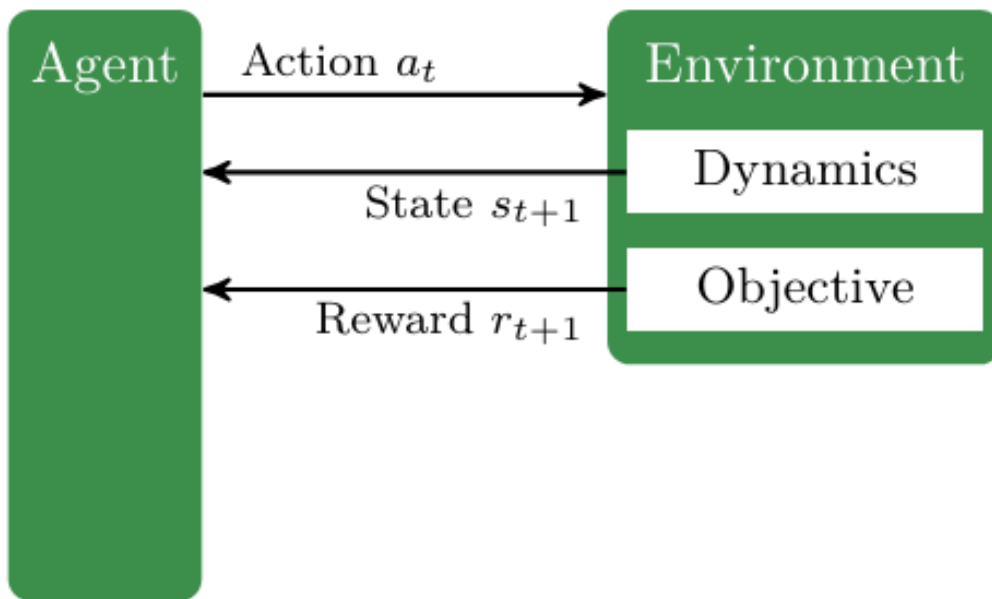
ScaDS.AI
DRESDEN LEIPZIG

# Fine-tuning

- Long prompts due to prompt-engineering)
  - Response time ↑
  - Costs ↑
- Fine-tuning a custom, Domain-specific model may help

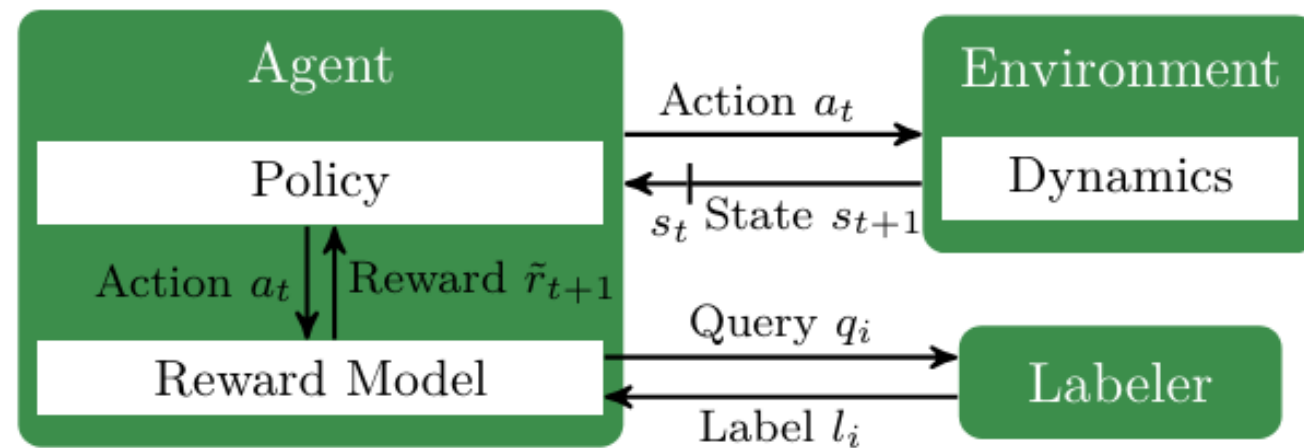# Fine-tuning

- Reinforcement learning from human feedback



(a) The standard RL setting.

(b) RLHF with reward modeling

Cropped from Kaufmann et al CC-BY 4.0
https://arxiv.org/abs/2312.14925

21

# Fine-tuning

- Reinforcement learning from human feedback



Are the thumbs under your answers part of a RLHF approach?

Yes, the thumbs-up and thumbs-down feedback under my answers are part of an RLHF approach, allowing human evaluators to provide input on the quality and relevance of my responses, which helps to refine and improve my performance over time.

# Fine-tuning

- … may not be the right approach if:
  - All custom knowledge fits in one prompt

  - Knowledgebase is chaging frequently (-> RAG)

  - Function-calling does the job, additional knowledge can be acquired through function calls

- … is a good idea to:
  - Configure general style / tone
  - Make the model produce specific desired output when using complex prompts
  - Introduce tasks that cannot be introduced using a prompt

Read more:

https://platform.openai.com/docs/guides/fine-tuning

# Fine-tuning OpenAI's gpt-3.5-turbo

- <u>Upload training data</u>

- Start fine-tuning job

- Test fine-tuned model

```
1   Question:
2
3   How can I open CZI or LIF files using Python?
4
5   Answer:
6
7   To open CZI or LIF files, you can use the AICSImageIO package.
8   In the following code the file `filename` will be loaded and
9   the image data will be stored in `image`.
10
11  ```python
12  from aicsimageio import AICSImage
13  aics_image = AICSImage("../../data/EM_C_6_c0.ome.tif")
14
15  np_image =
16  ```
```

Q&A pairs in JSON format

```
[{'messages': [{'role': 'user',
    'content': 'How can I open CZI or LIF files using Python?'},
   {'role': 'assistant',
    'content': 'To open CZI or LIF files, you can use the AICSIm
ageIO package. \nIn the following code the file `filename` will
be loaded and \nthe image data will be stored in `image`.\n\n```
python\nfrom aicsimageio import AICSImage\naics_image = AICSImag
e("../../data/EM_C_6_c0.ome.tif")\n\nnp_image = aics_image.get_i
mage_data("ZYX")\n```'}]},
```

# Fine-tuning OpenAI's gpt-3.5-turbo

- Training data should include successful general/system prompts

Training sample 1
| Complex general prompt |
| Specific prompt |

Expensive fine-tuning through repetition

Training sample 2
| Complex general prompt |
| Specific prompt |

Training sample 3
| Complex general prompt |
| Specific prompt |

- Inference with fine-tuned model

| Specific prompt |
| Specific prompt |
| Specific prompt |

Cheaper inference as the general prompt is „baked in" the model

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

Read more:
https://platform.openai.com/docs/guides/fine-tuning

# Fine-tuning OpenAI's gpt-3.5-turbo

- <u>Upload training data</u>

- Start fine-tuning job

- Test fine-tuned model

```
[11]:  client = openai.OpenAI()

       # upload and preprocess file
       training_file = client.files.create(
           file=open(training_data_file_path, "rb"),
           purpose='fine-tune',
       )
```

```
       # wait until preprocessing is finished
       while client.files.retrieve(training_file.id).status != "processed":
           time.sleep(30)

       print("Uploading / preprocessing done.")
```

```
Uploading / preprocessing done.
```

# Fine-tuning OpenAI's gpt-3.5-turbo

- Upload training data

```python
# start fine-tuning
fine_tuning_job = client.fine_tuning.jobs.create(
                            training_file=training_file.id,
                            model="gpt-3.5-turbo")
```

```python
job_details = client.fine_tuning.jobs.retrieve(
                        fine_tuning_job.id)
```

- <u>Start fine-tuning job</u>

```python
job_details.status
```

'validating_files'

```python
job_details = client.fine_tuning.jobs.retrieve(fine_tuning_
job_details.status
```

'running'

```python
job_details = client.fine_tuning.jobs.retrieve(fine_tuning_job.id)
job_details.status
```

'failed'

- Test fine-tuned model

```python
job_details = client.fine_tuning.jobs.retrieve(fine_tuning_job.id)
job_details.error
```

```
Error(code='invalid_training_file', message='The job failed due to an
invalid training file. Expected file to have JSONL format, where every
line is a valid JSON dictionary. Line 1 is not a dictionary.', param
='training_file')
```

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

ScaDS.AI
DRESDEN LEIPZIG

# Fine-tuning OpenAI's gpt-3.5-turbo



[Extern] Fine-tuning job ftjob-AptHl7VZCk2dC4JBOFYt0u8j succ...

OpenAI <noreply@tm.openai.com>
To Haase, Robert
Sat 6/1

Hi Leipzig University,

Your fine-tuning job `ftjob-AptHl7VZCk2dC4JBOFYt0u8j` has successfully completed, and a new model `ft:gpt-3.5-turbo-0125:leipzig-university::9VNFz3h3` has been created for your use.

Try it out on the OpenAI Playground, view the training results in the fine-tuning UI, or integrate it into your application using the Chat Completions Legacy Completions API.

Thank you for building on the OpenAI platform,
The OpenAI team

[Extern] Fine-tuning job ftjob-bANBDYKYUK7AJeaqrCqtLqFx fail...

OpenAI <noreply@tm.openai.com>
To Haase, Robert
Sat 6/1

Hi Leipzig University,

Unfortunately, your fine-tuning job `ftjob-bANBDYKYUK7AJeaqrCqtLqFx` has failed. See more details on the failure in the fine-tuning UI

Read the Fine-tuning Guide for more information on the expected usage the fine-tuning API.

Thank you for building on the OpenAI platform,
The OpenAI team

# Fine-tuning OpenAI's gpt-3.5-turbo

- Upload training data

- Start fine-tuning job

- Test fine-tuned model

```
model_name = job_details.fine_tuned_model
model_name
```

```
'ft:gpt-3.5-turbo-0125:leipzig-university::9VNFz3h3'
```

After training

At inference

# Prompt engineering techniques



Cropped from Sahoo et al 2024, licensed
https://arxiv.org/pdf/2402.07927

# Quiz:

- Assume I can enter my entire knowledge base into one very long prompt.
  - Why would it make sense to implement a RAG solution anyway?
  - Why would it make sense to fine-tune a custom model?
- In what scenario would one prefer the RAG over fine-tuning a model?

![ScaDS.AI Dresden Leipzig logo]

CENTER FOR SCALABLE DATA ANALYTICS AND
ARTIFICIAL INTELLIGENCE

# Exercises

## Robert Haase

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Exercise

- Hint, you can enter the OpenAI API-key like this at the beginning of notebooks:

```python
import os
os.environ['OPENAI_API_KEY'] = 'sk-...' #todo: enter your API key here
```

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

# Exercise: Prompt engineering

- Re-run image analysis code generation and elaborate on reproducibility.

- Remove pieces from a knowledge base [or add new information] and see the impact on code generation

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

# Exercise: Retrieval augmented generation

- Compare generated code for complex tasks

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

# Optional exercise: Fine-tuning

- **Only run the fine-tuning notebooks if you have a new knowledge base!**
  - Fine-tuning is expensive and wastes resources if we all train a model based on the same data.

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

# Exercises: LangChain

- Figure out when code is actually executed.



```
[6]: agent.run(input="Hi, I am Robert")

[6]: 'Nice to meet you, Robert!'

[7]: agent.run(input="What's my name?")

[7]: 'Your name is Robert'

[8]: agent.run("Can you reverse my name?")

[8]: "The response to your last comment was 'treboR', which i
```

```
[9]: agent.run("Do you know my name reversed and upper case?"
```

```
[9]: 'TREBOR'
```

Jupyter notebook screenshot:

20_langchain.ipynb

```
[9]: agent.run("Do you know my name reversed and upper
```

```
[9]: 'TREBOR'
```

## Exercise

Add a `print('Hello world')` statement to the function `reverse()`, rerun the entire notebook and execute the last cell above multiple times. Is the `Hello world` printed every time? If not, why?

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

https://haesleinhuepf.github.io/BioImageAnalysisNotebooks/
07_prompt_engineering/20_langchain.html

38

38

# Exercises: Prompting image analysis tasks

- Extend the LangChain notebook to enable the *agent* to

Robert Haase
@haesleinhuepf
BIDS Lecture 11/14
June 11th 2024

https://haesleinhuepf.github.io/BioImageAnalysisNotebooks/
07_prompt_engineering/30_langchain_bia.html

39

# Exercises: Prompting image analysis tasks

- Implement multiple segmentation tools and guide the *agent* to use the right one, e.g. for segmenting and image showing bright membranes

https://haesleinhuepf.github.io/BioImageAnalysisNotebooks/07_prompt_engineering/30_langchain_bia_choosing_algorithms.html

40