# Continue Dev

- „Local" Copilot-like extension
- Pros:
  - Privacy
  - Cost(?)
  - Customizable models
- Cons:
  - Speed(?)
  - Accuracy

ELTE | FACULTY OF INFORMATICS

# Continue Dev

- Install from extensions
- Use continue tab
  - Try with OpenAI API key
- Functionality similar to Copilot Chat-mode
- Custom commands/prompts available

# Continue Dev

- Local LLM-s with Huggingface
- https://huggingface.co/docs/text-generation-inference - Huggingface TGI
- Models at: https://huggingface.co/
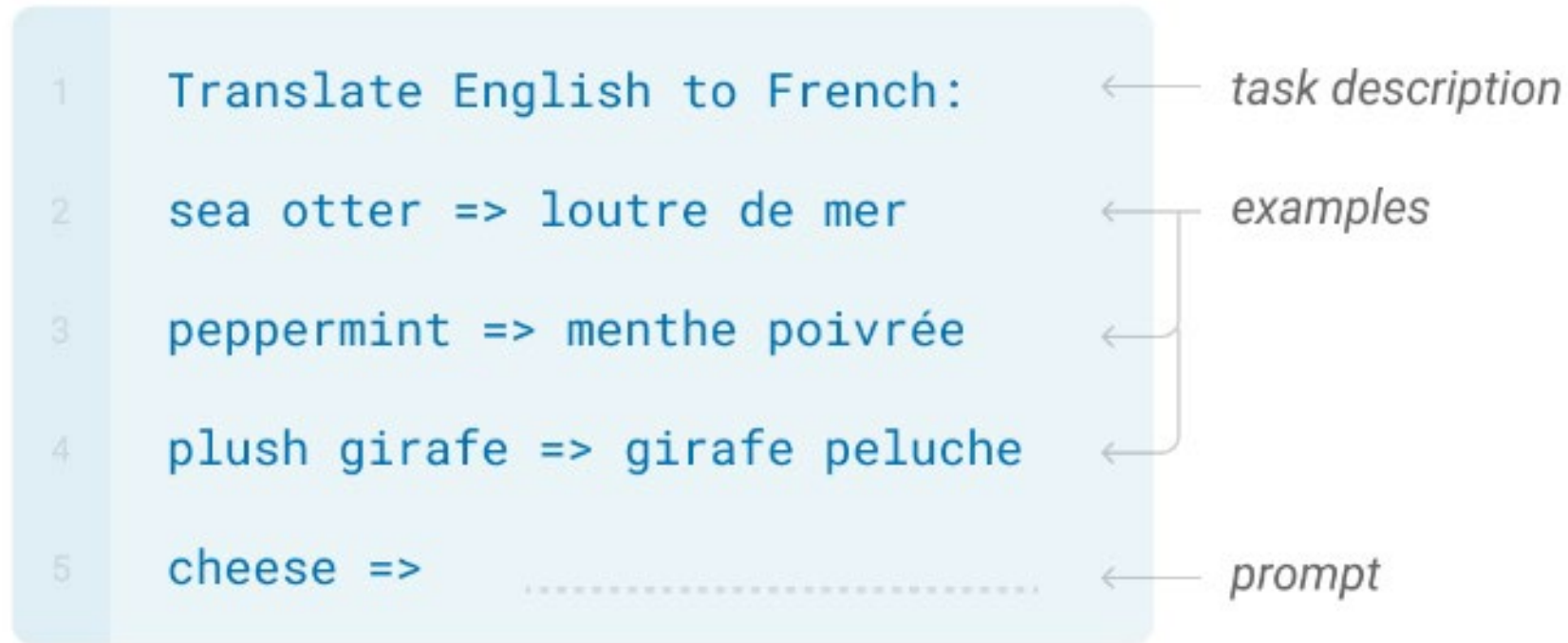- Supported models: https://docs.continue.dev/setup/select-model

# Prompting chat models

- System prompt
  - General definition for the whole conversation, does not get cut out of history
- User-Assistant prompts form the history which is limited. Breaking alternation is often allowed.
- Textual tool outputs might also be included in the conversation

System Prompt

User Prompt

Assistant response

User Prompt

Tool response

Assistant response

# In-Context Learning

- Few-shot prompts include the task definition and a few examples.



```
1   Translate English to French:        ←——  task description

2   sea otter => loutre de mer          ←——┐
                                           │  examples
3   peppermint => menthe poivrée        ←——┤

4   plush girafe => girafe peluche      ←——┘

5   cheese =>  ······················   ←——  prompt
```

# Problem Decomposition

- LLMs have limited single-step reasoning capabilities
- They are good at small common sense reasoning steps
- Hard problems should be decomposed into a Chain of Thoughts (CoT), where each thought is just a small step towards the solution.
- Latest models have built-in (by alignment) CoT capabilities

# CoT Example

Which one is older Bosch or the Eiffel tower?

The Eiffel Tower is older than Bosch. The Eiffel Tower was completed in 1889, while Bosch was founded in 1886.

Which one is older Bosch or the Eiffel tower?
Give your answer step-by-step!

1. The Eiffel Tower was completed in 1889.
2. Bosch is a company that was founded in 1886.
3. Therefore, Bosch is older than the Eiffel Tower by 3 years.

ELTE | FACULTY OF INFORMATICS

# Provide external knowledge

- Specific knowledge is missing from LLMs
- Models frequently hallucinate

# Retrieval Augmented Generation (RAG)

- Using embedding models to encode the meaning of text chunks (1-3k dimensional vectors)
  - Usually encoder-style LLMs finetuned for similarity search
- VectorDB – similarity search
- Hybrid search possible
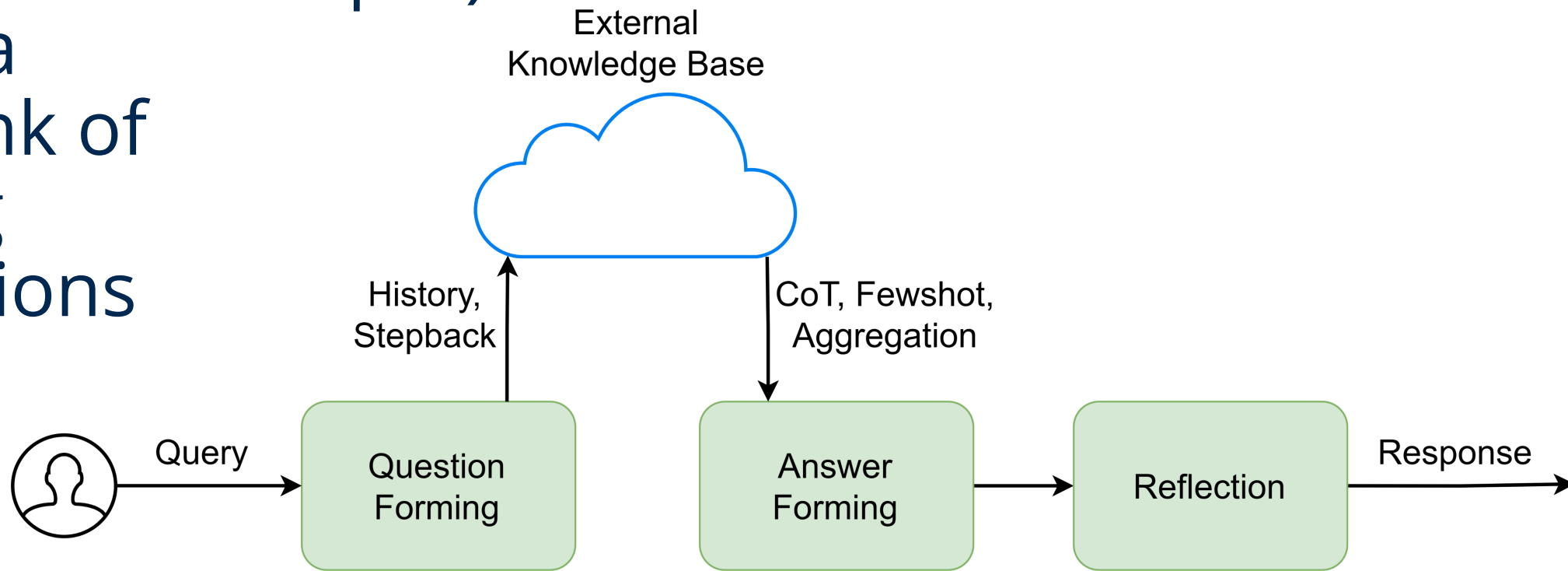- Adding knowledge from the same topic to the prompt

External Knowledge Base

Query

Response

LLM

# 2-Step RAG

- Questions are not always similar to the relevant knowledge:
  - References in the chat
  - Background information is better than specific
- Answers are assembled in a separate step

External Knowledge Base

History, Stepback

CoT, Fewshot, Aggregation

Query

Question Forming

Answer Forming

Response

ELTE | FACULTY OF INFORMATICS

# Reflection

- An extra reflection node (which reflects on the suitability of the output)
can filter a
large chunk of
remaining
hallucinations

External
Knowledge Base

History,
Stepback

CoT, Fewshot,
Aggregation

Query

Question
Forming

Answer
Forming

Reflection

Response

# Tooling

- Tools that can be called via text parameters and return text input are suitable for LLM chains
  - Programming kernel
  - Search engines
  - Knowledge graphs
  - Action handlers
  - Application APIs
- JSON is a popular format

Various Tools

Query

LLM

Response

# Agents

# Agents

- Specialist Agents are LLMs with external knowledge/tools that specialize in performing a single task
- Controller is an LLM, that uses a loop of planning, task execution and evaluation

**Action**
*External tool*

**Thought**
*Observation*

**Reasoning**
*CoT*

**Plan**
*Candidate resp.*

**Criticism**
*Reflection*

# Example: Design Assistant

- Integration of scientific tools and LLM pipelines

- Immature stage:
1st conference on LLM-aided design this June

- Agentified chat systems are viable

# Example: Design Assistant

# Example: Design Assistant



Engineer

Chat Agent

Controller

1. Find fan in model.
2. Calculate load
3. Find suitable bearing

Best Practice
Advisor

Query

Response

CAD/CAE Query
Composer

Query

Response

Catalogue Reader
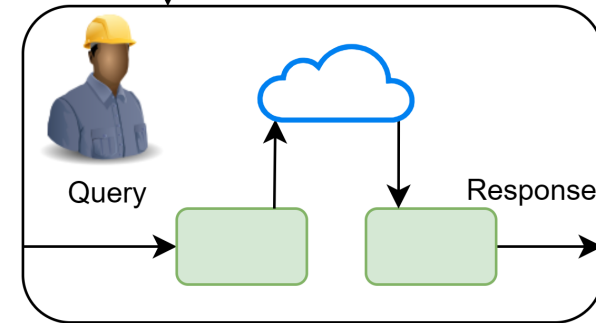
Query

Response

# Example: Design Assistant



Engineer → Chat Agent → Controller

{mode: "search", query: *"SELECT * WHERE name='fan'"*}

**Best Practice Advisor** — Query / Response

**CAD/CAE Query Composer** — Query / Response

**Catalogue Reader** — Query / Response

# Example: Design Assistant



Engineer

Chat Agent

Controller

PartID = 23
Type = Cooling Fan
ParentModel = EngineA

Best Practice Advisor

Query  Response

CAD/CAE Query Composer

Query  Response

Catalogue Reader

Query  Response

# Example: Design Assistant



Engineer ↔ Chat Agent → Controller

Find maximal RPM for "EngineA" cooling fans

Best Practice Advisor
Query → Response

CAD/CAE Query Composer
Query → Response

Catalogue Reader
Query → Response

# Example: Design Assistant

# Example: Design Assistant



Find maximal load for 1000 RPM!

Engineer

Chat Agent

Controller

Best Practice Advisor

Query | Response

CAD/CAE Query Composer

Query | Response

Catalogue Reader

Query | Response

# Example: Design Assistant

Engineer

Chat Agent

Controller

{mode: "simulate", params: *"partID: 23; setRPM: 1000"* }

Best Practice Advisor

Query — Response

CAD/CAE Query Composer

Query — Response

Catalogue Reader

Query — Response

# Example: Design Assistant

# Example: Design Assistant



{DB: "bearings", params: "{MAX_LOAD}, {DIM}"}

# Example: Design Assistant

# Example: Design Assistant
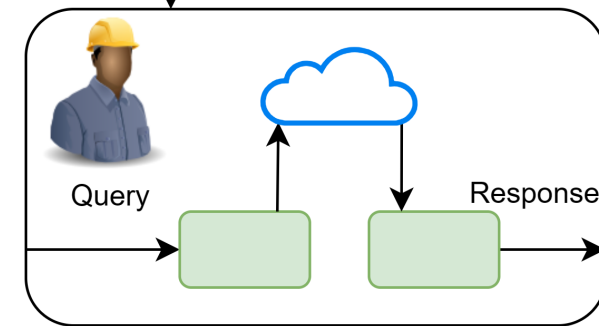
# Example: Design Assistant

# Coding Agents

- Planning included
- Tools are mostly code-related
- Might perform RAG on the repository

- Example:
- Pythagora.io

# Get an OpenAI API key
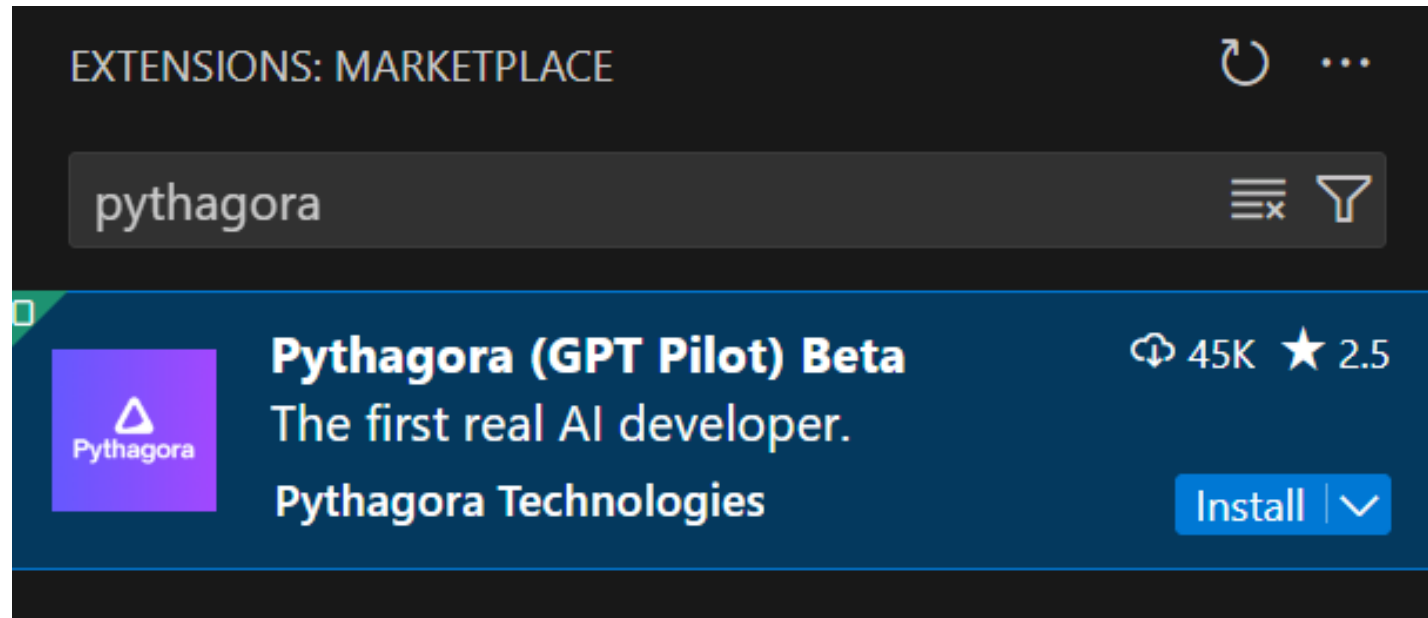
- Platform.openai.com
- Dashboard
- API keys
- Create new
- Save it for later

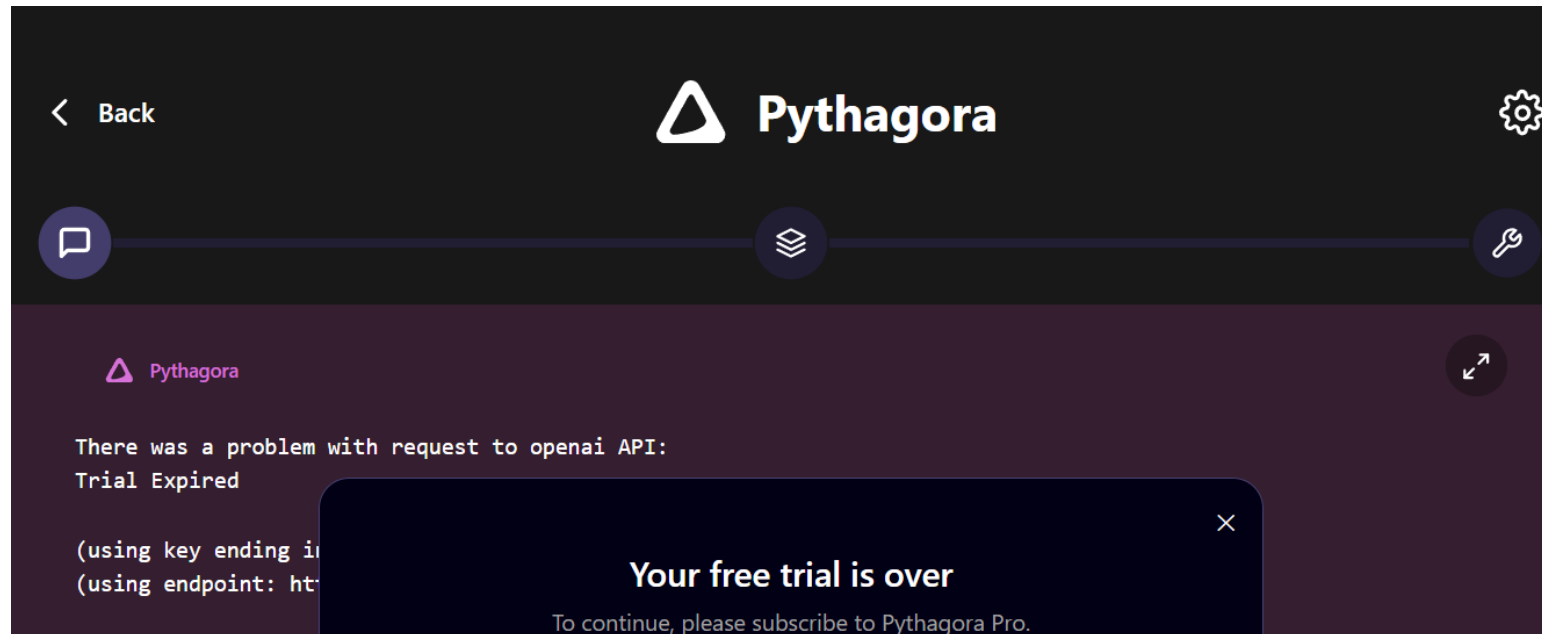**OpenAI**

ELTE | FACULTY OF INFORMATICS

# Get the Pythagora Extension

- Install from VSCode Extensions

# Get the Pythagora Extension

- Open the Pythagora tab, wait for it to setup
- Click on the Settings cogwheel

# Get the Pythagora Extension

- Set and copy your GPT Pilot path to **any value**
- Navigate to that folder

**GPT Pilot path**

e:\Code\gptpilot | **Change**

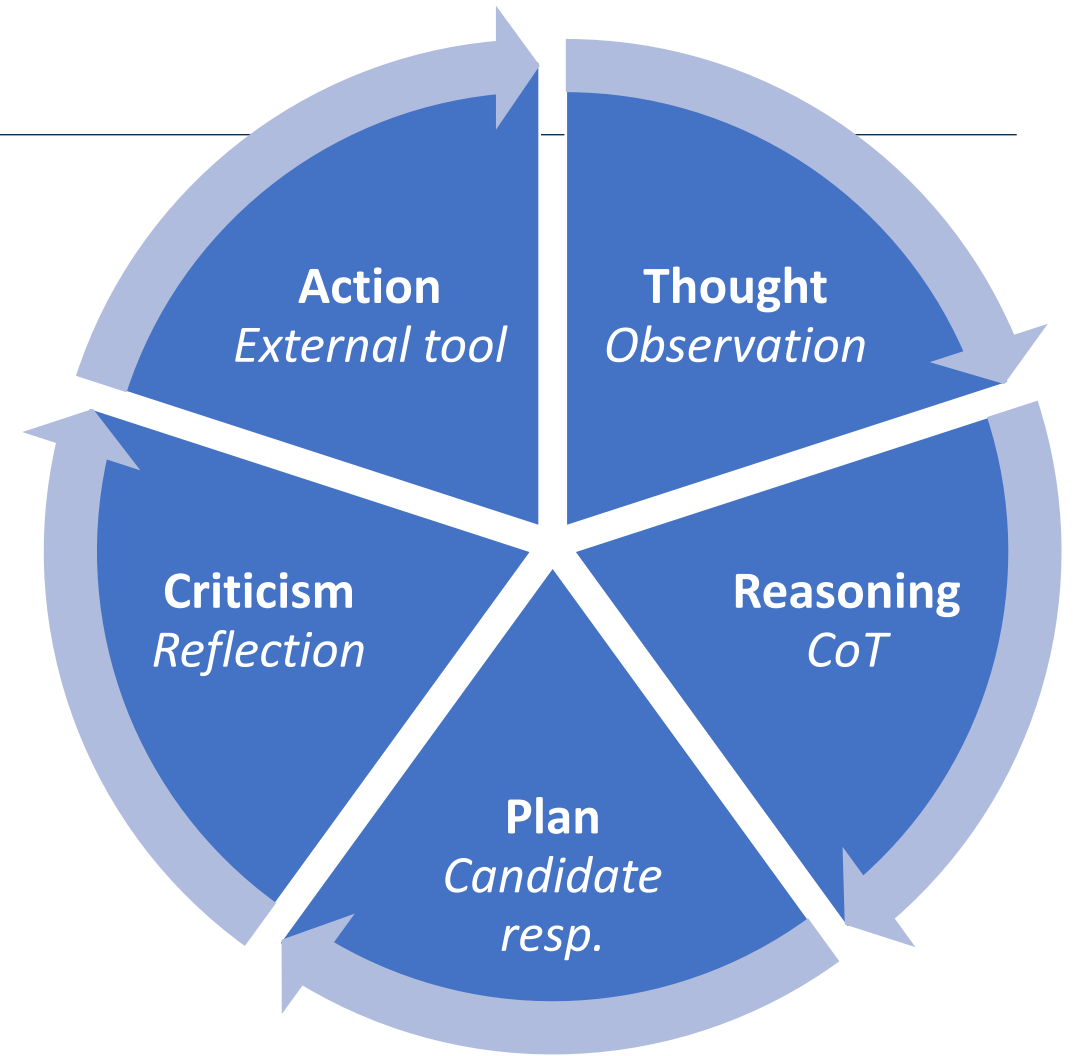ELTE | FACULTY OF INFORMATICS

# Get the Pythagora Extension

- Edit ./pilot/.env
- **OPENAI_ENDPOINT=https://api.openai.com/v1/chat/completions**
- **OPENAI_API_KEY=<YOUR KEY>**
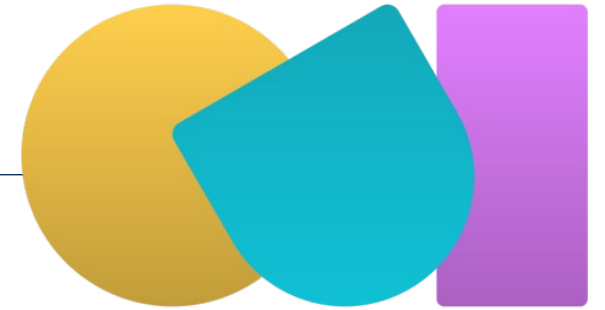- **MODEL_NAME=gpt-4o**
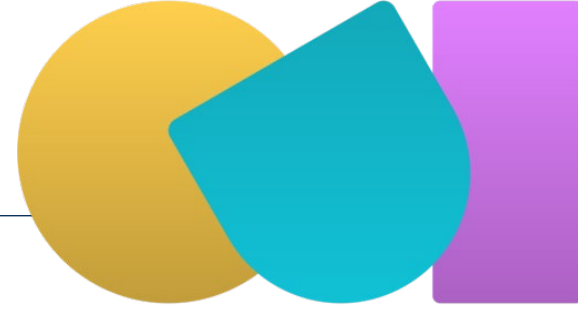
- Restart VSCode

# Enjoy the interview

- Create an app answering the questions by Pythagora
- Watch it plan, act, revise!

# UI Dev Tool

- OpenUI
- Clone from GitHub
  - git clone https://github.com/wandb/openui
  - $env:OPENAI_API_KEY = <your key>
  - cd openui/backend
  - pip install .
  - python -m openui

# UI Dev Tool

- Set the model to gpt-4o in the settings
- Draw or screenshot a UI!
- Define local modifications
- Dynamic connections
- Basic animations
- etc