# Generative AI for software Development:

**Authors:**
**EL ALOUT Ismail, EL HABTI Ouassel, KARMAOUI Oussama**
**LAHMOUZ Zakaria, HALILOUA Othmane**

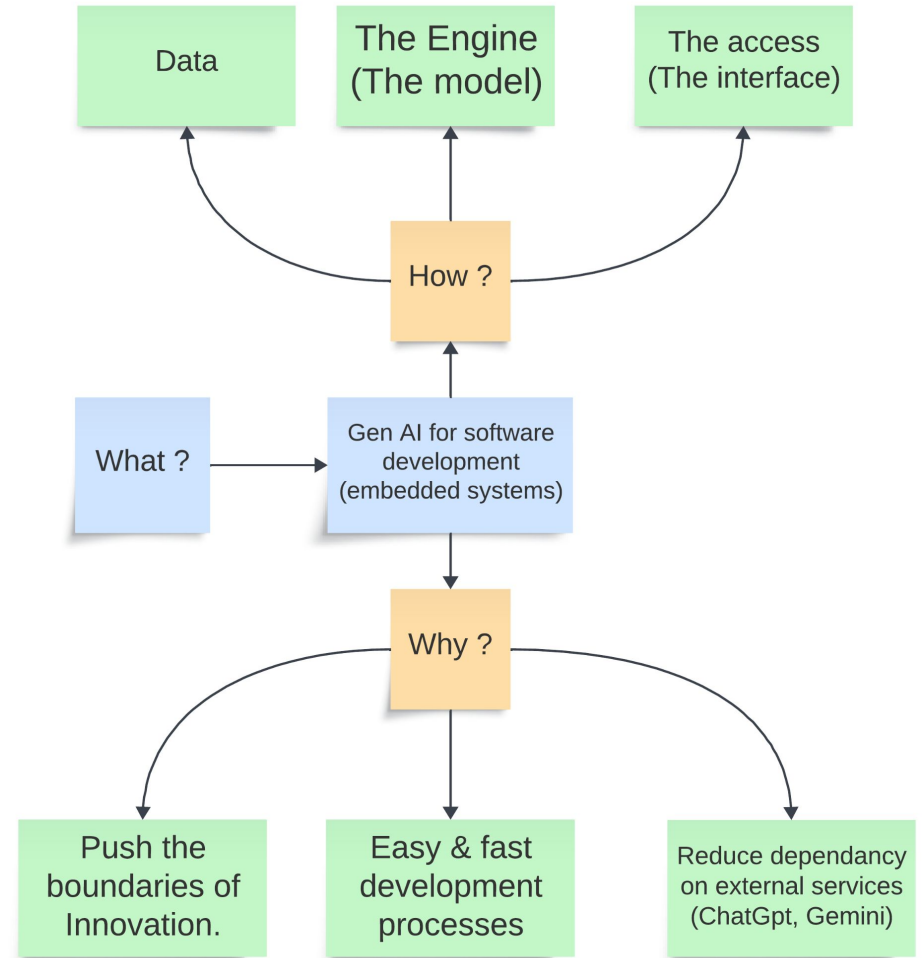**Supervisor:**
**Mr CHAARI Lotfi**

# PLAN

- Introduction
- Work Organization and task planning
- Development and versioning tools
- Dataset building
- File System Architecture and Data Pipelines
- Fine-tuning
- Deployment
- Demonstration
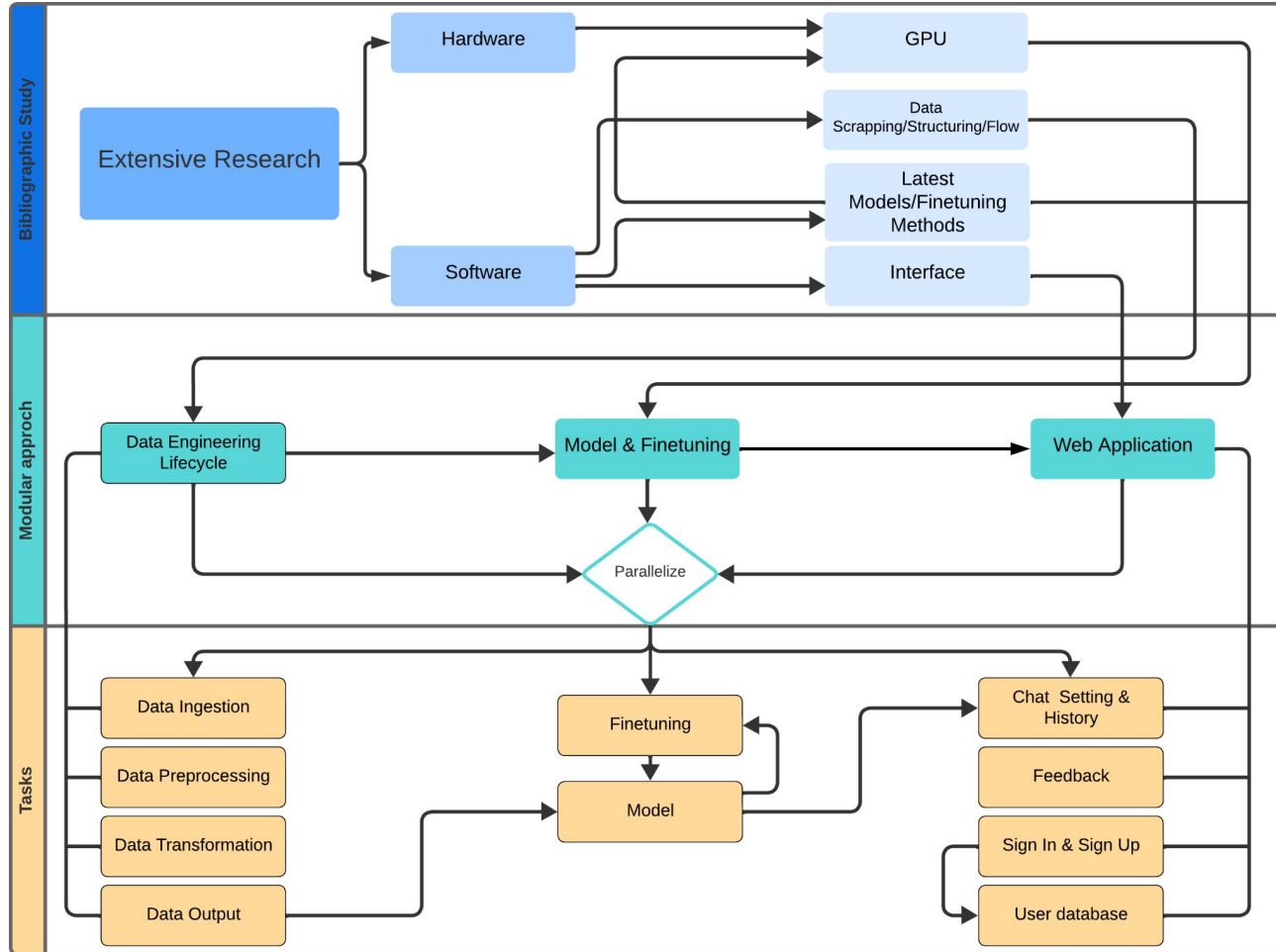- Conclusion

# INTRODUCTION

The What ? The why ? And The How ?

- AI has pushed the boundaries of what can be achieved.

- One would like to live up to the standards.

- One would want to be a protagonist in the innovation act, to better his work, and reduce his dependence.

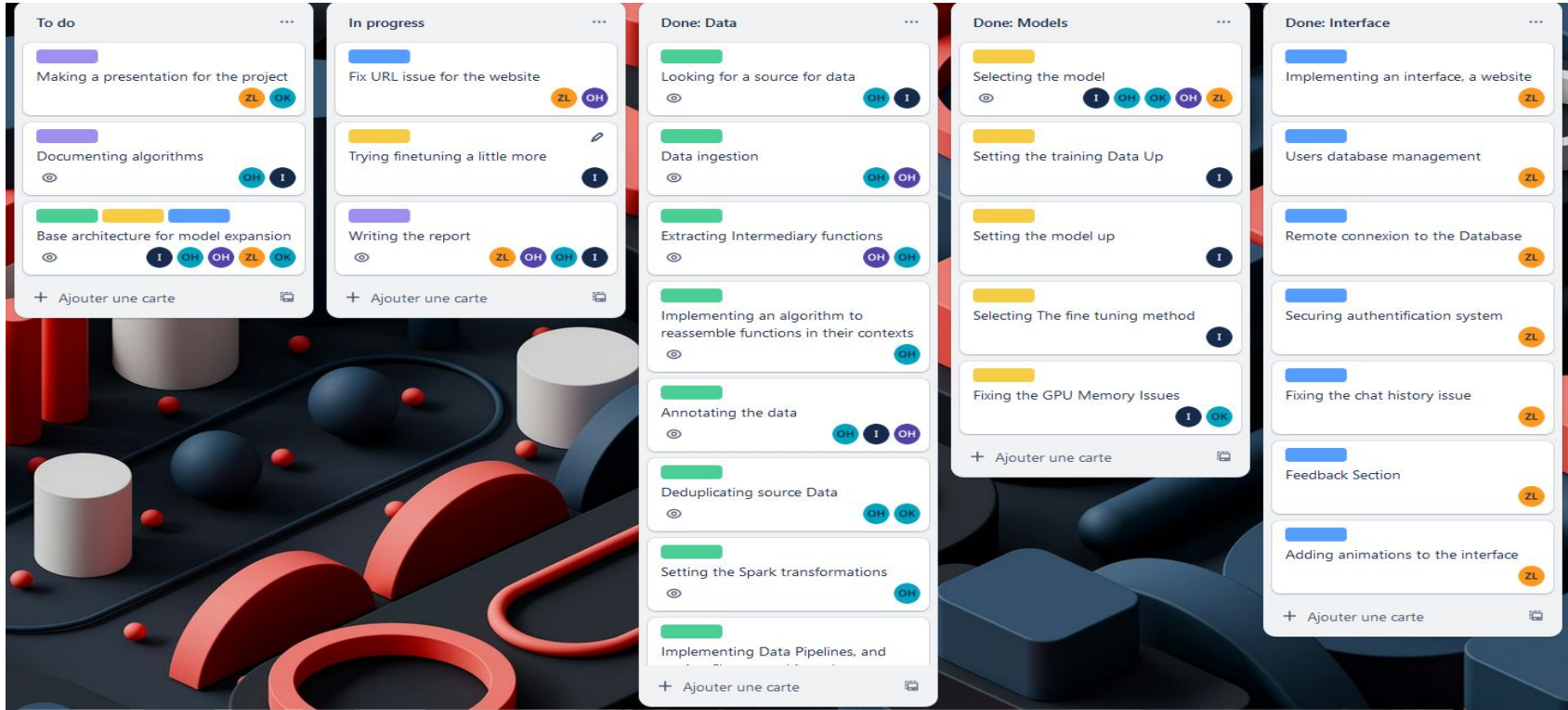- To that end, three pivotal components: Data, the model, the interface.

# Work organization and task planning

# Task planning & Thought Process:

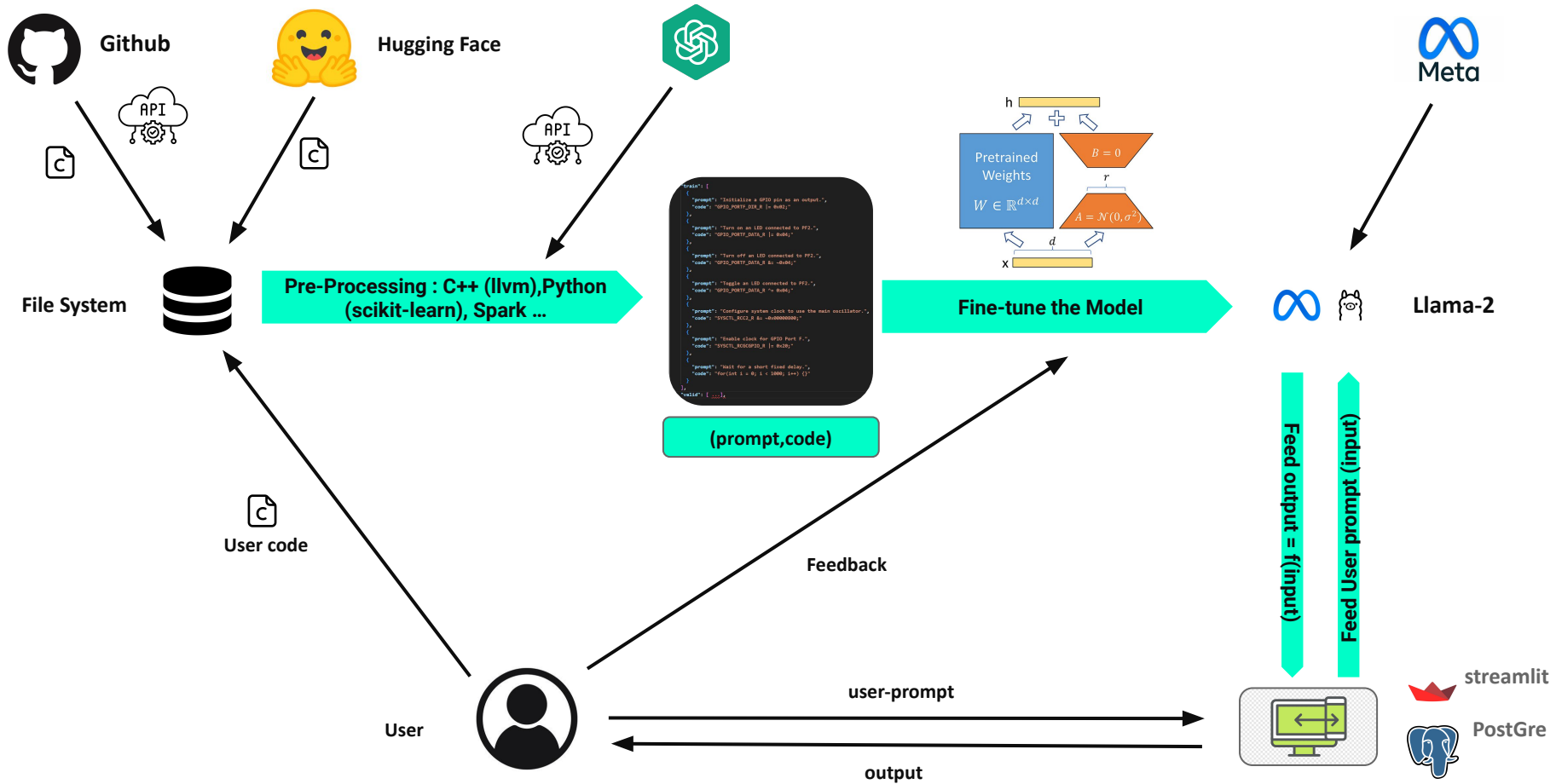# Project Management:

The use of Trello allowed us to better depict the tasks that require tackling, monitoring, refining or initiating. We used our modular development approach so that member can be assigned easily to corresponding tasks, with brackets to evaluate the completion rate.

# Development & Versioning Tools

# Development:

| | | |
|---|---|---|
| Python | | Versatility, readability, and a rich ecosystem of libraries for all development modules. |
| C++ | | Efficient thanks to its low-level control to parse and manipulate code snippets effectively. |
| Spark | | Data preprocessing and transformation tasks, harnessing its distributed computing capabilities to handle massive datasets efficiently. |
| Pytorch | PyTorch | Dynamic computational graph construction and a high degree of flexibility. |
| Streamlit | | Hosting the chatbot interface, providing an intuitive and interactive platform. |
| PostGre | | Robust database management system, facilitating efficient storage and retrieval of user data within the application. |
| CSS | | Used to style the web application, enhancing user experience and interface aesthetics. |

Github

Hugging Face

Meta

File System

Pre-Processing : C++ (llvm),Python (scikit-learn), Spark …

Fine-tune the Model

Llama-2

Pretrained Weights

$W \in \mathbb{R}^{d \times d}$

$B = 0$

$A = \mathcal{N}(0, \sigma^2)$

h

x

r

d

(prompt,code)

User code

Feedback

Feed output = f(input)

Feed User prompt (input)

User

user-prompt

output

streamlit

PostGre

10

# Versioning:
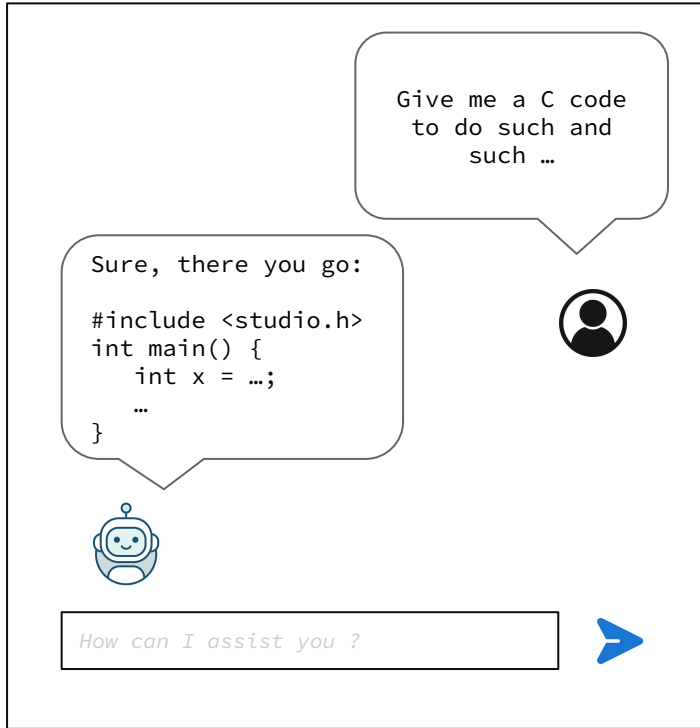
| | |
|---|---|
| **Github** | ❏     Allows for parallel development and smoother integration of changes into the main project. |
| **Google Drive** | ❏     Flexibility with Google Colab which was used during Fine-tuning.<br>❏     Allows to manage and transfer large files. (the model's configuration files) |

# Dataset building

# Data Scraping:

## First of all, what do we need ?

Give me a C code to do such and such …

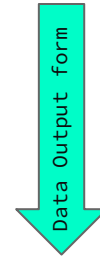Sure, there you go:

```
#include <studio.h>
int main() {
    int x = …;
    …
}
```

How can I assist you ?

The AI must understand the prompt (the user's request), and the code provided (the chatbot answer).
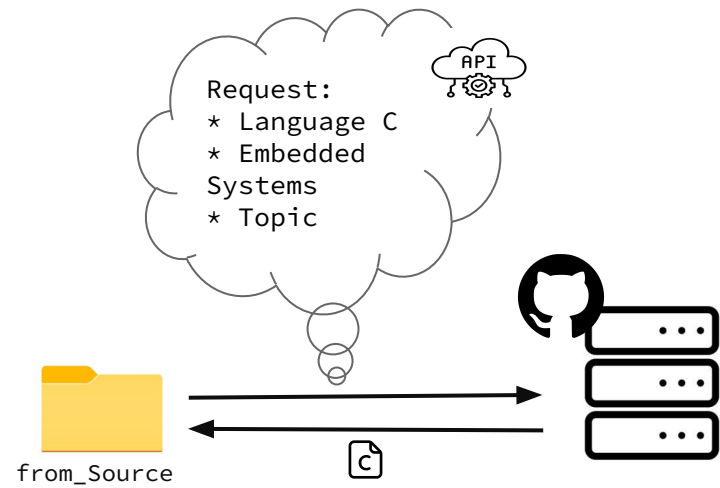
Our Dataset must be a set of coupled Prompts and corresponding codes.

Data Output form

("write a C code that does such…","#include<studio.h>\nint main…")

13

## Code Snippets:

- Go for Open Source platforms: rich in code, no cost.
- Automating the downloading process: The use of GitHub's Python API.
- We need a Storage System: The exploitation of an external Hard Drive.

Request:
* Language C
* Embedded Systems
* Topic

API

from_Source

C

## Topics:

Embedded Networking

Real Time Operating Systems …

## Prompts:

We need the prompt that would make an appropriate request for each code snippet. ChatGpt was of great assistance, we have no expertise in Embedded Systems. Again requires automation, can't be done manually.

# Data Pipelines & File System Architecture

# Data Ingestion

Data ingestion is the process of integrating data into the service upon which the development in ongoing.

The initial was simple. Use a physical Storage system for data ingestion.

The development was carried on independently of the specifics of physical storage systems, all that was needed was the location of the input.

Fertile ground to migrate to cloud solutions, for more, speed flexibility, and fault tolerance.

# Data Preprocessing:

Having thousands of raw C code snippets, how can we clean up ?

- Deduplication: No need to have two codes that do the same thing.
- Code cleaning: Comments removal, empty lines removal… (comments can be used ?)
- Deriving more insights from each code: For more specific requests.
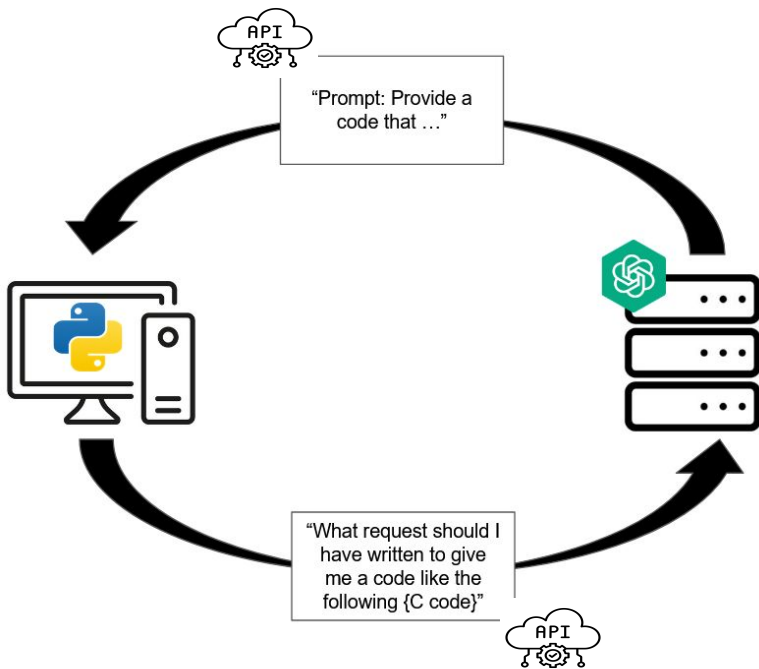- Reconstruction of derived codes in the context of the father code.

## Data Transformation

- C code:"""//…

```
unsigned int fec_golay2412_decode_symbol(unsigned int _sym_enc){

        if (_sym_enc >= (1<<15)) {

                fprintf(stderr,"error, fec_golay2412_decode(), input
        symbol too large\\n");

                exit(1);

        }

        return 0;}

//…"""
```

- Prompt: (provided by chatGPT)

"Can you provide me with a function to decode a received symbol encoded with this code? I need to ensure that the input symbol size is within the correct range for processing."



"Prompt: Provide a code that …"

"What request should I have written to give me a code like the following {C code}"

18

```
#include <avr/io.h>
#include <util/delay.h>

#define LED_PIN 0

void init_LED(void) {
    DDRB |= (1 << LED_PIN); // Set LED pin as output
}

void toggle_LED(void) {
    PORTB ^= (1 << LED_PIN); // Toggle LED
}

int main(void) {
    init_LED(); // Initialize LED pin

    while (1) {
        toggle_LED(); // Toggle LED
        _delay_ms(500); // Delay 500 milliseconds
    }

    return 0;
}
```

```
#include <avr/io.h>
#include <util/delay.h>

#define LED_PIN_MASK (1 << 0)

void init_LED(void) {
    DDRB |= LED_PIN_MASK; // Set LED pin as output
}

void toggle_LED(void) {
    PORTB ^= LED_PIN_MASK; // Toggle LED
}

int main(void) {
    init_LED(); // Initialize LED pin

    while (1) {
        toggle_LED(); // Toggle LED
        _delay_ms(500); // Delay 500 milliseconds
    }

    return 0;
}
```
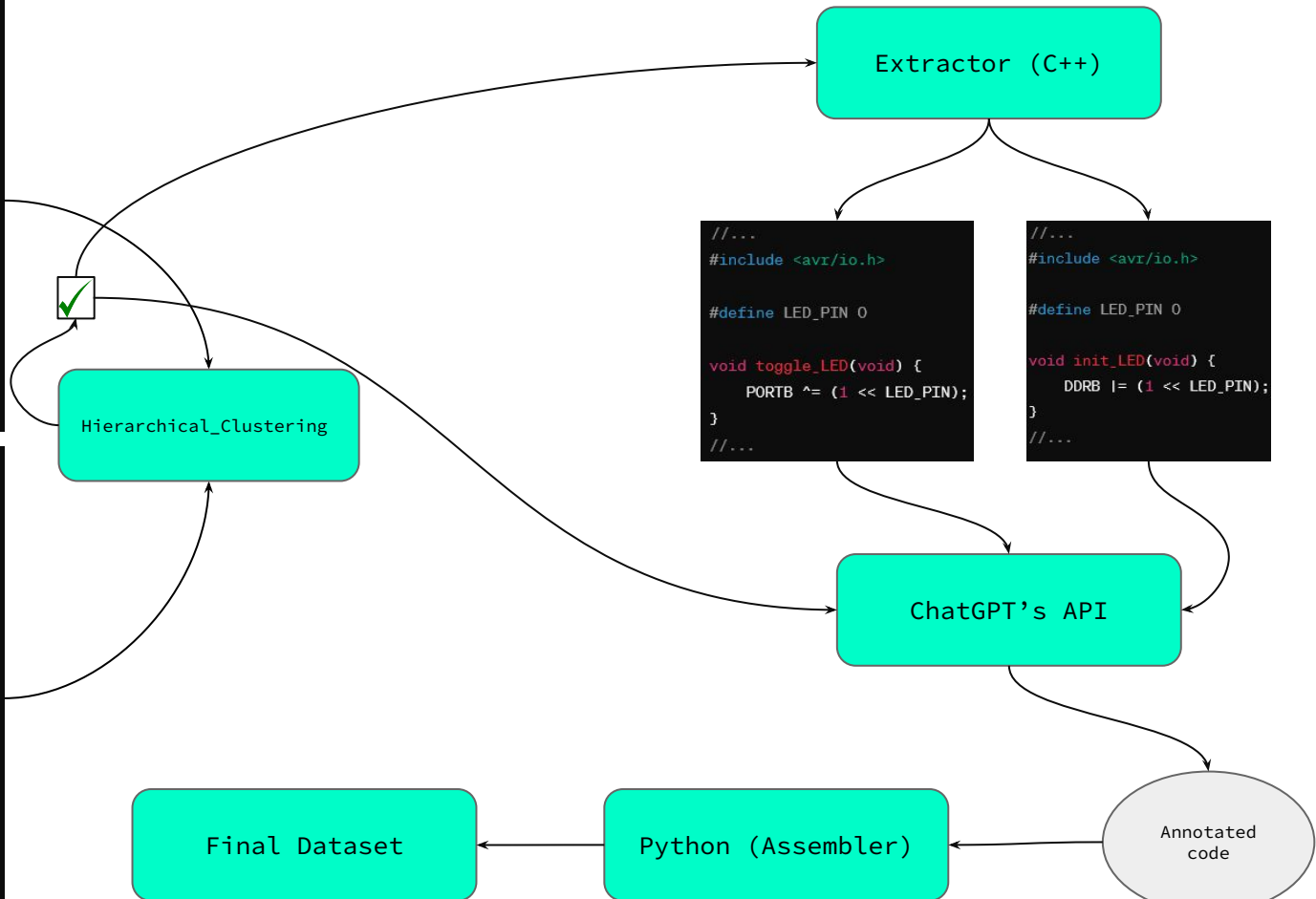
Extractor (C++)

Hierarchical_Clustering

```
//...
#include <avr/io.h>

#define LED_PIN 0

void toggle_LED(void) {
    PORTB ^= (1 << LED_PIN);
}
//...
```

```
//...
#include <avr/io.h>

#define LED_PIN 0

void init_LED(void) {
    DDRB |= (1 << LED_PIN);
}
//...
```

ChatGPT's API

Annotated code

Final Dataset ← Python (Assembler) ← Annotated code

# Data output



3879 instances of prompt code

# Fine-tuning

# Model selection: Llama-2



- ❏ Llama-2 is a family of large language model released by Meta, based on the Transformer Decoder-only architecture.
- ❏ Decoder-only means that the model utilizes only the decoder part of the Transformer architecture.
- ❏ There are 3 released versions by Meta with different parameter sizes: 7B, 13B and 70B.

  *Selected Model:* Llama-2 7B because of limited Hardware resources.



The Transformer architecture

# What is Fine-tuning ?

**MAIN IDEA:** Customize a pre-trained large language model for a particular task.

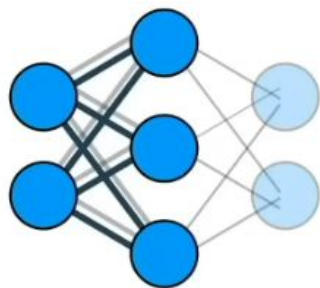| *Full Fine-tuning:*<br>Train each layer of the model on the target dataset. | *Parameter efficient Fine-tuning (PEFT):*<br>Freeze some layers and only adjust the weights of a smaller set of the last layers. |
|---|---|
| ❌ Very expensive computationally.<br>(Needs to allocate around 112GB for a 7B model).<br><br>❌ Long training time. | ✓ Less computational and storage costs.<br><br>✓ Less training time.<br><br>✓ Reach the same accuracy as a full fine-tuning. |

# PEFT approach: LoRA adapter

Fine-tunes the model by incorporating trainable rank decomposition matrices into each layer (specifically Feed-Forward layers).

$$h(x) = W_0 x + \Delta W x \qquad \Delta W = BA$$
$$= W_0 x + BAx$$

$x \rightarrow h(x)$

$$\left( \; \boxed{W_0} \; + \; \boxed{B \; A} \; \right) x = h(x)$$

**Frozen**      **Trainable**

$d = 1,000$
$k = 1,000$ $\implies$ $(d \times r) + (r \times k) = \textbf{4,000}$
$r = 2$ **trainable parameters**

$$W_0, \Delta W \in R^{d \times k}$$
$$B \in R^{d \times r}$$
$$A \in R^{r \times k}$$
$$h(x) \in R^{d \times 1}$$

# Fine-tuning Results:

❏   Number of Epochs: 20 epochs.

❏   Batch size: 1

❏   Training time: approximately 19 hours.

❏   Training Loss: 0.045

**Note:** This fine-tuning requires at least a GPU with 24GB VRAM. It cannot be done with the free Google Colab version.

train/loss

| Run ↑ | Smoothed | Value | Step | Relative |
|-------|----------|-------|------|----------|
| ● runs\Mar12_10-00-33_pop-os | 0.05 | 0.0453 | 29,084 | 7.474 hr |

*Interpretation:* The fluctuations are common for small batch size choices. The overall trend which shows that the loss is decreasing over time means that the model is still learning.

# Deployment

Chatbot for C code generation

How may I assist you today?

Your message

Creating an authentication system is essential:

- Enhance security

- Communicate with users

- Multiaccess of the app

- Subscription management

Navigation
- ◯ Login
- 🔴 Signup

username

email

password
👁

Signup

# Improve the user experience

# Demonstration

# Reinforcement learning integration

prompt

LLM

Generate samples

Go trough

COMPILER
+ UNIT
TESTS

Result

Feedback signal

REWARD
FUNCTION

Use the samples as a
seed (a base) for next
code generations

All
samples
failed ?

no

yes

REPAIR
NETWORK

# Conclusion

- Realised the ability to build the required dataset for our use case from scratch.
- Successful adaptation of the LoRA method to our model.
- Interactive Web Interface.
- Future improvement through Reinforcement learning.
- Soft skills: Communication, Adaptability, Flexibility, Initiative taking.

# References

- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen: "LoRA: Low-Rank Adaptation of Large Language Models" https://arxiv.org/abs/2106.09685
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, Colin Raffel: "Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning" https://arxiv.org/abs/2205.05638
- https://www.entrypointai.com/blog/lora-fine-tuning/
- https://www.leewayhertz.com/parameter-efficient-fine-tuning/
- https://arxiv.org/pdf/2207.01780.pdf
- Getting Started with Streamlit for Data Science: Create and deploy Streamlit web applications from scratch in Python:Tyler Richards