

# Knowledge Graph Generation

Project S9 Term 2023 - 2024



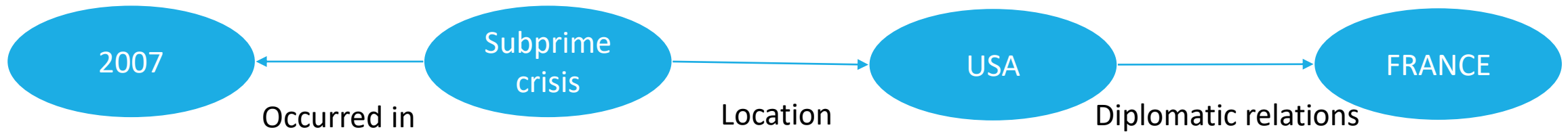
# PLAN :

- INTRODUCTION – MOTIVATION
- PIPELINE
- M\_REBEL
- R&D MERGE
- ALL\_MINI
- FINETUNING + METRICS
- IMPLEMENTATION
- STORING KB
- USER/ADMIN INTERFACE
- PROJECT MANAGEMENT
- CODE CARBON
- CONCLUSION

# Introduction – Motivation

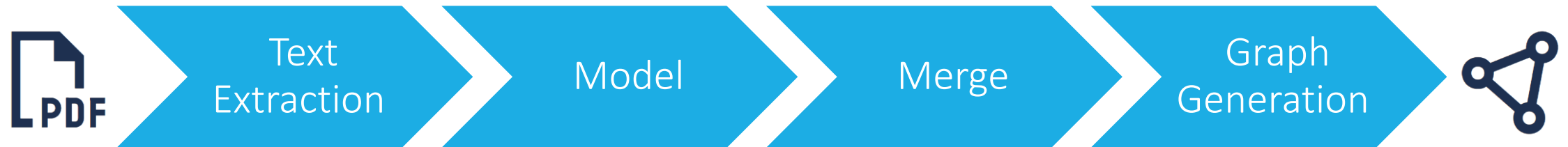
---

## What is a Knowledge Graph?

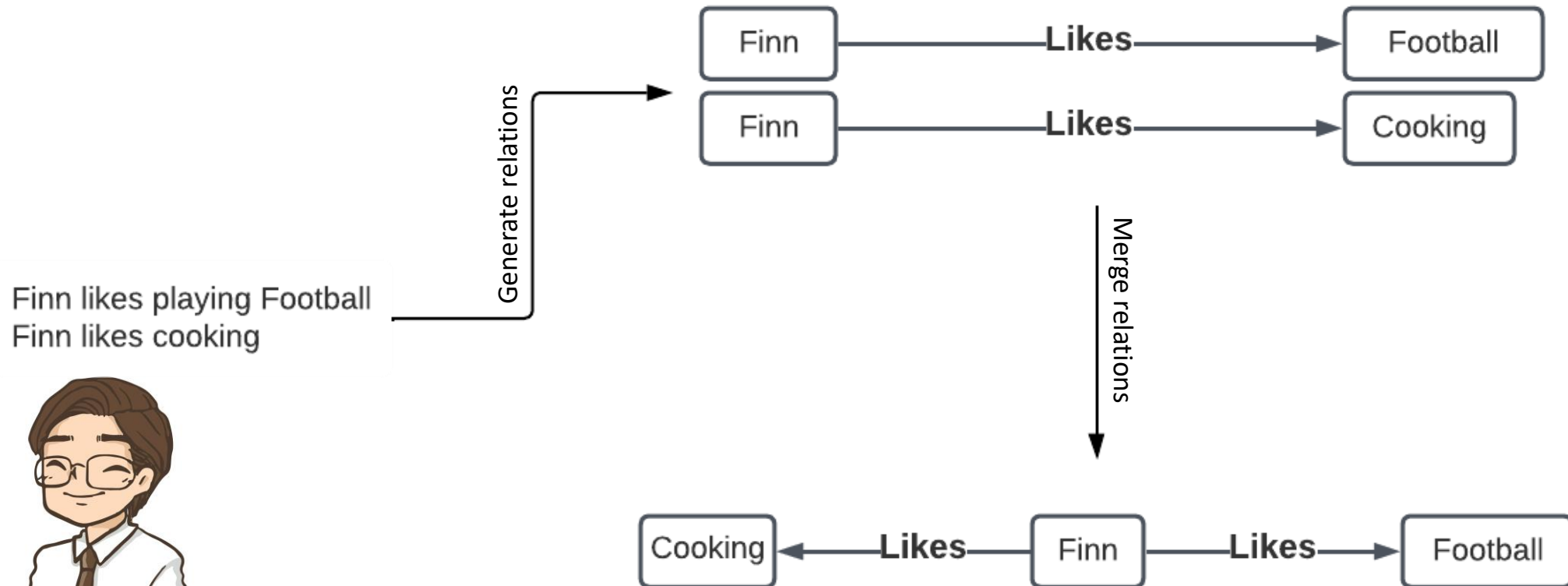


# Pipeline

---



# Example



# MRebel

---

## **MREBEL'S Foundation:**

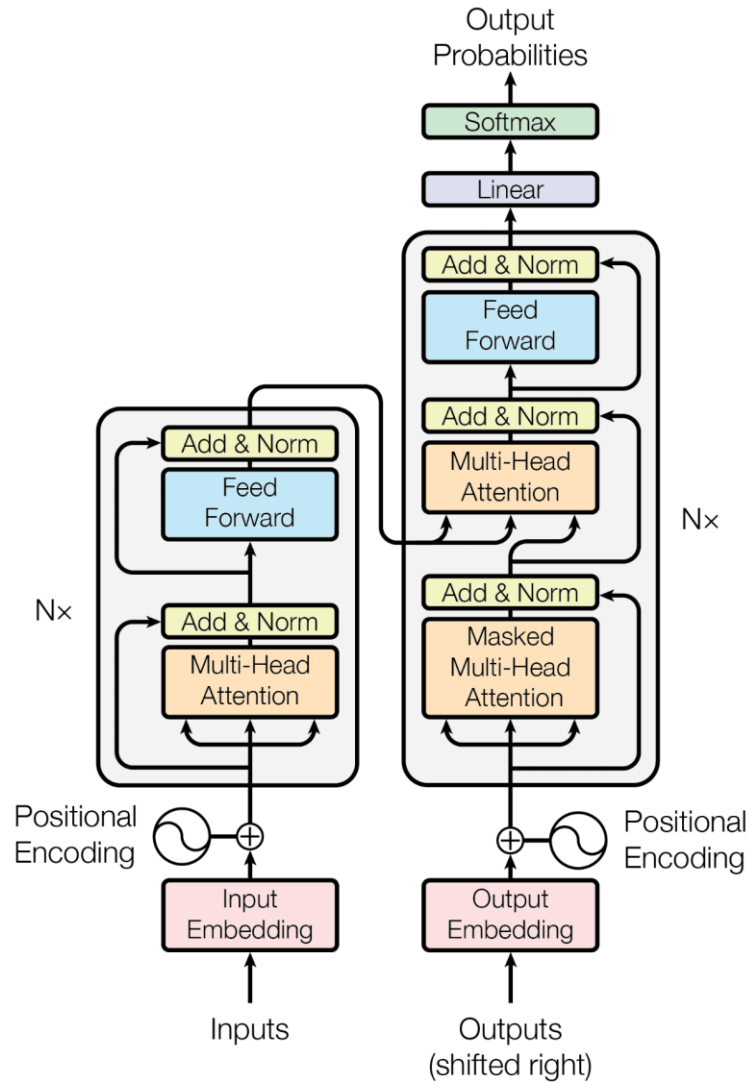
Built on BART: A hybrid model combining bidirectional context and autoregressive generation.

## **Relation Extraction as seq2seq:**

Reframes Relation Extraction: Converts the task into a sequence-to-sequence language generation problem

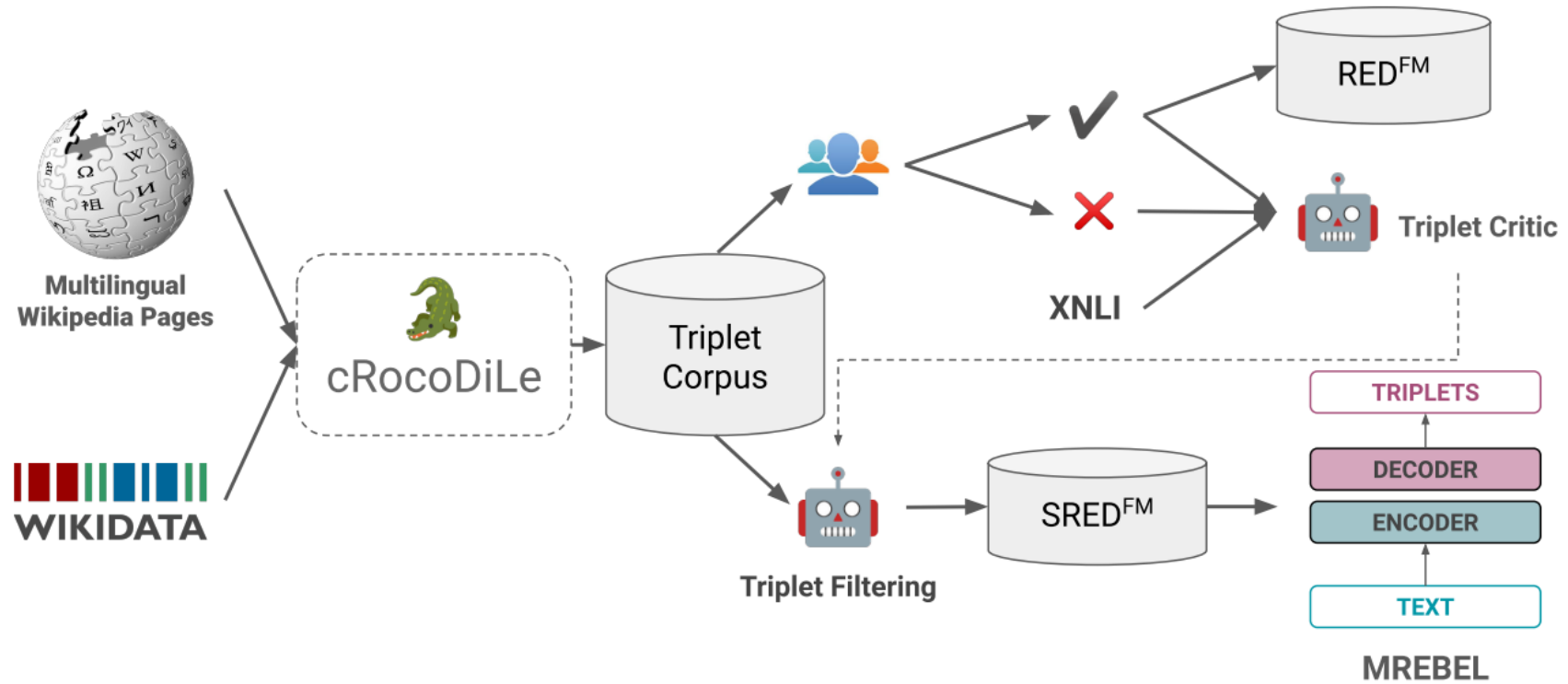


# The Transformer Architecture



*The encoder-decoder structure of the Transformer architecture  
Taken from "Attention Is All You Need"*

# MRebel



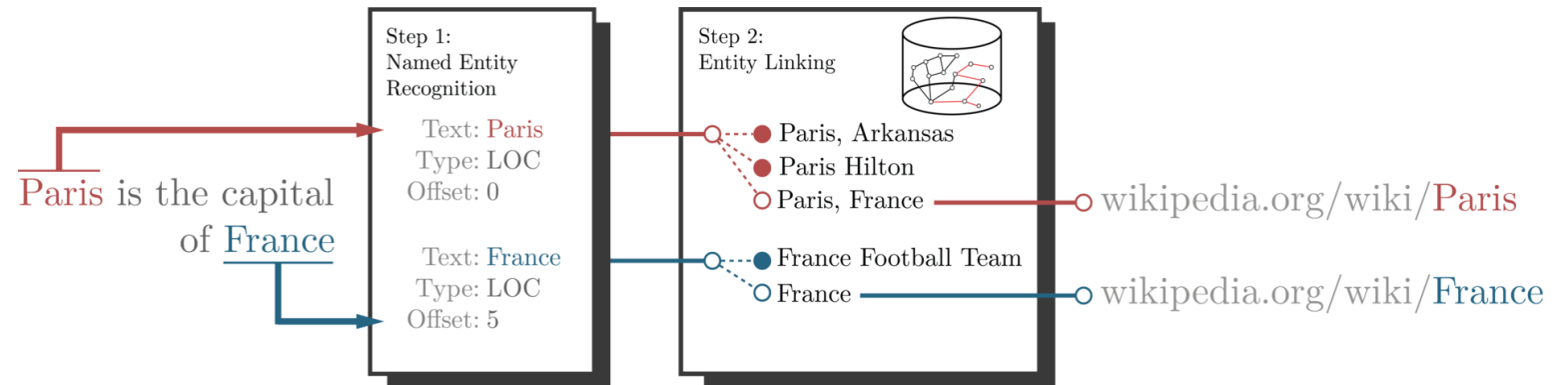
Complete pipeline for the creation of REDFM, SREFM and mREBEL.  
Retrieved from <https://doi.org/10.48550/arXiv.2306.09802>



# Wikipedia-Based Entity Merging Approach

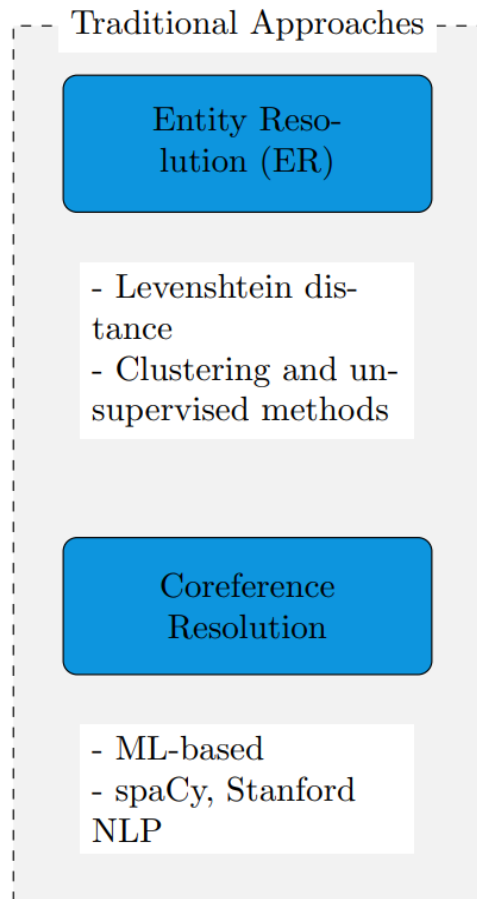
## Approach Evaluation:

- Time consuming page fetching.
- Repetitive network requests to external servers.
- Scalability issues



*Illustration of entity alignment using WikiPedia pages*

# Exploring Solutions For Entity Merging Challenges



## Named Entity Disambiguation (NED)

- Wikipedia/DB for disambiguation
- Linking to unique IDs

## Semantic Similarity

- word2vec, GloVE
- Relies on text but not context

## Graph-based Methods

- Effective in concrete cases
- Cluster & community creation

## Sentence Similarity Tasks

- All-miniLM
- Contrastive learning
- Promising Approach

# Introduction to the All-MiniLM Model

## Sentence Transformers family

Sentences --> 384-dimensional vector space.

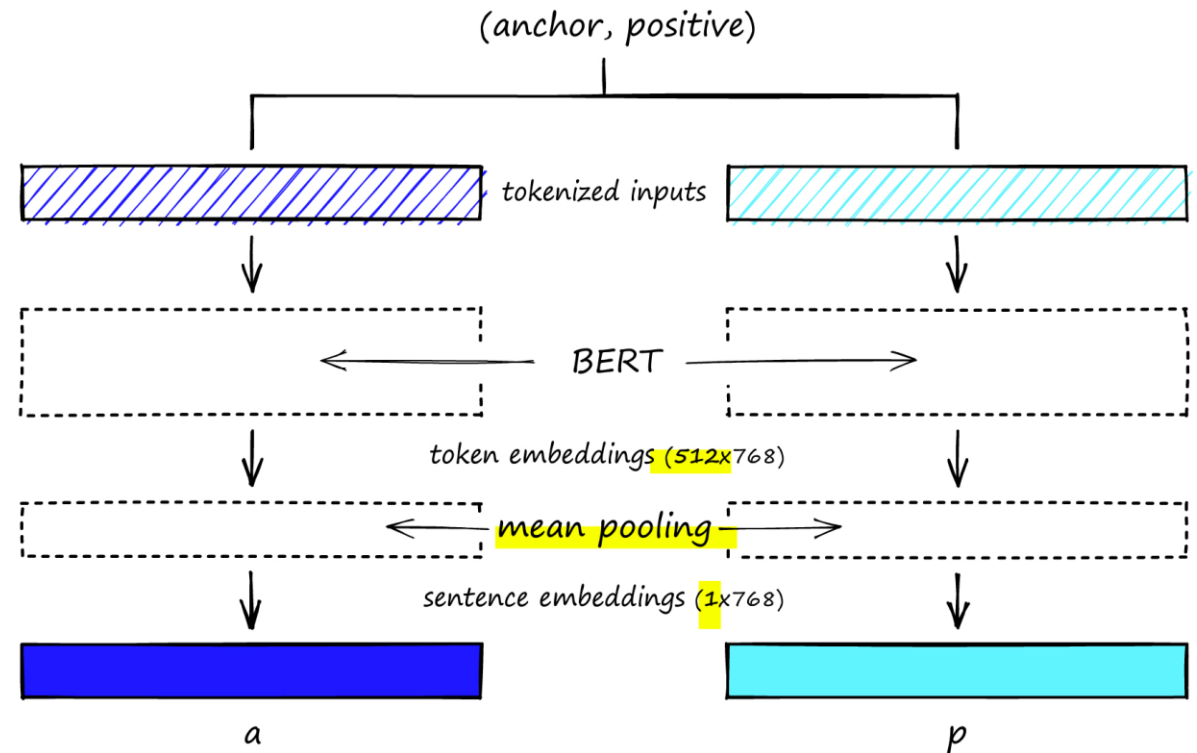
## Semantically Encode meaningful sentence embeddings

Semantic similarity comparison...

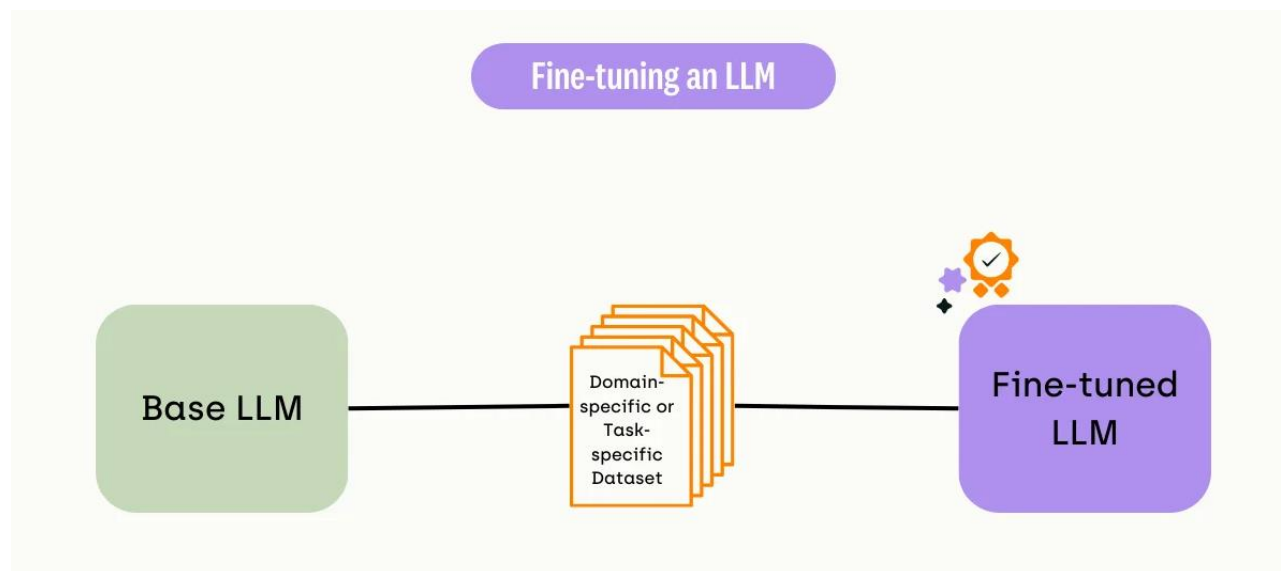
## Advantages:

Fast inference, small size

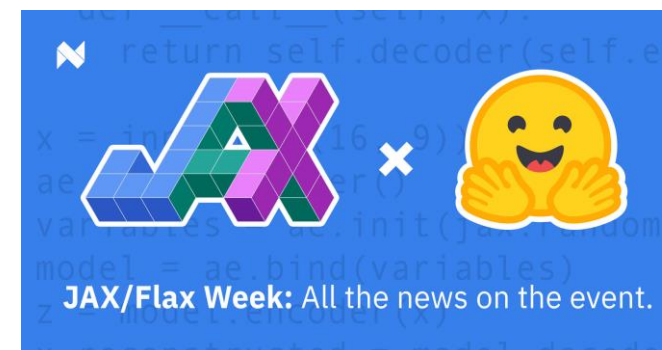
Great balance between performance and efficiency



# Introduction to the All-MiniLM Model



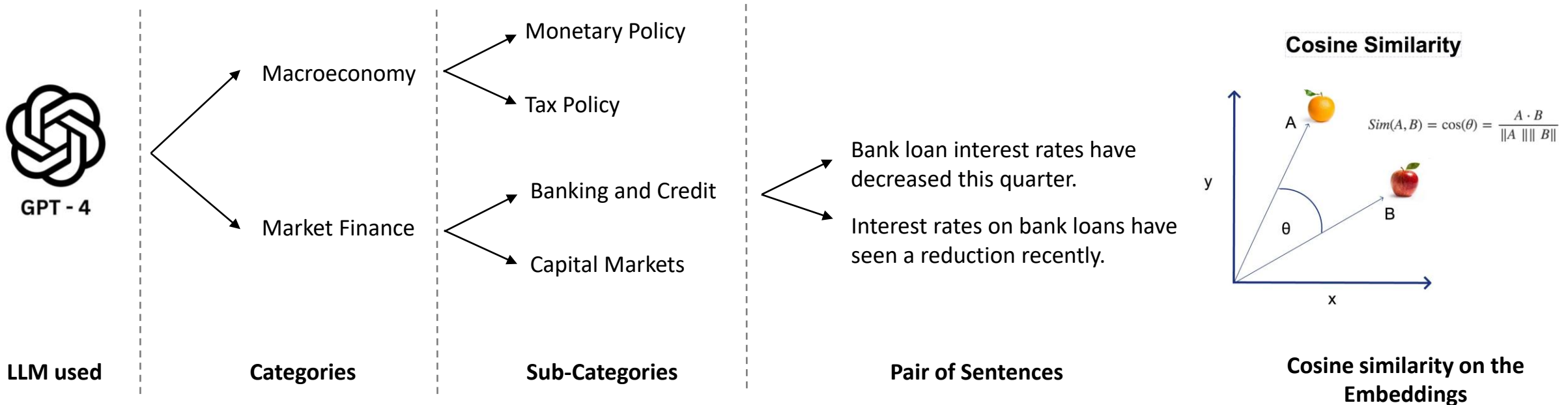
Fine-Tuning Scheme



Fine-tuned on a 1 billion sentence pair dataset during an HF event

The fine-tuning process involved a **contrastive learning** objective.

# Custom Fine-tuning of all-MiniLM



# Custom Fine-tuning of all-MiniLM

---

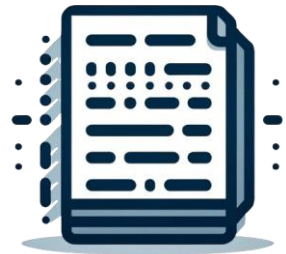
## Sample:

The Federal Reserve cut interest rates  
Interest rates were reduced by the central bank

## Logits outputs

Base Model : **0,765**

FT Model : **0,964**



## Assessing the results:

350 Test Sentences (10% of the train set's size)

Balanced Class Distribution

Merge Threshold: **0,8** *(if above then merge)*

## Model accuracy

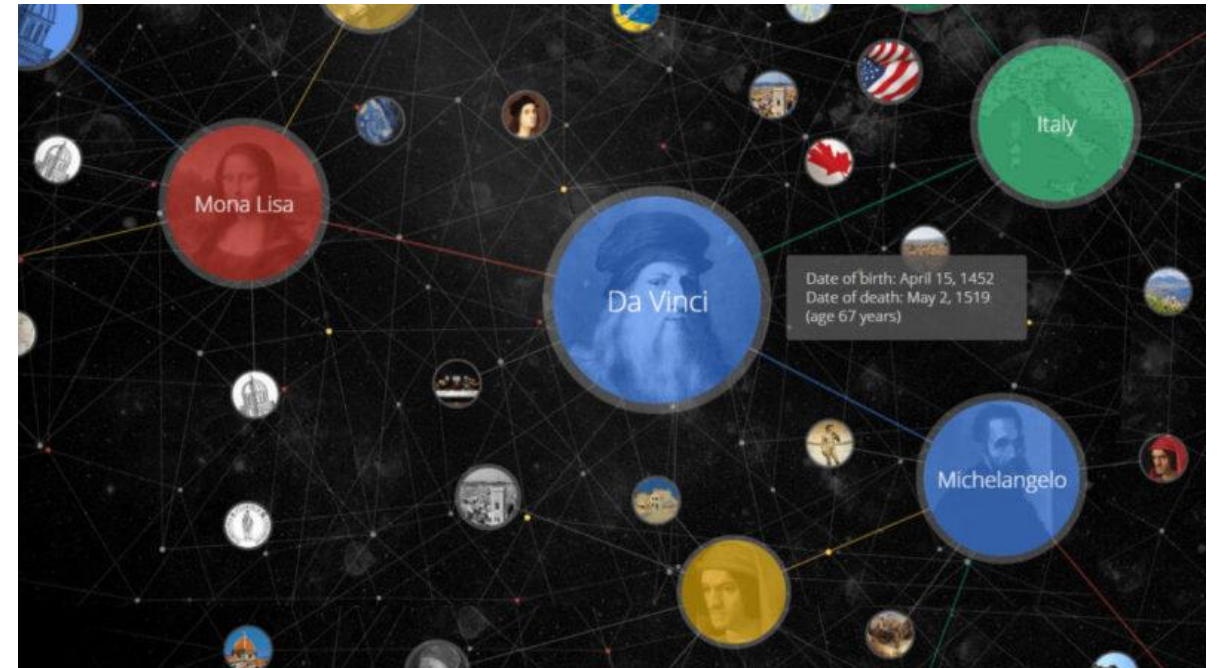
Base Model : **76,5%**

FT Model : **96,4%**

# Evaluating our Pipeline: Finding efficient Metrics

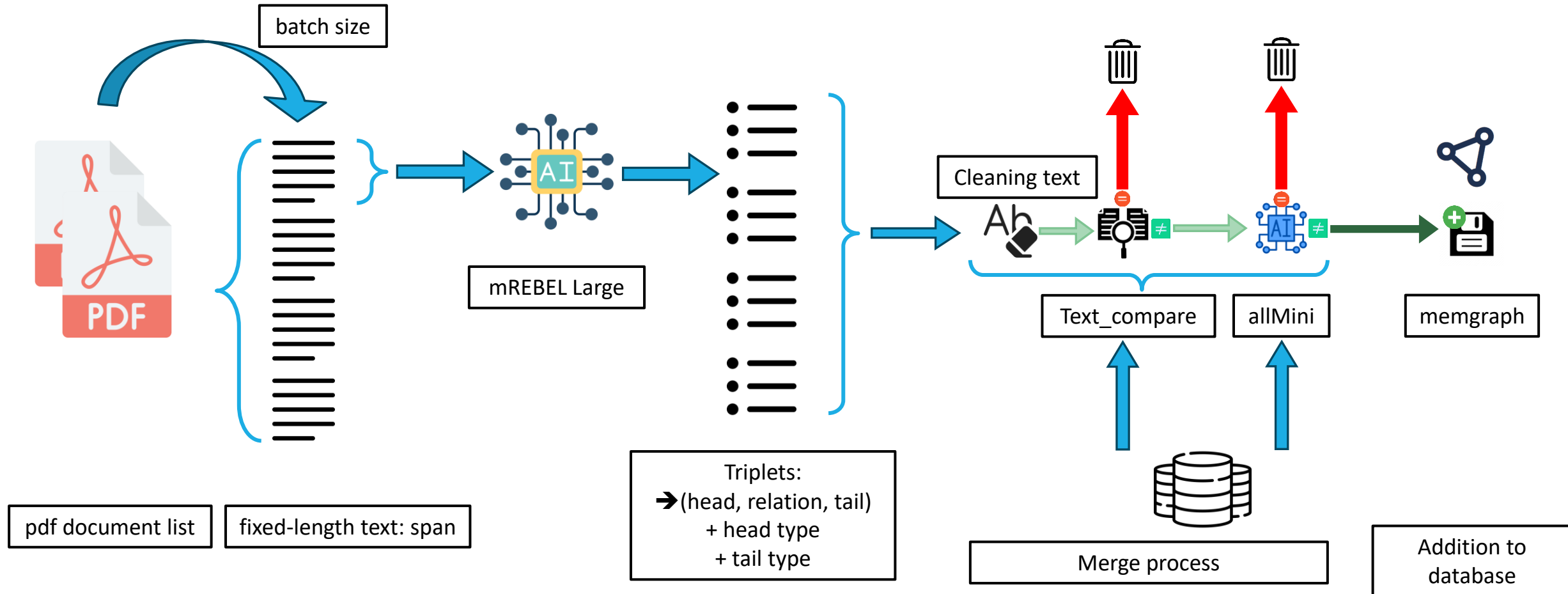
## Unsupervised Metrics:

- **Mean Rank (MR):** Reflects the average ranking of correct semantic triples in a list sorted by their likelihood. *A lower number indicates better performance.*
- **Hits@K:** Measures how often the correct semantic triples are ranked within the top K positions in the list. *Higher values mean more accurate predictions.*
- **Mean Reciprocal Rank (MRR):** Focuses on the top results and is less dependant to misleading information.



Google Knowledge Graph

# Detailed pipeline implementation





# Storing the Knowledge Base

Multiple ways



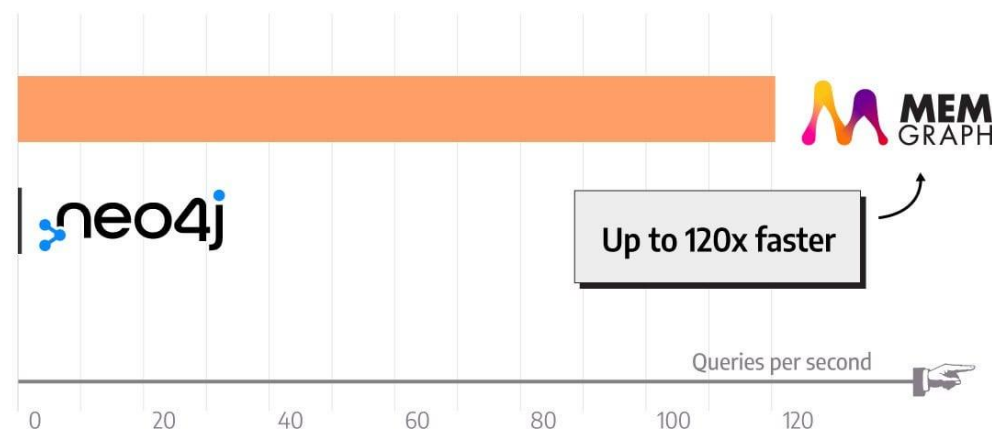
## PROS :

- Fast write to file
- Accessibility for file edit

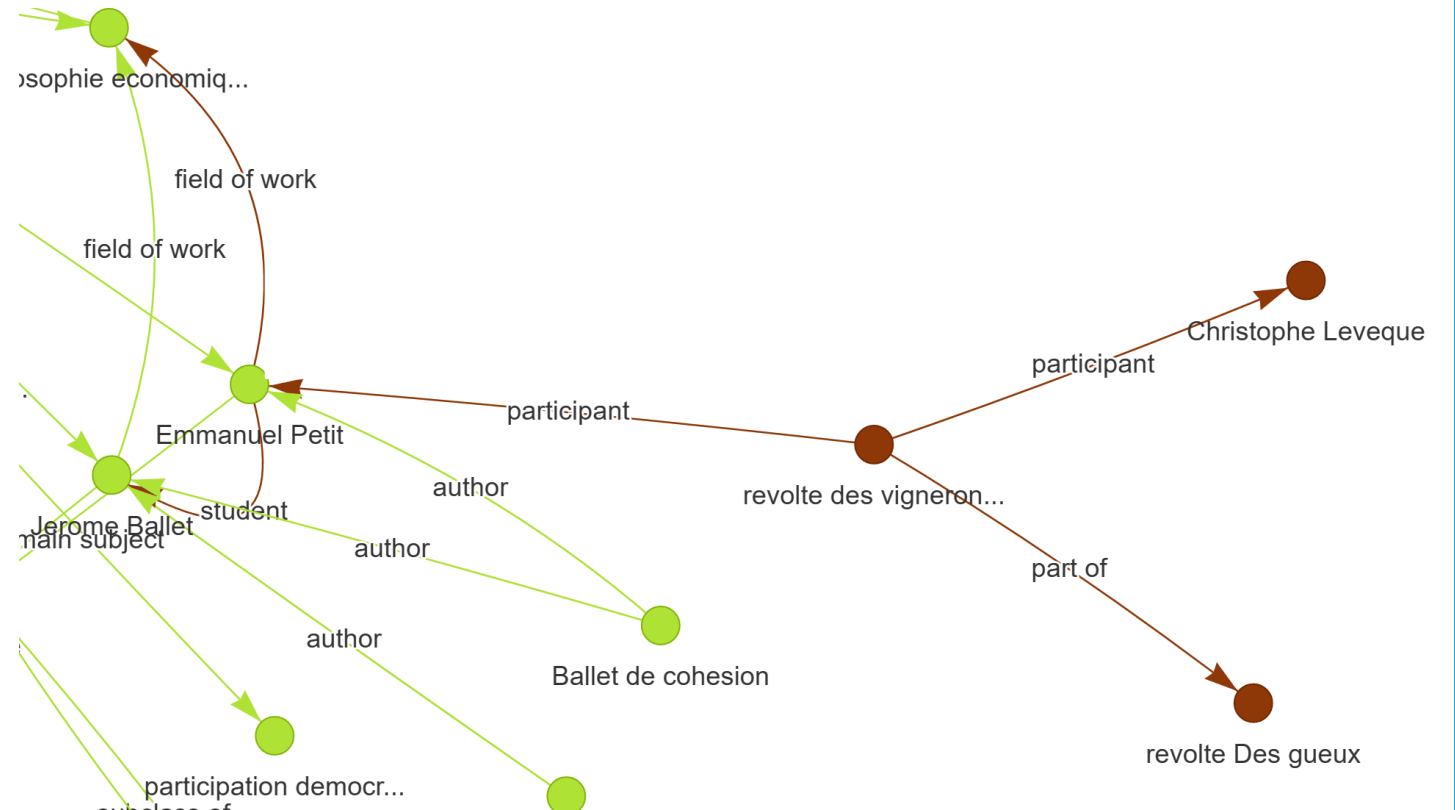
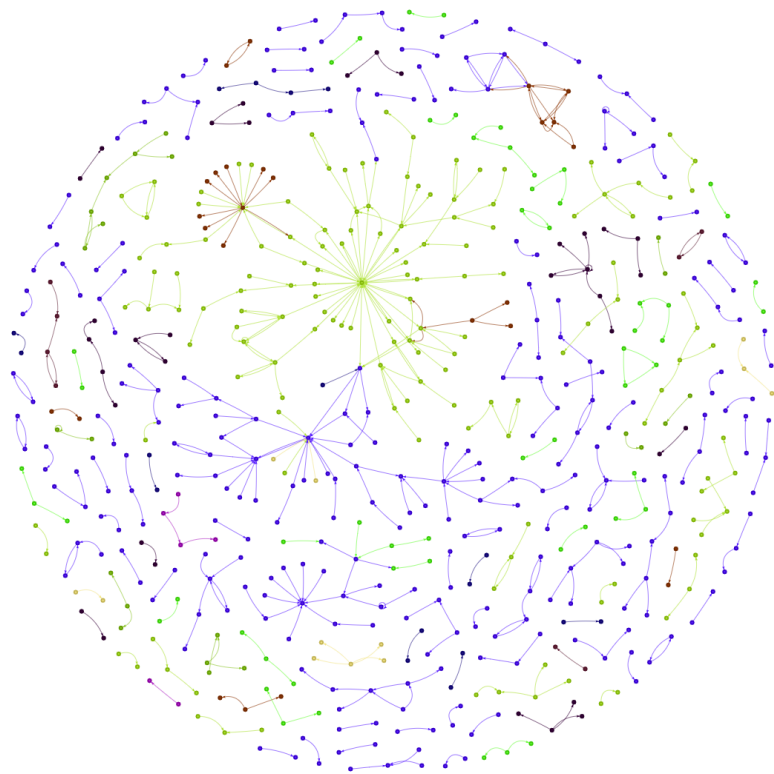
## CONS :

- Slow read of data
- Security Issues
- Stability Issues
- Computationally expensive

Throughput



memgraph-vs-neo4j-performance-benchmark-comparison

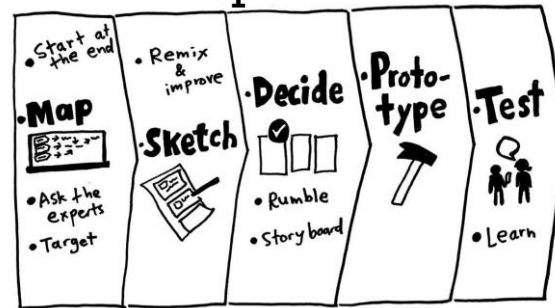


# Project Management

## Kanban



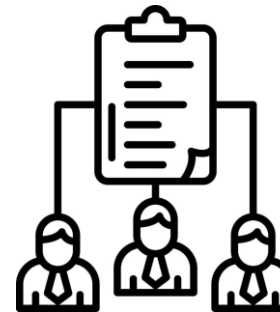
## sprint



## Weekly meeting - Friday



Discord

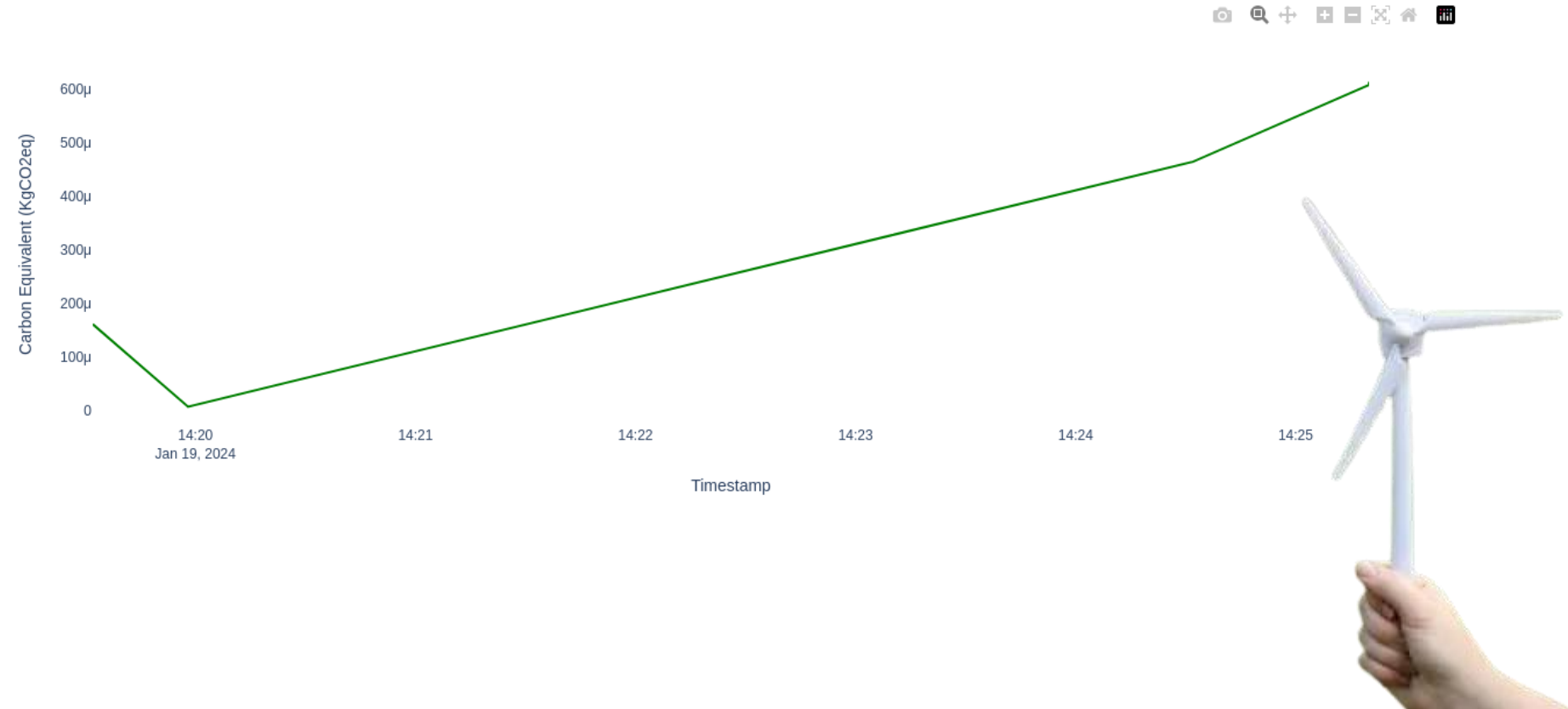
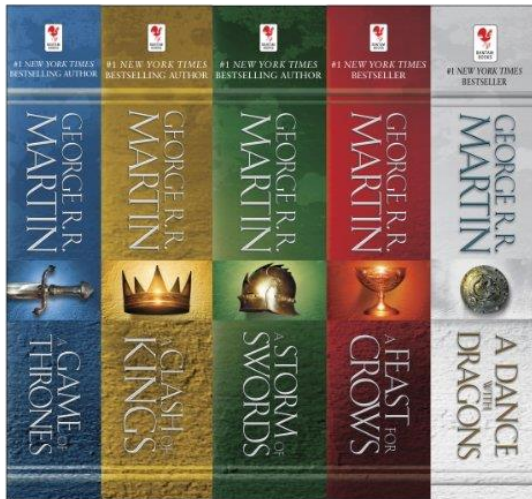


- Optimize and match our work with the client's wishes.
- Quality over quantity concept and Divide and conquer model

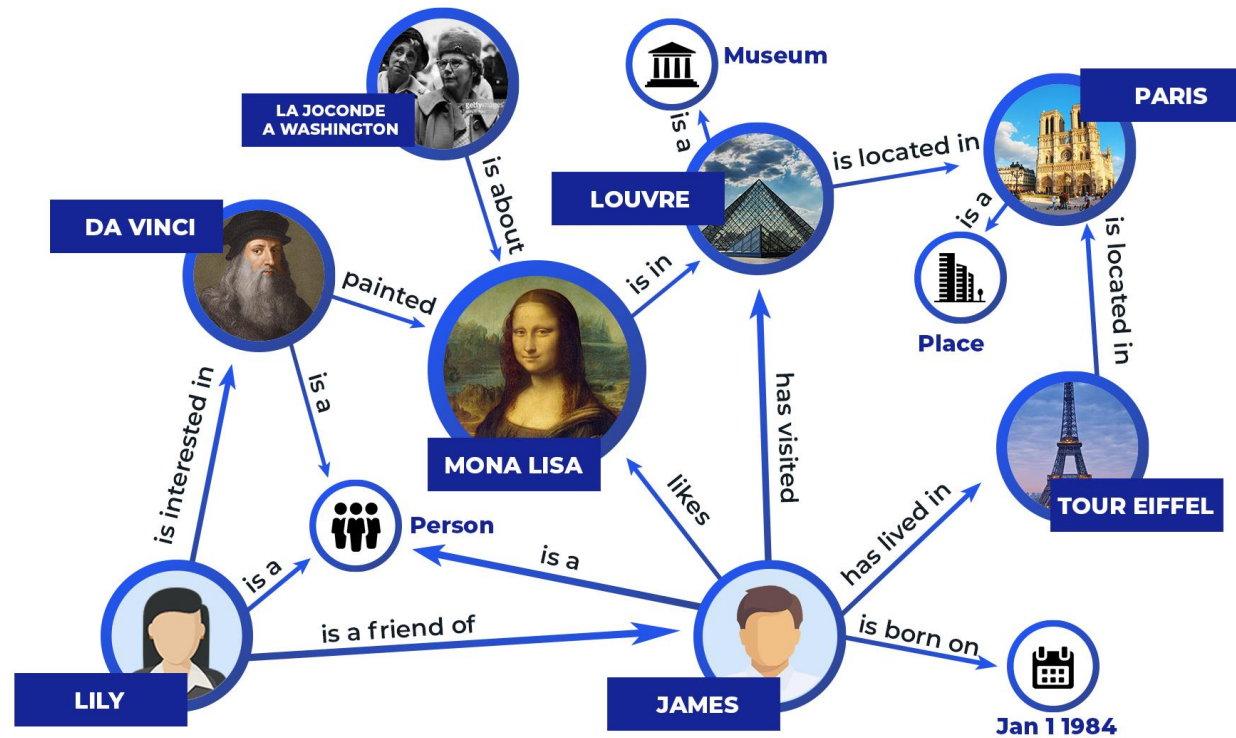
# Environmental study



## Emissions Timeline

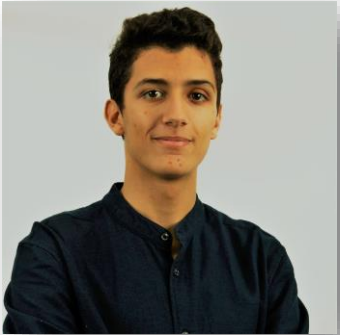


# Conclusion



# Team

---



**Jad El Karchi**  
AI & LLM  
student-engineer  
@ Enseirb-Matmeca



**Félicien Fichet**  
AI & OPTIMIZATION  
student-engineer  
@ Enseirb-Matmeca



**Olha Nahorna**  
Research Engineer &  
Project Leader  
@ BSE



**Mohamed Seddiq Elalaoui**  
AI & Data Science  
Student-engineer  
@ Enseirb-Matmeca



**Joseph Beasse**  
AI & NLP  
Student-engineer  
@ Enseirb-Matmeca

# Q & A

---

Thank you for your attention !

