Final presentation

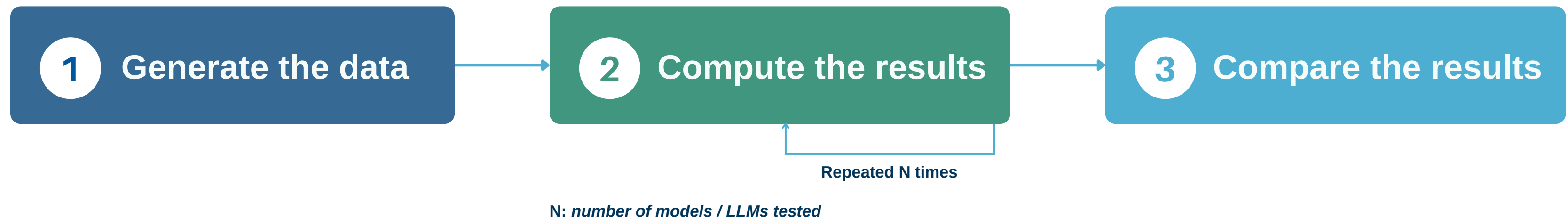# Outils Formels Avancés 2024

## AI fact checker

28 Mai 2024

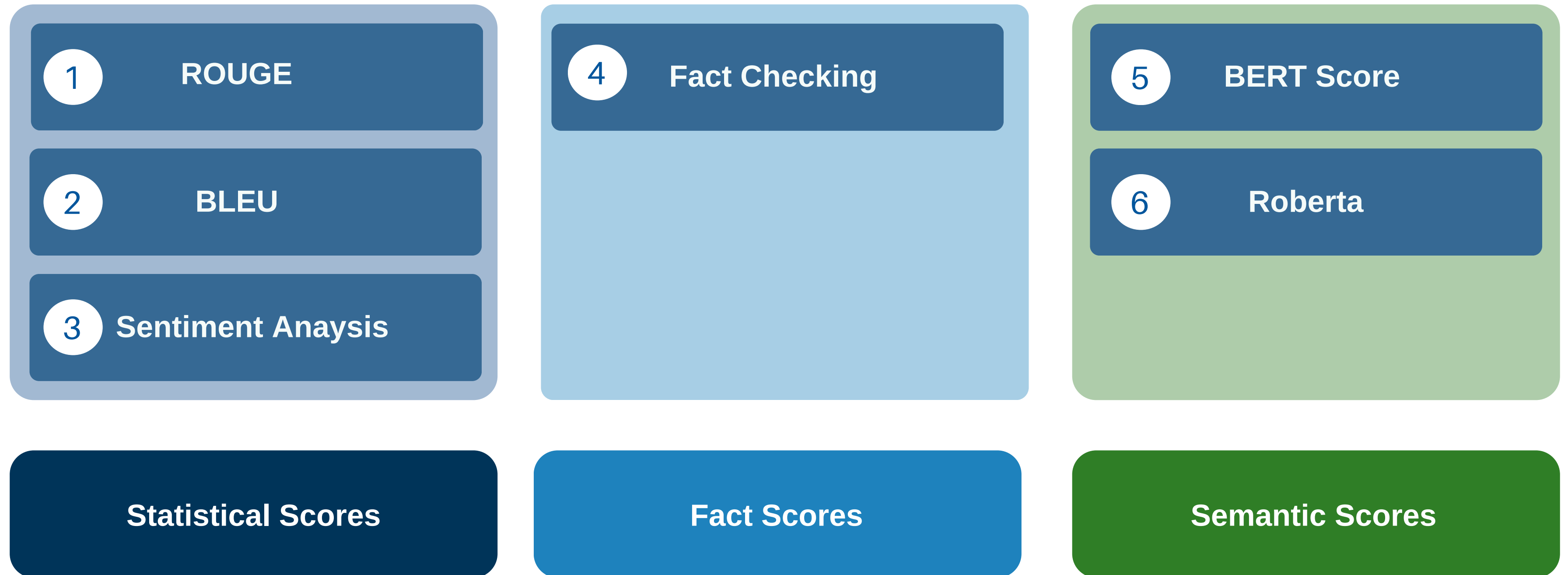# GOAL of the project

## AI Fact Checker

The main goal is to use the programming logic to verify the truth of the LLM's answers.

- Find references

- Find non-factual answers from LLMs

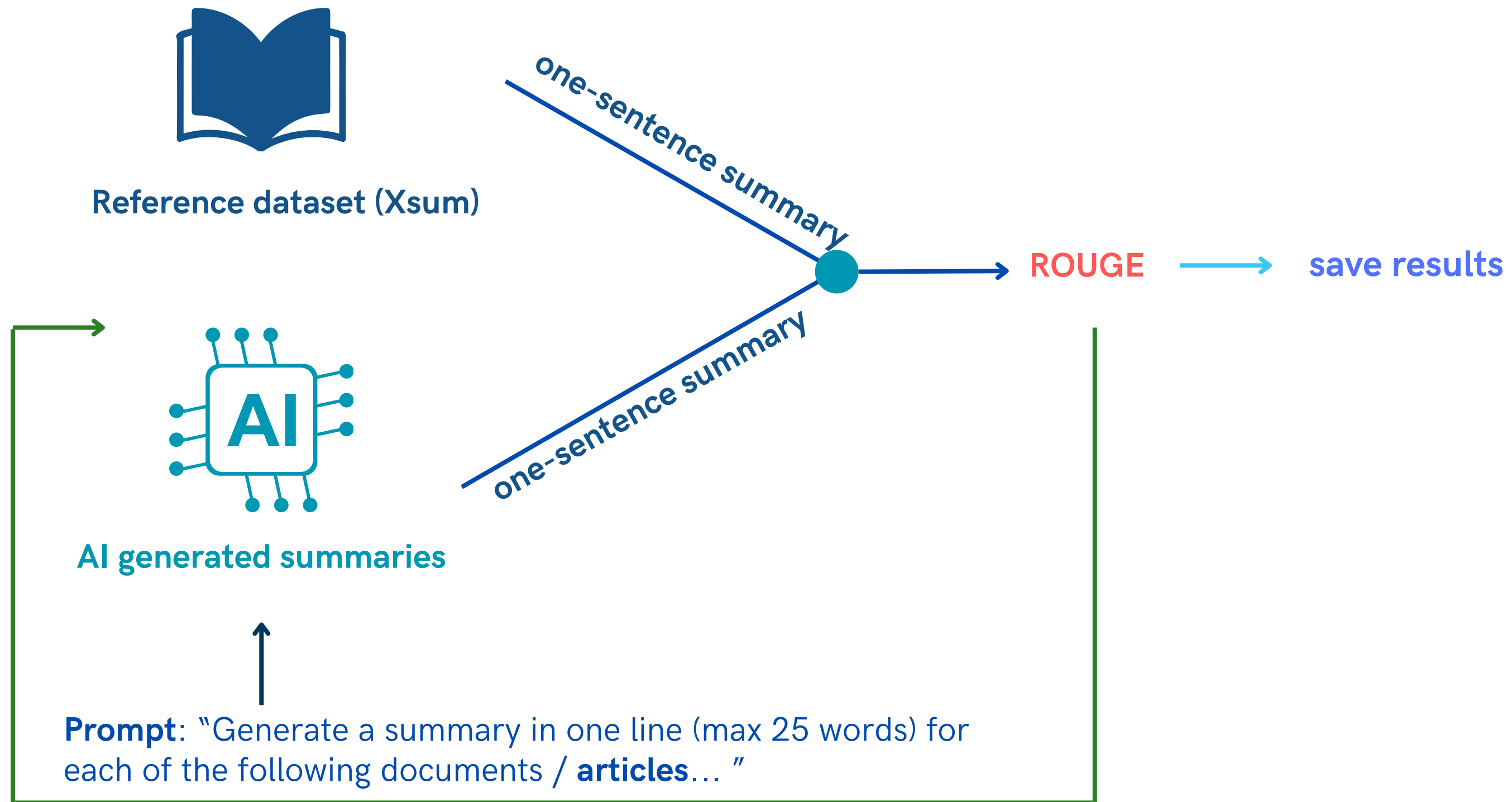- Define tests that measure the performance of LLMs

# Pipeline



| 1 | **Generate the data** |

| 2 | **Compute the results** |

| 3 | **Compare the results** |

**Repeated N times**

**N:** *number of models / LLMs tested*

# Evaluation Framework: Metrics

| 1 | ROUGE |
|---|-------|

| 2 | BLEU |
|---|------|

| 3 | Sentiment Anaysis |
|---|-------------------|

| 4 | Fact Checking |
|---|---------------|

| 5 | BERT Score |
|---|------------|

| 6 | Roberta |
|---|---------|

**Statistical Scores**

**Fact Scores**

**Semantic Scores**

Reference dataset (Xsum)

one-sentence summary

ROUGE → save results

AI generated summaries

one-sentence summary

**Prompt**: "Generate a summary in one line (max 25 words) for each of the following documents / **articles**… "

Evaluating 10 articles
Unigram and LCS-gram

Reference dataset

one-sentence translation

AI

Generated Translation

one-sentence translated

BLUE → save results

**Prompt**: "Generate a summary in one line (max 25 words) for each of the following documents / **articles**… "

## Evaluating 10 articles

Average of precision of 1-gram, 2-gram,3-gram and 4-gram

# 3. Sentiment Analysis

IMDB review dataset

Reference labels

Review classification

Review classification

Similarity check

$$\frac{\text{Num correct}}{\text{Total \# Review}}$$

save results

AI predicted labels

**Prompt**: "classify the following 10 sentences: positive, negative"

**Evaluating 10 random reviews**

# 4. Fact Checker

**3** Logic fact verification

**4** Answer output

SPARQL and wikidata

True = +1

False = 0

**Repeated with 10 different questions**

# 5. BERT Score

Reference dataset (Xsum)

one-sentence summary

AI generated summaries

one-sentence summary

BERT SCORE → save results

**Prompt**: "Generate a summary in one line (max 25 words) for each of the following documents / **articles**… "

**Evaluating 10 articles**

# 6. Roberta

Semantic Scores

Reference dataset (Xsum)

one-sentence summary

Roberta

save results

AI generated summaries

one-sentence summary

**Prompt**: "Generate a summary in one line (max 25 words) for each of the following documents / **articles**… "

**Evaluating 10 articles**

# Computing results

| | | | | | |
|---|---|---|---|---|---|
| **1** ROUGE | **2** BLEU | **3** Sentiment Anaysis | **4** Fact Checking | **5** BERT Score | **6** Roberta |

# Comparing LLMS



GPT-4    Large    Claude 2    Gemini Pro    GPT-3.5    LLaMA 2 70B

| Model | MMLU | Common sense and reasoning | | | | Knowledge | |
| | | HellaS | WinoG | Arc C (5) | Arc C (25) | TriQA | TruthfulQA |
| --- | --- | --- | --- | --- | --- | --- | --- |
| Mistral Large | 81.2% | 89.2% | 86.7% | **94.2%** | 94.0% | 82.7% | **50.5%** |
| LLaMA 2 70B | 69.9% | 87.1% | 83.2% | 86.0% | 85.1% | 77.6% | 44.7% |
| GPT 3.5 | 70.0% | 85.5% | 81.6% | 85.2% | 85.2% | – | – |
| GPT 4 | **86.4%** | **95.3%** | **87.5%** | – | **96.3%** | – | – |
| Claude 2 | 78.5% | – | – | 91.0% | – | **87.5%** | – |
| Gemini Pro 1.0 | 71.8% | 84.7% | – | – | – | – | – |

# Comparing LLMS



Scores Per Metric for Each Model

| Models | Rouge1 | RougeL | Bleu | Sentiment | Fact checking | Bert Score | RoberTa |
|---|---|---|---|---|---|---|---|
| ChatGPT_3.5 | 0.32 | 0.20 | 0.09 | 0.90 | 1.00 | 0.88 | 0.62 |
| ChatGPT_4-o | 0.33 | 0.21 | 0.09 | 0.90 | 0.80 | 0.89 | 0.65 |
| ChatGPT_4 | 0.27 | 0.16 | 0.09 | 0.90 | 0.80 | 0.88 | 0.51 |
| Meta_Lama_3_70B | 0.30 | 0.20 | 0.09 | 0.90 | 0.80 | 0.88 | 0.49 |
| Mistral_Large | 0.21 | 0.15 | 0.08 | 0.80 | 0.60 | 0.86 | 0.56 |
| Mistral_small | 0.26 | 0.17 | 0.08 | 0.80 | 0.70 | 0.85 | 0.54 |
| Mistral_Next | 0.32 | 0.18 | 0.08 | 0.80 | 0.80 | 0.87 | 0.62 |
| Blackbox_ai | 0.24 | 0.19 | 0.09 | 0.90 | 0.80 | 0.87 | 0.60 |

Metrics

# Findings and improvements

## Findings

- Different metrics are interesting

- Weighting the metrics was not useful in the end

- Variability in Prompt Responses

## Improvements for testing

- Adaptive Prompt Engineering
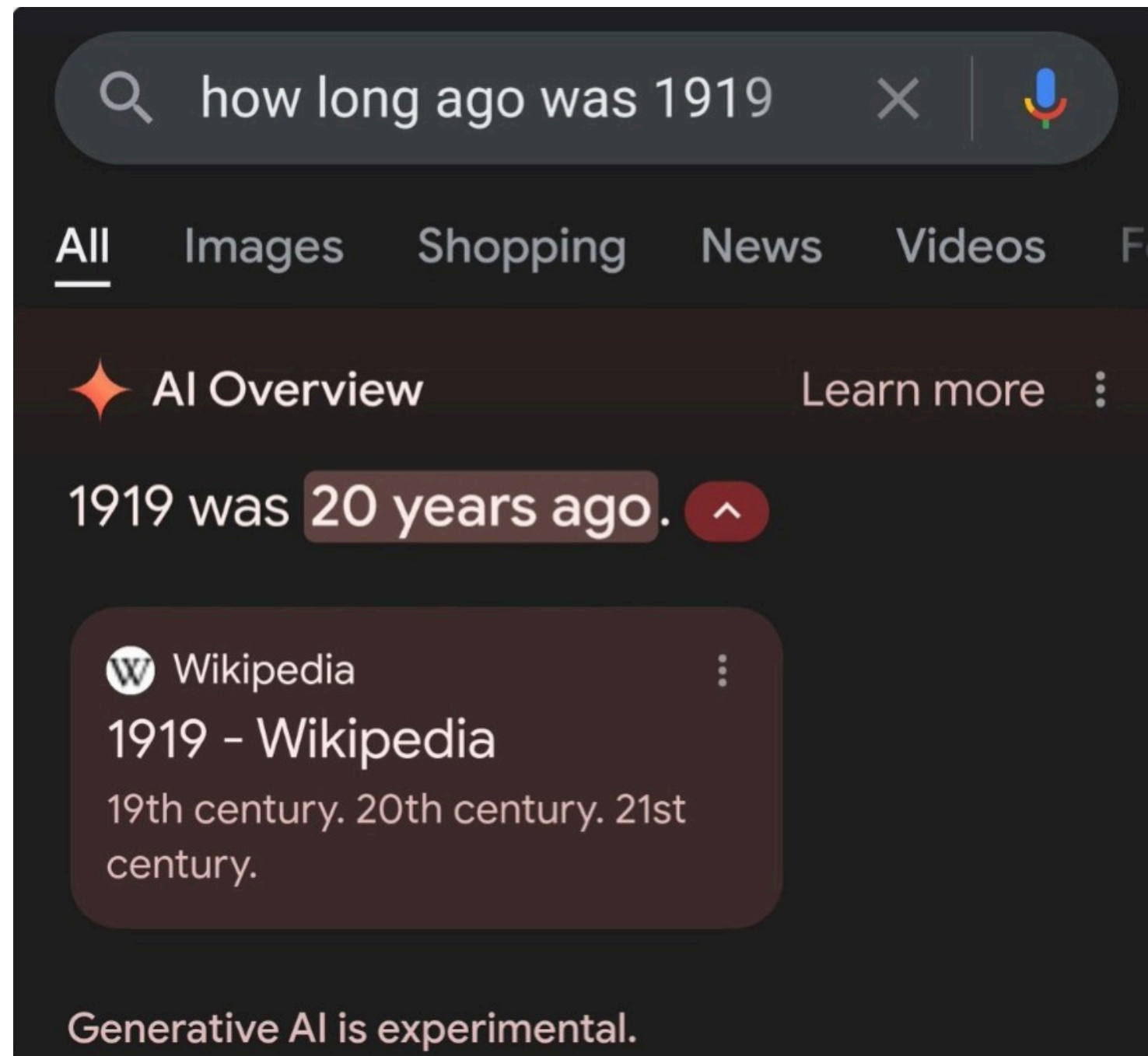
# Findings and improvements

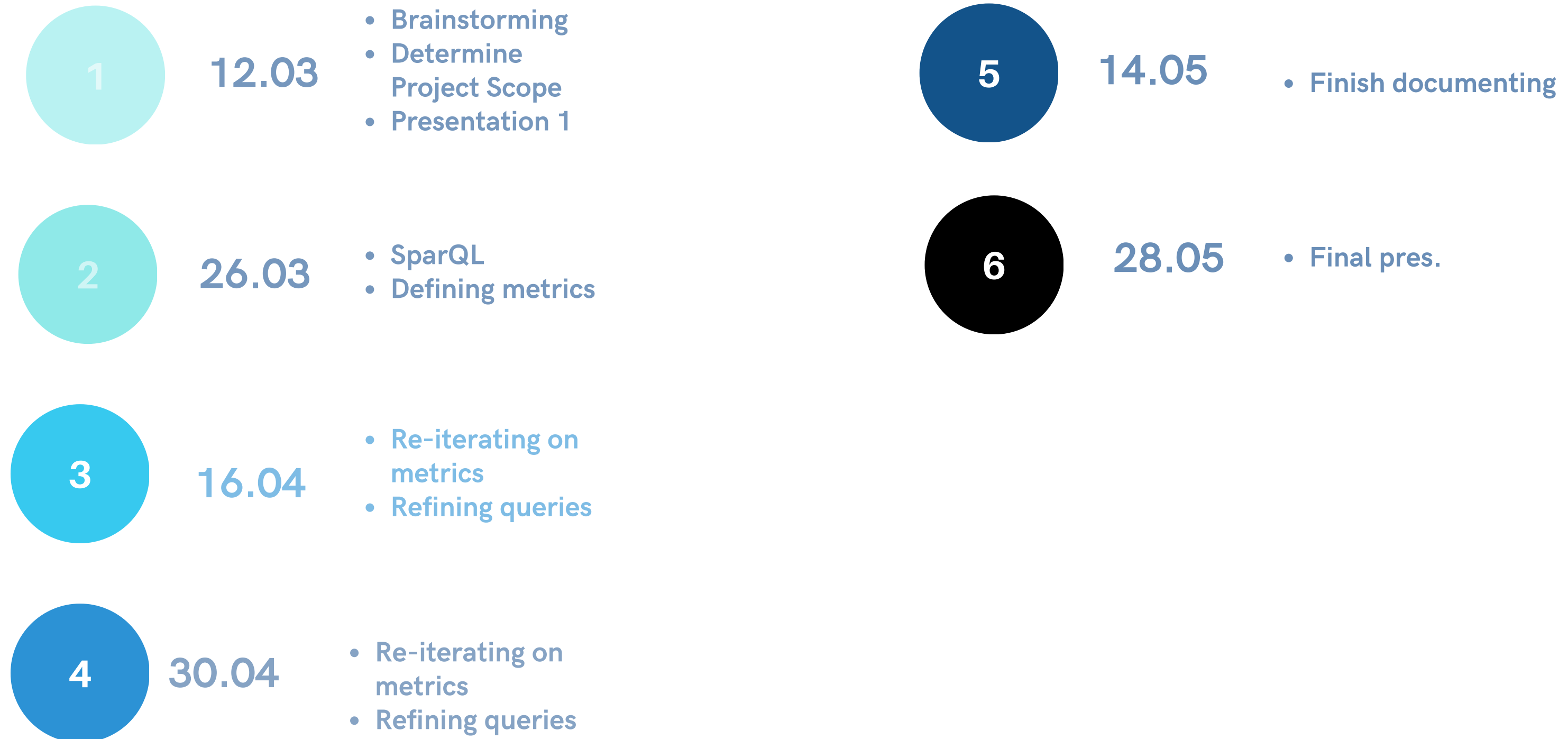# Improvement suggestion for LLMs

## Improvements

- Post-process on answers

- Iterative answer generation

- Fine-tune on references

- Get them access to do calculation

# State of the art :)

# Tentative Planning

**1**    12.03
- Brainstorming
- Determine Project Scope
- Presentation 1

**2**    26.03
- SparQL
- Defining metrics

**3**    16.04
- Re-iterating on metrics
- Refining queries

**4**    30.04
- Re-iterating on metrics
- Refining queries

**5**    14.05
- Finish documenting

**6**    28.05
- Final pres.

# Final Presentation

End for semester!