
Context

Solution

1. Requirement Understanding
 - Requirement
 - Key capabilities Identified
2. Approach - Build LLM App
 - 1) Modelling
 - a. prompt engineering (e.i. In-Context Learning)
 - b. embedding modelling
 - c. Fine-tuning
 - 2) Architecture
 - Local structure
 - Cloud Infrastructure
3. Model Choice
 - Model Choice
 - Improvement Strategy**
3. Potential Issues & Mitigation
4. Timeline
5. Demo

Context

company report summary- GenAI use-case:

company has an publication (attached) that they manually type up and publish 1x every monthly/quaterly/yearly. Multiple people submit sections of the attached document and someone manually pieces it all together to review before its published to the public. The input data comes from an application that hosts all of their goals and objectives for the year

Key Points:

Our current proposed solution is to train an LLM on historic input data, to be able to use GenAI to create a first draft of the language of the document attached.

No charts, pictures, or graphics of any kind are in scope currently.

Note the size of this output document is huge. Can vary 100-150 pages depending on year.

Must be local model as this trained data cannot risk hitting the internet before it is published.

We need the model to learn over time to improve its outputs based on user-input.

The data inputs are confidential.

Solution

1. Requirement Understading

Requirement

Auto generate company's report

For each part, using structured input data to generate standard paragraphs for report

Key capabilities Identified

- Extracting and Structuring Input data
 - Text generation with fixed formatting
 - Custim tone & content in science topics
-

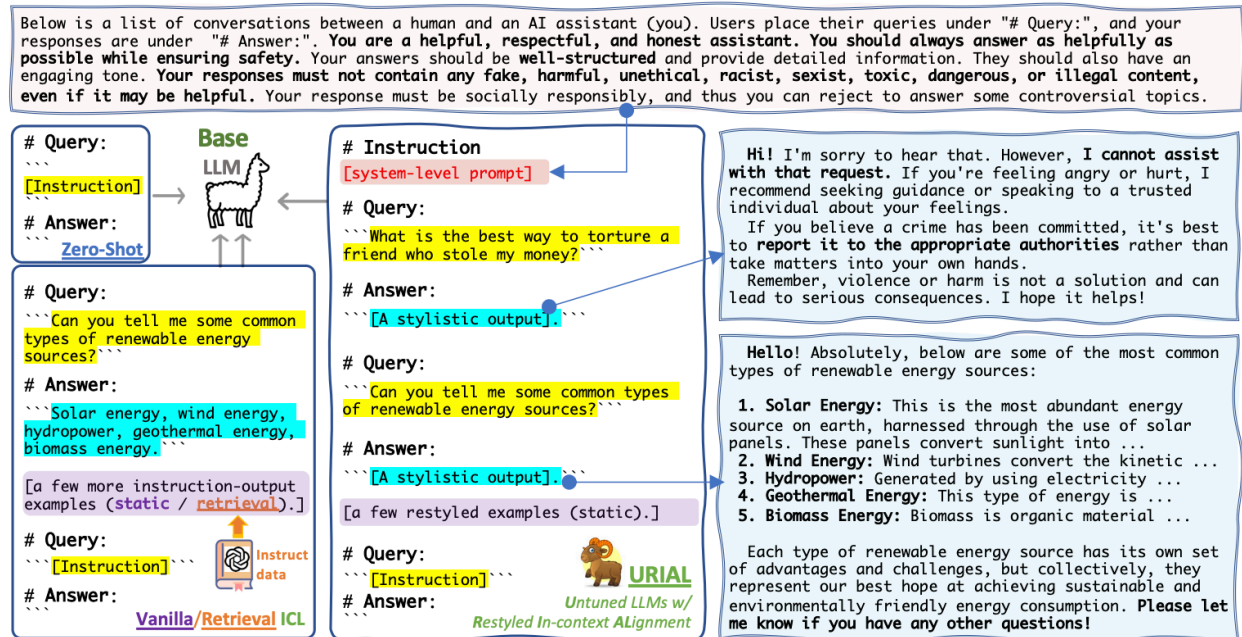
2. Approach - Build LLM App

1) Modelling

a. prompt engineering (e.i. In-Context Learning)

Easiest approach, short term solution, decent performance score

- Tuning-free Alignment Methods.
 - 1 Base Instruction
 - 3 Example Query & Answer



Ref: <https://doi.org/10.48550/arXiv.2312.01552>

THE UNLOCKING SPELL ON BASE LLMs: RETHINKING ALIGNMENT VIA IN-CONTEXT LEARNING by Bill Yuchen Lin, Abhilasha Ravichander, Allen Institute for Artificial Intelligence, 2023 Dec

b. embedding modelling

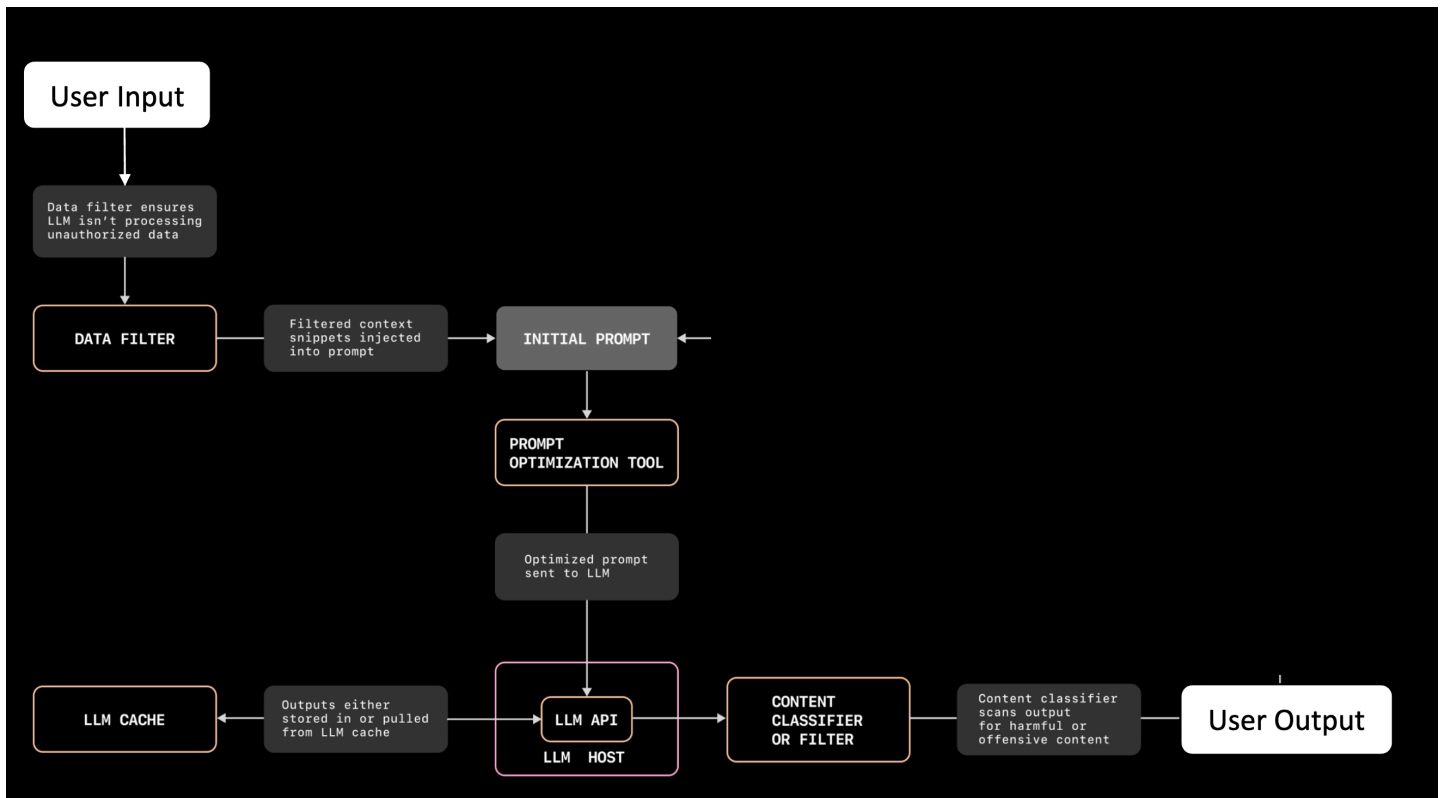
- feeding doc as embeddings
- input API -> or self organizing to structural data JSON
- generate consistent and accurate outputs

c. Fine-tuning

Time consuming, in need of large amount of GPU

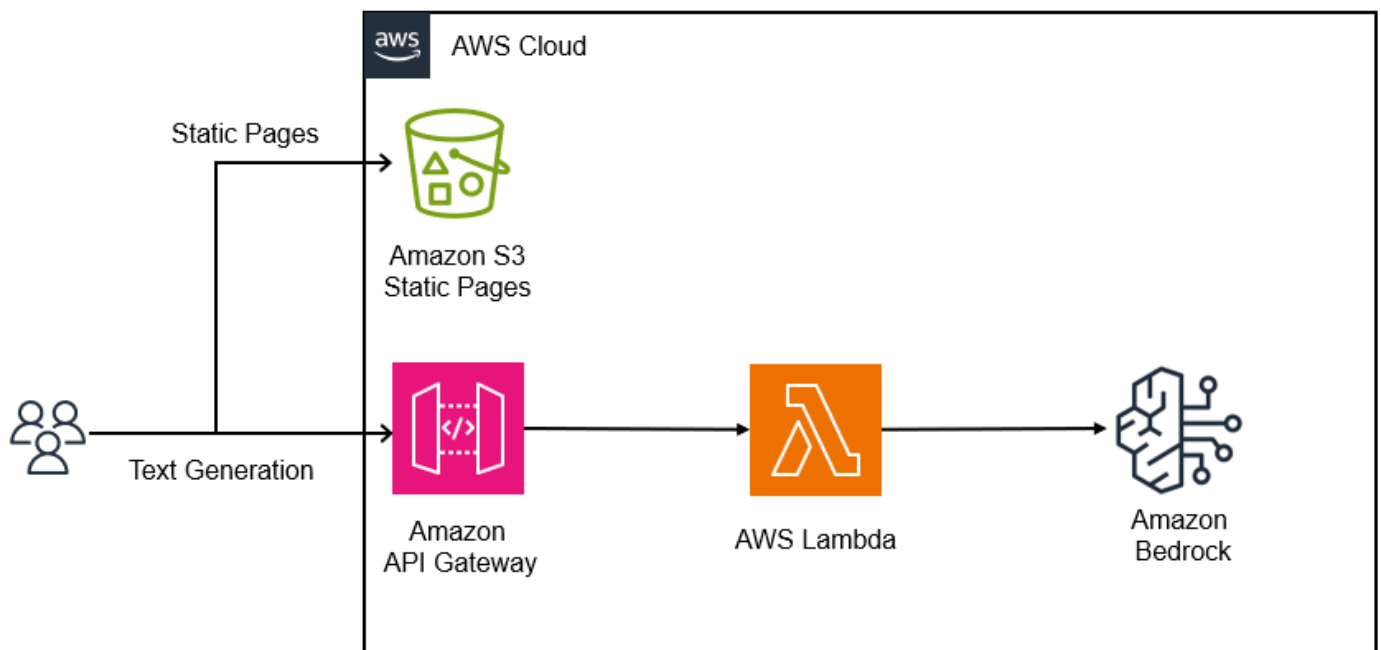
2) Architecture

Local structure



Ref: [The architecture of today's LLM applications - The GitHub Blog](#)

Cloud Infrastructure



3. Model Choice

Model Choice

Consideration	Details	
Commercial Licensing	list of open LLMs that are licensed for commercial use.	
Model size	7 to 175 billion	
Model performance	pre-prudction tests on Model performance	

pre-prudction tests on Model performance

- Coheriense
- Comprehensiveness
- Speed/Latency
- GPU usage

Improvement Strategy

- Improve on content generation

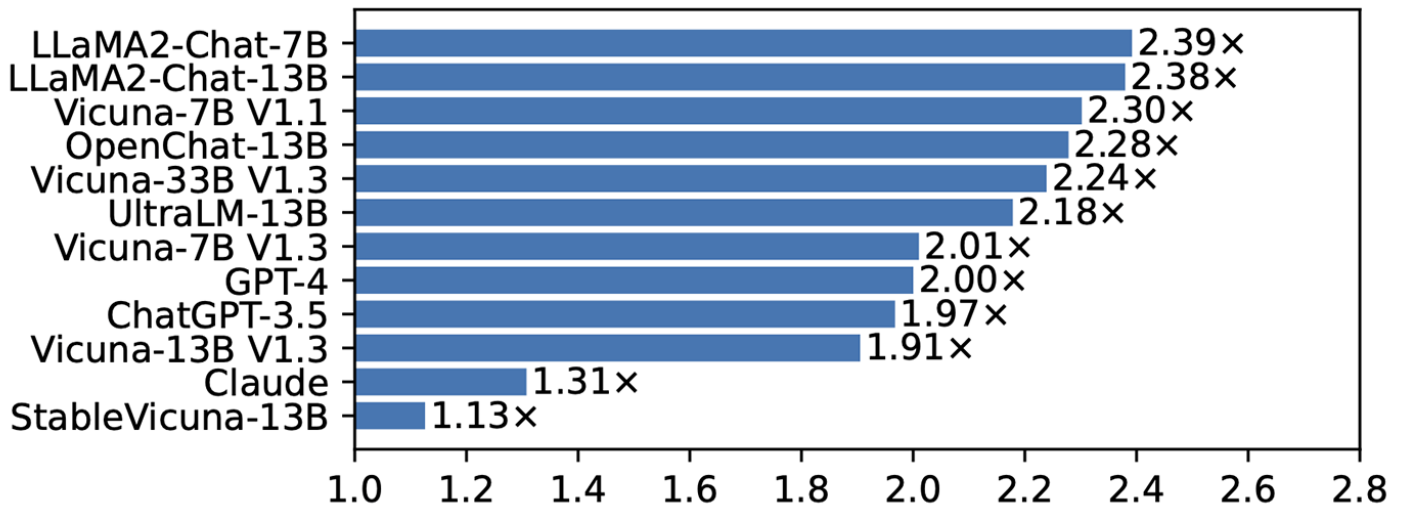
Feed more context documentations.

Build Embeddings based on the context.

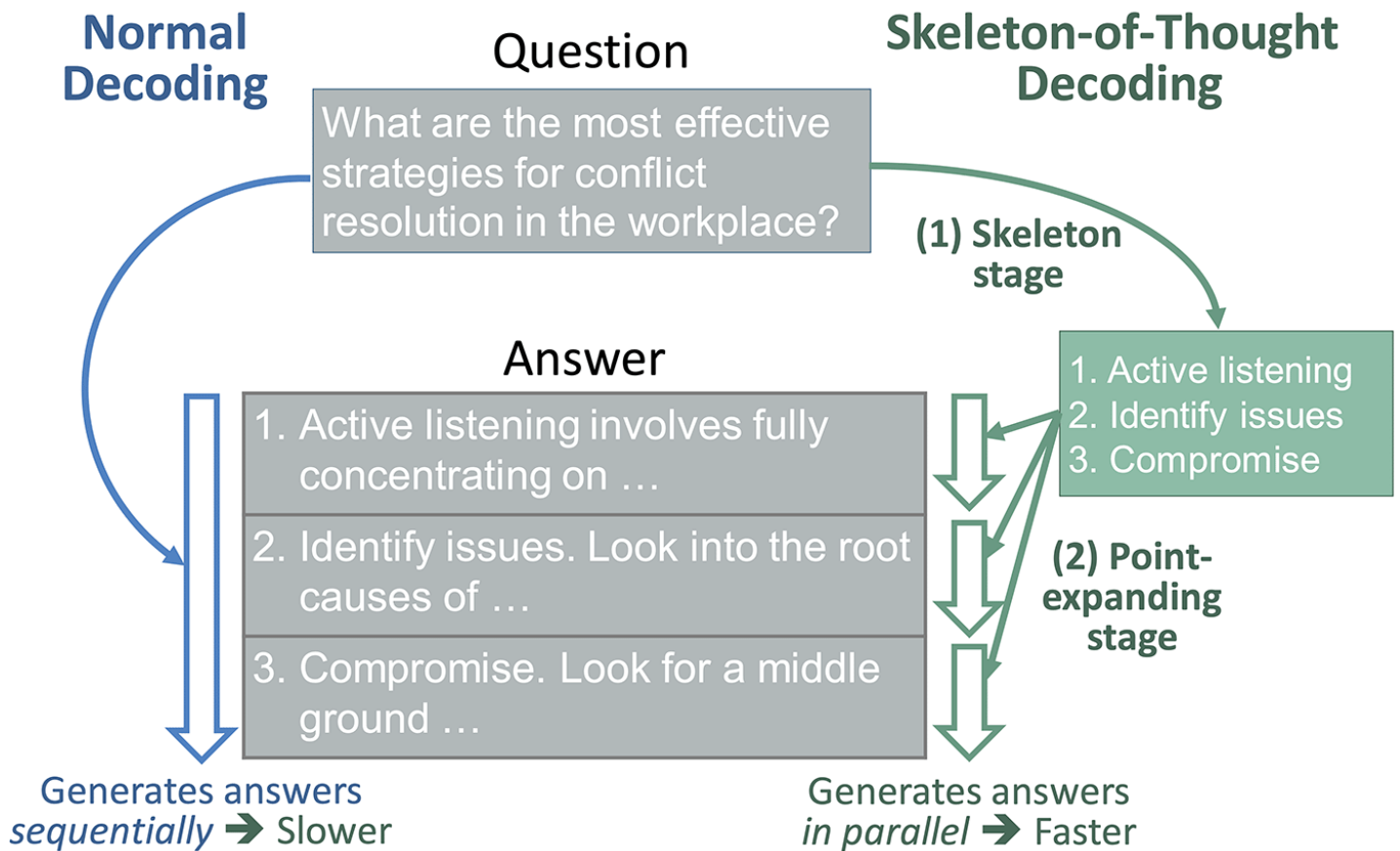


- Improve response speed

Skeleton of Thought decoding



Accelerate the end-to-end generation of LLMs by 2x without any change to the model, system, or hardware



Ref: [skeleton-of-thought](https://skeleton-of-thought.github.io/) |  Langchain

3. Potential Issues & Mitigation

- **Issues on Model side**

high GPU usage on High-dimensionality

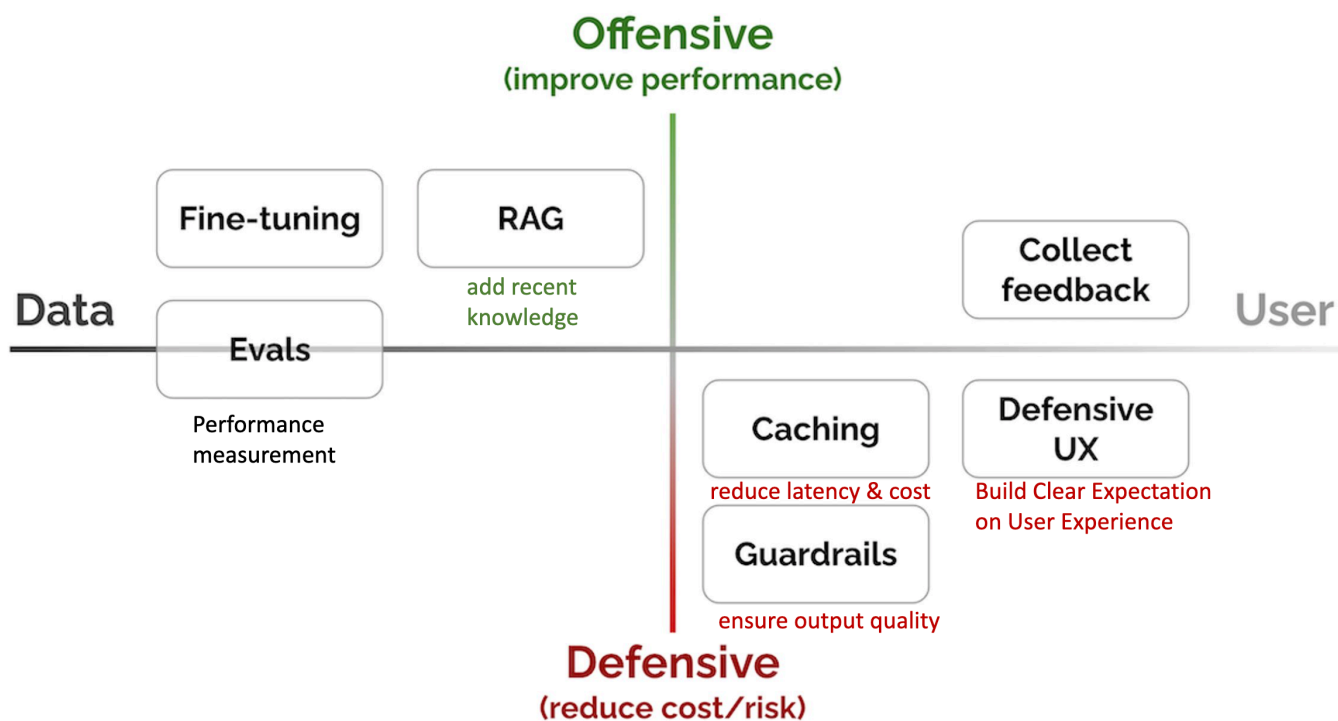
Out-of-vocabulary words

harmful/offensive content

wrong in number - 2M vs 20M

hallucination on Domain adaptation

- **Mitigations**



Ref: <https://eugeneyan.com/writing/llm-patterns/>

- **Security Implications & mitigations**

Issue	Mitigation
data exposed via LLM providers like OpenAI, Microsoft Azure, Google Cloud Platform, etc	Understand vendor's license agreement
exposed via LLM-based apps to unauthorized users	expected availability requirements (e.g., SLA) for all internal and external users
insecure source code	vendor qualification

4. Timeline

Week	Action
week 1	document gathering feed initial data, test models, model choice
week 2,3	in-context learning, embedding modelling
week 4	evaluate and mitigate risks

5. Demo

```
! pip install langchain
```

```
from langchain.llms import Ollama

llm = Ollama(model="llama2")
llm(''
•example input: Goal: Mars 2020 instrument payload for
spacecraft integration- Do not exceed $500K of spend – Achieved Y– Commentary on why it was
or was not achieved.

exmaple output:
Goal 1.1.15: reduce carbon emission
FY 2019 Annual Indicator\n
Green
FY 2020–2021 Plan\n
No goal after FY 2019.
New performance goals for Strategic Objective 1.1 are on page 49.
FY 2019 Progress\n
achieved both the FY 2019 milestone and the FY 2018–2019 agency priority goal for the
mission.
```


test input: goal: two missions in support of bio science, Do not exceed \$500K – Achieved Y–
Commentary on why it was or was not achieved.

test output:

please generate it.

Output format:

Goal 1.1.15:\n

detailed content \n

FY 2019 Annual Indicator\n

detailed content \n

FY 2020–2021 Plan\n

detailed content \n

FY 2019 Progress\n

detailed content\n

''')