# Symphony: Towards Natural Language Query Answering over Multi-modal Data Lakes (Vision)

**Zui Chen**
*Tsinghua*

**Zihui Gu**
*Renmin*

**Lei Cao**
*University of Arizona*
*MIT*

**Ju Fan**
*Renmin*

**Sam Madden**
*MIT*

**Nan Tang**
*QCRI*

I ♥ AMSTERDAM

**Jan 10, 2023 @ CIDR,**

# Multi-modal Data Lakes

**Text files, CSV files, databases, log files, knowledge graphs …**

# One Driving Application (PSA)

FIFA 2022
population
education
economic activities
disabilities
buildings
…
…

# One Driving Application (PSA)



FIFA 2022
population
education
economic activities
disabilities
buildings
…
…

# One Driving Application (PSA)

FIFA 2022
population
education
economic activities
disabilities
buildings
…
…

# Data Management Perspective

**Structured Queries (SQL)**

Data
Preparation

Extraction
Transformation
Loading
Cleaning
Schema matching
Deduplication
…

**High human-cost**

# Data Management Perspective

**Structured Queries (SQL)**

Data
Preparation

Extraction
Transformation
Loading
Cleaning
Schema matching
Deduplication

…

**High human-cost** ⟶ **No Data Preparation**

**Data lake
(raw format)**

# Data Management Perspective

**Structured Queries (SQL)** $\longrightarrow$ **Natural Language Queries**

Data
Preparation

Extraction
Transformation
Loading
Cleaning
Schema matching
Deduplication
…

**High human-cost** $\longrightarrow$ **No Data Preparation**

Data lake
(raw format)

# Data Management Perspective

**Structured Queries (SQL)** $\longrightarrow$ **Natural Language Queries**

**Symphony**

Reasoning 🤖

Retrieval

**Data Preparation**

Extraction
Transformation
Loading
Cleaning
Schema matching
Deduplication

...

**Data lake (raw format)**

**High human-cost** $\longrightarrow$ **No Data Preparation**

# Foundation Models?

**Foundation Models**
(GPT3, ChatGPT)

**Implicit knowledge**

**Symphony**

Reasoning

Retrieval

**Explicit knowledge**

# Foundation Models?

**Foundation Models**
(GPT3, ChatGPT)

**Implicit knowledge**

**(1) Poor (DB) reasoning**

**Symphony**

Reasoning

Retrieval

**Explicit knowledge**

# Foundation Models?

**Foundation Models**
(GPT3, ChatGPT)

**Implicit knowledge**

**(1) Poor (DB) reasoning**
**(2) Not up-to-date**

**Symphony**

Reasoning

Retrieval

**Explicit knowledge**

# Foundation Models?

**Foundation Models**
(GPT3, ChatGPT)

**Implicit knowledge**

**(1) Poor (DB) reasoning**
**(2) Not up-to-date**

**Question:** Who gave a talk about "Dutch tomatoes in CIDR 2023?"

**File**

The Netherlands has basically reinvented the tomato

By **Netherlands Foreign Investment Agency** | Published February 17, 2017

Gustavo Alonso keynote at CIDR'23

**Symphony**

Reasoning

Retrieval

**Explicit knowledge**

# A Running Example



(1) Index

# A Running Example



## (1) Index

**Q** Which **songs** appeared in a **film produced by Alankar Chitra** and **directed by Shanker Mukherjee?**

# A Running Example



**(1) Index**

**Q** Which **songs** appeared in a **film produced by Alankar Chitra** and **directed by Shanker Mukherjee?**

**(2) Data Retrieval**

Faraar (transl.Absconding) is a 1975 Bollywood crime film drama. The film is produced by Alankar Chitra and directed by Shanker Mukherjee. The film stars Amitabh Bachchan, Sharmila Tagore, Sanjeev Kumar, Sulochna, Sajjan, Agha and Bhagwan Dada…

Source: https://en.wikipedia.org/wiki/Faraar **P1**

| Year | Song | Film | **T1** |
|------|------|------|---|
| 1971 | Zindagi Ek Safer | Andaz | … |
| 1971 | Yeh Jo Mohabbat | Kati Patang | … |
| 1975 | Main Pyaasa tum | Faraar | … |
| … | … | … | … |

Source: https://en.wikipedia.org/wiki/Kishore_Kumar

6

# A Running Example



## (1) Index

**Q** Which **songs** appeared in a **film produced by Alankar Chitra** and **directed by Shanker Mukherjee?**

## (2) Data Retrieval

Faraar (transl.Absconding) is a 1975 Bollywood crime film drama. The film is produced by Alankar Chitra and directed by Shanker Mukherjee. The film stars Amitabh Bachchan, Sharmila Tagore, Sanjeev Kumar, Sulochna, Sajjan, Agha and Bhagwan Dada…
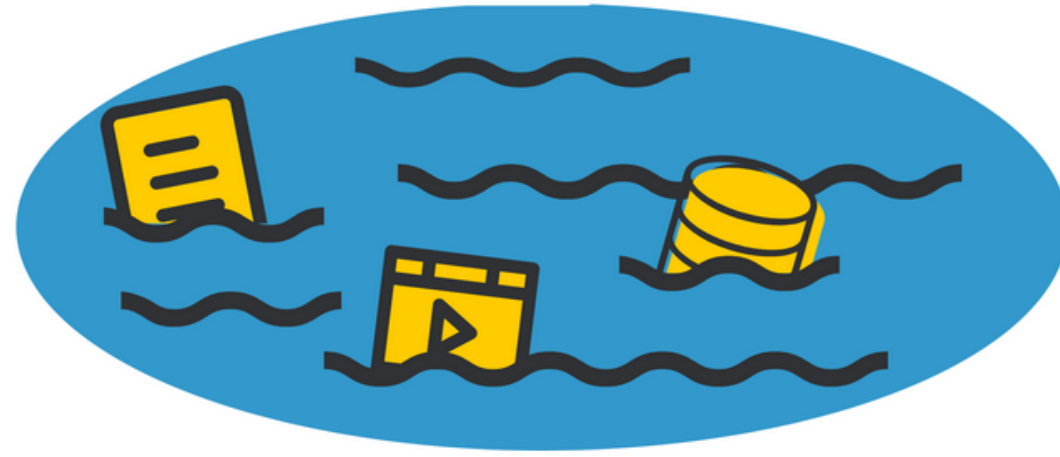
**P1**

Source: https://en.wikipedia.org/wiki/Faraar

**T1**

| Year | Song | Film | … |
|------|------|------|---|
| 1971 | Zindagi Ek Safer | Andaz | … |
| 1971 | Yeh Jo Mohabbat | Kati Patang | … |
| 1975 | Main Pyaasa tum | Faraar | … |
| … | … | … | … |

Source: https://en.wikipedia.org/wiki/Kishore_Kumar

## (3) Query Decomposition

**Prompt 1:** **The passage P1** has the following content: …
**The table T1** has the following columns: Year, Song, Film, Music Director, Lyricist.
**Based on P1 and T1**, the question is "Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?".
**What sub-questions** can it be broken down into?

**GPT-3:** **q1** What is the name of the film produced by Alankar Chitra and directed by Shanker Mukherjee. It can be answered by P1.

**Prompt 2:** The first sub-question is "What is the name of the film produced by Alankar Chitra and directed by Shanker Mukherjee?", it can be answered by P1.

**GPT-3:** **q2** the second sub-question is "What is the name of the song in the film?", it can be answered by T1.

6

# A Running Example



**(1) Index**

**Q** Which **songs** appeared in a **film produced by Alankar Chitra** and **directed by Shanker Mukherjee?**

**(2) Data Retrieval**

Faraar (transl.Absconding) is a 1975 Bollywood crime film drama. The film is produced by Alankar Chitra and directed by Shanker Mukherjee. The film stars Amitabh Bachchan, Sharmila Tagore, Sanjeev Kumar, Sulochna, Sajjan, Agha and Bhagwan Dada...

Source: https://en.wikipedia.org/wiki/Faraar

**P1**

| Year | Song | Film | **T1** |
|------|------|------|------|
| 1971 | Zindagi Ek Safer | Andaz | ... |
| 1971 | Yeh Jo Mohabbat | Kati Patang | ... |
| 1975 | Main Pyaasa tum | Faraar | ... |
| ... | ... | ... | ... |

Source: https://en.wikipedia.org/wiki/Kishore_Kumar

**(3) Query Decomposition**

**Prompt 1:** **The passage P1** has the following content: …
**The table T1** has the following columns: Year, Song, Film, Music Director, Lyricist.
**Based on P1 and T1**, the question is "Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?".
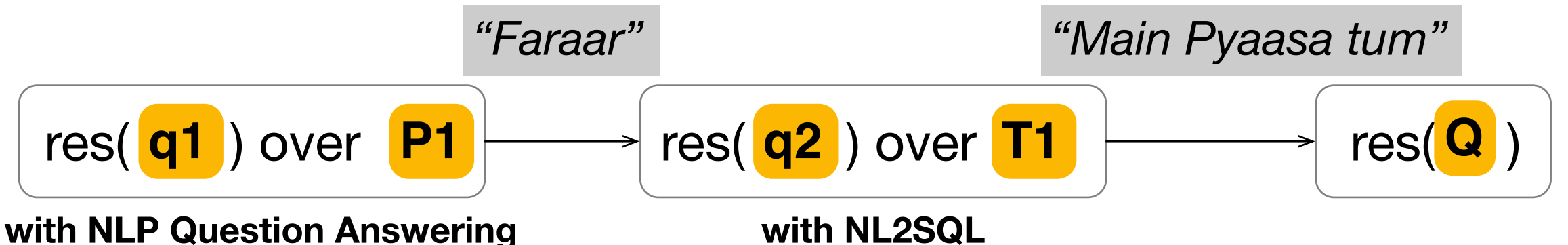**What sub-questions** can it be broken down into?

**GPT-3:** **q1** What is the name of the film produced by Alankar Chitra and directed by Shanker Mukherjee. It can be answered by **P1**.

**Prompt 2:** The first sub-question is "What is the name of the film produced by Alankar Chitra and directed by Shanker Mukherjee?", it can be answered by P1.

**GPT-3:** **q2** the second sub-question is "What is the name of the song in the film?", it can be answered by **T1**.

**(4) Sub-query Evaluation**

*"Faraar"*      *"Main Pyaasa tum"*

res( **q1** ) over **P1** → res( **q2** ) over **T1** → res( **Q** )

**with NLP Question Answering**      **with NL2SQL**

6

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |

X
(text, tables, databases, …)

Vectors

Encoder

**Index:**
FAISS
∞ Meta

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
| --- | --- | --- | --- |

X
(text, tables, databases, …)

Vectors

Encoder

**Index:**
FAISS
∞ Meta

**AutoEncoders**

**Fixed-size vector**
**vec(X)**

serialize(X) → Transformer-based Encoder → Transformer-based Decoder → X

**Compress**          **Reconstruct**

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
| --- | --- | --- | --- |



**Q**

Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?

Transformer-based Encoder

**vec(Q)**

Vector-based similarity search

**Index:** FAISS
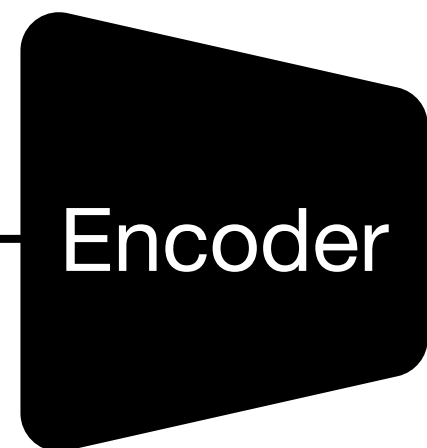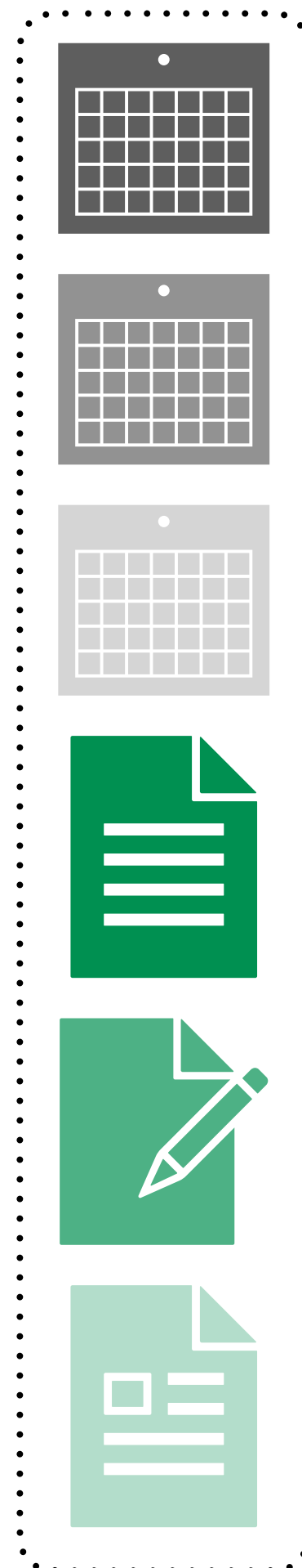
∞ Meta

Top-*K*

**T1**

**P1**

(1) Index (X-to-Vec)  (2) Data Retrieval  (3) Query Decomposition  (4) Sub-query execution
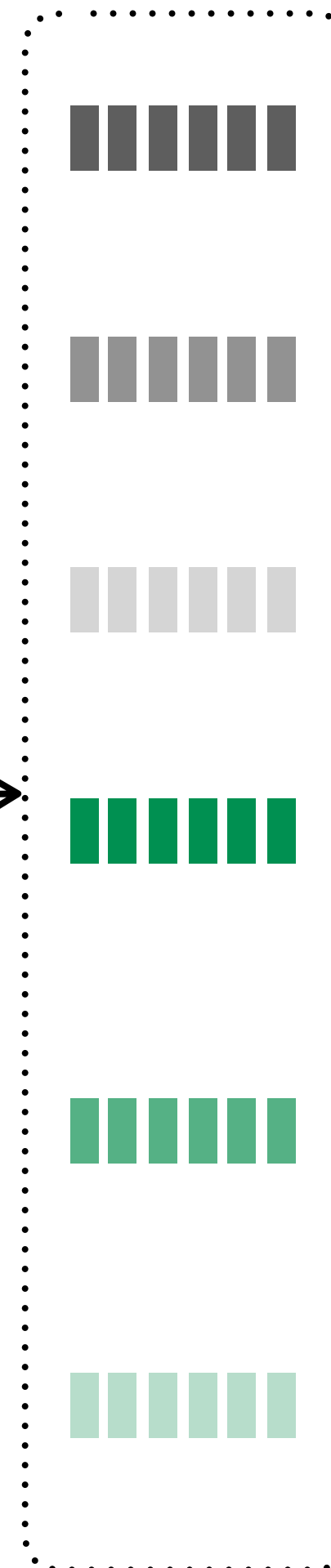
**Q**

Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?

Transformer-based Encoder

**vec(Q)**

Vector-based similarity search

**Index:** FAISS ∞ Meta

Top-*K*

T1

P1

Large *K*

**Human-in-the-loop**

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
|---|---|---|---|

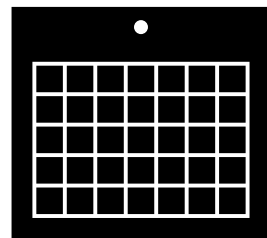**Q** Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?
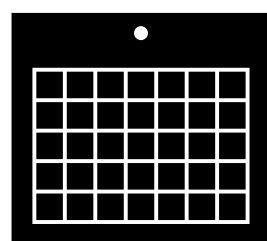
**T1**

**P1**

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
| --- | --- | --- | --- |

**GPT-3**

**Q** Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?
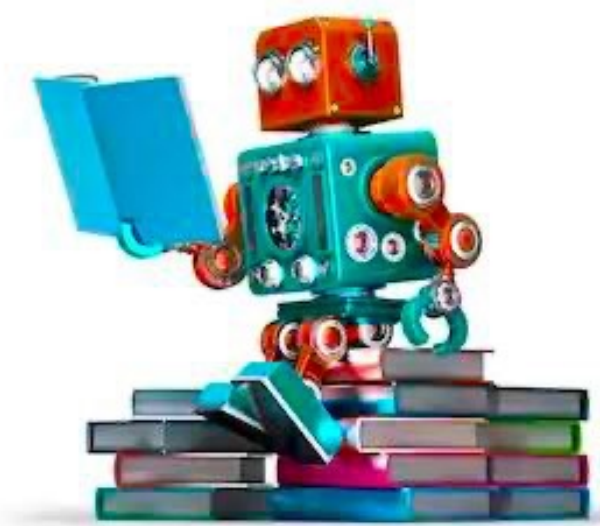
**T1**

**P1**

**Prompt 1: The passage P1** has the following content: …
**The table T1** has the following columns: Year, Song, Film, Music Director, Lyricist.
**Based on P1 and T1**, the question is "Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?".
**What sub-questions** can it be broken down into?

**GPT-3:** **q1** What is the name of the film produced by Alankar Chitra and directed by Shanker Mukherjee. It can be answered by **P1**.

**Prompt 2:** The first sub-question is "What is the name of the film produced by Alankar Chitra and directed by Shanker Mukherjee?", it can be answered by P1.

**GPT-3:** **q2** the second sub-question is "What is the name of the song in the film?", it can be answered by **T1**.

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
| --- | --- | --- | --- |

$$\mathbf{prompt}_1 = \begin{array}{l} \text{Serialize}(d'_1); \text{Serialize}(d'_2); \ldots \text{Serialize}(d'_m) \\ \text{Based on } d'_1, d'_2, \ldots \text{ and } d'_m, \text{ the question is } Q, \\ \text{what sub-questions can it be broken down into?} \end{array}$$

**Q** Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?

**T1**

**P1**

**Prompt Generation**

**Initial Prompt**

**prompt₁**

**GPT-3**

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
|---|---|---|---|

$$\textbf{prompt}_1 = \quad \text{Serialize}(d_1'); \text{Serialize}(d_2'); \ldots \text{Serialize}(d_m')$$
$$\text{Based on } d_1', d_2', \ldots \text{ and } d_m', \text{ the question is } Q,$$
$$\text{what sub-questions can it be broken down into?}$$

**Q** Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?

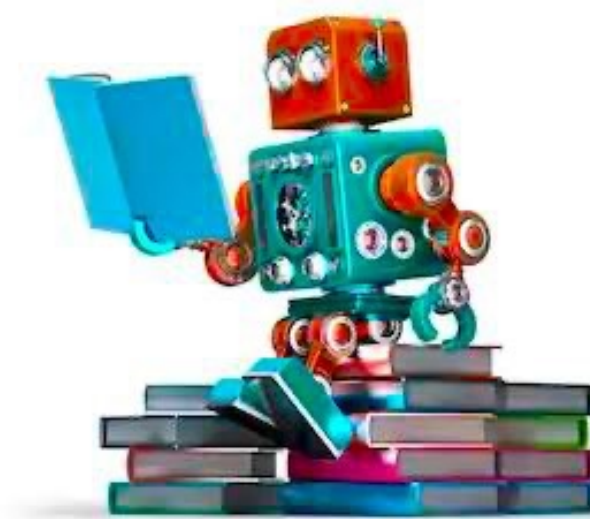**T1**

**P1**

**Prompt Generation**

**Initial Prompt**

**prompt₁**

**(qᵢ, dᵢ)**

**GPT-3**

10

$$\textbf{prompt}_1 = \quad \text{Serialize}(d'_1); \text{Serialize}(d'_2); \ldots \text{Serialize}(d'_m)$$
$$\text{Based on } d'_1, d'_2, \ldots \text{ and } d'_m, \text{ the question is } Q,$$
$$\text{what sub-questions can it be broken down into?}$$

**Q** Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?

**T1**

**P1**

**Initial Prompt**

**prompt₁**

**Prompt Generation**

**(qᵢ, dᵢ)**

**Next Prompt**

**promptᵢ**

**GPT-3**

$$\textbf{prompt}_i = \quad \text{The [N] sub-query is [Q], it can be answered by [D]}$$

10

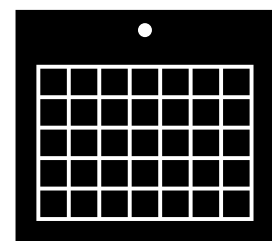(1) Index (X-to-Vec)  (2) Data Retrieval  (3) Query Decomposition  (4) Sub-query execution

$$\textbf{prompt}_1 = \text{Serialize}(d'_1); \text{Serialize}(d'_2); \ldots \text{Serialize}(d'_m)$$

Based on $d'_1, d'_2, \ldots$ and $d'_m$, the question is $Q$, what sub-questions can it be broken down into?

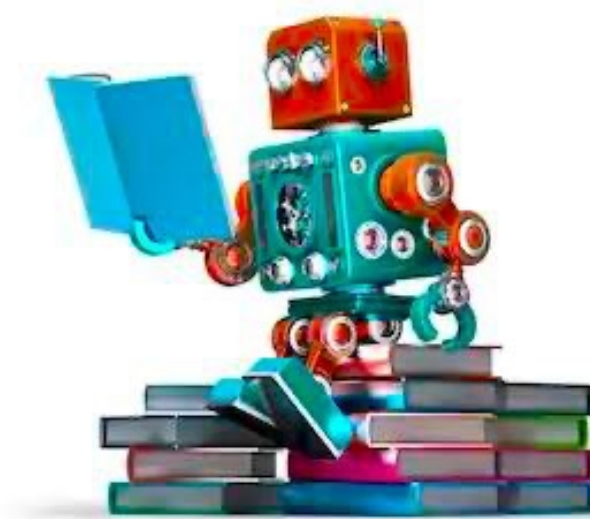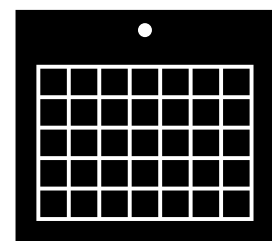Q Which songs appeared in a film produced by Alankar Chitra and directed by Shanker Mukherjee?

T1

P1

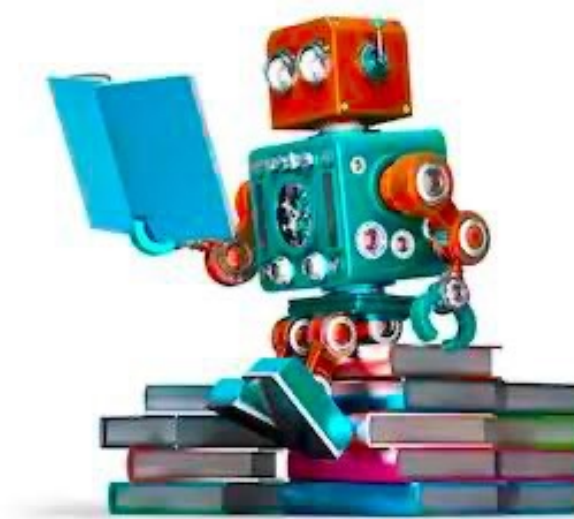Initial Prompt

prompt$_1$

Prompt Generation

$(q_i, d_i)$

Next Prompt

prompt$_i$

GPT-3

Stop

$$\textbf{prompt}_i = \text{The [N] sub-query is [Q], it can be answered by [D]}$$

10

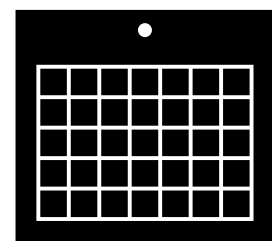| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
|---|---|---|---|

**Natural Language Query over Text**

NLP
Question
Answering

**Natural Language Query over Table(s)**

NL2SQL
(SIGMOD'23)

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
|---|---|---|---|

**Natural Language Query over Text**

NLP Question Answering

NLDBs "Database reasoning over multiple text files"

**Natural Language Query over Table(s)**

NL2SQL (SIGMOD'23)

TableQA (EMNLP'22)

| (1) Index (X-to-Vec) | (2) Data Retrieval | (3) Query Decomposition | (4) Sub-query execution |
| --- | --- | --- | --- |

**Natural Language Query over Text**

**Natural Language Query over Table(s)**

| NLP Question Answering | NLDBs "Database reasoning over multiple text files" | NL2SQL (SIGMOD'23) | TableQA (EMNLP'22) |
| --- | --- | --- | --- |

**Query Optimizer**

# Preliminary Result

- Data lake
  - 400K web tables
  - 6M passages (text files)
- Queries: 18, manually designed
  - Goal: ensuring each query can be answered by the data lake
- Test:
  - Is the decomposition correct?
  - Is the sub-query evaluation correct?

# Sample Successful Case (12/18)

- Question: Which is taller, the tallest building in the UK or the tallest building in South Korea?
- Retrieved datasets:

Table T1

| Rank | Official Name | Height (m) | ... |
|------|---------------|------------|-----|
| 1 | The Shard | 310 | ... |
| 2 | 22 Bishopsgate | 278 | ... |
| 3 | One Canada Square | 235 | ... |
| ... | ... | ... | ... |

*Source: https://en.wikipedia.org/wiki/*
*List_of_tallest_buildings_in_the_United_Kingdom*

Table T2

| Rank | Name | Height (m) | ... |
|------|------|------------|-----|
| 1 | Lotte World Tower | 550 | ... |
| 2 | Landmark Tower | 412 | ... |
| 3 | Tower A | 339 | ... |
| ... | ... | ... | ... |

*Source: https://en.wikipedia.org/wiki/*
*List_of_tallest_buildings_in_South_Korea*

- Decomposed queries:
  1. "What is the height of the tallest building in the UK?", it can be answered by T1;
  2. "What is the height of the tallest building in the South Korea?", it can be answered by T2.
- GPT-3 for aggregation: Question, result of subquery-1, result of subquery-2

# Sample Failure Case (6/18)

- Question: What year was the first German film that won the Academy Award for Best Foreign language Film released?

- Retrieved datasets:

Passage P1

> ***The Tin Drum*** *(German: Die Blechtrommel) is a **1979** film adaptation of Günter Grass' novel of the same title, directed by Volker Schlöndorff from a screenplay co-written with Jean-Claude Carrière and Franz Seitz.*
> …

Source: https://en.wikipedia.org/wiki/The_Tin_Drum_(film)

Table T1

| ceremony | film title used in nomination | result | ... |
|---|---|---|---|
| 51st | the glass cell | nominee | ... |
| 52nd | **the tin drum** | **won academy award** | |
| 53rd | fabian | not nominated | ... |
| ... | ... | | ... |

Source: https://en.wikipedia.org/wiki/
List_of_German_submissions_for_the_Academy_Award
_for_Best_Foreign_Language_Film

- Decomposed queries:

  1. "What year was the film The Tin Drum released?", it can be answered by T1;

  2. "What was the first German film that won the Academy Award for Best Foreign language Film?", it can be answered by T1.

# Sample Failure Case (6/18)

- Question: What year was the first German film that won the Academy Award for Best Foreign language Film released?

- Retrieved datasets:

Passage P1

> ***The Tin Drum*** *(German: Die Blechtrommel) is a* **1979** *film adaptation of Günter Grass' novel of the same title, directed by Volker Schlöndorff from a screenplay co-written with Jean-Claude Carrière and Franz Seitz.* …

Source: https://en.wikipedia.org/wiki/The_Tin_Drum_(film)

Table T1

| ceremony | film title used in nomination | result | ... |
|----------|-------------------------------|--------|-----|
| 51st | the glass cell | nominee | ... |
| 52nd | **the tin drum** | **won academy award** | |
| 53rd | fabian | not nominated | ... |
| ... | ... | | ... |

Source: https://en.wikipedia.org/wiki/
List_of_German_submissions_for_the_Academy_Award
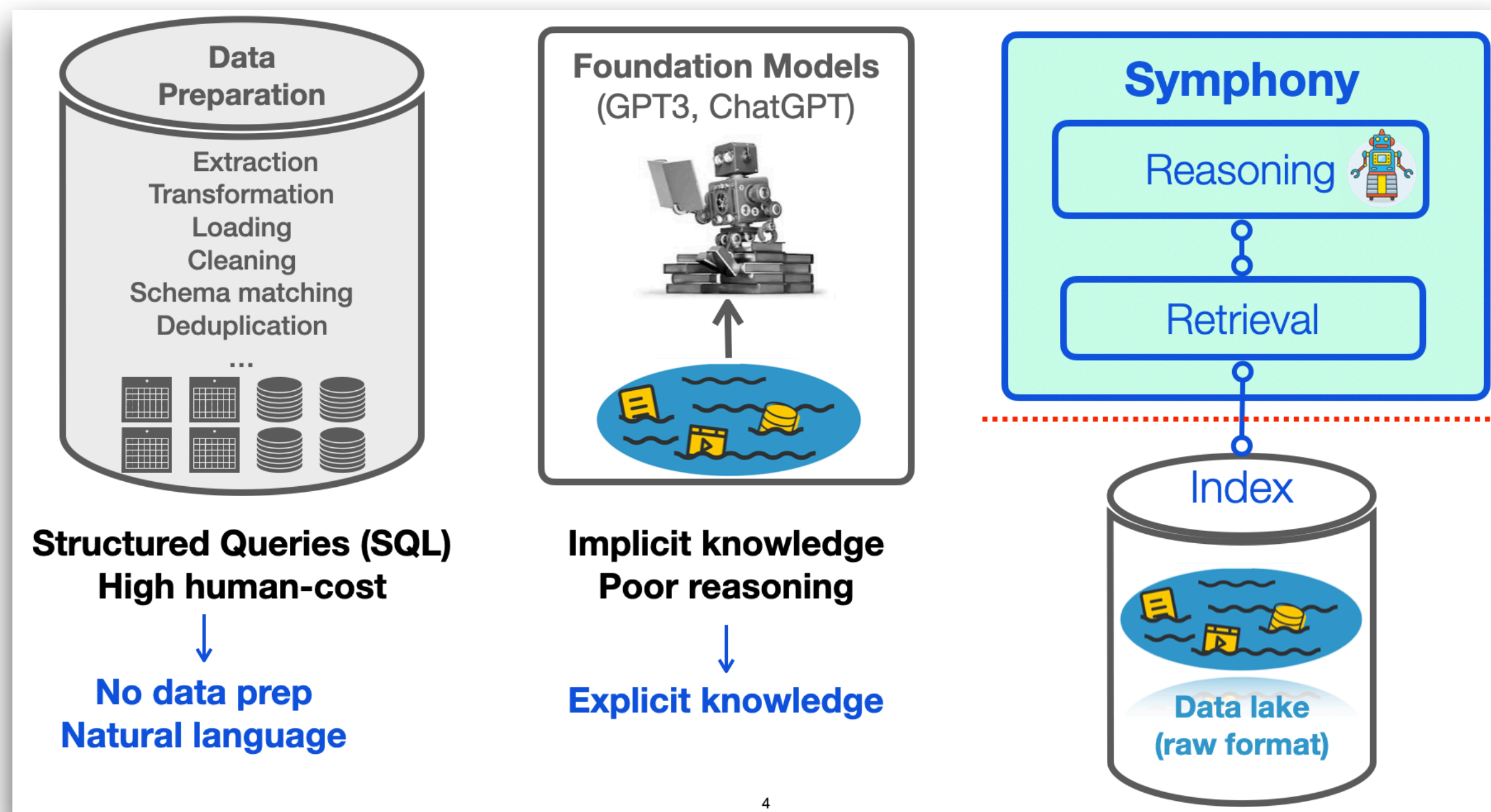_for_Best_Foreign_Language_Film

- Decomposed queries:

1. "What year was the film The Tin Drum released?", it can be answered by ~~T1,~~ **P1**

2. ~~"What was the first German film that won the Academy Award for Best Foreign language Film?", it can be answered by T1.~~

# Conclusion

- **New way** of exploring multi-modal data lakes

- **AI-Assistant**: Political/business leaders surrounded by advisors



# Future Work

- Improving the index (X-to-Vec)

  - Contrastive learning

  - Combining with traditional string similarity search

- Sub-query execution module

  - NL2SQL