

# Seeking Order in Disorder: Towards Accurate and Efficient Document Analytics



Yiming Lin<sup>1</sup>, Madelon Hulsebos<sup>1</sup>, Ruiying Ma<sup>2</sup>, Shreya Shankar<sup>1</sup>,  
Sepanta Zeighami<sup>1</sup>, Eugene Wu<sup>3</sup>, Aditya Parameswaran<sup>1</sup>  
<sup>1</sup>UC, Berkeley, <sup>2</sup>Tsinghua University, <sup>3</sup>Columbia University

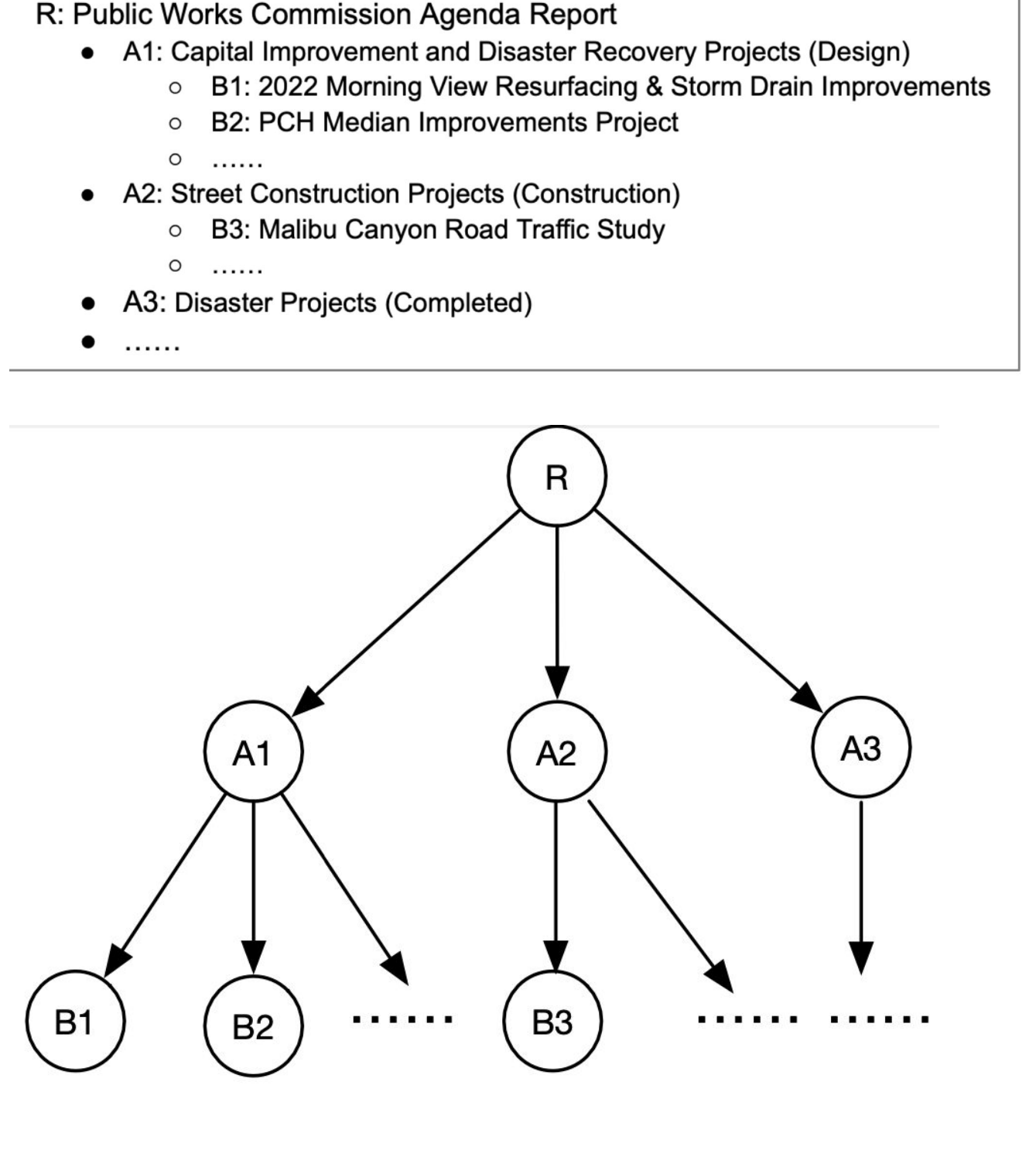


COLUMBIA UNIVERSITY

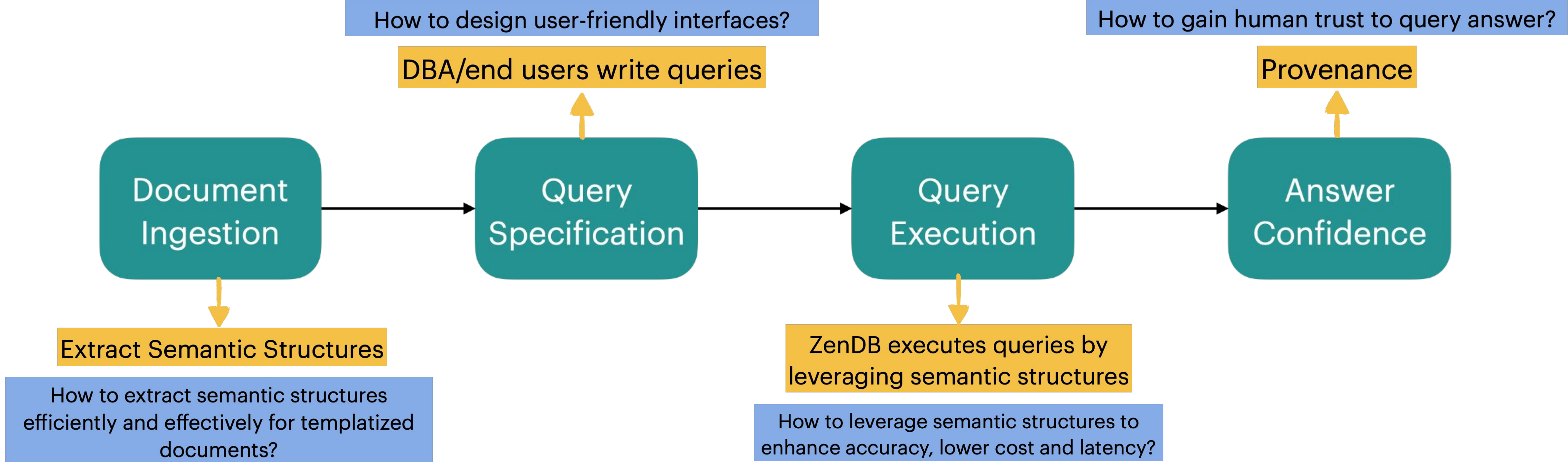
## Background & Question

- Over 80% of data exists in unstructured formats, and extracting values from unstructured documents remains a considerable challenge.
- LLMs provide us the ability to extract semantic objects, labels, and relationships from text much easier than before
- Q: can we build a data management system for text data and expose a SQL-like query language for perform advanced analytics beyond simple retrieval?

## Templatized Documents



## ZenDB's Workflow



## Query Interfaces

- Require minimal domain knowledge:
  - table name and description
  - Attribute name/type/description

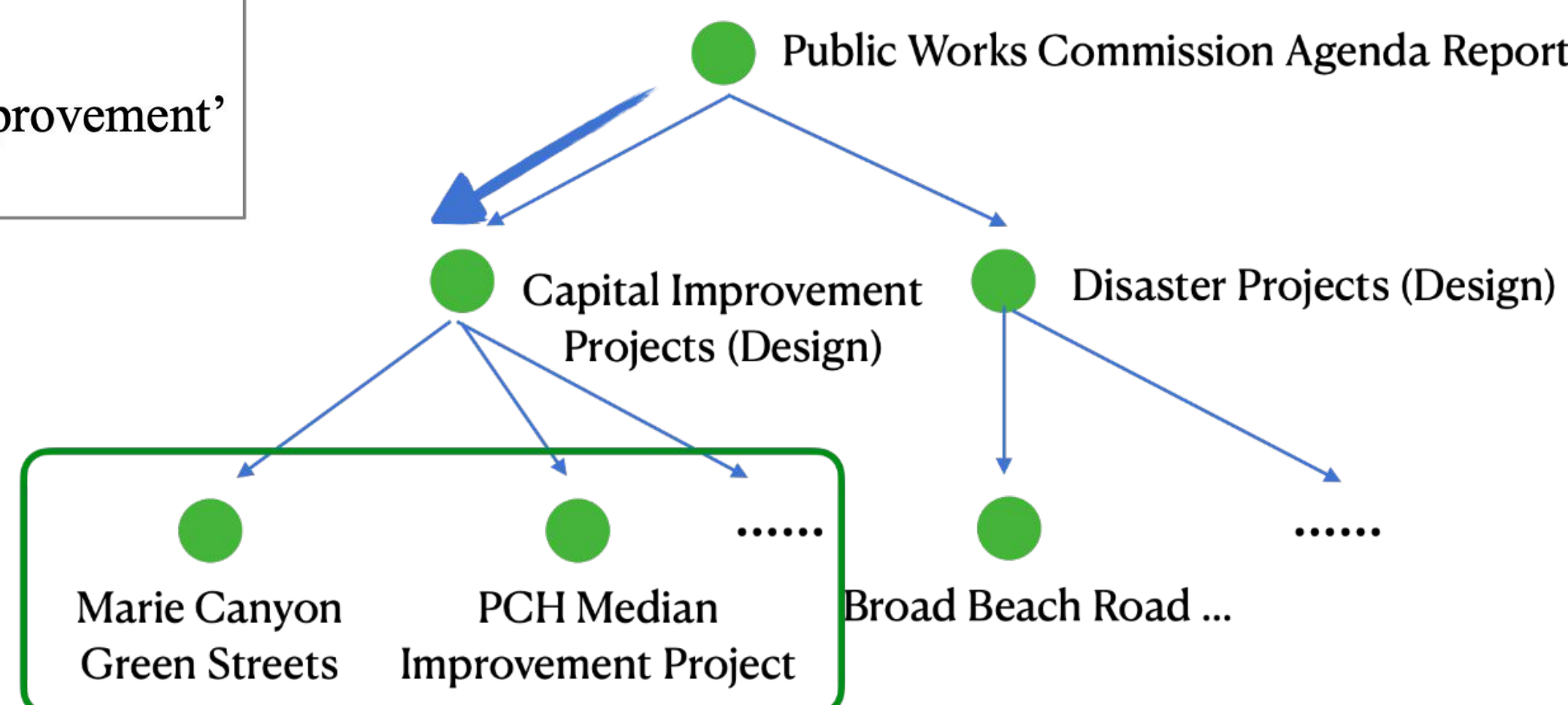
```
CREATE TABLE Projects AS (table_desc = 'The projects table contains the description of a set of civic agenda projects.')
CREATE ATTRIBUTE name ON Projects AS (attr_type = text, attr_desc = 'name of project')
CREATE ATTRIBUTE type ON Projects AS (attr_type = text, attr_desc = 'type of project')
CREATE ATTRIBUTE begin_time ON Projects AS (attr_type = date, attr_desc = 'begin time of project')
```

Q1: 'What is the number of capital improvement projects that start later than 2022'

```
Q2: SELECT COUNT(Projects.name)
FROM Projects,
WHERE Projects.type = 'Capital Improvement'
AND Projects.begin_time > '2022'
```

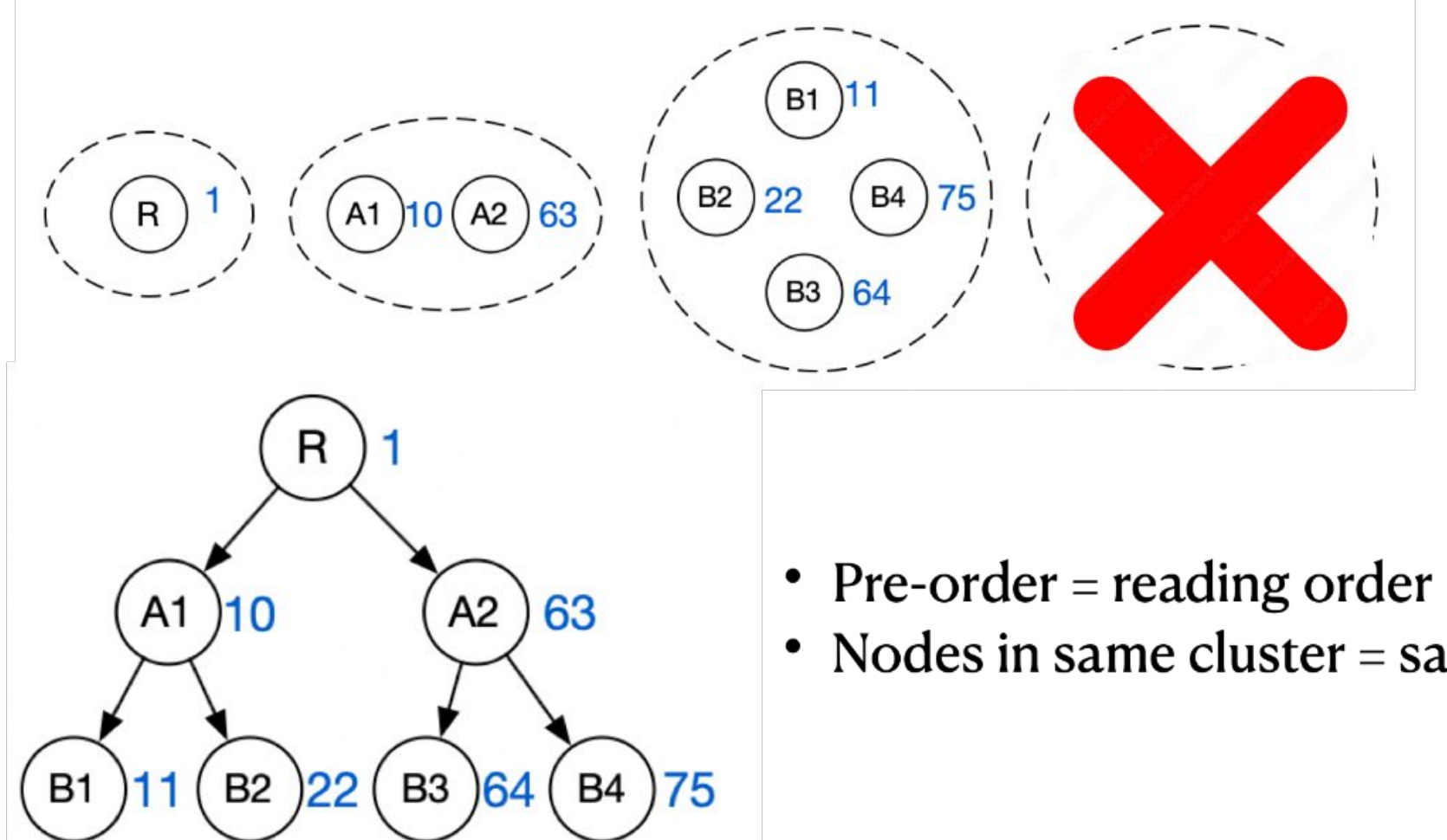
### Sketch for each node:

- Name of current node and ancestors
- Summary
- Top-1 sentence similar to query



## Semantic Structure

### Tree Construction



Documents sharing templates - by no-LLM visual patterns matching

## Experiment Takeaways

- ZenDB VS LLM** (all GPT-4)
  - Up to +29% precision, +31% recall, 30x saving of costs, 4x latency saving
- ZenDB VS RAG** (all GPT-4)
  - Up to +61% precision, +80% recall, 1.7x higher cost and 1.3x higher latency
- ZenDB + GPT-3.5-Turbo** (100x cheaper)
  - Paying **one dollar**, you can run **5.5k SQLs** on average on single document, with usable quality (-7% precision and -5% recall VS ZenDB + GPT-4)