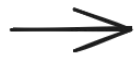


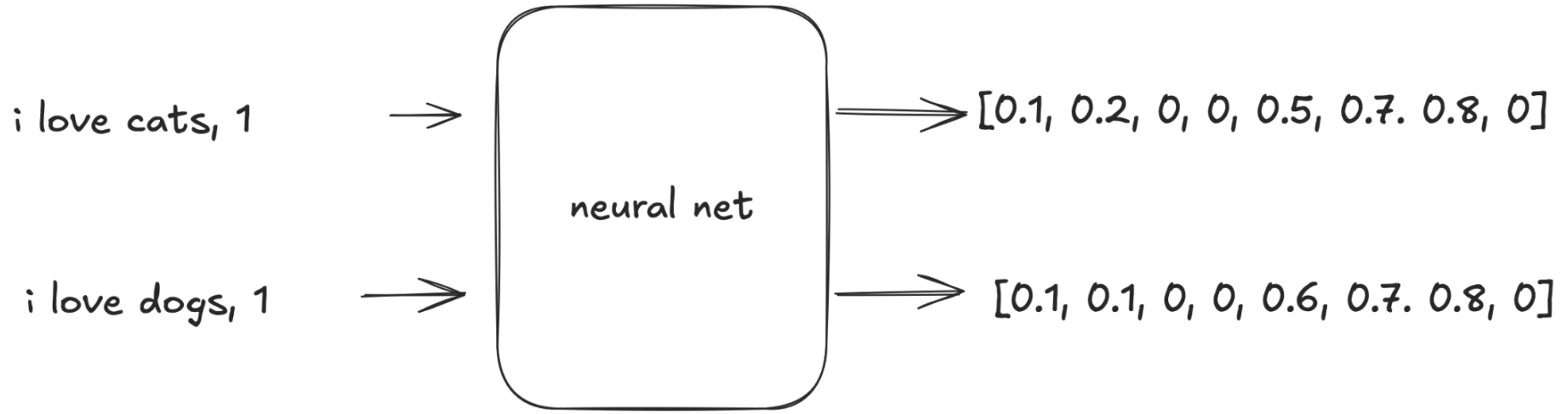
RAG, A2A, Resource-Aware Optimization

made with  for “Little ML book club”

i love cats

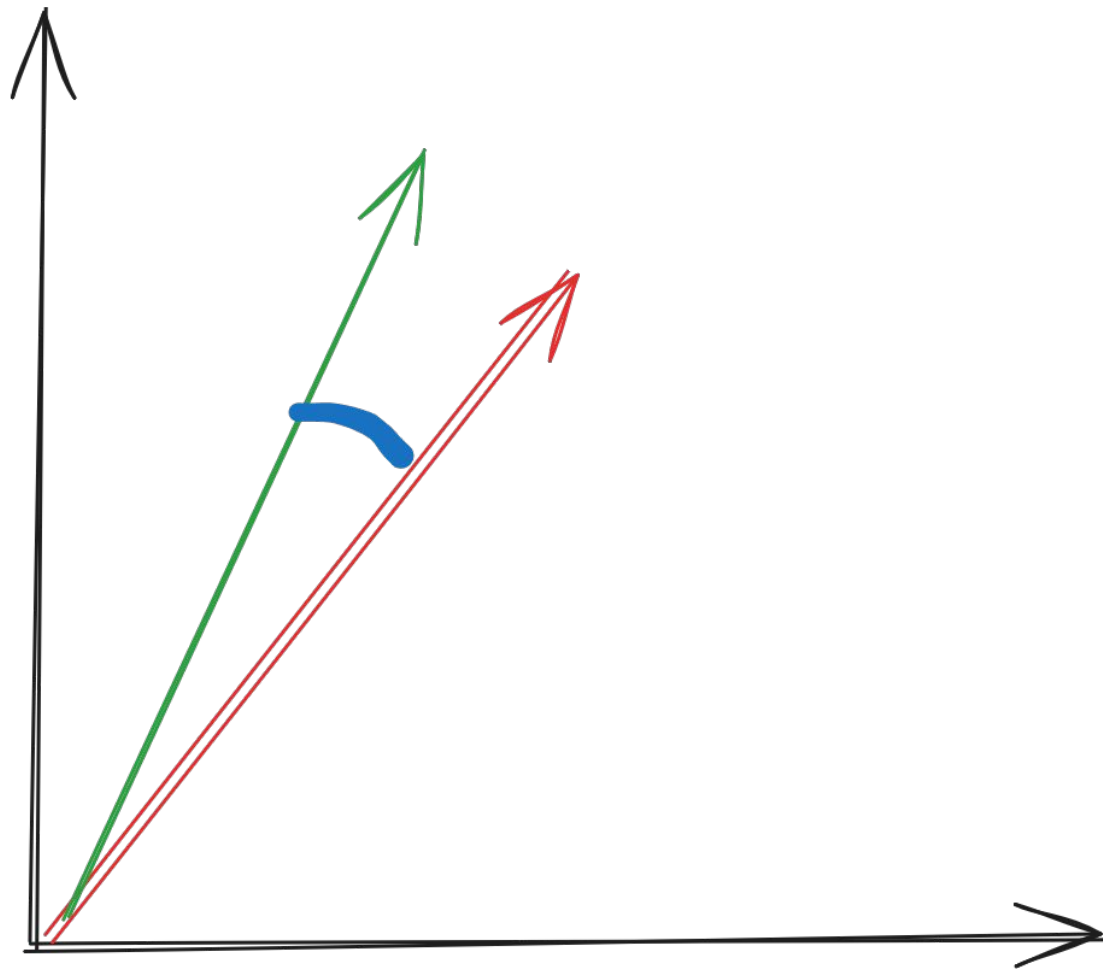


[0.1, 0.2, 0, 0, 0.5, 0.7, 0.8, 0]



[0.1, 0.2, 0, 0, 0.5, 0.7, 0.8, 0]

[0.1, 0.1, 0, 0, 0.6, 0.7, 0.8, 0]





problem - exact word search

TF-IDF Working Scheme

Step 1: Document Collection

The cat
sat on
the mat.

Doc 1

A dog
sat on
the log.

Doc 2

Cats and
dogs are
pets.

Doc 3

Step 2: Tokenization & Preprocessing

Preprocessing

Tokenize,
Lowercase,
Remove
Stop Words
(*'the', 'a', 'on',
'and', 'are'*)

"cat"
"sat"
"mat"
"dog"
"log"
"cats"
"dogs"
"pets"

How often in *this* doc?

Term Frequency (TF) (per document)

Doc 1:

- $TF("cat") = \frac{1}{3}$
- $TF("sat") = \frac{1}{3}$
- $TF("mat") = \frac{1}{3}$

How rare *across* docs?

Inverse Document Frequency (IDF) (across all docs)

Total Docs (N) = 3

$DF("cat") = 1$ (in Doc 1) $\rightarrow IDF("cat") = \log\left(\frac{3}{1}\right)$

$DF("sat") = 2$ (in Doc 1, Doc 2) $\rightarrow IDF("sat") = \log\left(\frac{3}{2}\right)$

$DF("dog") = 2$ (in Doc 2, Doc 3) $\rightarrow IDF("dog") = \log\left(\frac{3}{2}\right)$

Step 5: TF-IDF Calculation

TF-IDF = TF * IDF

"cat" in Doc 1:

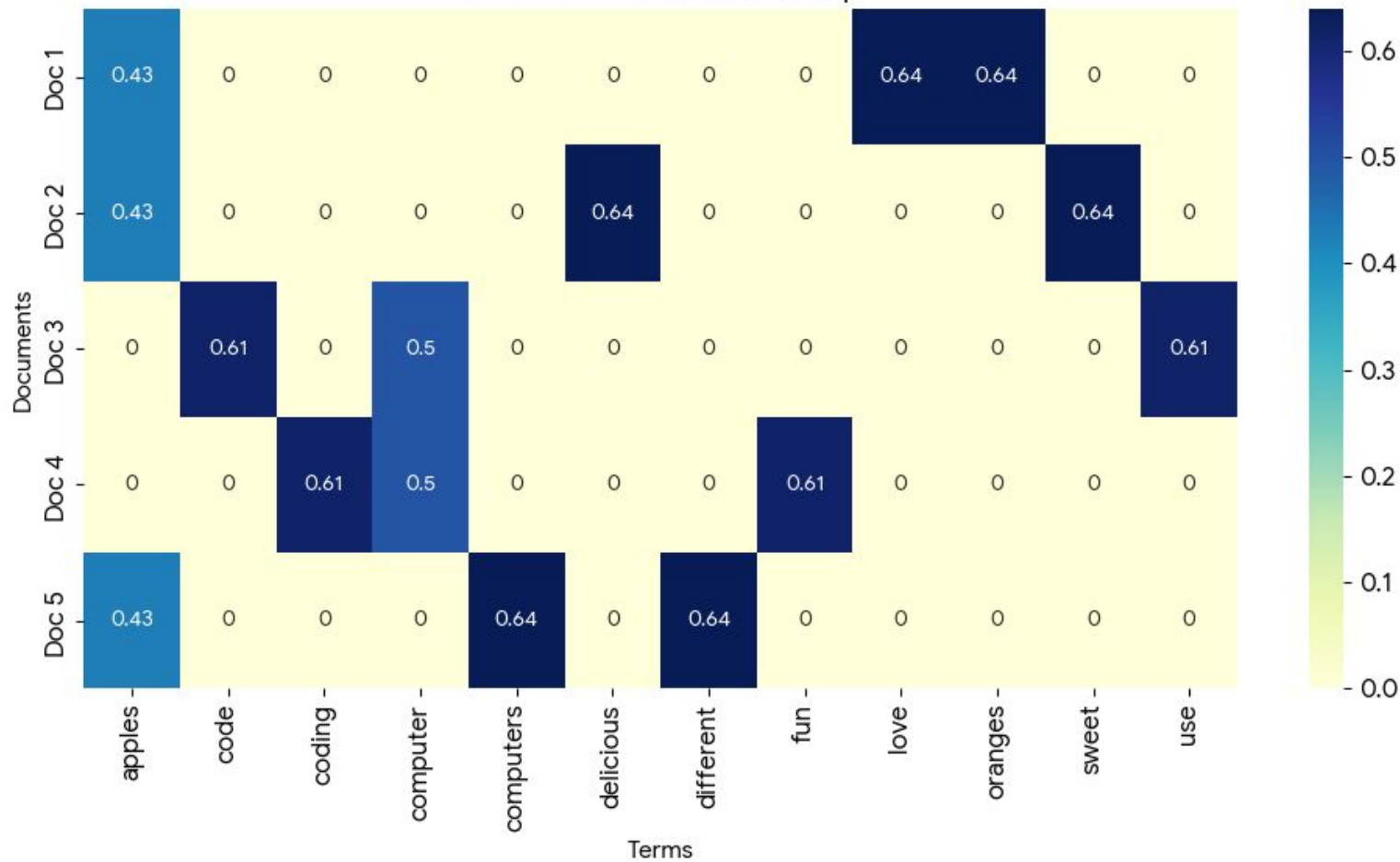
$$TF-IDF = \frac{1}{3} * \log\left(\frac{3}{1}\right)$$

Final Importance Score

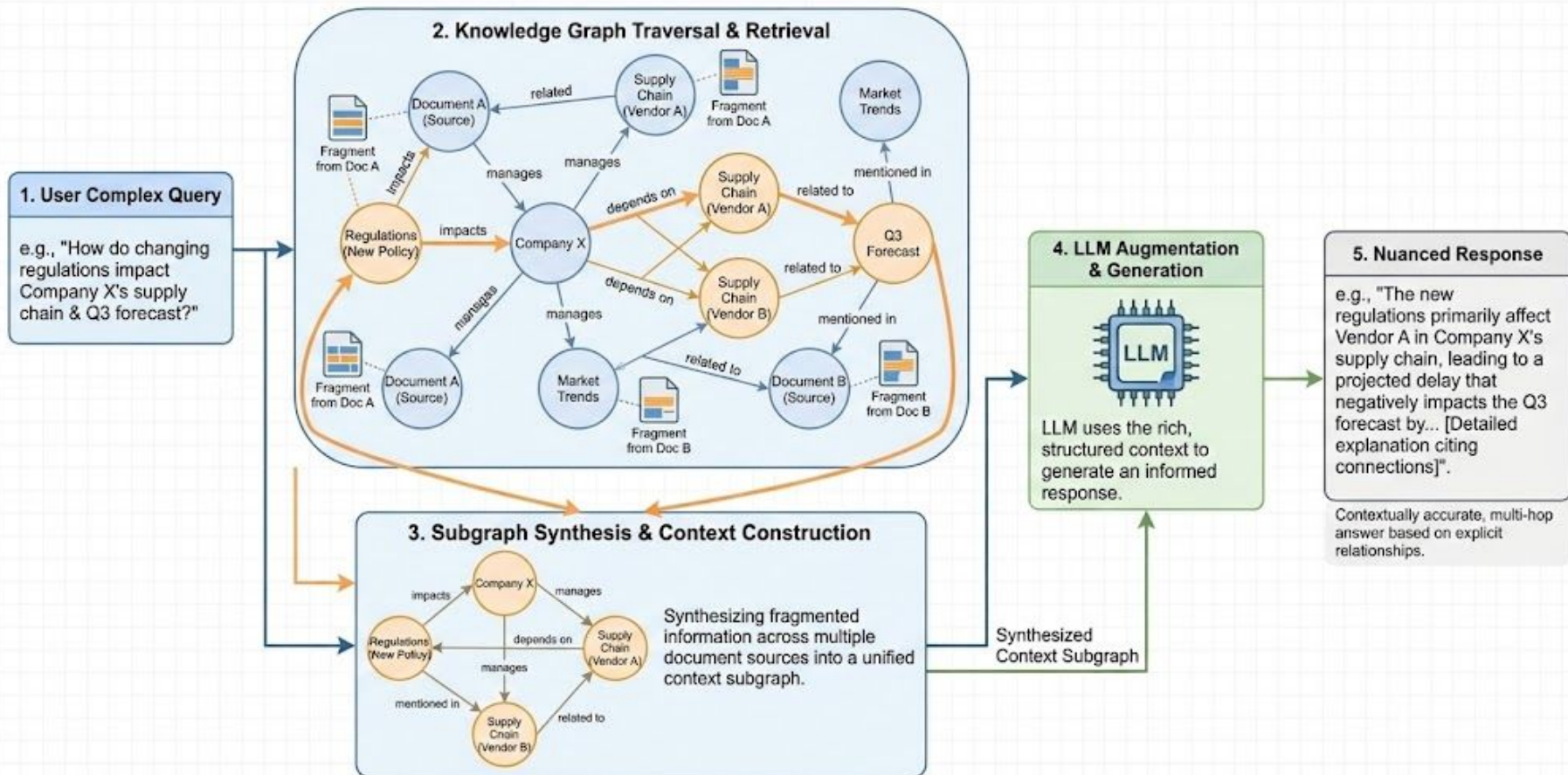
Step 6: TF-IDF Matrix

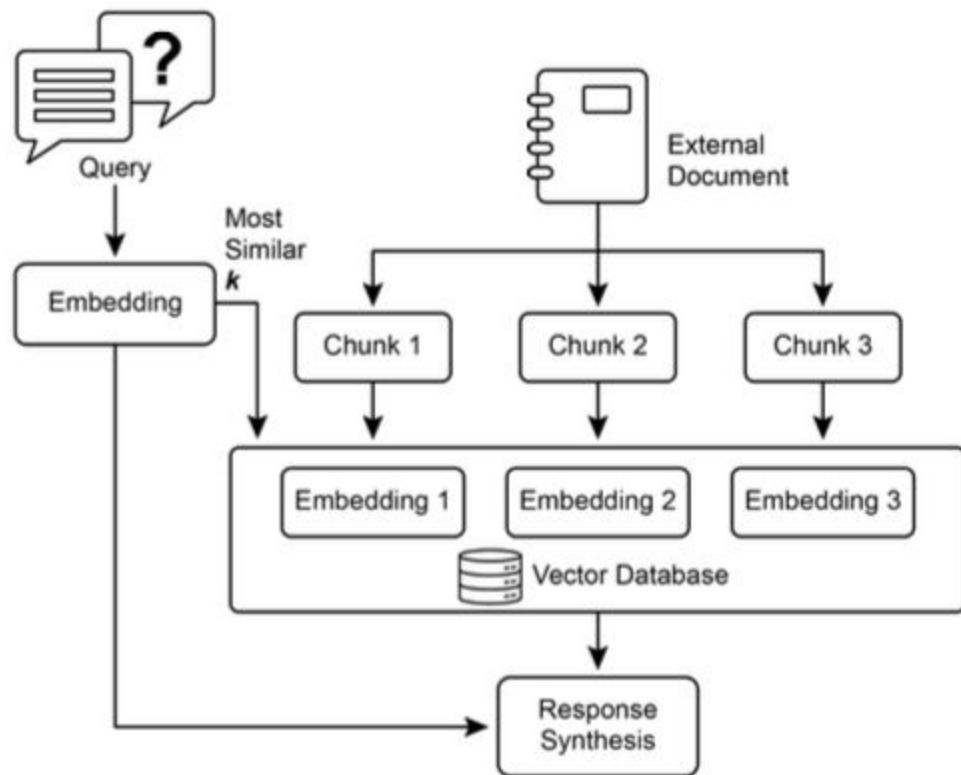
	cat	sat	mat	dog	log	cats	dogs	pets
Doc 1	0.3	0.1	0.04	0.09	0.02	0.03	0.03	0.08
Doc 2	0.11	0.09	0.02	0.09	0.04	0.05	0.03	0.10
Doc 3	0.2	0.03	0.05	0.1	0.02	0.03	0.03	0.07

TF-IDF Visualization Heatmap



Graph RAG Working Scheme: Knowledge Graph-Based Retrieval & Synthesis





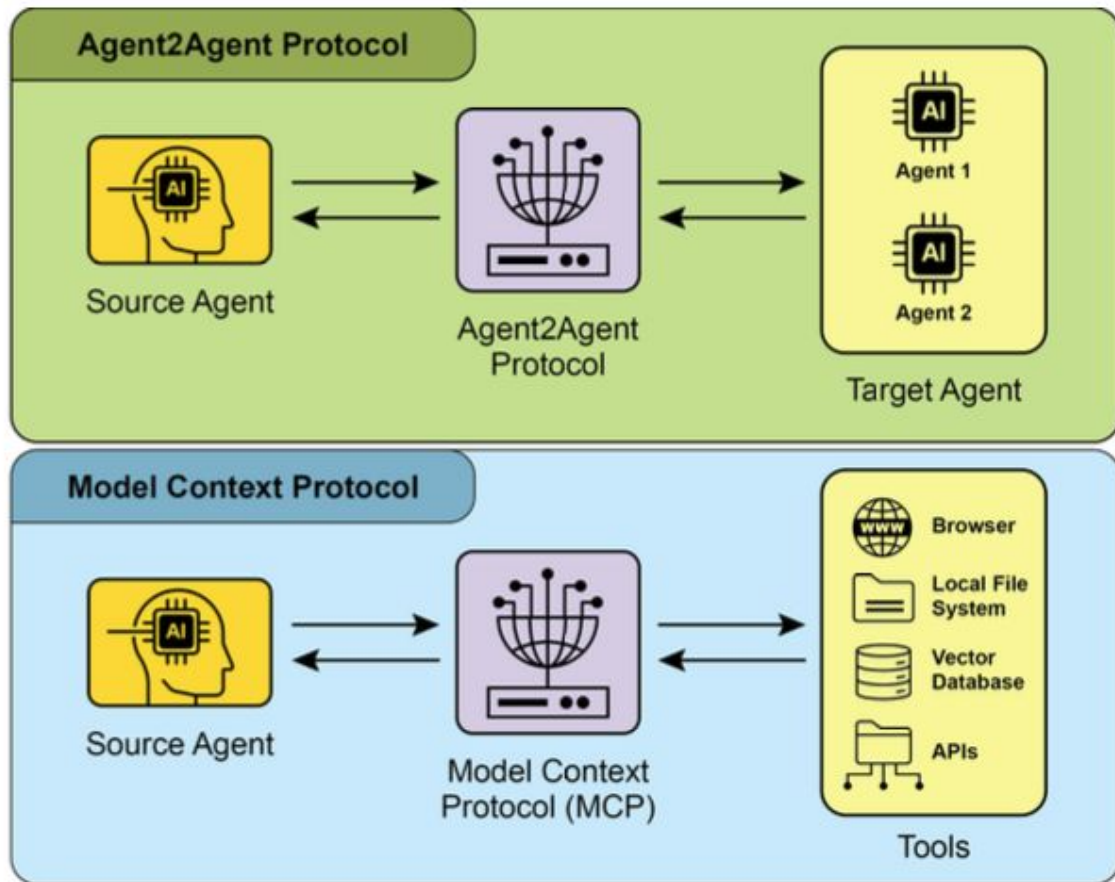


Fig.1: Comparison A2A and MCP Protocols

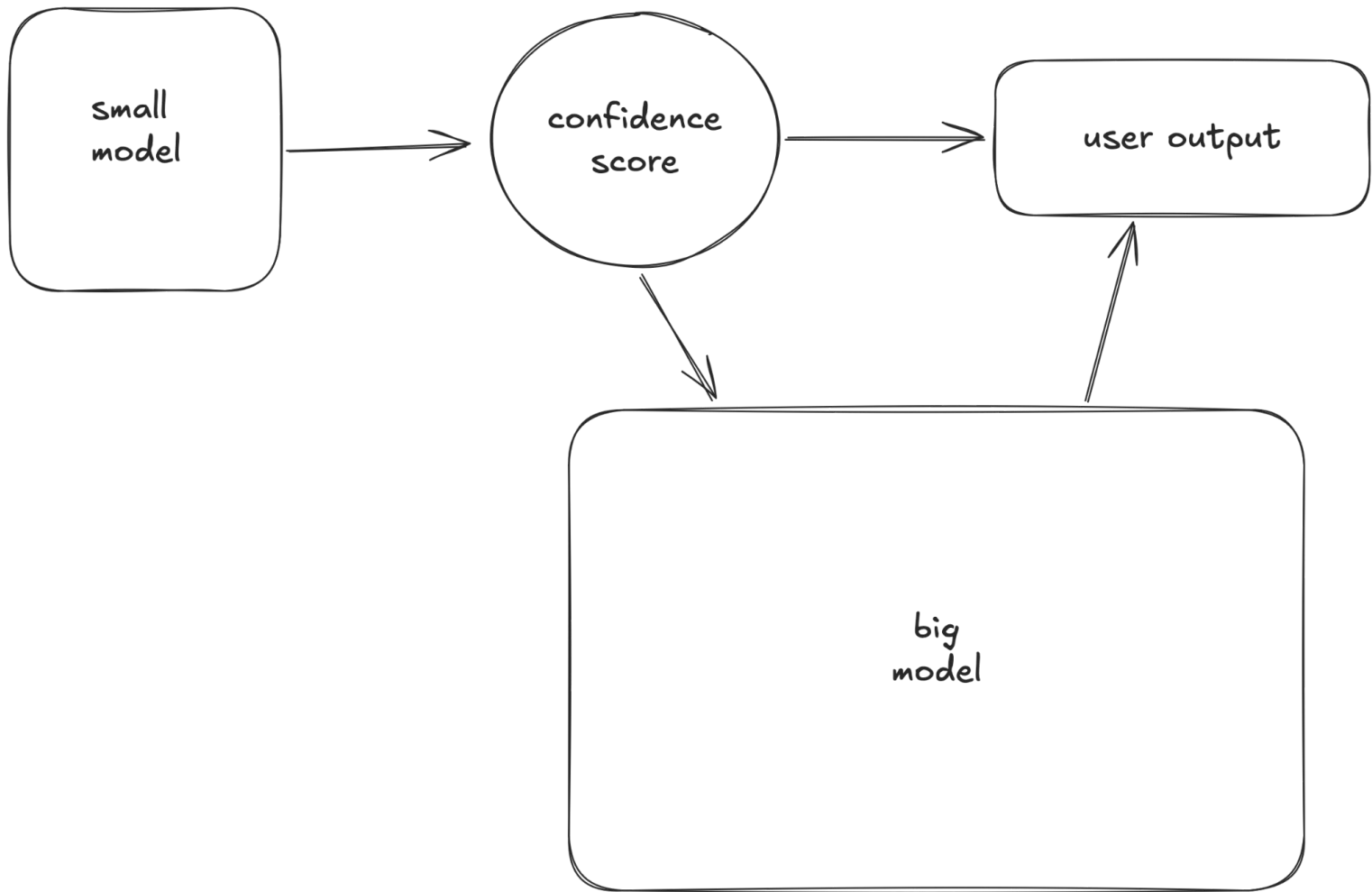
#Synchronous Request Example

```
{
  "jsonrpc": "2.0",
  "id": "1",
  "method": "sendTask",
  "params": {
    "id": "task-001",
    "sessionId": "session-001",
    "message": {
      "role": "user",
      "parts": [
        {
          "type": "text",
          "text": "What is the exchange rate from USD to EUR?"
        }
      ]
    }
  },
  "acceptedOutputModes": ["text/plain"],
  "historyLength": 5
}
```

```
# Streaming Request Example
{
  "jsonrpc": "2.0",
  "id": "2",
  "method": "sendTaskSubscribe",
  "params": {
    "id": "task-002",
    "sessionId": "session-001",
    "message": {
      "role": "user",
      "parts": [
        {
          "type": "text",
          "text": "What's the exchange rate for JPY to GBP today?"
        }
      ]
    },
    "acceptedOutputModes": ["text/plain"],
    "historyLength": 5
  }
}
```

- Synchronous Request/Response
- Asynchronous Polling
- Streaming Updates (Server-Sent Events - SSE)
- Push Notifications (Webhooks)

- use less data
- use smaller models when possible



see you next time
for

“Reasoning Technique, Evaluation and
Monitoring, Prioritization, Exploration and
Discovery”