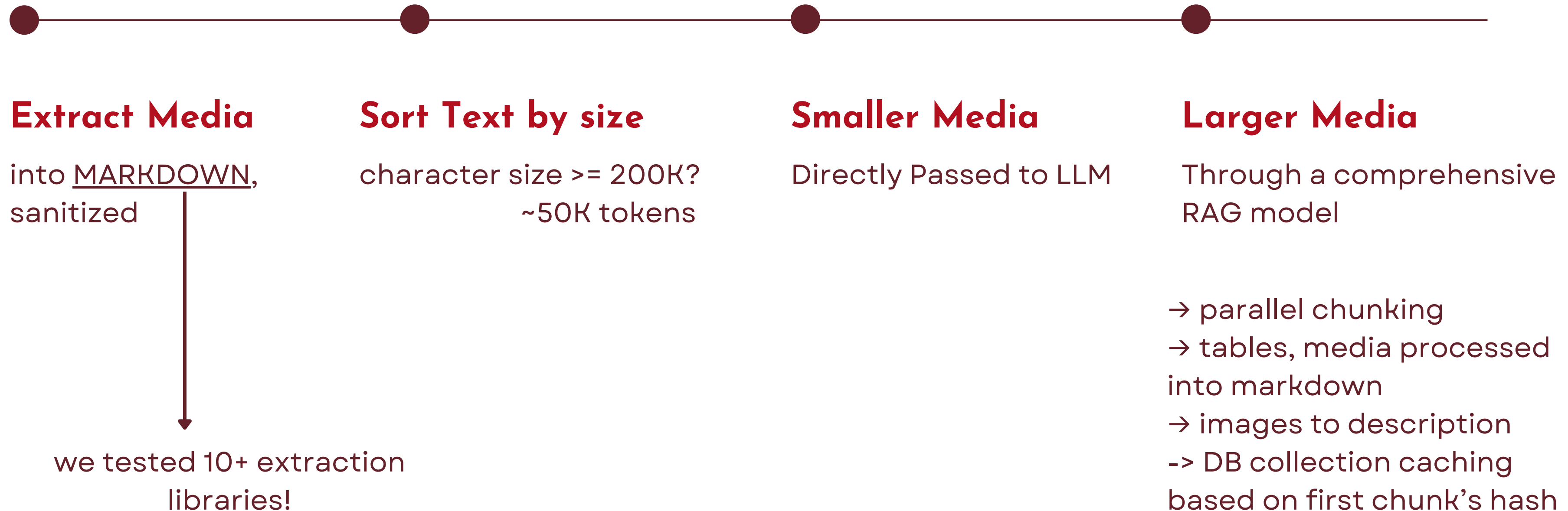Team

# Laal Tamatar

# Overall Round Progress

| 1 | Basic RAG (late entry) |
|---|---|
| 2 | RAG with optimized retrieval ← local Milvus DB, t-e-3 large |
| 3 | Sanitization, optimized chunking ← MD introduced |
| 4 | added single tool calling, agentic |
| 5 | agentic interaction, using a cURL tool |
| 6 | more tools introduced <---- cURL v/s headless browser debate |
| 7A | Agentic workflow with DOM snapshots v/s looped cURL approach |
| 7B | Fixated to cURL + tools agentic workflow |

# Rounds 1-4 (prelims) overview

**Extract Media**

into MARKDOWN, sanitized

we tested 10+ extraction libraries!

**Sort Text by size**

character size >= 200K? ~50K tokens

**Smaller Media**

Directly Passed to LLM

**Larger Media**

Through a comprehensive RAG model

→ parallel chunking
→ tables, media processed into markdown
→ images to description
-> DB collection caching based on first chunk's hash

HTTP Requests

{ documentUrl, questions[] }

{ documentUrl, questions[] }

{ documentUrl, questions[] }

'transcribe' router

API Entrypoint
routes requests gracefully

nodejs

questions[]

documentUrl

Multi-format
document
parsing & sorting

python microservice

files > 120MB
REJECT

nodejs

⟹ Response

markdown

pdf parser
(nodejs)

pdf, img

text
with image
descriptions

processingStrategy = 'batch'

questions[]

Google Gemini

2.5 flash / lite

⟹ Response

page length <= 200
~75K tokens

chunks[]

questions[]

processingStrategy
='rag'

text-emedding-3-large

questions[] || chunks[]

Generate embeddings in parallel!

vectors[][]

milvus

local vector DB

(Dockerized)

in memory
chunk-vector
map

attach context

⟹ Response

semantic search

# Finale Approach

we had two!

# The dilemma

## the 'simple' approach

### master cURL tool + util tools

- Works for both, server and client side rendering
- after the context limit approaches 50K tokens → switch to RAG
- **(chunks are created, stored in a vector DB, semantically searched and returned)**
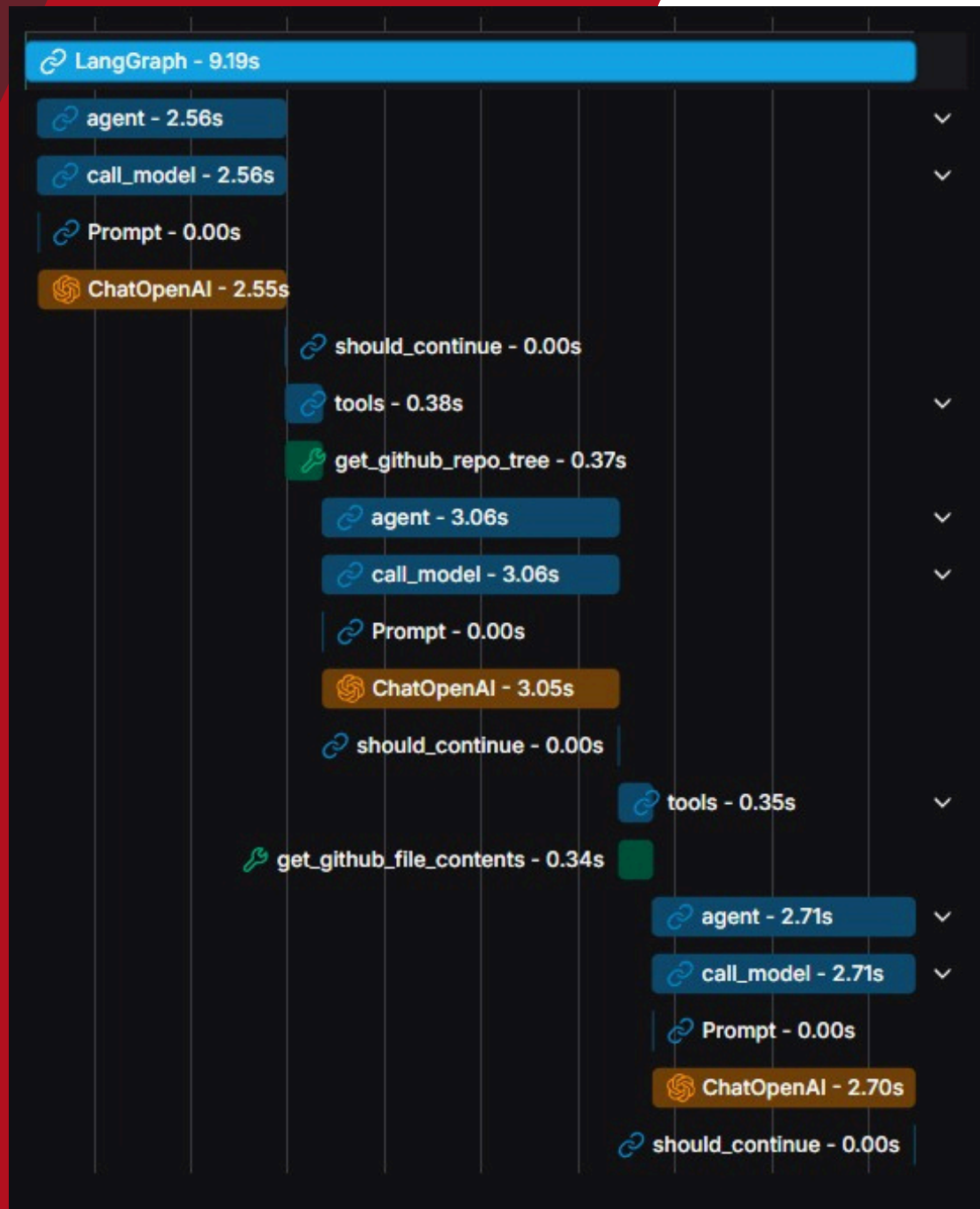- Faster
- pretty generic

## the 'optimal' approach

### headless browser + DOM snapshots on change

- Headless browser + DOM snapshots on every change
- A tool for the agent to read sanitized DOM snapshots, and inject JS into the latest page.
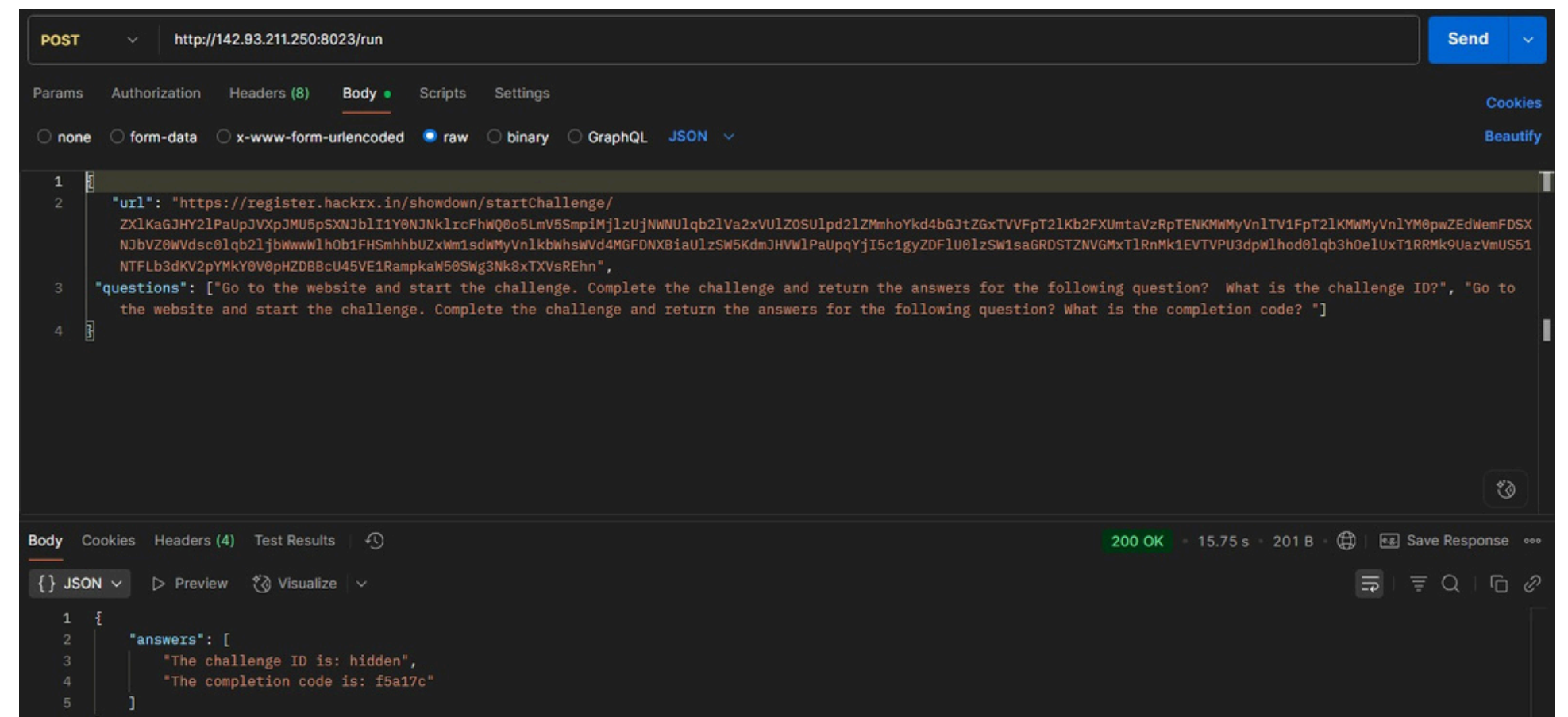- Memory overhead, less scalable
- cleaner

# Let's talk numbers

what worked out for us!



V/S

9s

15s

# Thank you!