


Agenda

→ Choosing the right LLM for our Appⁿ.

→ Cost  Model-A : 0.1 \$ per million token
Model-B : 10 \$ per million token

→ latency : Response time

→ Very important parameter for
Chat / Voice agents.

→ Any real-time appⁿ.

→ Intelligence.

↳ Reasoning

⇒ While choosing the LLM for the use-case,
we should try to find out the right balance
b/w these 3 parameters.

Ex: Voice Assistant agent.

latency ~ 800-900 ms (< 1sec)

Gemini Flash 2.5.

→ 250 tokens/sec
→ 0.25 sec time-to-first-token.
→ 0.15 \$ / Million tokens.

GPT-4o.

→ Better reasoning
→ latency & cost will be on the higher side

Ex Chat Agent with complex reasoning

⇒ Claude Sonnet x

→ Most predictable outputs ⇒ Accuracy ↑↑.
→ Code Generation / Structured output.

Use-cases: Legal analysis | Financial analysis | ...

⇒ 15 \$ / M tokens.

Ex: High Volume & low Complexity Chat Agent.

Claude Sonnet ~~X~~

Gemini Flash 2.5.

→ 250 tokens/sec
→ 0.25 sec time-to-first-token.
→ 0.15 \$ / Million tokens.
→ 100x Cheaper than Claude.

Model selection shouldn't be hard-coded.

⇒ We should be able switch our models based on the requirement.

⇒ Model selection should be configurable.

1. Router Based Model Selection.

→ Route the request to the right model based on the complexity of request.

Simple query \Rightarrow Gemini Flash \hookrightarrow Fast + Cheap.

Complex query \Rightarrow Claude Sonnet \hookrightarrow Smart + Consistent.

Image query \Rightarrow Gemini Pro, . . .

2. Agent selection \Rightarrow Configurable

\hookrightarrow We'll be able to change the underlying model without making the code changes.

anthropic.messages.create(model="Claude-sonnet-4.5")
 \hookrightarrow Hard Coding

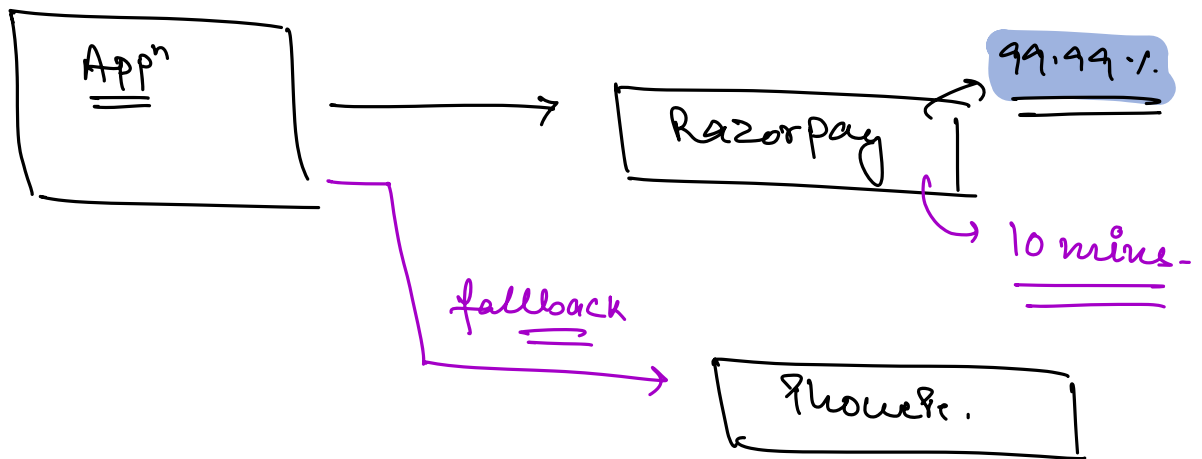
anthropic.messages.create(model=\$model-config)
 \hookrightarrow Config Variable

3. Fallback chain for Resilient Systems.

SLA

\hookrightarrow Service Level Agreement

99.9%



LLM Optimization Strategies.

1. Semantic Caching.
2. Prompt Optimization.
3. Batch Processing.
4. Two-Tier Processing