

# LLM Load Balancer

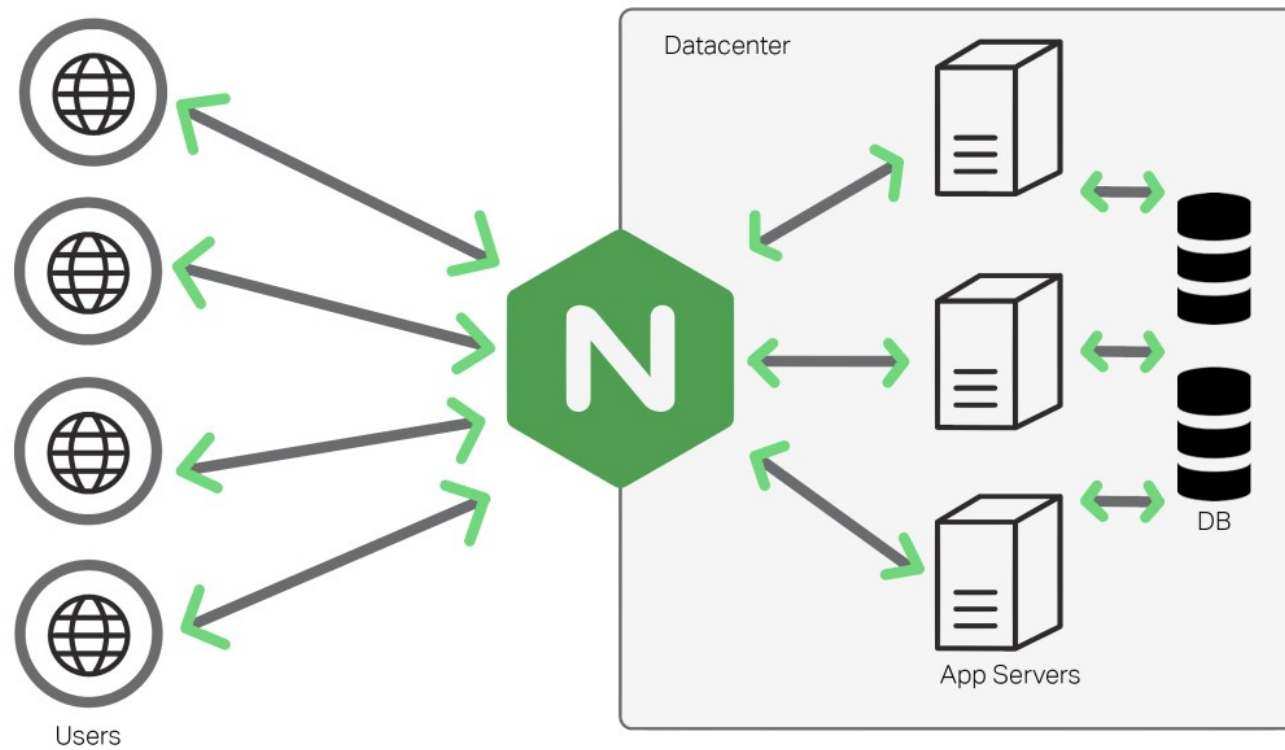
Blaise Swartwood

Jared Kagay

CSSE432

05/22/25

# Load Balancer



# Implementation



Implemented in Python



ChatGPT-2 backend model



Key features:

Round robin

Least Connections

Semantic LRU Caching

# Load Balancer Basics

1

Set up simple server backend that is multithreaded to just echo back

2

Set up client to send data

3

Create load balancer server to forward requests

4

Load balancing algorithms implementation

# Enhancements

---

Async/await implementation

---

Server LLM hosting

---

LLM Caching (basic, semantic, LRU)

---

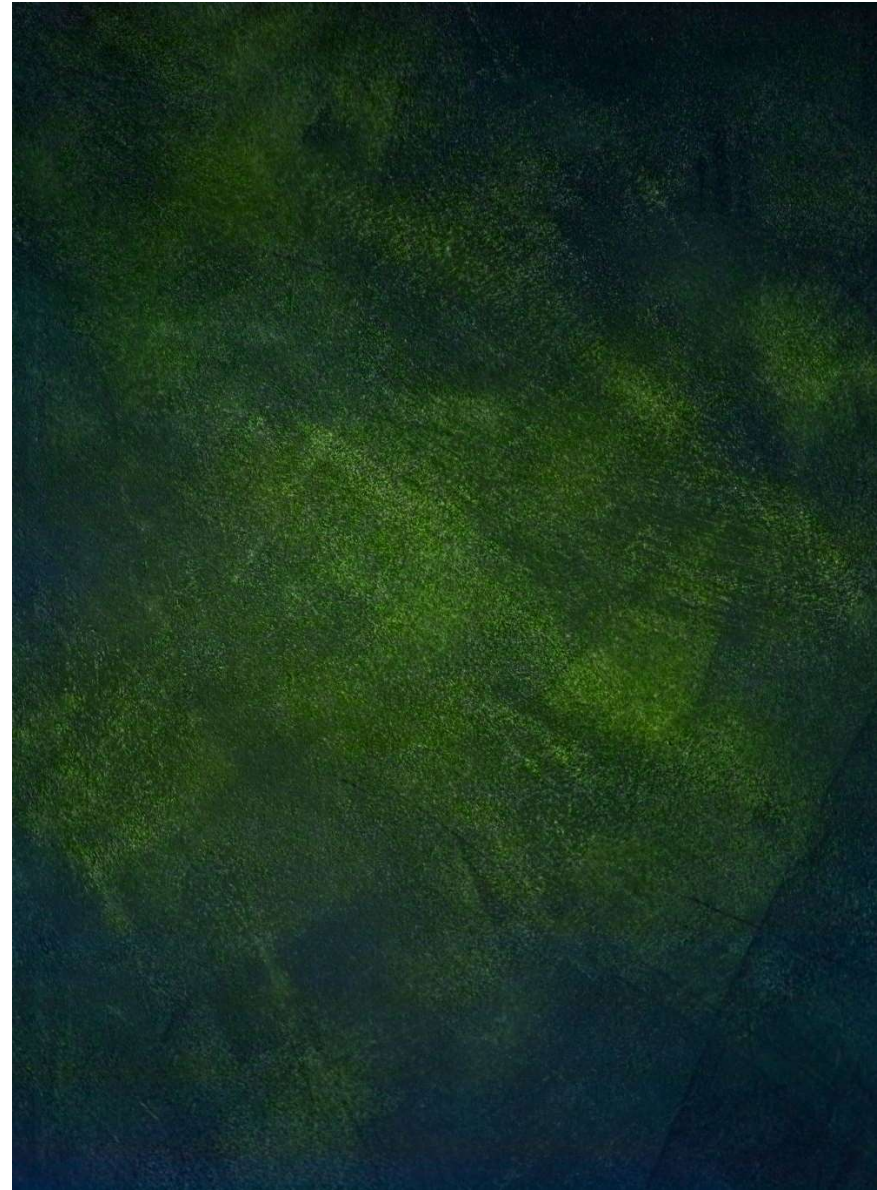
Dynamic Server Connection

---

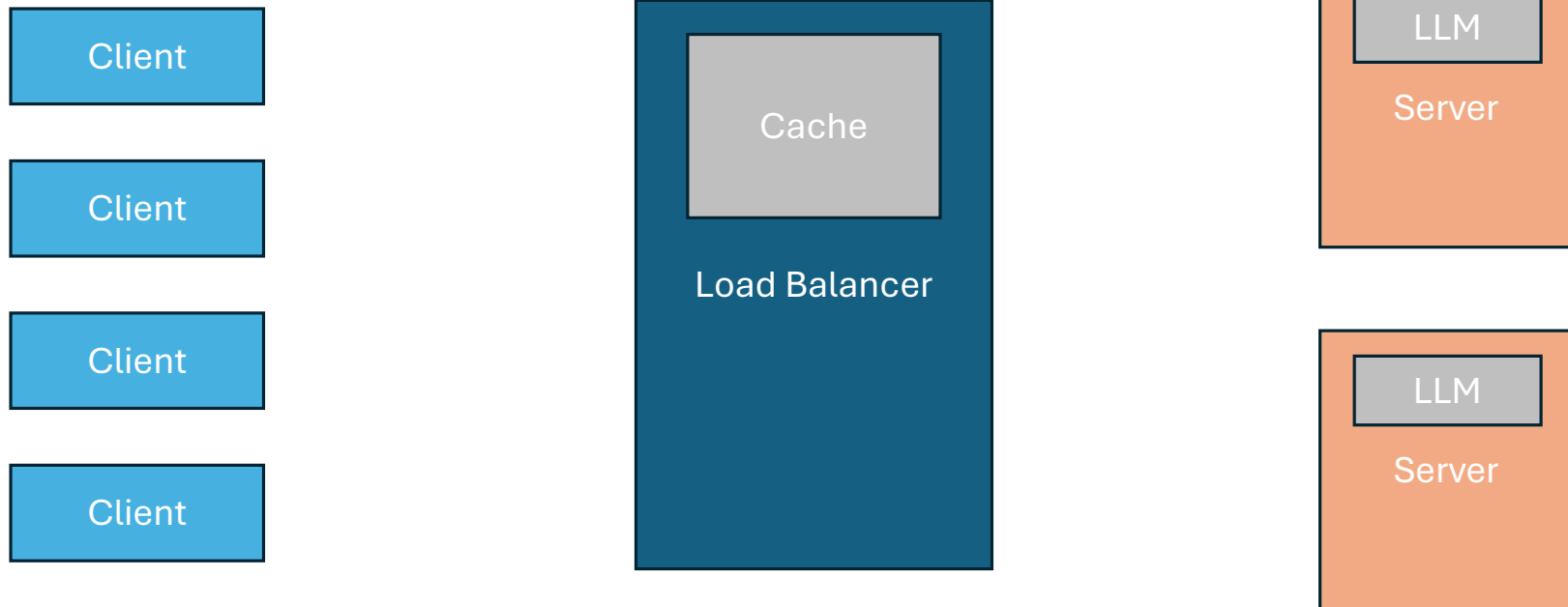
Health Checks

---

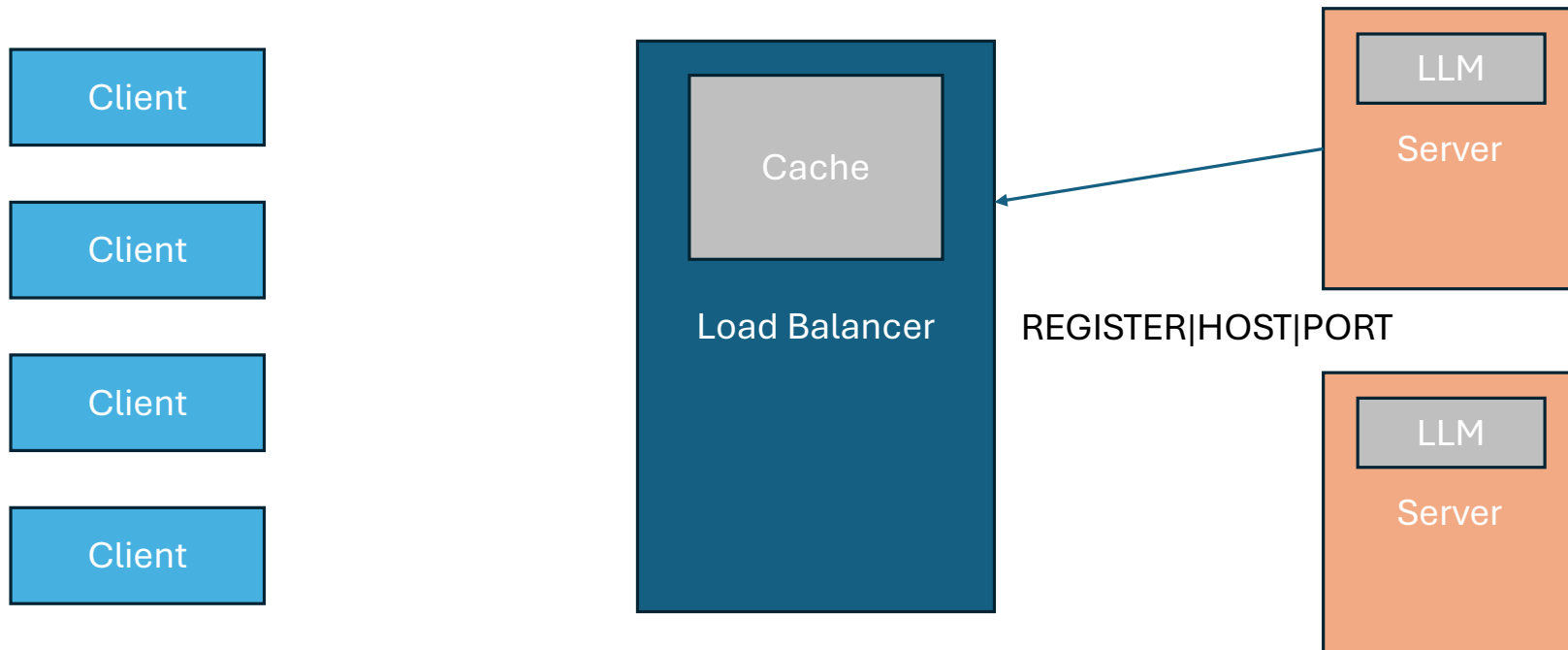
Add super simple UI



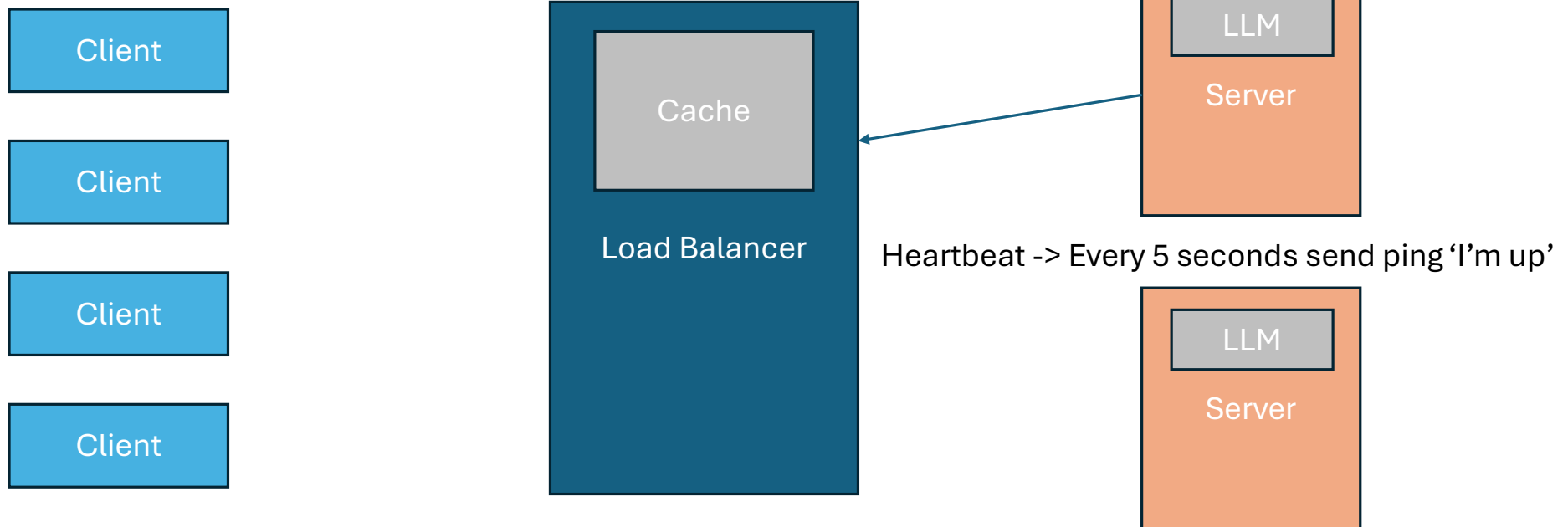
# Protocols:



# Protocols:

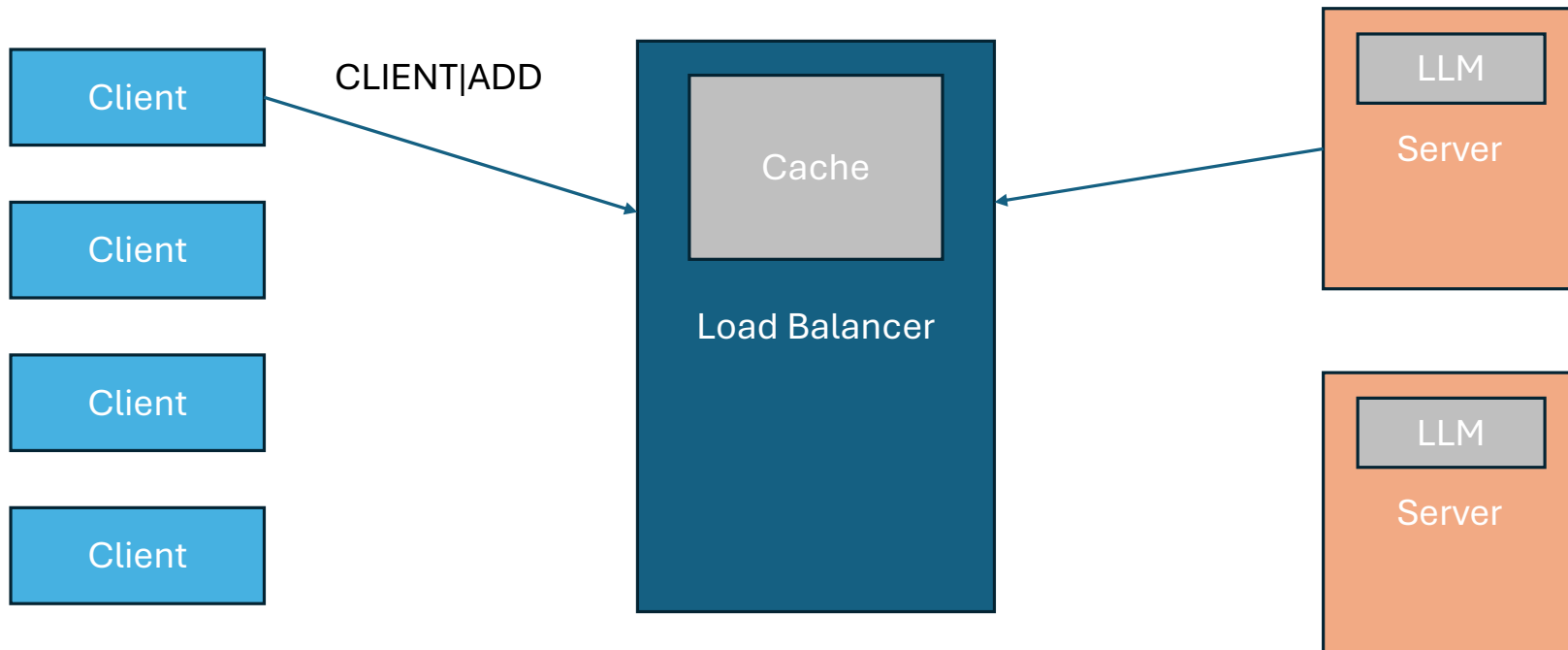


# Protocols:

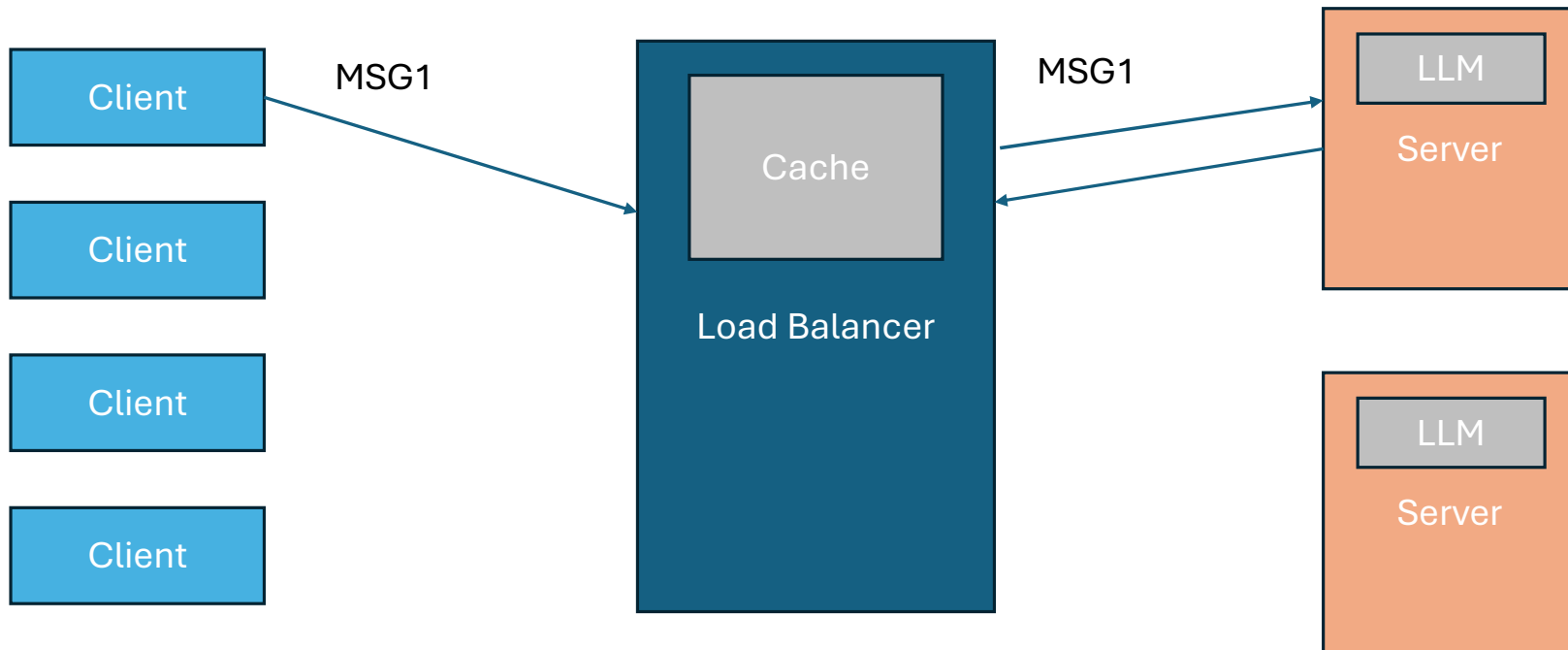




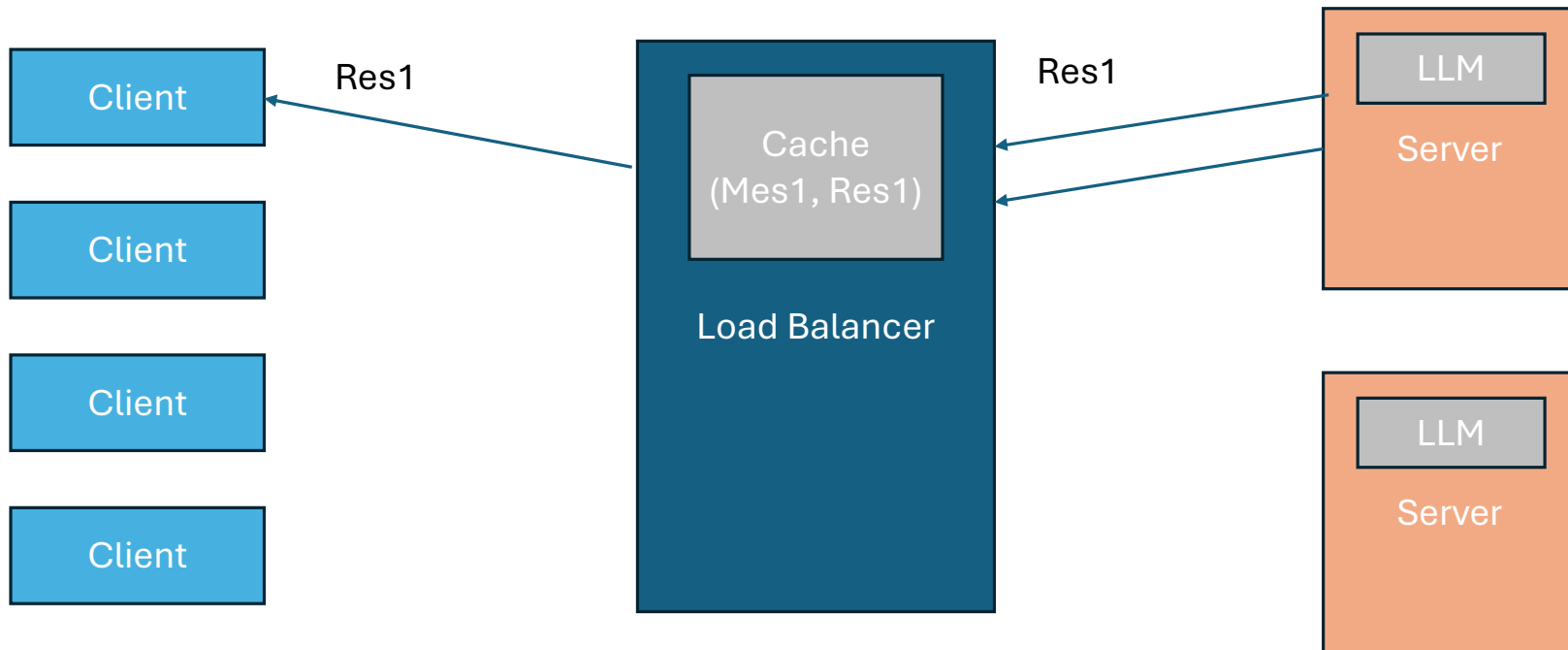
# Protocols:



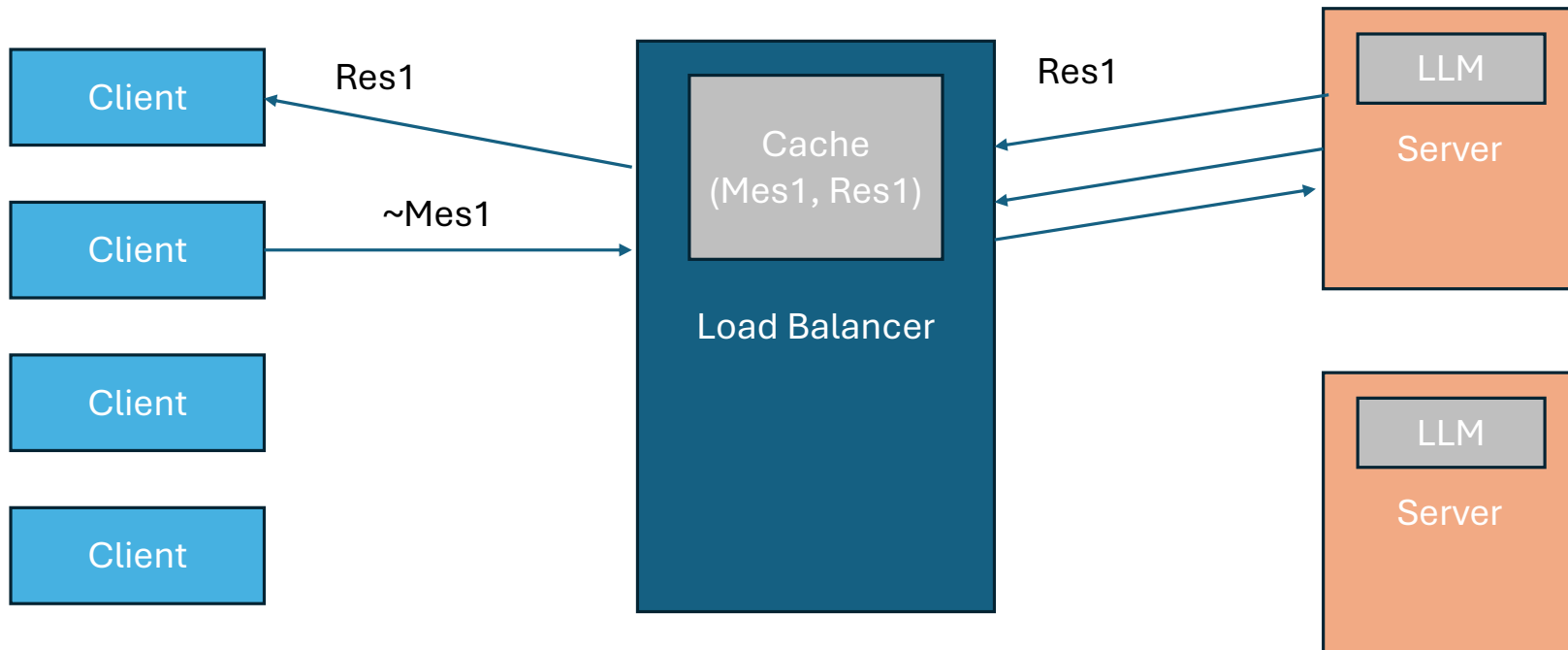
# Protocols:



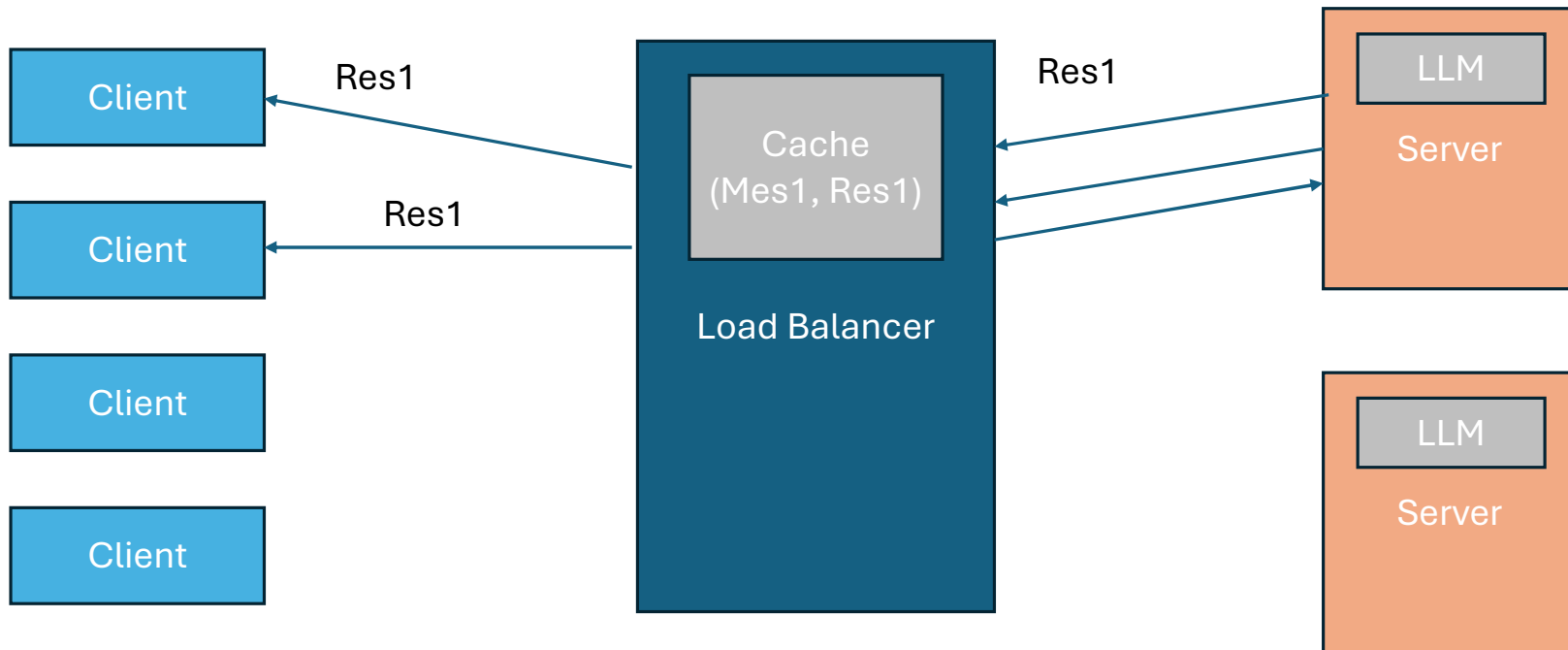
# Protocols:



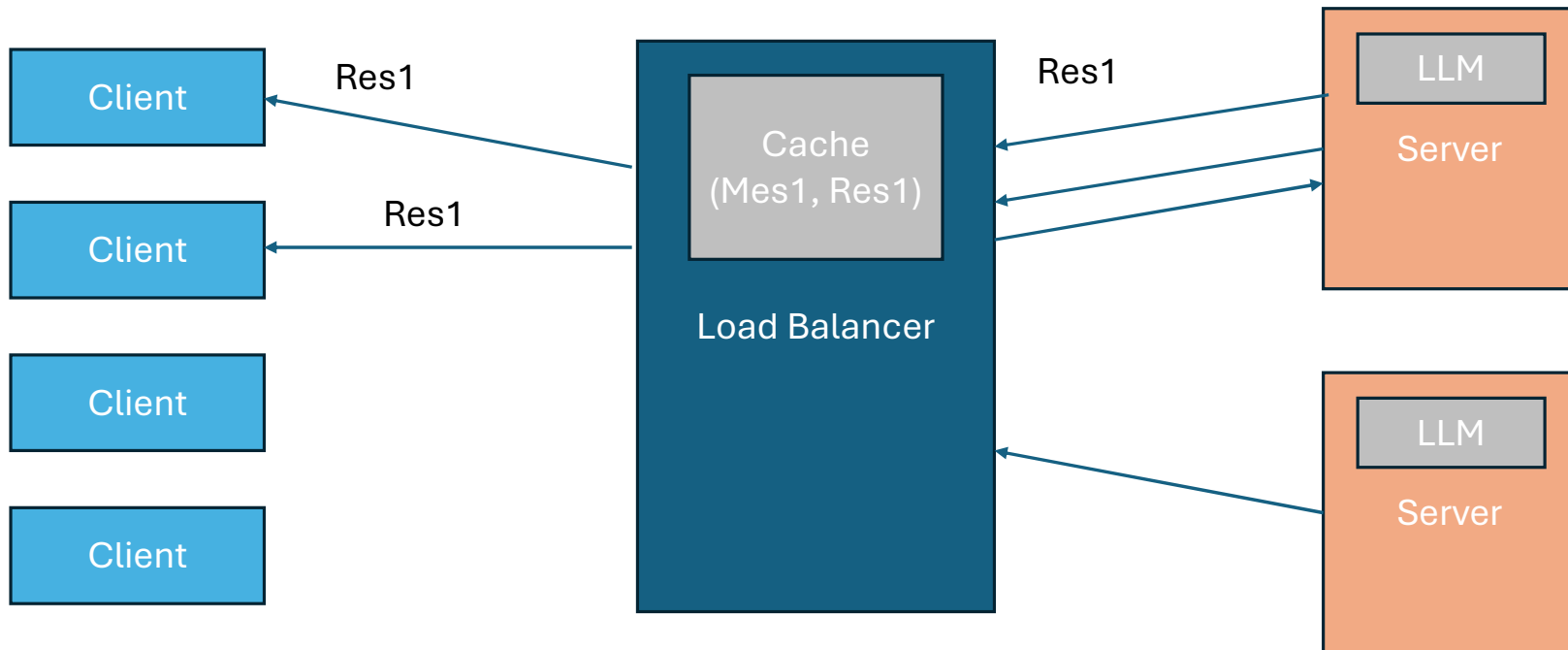
# Protocols:



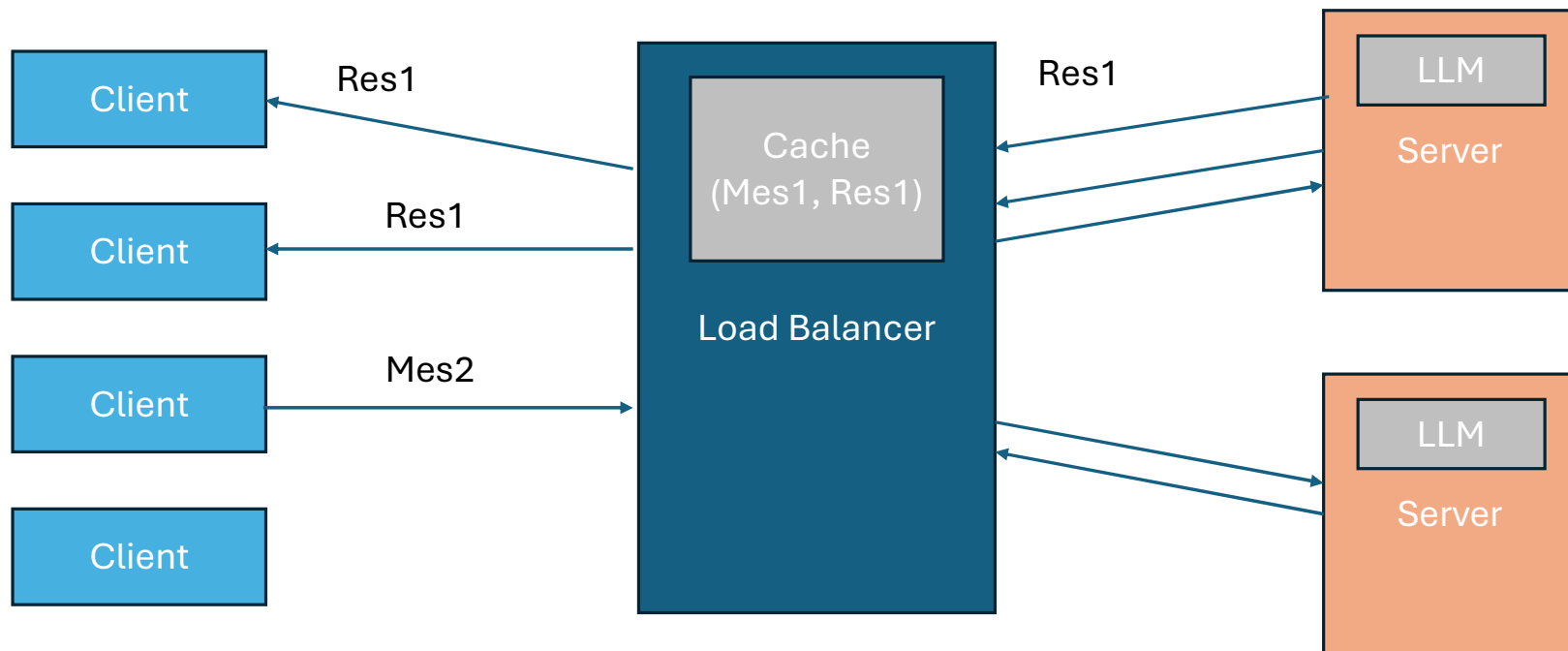
# Protocols:



# Protocols:



# Protocols:





Demo



# References

[Build Your Own  
Load Balancer |  
Coding  
Challenges](#)

[How to design a  
load balancer  
from scratch?](#)

[Let's Build! A  
Simple Load  
Balancer with  
Golang](#)