

Agenda

- Problem Statement
- Functional & Non-functional requirements
- API design
- Scale Estimation
- Architecture
- Sharding
- Bonus Material

High Level Design → *software*
birds eye view of the system
focus on the distributed architecture

facebook → laptop *X*
2B+ user 10B+ msg/day

software crash
power outage
n/w fails
n/w outage
diskfn.

10,000 computers → distributed across data centres → around the globe

communicate over network
↓
outage congested/slow

Staff Engg @ Google. → given a list of strings, sort them in acc order.
50 PB

SDE 2+ → top startup
→ MNC } at least 1 HLD round

HLD involves → conversational-

- ① Problem is vaguely defined
- ② No one correct solution → focus on tradeoffs
- ③ Conversation could go in any direction-

~45 mins

- ✓ → Problem Statement
- ✓ → Functional & Non-functional requirements
- ✓ → API design
- ✓ → Scale Estimation
- ✓ → Sharding
- Architecture

Problem Statement

- X Stock Broker → Zerodha / ICICI direct
- X Stock analysis → Trading View
- ✓ Stock Exchange → NSE/BSE NYSE/NASDAQ

Functional Requirements (MVP)

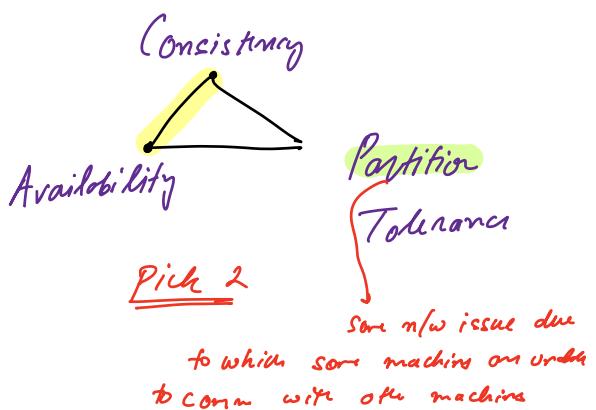
→ minimal set of features so that the product is usable

- ✓ → trade stocks : buy / sell / cancel / modify
 - ↳ various order types
 - ↳ matching Order
- ✓ → Real time price data
- ✓ → Market Reports
 - ↳ past orders
 - ↳ price history
 - ↳ overall market stats

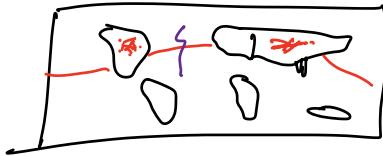
Non-functional Requirements

- NSE down → Availability !!
- Orders must be correctly created without errors → Consistency !!
- multiple servers → all servers in one place to distribute load

CAP theorem



Connected via v. high "speed
LAN (100Gbps)



giving up partition tolerance

- Fairness → if UI's order was placed before UI's order & both have same price then UI's Order should be fulfilled first.
- v.v. low latency
- Process thousands of trades / sec



API design → specify for endpoint name, input, output

↳ how does the external world use our system
brokers, traders stock exchange

Order processing

① POST /orders
↳ Create

Request
 $\left\{ \begin{array}{l} \text{Broker-order-id: } \text{---} \\ \text{stock: Reliance} \\ \text{price: } \text{---} \\ \text{quantity: } \text{---} \end{array} \right. \quad \left. \begin{array}{l} \text{type: } \text{---} \\ \text{timestamp: } \text{---} \\ \text{user-id: } \text{---} \end{array} \right\}$

Response

Success / Failure

$\left\{ \begin{array}{l} \text{type: } \text{---} \\ \text{price: } \text{---} \\ \text{quantity: } \text{---} \end{array} \right\}$

② PATCH /order/{id}
↳ update

③ DELETE /order/{id}
↳ cancel order

④ GET /orden/*id*
Orden details

(5) GET /orders → order book details

Market Data

① GET /stock / <symbol>

↳ current price

Reports

① GET /report/<date>

↳ total order value & amount
transacted on that day.

Scale Estimation

Back of the Envelope calculations

→ # users

total: 70M

India → 1.42 people

5% of Indians have a demat account.

DAU: 14 M

20% people who have an account trade on a given day.

→ # orders per day → Pareto principle
80-20 rule 20% people → plan trades
80% → monitoring.

3 million people actively trading on a day

↳ each active trade places $\underline{\underline{5-10}}$ orders

Orders per day $\approx 15M - 30M$

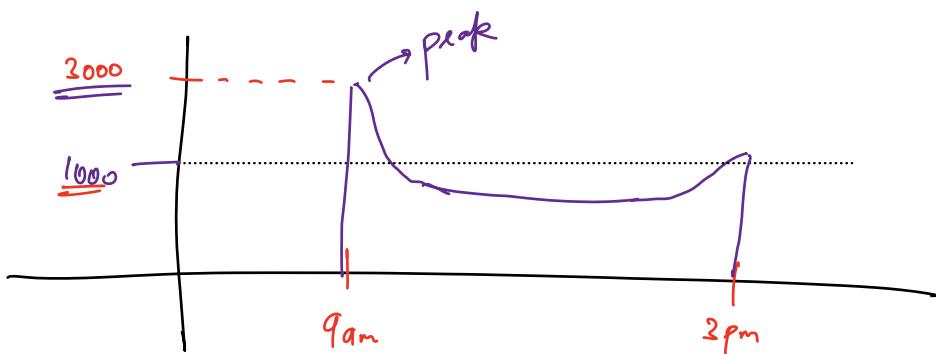
NSE from data $\approx 20M$ trades per day
↓
orders

daily volume = 2.6 trillion

20M orders / day

→ # orders / second ?

$$\frac{20M}{\cancel{60} \times \cancel{60} \times \cancel{24} \overline{6}} = \frac{20 \times 10^6}{3600 \times 6 \approx 20,000} = \frac{2 \times 10^7}{2 \times 10^4} \approx \underline{\underline{1000 \text{ orders/sec}}}$$



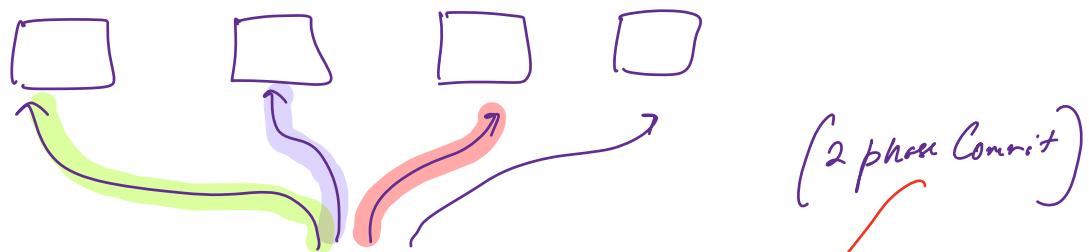
System should be able to handle $\underline{\underline{5000 \text{ orders/sec}}}$

Sharding

partition data across servers

→ multiple computers

→ v.v. high throughput (5000 req/sec)

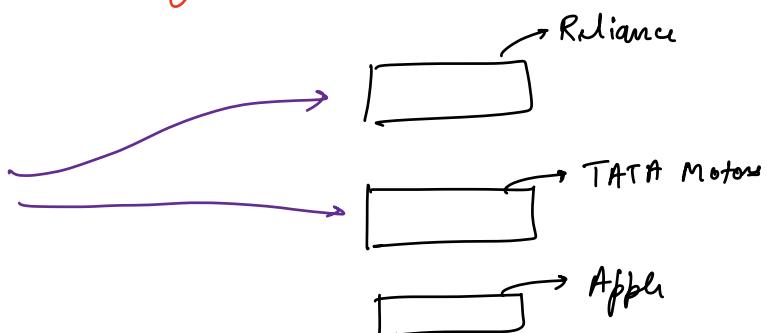


if we need to talk to multiple shards for 1 req
↳ v. high latency.

(all data)

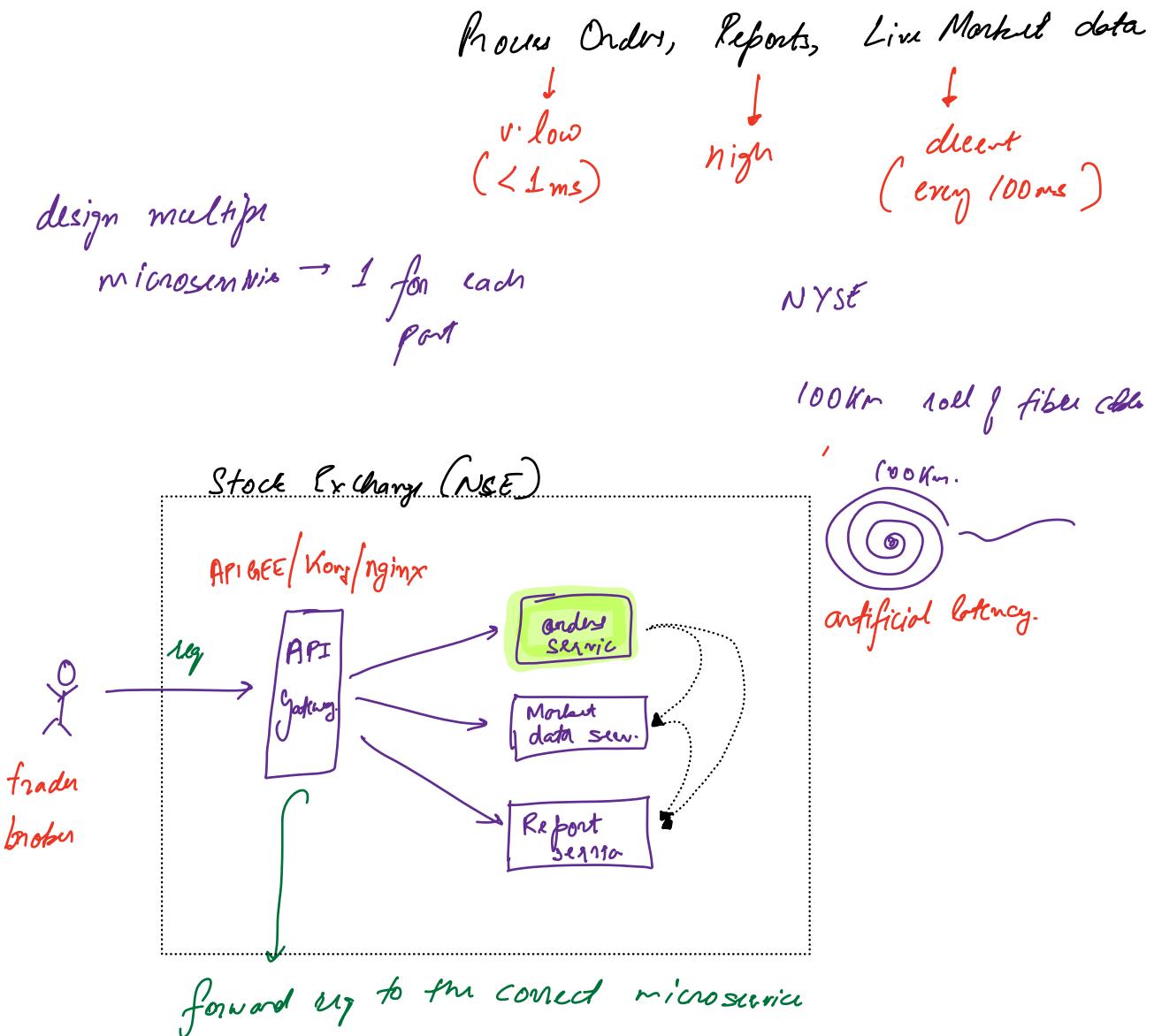
goal → everything that an order needs to be processed
must be on a single shard

Shard by stock symbol.

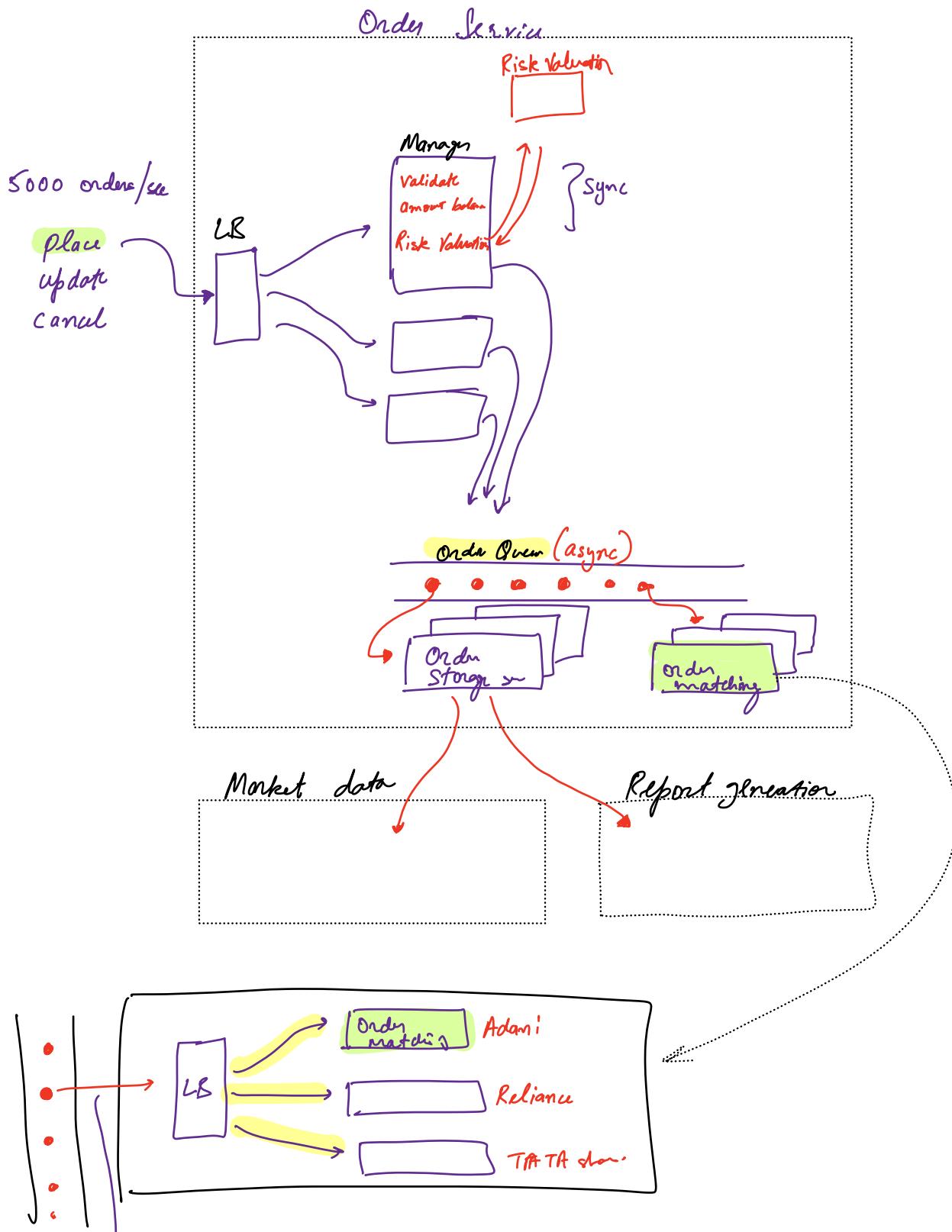


6:31 → 6:40

Architecture



Order Service



↓
 {Symbol:
 amount:
 price:
 type:
 }

designed for comm
over n/w.

Kafka / RabbitMQ / SQS

optimized for distribute
cloud architec.

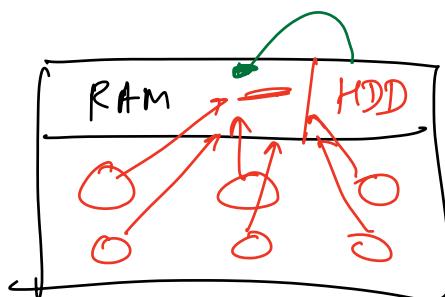
any stock exchange
will develop all these
solⁿs in hours & deploy

on same server

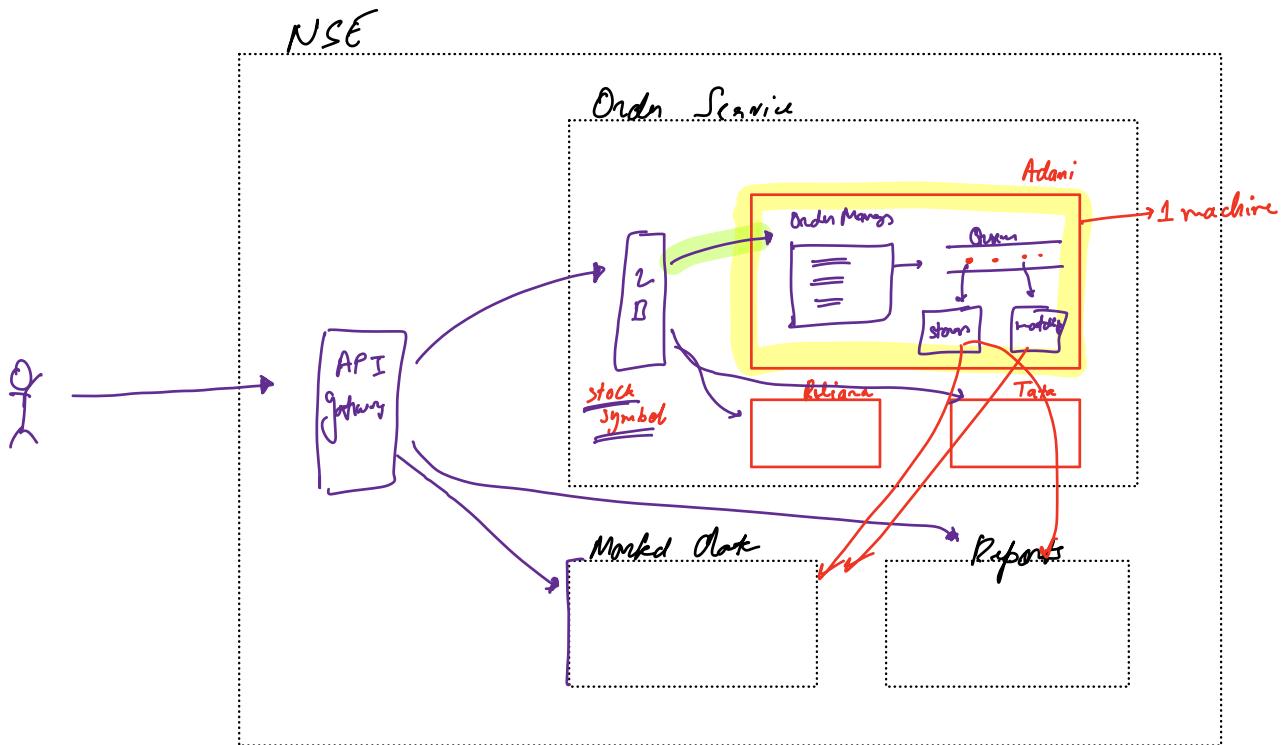
↳ services communicate via

v. high latency
from the perspective of
Stock Exchange !!

IPC (OS provided utilities)



Memory Mapped files



Robustness

Stock exchanges → 9's availability.

99.9999% time

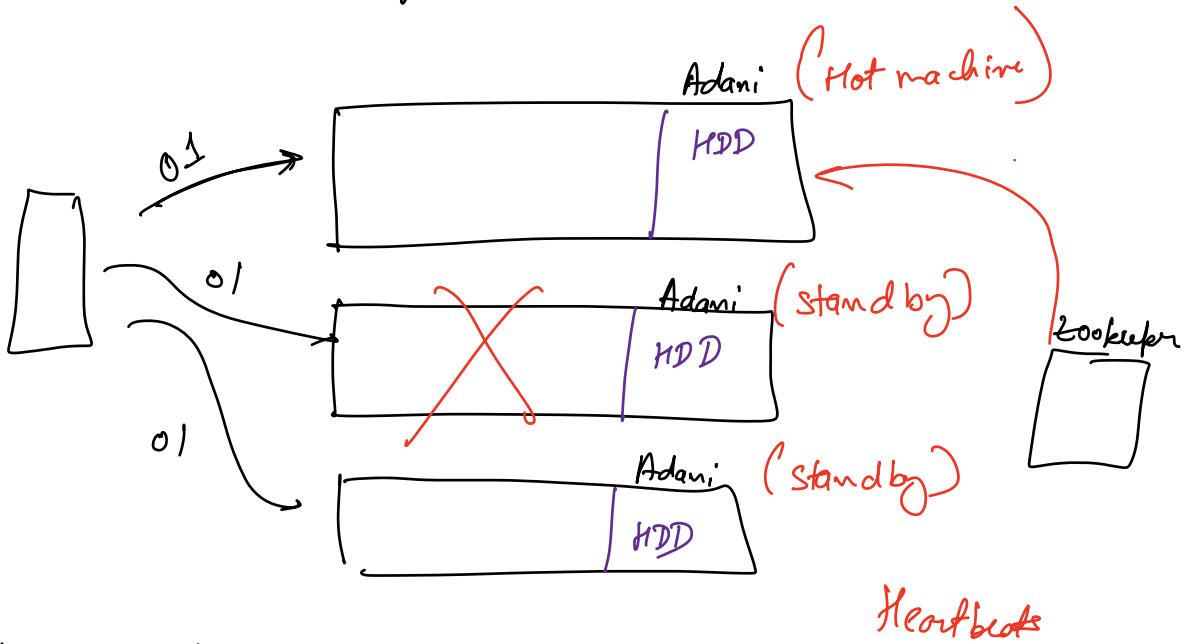
0.0001% time

0.0001% of 24 hours =

$$10^{-6} * 3600 * 24 = 0.08 \text{ seconds.}$$

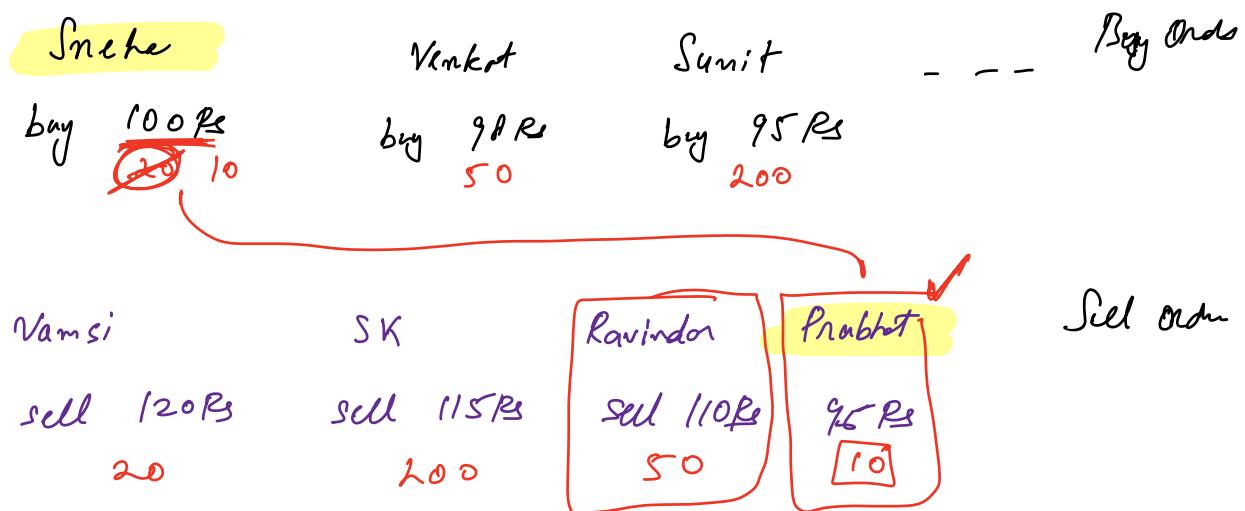
~6 seconds of downtime per day (NSE)

Hot machine & replicas

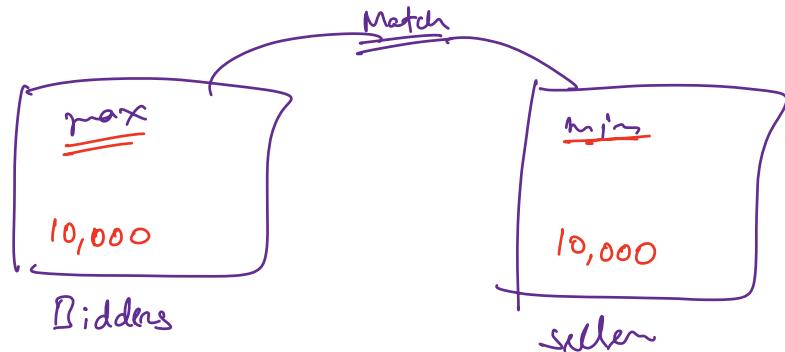


at end of day → persist in long storage
 Heartbeats → Cassandra.

Matching Algo → hidden DSA problem



match highest bidder with the lowest seller



linear scan $O(n)$



$O(fg^m)$

LRU cache

Hashmap + DLL

$O(1)$