

Graphite

- ↳ Git is not only the tool to manage code changes
- ↳ gregory sroc \Rightarrow mesherical best

V1

All Basic

Juno \rightarrow Distributed Key Value DB

- ↳ used in so key services at PayPal

↳ core backend service \rightarrow auth-backend

↳ efficient store & cache data

\rightarrow it's not in memory all the time

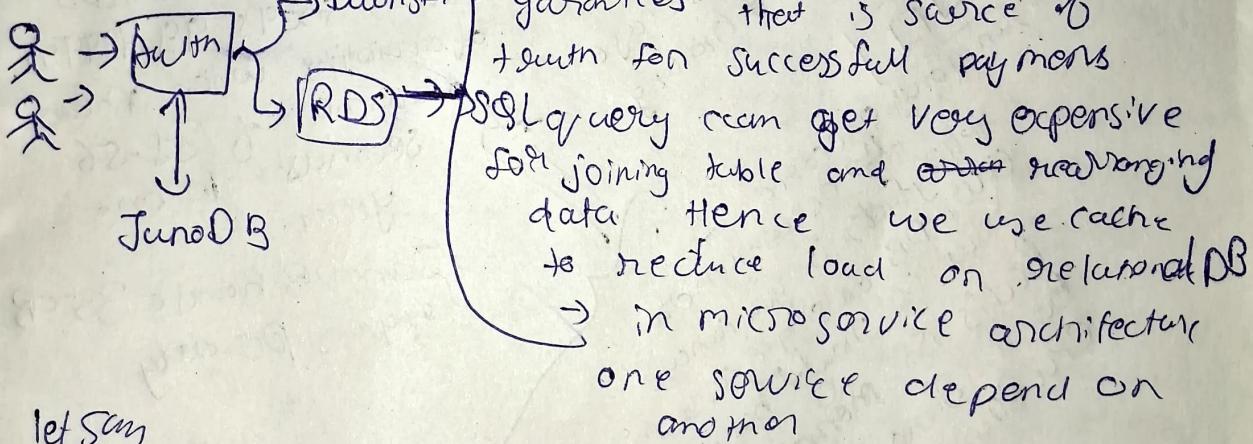
\rightarrow still it can cache data for you
it has persistent storage as well

↳ it reduces load of Relational

Database because of its assets

↳ guarantees that is source of

truth for successful payments



let say

Auth depend on another

microservice (Downstream) - - - Making large no. of call can make it go down - - - Juno DB can cache (But Doubt Downstream can have feature of caching) - - - why to increase load on server - - - Juno DB can help in heavy caching system

\rightarrow reduce help reduce SQL expensive like aggregation and joins.

Redis was written in single threaded C program

Juno DB

Started in single threaded C++ program

Rewritten in golang

↳ Highly efficient

↳ Highly concurrent (Goroutines)

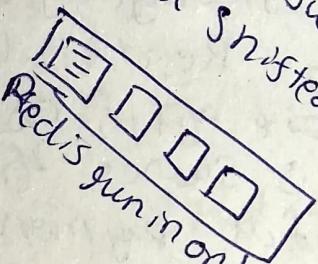
Goroutines are lightweight

threads ↳ Fibers

→ we can create many fibers to get things done

Parallelly

→ far more efficient than traditional POSIX threads which then multiplex over kernel threads.



Because there is one core per node, Redis is memory bound. But not all Redis operations are memory bound. Some are CPU bound. Hence Juno DB handles 350B requests per day.

Because of this JunoDB provides 99.95 99.9999% of availability which means downtime of 31.56 seconds in one year.

JunoDB handles 350B requests per day

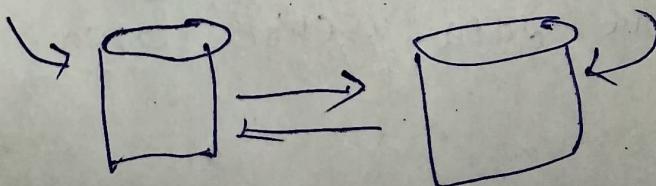
Juno DB use Both \rightarrow in memory & on disk cache
 \hookrightarrow TTL
 use cases \rightarrow from few seconds to few days.
 \hookrightarrow user token
 account details \rightarrow no need to hit auth on every request why to go to db JWT then do verify then do DB across check why not do cache of user details
 API responses \rightarrow no need to fetch account detail every time
 \rightarrow no need to give complete responses every time

user preferences eg transaction history, recent activity
 \rightarrow notification preferences

+ idempotency \rightarrow avoid duplicate processing
 no matter how many time the request come
 eg suppose while transaction state did not get updated B that why another request and same time you done processing
 But you want amount should get transferred hence no matter how many time request come you do only one processing

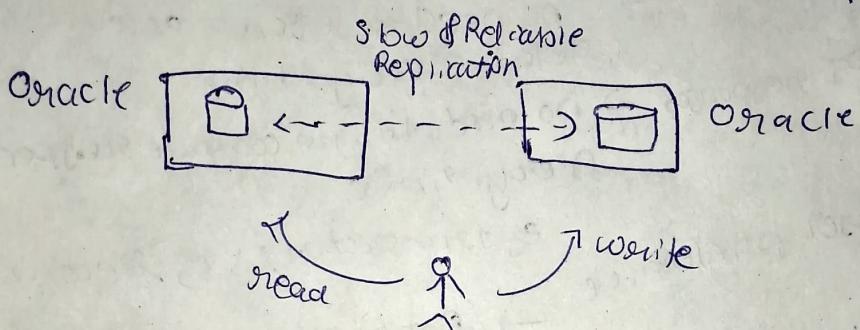
eg you can't drop notification twice though will confuse user that payment got transferred twice

+ latency Bridging : JunoDB has really fast inter cluster replication
 near instant Backup in another cluster
 Because

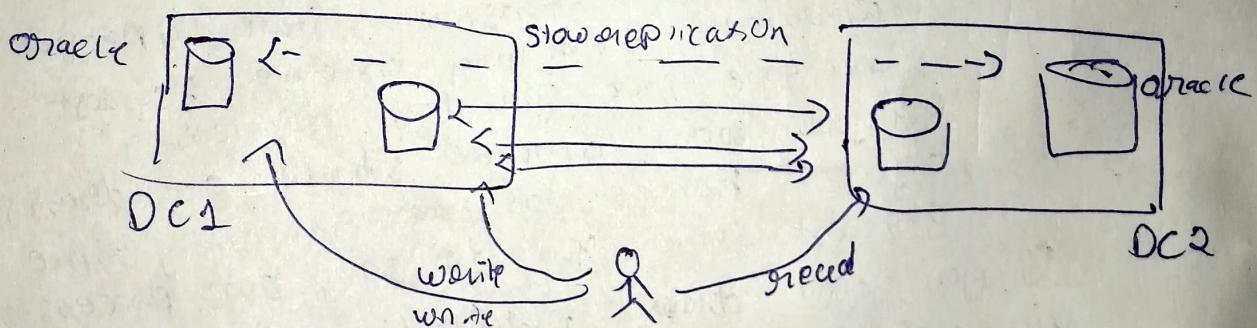


fast replication requirement when suppose
 1 read go to DB Cluster 1
 write go to DB Cluster 2 near consistent
 we will have inconsistent Broken cluster

PayPal uses Oracle as Primary DB
 they have active active setup across datacenters



replication lag Btw 2 Slow But reliable Oracle's DB
 is pain point so they do



Video 2 System Design

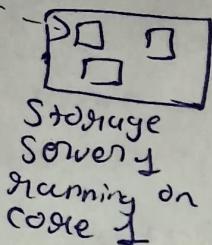
High level architecture & System Design of Storage Server

Storage Server → Inside DB --- Inside the DB the server which is used to do get put Post Delete operation Hence Juto DB is not purely implemented if they do this can store data in disk or on memory

Data is split into Partition (also called Shards).

→ this Shards converted in → Rock DB

Shards

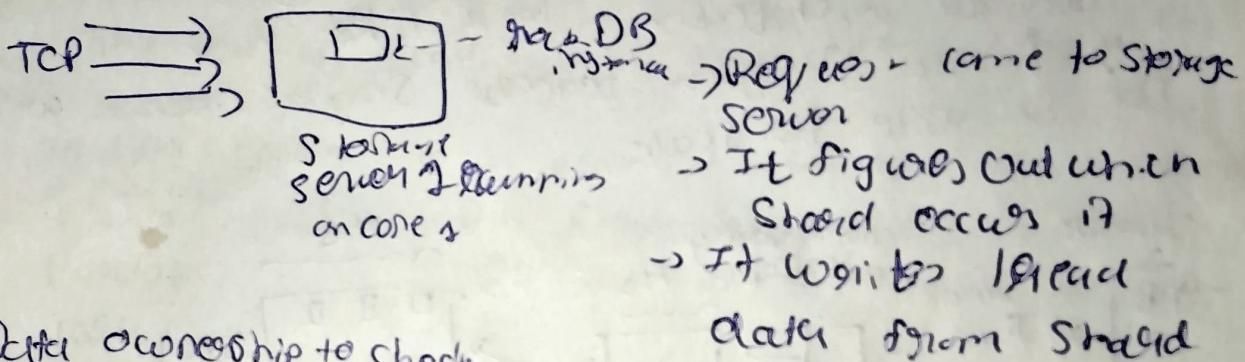


SS 2
C 2

SS 3
C 3

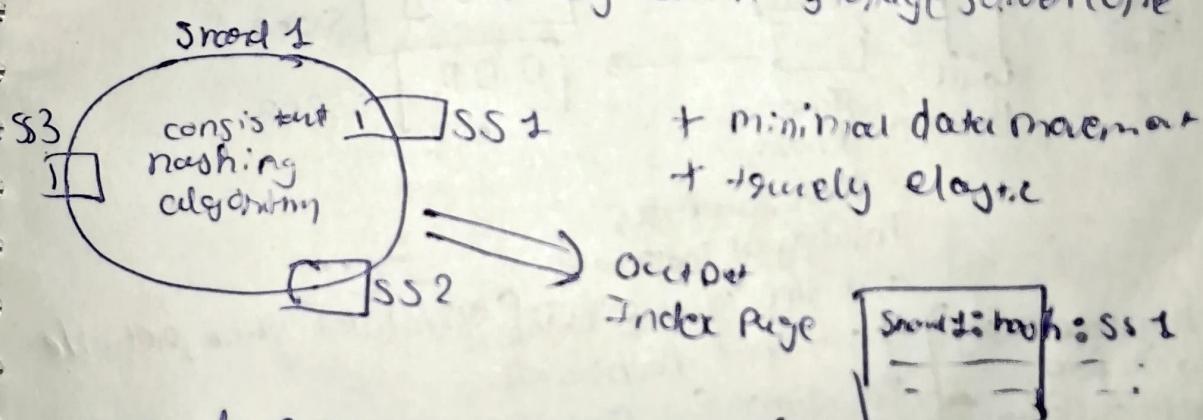
(very famous embedded DB used to embed data)

and given for Get, Put, Post, Delete operation on server

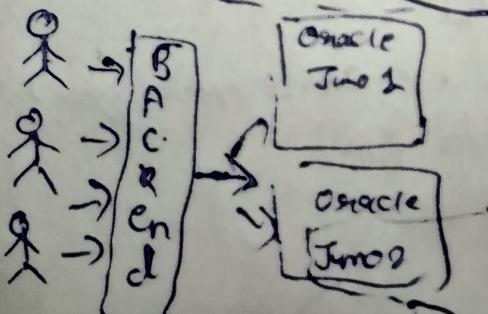
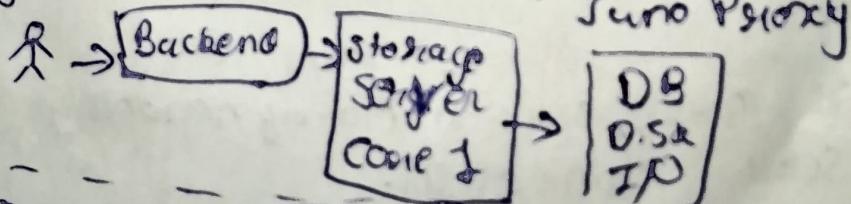


Diff ownership to check

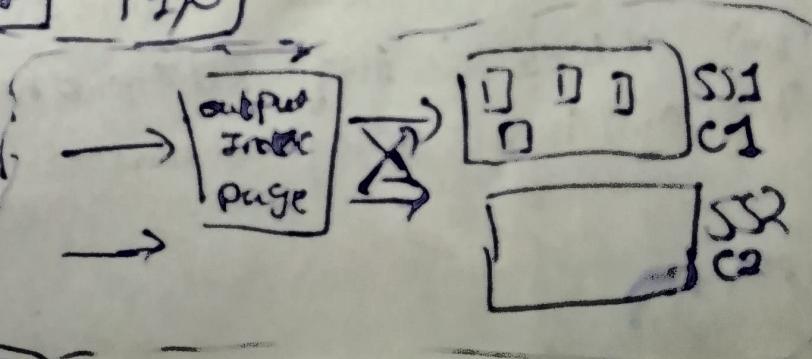
Which Shard owned by which Storage server core



current Redis System



Oracle Junos

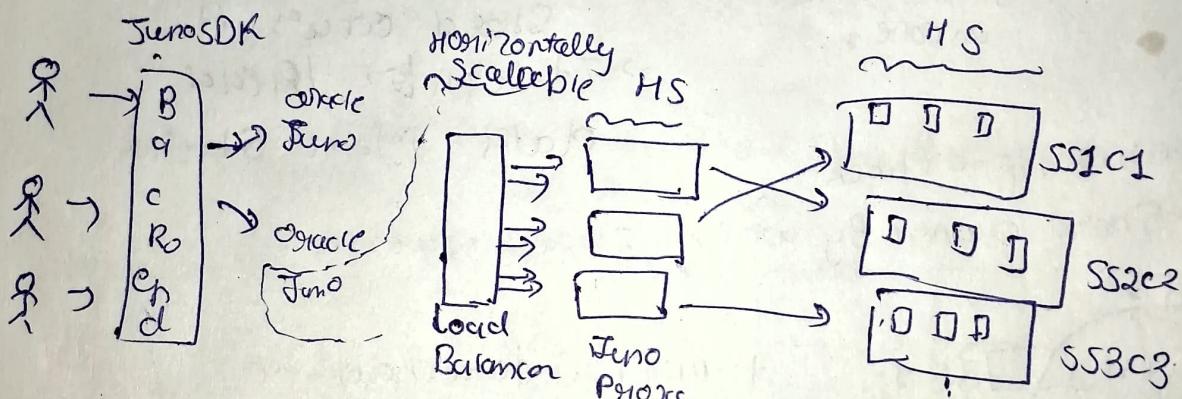


Advantage

- Storage topology is abstracted
- Client need not connect to all storage servers
- minimal config overhead

is proxy at single machine?

No Juno can have multiple JUNO proxy which can be put behind LB. Overall architecture look like this. This feature is given because single machine will not be able to scale.



to store key value mostly certain data for distributed systems

Each Juno Proxy has copy of index page and they run ETCD [distributed store, reliable & consistent hashing all the time for each request]

Juno Proxy is Stateless, equal queue in no. of handling connection they can scale seamlessly.

it is Stateless \Rightarrow they all should be same replicas of each other and have same information all the time.

V3 Scrolling

why Juno DB needed to scale

as no. of microservice increase

↳ the no. of inbound connection / ^{network} increase

to get more performance there is a persistent connection maintaining between client to Juno Proxy & Juno proxy to storage server

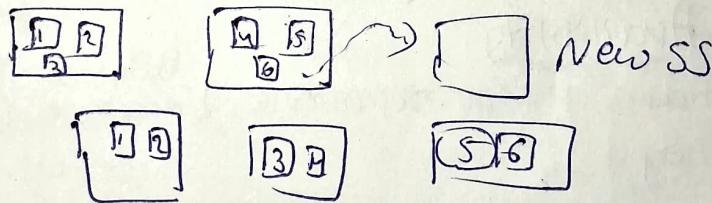
as no. of microservice ↑ no. of concurrent connection ↑

↳ the amount of data to be stored increase ⇒ storage

↳ the no. of read write increase ⇒ compute

as data grows it is essential to add more storage servers Juno DB since the number of partition (say 1024 and we can't change after that) 1024 shards

this 1024 shard are distributed to other servers by consistent hashing



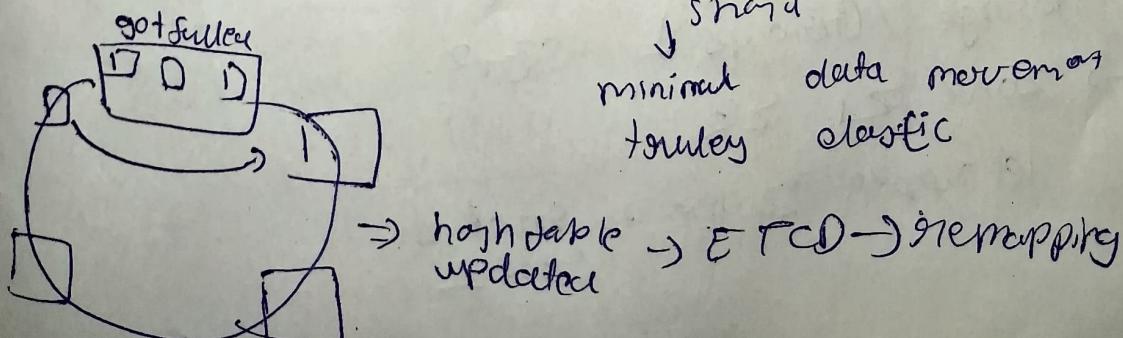
→ this is done by consistent hashing

→ when distribution of shard done control reaches to ETCO and updates the index page, mapping

ETCO shard i → storage server

↓ shard

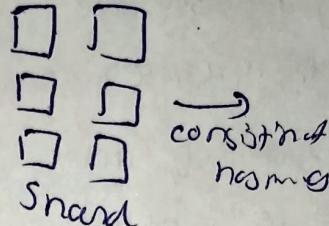
minimal data movement
fully elastic



Each shard has a bunch of microshard

when a shard is moved from one node to another
all micro shards are moved as a shard moves
micro shard are building Block of data

Key \xrightarrow{f} high \rightarrow %Shard
Minimum Shard



because
of no. of shard
is fixed we can
do mod shard

Total no. of shard are fixed and specified
cluster creation

\rightarrow Number should be large enough
these shard will be distributed across
storage nodes

Paypal has 200 storage nodes to process
100B req every day

V4 99.999% Availability

Paypal can't have poor experience because it deals
with money

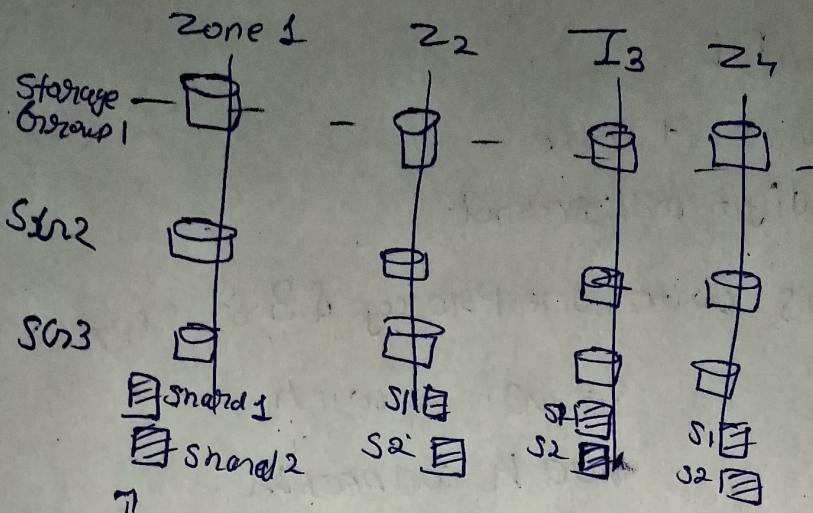
dependency
(compute +
storage) \rightarrow highly
available

It arranges the storage nodes in a grid

rows \rightarrow Storage Groups 1 Sgr1 2 Sgr2

column \rightarrow Zone Z1 Z2

Let suppose Zone as Rack 1
Rack 2



this for
snapshots
are replicated

swrite

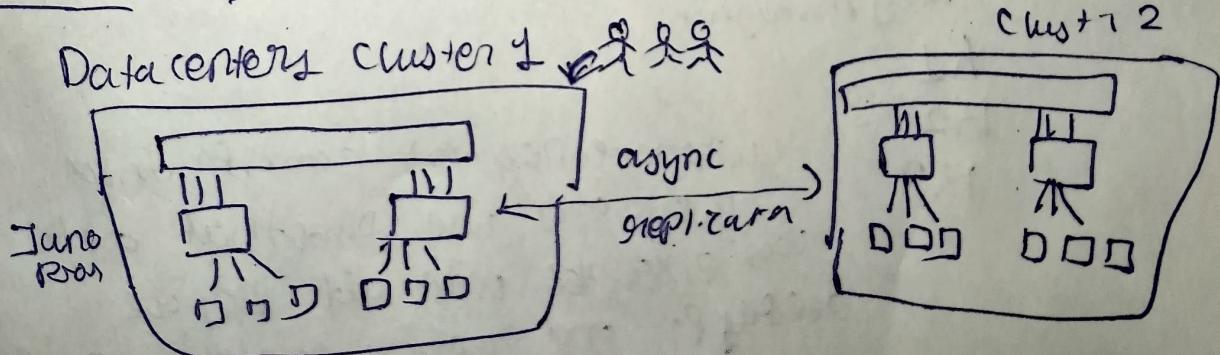
write ← Juno proxy will
successful wait for acknowledgement
from majority

read → SG3 → 1 2 3 4 5

because of
this beautiful
architecture PayPal
can take whole
zone for routine
health check, maintenance

always ← Juno proxy will
return wait from majority
dates + value

New Cross Data center Replication



Performance of JunoDB

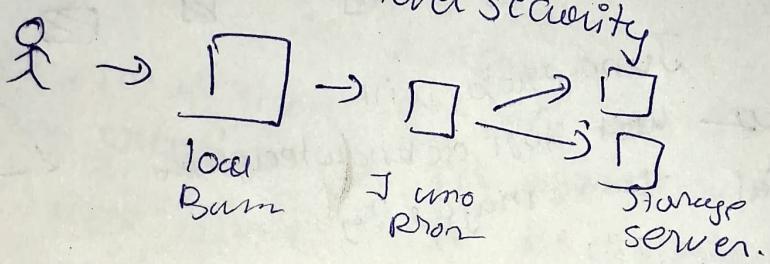
↳ performance + Scale

Single digit millisecond.

P99 of 4 ms with one proxy & 3 storage servers with 100 K connection & 5 RPS

JunoDB is secure for transit & at rest

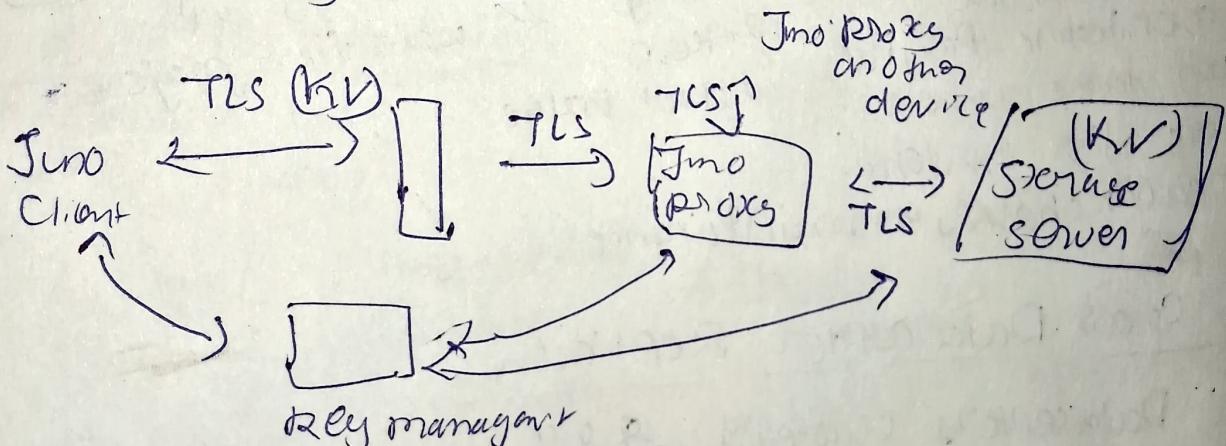
TLS transport level security



Ideally Client should do encryption

If client do not do encryption

Juno Proxy does it



K1

K2

K3

Whenever newer comes to that data and we find that data was encrypted with old key we decrypt and re-encrypt with new key