

Design an Online - Streaming Application Like HOTSTAR / NETFLIX

Step-1:- Requirement Gathering

① Functional Requirement

- (a) User are able to search his Movie / serial / Web series etc.
- (b) User able to add their favourite thing in watchlist.
- (c) Get details of the serial / Movies / Web - Browser
- (d) Watch History Maintenance
- (e) User can give rating to movie / serial / webservices etc.
- (f) Personalized Recommendations
- (g) User can buy Premium Membership and can cancel membership.
- (h) User able to register / login.

OR

- (a) Search Functionality
- (b) Watchlist Management
- (c) Detailed Information
- (d) Watch History
- (e) User Rating System
- (f) Personalized Recommendations
- (g) Buy / Cancel Premium Membership.
- (h) Registered / Login System.

② Non-Functional Requirement

- ① My System should be Highly Available.
- ② My System should be Highly reliable.
- ③ My System should be secure
- ④ My system should give high performance.
- ⑤ My System should be Scalable

Note:- High Performance means Low latency Time.

Why High Performance or what High performance mean here in Online streaming Application:-

- (a) Streaming should be fast
- (b) Instant Recommendation feature
- (c) Smooth / Interactive User Interface
- (d) System perform at regular speed even in large-scale traffic.
- (e) Fast Response Time
- (f) Immediate Action Execution
- (g) No waiting time.

Scalable:- App does not crash even user increase at large rate on App. Means if daily active user is increased by x% than normal numbers, my app must handle the load properly.

Reliable:- Means Minimizing the failures.

Step -2:- Prioritization

(a) Database Population:-

1.) **Primary Content Addition:** Data of Movies, Web series, TV

Shows are stored in Database.

2.) **Data structuring:** Schema Design to check and track metadata, genres, actors, ratings and availability.

3.) **Dependency Resolution:** Our this content data will be foundation for Search Recommendations and Subscriptions features. So their dependencies should be defined.

4.) **Scalability Consideration:-** Optimized storage and indexing is used in huge data maintenance and storage.

(b) User Management System:-

1.) **User Account Service:-** Maintain User Profile Data and their actions.

2.) **User Authentication Service:-** Authenticate User

login and registration

3.) User Authorization Service ← optional
not much required (But if we want to make thing more systematic, we can add this feature).

(c) Search Feature

Efficient content discovery is ensured in this step.

(d) Watch History and Favorite Content:-
User engagement tracking here is ensured in this phase.

(e) Personalized Recommendation :-

AI/ML-based suggestions are involved in this phase to build this feature.

(f) Premium Membership & Cancellation:-
Monetization and user control features are ensured in this phase

In prioritization, I focussed to build independent and most essential feature to build any application i.e from Top priority to the lowest priority, I define my phase work

Step-3:- Infrastructure Estimation / Back-of-the-Envelope Calculation / Feature-in-Hand Planning

(a) Storage Requirements:-

Let, Daily Active User = ± Billion

Next Year Estimated User = 0.5 Billion

Total Users = 4 Billion

75% users are assumed to be inactive users in Online

streaming application. (95 is assumed)

Assume data of user be \pm MB

Total Data for user in Next year

$$= 4.5 \times 10^9 \times 10^6 \text{ Byte}$$

$$= 4.5 \times 10^{15} \text{ Byte}$$

$$= 4.5 \text{ PB}$$

Data for content:-

Let number of movies (expected upto next year) on app be 200000 . Let 36B on on average each movie must contain the data.

$$\text{Estimated Data for movie} = 36B \times 2 \times 10^9$$

$$= 3 \times 10^9 \times 2 \times 10^9 \text{ Bytes}$$

$$= 6 \times 10^{14} \text{ Bytes}$$

$$= 600 \times 10^{12} \text{ Bytes}$$

$$= 600 \text{ TB } (\because 1 \text{ TB} = 10^{12} \text{ Bytes})$$

Let, number of serials (expected upto next year) on app be 10,000 . And each serial on on average consist of 750 episodes then

$$\text{Estimated Data for serial} = 750 \times 10^9 \times 10^9$$

$$= 7.5 \times 10^{15} \text{ Bytes}$$

$$= 7.5 \text{ PB}$$

$$[\because 1 \text{ PB} = 10^{15} \text{ Bytes}]$$

Note:- Average size of each episode is taken as 1GB i.e. of 45 minutes. each on average

Now, let web series upto next year be 5000, each web series of each episode containing content for 60 minute on on average

On on average each web series contain 10 episodes each web series contain 3 seasons on on average

$$\text{Each web series size} = 30 \times 2 \text{ GB} = 60 \text{ GB per Web Series}$$

$$\text{Total data required for web series} = 60 \times 10^9 \times 5 \times 10^3$$

$$= 3 \times 10^{14} \text{ Bytes} = 300 \text{ TB}$$

Final Approximate Storage Estimation

for Data Content Estimation for storing video :-

$$8.5 \text{ PB} (7.5 \text{ PB} + 0.6 \text{ TB} + 0.3 \text{ PB})$$

But, we apply some compression algorithm like lossless decomposition that reduce data need by 20% easily without quality loss.

Final storage requirement for storing content

$$= 80\% \text{ of } 8.5 \text{ PB}$$

$$= 0.8 \times 8.5 \text{ PB}$$

$$= 6.8 \text{ PB} (\text{approx})$$

$$\approx 7 \text{ PB}$$

for new content, we have taken $\pm \text{PB}$ as buffer storage

So, Total data for Distributed File System = 8 PB

Metadata storage Estimation :-

Let $\pm \text{MB}$ is required to store basic detail about movie, serial, web series i.e:-

- ① Name ② Release / Ending Date
- ③ Actor Cast ④ Producer / Director Information
- ⑤ Channel Detail (for Serial)
- ⑥ Link to the path where content is stored

for Movies:- $200 \text{ K} \times \pm \text{MB} = 200 \text{ GB}$

for TV serials:- $10 \text{ K} \times \pm \text{MB} = 10 \text{ GB}$

for Web Series:- $5 \text{ K} \times \pm \text{MB} = 5 \text{ GB}$

Total Metadata Storage Required = $(200 + 10 + 5) \text{ GB}$
= 215 GB

Caching & Extra indexing can take up to extra data of 35 GB.

so, Total Metadata Storage Required = 250 GB

Data for storing User interactions & Activity logs for storing Watch History, Search Queries, Click events, personalized recommendations etc.

Estimation:-

Daily Active Users : \pm Billion

Every User :- 10KB /day activity generation

Daily storage :- \pm Billion \times 10KB = 10TB /day

Annual storage :- 10TB \times 365 = 3.65PB /year

Data for Cache storage (Fast Access data)

10-20% of total content data is stored in cache on average.

Now, Let 20% of data (Total data) is stored in cache.

Rough estimate :- 2PB for caching

Final Data Storage Summary:-

User Data :- 4.5PB

Content storage :- 10PB

Metadata storage :- 215GB (\approx 0.2TB)

User Activity logs :- 3.65 PB /year

Cache storage :- 2PB

Total storage :- 20.35PB (including all components)

Let , Total storage :- 21PB

As Buffer for future growth $\sim \pm$ PB
which is good for our system to be made scalable.

② Traffic Estimation

Major factors to measure it :-

- (a) DAU (Daily Active Users) \rightarrow \pm Billion
- (b) Peak vs Avg. Load
- (c) Requests per second
- (d) Read vs Write Ratio
- (e) Concurrent Users (f) Session Duration

Traffic Estimation is basically estimating the load of upcoming requests to our system.

(b) Peak vs Avg. Load Estimation :-

Let \pm hours /day is average time for each daily active user for streaming.

1 hour = 3600 seconds

1 Billion Users / 24 hours = 41.6 Million Users per hour

41.6 Million / 3600 sec = 11500 users per second

Avg Load = 11500 users per second

Peak Load = 4 times Avg Load

$$\begin{aligned} &= 4 \times 11500 \text{ users per second} \\ &= 46000 \text{ users per second.} \end{aligned}$$

Request per second depend upon:-

- (a) Homepage requests (Searching, Recommendations)
- (b) Video Requests (Streaming, Buffering)
- (c) Database Queries (Watch History, Favourites)

Using Traffic Estimation to find:-

- (a) Server Capacity Plan
- (b) Scalability Planning
- (c) Caching Strategy

Concurrent user & RPS (Request per second) Load estimation:-

Peak concurrency factor ≈ 1 to 5% of DAU
for \pm Billion DAU, Let ± 1 be the peak concurrency factor i.e 10 million users.

Let, each user do 5 operations per minute

So, Per Minute Request = $10M \times 5$ Request
 $= 50M$ Request per minute

Per Second Request or RPS = $\frac{50M}{60}$
 $= \sim 833K$

Estimated RPS for Worst Case

Read Vs. Write Ratio

By Analysis, Application will be ready-heavy.

Major Read Operations :-

- ① Homepage Load :- Recommendations, Trending Videos
 - ② Search Query ④ Watch History Fetch
 - ③ Streaming Request ⑤ Favorite List Fetch.
- Minor Write Operations:-

- ① Give Rating ② Favorite List add/delete
- ③ Updating Watch History

Estimated Ratio: 95% Reads & 5% Writes

Throughput Calculation:-

Throughput means how much request a system can successfull process in a second.

Now,
Estimated RPS = 833K

Throughput = RPS × Average Request Size.

Let Avg Request Size = 10KB

$$\text{Then, Throughput} = 833K \times 10KB \\ = 8.3 \text{ GB/sec.}$$

If request is of video stream then request size can also be larger

Bandwidth estimation:-

Mostly used for video streaming estimation
In streaming, Bit rate matters a lot i.e

Video Quality	Bitrate (Mbps)	Size per hour (GB)
480p (SD)	1 Mbps	0.456 GB
720p (HD)	3 Mbps	1.35 GB
1080p (Full HD)	5 Mbps	2.25 GB
4K (Ultra HD)	20 Mbps	9 GB

Now, If 10million users are streaming at avg 1080p quality (5Mbps)

$$\text{Total Bandwidth Required} = 10M \times 5 \text{ Mbps} \\ = 50 \text{ Tbps (Terabits per second)}$$

The analysis done for Bandwidth estimation is for Worst-Case Analysis But in Reality by CDN i.e Content delivery, Network bandwidth requirement can be reduced 10-20 times.

Infrastructure Estimation Steps:- (Followed)

① Request Handling:-

- (a) Calculating RPS
- (b) Estimated 1000 CPU's to process 100 RPS per core (Theoretically a High End Machine in an ideal case of RAM 16-32 GB can handle 8.3 GB/sec Throughput single handedly, but in practicality it is not possible)
- (c) Auto Scaling required for peak loads
- (d) Video Processing (Encoding / Transcoding)
 - GPU-based Cloud instances needed
 - Parallel Processing to optimize Encoding Speed

② Network Bandwidth Estimation:-

- (a) 50Tbps for Peak load
- (b) CDN & Edge Caching to reduce direct load.
- (c) 8.3 GB/sec → Throughput Calculation for API requests & User Actions.

③ Scaling Strategy:-

Database Sharding, Horizontal scaling (Microservices + Load Balancers), Auto Scaling (Real-Time Resource Optimization based on Load), Vertical scaling (Efficient Resource Allocation on cloud instances), CDN & Edge Computing (Reduce global latency for video delivery).

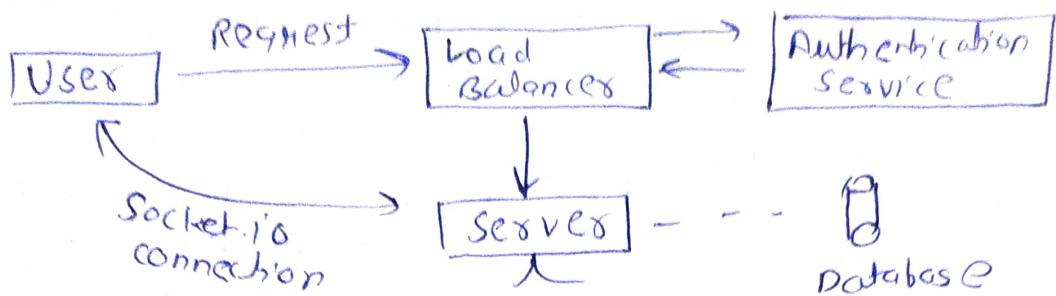
④ Traffic Estimation:-

- Estimate Daily Active Users → 1 Billion
- Estimate Concurrent Users at Peak Time → 10M
- Find out Application is Read-heavy or Write-heavy. → Read Heavy

In this step, Infrastructure Estimation is done.

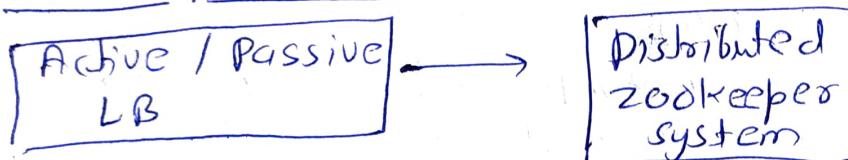
Step 4:- Component Identification

(a) starting with very basic design:-



Introducing Zookeeper so as to choose right server as it monitors health check-ups
 Maintain Zookeeper in distributed manner which are given some servers in which they send their request by their Health Monitor checkup

More Optimization:-



First depending on Load, Algorithm used in LB will be:-

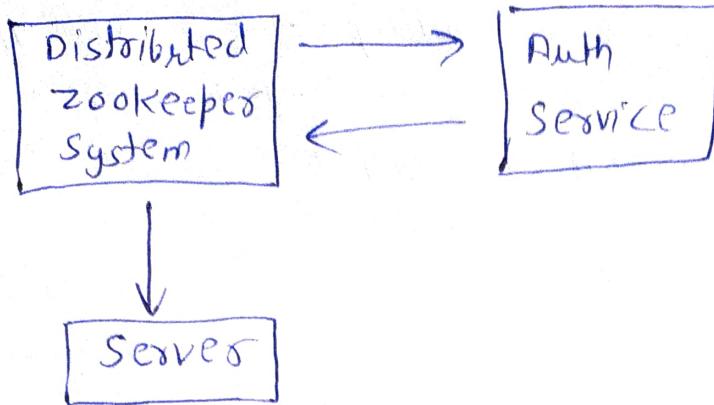
- (a) Round Robin Algorithm
- (b) Least Connection Algorithm
- (c) Weighted Round Robin Algorithm

↓
Request size Tses

Advantage of this Design:-

- (a) High Availability:- LB & Zookeeper both handle failure.
- (b) Scalability:- Round Robin approach (followed by Least Connection & Weighted Round Robin Algorithm if load Tses) evenly distribute load.
- (c) Fault Tolerance:- If Zookeeper failed, Backup instance will handle it.

Now,

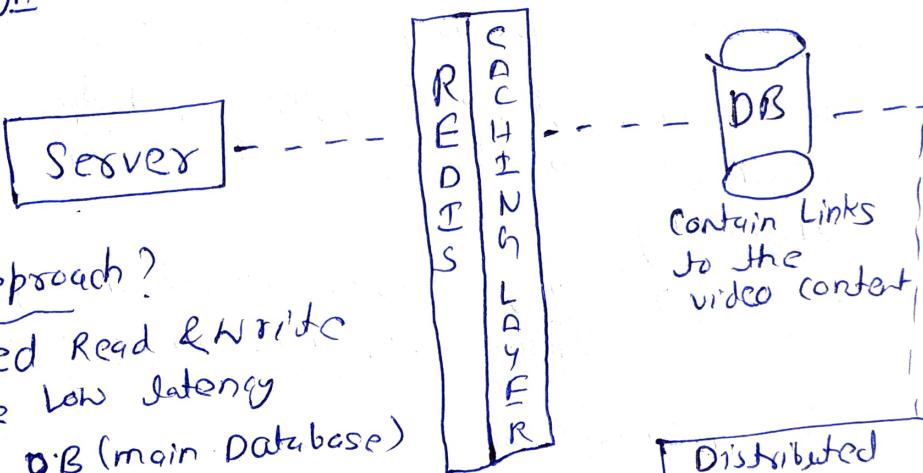


Advantage of This Approach :-

- (a) Secure Authentication
- (b) Efficient Communication
- (c) Low Latency
(socket.io bring low latency to the system)

zookeeper exit itself when communication is established b/w user and server, this optimizes the load

More Design :-



Why this Approach?

- ① High speed Read & Write
- ↳ Provide low latency
- ② Reduced DB (main Database) load.
- ③ Better Scalability, Multiple servers can connect / efficiently scaled easily.

Server taken here is basically a distributed server system

Distributed File System like HDFS
Efficient for storing videos/audios.

Note:- Video is streamed on server to user app in the form of small video packets. We can use CDN or Web socket connection for it

Load Balancing Chosen → L4 LB

Why Chosen L4 Software LB?

- (a) Cost effective
- (b) Faster Routing due to TCP/UDP
- (c) Better Scalability (Horizontal Scaling Possible)
- (d) High performance (High Throughput)

More Design :-



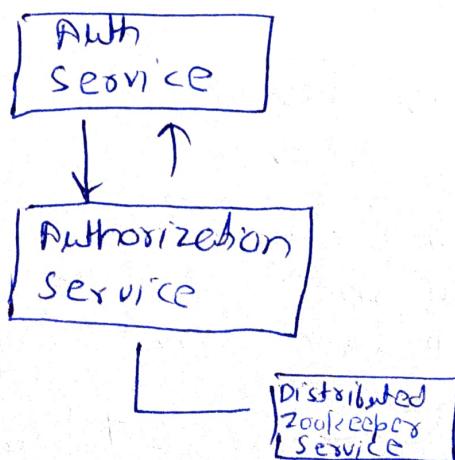
We need to store 4.5PB of data which is very large. Due to this we prefer PostgreSQL as its scalability, performance & reliability is much higher than MySQL. So PostgreSQL is preferred for Big Data Workloads.

You cannot store 4.5PB of Data in single machine. To store it efficiently, we use Sharding.

Range-based Partitioning
(Each user has Unique ID)

Advantage

- ① Linear scalability
- ② Better Performance
- ③ Fault Tolerance
- ④ Write optimization



Main DB is in Cassandra Data Type and Redis is used where caching is used and activity logs of 3-65 PB data can be stored in Cassandra.

Why Cassandra :-

- (a) Scalability :- It is horizontally scalable, can efficiently handle huge datasets.
- (b) High Write Throughput :- 5-16 PB of data can be efficiently stored in Cassandra. For storing 16PB of data store, high write speed is required that is efficiently managed by Cassandra log-structured storage engine.
- (c) No Single Point of Failure :- Cassandra uses Peer-to-peer Architecture so that system will not crash if a node fails.
- (d) Range Queries :- (Partition Key + Clustering Columns) can be used for optimizing user-specific queries.
- (e) Time-Series Data Support :- If video access logs or user activity store is required then Cassandra is optimized for time-series data also.
- In Cassandra DB, sharding can be done to increase DB query optimization.
Caching layer is put over the Cassandra database.

Id Generator Service :- Produce unique ID based on the content given to it for users, transactions etc.

To produce Unique id

Since Our application is read-heavy - Master-Distributed DB architecture is preferred.

Strategies for Distributed Database :-

- (a) Master-Slave for User-DB because of high read operations
- (b) Master-Master for high-write components like for user generated content like comments, reviews etc.

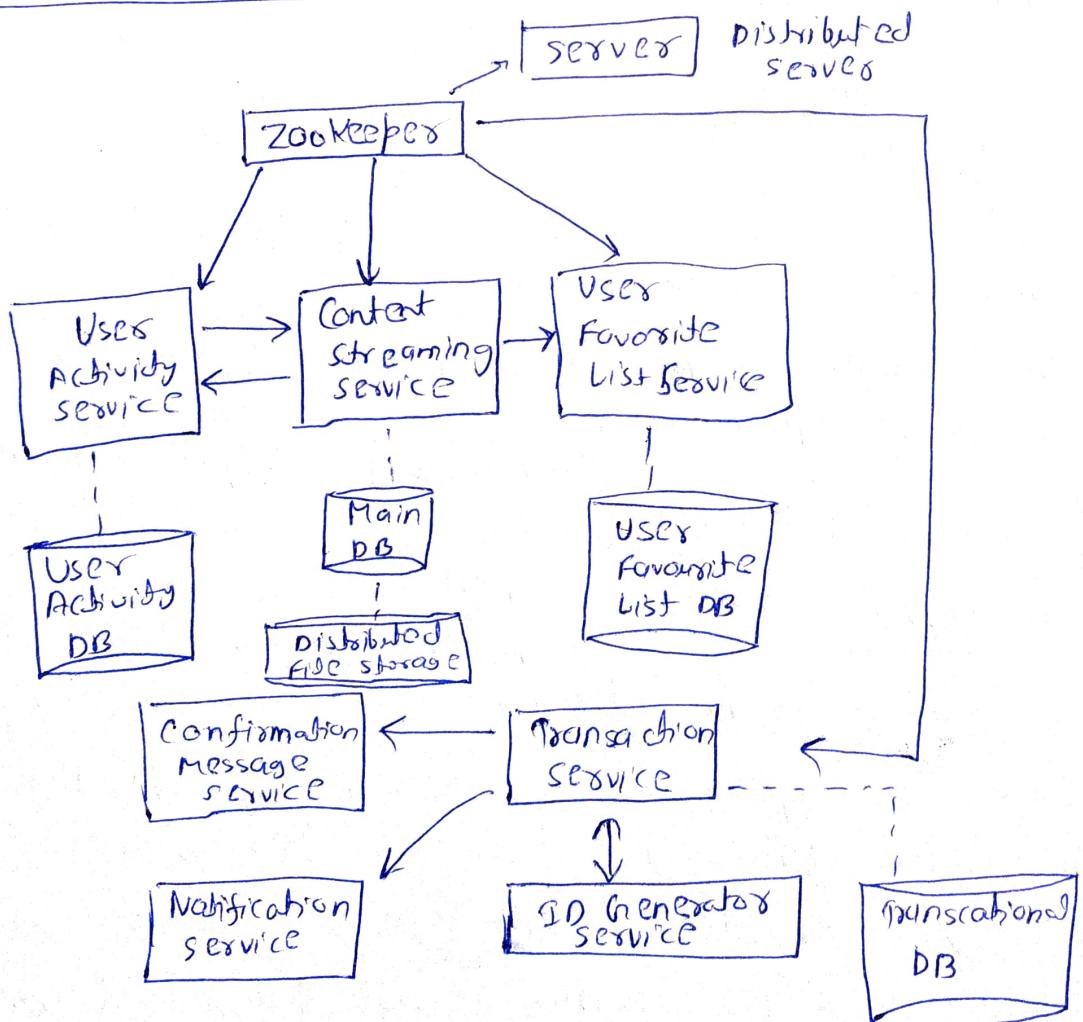
stateless or statefull System! -

stateless is preferred because:-

- (a) scalability :- stateless are easily horizontally scale.
- (b) Microservices Architecture:- Best, because services works independent of each other in it.
- (c) Load Balancing:- Can send request to another system as no session dependency will be there.
- (d) Failure Handling:- If a server failed , request go to another server without any data loss.

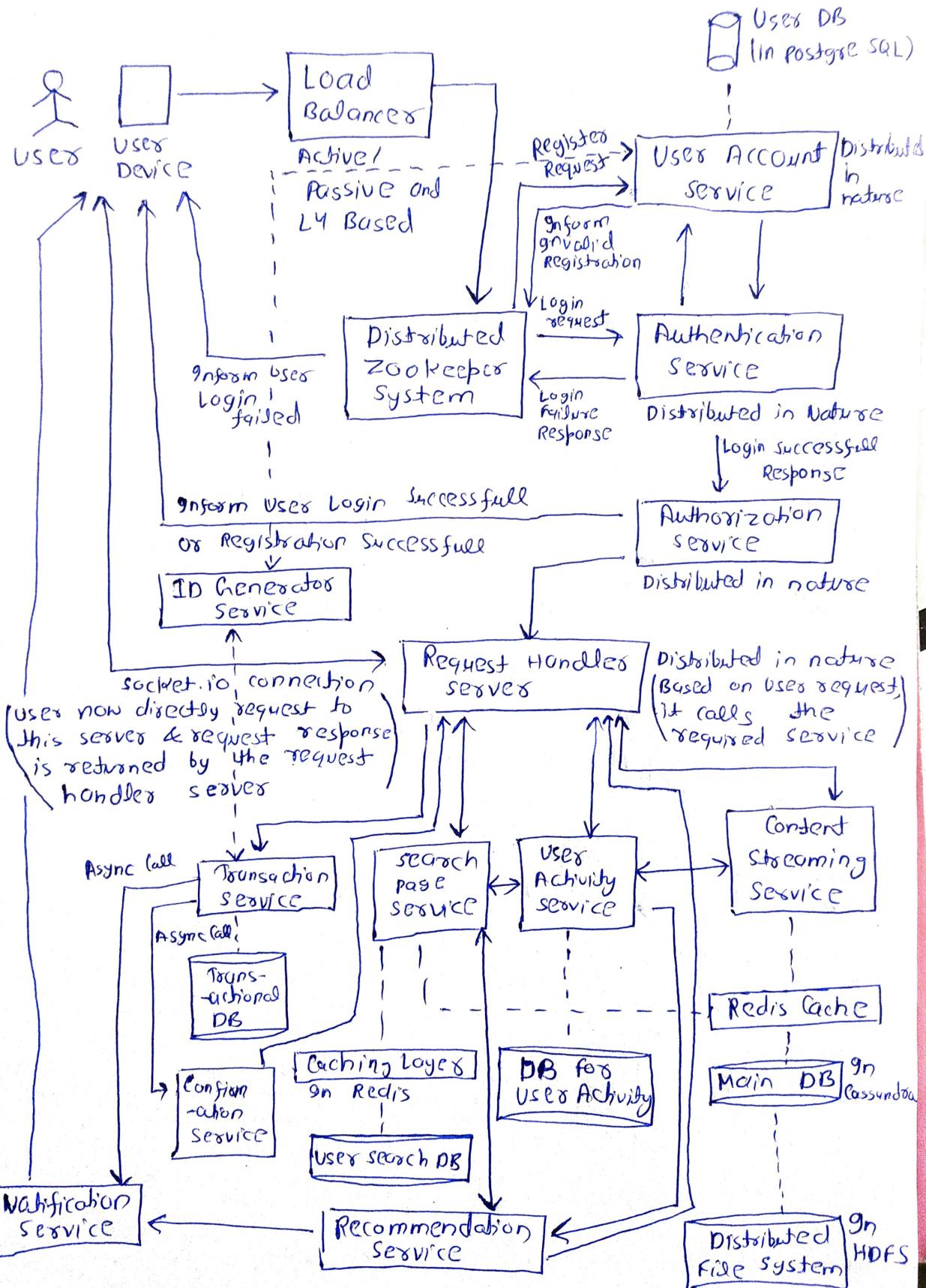
Service to Service Communication via REST (JSON over HTTP)

Introducing More Service to Complete Our All Functional And Non Functional Requirement



Search Page service and Recommendation service are also required.

Step-5 :- Final Design and Communication Protocol Identification



Full Communication Protocol :-

From	To	Protocol	Reason
User Device	Load Balancer	REST (HTTPS)	Secure Request Routing.
Load Balancer	Request Handler service	REST (HTTP)	Request Forwarding
Request Handler Service	Authentication service	REST (HTTP)	User Login Authentication
Authentication Service	User Account service	SQL Query (PostgreSQL)	User details fetch.
Authentication Service	Distributed zookeeper system	Zookeeper API (TCP)	Distributed Instance Coordination
Authentication Service	Authorization Service	REST (HTTP)	Role-based access control.
Authorization Service	Request Handler Service	REST (HTTP)	Access Approval
User Device	Request Handler Service	Web socket	Real-time updates (Login session, etc.)
Request Handler Service	Transaction service	REST (HTTPS)	Secure Payments
Transaction Service	Transactional database	SQL Query (PostgreSQL)	Payment Records
Transaction Service	Notification Service	Async Messaging (Kafka/Rabbit MQ)	Payment Confirmation Notification.
User Device	Search page service	REST (HTTPS)	Search Queries
Search page Service	Caching Layer (Redis)	Redis Binary Protocol	Fast Search Results
Search page Service	User search DB	SQL Query (PostgreSQL)	Search History
User Device	User Activity Service	Web socket	Real-time user tracking
User Activity Service	User Activity Database	SQL/Nosql Query (Cassandra)	User actions log

Full Communication Protocol :-

From	To	Protocol	Reason
User Activity Service	Recommendation Service	REST (HTTP)	Behavior data sharing
Recommendation service	DB for user Activity	SQL / NoSQL Query (Cassandra)	Analysis input
Recommendation service	Notification Service	Async Messaging (Kafka / Rabbit MQ)	Personalized Recommendations
User Device	Content Streaming Service	REST (HTTPS)	Video Requests
Content streaming service	Redis Cache	Redis Binary Protocol	Quick video loading
Content streaming service	Main Database (Cassandra)	SQL / NoSQL Query	Video Metadata
Content streaming service	Distributed File storage (HDFS)	HDFS API (HTTP/TCP)	High volume video storage
Content streaming service	User Device	UDP	High speed video streaming
User Device	Notification Service	WebSocket / Firebase Cloud Messaging (FCM)	Real-time Push Notification
Notification service	Database for Notifications	NoSQL (Cassandra)	Notification logs storage
Distributed Zookeeper System	All Distributed Services	Zookeeper API (TCP)	Cluster Coordination