

Notification Systems

Primary → High level design
Scaling for massive load
infrastructure

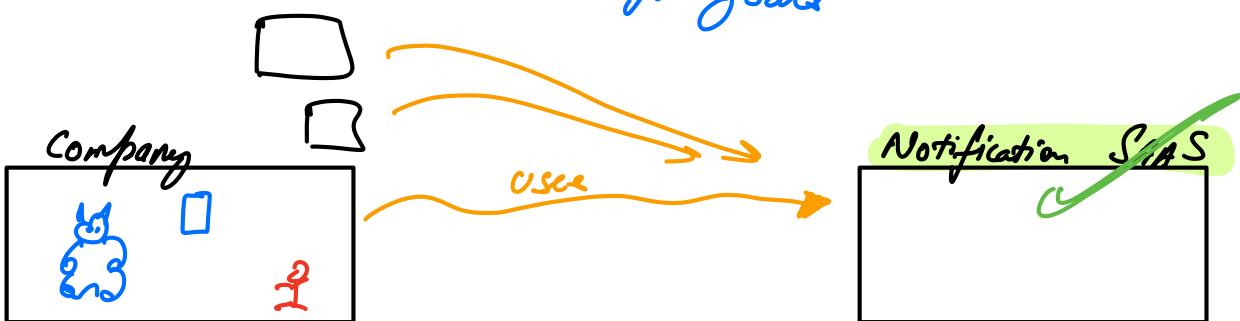
(go from 1000 → 1B users!
NoSQL/SQL, Caching, MQ &..)

V.Little → low level design (good, extensible, maintainable code)
db schema / tables / attributes
code

No → Code

Requirements

Functional Requirements
with code for / expose APIs for
Non functional Requirements
other design goals



- not a notification SaaS
- they have a separate (not related to product) notification
- they still need notifications
- enable notifications from Business to user

You

- build your own notification SaaS
- a senior architect / CTO considering to build notifications in-house

do NOT → plan on a 3rd Party!

- just do your HLD

5% of your entire company's productivity.

Functional Requirements

- ✓ client should be able to send notifications to users
company / backend code
- ✓ client should be able to send multiple types of notifications
 - SMS
 - Email
 - In-app
 - WhatsApp

pluggable / extensible
we should be able to very easily add a new notification type / channel
in the future!
- ✓ client should be able to send bulk notifications
 - examples ↗ Email campaign
 - ↗ celebrity problem (Twitter)
- X Should support multiple clients (multi-tenancy)

X^o client base limits

∴ it's a same

- how many notifications in total } • VIP clients (with more notifications)
• Pay as you go (0.001\$ / notification)

also, rate limiters

✓^o users should be able to receive notifications

✓^o take into account user preferences

- Max 10 of notifications / day

- Select channels
 - Email ✓
 - SMS X
 - In-app ✓
 - WhatsApp X

- select type

- turn off promotional
 - turn off notification about delivery
 - payment

✓^o Observability

- how many notifications are being delivered

delivery confirmation

Google → $\approx 10^7$ servers

- how many errors are raised per channel

- backups & logs

✓ channel restrictions
Compliance

client

user
channel
SMS/
Email/
...

✓ priority (Categories of notification)
OTP / recovery / transactional / promotional / ...

Non-functional

✓ low latency (≈ 10 sec)

✓ reliability

guaranteed (almost) delivery of notifications
retries / observability / success callbacks

✓ massive scale

1 hour latency is okay!

- ~~Security~~ Not discussed!

: you should not build security systems unless
you are an expert!

Elon Musk → 100M followers

100M

Youtube

any user has 100 subscribers.

2 RT users on YouTube

100M videos / day.

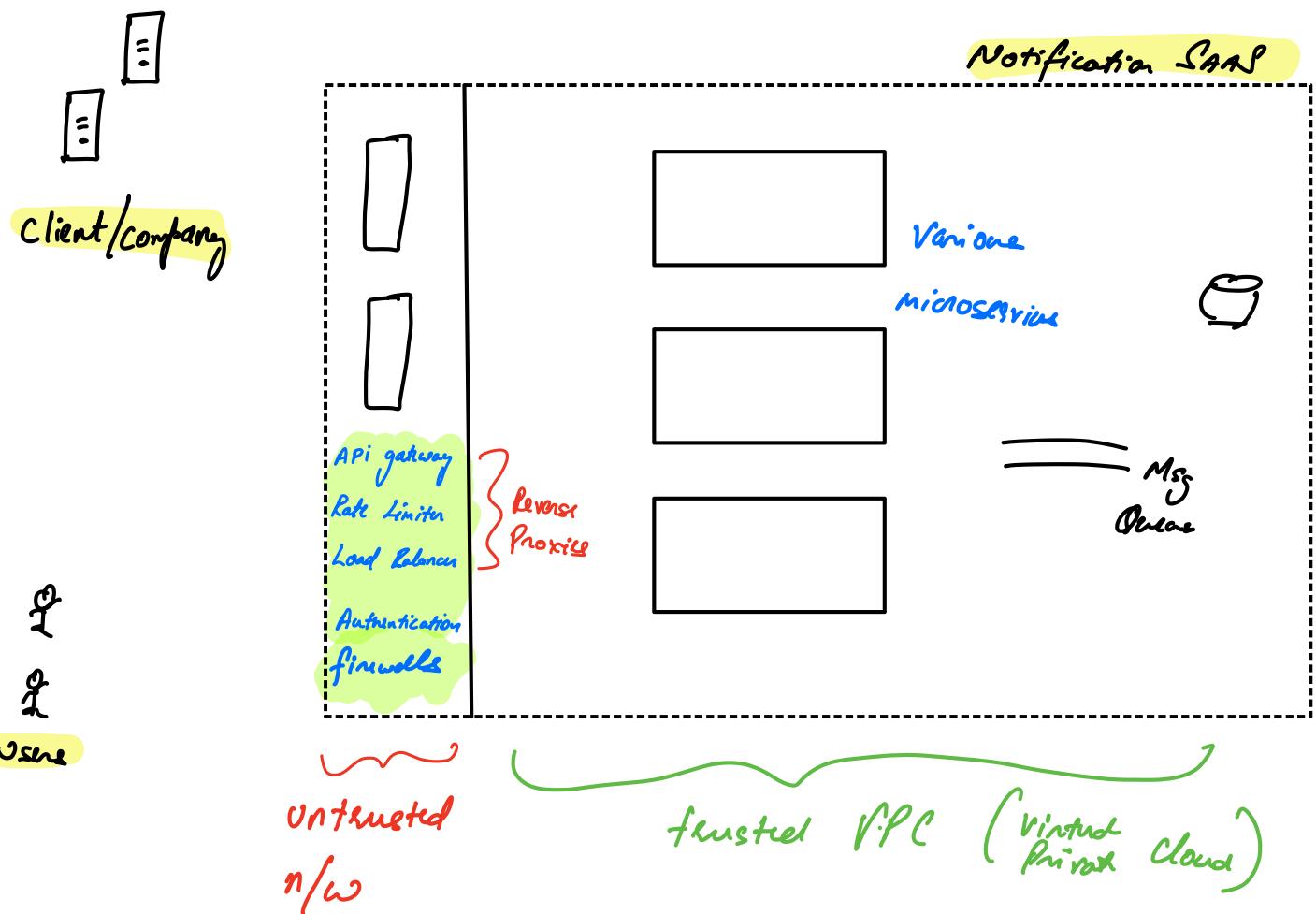
You will need to send trillions/quadrillions of notifications/day.

Now you must

subscribe + Bell icon!

Instagram / ...

Basics → typical infrastructure



HLI Basic

- Networking basics → DNS / CDN / IP / OSI
- Database basics → ACID / Normalization / Schema design / indices
- Load Balancing & consistent Hashing
- Sharding (Horizontal vs Vertical Partition)
- Replication (Master slave / multi-master)
- CAP & PACELC theorem
 - immediate consistency
 - eventual consistency
 - data loss
 - funable consistency
- Caching
 - CDNs / browser cache
 - global / local / distributed
 - in-memory (Redis)
 - Invalidation (TTL / write through / write around / write back)
 - Eviction (LRU, FIFO, LFU, ...)
- API gateway / LB / Auth / Rate Limits / reverse proxy / VPC / subnetting
- SQL vs NoSQL
 - Key-value (Redis / memcached / DynamoDB)
 - document (Mongo / Cockroach)
 - Column (Cassandra / Hbase / ... Systa)

- File (LS3 / git LFs / HDFs)
 - Internals of db → LSM Trees / bloom filters / sparse index / sharding / db orchestration / data replication
 - Message Queues (RabbitMQ / SQS / ...)
Event driven architecture
Pub / sub
 - Search → Elastic search
 - Microservices
 - distributed transactions
 - Observability
- (2 Phase Commit / Saga pattern / choreography / orchestration / correlation ID / compensating transactions)

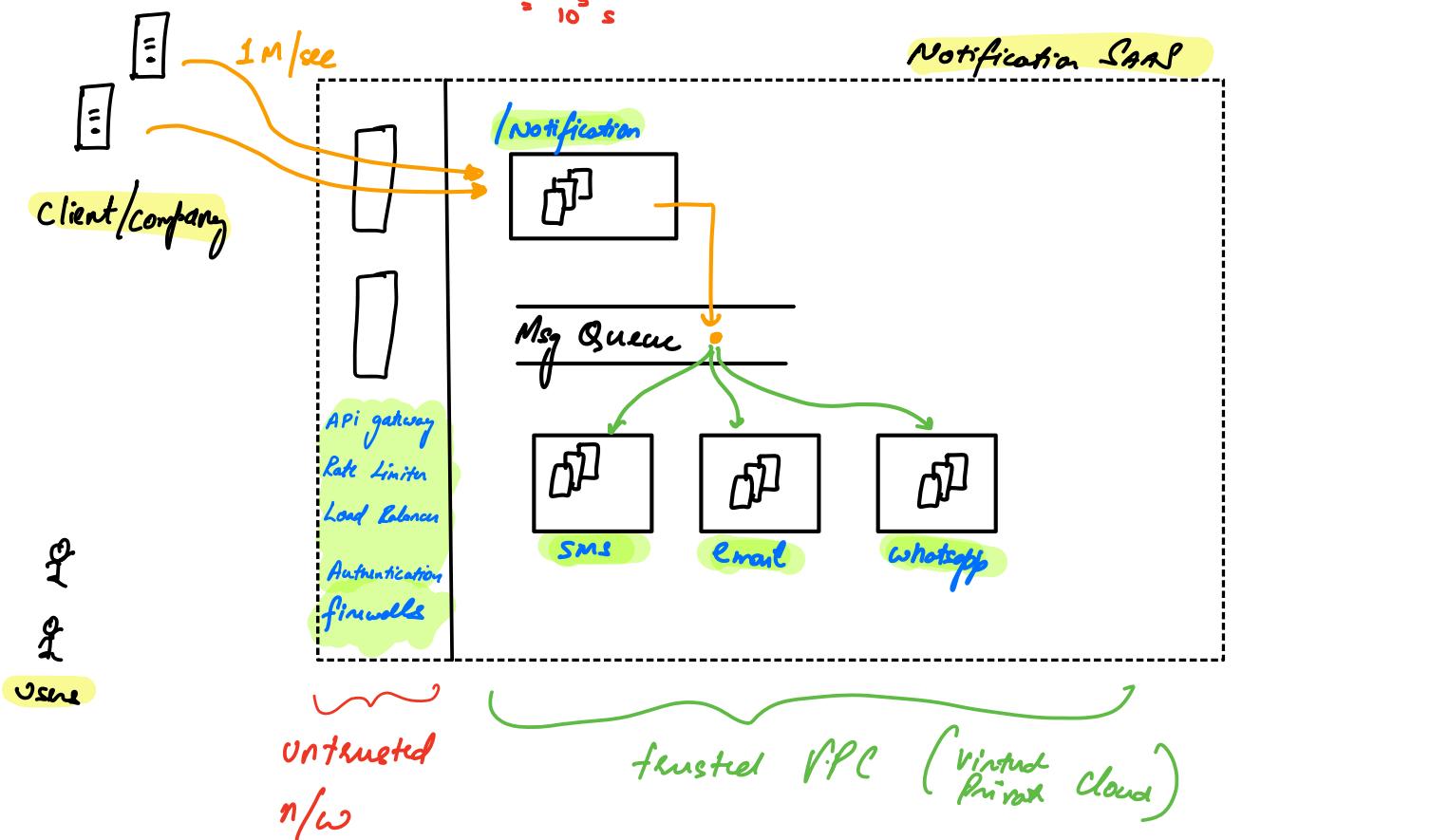
Sending Notifications + Incorporate various channels

total (normal + bulk) we have 2B users
each receives 50 notifications / day

$$\text{Rate} = 2 \text{B users} * \frac{50 \text{ notifi.}}{\text{user/day}} = \frac{2 \times 10^9 * 50}{10^5} \frac{\text{notifications}}{\text{sec}} = 10^6$$

$\frac{0.64 \times 10^9}{60 \times 60 \times 24} = 10^5 \text{ s}$

= 1M notifications/s



Client API request to send Notification

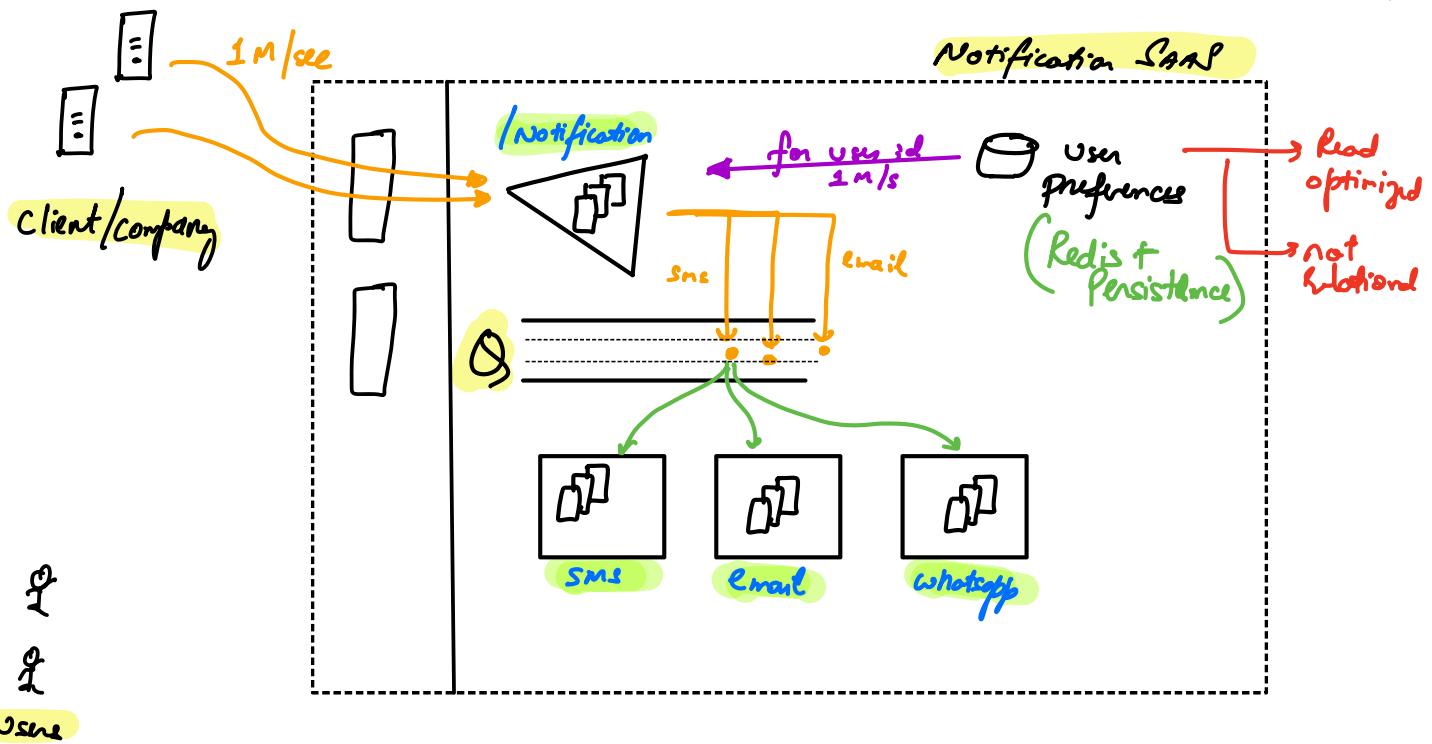
```

  user-id : 1136,
  type : 'OTP'
  channels : {
    sms : { msg: 'Hi there', attached: 'URL' }
    email : { subject: '...', body: '...', attach: '...' }
    inapp : { msg: '...', action: 'open(navigation)' }
  }
  promotional / urgent / ...
  
```

In db

maintain

User id → Phone / email / ...

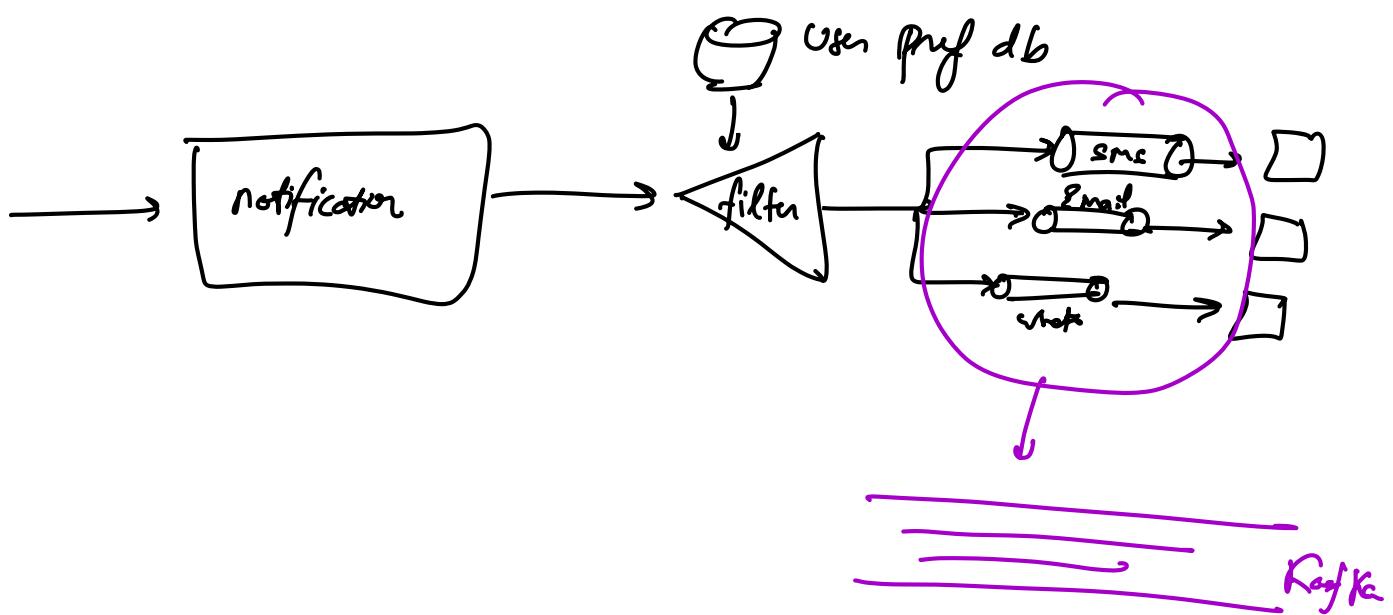


notification fan out service

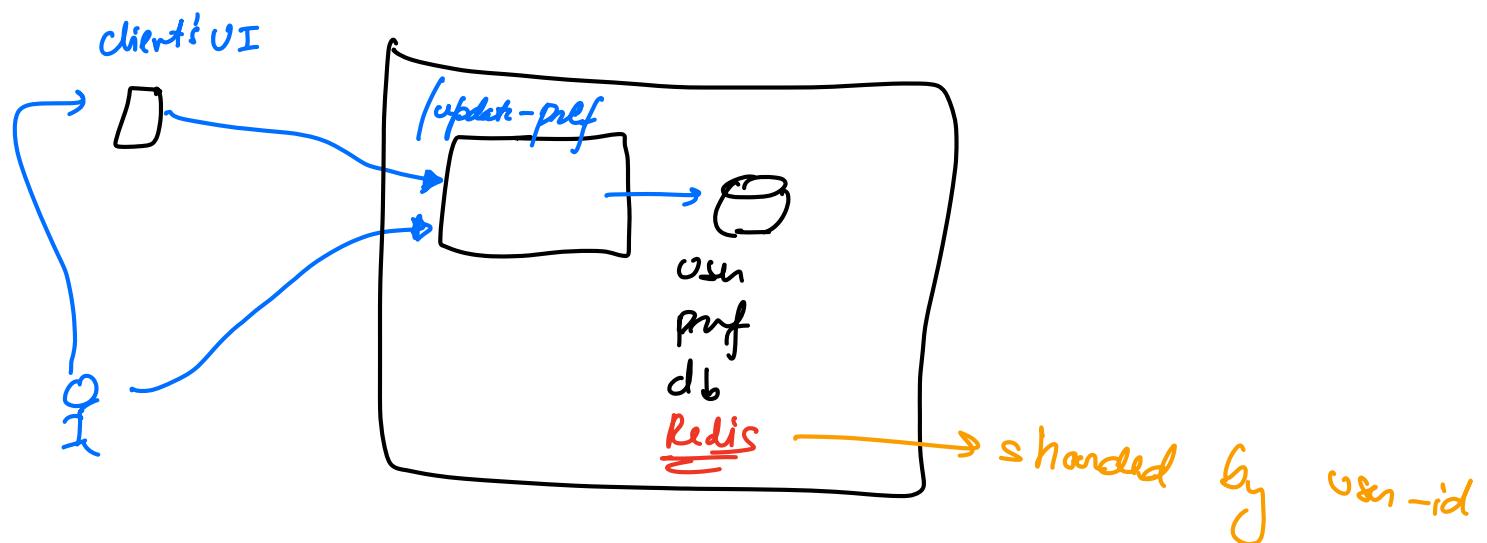
- takes multi-channel API req from client
- check user preference
- drops a msg in \mathbb{Q} for each enabled channel.

Msg \mathbb{Q} is partitioned primarily by channel / topic
sms / email / whatsapp

Persisted \mathbb{Q} → Should not lose msgs



User Preferences



2 Ω users

10 channels

User : id : channel : type : Yes / No / max per day
 10 R 4B

\rightarrow 14 B / entry

$$2R * 10 * 14 = 280 \text{ R bytes} = 280 \text{ GB of data}$$

as of 2024, max system config?

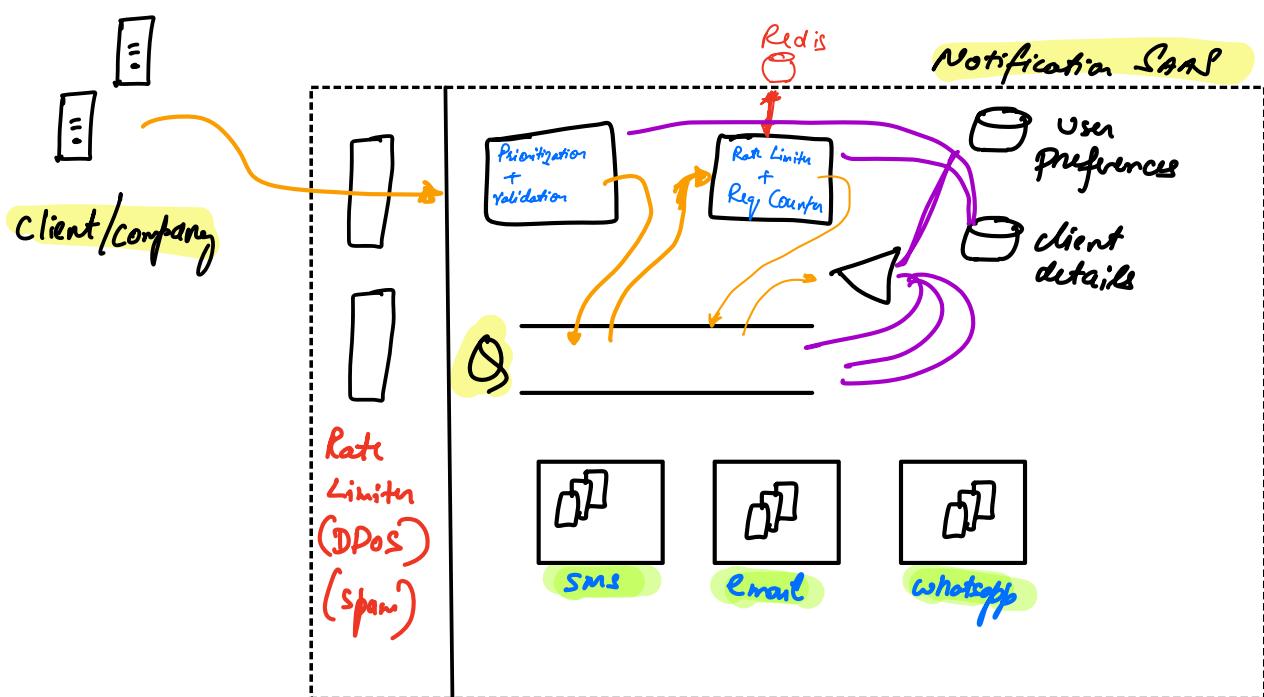
RAM 12TB

CPU 300 core / 600 threads

Disk 2PB

N/w 10Tbps

Including Limits



Notification service also applies

Rate limiting based on

- client profile (VIP/amount/...)
- user preferences.

Notification SaaS

