

# Scaling: 10K RPM vs Million RPM

## Exploring the Differences



### Introduction

Writing an app to run at 10K RPM (requests per minute) isn't too difficult. 10K RPM equals about 170 RPS (Requests per second). A basic MVP (minimum viable product) with key features and good code can achieve this easily.

You will need around 10 core for this all of this very efficiently. Obviously it depends on task to task, here I am discussing about some basic tasks like login, listing, updates and basic calculations.



That's it for now.

Stay away from over engineering. If you engineering is not aligned with your business you will be wasting a lot of resource which can actually create impact.

Everyline of code or architecture change that you do. Always keep in mind how will this impact my CPU, memory, disk and network.

Questions or Suggestions

# Who am I?

## Gaurav Yadav

I love sports and love to dive deeper in lower level implementation in computer science. I have spent a lot of time of my career in tuning systems for performance and efficiency in terms of cost and utilization

I was a Principle Engineer at OlaCabs and currently working as Principal Architect with Physicswallah



# Introduction

Writing an app to run at 10K RPM (requests per minute) isn't too difficult. 10K RPM equals about 170 RPS (Requests per second). A basic MVP (minimum viable product) with key features and good code can achieve this easily.

You will need around 10 core for this all of this very efficiently. [Obviously it depends on task to task, here I am discussing about very basic tasks like login, listing, updates and basic calculations].



# How can you scale your applications for Million RPM?

We are going to talk about different aspects that can be very useful when you strive to go for next step in your scaling journey.

Assumptions for Distributed systems:

Network is prone to disturbance. It can go down or become slow.  
Your machine/VMs can go down.

# Scaling: 10K RPM vs Million RPM

## Exploring the Differences



How to choose right database?  
There are many different databases, so before selecting one you need to understand what you want to do with it. There are many factors to consider, such as performance, scalability, and cost.



Do you really understand your application?  
It's important to have a deep understanding of your application's requirements and constraints. This will help you make informed decisions about how to scale it effectively.



Caching  
Caching is a technique used to store frequently accessed data in memory to reduce the time taken to retrieve it from a slower storage system.



Connection Pooling  
Connection pooling is a technique used to reuse database connections instead of creating new ones every time a request is made.



Stay alive or die: when your services are started taking: Keep Alive  
Keep Alive is a technique used to keep a connection alive even after it has been idle for a certain amount of time.



Polling vs Pushing  
Polling is a technique where a client repeatedly sends requests to a server to check for updates. Pushing is a technique where the server sends updates to the client as soon as they are available.



DNS Caching  
DNS caching is a technique used to store the results of previous DNS queries in memory to reduce the time taken to resolve them in the future.



CDN will save your internal systems  
A Content Delivery Network (CDN) is a network of servers located around the world that cache and deliver content to users based on their location.



Do you know your fail?  
It's important to understand what can go wrong with your application and how to prevent it.



Let's talk about about database  
Databases are the heart of your application. It's important to understand how they work and how to optimize them for performance.



08

09

Alerts, Metrics and APMs: Become necessary



10

Do you know your fail?

### Introduction

Writing an app to run at 10K RPM (requests per minute) isn't too difficult. 10K RPM equals about 170 RPS (Requests per second). A basic MVP (minimum viable product) with key features and good code can achieve this easily.

You will need around 10 core for this all of this very efficiently. Obviously it depends on task to task, here I am discussing about some basic tasks like login, listing, updates and basic calculations.



How to choose right Database?  
There are many different databases, so before selecting one you need to understand what you want to do with it. There are many factors to consider, such as performance, scalability, and cost.



Caching  
Caching is a technique used to store frequently accessed data in memory to reduce the time taken to retrieve it from a slower storage system.



Connection Pooling  
Connection pooling is a technique used to reuse database connections instead of creating new ones every time a request is made.



DNS Caching  
DNS caching is a technique used to store the results of previous DNS queries in memory to reduce the time taken to resolve them in the future.



That's it for now.

Stay away from over engineering. If you engineering is not aligned with your business you will be wasting a lot of resource which can actually create impact.

Everyline of code or architecture change that you do. Always keep in mind how will this impact my CPU, memory, disk and network.

Questions or Suggestions



# Q1



## Do you really understanding your application

Try answering below 2 questions:

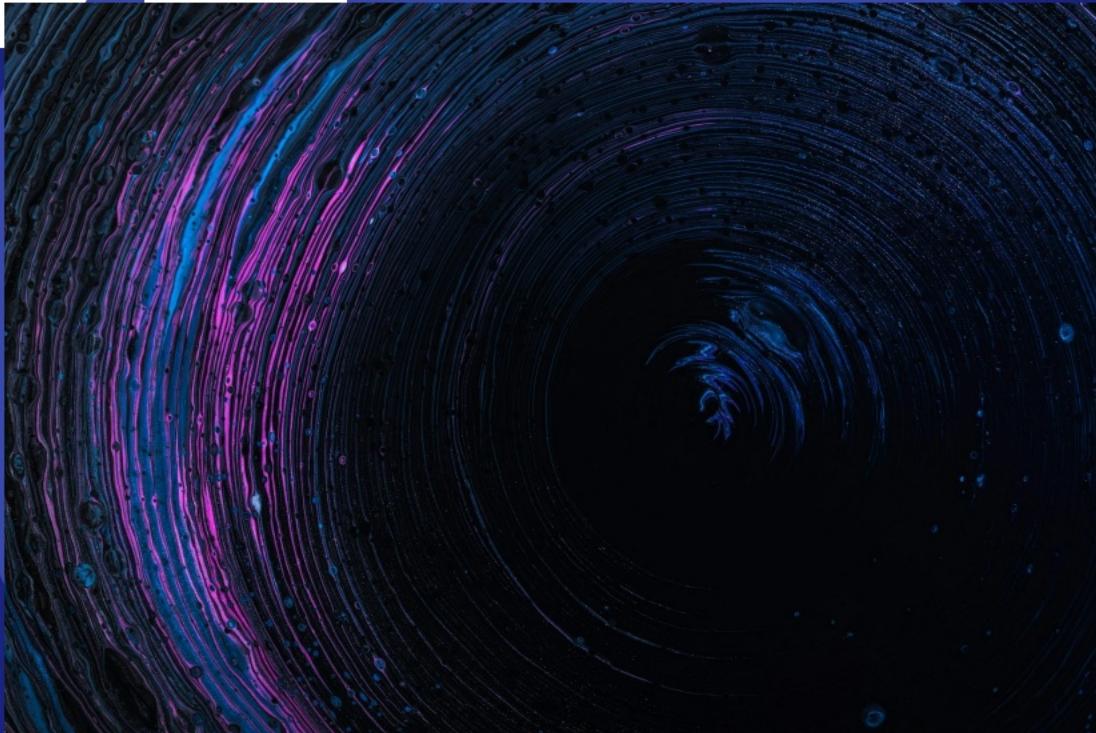
1. Is your application read heavy or write heavy
2. Is your application compute heavy or IO heavy

Answer these and write it down as it will help you in further slides.

Facts:

Most of the applications in the world are IO heavy. Lets see how?

# 02



## Caching

Caching is the first thing that comes to my mind if anyone asks how to scale my application.

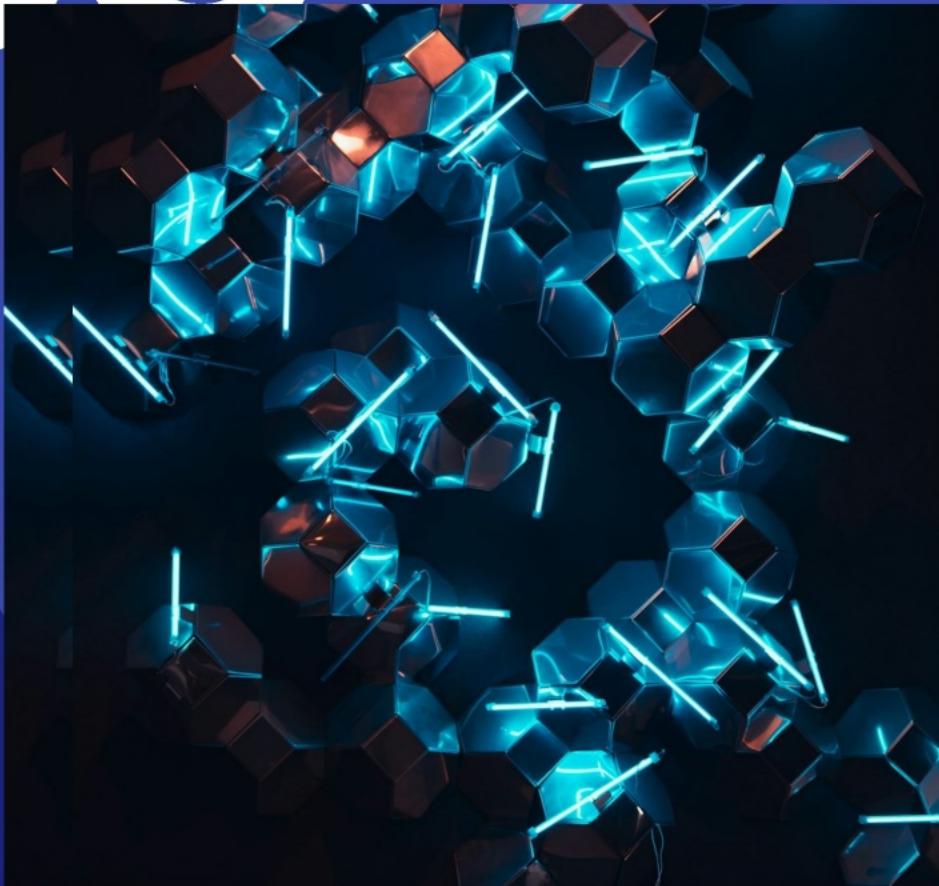
Caching saves your databases which are the most critical stateful part of your application.

The most problematic part of databases is that they are stateful.

### **Why is the stateful nature a problem?**

This is the very thing which makes horizontal scaling tricky?

# 03

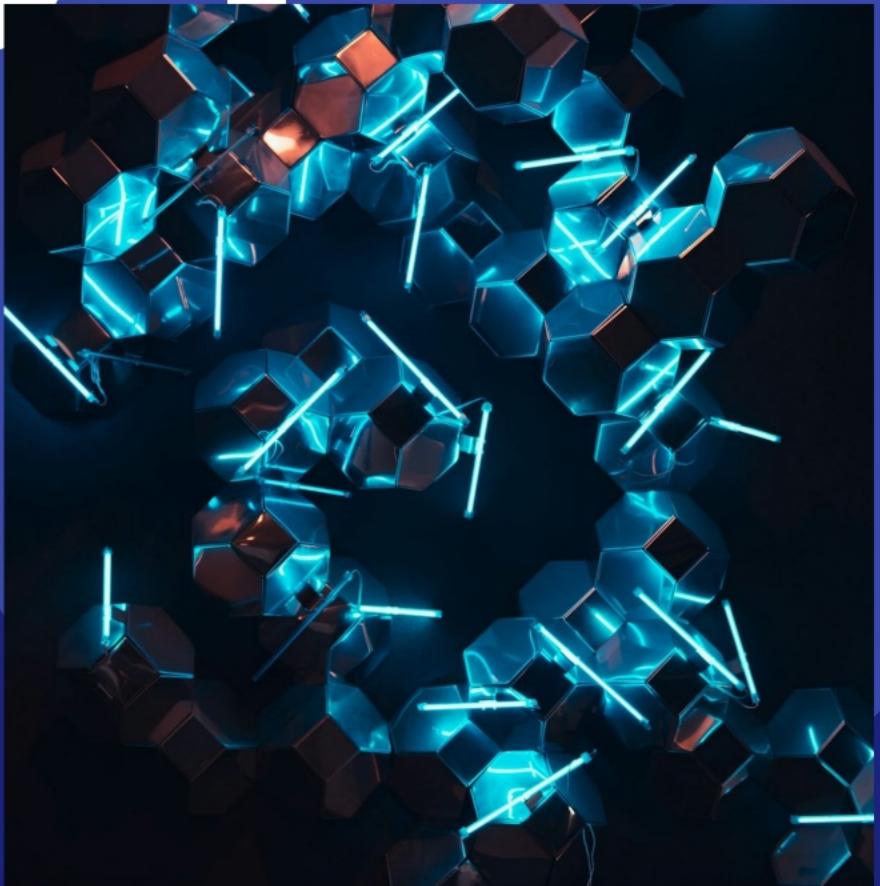


## DNS Caching

Every connection that you try to make to any dependency involves this step. Before start the three way handshake, DNS queries are fired. When working at Million RPM scale this can cause trouble at few places.

You can run DNS cache servers or add DNS caching in your applications.

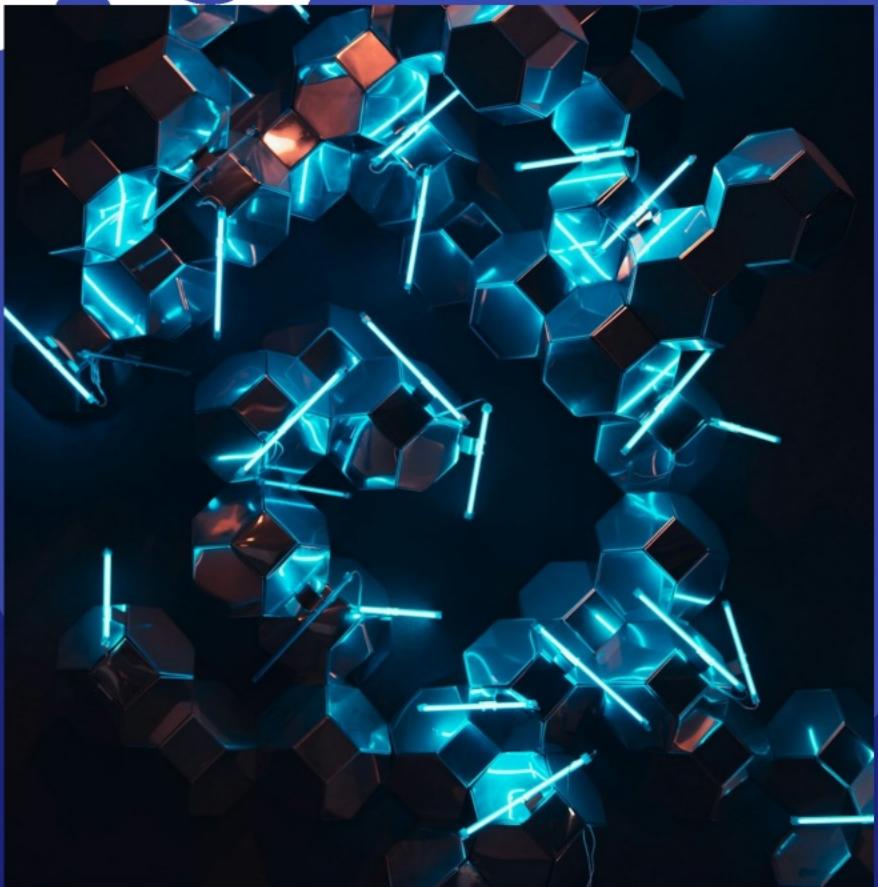
# 04



## **Stay alive no need to do handshake once you started talking: Keep Alive**

This is kind of an extension of connection pooling. When you are making request to dependencies, you can request the server to keep this connection open.

# 05



## Connection Pooling

What connection pooling means is you create a set of connection with the dependency and keep on using those instead of creating the connections for each requests.

How this helps?

Each connection use some memory and CPU, CPU usage is high when lot of connections are getting created. So creating connections at the start of the application and using them helps.

# 06



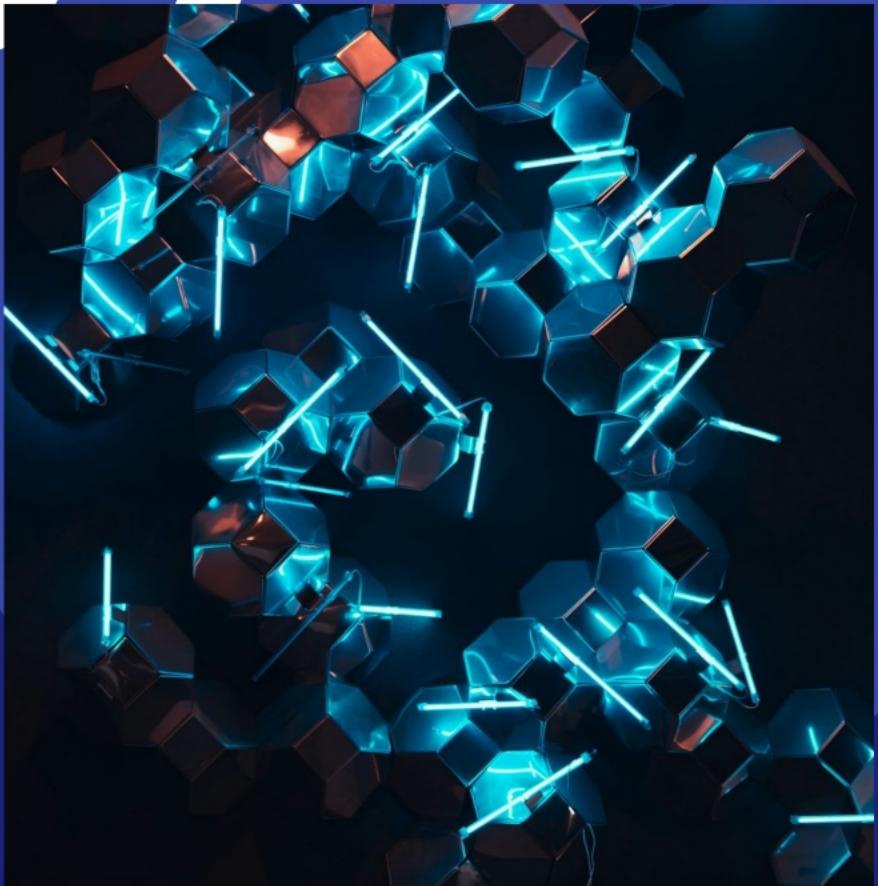
## **CDN will save your internal systems**

This can bring you major win if you are able to tackle the problem of cache invalidation.

How?

CDN serves the data from the POPs and requests never actually hit your gateways.

# 07



## Polling vs Pushing

This can be a very important decision and will depend on your application. Application which needs fresh data should focus on pushing[means socket connections] instead of polling data every second. For example: Stock market tickers

# 08



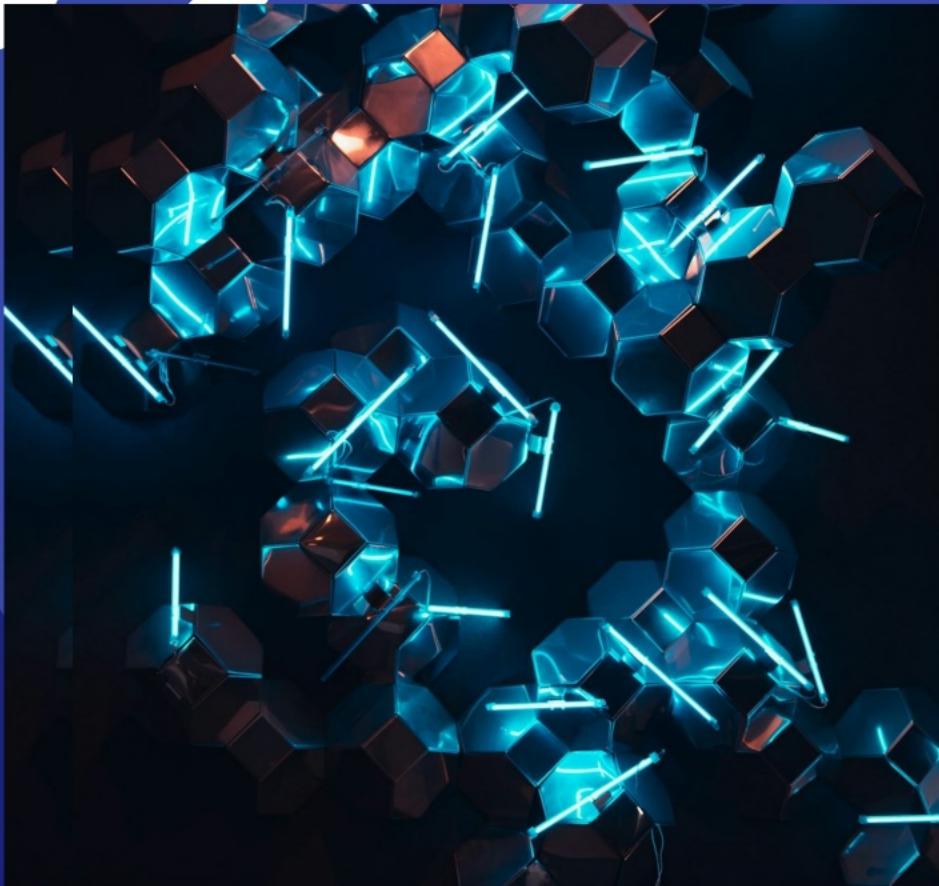
## **Rate Limiting, retries and circuit breakers will become new normal**

Rate limiting is to make sure that your apps are not overloaded during spikes and shed load in case of unwanted traffic.

Circuit breaker is to make sure that in case of any issue with services, there is a graceful way of handling the error.

Retries are important so that in any of the above cases you can request the resource back again and get the intended response. This can also lead to request amplification if not tackled properly.

# 09



## Alerts Metrics and APMs become necessity

This is nothing new but to reiterate again. Anything that you cannot measure, you will not be able to improve. As you grow in scale, there are some metrics which becomes relevant which were not earlier like p99.

Also APM will help you in finding out the hotspot in your application that need immediate attention.

You want to be informed for any CPU or Memory increase because in large system the chances of cascading is very high. For example: let's say you are running 100 instances and 20 instances go down due to load. This can put load on existing 80 instances and bring them down and can stuck in loop.

# 10



## Do you know your tail?

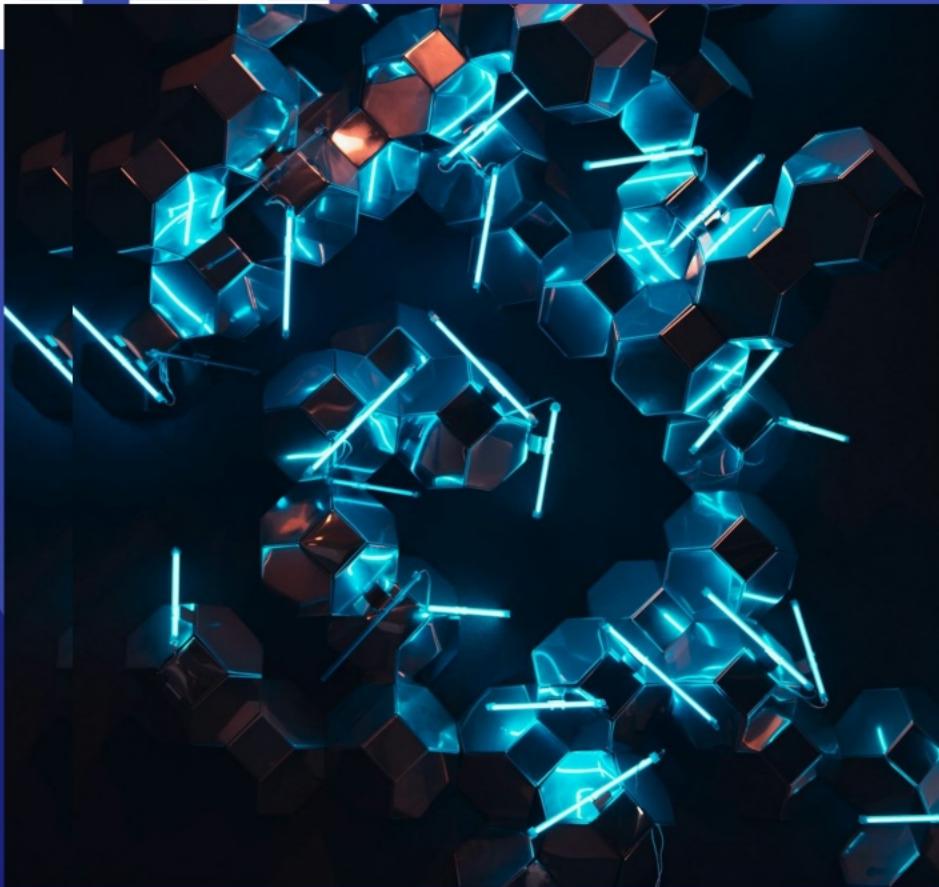
Tail here means tail latencies. Tail latencies are p99 or p99.9. You can define your own P9\* to watch.

What exactly does p99 mean?

Why is it important?

Because at scale p99 will impact a larger audience and you should always keep a check on that.

# 11



## Let's talk about database

Databases are the most critical part of your applications. When we talk about scaling there few things that you have to do at application level like using connection pooling and proper indexing and writing perfromant queries. With all this in place we will still be needing to scale databases to serve.

How to scale reads?  
Relication

How to scale writes?  
Sharding

What else to consider?  
Should you run a proxy in front of your database to handle connection and limit connections to database

# 12



## How to choose right Database?

This can be a very critical decision. So before taking this decision you should know about few things.

1. Do you know about CAP? What do you want CA, AP or CP?
2. Do you understand the nature of your data?
3. Do you understand the nature of your queries?
4. OLTP vs OLAP
5. Columnar vs Row based

# 13



## How to choose right Database?

Lets see few example. There is no direct roadmap for choosing the right database.

1. Very low Latency: You want low latency, stay away from disk. Inmemory databases.
2. CP vs AP: CA you can prefer MYSQL for AP you can opt for MongoDB or similar DB
3. For analytics you should opt for AP and columnar databases. Do you know why columnar databases for analytics?
4. Data type: Structured and unstructured, relational data.
5. Write Intensive: LSM based databases like cassandra, scylladb.
6. Do read about time series databases and vector databases. Trending right now. Why?

# Scaling: 10K RPM vs Million RPM

## Exploring the Differences



How to choose right database?  
There are many different databases, so before selecting one you need to understand what you want to do with it. There are many factors to consider, such as performance, scalability, and cost.



Do you really understand your application?  
It's important to understand what your application does and how it works. This will help you identify potential bottlenecks and optimize your code.



Caching  
Caching is a technique used to store frequently accessed data in memory to reduce the time it takes to retrieve it from storage.



Connection Pooling  
Connection pooling is a technique used to reuse database connections instead of creating a new one every time a request is made.



Stay alive or die trying! Start taking: Keep Alive  
Keep Alive is a technique used to keep a connection alive by sending small requests at regular intervals.



Polling vs Pushing  
Polling is a technique used to check for updates by periodically sending requests to a server. Pushing is a technique used to send updates to clients as soon as they are available.



DNS Caching  
DNS caching is a technique used to store the results of previous DNS queries in memory to reduce the time it takes to resolve them in the future.



CDN will save your internal systems  
A Content Delivery Network (CDN) is a network of servers located around the world that cache and deliver content to users based on their location.



Do you know your fail?  
It's important to understand what can go wrong with your application and how to fix it.



Let's talk about about database  
Databases are the heart of your application. It's important to understand how they work and how to optimize them.



Data Logging, retries and circuit breakers will handle your errors  
Data logging, retries, and circuit breakers are techniques used to handle errors and ensure your application remains available.



Alerts, Metrics and APMs: Become necessary  
Alerts, metrics, and Application Performance Monitoring (APM) are tools used to monitor your application and detect issues early.

### Introduction

Writing an app to run at 10K RPM (requests per minute) isn't too difficult. 10K RPM equals about 170 RPS (Requests per second). A basic MVP (minimum viable product) with key features and good code can achieve this easily.

You will need around 10 core for this all of this very efficiently. Obviously it depends on task to task, here I am discussing about some basic tasks like login, listing, updates and basic calculations.



That's it for now.

Stay away from over engineering. If you engineering is not aligned with your business you will be wasting a lot of resource which can actually create impact.

Everyline of code or architecture change that you do. Always keep in mind how will this impact my CPU, memory, disk and network.

Questions or Suggestions

# **Thats it for now.**

Stay away from over engineering, if your engineering is not aligned with your business you will be wasting a lot of resource which can actually create impact.

Everyline of code or architecture change that you do. Always keep in mind how will this imapct my CPU, memory, disk and network.

Questions or Suggestions

# Scaling: 10K RPM vs Million RPM

## **Exploring the Differences**

