

[SYSTEMS DESIGN]

1.Design a global and fast code-deployment system.

Summary :

- Many systems design questions are intentionally left very vague and are literally given in the form of `Design FooBar`. It's your job to ask clarifying questions to better understand the system that you have to build.
- We've laid out some of these questions below; their answers should give you some guidance on the problem. Before looking at them, we encourage you to take few minutes to think about what questions you'd ask in a real interview.

Clarifying questions to ask :

Q: What exactly do we mean by a code-deployment system? Are we talking about building, testing, and shipping code?

A: We want to design a system that takes code, builds it into a binary (an opaque blob of data—the compiled code), and deploys the result globally in an efficient and scalable way. We don't need to worry about testing code; let's assume that's already covered.

Q: What part of the software-development lifecycle, so to speak, are we designing this for? Is this process of building and deploying code happening when code is being submitted for code review, when code is being merged into a codebase, or when code is being shipped?

A: Once code is merged into the trunk or master branch of a central code repository, engineers should be able to trigger a build and deploy that build (through a UI, which we're not designing). At that point, the code has already been reviewed and is ready to ship. So, to clarify, we're not designing the system that handles code being submitted for review or being merged into a master branch—just the system that takes merged code, builds it, and deploys it.

Q: Are we essentially trying to ship code to production by sending it to, presumably, all of our application servers around the world?

A: Yes, exactly.

Q: How many machines are we deploying to? Are they located all over the world?

A: We want this system to scale massively to hundreds of thousands of machines spread across 5-10 regions throughout the world.

Q: This sounds like an internal system. Is there any sense of urgency in deploying this code? Can we afford failures in the deployment process? How fast do we want a single deployment to take?

A: This is an internal system, but we'll want to have decent availability, because many outages are resolved by rolling forward or rolling back buggy code, so this part of the infrastructure may be necessary to avoid certain terrible situations. In terms of failure tolerance, any build should eventually reach a SUCCESS or FAILURE state. Once a binary has been successfully built, it should be shippable to all machines globally within 30 minutes.

Q: So, it sounds like we want our system to be available, but not necessarily highly available, we want a clear end-state for builds, and we want the entire process of building and deploying code to take roughly 30 minutes. Is that correct?

A: Yes, that's correct.

Q: How often will we be building and deploying code, how long does it take to build code, and how big can the binaries that we'll be deploying get?

A: Engineering teams deploy hundreds of services or web applications, thousands of times per day; building code can take up to 15 minutes; and the final binaries can reach sizes of up to 10 GB. The fact that we might be dealing with hundreds of different applications shouldn't matter though; you're just designing the build pipeline and deployment system, which are agnostic to the types of applications that are getting deployed.

Q: When building code, how do we have access to the actual code? Is there some sort of reference that we can use to grab code to build?

A: Yes; you can assume that you'll be building code from commits that have been merged into a master branch. These commits have SHA identifiers (effectively arbitrary strings) that you can use to download the code that needs to be built.

SOLUTION WALKTHROUGH :

Our solution walkthroughs are meant to supplement our video solutions. We recommend starting with the video solution and using the walkthrough either as a point of reference while you watch the video or as a review tool if you need to brush up on this question's solution later on.

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

From the answers we were given to our clarifying questions (see Prompt Box), we're building a system that involves repeatedly (in the order of thousands of times per day) building and deploying code to hundreds of thousands of machines spread out across 5-10 regions around the world.

Building code will involve grabbing snapshots of source code using commit SHA identifiers; beyond that, we can assume that the actual implementation details of the building action are taken care of. In other words, we don't need to worry about how we would build JavaScript code or C++ code; we just need to design the system that enables the repeated building of code.

Building code will take up to 15 minutes, it'll result in a binary file of up to 10GB, and we want to have the entire deployment process (building and deploying code to our target machines) take at most 30 minutes.

Each build will need a clear end-state (SUCCESS or FAILURE), and though we care about availability (2 to 3 nines), we don't need to optimize too much on this dimension.

2. Coming Up with A Plan

It's important to organize ourselves and to lay out a clear plan regarding how we're going to tackle our design. What are the major, distinguishable components of our how system?

It seems like this system can actually very simply be divided into two clear subsystems:
the Build System that builds code into binaries

the Deployment System that deploys binaries to our machines across the world

Note that these subsystems will of course have many components themselves, but this is a very straightforward initial way to approach our problem.

3. Build System -- General Overview

From a high-level perspective, we can call the process of building code into a binary a job, and we can design our build system as a queue of jobs. Jobs get added to the queue, and each job has a commit identifier (the commit SHA) for what version of the code it should build and the name of the artifact that will be created (the name of the resulting binary). Since we're agnostic to the type of the code being built, we can assume that all languages are handled automatically here.

We can have a pool of servers (workers) that are going to handle all of these jobs. Each worker will repeatedly take jobs off the queue (in a FIFO manner—no prioritization for now), build the relevant binaries (again, we're assuming that the actual implementation details of building code are given to us), and write the resulting binaries to blob storage (Google Cloud Storage or S3 for instance). Blob storage makes sense here, because binaries are literally blobs of data.

4. Build System -- Job Queue

A naive design of the job queue would have us implement it in memory (just as we would implement a queue in coding interviews), but this implementation is very problematic; if there's a failure in our servers that hold this queue, we lose the entire state of our jobs: queued jobs and past jobs.

It seems like we would be unnecessarily complicating matters by trying to optimize around this in-memory type of storage, so we're likely better off implementing the queue using a SQL database.

5. Build System -- SQL Job Queue

We can have a jobs table in our SQL database where every record in the database represents a job, and we can use record-creation timestamps as the queue's ordering mechanism.

Our table will be:

id: string, the ID of the job, auto generated

created_at: timestamp

commit_sha: string

name: string, the pointer to the job's eventual binary in blob storage

status: string, QUEUED, RUNNING, SUCCEEDED, FAILED

We can implement the actual dequeuing mechanism by looking at the oldest creation_timestamp with a QUEUED status. This means that we'll likely want to index our table on both created_at and status.

6. Build System -- Concurrency

ACID transactions will make it safe for potentially hundreds of workers to grab jobs off the queue without unintentionally running the same job twice (we'll avoid race conditions). Our actual transaction will look like this:

```
BEGIN TRANSACTION;
```

```
SELECT * FROM jobs_table WHERE status = 'QUEUED' ORDER BY created_at  
ASC LIMIT 1;
```

```
// if there's none, we ROLLBACK;
```

```
UPDATE jobs_table SET status = 'RUNNING' WHERE id = id from previous query;
```

```
COMMIT;
```

All of the workers will be running this transaction every so often to dequeue the next job; let's say every 5 seconds. If we arbitrarily assume that we'll have 100 workers sharing the same queue, we'll have $100/5 = 20$ reads per second, which is very easy to handle for a SQL database.

7. Build System -- Lost Jobs

Since we're designing a large-scale system, we have to expect and handle edge cases. Here, what if there's a network partition with our workers or one of our workers dies mid-build? Since builds last around 15 minutes on average, this will very likely happen. In this case, we want to avoid having a "lost job" that we were never made

aware of, and with our current design, the job will remain RUNNING forever. How do we handle this?

We could have an extra column on our jobs table called last_heartbeat. This will be updated in a heartbeat fashion by the worker running a particular job, where that worker will update the relevant row in the table every 3-5 minutes to just let us know that it's still running the job.

We can then have a completely separate service that polls the table every so often (say, every 5 minutes, depending on how responsive we want this build system to be), checks all of the RUNNING jobs, and if their last_heartbeat was last modified longer than 2 heartbeats ago (we need some margin of error here), then something's likely wrong, and this service can reset the status of the relevant jobs to QUEUED, which would effectively bring them back to the front of the queue.

The transaction that this auxiliary service will perform will look something like this:

```
UPDATE jobs_table SET status = 'QUEUED' WHERE  
    status = 'RUNNING' AND  
    last_heartbeat < NOW() - 10 minutes.
```

8. Build System -- Scale Estimation

We previously arbitrarily assumed that we would have 100 workers, which made our SQL-database queue able to handle the expected load. We should try to estimate if this number of workers is actually realistic.

With some back-of-the-envelope math, we can see that, since a build can take up to 15 minutes, a single worker can run 4 jobs per hour, or ~100 (96) jobs per day. Given thousands of builds per day (say, 5000-10000), this means that we would need 50-100 workers (5000 / 100). So our arbitrary figure was accurate.

Even if the builds aren't uniformly spread out (in other words, they peak during work hours), our system scales horizontally very easily. We can automatically add or remove workers whenever the load warrants it. We can also scale our system vertically by making our workers more powerful, thereby reducing the build time.

9. Build System -- Storage

We previously mentioned that we would store binaries in blob storage (GCS). Where does this storage fit into our queueing system exactly?

When a worker completes a build, it can store the binary in GCS before updating the relevant row in the jobs table. This will ensure that a binary has been persisted before its relevant job is marked as SUCCEEDED.

Since we're going to be deploying our binaries to machines spread across the world, it'll likely make sense to have regional storage rather than just a single global blob store.

We can design our system based on regional clusters around the world (in our 5-10 global regions). Each region can have a blob store (a regional GCS bucket). Once a worker successfully stores a binary in our main blob store, the worker is released and can run another job, while the main blob store performs some asynchronous replication to store the binary in all of the regional GCS buckets. Given 5-10 regions and 10GB files, this step should take no more than 5-10 minutes, bringing our total build-and-deploy duration so far to roughly 20-25 minutes (15 minutes for a build and 5-10 minutes for global replication of the binary).

10. Deployment System -- General Overview

From a high-level perspective, our actual deployment system will need to allow for the very fast distribution of 10GB binaries to hundreds of thousands of machines across all of our global regions. We're likely going to want some service that tells us when a binary has been replicated in all regions, another service that can serve as the source of truth for what binary should currently be run on all machines, and finally a peer-to-peer-network design for our actual machines across the world.

11. Deployment System -- Replication-Status Service

We can have a global service that continuously checks all regional GCS buckets and aggregates the replication status for successful builds (in other words, checks that a given binary in the main blob store has been replicated across all regions). Once a binary has been replicated across all regions, this service updates a separate SQL database with rows containing the name of a binary and a `replication_status`. Once a binary has a "complete" `replication_status`, it's officially deployable.

12. Deployment System -- Blob Distribution

Since we're going to deploy 10 GBs to hundreds of thousands of machines, even with our regional clusters, having each machine download a 10GB file one after the other from a regional blob store is going to be extremely slow. A peer-to-peer-network approach will be much faster and will allow us to hit our 30-minute time frame for deployments. All of our regional clusters will behave as peer-to-peer networks.

13. Deployment System -- Trigger

Let's describe what happens when an engineer presses a button on some internal UI that says, "Deploy build/binary B1 to every machine globally". This is the action that triggers the binary downloads on all the regional peer-to-peer networks.

To simplify this process and to support having multiple builds getting deployed concurrently, we can design this in a goal-state oriented manner.

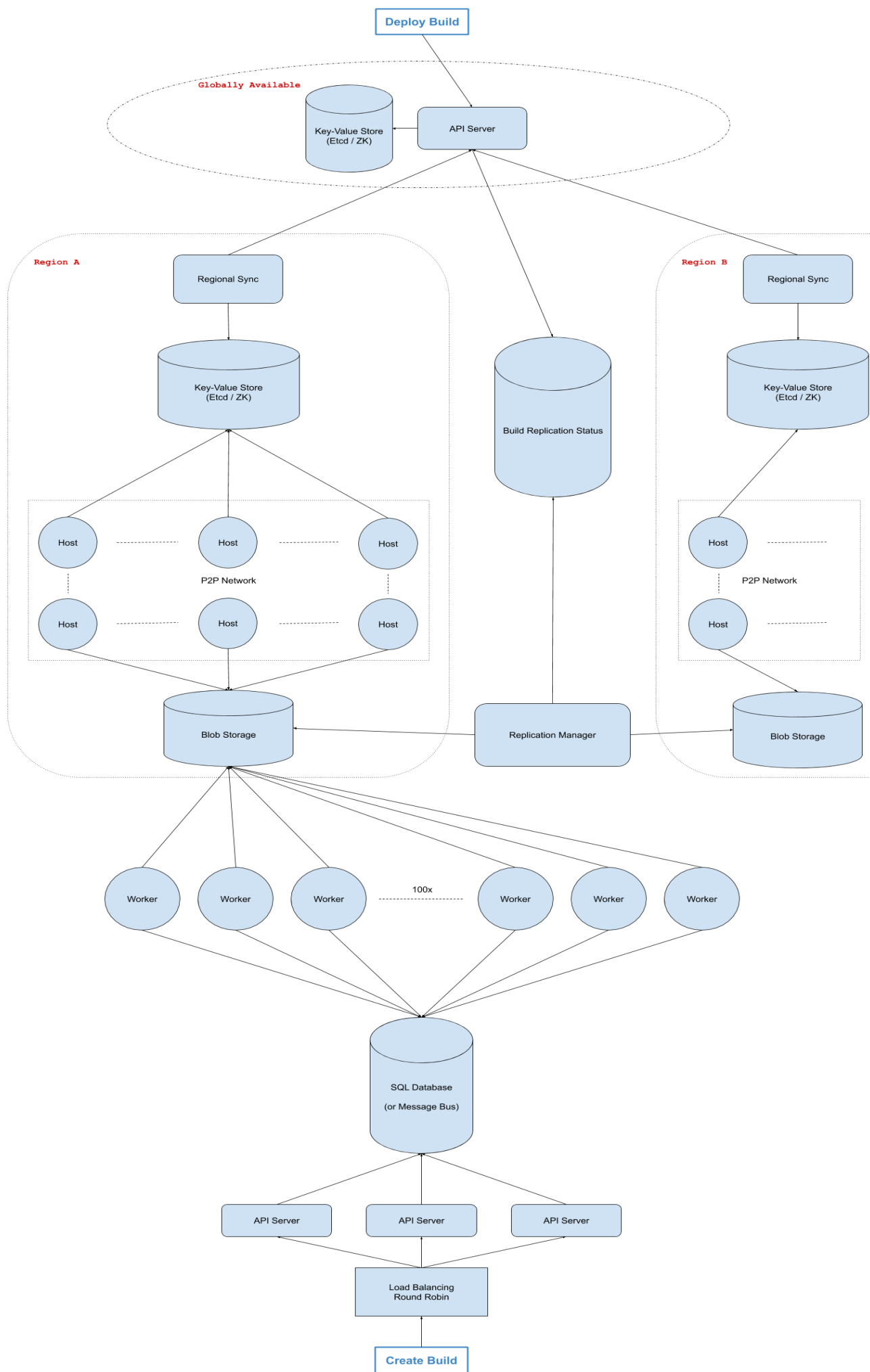
The goal-state will be the desired build version at any point in time and will look something like: "current_build: B1", and this can be stored in some dynamic configuration service (a key-value store like Etcd or ZooKeeper). We'll have a global goal-state as well as regional goal-states.

Each regional cluster will have a K-V store that holds configuration for that cluster about what builds should be running on that cluster, and we'll also have a global K-V store.

When an engineer clicks the "Deploy build/binary B1" button, our global K-V store's build_version will get updated. Regional K-V stores will be continuously polling the global K-V store (say, every 10 seconds) for updates to the build_version and will update themselves accordingly.

Machines in the clusters/regions will be polling the relevant regional K-V store, and when the build_version changes, they'll try to fetch that build from the P2P network and run the binary.

System diagram is given on next page.



2. Design The Uber API.

Summary :

- Many systems design questions are intentionally left very vague and are literally given in the form of `Design FooBar`. It's your job to ask clarifying questions to better understand the system that you have to build.
- We've laid out some of these questions below; their answers should give you some guidance on the problem. Before looking at them, we encourage you to take few minutes to think about what questions you'd ask in a real interview.

Clarifying Questions to Ask :

Q: Uber has a lot of different services: there's the core ride-hailing Uber service, there's UberEats, there's UberPool--are we designing the API for all of these services, or just for one of them?

A: Let's just design the core rides API -- not UberEats or UberPool.

Q: At first thought, it seems like we're going to need both a passenger-facing API and a driver-facing API--does that make sense, and if yes, should we design both?

A: Yes, that totally makes sense. And yes, let's design both, starting with the passenger-facing API.

Q: To make sure we're on the same page, this is the functionality that I'm envisioning this API will support: A user (a passenger) goes on their phone and hails a ride; they get matched with a driver; then they can track their ride as it's in progress, until they reach their destination, at which point the ride is complete. Throughout this process, there are a few more features to support, like being able to track where the passenger's driver is before the passenger gets picked up, maybe being able to cancel rides, etc.. Does this capture most of what you had in mind?

A: Yes, this is precisely what I had in mind. And you can work out the details as you start designing the API.

Q: Do we need to handle things like creating an Uber account, setting up payment preferences, contacting Uber, etc..? What about things like rating a driver, tipping a driver, etc.?

A: For now, let's skip those and really focus on the core taxiing service.

Q: Just to confirm, you want me to write out function signatures for various API endpoints, including parameters, their types, return values, etc., right?

A: Yup, exactly.

SOLUTION WALKTHROUGH :

1. Gathering Requirements

As with any API design interview question, the first thing that we want to do is to gather API requirements; we need to figure out what API we're building exactly.

We're designing the core ride-hailing service that Uber offers. Passengers can book a ride from their phone, at which point they're matched with a driver; they can track their driver's location throughout the ride, up until the ride is finished or canceled; and they can also see the price of the ride as well as the estimated time to destination throughout the trip, amongst other things.

The core taxiing service that Uber offers has a passenger-facing side and a driver-facing side; we're going to be designing the API for both sides.

2. Coming Up with A Plan

It's important to organize ourselves and to lay out a clear plan regarding how we're going to tackle our design. What are the major, potentially contentious parts of our API? Why are we making certain design decisions?

We're going to center our API around a Ride entity; every Uber ride will have an associated Ride containing information about the ride, including information about its passenger and driver.

Since a normal Uber ride can only have one passenger (one passenger account--the one that hails the ride) and one driver, we're going to cleverly handle all permissioning related to ride operations through passenger and driver IDs. In other words, operations like GetRide and EditRide will purely rely on a passed userId, the userId of the

passenger or driver calling them, to return the appropriate ride tied to that passenger or driver.

We'll start by defining the Ride entity before designing the passenger-facing API and then the driver-facing API.

3. Entities

Ride :

The Ride entity will have a unique id, info about its passenger and its driver, a status, and other details about the ride.

- rideId: string
- passengerInfo: PassengerInfo
- driverInfo?: DriverInfo
- rideStatus: RideStatus - enum CREATED/MATCHED/STARTED/FINISHED/CANCELED
- start: GeoLocation
- destination: GeoLocation
- createdAt: timestamp
- startTime: timestamp
- estimatedPrice: int
- timeToDestination: int

We'll explain why the driverInfo is optional when we get to the API endpoints.

PassengerInfo :

id: string

name: string

rating: int

DriverInfo :

id: string

name: string

rating: int

ridesCount: int

vehicleInfo: VehicleInfo

VehicleInfo :

licensePlate: string

description: string

4. Passenger API

The passenger-facing API will be fairly straightforward. It'll consist of simple CRUD operations around the Ride entity, as well as an endpoint to stream a driver's location throughout a ride.

CreateRide(userId: string, pickup: Geolocation, destination: Geolocation)

=> Ride

Usage: called when a passenger books a ride; a Ride is created with no DriverInfo and with a CREATED RideStatus; the Uber backend calls an internal FindDriver API that uses an algorithm to find the most appropriate driver; once a driver is found and accepts the ride, the backend calls EditRide with the driver's info and with a MATCHED RideStatus.

GetRide(userId: string)

=> Ride

Usage: polled every couple of seconds after a ride has been created and until the ride has a status of MATCHED; afterwards, polled every 20-90 seconds throughout the trip to update the ride's estimated price, it's time to destination, its RideStatus if it's been canceled by the driver, etc..

EditRide(userId: string, [...params?: all properties on the Ride object that need to be edited])

=> Ride

CancelRide(userId: string)

=> void

Wrapper around EditRide -- effectively calls EditRide(userId: string, rideStatus: CANCELLED).

StreamDriverLocation(userId: string)

=> Geolocation

Usage: continuously streams the location of a driver over a long-lived websocket connection; the driver whose location is streamed is the one associated with the Ride tied to the passed userId.

5. Driver API

The driver-facing API will rely on some of the same CRUD operations around the Ride entity, and it'll also have a SetDriverStatus endpoint as well as an endpoint to push the driver's location to passengers who are streaming it.

SetDriverStatus(userId: string, driverStatus: DriverStatus)

=> void

DriverStatus: enum UNAVAILABLE/IN_RIDE/STANDBY

Usage: called when a driver wants to look for a ride, is starting a ride, or is done for the day; when called with STANDBY, the Uber backend calls an internal FindRide API that uses an algorithm to enqueue the driver in a queue of drivers waiting for rides and to find the most appropriate ride; once a ride is found, the ride is internally locked to the driver for 30 seconds, during which the driver can accept or reject the ride; once the driver accepts the ride, the internal backend calls EditRide with the driver's info and with a MATCHED RideStatus.

GetRide(userId: string)

=> Ride

Usage: polled every 20-90 seconds throughout the trip to update the ride's estimated price, it's time to destination, whether it's been canceled, etc..

EditRide(userId: string, [...params?: all properties on the Ride object that need to be edited])

=> Ride

AcceptRide(userId: string)

=> void

Calls EditRide(userId, MATCHED) and SetDriverStatus(userId, IN_RIDE).

CancelRide(userId: string)

=> void

Wrapper around EditRide -- effectively calls EditRide(userId, CANCELLED).

PushLocation(userId: string, location: Geolocation)

=> void

Usage: continuously called by a driver's phone throughout a ride; pushes the driver's location to the relevant passenger who's streaming the location; the passenger is the one associated with the Ride tied to the passed userId.

6. UberPool

As a stretch goal, your interviewer might ask you to think about how you'd expand your design to handle UberPool rides.

UberPool rides allow multiple passengers (different Uber accounts) to share an Uber ride for a cheaper price.

One way to handle UberPool rides would be to allow Ride objects to have multiple passengerInfos. In this case, Rides would also have to maintain a list of all destinations that the ride will stop at, as well as the relevant final destinations for individual passengers.

Perhaps a cleaner way to handle UberPool rides would be to introduce an entirely new entity, a PoolRide entity, which would have a list of Rides attached to it. Passengers would still call the CreateRide endpoint when booking an UberPool ride, and so they would still have their own, normal Ride entity, but this entity would be attached to a PoolRide entity with the rest of the UberPool ride information.

Drivers would likely have an extra DriverStatus value to indicate if they were in a ride but still accepting new UberPool passengers.

Most of the other functionality would remain the same; passengers and drivers would still continuously poll the GetRide endpoint for updated information about the ride, passengers would still stream their driver's location, passengers would still be able to cancel their individual rides, etc..

3. Design Google Drive.

Clarifying Questions to Ask :

Q: Are we just designing the storage aspect of Google Drive, or are we also designing some of the related products like Google Docs, Sheets, Slides, Drawings, etc.?

A: We're just designing the core Google Drive product, which is indeed the storage product. In other words, users can create folders and upload files, which effectively stores them in the cloud. Also, for simplicity, we can refer to folders and files as "entities".

Q: There are a lot of features on Google Drive, like shared company drives vs. personal drives, permissions on entities (ACLs), starred files, recently-accessed files, etc.. Are we designing all of these features or just some of them?

A: Let's keep things narrow and imagine that we're designing a personal Google Drive (so you can forget about shared company drives). In a personal Google Drive, users can store entities, and that's all that you should take care of. Ignore any feature that isn't core to the storage aspect of Google Drive; ignore things like starred files, recently-accessed files, etc.. You can even ignore sharing entities for this design.

Q: Since we're primarily concerned with storing entities, are we supporting all basic CRUD operations like creating, deleting, renaming, and moving entities?

A: Yes, but to clarify, creating a file is actually uploading a file, folders have to be created (they can't be uploaded), and we also want to support downloading files.

Q: Are we just designing the Google Drive web application, or are we also designing a desktop client for Google drive?

A: We're just designing the functionality of the Google Drive web application.

Q: Since we're not dealing with sharing entities, should we handle multiple users in a single folder at the same time, or can we assume that this will never happen?

A: While we're not designing the sharing feature, let's still handle what would happen if multiple clients were in a single folder at the same time (two tabs from the same browser, for example). In this case, we would want changes made in that folder to be reflected to all clients within 10 seconds. But for the purpose of this question, let's not worry about conflicts or anything like that (i.e., assume that two clients won't make changes to the same file or folder at the same time).

Q: How many people are we building this system for?

A: This system should serve about a billion users and handle 15GB per user on average.

Q: What kind of reliability or guarantees does this Google Drive service give to its users?

A: First and foremost, data loss isn't tolerated at all; we need to make sure that once a file is uploaded or a folder is created, it won't disappear until the user deletes it. As for availability, we need this system to be highly available.

SOLUTION WALKTHROUGH :

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

We're designing the core user flow of the Google Drive web application. This consists of storing two main entities: folders and files. More specifically, the system should allow users to create folders, upload and download files, and rename and move entities once they're stored. We don't have to worry about ACLs, sharing entities, or any other auxiliary Google Drive features.

We're going to be building this system at a very large scale, assuming 1 billion users, each with 15GB of data stored in Google Drive on average. This adds up to approximately 15,000 PB of data in total, without counting any metadata that we might store for each entity, like its name or its type.

We need this service to be Highly Available and also very redundant. No data that's successfully stored in Google Drive can ever be lost, even through catastrophic failures in an entire region of the world.

2. Coming Up With A Plan

It's important to organize ourselves and to lay out a clear plan regarding how we're going to tackle our design. What are the major, distinguishable components of our how system?

First of all, we'll need to support the following operations:

For Files

- UploadFile
- DownloadFile
- DeleteFile
- RenameFile
- MoveFile

For Folders

- CreateFolder
- GetFolder
- DeleteFolder
- RenameFolder
- MoveFolder

Secondly, we'll have to come up with a proper storage solution for two types of data:

File Contents: The contents of the files uploaded to Google Drive. These are opaque bytes with no particular structure or format.

Entity Info: The metadata for each entity. This might include fields like entityID, ownerID, lastModified, entityName, entityType. This list is non-exhaustive, and we'll most likely add to it later on.

Let's start by going over the storage solutions that we want to use, and then we'll go through what happens when each of the operations outlined above is performed.

3. Storing Entity Info

To store entity information, we can use key-value stores. Since we need high availability and data replication, we need to use something like Etcd, Zookeeper, or Google Cloud Spanner (as a K-V store) that gives us both of those guarantees as well as consistency (as opposed to DynamoDB, for instance, which would give us only eventual consistency).

Since we're going to be dealing with many gigabytes of entity information (given that we're serving a billion users), we'll need to shard this data across multiple clusters of these K-V stores. Sharding on entityID means that we'll lose the ability to perform batch operations, which these key-value stores give us out of the box and which we'll need when we move entities around (for instance, moving a file from one folder to another would involve editing the metadata of 3 entities; if they were located in 3 different shards that wouldn't be great). Instead, we can shard based on the ownerID of the entity, which means that we can edit the metadata of multiple entities atomically with a transaction, so long as the entities belong to the same user.

Given the traffic that this website needs to serve, we can have a layer of proxies for entity information, load balanced on a hash of the ownerID. The proxies could have some caching, as well as perform ACL checks when we eventually decide to support them. The proxies would live at the regional level, whereas the source-of-truth key-value stores would be accessed globally.

4. Storing File Data

When dealing with potentially very large uploads and data storage, it's often advantageous to split up data into blobs that can be pieced back together to form the original data. When uploading a file, the request will be load balanced across multiple servers that we'll call "blob splitters", and these blob splitters will have the job of splitting files into blobs and storing these blobs in some global blob-storage solution like GCS or S3 (since we're designing Google Drive, it might not be a great idea to pick S3 over GCS :P).

One thing to keep in mind is that we need a lot of redundancy for the data that we're uploading in order to prevent data loss. So we'll probably want to adopt a strategy like: try pushing to 3 different GCS buckets and consider a write successful only if it went through in at least 2 buckets. This way we always have redundancy without necessarily sacrificing availability. In the background, we can have an extra service in charge of further replicating the data to other buckets in an async manner. For our main 3 buckets, we'll want to pick buckets in 3 different availability zones to avoid having all of our redundant storage get wiped out by potential catastrophic failures in the event of a natural disaster or huge power outage.

In order to avoid having multiple identical blobs stored in our blob stores, we'll name the blobs after a hash of their content. This technique is called Content-Addressable Storage, and by using it, we essentially make all blobs immutable in storage. When a file changes, we simply upload the entire new resulting blobs under their new names computed by hashing their new contents.

This immutability is very powerful, in part because it means that we can very easily introduce a caching layer between the blob splitters and the buckets, without worrying about keeping caches in sync with the main source of truth when edits are made--an edit just means that we're dealing with a completely different blob.

5. Entity Info Structure

Since folders and files will both have common bits of metadata, we can have them share the same structure. The difference will be that folders will have an `is_folder` flag set to true and a list of `children_ids`, which will point to the entity information for the folders and files within the folder in question. Files will have an `is_folder` flag set to false and a `blobs` field, which will have the IDs of all of the blobs that make up the data within the relevant file. Both entities can also have a `parent_id` field, which will point to the entity information of the entity's parent folder. This will help us quickly find parents when moving files and folders.

File Info

```
{
  blobs: ['blob_content_hash_0', 'blob_content_hash_1'],
  id: 'some_unique_entity_id'
  is_folder: false,
  name: 'some_file_name',
  owner_id: 'id_of_owner',
  parent_id: 'id_of_parent',
}
```

Folder Info

```
{
  children_ids: ['id_of_child_0', 'id_of_child_1'],
  id: 'some_unique_entity_id'
```

```
is_folder: true,  
name: 'some_folder_name',  
owner_id: 'id_of_owner',  
parent_id: 'id_of_parent',  
}
```

6. Garbage Collection

Any change to an existing file will create a whole new blob and de-reference the old one. Furthermore, any deleted file will also de-reference the file's blobs. This means that we'll eventually end up with a lot of orphaned blobs that are basically unused and taking up storage for no reason. We'll need a way to get rid of these blobs to free some space.

We can have a Garbage Collection service that watches the entity-info K-V stores and keeps counts of the number of times every blob is referenced by files; these counts can be stored in a SQL table.

Reference counts will get updated whenever files are uploaded and deleted. When the reference count for a particular blob reaches 0, the Garbage Collector can mark the blob in question as orphaned in the relevant blob stores, and the blob will be safely deleted after some time if it hasn't been accessed.

7. End To End API Flow

Now that we've designed the entire system, we can walk through what happens when a user performs any of the operations we listed above.

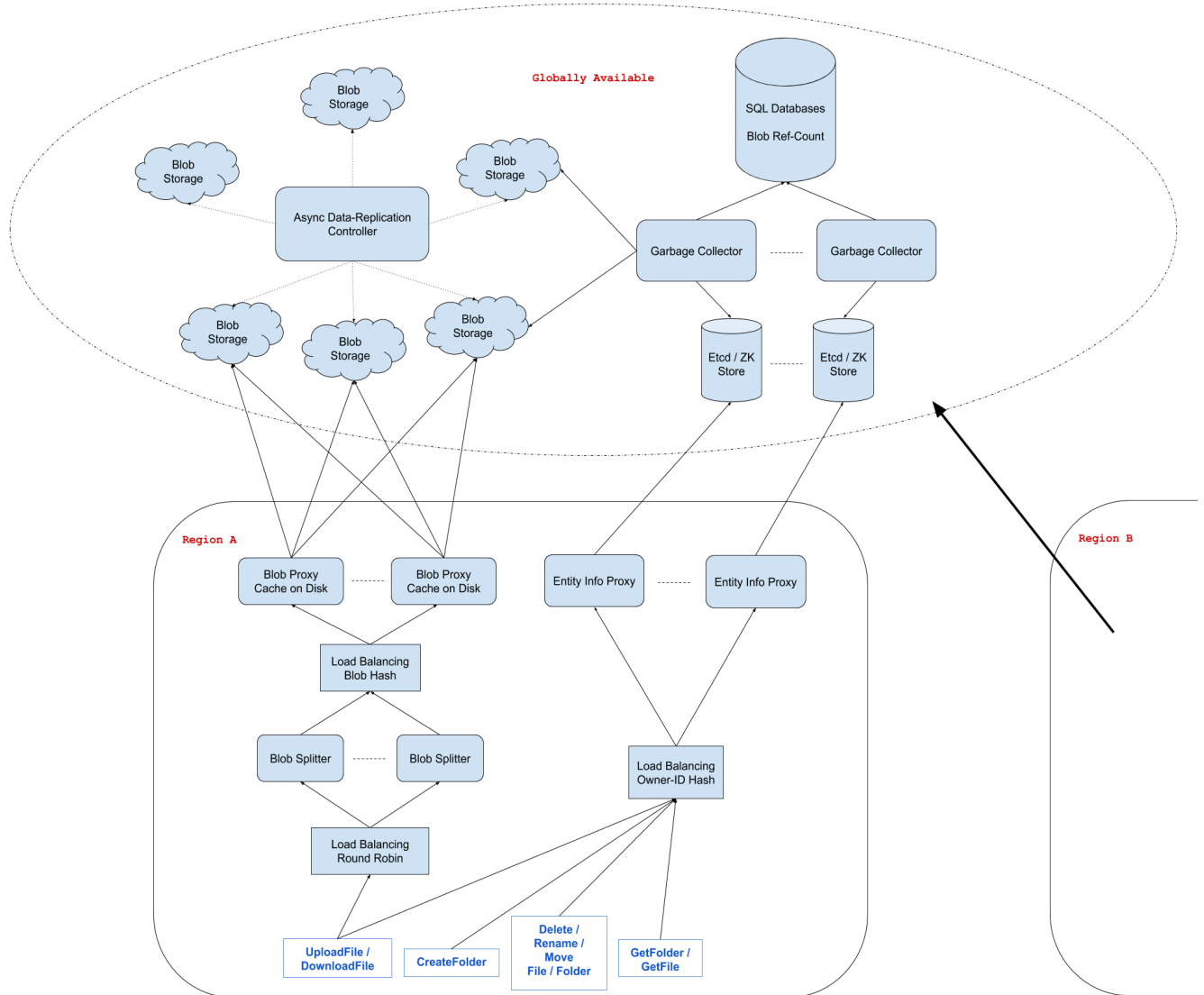
CreateFolder is simple; since folders don't have a blob-storage component, creating a folder just involves storing some metadata in our key-value stores.

UploadFile works in two steps. The first is to store the blobs that make up the file in the blob storage. Once the blobs are persisted, we can create the file-info object, store the blob-content hashes inside its blobs field, and write this metadata to our key-value stores.

DownloadFile fetches the file's metadata from our key-value stores given the file's ID. The metadata contains the hashes of all of the blobs that make up the content of the file, which we can use to fetch all of the blobs from blob storage. We can then assemble them into the file and save it onto local disk.

All of the Get, Rename, Move, and Delete operations atomically change the metadata of one or several entities within our key-value stores using the transaction guarantees that they give us.

8. System Diagram



4. Design The Reddit API

Summary :

Design an API for Reddit subreddits given the following information.

The API includes these 2 entities:

- User | userId: string, ...
- SubReddit | subredditId: string, ...

Both of these entities likely have other fields, but for the purpose of this question, those other fields aren't needed.

Your API should support the basic functionality of a subreddit on Reddit.

Many systems design questions are intentionally left very vague and are literally given in the form of `Design Foobar`. It's your job to ask clarifying questions to better understand the system that you have to build.

We've laid out some of these questions below; their answers should give you some guidance on the problem. Before looking at them, we encourage you to take few minutes to think about what questions you'd ask in a real interview.

Clarifying Questions to Ask :

Question 1

Q: To make sure that we're on the same page: a subreddit is an online community where users can write posts, comment on posts, upvote / downvote

posts, share posts, report posts, become moderators, etc.--is this correct, and are we designing all of this functionality?

A: Yes, that's correct, but let's keep things simple and focus only on writing posts, writing comments, and upvoting / downvoting. You can forget about all of the auxiliary features like sharing, reporting, moderating, etc..

Question 2

Q: So we're really focusing on the very narrow but core aspect of a subreddit: writing posts, commenting on them, and voting on them.

A: Yes.

Question 3

Q: I'm thinking of defining the schemas for the main entities that live within a subreddit and then defining their CRUD operations -- methods like Create/Get/Edit/Delete/List<Entity> -- is this in line with what you're asking me to do?

A: Yes, and make sure to include method signatures -- what each method takes in and what each method returns. Also include the types of each argument.

Question 4

Q: The entities that I've identified are Posts, Comments, and Votes (upvotes and downvotes). Does this seem accurate?

A: Yes. These are the 3 core entities that you should be defining and whose APIs you're designing.

Question 5

Q: Is there any other functionality of a subreddit that we should design?

A: Yes, you should also allow people to award posts. Awards are a special currency that can be bought for real money and gifted to comments and posts. Users can buy some quantity of awards in exchange for real money, and they can give awards to posts and comments (one award per post / comment).

SOLUTION WALKTHROUGH :

1. Gathering Requirements

As with any API design interview question, the first thing that we want to do is to gather API requirements; we need to figure out what API we're building exactly.

We're designing the core user flow of the subreddit functionality on Reddit. Users can write posts on subreddits, they can comment on posts, and they can upvote / downvote posts and comments.

We're going to be defining three primary entities: Posts, Comments, and Votes, as well as their respective CRUD operations.

We're also going to be designing an API for buying and giving awards on Reddit.

2. Coming Up With A Plan

It's important to organize ourselves and to lay out a clear plan regarding how we're going to tackle our design. What are the major, potentially contentious parts of our API? Why are we making certain design decisions?

The first major point of contention is whether to store votes only on Comments and Posts, and to cast votes by calling the EditComment and EditPost methods, or whether to store them as entirely separate entities--siblings of Posts and Comments, so to speak. Storing them as separate entities makes it much more straightforward to edit or remove a particular user's votes (by just calling EditVote, for instance), so we'll go with this approach.

We can then plan to tackle Posts, Comments, and Votes in this order, since they'll likely share some common structure.

3. Posts

Posts will have an id, the id of their creator (i.e., the user who writes them), the id of the subreddit that they're on, a description and a title, and a timestamp of when they're created.

Posts will also have a count of their votes, comments, and awards, to be displayed on the UI. We can imagine that some backend service will be calculating or updating these numbers when some of the Comment, Vote, and Award CRUD operations are performed.

Lastly, Posts will have optional `deletedAt` and `currentVote` fields. Subreddits display posts that have been removed with a special message; we can use the `deletedAt` field to accomplish this. The `currentVote` field will be used to display to a user whether or not they've cast a vote on a post. This field will likely be populated by the backend upon fetching Posts or when casting Votes.

- `postId`: string
- `creatorId`: string
- `subredditId`: string
- `title`: string
- `description`: string
- `createdAt`: timestamp
- `votesCount`: int
- `commentsCount`: int
- `awardsCount`: int
- `deletedAt?`: timestamp
- `currentVote?`: enum UP/DOWN

Our `CreatePost`, `EditPost`, `GetPost`, and `DeletePost` methods will be very straightforward. One thing to note, however, is that all of these operations will take in the `userId` of the user performing them; this id, which will likely contain authentication information, will be used for ACL checks to see if the user performing the operations has the necessary permission(s) to do so.

`CreatePost(userId: string, subredditId: string, title: string, description: string)`

`=> Post`

`EditPost(userId: string, postId: string, title: string, description: string)`

=> Post

GetPost(userId: string, postId: string)

=> Post

DeletePost(userId: string, postId: string)

=> Post

Since we can expect to have hundreds, if not thousands, of posts on a given subreddit, our ListPosts method will have to be paginated. The method will take in optional pageSize and pageToken parameters and will return a list of posts of at most length pageSize as well as a nextPageToken--the token to be fed to the method to retrieve the next page of posts.

ListPosts(userId: string, subredditId: string, pageSize?: int, pageToken?: string)

=> (Post[], nextPageToken?)

4. Comments

Comments will be similar to Posts. They'll have an id, the id of their creator (i.e., the user who writes them), the id of the post that they're on, a content string, and the same other fields as Posts have. The only difference is that Comments will also have an optional parentId pointing to the parent post or parent comment of the comment. This id will allow the Reddit UI to reconstruct Comment trees to properly display (indent) replies. The UI can also sort comments within a reply thread by their createdAt timestamps or by their votesCount.

- commentId: string
- creatorId: string
- postId: string
- createdAt: timestamp
- content: string
- votesCount: int

- awardsCount: int
- parentId?: string
- deletedAt?: timestamp
- currentVote?: enum UP/DOWN

Our CRUD operations for Comments will be very similar to those for Posts, except that the CreateComment method will also take in an optional parentId pointing to the comment that it's replying to, if relevant.

CreateComment(userId: string, postId: string, content: string, parentId?: string)

=> Comment

EditComment(userId: string, commentId: string, content: string)

=> Comment

GetComment(userId: string, commentId: string)

=> Comment

DeleteComment(userId: string, commentId: string)

=> Comment

ListComments(userId: string, postId: string, pageSize?: int, pageToken?: string)

=> (Comment[], nextPageToken?)

5. Votes

Votes will have an id, the id of their creator (i.e., user who casts them), the id of their target (i.e., the post or comment that they're on), and a type, which will be a simple UP/DOWN enum. They could also have a createdAt timestamp for good measure.

- `voteId`: string
- `creatorId`: string
- `targetId`: string
- `type`: enum UP/DOWN
- `createdAt`: timestamp

Since it doesn't seem like getting a single vote or listing votes would be very useful for our feature, we'll skip designing those endpoints (though they would be straightforward).

Our `CreateVote`, `EditVote`, and `DeleteVote` methods will be simple and useful. The `CreateVote` method will be used when a user casts a new vote on a post or comment; the `EditVote` method will be used when a user has already cast a vote on a post or comment and casts the opposite vote on that same post or comment; and the `DeleteVote` method will be used when a user has already cast a vote on a post or comment and just removes that same vote.

`CreateVote(userId: string, targetId: string, type: enum UP/DOWN)`

=> `Vote`

`EditVote(userId: string, voteId: string, type: enum UP/DOWN)`

=> `Vote`

`DeleteVote(userId: string, voteId: string)`

=> `Vote`

6. Awards

We can define two simple endpoints to handle buying and giving awards. The endpoint to buy awards will take in a paymentToken, which will be a string that contains all of the necessary information to process a payment. The endpoint to give an award will take in a targetId, which will be the id of the post or comment that the award is being given to.

BuyAwards(userId: string, paymentToken: string, quantity: int)

GiveAward(userId: string, targetId: string)

5. Design A StockBroker

Design a stockbroker: a platform that acts as the intermediary between end-customers and some central stock exchange.

Many systems design questions are intentionally left very vague and are literally given in the form of `Design Foobar`. It's your job to ask clarifying questions to better understand the system that you have to build.

We've laid out some of these questions below; their answers should give you some guidance on the problem. Before looking at them, we encourage you to take few minutes to think about what questions you'd ask in a real interview.

Clarifying Questions to Ask :

Question 1

Q: What do we mean exactly by a stock broker? Is this something like Robinhood or Etrade?

A: Yes, exactly.

Question 2

Q: What is the platform supposed to support exactly? Are we just supporting the ability for customers to buy and sell stocks, or are we supporting more? For instance, are we allowing other types of securities like options and futures to be traded on our platform? Are we supporting special types of orders like limit orders and stop losses?

A: We're only supporting market orders on stocks in this design. A market order means that, given a placed order to buy or sell a stock, we should try to execute the order as soon as possible regardless of the stock price. We also aren't designing any "margin" system, so the available balance is the source of truth for what can be bought.

Question 3

Q: Are we designing any of the auxiliary aspects of the stock brokerage, like depositing and withdrawing funds, downloading tax documents, etc.?

A: No -- we're just designing the core trading aspect of the platform.

Question 4

Q: Are we just designing the system to place trades? Do we want to support other trade-related operations like getting trade statuses? In other words, how comprehensive should the API that's going to support this platform be?

A: In essence, you're only designing a system around a PlaceTrade API call from the user, but you should define that API call (inputs, response, etc.).

Question 5

Q: Where does a customer's balance live? Is the platform pulling a customer's money directly from their bank account, or are we expecting that customers will have already deposited funds into the platform somehow? In other words, are we ever directly interacting with banks?

A: No, you won't be interacting with banks. You can assume that customers have already deposited funds into the platform, and you can further assume that you have a SQL table with the balance for each customer who wants to make a trade.

Question 6

Q: How many customers are we building this for? And is our customer-base a global one?

A: Millions of customers, millions of trades a day. Let's assume that our customers are only located in 1 region -- the U.S., for instance.

Question 7

Q: What kind of availability are we looking for?

A: As high as possible, with this kind of service people can lose a lot of money if the system is down even for a few minutes.

Question 8

Q: Are we also designing the UI for this platform? What kinds of clients can we assume we have to support?

A: You don't have to design the UI, but you should design the PlaceTrade API call that a UI would be making to your backend. Clients would be either a mobile app or a webapp.

Question 9

Q: So we want to design the API for the actual brokerage, that itself interacts with some central stock exchange on behalf of customers. Does this exchange have an API? If yes, do we know what it looks like, and do we have any guarantees about it?

A: Yes, the exchange has an API, and your platform's API (the PlaceTrade call) will have to interact with the exchange's API. As far as that's concerned, you can assume that the call to the exchange to make an actual trade will take in a callback (in addition to the info about the trade) that will get executed when that trade completes at the exchange level (meaning, when the trade either gets FILLED or REJECTED, this callback will be executed). You can also assume that the exchange's system is highly available--your callback will always get executed at least once.

SOLUTION WALKTHROUGH :

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

We're building a stock-brokerage platform like Robinhood that functions as the intermediary between end-customers and some central stock exchange. The idea is that the central stock exchange is the platform that actually executes stock trades, whereas the stockbroker is just the platform that customers talk to when they want to place a trade--the stock brokerage is "simpler" and more "human-readable", so to speak.

We only care about supporting market trades--trades that are executed at the current stock price--and we can assume that our system stores customer balances (i.e., funds that customers may have previously deposited) in a SQL table.

We need to design a PlaceTrade API call, and we know that the central exchange's equivalent API method will take in a callback that's guaranteed to be executed upon completion of a call to that API method.

We're designing this system to support millions of trades per day coming from millions of customers in a single region (the U.S., for example). We want the system to be highly available.

2. Coming Up With A Plan

It's important to organize ourselves and to lay out a clear plan regarding how we're going to tackle our design. What are the major, distinguishable components of our system?

We'll approach the design front to back:

the PlaceTrade API call that clients will make

the API server(s) handling client API calls

the system in charge of executing orders for each customer

We'll need to make sure that the following hold:

trades can never be stuck forever without either succeeding or failing to be executed

a single customer's trades have to be executed in the order in which they were placed

balances can never go in the negatives

3. API Call

The core API call that we have to implement is PlaceTrade.

We'll define its signature as:

```
PlaceTrade(  
    customerId: string,  
    stockTicker: string,  
    type: string (BUY/SELL),  
    quantity: integer,
```

```
) => (  
  tradeId: string,  
  stockTicker: string,  
  type: string (BUY/SELL),  
  quantity: integer,  
  createdAt: timestamp,  
  status: string (PLACED),  
  reason: string,  
)
```

The customer ID can be derived from an authentication token that's only known to the user and that's passed into the API call.

The status can be one of:

- PLACED
- IN PROGRESS
- FILLED
- REJECTED

That being said, PLACED will actually be the defacto status here, because the other statuses will be asynchronously set once the exchange executes our callback. In other words, the trade status will always be PLACED when the PlaceTrade API call returns, but we can imagine that a GetTrade API call could return statuses other than PLACED.

Potential reasons for a REJECTED trade might be:

insufficient funds

random error

past market hours

4. API Server(s)

We'll need multiple API servers to handle all of the incoming requests. Since we don't need any caching when making trades, we don't need any server stickiness, and we can just use some round-robin load balancing to distribute incoming requests between our API servers.

Once API servers receive a PlaceTrade call, they'll store the trade in a SQL table. This table needs to be in the same SQL database as the one that the balances table is in, because we'll need to use ACID transactions to alter both tables in an atomic way.

The SQL table for trades will look like this:

- id: string, a random, auto-generated string
- customer_id: string, the id of the customer making the trade
- stockTicker: string, the ticker symbol of the stock being traded
- type: string, either BUY or SELL
- quantity: integer (no fractional shares), the number of shares to trade
- status: string, the status of the trade; starts as PLACED
- created_at: timestamp, the time when the trade was created
- reason: string, the human-readable justification of the trade's status

The SQL table for balances will look like this:

- id: string, a random, auto-generated string
- customer_id: string, the id of the customer related to the balance
- amount: float, the amount of money that the customer has in USD
- last_modified: timestamp, the time when the balance was last modified

5. Trade-Execution Queue

With hundreds of orders placed every second, the trades table will be pretty massive. We'll need to figure out a robust way to actually execute our trades and to update our table, all the while making sure of a couple of things:

We want to make sure that for a single customer, we only process a single BUY trade at any time, because we need to prevent the customer's balance from ever reaching negative values.

Given the nature of market orders, we never know the exact dollar value that a trade will get executed at in the exchange until we get a response from the exchange, so we have to speak to the exchange in order to know whether the trade can go through.

We can design this part of our system with a Publish/Subscribe pattern. The idea is to use a message queue like Apache Kafka or Google Cloud Pub/Sub and to have a set of topics that customer ids map to. This gives us at-least-once delivery semantics to make sure that we don't miss new trades. When a customer makes a trade, the API server writes a row to the database and also creates a message that gets routed to a topic for that customer (using hashing), notifying the topic's subscriber that there's a new trade.

This gives us a guarantee that for a single customer, we only have a single thread trying to execute their trades at any time.

Subscribers of topics can be rings of 3 workers (clusters of servers, essentially) that use leader election to have 1 master worker do the work for the cluster (this is for our system's high availability)--the leader grabs messages as they get pushed to the topic and executes the trades for the customers contained in the messages by calling the exchange. As mentioned above, a single customer's trades are only ever handled by the same cluster of workers, which makes our logic and our SQL queries cleaner.

As far as how many topics and clusters of workers we'll need, we can do some rough estimation. If we plan to execute millions of trades per day, that comes down to about 10-100 trades per second given open trading hours during a third of a day and non-uniform trading patterns. If we assume that the core execution logic lasts about a second, then we should have roughly 10-100 topics and clusters of workers to process trades in parallel.

~100,000 seconds per day ($3600 * 24$)

~1,000,000 trades per day

trades bunched in 1/3rd of the day

--> $(1,000,000 / 100,000) * 3 = \sim 30$ trades per second

6. Trade-Execution Logic

The subscribers (our workers) are streaming / waiting for messages. Imagine the following message were to arrive in the topic queue:

```
{"customerId": "c1"}
```

The following would be pseudo-code for the worker logic:

```
// We get the oldest trade that isn't in a terminal state.
trade = SELECT * FROM trades WHERE
    customer_id = 'c1' AND
    (status = 'PLACED' OR status = 'IN PROGRESS')
    ORDER BY created_at ASC LIMIT 1;

// If the trade is PLACED, we know that it's effectively
// ready to be executed. We set it as IN PROGRESS.
if trade.status == "PLACED" {
    UPDATE trades SET status = 'IN PROGRESS' WHERE id = trade.id;
}

// In the event that the trade somehow already exists in the
// exchange, the callback will do the work for us.
if exchange.TradeExists(trade.id) {
    return;
}

// We get the balance for the customer.
balance = SELECT amount FROM balances WHERE
    customer_id = 'c1';

// This is the callback that the exchange will execute once
// the trade actually completes. We'll define it further down
// in the walkthrough.
callback = ...

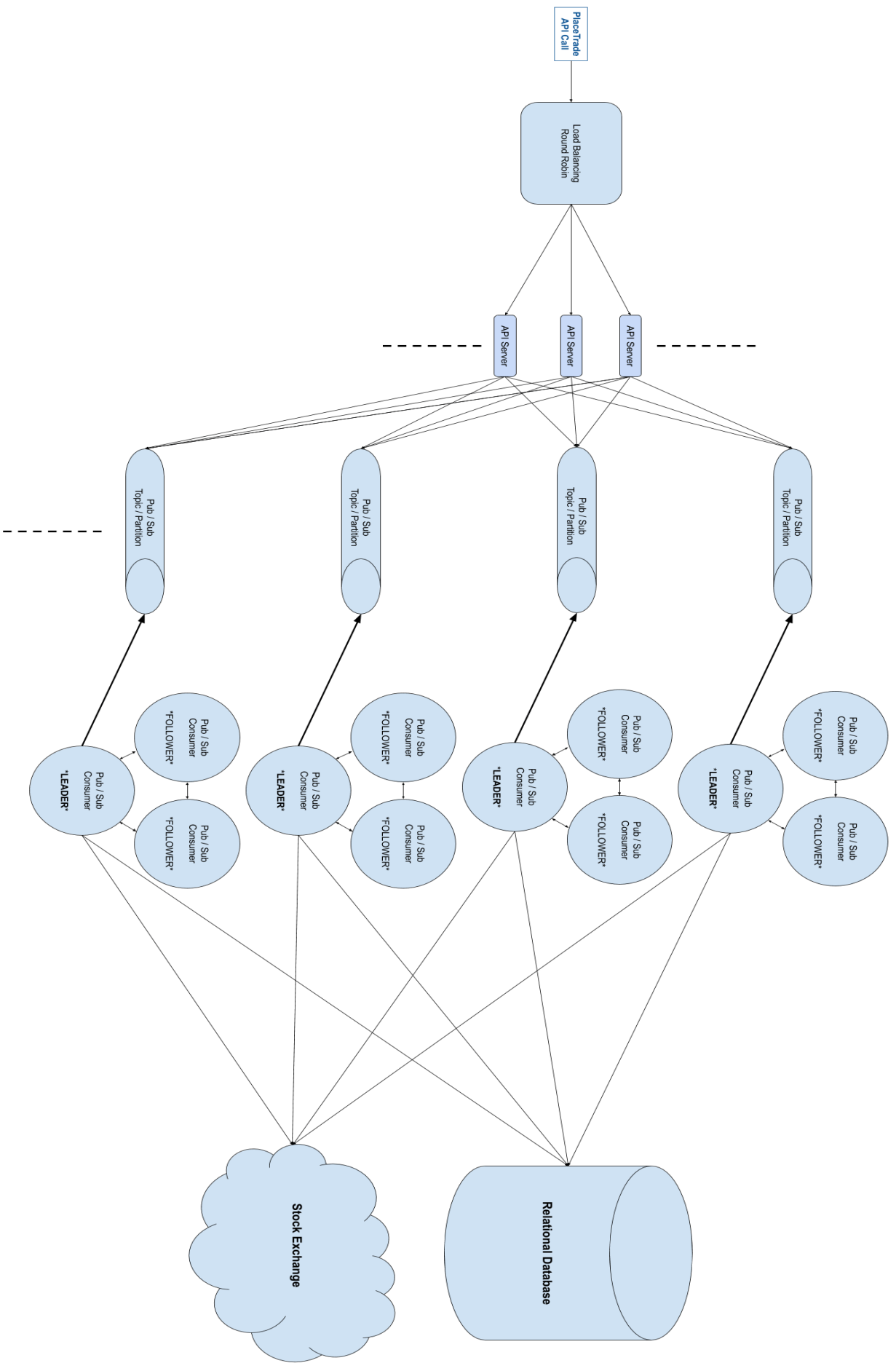
exchange.Execute(
    trade.stockTicker,
    trade.type,
    trade.quantity,
    max_price = balance,
    callback,
)
```

7. Exchange Callback

Below is some pseudo code for the exchange callback:

```
function exchange_callback(exchange_trade) {
  if exchange_trade.status == 'FILLED' {
    BEGIN TRANSACTION;
    trade = SELECT * FROM trades WHERE id = database_trade.id;
    if trade.status <> 'IN PROGRESS' {
      ROLLBACK;
      pubsub.send({customer_id: database_trade.customer_id});
      return;
    }
    UPDATE balances SET amount -= exchange_trade.amount WHERE customer_id = database_trade.customer_id;
    UPDATE trades SET status = 'FILLED' WHERE id = database_trade.id;
    COMMIT;
  } else if exchange_trade.status == 'REJECTED' {
    BEGIN TRANSACTION;
    UPDATE trades SET status = 'REJECTED' WHERE id = database_trade.id;
    UPDATE trades SET reason = exchange_trade.reason WHERE id = database_trade.id;
    COMMIT;
  }
  pubsub.send({customer_id: database_trade.customer_id});
  return http.status(200);
}
```


8. System Diagram



6. Design A AlgoExpert

Clarifying Questions to Ask :

Question 1

Q: Are we designing the entire AlgoExpert platform or just a specific part of it, like the coding workspace?

A: Since we only have about 45 minutes, you should just design the core user flow of the AlgoExpert platform. The core user flow includes users landing on the home page of the website, going to the questions list, marking questions as complete or in progress, and then writing and running code in various languages for each language. Don't worry about payments or authentication; you can just assume that you have these services working already (by the way, we mainly rely on third-party services here, like Stripe, PayPal, and OAuth2).

Question 2

Q: AlgoExpert doesn't seem like a system of utmost criticality (like a hospital system or airplane software); are we okay with 2 to 3 nines of availability for the system?

A: Yes, this seems fine--no need to focus too much on making the system highly available.

Question 3

Q: How many customers should we be building this for? Is AlgoExpert's audience global or limited to one country?

A: AlgoExpert's website receives hundreds of thousands of users every month, and tens of thousands of users may be on the website at any point in time. We want the website to feel very responsive to people everywhere in the world, and the U.S. and India are the platform's top 2 markets that we especially want to cater to.

Question 4

Q: Does AlgoExpert make changes to its content (questions list and question solutions) often?

A: Yes--every couple of days on average. And we like to have our changes reflected in production globally within the hour.

Question 5



Q: How much of the code-execution engine behind the coding workspace should we be designing? Do we have to worry about the security aspect of running random user code on our servers?

A: You can disregard the security aspects of the code-execution engine and just focus on its core functionality--the ability to run code in various languages at any given time with acceptable latency.

Question 6



Q: While we'll care about latency across the entire system, the code-execution engine seems like the place where we'll care about it most, since it's very interactive, and it also seems like the toughest part of our system to support low latencies; are we okay with anywhere between 1 and 3 seconds for the average run-code latency?

A: Yes--this seems reasonable and acceptable from a product point of view.

SOLUTION WALKTHROUGH :

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

From the answers we were given to our clarifying questions (see Prompt Box), we're building the core AlgoExpert user flow, which includes users landing on the website, accessing questions, marking them as complete, writing code, running code, and having their code saved.

We don't need to worry about payments or authentication, and we don't need to go too deep into the code-execution engine.

We're building this platform for a global audience, with an emphasis on U.S. and India users, and we don't need to overly optimize our system's availability. We probably don't need more than two or three nines, because we're not building a health or security system, and this gets us somewhere between 8 hours and 3 days of downtime per year, which is reasonable. All in all, this means that we don't need to worry too much about availability.

We care about latency and throughput within reason, but apart from the code-execution engine, this doesn't seem like a particularly difficult aspect of our system.

2. Coming Up With A Plan

It's important to organize ourselves and to lay out a clear plan regarding how we're going to tackle our design. What are the major, distinguishable components of our system?

On the one hand, AlgoExpert has a lot of static content; the entire home page, for instance, is static, and it has a lot of images. On the other hand, AlgoExpert isn't just a static website; it clearly has a lot of dynamic content that users themselves can generate (code that they can write, for example). So we'll need to have a robust API backing our UI, and given that user content gets saved on the website, we'll also need a database backing our API.

We can divide our system into 3 core components:

Static UI content

Accessing and interacting with questions (question completion status, saving solutions, etc.)

Ability to run code

Note that the second bullet point will likely get further divided.

3. Static UI Content

For the UI static content, we can put public assets like images and JavaScript bundles in a blob store: S3 or Google Cloud Storage. Since we're catering to a global audience and we care about having a responsive website (especially the home page of the website), we might want to use a Content Delivery Network (CDN) to serve that content. This is especially important for a better mobile experience because of the slow connections that phones use.

4. Main Clusters And Load Balancing

For our main backend servers, we can have 2 primary clusters in the 2 important regions: U.S. and India.

We can have some DNS load balancing to route API requests to the cluster closest to the user issuing the requests, and within a region, we can have some path-based load balancing to separate our services (payments, authentication, code execution, etc.), especially since the code execution platform will probably need to run on different kinds of servers compared to those of the rest of the API. Each service can probably have a set of servers, and we can do some round-robin load balancing at that level (this is probably handled directly at the path-based load balancing layer).

5. Static API Content

There's a lot of static API content on AlgoExpert: namely, the list of questions and all of their solutions. We can store all of this data in a blob store for simplicity.

6. Caching

We can implement 2 layers of caching for this static API content.

We can have client-side caching; this will improve the user experience on the platform (users will only need to load questions once per session), and this will reduce the load on our backend servers (this will probably save 2-3 network calls per session).

We can also have some in-memory caching on our servers. If we approximate 100 questions with 10 languages and 5KB per solution, this should be less than $100 * 10 * 5000 \text{ bytes} = 5\text{MB}$ of total data to keep in memory, which is perfectly fine.

Since we were told that we want to make changes to static API content every couple of days and that we want those changes to be reflected in production as soon as possible, we can invalidate, evict and replace the data in our server-side caches every 30 minutes or so.

7. Access Control

Whenever you're designing a system, it's important to think about any potential access control that needs to be implemented. In the case of AlgoExpert, there's straightforward access control with regards to question content: users who haven't purchased AlgoExpert can't access individual questions. We can implement this fairly easily by

just making some internal API call whenever a user requests our static API content to figure out if the user owns the product before returning the full content for questions.

8. User Data Storage

For user data, we have to design the storage of question completion status and of user solutions to questions. Since this data will have to be queried a lot, a SQL database like Postgres or MySQL seems like a good choice.

We can have 2 tables. The first table might be `question_completion_status`, which would probably have the following columns:

- `id`: integer, primary key (an auto-incremented integer for instance)
- `user_id`: string, references the id of the user (can be obtained from auth)
- `question_id`: string, references the id of the question
- `completion_status`: string, enum to represent the completion status of the question

We could have a uniqueness constraint on (`user_id`, `question_id`) and an index on `user_id` for fast querying.

The second table might be `user_solutions`:

- `id`: integer, primary key (an auto-incremented integer for instance)
- `user_id`: string, references the id of the user (can be obtained from auth)
- `question_id`: string, references the id of the question
- `language`: string, references the language of the solution
- `solution`: string, contains the user's solution

We could have a uniqueness constraint on (`user_id`, `question_id`, `language`) and an index on `user_id` as well as one on `question_id`. If the number of languages goes up significantly, we might also want to index on `language` to allow for fast per-language querying so that the UI doesn't fetch all of a user's solutions at the same time (this might be a lot of data for slow connections).

9. Storage Performance

Marking questions as complete and typing code in the coding workspace (with a 1-3 second debounce for performance reasons) will issue API calls that write to the database. We likely won't get more than 1000 writes per second given our user numbers (assuming roughly 10,000 users on the platform at any given point in time), which SQL databases can definitely handle.

We can have 2 major database servers, each serving our 2 main regions: 1 in North America and 1 in India (perhaps serving Southeast Asia). If need be, we can add a 3rd cluster serving Europe exclusively (or other parts of the world, as our platform grows).

10. Inter-Region Replication

Since we'll have 2 primary database servers, we'll need to keep them up to date with each other. Fortunately, users on AlgoExpert don't share user-generated content; this means that we don't need data that's written to 1 database server to immediately be written to the other database server (this would likely have eliminated the latency improvements we got by having regional databases).

That being said, we do need to keep our databases up to date with each other, since users might travel around the world and hit a different database server than their typical one.

For this, we can have some async replication between our database servers. The replication can occur every 12 hours, and we can adjust this according to behavior in the system and amount of data that gets replicated across continents.

11. Code Execution

First of all, we should implement some rate limiting. A service like code execution lends itself perfectly to rate limiting, and we can implement some tier-based rate limiting using a K-V Store like Redis to easily prevent DoS attacks. We can limit the number of code runs to once every second, 3 times per 5 seconds, and 5 times per minute. This will prevent DoS attacks through the code-execution API, but it'll still allow for a good user experience when running code.

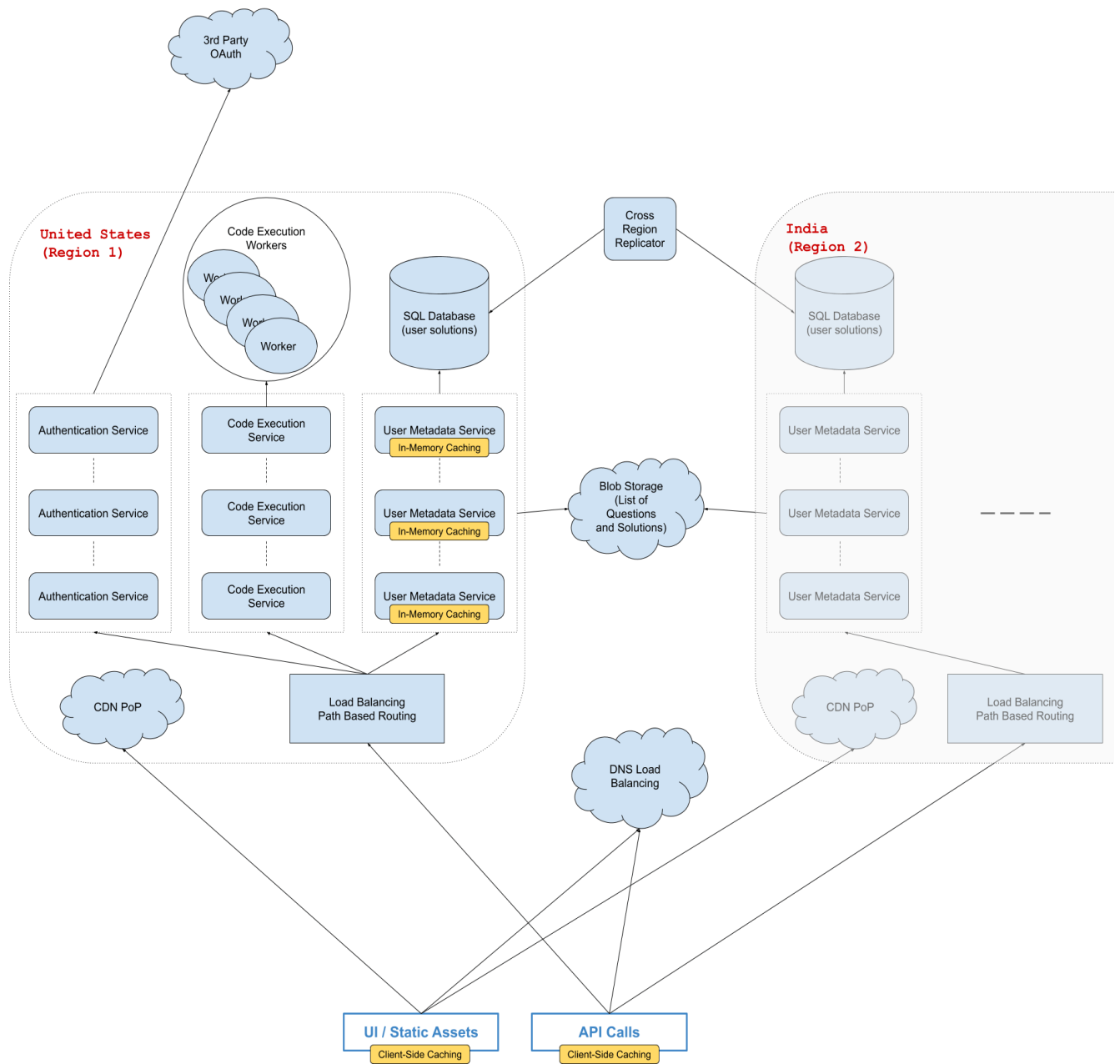
Since we want 1-3 seconds of latency for running code, we need to keep a set of special servers--our "workers"-- ready to run code at all times. They can each clean up after running user code (remove extra generated files as a result of compilation, for example) so that they don't need to be killed at any point. Our backend servers can contact a free worker and get the response from that worker when it's done running code (or if the code timed out), and our servers can return that to the UI in the same request.

Given that certain languages need to be compiled, we can estimate that it would take on average 1 second to compile and run the code for each language. People don't run code that often, so we can expect 10 run-codes per second in total given roughly 10,000 users on the website at once, so we'll probably need 10-100 machines to satisfy our original latency requirement of 1-3 seconds per run-code (10 machines if 10 run-codes per second is accurate, more if we experience higher load).

This design scales horizontally with our number of users, and it can scale vertically to make running code even faster (more CPU == faster runs).

Lastly, we can have some logging and monitoring in our system, especially for running code (tracking run-code events per language, per user, per question, average response time, etc.). This will help us automatically scale our clusters when user demand goes up or down. This can also be useful to know if any malicious behavior is happening with the code-execution engine.

12. System Diagram



7. Design Airbnb

Clarifying questions to ask

Question 1

Q: Like a lot of other sharing-economy products out there, Airbnb has two sides: a host-facing side and a renter-facing side. Are we designing both of these sides or just one of them?

A: Let's design both of these sides of the product.

Question 2

Q: Okay. So we're probably designing the system for hosts to create and maybe delete listings, and the system for renters to browse through properties, book them, and manage their bookings afterwards. Is that correct?

A: Yes for hosts; but let's actually just focus on browsing through listings and booking them for renters. We can ignore everything that happens after booking on the renter-facing side.

Question 3

Q: Okay, but for booking, is the idea that, when a user is browsing a property for a specific date range, the property gets temporarily reserved for them if they start the booking process?

A: Yes. More specifically, multiple users should be allowed to look at the same property, for the same date range, concurrently without issues. But once a user starts the booking process for a property, it should be reflected that this property is no longer available for the dates in question if another user tries to book it.

Question 4

Q: I see. But so, let's say two users are looking at the exact same property for an overlapping date range, and one user presses "Book Now", at which point they have to enter credit card information. Should we immediately lock the property for the other user for some predetermined period of time, like maybe 15 minutes, and if the first person actually goes through with booking the property, then this "lock" becomes permanent?

A: Yes, that makes sense. In real life, there might be slight differences, but for the sake of this design, let's go with that.

Question 5

Q: Okay. And do we want to design any auxiliary features like being able to contact hosts, authentication and payment services, etc., or are we really just focusing on browsing and reserving?

A: Let's really just focus on browsing and booking. We can ignore the rest.

Question 6

Q: I see. So, since it sounds like we're designing a pretty targeted part of the entire Airbnb service, I want to make sure that I know exactly every functionality that we want to support. My understanding is that users can go on the main Airbnb website or app,

they can look up properties based on certain criteria, like location, available date range, pricing, property details, etc., and then they can decide to book a location. As for hosts, they can basically just create a listing and delete it like I said earlier. Is that correct?

A: Yes. But actually, for this design, let's purely filter based on location and available date range as far as listing characteristics are concerned; let's not worry about other criteria like pricing and property details.

Question 7

Q: What is the scale that we're designing this for? Specifically, roughly how many listings and renters do we expect to cater to?

A: Let's only consider Airbnb's U.S. operations. So let's say 50 million users and 1 million listings.

Question 8

Q: Regarding systems characteristics like availability and latency, I'm assuming that, even if a renter looks up properties in a densely-populated area like NYC, where there might be a lot of listings, we care about serving these listings fast, accurately, and reliably. Is that correct?

A: Yes, that's correct. Ideally, we don't want any downtime for renters browsing listings.

Solution walkthrough

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

We're designing the core system behind Airbnb, which allows hosts to create property listings and renters to browse through these listings and book them.

Specifically, we'll want to support:

On the host side, creating and deleting listings.

On the renter side, browsing through listings, getting individual listings, and "reserving" listings.

"Reserving" listings should happen when a renter presses some "Book Now" button and should effectively lock (or reserve) the listing for some predetermined period of time (say, 15 minutes), during which any other renter shouldn't be able to reserve the listing or to even browse it (unless they were already browsing it).

We don't need to support anything that happens after a reservation is made, except for freeing up the reservation after 15 minutes if the renter doesn't follow through with the booking and making the reservation permanent if the renter does actually book the listing in question.

Regarding listings, we should focus on browsing and reserving them based on location and available date range; we can ignore any other property characteristics like price, number of bedrooms, etc.. Browsing listings should be as quick as possible, and it should reflect newly created listings as fast as possible. Lastly, reserved and booked listings shouldn't be browsable by renters.

Our system should serve a U.S.-based audience with approximately 50 million users and 1 million listings.

2. Coming Up With A Plan

We'll tackle this system by dividing it into two main sections:

The host side.

The renter side.

We can further divide the renter side as follows:

Browsing (listing) listings.

Getting a single listing.

Reserving a listing.

3. Listings Storage & Quadtree

First and foremost, we can expect to store all of our listings in a SQL table. This will be our primary source of truth for listings on Airbnb, and whenever a host creates or deletes a listing, this SQL table will be written to.

Then, since we care about the latency of browsing listings on Airbnb, and since this browsing will require querying listings based on their location, we can store our listings in a region quadtree, to be traversed for all browsing functionality.

Since we're optimizing for speed, it'll make sense to store this quadtree in memory on some auxiliary machine, which we can call a "geo index," but we need to make sure that we can actually fit this quadtree in memory.

In this quadtree, we'll need to store all the information about listings that needs to be displayed on the UI when a renter is browsing through listings: a title, a description, a link pointing to a property image, a unique listing ID, etc..

Assuming a single listing takes up roughly 10 KB of space (as an upper bound), some simple math confirms that we can store everything we need about listings in memory.

~10 KB per listing

~1 million listings

$\sim 10 \text{ KB} * 1000^2 = 10 \text{ GB}$

Since we'll be storing our quadtree in memory, we'll want to make sure that a single machine failure doesn't bring down the entire browsing functionality. To ensure this, we can set up a cluster of machines, each holding an instance of our quadtree in memory, and these machines can use leader election to safeguard us from machine failures.

Our quadtree solution works as follows: when our system boots up, the geo-index machines create the quadtree by querying our SQL table of listings. When listings are created or deleted, hosts first write to the SQL table, and then they synchronously update the geo-index leader's quadtree. Then, on an interval of say, 10 minutes, the geo-index leader and followers all recreate the quadtree from the SQL table, which allows them to stay up to date with new listings.

If the leader dies at any point, one of the followers takes its place, and data in the new leader's quadtree will be stale for at most a few minutes until the interval forces the quadtree to be recreated.

4. Listing Listings

When renters browse through listings, they'll have to hit some ListListings API endpoint. This API call will search through the geo-index leader's quadtree for relevant listings based on the location that the renter passes.

Finding relevant locations should be fairly straightforward and very fast, especially since we can estimate that our quadtree will have a depth of approximately 10, since 4^{10} is greater than 1 million.

That being said, we'll have to make sure that we don't return listings that are unavailable during the date range specified by the renter. In order to handle this, each listing in the quad tree will contain a list of unavailable date ranges, and we can perform a simple binary search on this list for each listing, in order to determine if the listing in question is available and therefore browsable by the renter.

We can also make sure that our quadtree returns only a subset of relevant listings for pagination purposes, and we can determine this subset by using an offset: the first page of relevant listings would have an offset of 0, the second page would have an offset of 50 (if we wanted pages to have a size of 50), the third page would have an offset of 100, and so on and so forth.

5. Getting Individual Listings

This API call should be extremely simple; we can expect to have listing IDs from the list of listings that a renter is browsing through, and we can simply query our SQL table of listings for the given ID.

6. Reserving Listings

Reserved listings will need to be reflected both in our quadtree and in our persistent storage solution. In our quadtree, because they'll have to be excluded from the list of

browsable listings; in our persistent storage solution, because if our quadtree needs to have them, then the main source of truth also needs to have them.

We can have a second SQL table for reservations, holding listing IDs as well as date ranges and timestamps for when their reservations expire. When a renter tries to start the booking process of a listing, the reservation table will first be checked to see if there's currently a reservation for the given listing during the specified date range; if there is, an error is returned to the renter; if there isn't, a reservation is made with an expiration timestamp 15 minutes into the future.

Following the write to the reservation table, we synchronously update the geo-index leader's quadtree with the new reservation. This new reservation will simply be an unavailability interval in the list of unavailabilities on the relevant listing, but we'll also specify an expiration for this unavailability, since it's a reservation.

A listing in our quadtree might look something like this:

```
{
  "unavailabilities": [
    {
      "range": ["2020-09-22T12:00:00-05:00", "2020-09-28T12:00:00-05:00"],
      "expiration": "2020-09-16T12:00:00-04:00"
    }
    {
      "range": ["2020-10-02T12:00:00-05:00", "2020-10-10T12:00:00-05:00"],
      "expiration": null
    },
  ],
  "title": "Listing Title",
  "description": "Listing Description",
  "thumbnailUrl": "Listing Thumbnail URL",
  "id": "Listing ID"
}
```

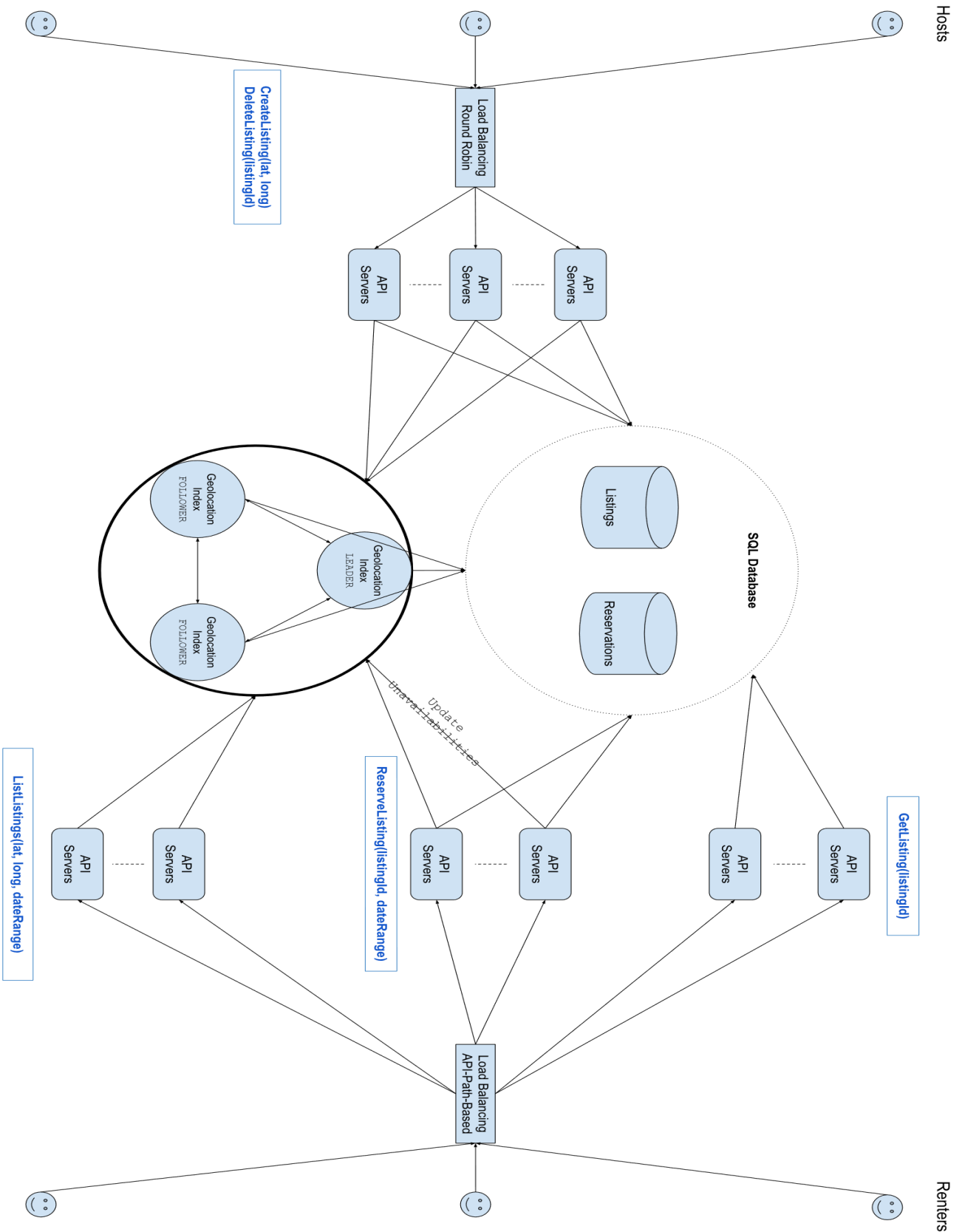
7. Load Balancing

On the host side, we can load balance requests to create and delete listings across a set of API servers using a simple round-robin approach. The API servers will then be in charge of writing to the SQL database and of communicating with the geo-index leader.

On the renter side, we can load balance requests to list, get, and reserve listings across a set of API servers using an API-path-based server-selection strategy. Since workloads for these three API calls will be considerably different from one another, it makes sense to separate these calls across different sets of API servers.

Of note is that we don't want any caching done at our API servers, because otherwise we'll naturally run into stale data as reservations, bookings, and new listings appear.

8. System Diagram



8. Design Netflix

Clarifying questions to ask

Question 1

Q: From a high-level point of view, Netflix is a fairly straightforward service: users go on the platform, they're served movies and shows, and they watch them. Are we designing this high-level system entirely, or would you like me to focus on a particular subsystem, like the Netflix home page?

A: We're just designing the core Netflix product--so the overarching system / product that you described.

Question 2

Q: Should we worry about auxiliary services like authentication and payments?

A: You can ignore those auxiliary services; focus on the primary user flow. That being said, one thing to note is that, by nature of the product, we're going to have access to a lot of user-activity data that's going to need to be processed in order to enable Netflix's recommendation system. You'll need to come up with a way to aggregate and process user-activity data on the website.

Question 3

Q: For this recommendation system, should I think about the kind of algorithm that'll fuel it?

A: No, you don't need to worry about implementing any algorithm or formula for the recommendation engine. You just need to think about how user-activity data will be gathered and processed.

Question 4

Q: It sounds like there are 2 main points of focus in this system: the video-serving service and the recommendation engine. Regarding the video-serving service, I'm assuming that we're looking for high availability and fast latencies globally; is this correct?

A: Yes, but just to clarify, the video-streaming service is actually the only part of the system for which we care about fast latencies.

Question 5

Q: So is the recommendation engine a system that consumes the user-activity data you mentioned and operates asynchronously in the background?

A: Yes.

Question 6

Q: How many users do we expect to be building this for?

A: Netflix has about 100M to 200M users, so let's go with 200M.

Question 7

Q: Should we worry about designing this for various clients, like desktop clients, mobile clients, etc.?

A: Even though we're indeed designing Netflix to be used by all sorts of clients, let's focus purely on the distributed-system component--so no need to get into details about clients or to optimize for certain clients.

Solution walkthrough

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

We're designing the core Netflix service, which allows users to stream movies and shows from the Netflix website.

Specifically, we'll want to focus on:

Delivering large amounts of high-definition video content to hundreds of millions of users around the globe without too much buffering.

Processing large amounts of user-activity data to support Netflix's recommendation engine.

2. Coming Up With A Plan

We'll tackle this system by dividing it into four main sections:

Storage (Video Content, Static Content, and User Metadata)

General Client-Server Interaction (i.e., the life of a query)

Video Content Delivery

User-Activity Data Processing

3. Video-Content Storage

Since Netflix's service, which caters to millions of customers, is centered around video content, we might need a lot of storage space and a complex storage solution. Let's start by estimating how much space we'll need.

We were told that Netflix has about 200 million users; we can make a few assumptions about other Netflix metrics (alternatively, we can ask our interviewer for guidance here):

Netflix offers roughly 10 thousand movies and shows at any given time

Since movies can be up to 2+ hours in length and shows tend to be between 20 and 40 minutes per episode, we can assume an average video length of 1 hour

Each movie / show will have a Standard Definition version and a High Definition version. Per hour, SD will take up about 10GB of space, while HD will take about 20GB.

~10K videos (stored in SD & HD)

~1 hour average video length

~10 GB/h for SD + ~20 GB/h for HD = 30 GB/h per video

~30 GB/h * 10K videos = 300,000 GB = 300 TB

This number highlights the importance of estimations. Naively, one might think that Netflix stores many petabytes of video, since its core product revolves around video content; but a simple back-of-the-napkin estimation shows us that it actually stores a very modest amount of video.

This is because Netflix, unlike other services like YouTube, Google Drive, and Facebook, has a bounded amount of video content: the movies and shows that it offers; those other services allow users to upload unlimited amounts of video.

Since we're only dealing with a few hundred terabytes of data, we can use a simple blob storage solution like S3 or GCS to reliably handle the storage and replication of Netflix's video content; we don't need a more complex data-storage solution.

4. Static-Content Storage

Apart from video content, we'll want to store various pieces of static content for Netflix's movies and shows, including video titles, descriptions, and cast lists.

This content will be bounded in size by the size of the video content, since it'll be tied to the number of movies and shows, just like the video content, and since it'll naturally take up less space than the video data.

We can easily store all of this static content in a relational database or even in a document store, and we can cache most of it in our API servers.

5. User Metadata Storage

We can expect to store some user metadata for each video on the Netflix platform. For instance, we might want to store the timestamp that a user left a video at, a user's rating on a video, etc..

Just like the static content mentioned above, this user metadata will be tied to the number of videos on Netflix. However, unlike the static content, this user metadata will grow with the Netflix userbase, since each user will have user metadata.

We can quickly estimate how much space we'll need for this user metadata:

~200M users

~1K videos watched per user per lifetime (~10% of total content)

~100 bytes/video/user

~100 bytes * 1K videos * 200M users = 100 KB * 200M = 1 GB * 20K = 20 TB

Perhaps surprisingly, we'll be storing an amount of user metadata in the same ballpark as the amount of video content that we'll be storing. Once again, this is because of the bounded nature of Netflix's video content, which is in stark contrast with the unbounded nature of its userbase.

We'll likely need to query this metadata, so storing it in a classic relational database like Postgres makes sense.

Since Netflix users are effectively isolated from one another (they aren't connected like they would be on a social-media platform, for example), we can expect all of our latency-sensitive database operations to only relate to individual users. In other words, potential operations like `GetUserInfo` and `GetUserWatchedVideos`, which would require fast latencies, are specific to a particular users; on the other hand, complicated database operations involving multiple users' metadata will likely be part of background data-engineering jobs that don't care about latency.

Given this, we can split our user-metadata database into a handful of shards, each managing anywhere between 1 and 10 TB of indexed data. This will maintain very quick reads and writes for a given user.

6. General Client-Server Interaction

The part of the system that handles serving user metadata and static content to users shouldn't be too complicated.

We can use some simple round-robin load balancing to distribute end-user network requests across our API servers, which can then load-balance database requests according to `userId` (since our database will be sharded based on `userId`).

As mentioned above, we can cache our static content in our API servers, periodically updating it when new movies and shows are released, and we can even cache user metadata there, using a write-through caching mechanism.

7. Video Content Delivery

We need to figure out how we'll be delivering Netflix's video content across the globe with little latency. To start, we'll estimate the maximum amount of bandwidth consumption that we could expect at any point in time. We'll assume that, at peak traffic, like when a popular movie comes out, a fairly large number of Netflix users might be streaming video content concurrently.

~200M total users

~5% of total users streaming concurrently during peak hours

~20 GB/h of HD video \approx 5 MB/s of HD video

$\sim 5\% \text{ of } 200\text{M} * 5 \text{ MB/s} = 10\text{M} * 5 \text{ MB/s} = 50 \text{ TB/s}$

This level of bandwidth consumption means we can't just naively serve the video content out of a single data center or even dozens of data centers. We need many thousands of locations around the world to be distributing this content for us. Thankfully, CDNs solve this precise problem, since they have many thousands of

Points of Presence around the world. We can thus use a CDN like Cloudflare and serve our video content out of the CDN's PoPs.

Since the PoPs can't keep the entirety of Netflix's video content in cache, we can have an external service that periodically repopulates CDN PoPs with the most important content (the movies and shows most likely to be watched).

8. User-Activity Data Processing

We need to figure out how we'll process vast amounts of user-activity data to feed into Netflix's recommendation engine. We can imagine that this user-activity data will be gathered in the form of logs that are generated by all sorts of user actions; we can expect terabytes of these logs to be generated every day.

MapReduce can help us here. We can store the logs in a distributed file system like HDFS and run MapReduce jobs to process massive amounts of data in parallel. The results of these jobs can then be fed into some machine learning pipelines or simply stored in a database.

Map Inputs

Our Map inputs can be our raw logs, which might look like:

```
{"userId": "userId1", "videoId": "videoId1", "event": "CLICK"}  
{"userId": "userId2", "videoId": "videoId2", "event": "PAUSE"}  
{"userId": "userId3", "videoId": "videoId3", "event": "MOUSE_MOVE"}
```

Map Outputs / Reduce Inputs

Our Map function will aggregate logs based on `userId` and return intermediary key-value pairs indexed on each `userId`, pointing to lists of tuples with `videoId`s and relevant events.

These intermediary k/v pairs will be shuffled appropriately and fed into our Reduce functions.

```
{"userId1": [("CLICK", "videoId1"), ("CLICK", "videoId1"), ..., ("PAUSE", "videoId2")]}  
{"userId2": [("PLAY", "videoId1"), ("MOUSE_MOVE", "videoId2"), ..., ("MINIMIZE",  
"videoId3")]}
```

Reduce Outputs

Our Reduce functions could return many different outputs. They could return k/v pairs for each `userId`/`videoId` combination, pointing to a computed score for that user/video pair; they could return k/v pairs indexed at each `userId`, pointing to lists of (`videoId`,

score) tuples; or they could return k/v pairs also indexed at eacher userId but pointing to stack-rankings of videolds, based on their computed score.

("userId1|videold1", score)

("userId1|videold2", score)

OR

{"userId1": [("videold1", score), ("videold2", score), ..., ("videold3", score)]}

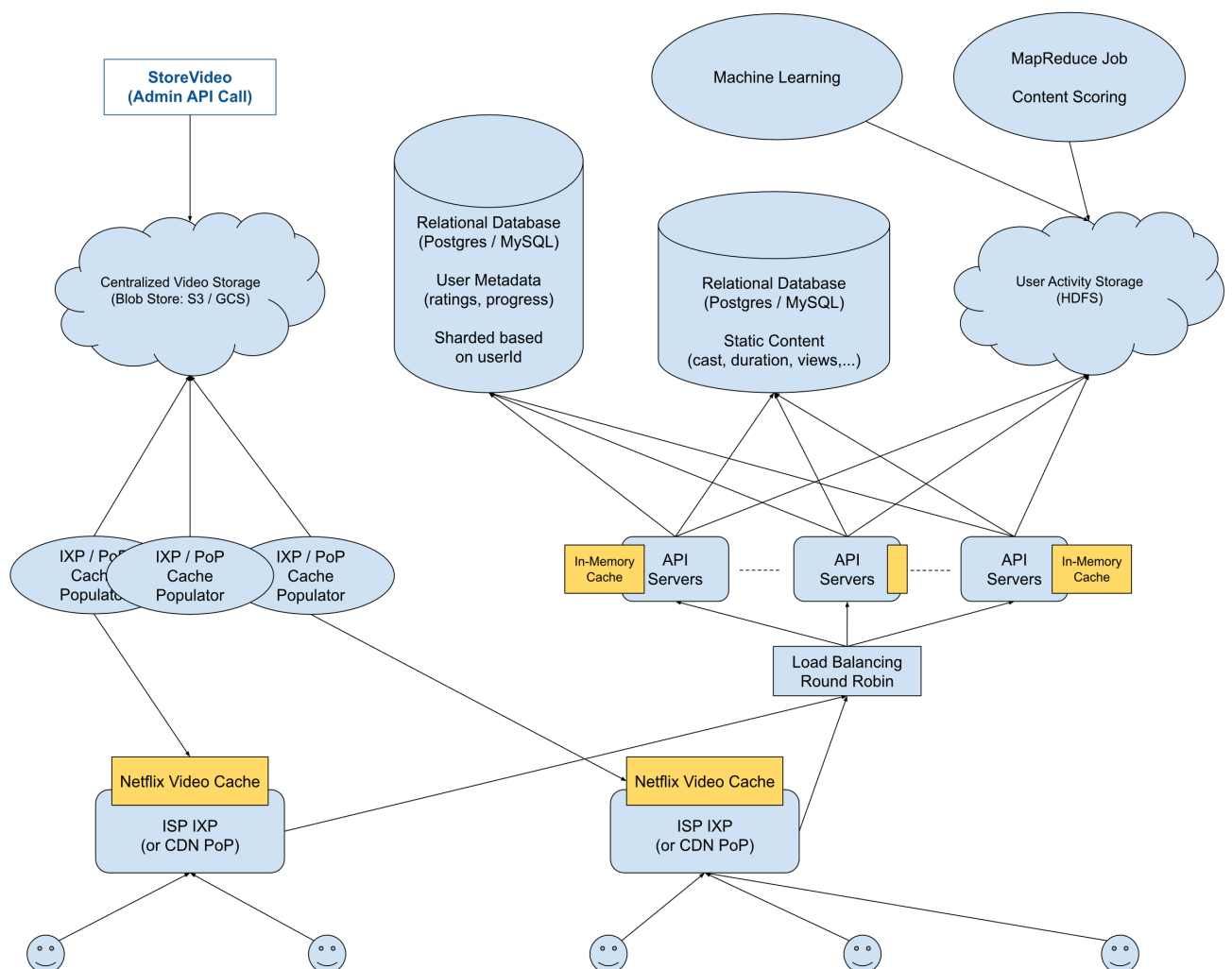
{"userId2": [("videold1", score), ("videold2", score), ..., ("videold3", score)]}

OR

("userId1", ["videold1", "videold2", ..., "videold3"])

("userId2", ["videold1", "videold2", ..., "videold3"])

9. System Diagram



8. Design Facebook News Feed

Clarifying questions to ask

Question 1

Q: Facebook News Feed consists of multiple major features, like loading a user's news feed, interacting with it (i.e., posting status updates, liking posts, etc.), and updating it in real time (i.e., adding new status updates that are being posted to the top of the feed, in real time). What part of Facebook News Feed are we designing exactly?

A: We're designing the core functionality of the feed itself, which we'll define as follows: loading a user's news feed and updating it in real time, as well as posting status updates. But for posting status updates, we don't need to worry about the actual API or the type of information that a user can post; we just want to design what happens once an API call to post a status update has been made. Ultimately, we primarily want to design the feed generation/refreshing piece of the data pipeline (i.e, how/when does it get constructed, and how/when does it get updated with new posts).

Question 2

Q: To clarify, posts on Facebook can be pretty complicated, with pictures, videos, special types of status updates, etc.. Are you saying that we're not concerned with this aspect of the system? For example, should we not focus on how we'll be storing this type of information?

A: That's correct. For the purpose of this question, we can treat posts as opaque entities that we'll certainly want to store, but without worrying about the details of the storage, the ramifications of storing and serving large files like videos, etc..

Question 3

Q: Are we designing the relevant-post curation system (i.e., the system that decides what posts will show up on a user's news feed)?

A: No. We're not designing this system or any ranking algorithms; you can assume that you have access to a ranking algorithm that you can simply feed a list of relevant posts to in order to generate an actual news feed to display.

Question 4

Q: Are we concerned with showing ads in a user's news feed at all? Ads seem like they would behave a little bit differently than posts, since they probably rely on a different ranking algorithm.

A: You can treat ads as a bonus part of the design; if you find a way to incorporate them in, great (and yes, you'd have some other ads-serving algorithm to determine what ads need to be shown to a user at any point in time). But don't focus on ads to start.

Question 5

Q: Are we serving a global audience, and how big is our audience?

A: Yes -- we're serving a global audience, and let's say that the news feed will be loaded in the order of 100 million times a day, by 100 million different users, with 1 million new status updates posted every day.

Question 6

Q: How many friends does a user have on average? This is important to know, since a user's status updates could theoretically have to show up on all of the user's friends' news feeds at once.

A: You can expect each user to have, on average, 500 friends on the social network. You can treat the number of friends per user as a bell-shaped distribution, with some users who have very few friends, and some users who have a lot more than 500 friends.

Question 7

Q: How quickly does a status update have to appear on a news feed once it's posted, and is it okay if this varies depending on user locations with respect to the location of the user submitting a post?

A: When a user posts something, you probably want it to show up on other news feeds fairly quickly. This speed can indeed vary depending on user locations. For instance, we'd probably want a local friend within the same region to see the new post within a few seconds, but we'd likely be okay with a user on the other side of the world seeing the same post within a minute.

Question 8

Q: What kind of availability are we aiming for?

A: Your design shouldn't be completely unavailable from a single machine failure, but this isn't a high availability requirement. However, posts shouldn't ever just disappear. Once the user's client gets confirmation that the post was created, you cannot lose it.

Solution walkthrough

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

We're designing the core user flow of the Facebook News Feed. This consists of loading a user's news feed, scrolling through the list of posts that are relevant to them, posting status updates, and having their friends' news feeds get updated in real time. We're Specifically designing the pipeline that generates and serves news feeds and the system that handles what happens when a user posts and news feeds have to be updated.

We're dealing with about 1 billion users, each with 500 friends on average.

Getting a news feed should feel fairly instant, and creating a post should update all of a user's friends' news feeds within a minute. We can have some variance with regards to feed updates depending on user locations.

Additionally, we can't be satisfied with a single cluster serving everyone on earth because of large latencies that would occur between that cluster and the user in some parts of the world, so we need a mechanism to make sure the feed gets updated within a minute in the regions other than the one the post was created in.

We can assume that the ranking algorithms used to generate news feeds with the most relevant posts is taken care of for us by some other system that we have access to.

2. Coming Up With A Plan

We'll start with the extremities of our system and work inward, first talking about the two API calls, CreatePost and GetNewsFeed, then, getting into the feed creation and storage strategy, our cross-region design, and finally tying everything together in a fast and scalable way.

3. CreatePost API

For the purpose of this design, the CreatePost API call will be very simple and look something like this:

```
CreatePost(  
    user_id: string,  
    post: data  
)
```

When a user creates a post, the API call goes through some load balancing before landing on one of many API servers (which are stateless). Those API servers then create a message on a Pub/Sub topic, notifying its subscribers of the new post that was just created. Those subscribers will do a few things, so let's call them S1 for future reference. Each of the subscribers S1 reads from the topic and is responsible for creating the facebook post inside a relational database.

4. Post Storage

We can have one main relational database to store most of our system's data, including posts and users. This database will have very large tables.

5. GetNewsFeed API

The GetNewsFeed API call will most likely look like this:

```
GetNewsFeed(  
    user_id: string,  
    pageSize: integer,  
    nextPageToken: integer,  
) => (  
    posts: [{  
        user_id: string,  
        post_id: string,  
        post: data,  
    }],  
    nextPageToken: string,  
)
```

The `pageSize` and `nextPageToken` fields are used to paginate the newsfeed; pagination is necessary when dealing with large amounts of listed data, and since we'll likely want each news feed to have up to 1000 posts, pagination is very appropriate here.

6. Feed Creation And Storage

Since our databases tables are going to be so large, with billions of millions of users and tens of millions of posts every week, fetching news feeds from our main database every time a `GetNewsFeed` call is made isn't going to be ideal. We can't expect low latencies when building news feeds from scratch because querying our huge tables takes time, and sharding the main database holding the posts wouldn't be particularly helpful since news feeds would likely need to aggregate posts across shards, which would require us to perform cross-shard joins when generating news feeds; we want to avoid this.

Instead, we can store news feeds separately from our main database across an array of shards. We can have a separate cluster of machines that can act as a proxy to the relational database and be in charge of aggregating posts, ranking them via the ranking algorithm that we're given, generating news feeds, and sending them to our shards every so often (every 5, 10, 60 minutes, depending on how often we want news feeds to be updated).

If we average each post at 10kB, and a newsfeed comprises of the top 1000 posts that are relevant to a user, that's 10MB per user, or 10 000TB of data total. We assume that it's loaded 10 times per day per user, which averages at 10k QPS for the newsfeed fetching.

Assuming 1 billion news feeds (for 1 billion users) containing 1000 posts of up to 10 KB each, we can estimate that we'll need 10 PB (petabytes) of storage to store all of our users' news feeds. We can use 1000 machines of 10 TB each as our news-feed shards.

~10 KB per post

~1000 posts per news feed

~1 billion news feeds

$\sim 10 \text{ KB} * 1000 * 1000^3 = 10 \text{ PB} = 1000 * 10 \text{ TB}$

To distribute the newsfeeds roughly evenly, we can shard based on the user id.

When a GetNewsFeed request comes in, it gets load balanced to the right news feed machine, which returns it by reading on local disk. If the newsfeed doesn't exist locally, we then go to the source of truth (the main database, but going through the proxy ranking service) to gather the relevant posts. This will lead to increased latency but shouldn't happen frequently.

7. Wiring Updates Into Feed Creation

We now need to have a notification mechanism that lets the feed shards know that a new relevant post was just created and that they should incorporate it into the feeds of impacted users.

We can once again use a Pub/Sub service for this. Each one of the shards will subscribe to its own topic--we'll call these topics the Feed Notification Topics (FNT)--and the original subscribers S1 will be the publishers for the FNT. When S1 gets a new message about a post creation, it searches the main database for all of the users for whom this post is relevant (i.e., it searches for all of the friends of the user who created the post), it filters out users from other regions who will be taken care of asynchronously, and it maps the remaining users to the FNT using the same hashing function that our GetNewsFeed load balancers rely on.

For posts that impact too many people, we can cap the number of FNT topics that get messaged to reduce the amount of internal traffic that gets generated from a single

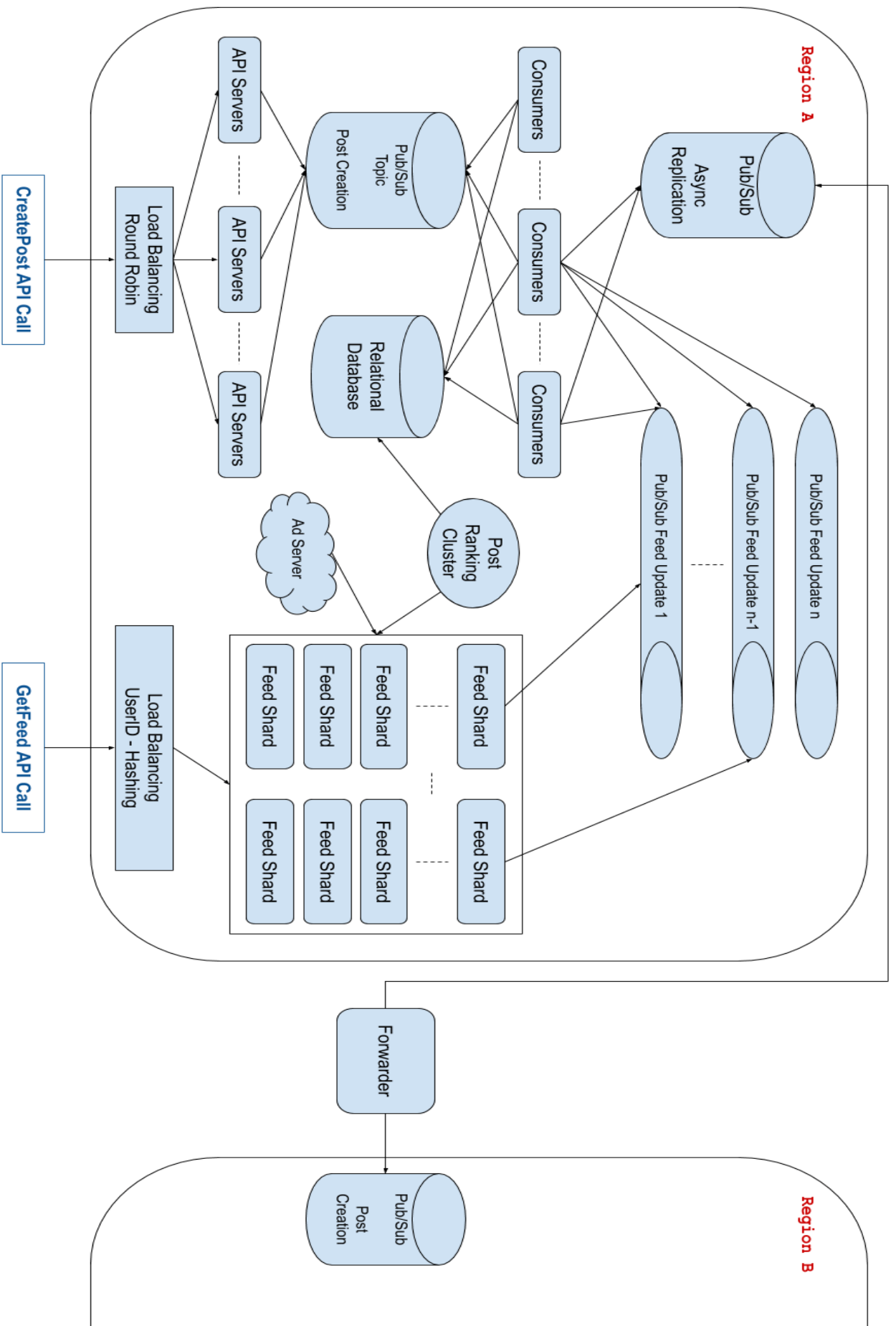
post. For those big users we can rely on the asynchronous feed creation to eventually kick in and let the post appear in feeds of users whom we've skipped when the feeds get refreshed manually.

8. Cross-Region Strategy

When CreatePost gets called and reaches our Pub/Sub subscribers, they'll send a message to another Pub/Sub topic that some forwarder service in between regions will subscribe to. The forwarder's job will be, as its name implies, to forward messages to other regions so as to replicate all of the CreatePost logic in other regions. Once the forwarder receives the message, it'll essentially mimic what would happen if that same CreatePost were called in another region, which will start the entire feed-update logic in those other regions. We can have some additional logic passed to the forwarder to prevent other regions being replicated to from notifying other regions about the CreatePost call in question, which would lead to an infinite chain of replications; in other words, we can make it such that only the region where the post originated from is in charge of notifying other regions.

Several open-source technologies from big companies like Uber and Confluent are designed in part for this kind of operation.

9. System Diagram



10. Design Slack

Clarifying questions to ask

Question 1

Q: There are a lot of things that you can do on Slack. Primarily, you use Slack to communicate with people in one-on-one channels, private channels, or public channels, all within an organization. But you can also do a bunch of other things on Slack, like create and delete channels, change channel settings, change Slack settings, invite people to channels, etc.. What exactly are we designing here?

A: We're designing the core messaging functionality, which involves communicating in both one-on-one channels and group channels in an organization. You don't have to worry about channel settings and all of those extra functionalities.

Question 2

Q: Okay. Do you want me to take care of the concept of private channels at all?

A: Let's just focus on users in a channel as far as access control is concerned; we can forget about the concept of a private channel.

Question 3

Q: Okay. And regarding communication, from my knowledge of Slack, when you load the web app or the desktop / mobile apps, you can obviously access all the messages of channels that you're in (including one-on-one channels), but you're also notified of channels that have unread messages for you and of the number of unread mentions that you have in each channel. Channels with unread messages are bold, if I remember correctly, and the number of unread mentions is simply visible next to channel names. Should we design our system to accommodate this?

A: Yes, we should take care of this. And on that note, one thing we'll want to handle is cross-device synchronization. In other words, if you have both the Slack desktop app and the Slack mobile app open, and both apps are showing that one channel is unread, and you read that channel on one of the apps, the other app should immediately get updated and should mark the channel as read. You'll have to handle this.

Question 4

Q: Hmm, okay. Speaking of different applications, by the way, are we designing the various device / software apps, or just the backend systems that the frontends / clients communicate with?

A: You'll only really focus on the backend systems for this question.

Question 5

Q: Okay. Also, there are a lot of different features in actual Slack messages. For example, adding custom emojis, pinning messages, saving messages, writing code snippets or text-blocks, etc.. Do you want me to handle all of this?

A: No, you can just treat messages as pure text for now. Of course, what you'll design will likely be extensible to different types of messages and will eventually be able to handle things like pinning or saving messages, but for this design, don't worry about that.

Question 6

Q: How many users do we expect to be building this for? And how large is the largest organization on slack? How many users does it have?

A: Slack has about 10 to 20 million users, so let's go with 20 million. And as for organizations, let's say that the largest single Slack customer has 50,000 people in the same organization. We can also approximate that the largest channel will be of that same size if all of an organization's employees are in the same channel (the typical #general channel, for example).

Question 7

Q: Since this is a chat application, I'm assuming that low latency is one of our top priorities, and also, since this service impacts millions of users, I'm assuming that we should design with high availability in mind. Are these correct assumptions?

A: Yes to both of those things, but for the sake of being a little more focused, don't worry about optimizing for availability. Let's focus primarily on latency and core functionality.

Question 8

Q: Okay. And are we building this for a global audience, or should we focus on a single region?

A: Let's handle a single region for this question, but just like with availability, don't focus too much on this aspect of the design.

Solution Walkthrough

1. Gathering System Requirements

As with any systems design interview question, the first thing that we want to do is to gather system requirements; we need to figure out what system we're building exactly.

We're designing the core communication system behind Slack, which allows users to send instant messages in Slack channels.

Specifically, we'll want to support:

Loading the most recent messages in a Slack channel when a user clicks on the channel.

Immediately seeing which channels have unread messages for a particular user when that user loads Slack.

Immediately seeing which channels have unread mentions of a particular user, for that particular user, when that user loads Slack, and more specifically, the number of these unread mentions in each relevant channel.

Sending and receiving Slack messages instantly, in real time.

Cross-device synchronization: if a user has both the Slack desktop app and the Slack mobile app open, with an unread channel in both, and if they read this channel on one device, the second device should immediately be updated and no longer display the channel as unread.

The system should have low latencies and high availability, catering to a single region of roughly 20 million users. The largest Slack organizations will have as many as 50,000 users, with channels of the same size within them.

That being said, for the purpose of this design, we should primarily focus on latency and core functionality; availability and regionality can be disregarded, within reason.

2. Coming Up With A Plan

We'll tackle this system by dividing it into two main sections:

Handling what happens when a Slack app loads.

Handling real-time messaging as well as cross-device synchronization.

We can further divide the first section as follows:

Seeing all of the channels that a user is a part of.

Seeing messages in a particular channel.

Seeing which channels have unread messages.

Seeing which channels have unread mentions and how many they have.

3. Persistent Storage Solution & App Load

While a large component of our design involves real-time communication, another large part of it involves retrieving data (channels, messages, etc.) at any given time when the Slack app loads. To support this, we'll need a persistent storage solution.

Specifically, we'll opt for a SQL database since we can expect this data to be structured and to be queried frequently.

We can start with a simple table that'll store every Slack channel.

Channels

id (channelid): <i>uuid</i>	orgld: <i>uuid</i>	name: <i>string</i>	description: <i>string</i>
...

Then, we can have another simple table representing channel-member pairs: each row in this table will correspond to a particular user who is in a particular channel. We'll use this table, along with the one above, to fetch a user's relevant when the app loads.

Channel Members

id: <i>uuid</i>	orgId: <i>uuid</i>	channelId: <i>uuid</i>	userId: <i>uuid</i>
...

We'll naturally need a table to store all historical messages sent on Slack. This will be our largest table, and it'll be queried every time a user fetches messages in a particular channel. The API endpoint that'll interact with this table will return a paginated response, since we'll typically only want the 50 or 100 most recent messages per channel.

Also, this table will only be queried when a user clicks on a channel; we don't want to fetch messages for all of a user's channels on app load, since users will likely never look at most of their channels.

Historical Messages

id: <i>uuid</i>	orgId: <i>uuid</i>	channelId: <i>uuid</i>	senderId: <i>uuid</i>	sentAt: <i>timestamp</i>	body: <i>string</i>	mentions: <i>List<uuid></i>
...

In order not to fetch recent messages for every channel on app load, all the while supporting the feature of showing which channels have unread messages, we'll need to store two extra tables: one for the latest activity in each channel (this table will be updated whenever a user sends a message in a channel), and one for the last time a particular user has read a channel (this table will be updated whenever a user opens a channel).

Latest Channel Timestamps

id: <i>uuid</i>	orgId: <i>uuid</i>	channelId: <i>uuid</i>	lastActive: <i>timestamp</i>
...

Channel Read Receipts

id: <i>uuid</i>	orgId: <i>uuid</i>	channelId: <i>uuid</i>	userId: <i>uuid</i>	lastSeen: <i>timestamp</i>
...

For the number of unread user mentions that we want to display next to channel names, we'll have another table similar to the read-receipts one, except this one will have a count of unread user mentions instead of a timestamp. This count will be updated (incremented) whenever a user tags another user in a channel message, and it'll also be updated (reset to 0) whenever a user opens a channel with unread mentions of themselves.

Unread Channel-User-Mention Counts

id: <i>uuid</i>	orgId: <i>uuid</i>	channelId: <i>uuid</i>	userId: <i>uuid</i>	count: <i>int</i>
...

4. Load Balancing

For all of the API calls that clients will issue on app load, including writes to our database (when sending a message or marking a channel as read), we're going to want to load balance.

We can have a simple round-robin load balancer, forwarding requests to a set of server clusters that will then handle passing requests to our database.

5. "Smart" Sharding

Since our tables will be very large, especially the messages table, we'll need to have some sharding in place.

The natural approach is to shard based on organization size: we can have the biggest organizations (with the biggest channels) in their individual shards, and we can have smaller organizations grouped together in other shards.

An important point to note here is that, over time, organization sizes and Slack activity within organizations will change. Some organizations might double in size overnight, others might experience seemingly random surges of activity, etc.. This means that, despite our relatively sound sharding strategy, we might still run into hot spots, which is very bad considering the fact that we care about latency so much.

To handle this, we can add a "smart" sharding solution: a subsystem of our system that'll asynchronously measure organization activity and "rebalance" shards accordingly. This service can be a strongly consistent key-value store like Etcd or ZooKeeper, mapping orgIds to shards. Our API servers will communicate with this service to know which shard to route requests to.

6. Pub/Sub System for Real-Time Behavior

There are two types of real-time behavior that we want to support:

Sending and receiving messages in real time.

Cross-device synchronization (instantly marking a channel as read if you have Slack open on two devices and read the channel on one of them).

For both of these functionalities, we can rely on a Pub/Sub messaging system, which itself will rely on our previously described "smart" sharding strategy.

Every Slack organization or group of organizations will be assigned to a Kafka topic, and whenever a user sends a message in a channel or marks a channel as read, our previously mentioned API servers, which handle speaking to our database, will also send a Pub/Sub message to the appropriate Kafka topic.

The Pub/Sub messages will look like:

```
{
  "type": "chat",
  "orgId": "AAA",
  "channelId": "BBB",
  "userId": "CCC",
  "messageId": "DDD",
  "timestamp": "2020-08-31T01:17:02",
  "body": "this is a message",
  "mentions": ["CCC", "EEE"]
},
{
  "type": "read-receipt",
  "orgId": "AAA",
  "channelId": "BBB",
  "userId": "CCC",
  "timestamp": "2020-08-31T01:17:02"
}
```

We'll then have a different set of API servers who subscribe to the various Kafka topics (probably one API server cluster per topic), and our clients (Slack users) will establish long-lived TCP connections with these API server clusters to receive Pub/Sub messages in real time.

We'll want a load balancer in between the clients and these API servers, which will also use the "smart" sharding strategy to match clients with the appropriate API servers, which will be listening to the appropriate Kafka topics.

When clients receive Pub/Sub messages, they'll handle them accordingly (mark a channel as unread, for example), and if the clients refresh their browser or their mobile app, they'll go through the entire "on app load" system that we described earlier.

Since each Pub/Sub message comes with a timestamp, and since reading a channel and sending Slack messages involve writing to our persistent storage, the Pub/Sub messages will effectively be idempotent operations.

7. System Diagram

