

Design UBER / OLA App which is a ride sharing app

Solution:- ① Find out the Functional Requirements

(a) User :-

- (i) Book a cab
- (ii) ETA of the cab arrival
- (iii) Cancel the cab
- (iv) Details of the ongoing trip.

(b) Driver :-

- (i) Notification for the new booking
- (ii) Accept / reject a new booking
- (iii) Update my location very frequently
- (iv) Details of the ongoing trip.

② Find out the Non-Functional Requirements

- (i) It must be real-time, no stale data.
- (ii) It should be scalable, in peak hours like evening hours, noon hours your app should work smoothly.
- (iii) It should be reliable.
- (iv) You keep your system highly available making your system eventually consistent.

③ Get Estimation / Back of the Envelope Calculation

① How much compute and storage and what will be network bandwidth, you must have numbers to make the infrastructure ready.

Let's assume:-

$$\text{Total customers / riders} = 500 \text{ million}$$

$$\text{Total number of drivers} = 5 \text{ million}$$

$$\text{Daily active riders} = 20 \text{ million}$$

Daily active drivers = 3 million (since most of the drivers are ready to give service)

Total storage needed for all the drivers if

Rider's data (i.e id, name, vehicle registration number) = 500 Million $\times 1 \text{ KB}$

$$\begin{aligned} \text{is } 1 \text{ KB} \\ &= 500 \times 10^6 \times 1 \text{ KB} \\ &= 500 \text{ GB} \end{aligned}$$

every day, users / riders are increasing at the rate of 50,000 / days

$\hookrightarrow 50 \text{ MB/day}$ is the storage requirement for the next $\frac{1}{2}$ year

$$\text{Current} \rightarrow 500 \text{ GB} \text{ (Storage Cap.)}$$

$$\begin{aligned} \text{Next 365 days} &= 365 \times 100/2 \text{ MB} \\ &= 182.5 \times 100 \text{ MB} \\ &= 18.2 \text{ GB} \end{aligned}$$

In $\frac{1}{2}$ year, we require $18.2 \text{ GB} \approx 520 \text{ GB}$ so as to make this application work smoothly for an year.

for an year, let drivers (new) come in the app everyday, let it be 5000

$$\begin{aligned} \text{In an year total new drivers} &= 365 \times 5000 \\ &= 365 \times \frac{10000}{2} \end{aligned}$$

$$= 182.5 \times 10000$$

$$= 1.825 \times 10^6$$

Now, Each driver meta data let say take 1 KB of space

$$\text{Then extra space required in next 365 days} = 1.825 \times 10^6 \times 10^3$$

$$\begin{aligned}
 &= 18.25 \text{ MB} \times 100 \\
 &= 1.825 \times 10^9 \text{ Bytes} \\
 &= 1.825 \text{ GB}
 \end{aligned}$$

$$\begin{aligned}
 \text{Total storage for drivers} &= 3 \text{ GB} + 1.825 \text{ GB} \\
 &= 4.825 \text{ GB}
 \end{aligned}$$

Trip data storage:-

since, we have daily active 20 million drivers,
so there will be 20 million trips per day

$$\text{Let, } \pm \text{Trip data} = 100 \text{ Bytes}$$

$$\begin{aligned}
 \text{1 day storage} &= 20 \times 10^6 \times 100 \text{ Bytes} \\
 &= 2 \times 10^9 \text{ Bytes} \\
 &= 2 \text{ GB}
 \end{aligned}$$

On 365 days,

$$\begin{aligned}
 \text{Storage required} &= 365 \times 2 \text{ GB} \\
 &= 730 \text{ GB}
 \end{aligned}$$

Noted Points:-

① Everytime a trip happens, it results in queries/APC calls
At what rate all these things are happening, how many transactions are happening, how many queries are happening, in order to understand the scale at which it is operating?

$$\begin{aligned}
 \text{In 24 hrs} &= 20 \text{ million} \\
 &= 20 \times 10^6
 \end{aligned}$$

$$\begin{aligned}
 \downarrow \text{second} &= \frac{20 \times 10^6}{24 \times 60 \times 60} = 230 \text{ trips/seconds}
 \end{aligned}$$

$$\begin{aligned}
 \text{Traffic Bandwidth of Trip} &= \text{Trips/seconds} \times \text{size of trip} \\
 / \text{Throughput} &= 230 \times 100 \text{ Bytes/second} \\
 &= 23 \text{ KB/seconds}
 \end{aligned}$$

② Driver needs to update its location very frequently:-
Let, driver location data take 36 Bytes

$$\begin{aligned}
 \text{Now, daily active drivers} &= 20 \text{ million} \\
 \text{So, Data to store location} &= 36 \times 20 \text{ million Bytes} \\
 &= 36 \times 20 \times 10^6 \text{ Bytes}
 \end{aligned}$$

$$\text{Throughput} = 720/4 \text{ location/second} = 180 \text{ MB/second}$$

How many servers are required?
Let, server can take 8000 requests in a day.

$$\text{Then, Number of Servers Required} = \frac{20 \times 10^6}{8000}$$

$$= \frac{20 \times 10^3}{8}$$

$$= \frac{1 \times 10^4}{4}$$

$$= 2500 \text{ SERVERS}$$

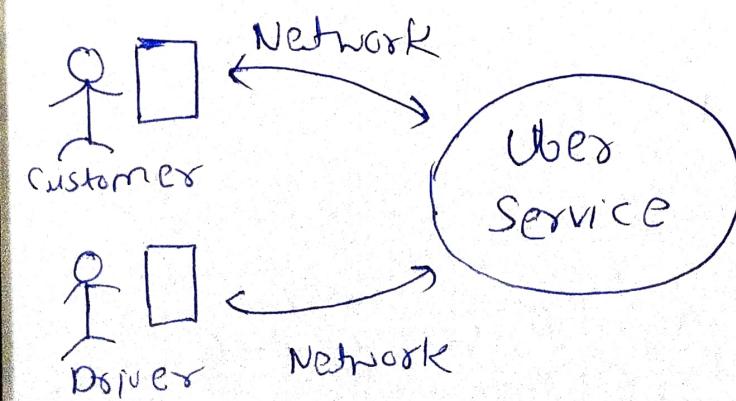
Why are companies trying to decrease their carbon footprint? What does it mean?
When you book a cab, it takes 5ml of water, so whenever you are booking a cab, you are depleting the natural environment resources.

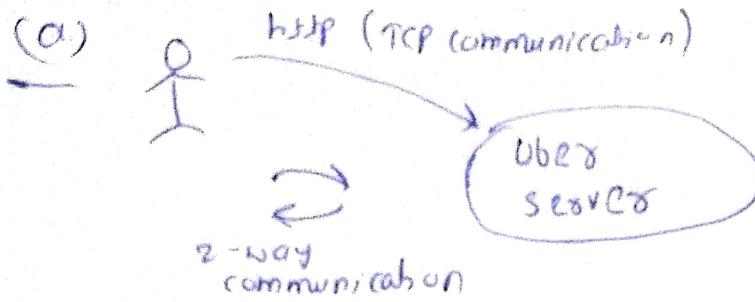
This is because water is used in machines in order to ensure machines run at the ambient temperature so that it can run all the time. Resources used are electricity and water, sometime for electricity petrol/diesel are used.

Because natural resources like Microsoft/Google/LinkedIn/Amazon etc. are focussing to minimize carbon footprint.

Note:- A prompt to chatbot also take 15-20 ml of water to keep servers machine temperature stable

③ How an overall system looks like:-





Approach (a) is resource intensive as http is short-lived.

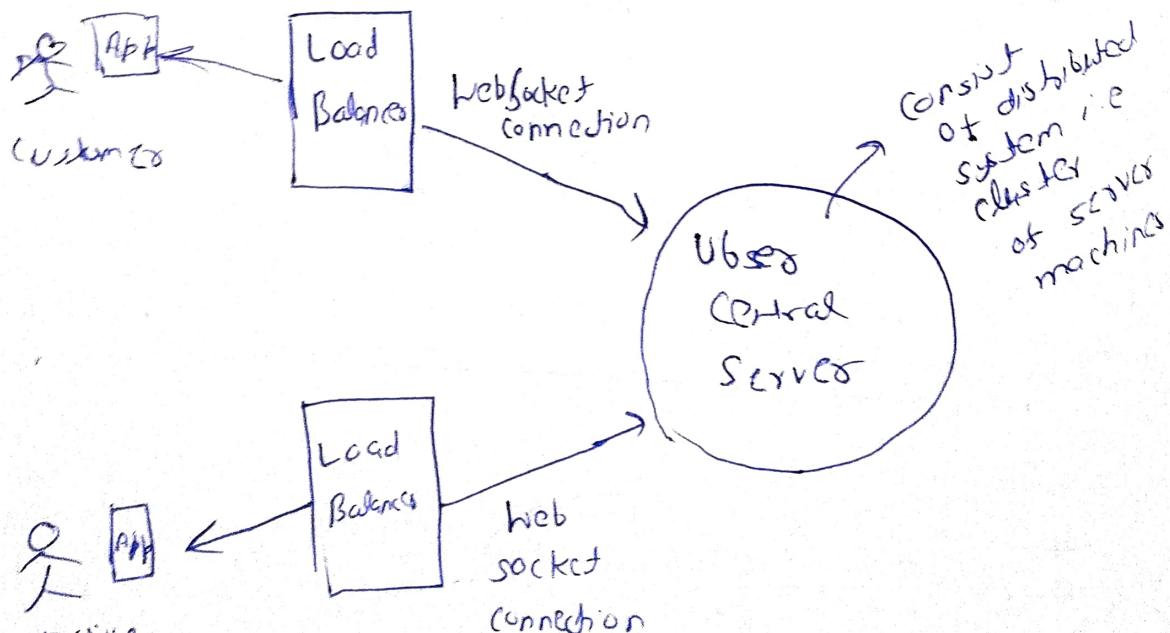
(b) So, to simplify approach (a), we can use socket I/O as it provides bidirectional communication but it is long-lived.

So, n/w for customer - central Uber system as well as for Drivers - central Uber system will be the web-socket connection.

Note:- 2500 Central

Uber system use these as we need 2500 servers.

(c) Since, there are so many servers and to have to balance the load, hence we need load balancing between customers & distributed system of Uber services, as well as with between drivers and distributed system.



(d) Let's design Uber Central system

uber central system consist of -

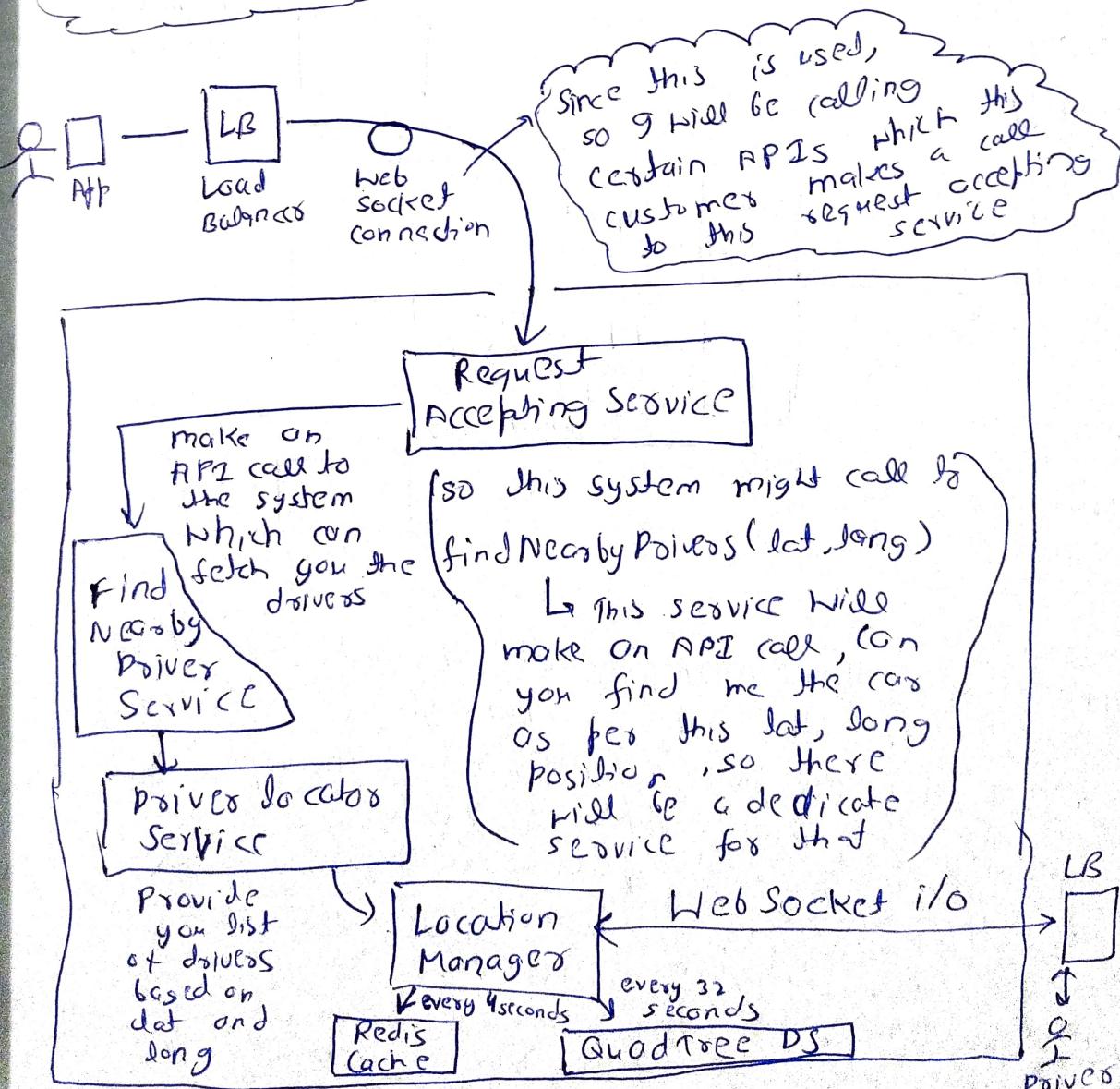
- (i) Driver location
- (ii) Trip Management
- (iii) Location Updation
- (iv) Drivers Matching
- (v) ETA calculation

} → If you add all these services in a single machine, then this machine is typically doing multiple tasks.

Hence, all these services are the specialized task, these are the individual task which should be ideally kept separately, so monolithic architecture is not a good approach, so we need to make it more granular.

↓ This architecture is called Monolithic architecture

→ So, better approach can be Service Oriented Architecture



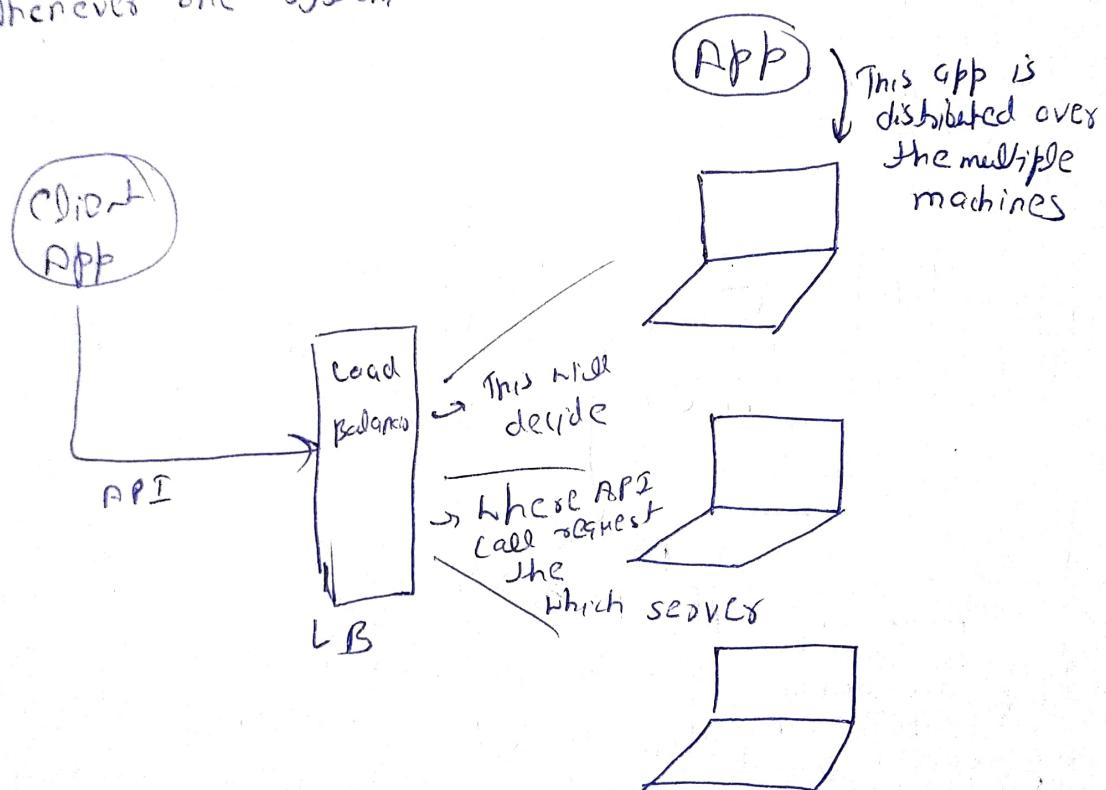
What are the APIs which the customers might be interested in calling? Some of the APIs called by customers:

- ① Request Ride (rideId, lat, long, destLat, destLong,
choose lat and longitude)
Type of Vehicle)

Why not chosen Microservices instead of SOA, because currently I am not sure about that every system is fully independent, so if you are not sure, you must start with SOA not directly Microservices

What is API

Whenever one system talks with another systems over a



In network, they normally do through the API itself, in the API when you are trying to hit the server / machine, in API you do not hard code the IP Address into it, you only go and provide the end points not IP address, so even if you are making a call to the API, you have to go through Load Balancers. In short, you say you have something called as API gateways which is similar to that we done in load balancing

How does location service works?
Let's say changes its location
very fast i.e. it updates its location every 4
seconds.

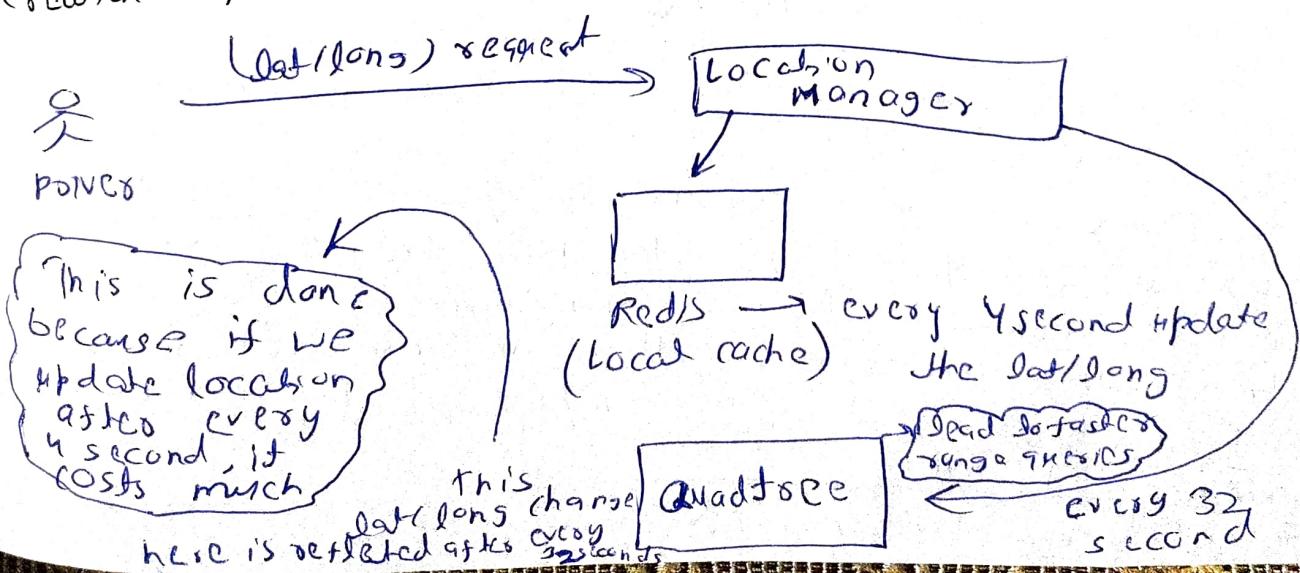
Now, every driver has some kind of lat & long
position which is updating. These basically we
are mapping the things on the different type of
locations, datastructure that we use here is
Quadtree Datastructure.

In Zomato, when you will
review the different restaurants, each restaurant
will have the location \rightarrow latitude / longitude. And
if you have to do it on the global scale, the
data structure which very efficiently stores the
locations of each of the entities over the
hemispheres is Quadtree Datastructure.

How drivers update their location

Let driver updates its
location after 4 seconds,
whenever driver changes
its location it makes a
call to the Location Manager
through socket i/o connection.

Location Manager will update the quadtree
datastructure. Now, whenever if you try to
update it is very slow as update leads to the
creation of the additional nodes in the Quadtree.



This is done
because if we
update location
after every
4 second, it
costs much

Lat/long change
here is reflected after every 32 seconds

Quadtree
Quadtree
range queries

every 32
second

Note:- Quadtree Data Structure is a persistent storage. Quadtree is updated by location manager as after 32 second, location manager pull data from Redis and put it on Quadtree, location manager updates Redis after every 4 second so sync between Redis and Quadtree is managed by the location manager.

Database used?

- (i) since I am dealing with huge data and it is scalable
- (ii) It provides you high read/write I/O operation
- (iii) So, database can be used here "Cassandra"
- (iv) In earlier times, users used Cassandra database but they replaced it by **Spanner** database. Spanner is a Google property database which also provide you high throughput in case of read and write operation.

Note:- If your all database lies in the single machine then you can use RDBMS preferably but if data cannot be stored on single machine RDBMS is not preferred.

Driver might use updateDriverLocation (DriverId, oldLat, oldLong, newLat, newLong) as an API call to location manager

Why can we not directly fetch data from Redis as Redis is fast?

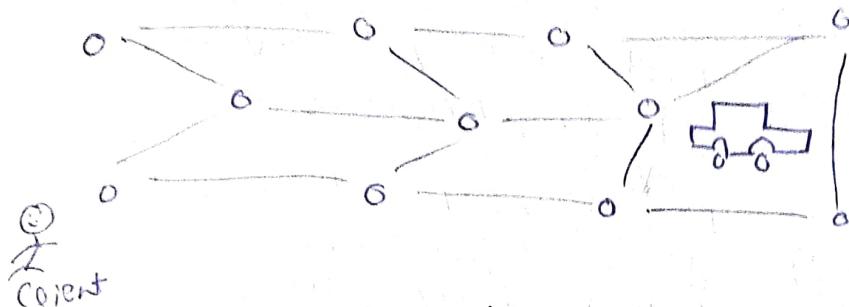
In Redis, we use key, value pairs. So, if you query over a key, you will only get one specific value. Now, if you try to do the range queries then you have to scan all the keys, it is complex as it takes a lot of time.

So, driver is doing write operation

When you book a cab
Uber provide you ETA
of where the cab arrives? How does
Uber do it?

Page No - 91

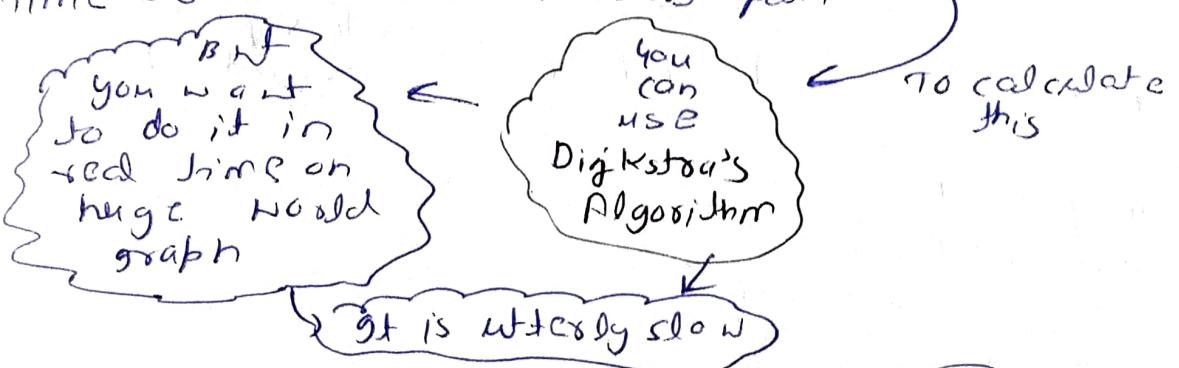
Any cab drivers make use of Google Maps. Google
Maps is nothing, it is just a complex graph.



orient

ETA arrival calculation?

- ① Find the shortest path
- ② Time to cover that shortest path



Note:- In order to speed up the shortest path stuff,
you always do the pre-calculations.
You always break the entire graph into the
smaller partitions and in the background, for
each two combination of the source &
the destination you always do the pre-
computation.

But this whatever result you get is
always stale data, in the real world, it is
not determined by distance only, it is
also determined by the traffic / accidents /
Adhoc traffic routes, so today's modern
ETA does not only require arithmetic /
numerical solution / predictions, you have to
be smart and you have to make use of
multiple data informations, so what happens in
these days, you are also going to use the

Machine Learning library
called as the Deep ETA
to predict the right ETAs

Noted Points -

Since, Uber is an open source so they have shared the information about them that how are they doing the work, they used

Traffic Data }
Map Data }

Both are put into
Routing Engines

Based on this data, It finds out the shortest path / distance b/w two places and whatever data comes from this routing engine is transferred to

Deep ETA

This predicts the right ETAs

So, ETA involves complex AI

too. What is

source/destination / time of request?