

lets design a bulletin-like social media centre, similar to facebook, Instagram etc.

- What is Twitter?
- social media service where users can read or post short messages (upto 280 characters) called tweets.
- available on web and mobile platforms such as android and iOS.

Requirements

- Functional
 - should be able to post new tweets (can be text, image, video etc.)
 - should be able to follow other users
 - should have a newsfeed feature consisting of tweets from the people the user is following
 - should be able to search tweets
- Non-functional
 - high availability with minimal latency
 - system should be scalable and efficient

Extended

- metrics and analytics
- reduced functionality
- favorite tweets

Estimation and Constraints

Traffic

There will be a read-heavy system.

- lets assume we have 1 billion users with 200 million daily active users (DAU), and on average each user tweets 5 times a day.

$$200 \text{ million} \times 5 \text{ tweets} = 1 \text{ billion/day}$$

- tweets can contain media such as images or videos. Assume 10% of tweets are media files shared by users.

$$10\% \times 1 \text{ billion} = 100 \text{ million/day files to store}$$

Request per second for our system?

1 billion requests per day translates into 12k requests per second.

$$\frac{1 \text{ billion}}{24 \text{ hrs} \times 3600 \text{ seconds}} = \sim 12k \text{ requests/second}$$

Storage

- assume each message on average is 100 bytes,

$$1 \text{ billion} \times 100 \text{ bytes} = \sim 100 \text{ GB/day}$$

- Also, 10% of daily messages (100million) are media files, each file to be 50KB on average

$$100 \text{ million} \times 50 \text{ KB/day} = \sim 5 \text{ TB/day}$$

- For 10 years, we require about 19PB of storage.

$$(5 \text{ TB} + 1 \text{ TB}) \times 365 \text{ days} \times 10 \text{ years} = \sim 19 \text{ PB}$$

Bandwidth

- our system is handling 5.1TB of ingests every day

$$\frac{5.1 \text{ TB}}{24 \text{ hrs} \times 3600 \text{ seconds}} = \sim 60 \text{ MB/second (minimum bandwidth)}$$

High level Database

Type	Estimate
Daily Active Users	100 million
Requests per second (RPS)	12k/s
Storage (per day)	~5.1TB
bandwidth	~60MB/s

Ranking Algorithm — as discussed we need a ranking algorithm to rank each tweet according to its relevance to each specific user.

For e.g. Facebook used to utilize an edgeRank algorithm. Here rank of each feed item is described by:

$$\text{Rank} = \text{Affinity} \times \text{Weight} \times \text{Decay}$$

where affinity → closeness of user to the creator of the edge, if a user frequently likes, comments or messages the edge creator, their value of affinity will be higher, resulting in a high rank for the post.

weight → value assigned according to each edge. A comment can have higher weightage than like, and thus a post with more comments is more likely to get a higher rank.

decay → is the measure of the creation of the edge. The older the edge, the lower will be the value of decay and eventually the rank.

Nowadays algorithms are much more complex and ranking is done using machine learning models which can take thousands of factors into consideration.

Retweets

- It is one of our extended requirements.

- To implement this feature, we can simply create a new tweet with the user ID of the user retweeting the original tweet and then modify the "type" enum and "content" property of the new tweet to link it with the original tweet.

For e.g. type enum property can be of type tweet, similar to text, video, image etc and content can be the ID of the original tweet.

→ how to very back implementation. To improve this, we can create a separate table itself to store retweets.

Search

- sometimes traditional DBMS are not performant enough, we need something which allows us to store, search and analyze huge volumes of data quickly, and in near real-time and give results with milliseconds. ElasticSearch can help us with this use case.

- ElasticSearch is a distributed, free and open search and analytics engines for all types of data, including textual, numerical, geospatial, structured, and unstructured. Built on top of Apache Lucene.

How do we identify trending topics?

- Trending functionality will be based on top of the search functionality.

- we can cache the most frequently searched queries, hashtags and topics in the last N seconds and update them every M seconds using some sort of batch job mechanism.

- our ranking algorithm can also be applied to the trending topics to give them more weight and prioritize them for the user.

Notifications

Push notifications are an integral part of any social media platform.

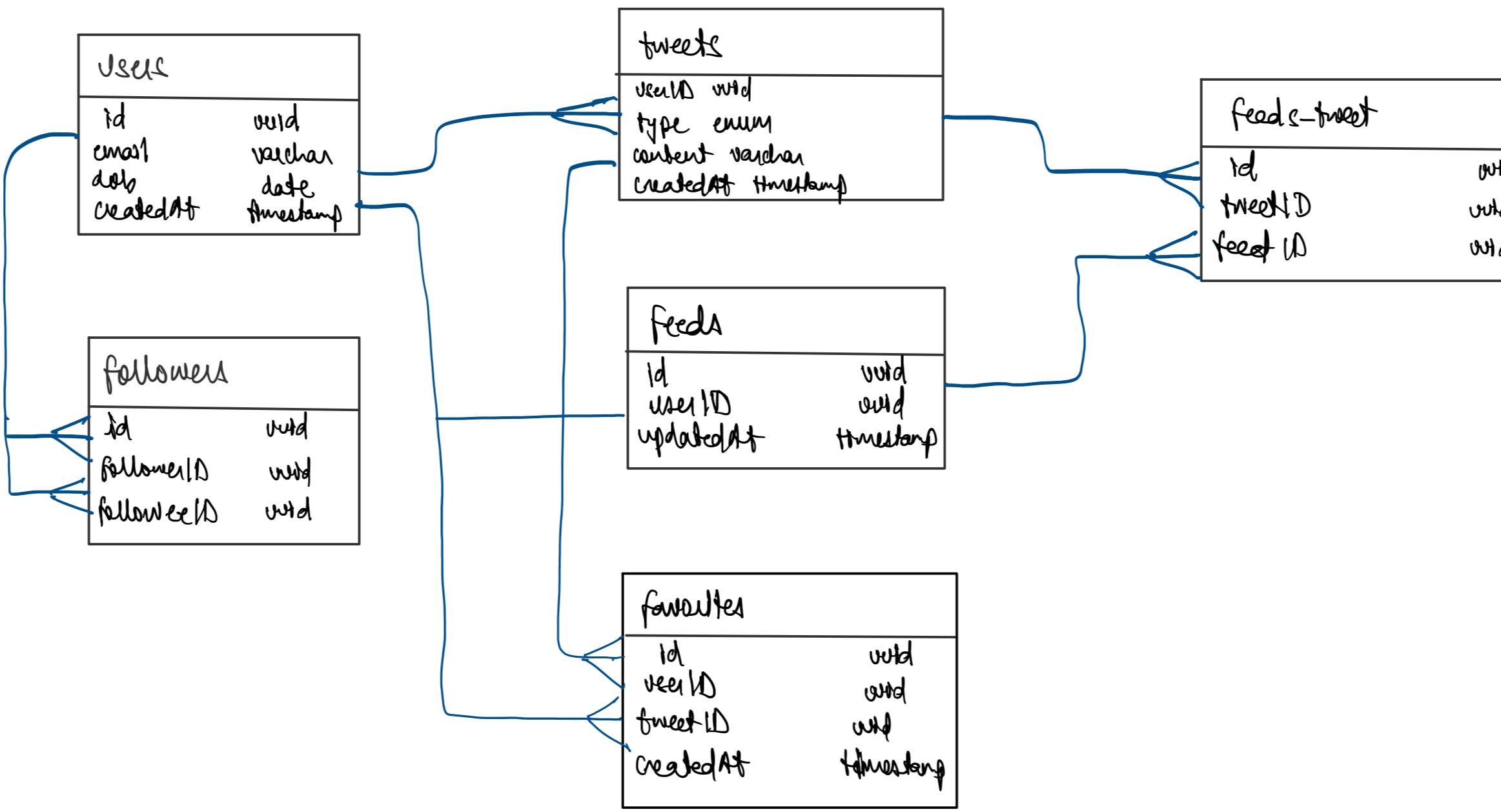
- we can use a messaging queue or a messaging broker such as Kafka with the notification service to dispatch requests to Pushbullet Cloud Messaging (FCM) or Apple Push Notification Service (APNS), which will handle the delivery of the push notifications to user devices.

Data Model Design

- Tables**
 - User** — this table will contain a user's information such as name, email, job and other details.
 - tweets** — this table will store tweets and their properties such as type (text, image, video etc.), content, context, etc. we will also store corresponding user.
 - favorites** — this table maps users with users for the favorite tweet functionality in our application.
 - followers** — this table maps the follower and followed users as users can follow each other (N:N relationship).
 - feeds** — this table stores feed properties with the corresponding user.
 - feed-tweet** — this table maps tweets and feed (N:M relationship).

Kind of database

- write data model seems quite relational, we don't necessarily need to store everything in a single db, as this can limit our scalability and quickly become a bottleneck.
- split data between different centers each having ownership over a particular table.
- we can use relational db such as PostgreSQL, or a distributed NoSQL db.



API Design

- Post a tweet — this API will allow user to post a tweet on the platform.
- `postTweet (userId: void, content: string, mediaURL?: string): boolean`
- Follow or unfollow a user — this API will allow user to follow or unfollow another user.
- `follow (followerId: void, followeeId: void): boolean`
- `unfollow (followerId: void, followeeId: void): boolean`
- Get newsfeed — this API will return all the tweets to be shown within a given newsfeed.
- `getNewsfeed (userId: void): Tweet[]`

High Level Design

- Architecture**
 - we will be using a microservices architecture since it will make it easier to horizontally scale and decouple our services.
 - each service will have ownership of its own data model.
- Core Services**
 - user service : handles user-related concerns such as authentication and user information.
 - newsfeed service : handles generation and publishing of user newsfeed. (discussed separately).
 - tweet service : handles tweet-related use cases such as posting a tweet, favorite etc.
 - search service : responsible for handling search related functionality. (discussed separately).
 - media service : handle the media (image, video, file etc.) uploads. (discussed separately).
 - notification service : simply push notification to the user.
 - analytics service : used for metrics and analytics use cases.

What about inter-service communication and service discovery?

- REST or HTTP performs well but we can further improve the performance using gRPC which is more lightweight and efficient.
- Service discovery is another thing we will have to take into account. We can use a service mesh that enables managed, observable and secure communication between individual services.

Detail Discussion

- Newsfeed** — when it comes to newsfeed, it seems easy enough to implement, but there are lots of things that can make or break this feature. Let's divide the problem into two parts:

Generation

- Let's assume we want to generate the feed for the user A, steps followed:
 1. Retrieve the IDs of all users and entities (friends, topics etc.) user A follows.
 2. Fetch the relevant tweets for each of the retrieved IDs.
 3. Use the ranking algorithm to rank the tweets based on parameters such as relevance, time, engagement etc.
 4. Deliver the ranked tweets data to the client in a paginated manner.

→ feed generation is an intensive process and can take quite a lot of time, especially for users following a lot of people.

- To improve the performance, the feed can be pre-generated and stored in the cache, then we can have a mechanism to periodically update the feed and apply our ranking algorithm to the new tweets.
- Publishing** — Publishing is the step where the feed data is pushed according to each specific user. This can be a quite heavy operation as a user may have millions of friends or followers. To deal, we have three approaches:
 - **Full Model (or fan-out on load)** — when a user creates a tweet, and a follower reloads their newsfeed, the feed is created and stored in memory. The most recent feed is only loaded when the user requests it. This approach reduces the no. of write operations on our database.
 - **downable** — user will not be able to view recent feeds unless they "pull" data from the server, which will increase the number of read operations on the server.
 - **Push Model (or fan-in on write)** — In this model, once a user creates a tweet, it is "pushed" to all follower's feeds immediately. This prevents the system from having to go through a user's entire follower list to check for updates. downside — it would increase the no. of write operations on the database.

Identity and resolve bottlenecks

- what if one of our services crashes?
- how will we distribute our traffic between our components?
- how can we reduce the load on our database?
- how to improve the availability of our cache?
- how can we make our notification system more robust?
- how can we reduce media storage cost?
- To make system more resilient:
 - running multiple instances of each of our services.
 - introducing load balancers between clients, servers, databases and cache servers.
 - using multiple read replicas for our distributed caches.
 - creating Once Delivery and message ordering is challenging in a distributed system, we can use a dedicated message broker such as NATS/Kafka to make our notification system more robust.
 - we can add media processing and compression capabilities to the media service to compress large files which will save a lot of storage space and reduce cost.

Detailed Design

Discussion of design decisions in detail:

Data Partitioning

- Horizontal partitioning/sharding can be a good first step. We can use partition schemes such as:
 - hash-based
 - UUID-based
 - range-based
 - composite partitioning

— Above approaches can still cause uneven data and load distribution, we can solve this using consistent hashing.

Mutual Friends

- For mutual friends, we can build a cocktail graph for each user.
- each node in graph will represent a user and directional edge will represent follower and followee. After that we can traverse the followers of a user to find and suggest a mutual friend.
- This will require a graph database such as Neo4j or Amazon Graph.

- This is a pretty simple algorithm, to improve our suggestion accuracy, we will need to incorporate a recommendation model which uses machine learning as part of our algorithm.

Metrics and Analytics

- Recording analytics and metrics is one of our extended requirements.
- As we will be using Apache Kafka to publish all sorts of events, we can process these events and run analytics on the data using Apache Spark which is an open-source unified analytics engine for large-scale data processing.

Caching

In a social media app, we have to be careful about using a cache as our users expect latest data.

— so to prevent usage spikes from our resources we can cache the top 20% of the tweets.

Cache-eviction policy

- we can use solutions like LRU or Memcached and cache 20% of the daily traffic but what kind of cache eviction policy?
- LRU can be a good choice. We discard least recently used key first.

How to handle cache miss

- whenever there is a cache miss, our servers can hit the database directly and update the cache with new entries.

Media access and storage

- our media service will be handling both access and storage of user media files.
- **Object storage** — object stores break data files up into pieces called objects. If then stores those objects in a single repository, which can be spread out across multiple networked systems.

— we can use distributed file storage such as HDFS or GlusterFS.

Content Delivery Network (CDN)

- CDN increases content availability and redundancy while reducing bandwidth costs. Generally, static files such as images, videos are served from CDN.

— we can use services like Amazon CloudFront or Cloudflare CDN for this use case.

