

SQL 13: SCHEMA DESIGN - 2

31/05/24

AGENDA

- ① Scaler Schema Design - Complete Today
- ② Represent Enums
- ③ Deciding PK of a mapping table
- ④ Represent FK
- ⑤ Represent Indexes
- ⑥ Case Study: Netflix Schema Design

Scaler Schema Design

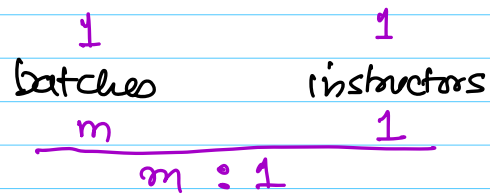
- ① batches
- ② Students
- ③ classes
- ④ instructors
- ⑤ mentors
- ⑥ mentor_sessions

Scaler Schema Design

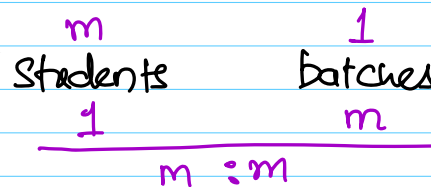
The requirements are as follows:

1. Scaler will have multiple batches.
2. For each batch, we need to store the name, start month and current instructor.
3. Each batch of Scaler will have multiple students.
4. Each batch has multiple classes.
5. For each class, store the name, date and time, instructor of the class.
6. For every student, we store their name, graduation year, University name, email, phone number.
7. Every student has a buddy, who is also a student.
8. A student may move from one batch to another.
9. For each batch a student moves to, the date of starting is stored.
10. Every student has a mentor.
11. For every mentor, we store their name and current company name.
12. Store information about all mentor sessions (time, duration, student, mentor, student rating, mentor rating).
13. For every batch, store if it is an Academy-batch or a DSML-batch.

- ① batches
- batch_id
 - batch_name
 - start_date
 - inst_id
 - bt_id

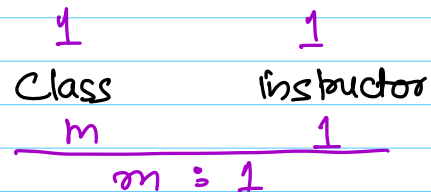


- ② Students
- st_id
 - st_name
 - email
 - phone-no
 - Univ-name
 - grad-year
 - buddy_id
 - mentor_id



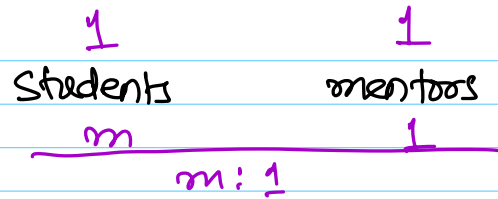
- ③ Student_batches
- st_id
 - b_id
 - date-of-start

- ④ classes
- c_id
 - c-name
 - c-time
 - inst_id



- ⑤ instructors
- inst_id
 - inst-name

- ⑥ mentors
- m_id
 - m_name
 - comp_name



- ⑦ mentor_sessions
- ms_id
 - ms_name
 - duration
 - st_rating
 - mentor_rating
 - start_time
 - st_id
 - mentor_id

Table: batches

①

b_id	b_name	b_type
1	A	Academy
2	B	DSML
3	C	DSML
4	D	Academy
5	E	Academy.
6	F	SST

- Pros
- ① Readability
 - ② Joins not required

- Cons
- ① Error prone while inserting data.
 - ② Memory - varchar(255) - Extra space.

(WIKES)

Query

```
SELECT * FROM BATCHES  
WHERE BTYPE = 'DSML';
```

③ string comparisons are slow.

②

Table: Batches

b.id	b.name	b.type
1	A	1
2	B	2
3	C	2
4	D	1
5	E	1
6	F	

Enum

Code / Documentation

Academy - 1

SST - 2

DSML - 3

Pros

- ① less space
- ② faster search/comparisons.

Cons

- ① No readability.
- ② Add/delete values in the middle
- can cause discrepancies
- ③ No info in DB about the mapping.

REPRESENT ENUMS

③ Look up table.

table: batch-types

bt-id	bt-name
1	Academy
2	DSML
3	SS T

Pros → All above problems are solved

Cons → Need to JOIN.

Mapping tables

① student-batches
- st-id
- b-id
- move-date

PK = (st-id, b-id)

② student-batches
- sb-id
- st-id
- b-id
- move-date

Pros → One extra col (less space)

② → index (st-id, b-id)
→ joins on st-id faster

Pros → ① size of index smaller.
(1 col).

Requirement

- Use Case ①:
- ① School has exams
 - ② For each a student joins, he/she has to take an exam (eg. contests, MBE etc)
 - ③ Each exam is associated with a batch.

Ques

Find all the exams taken by a student

Students

-st-id
-st-name

batches

-b-id
-b-name

Join

Student-batches

-st-id
-b-id
-more-date

Join

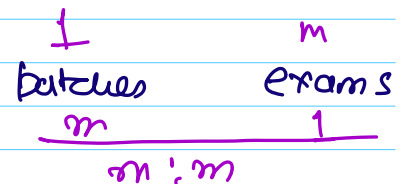
batch-exams

-be-id
-b-id
-e-id

Join

exams

-e-id
-e-name



d/p

st-name	exam-name
---------	-----------

Requirement:

Use Case ① Same exam will happen on different dates in different batches.

② If a student moves to another batch, he/she may have to give same exam again.

Ques

How to store the marks obtained in each exam for every student?

→ marks

→ date-of-exam.

✓ b-id

Student - exams

st-id	exam-id	marks	attempt-no.	date-of-exam
1	101	50		
1	101	70		

→ contest again in same batch.

Student - batch - exam

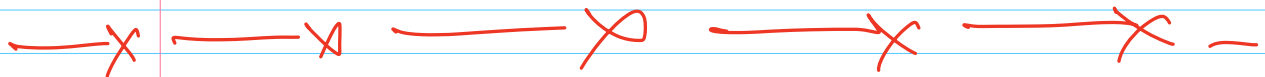
st-id	b-id	e-id	marks	doc	Attempt
1	1	101	50	-	
1	1	101	70	-	

PK = ? (st-id, b-id, exam-id)

st-id	<u>be-id</u>		marks	doe
	br-id	e-id		

st-id	be-id	marks	doe

twitter
Snowflake?



FK - All relations represented as FK

email → lot of queries

INDEXING

→ we will create an index on email. (Ans)

Note: Indexes will be created based on the UPPERCASE.

NOTE: DO NOT CREATE INDEXES PREMATURELY.

BREAK TILL - 8:33 AM.

Netflix Schema Design

Design Database Schema for a system like Netflix with following Use Cases.

Use Cases

1. Netflix has users.
2. Every user has an email and a password.
3. Users can create profiles to have separate independent environments.
4. Each profile has a name and a type. Type can be KID or ADULT.
5. There are multiple videos on netflix.
6. For each video, there will be a title, description and a cast.
7. A cast is a list of actors who were a part of the video. For each actor we need to know their name and list of videos they were a part of.
8. For every video, for any profile who watched that video, we need to know the status (COMPLETED/ IN PROGRESS).
9. For every profile for whom a video is in progress, we want to know their last watch timestamp.

Steps to do Schema Design

- ① Find all nouns/entities to create tables.
- ② Find all attributes (cols) that belong specifically to the tables. (PK)
- ③ Find the cardinalities of all relations

- add FK
- create mapping tables.

Tables.

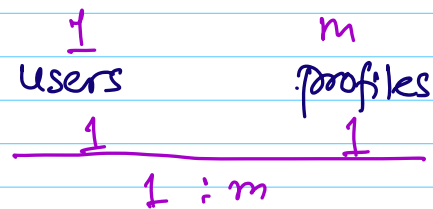
- ① users
- ② videos
- ③ profiles
- ④ profile-types (ENUM)
- ⑤ actors
- ⑥ watch-status (ENUM)

table : watch-type

pt-id	st-name

- ① users
 - u-id
 - u-name
 - email
 - password.

- ② profiles
 - p-id
 - p-name
 - u-id
 - pt-id



- ③ profile-types
 - pt-id
 - pt-name

ADULT	1
KID	2

④ videos

- v-id
- v-name
- title
- desc
- ~~wst-id~~

⑤ actors

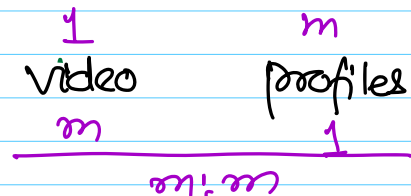
- a-id
- a-name

⑥ watch-status-type

- wst-id
- wst-name

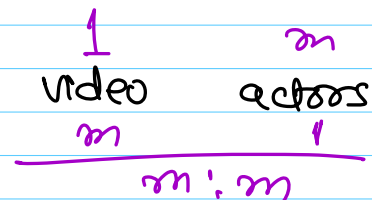
⑦ video-profiles

- vp-id
- video-id
- profile-id
- wst-id
- last-watch-timestamp.



⑧ video-actors

- video-id
- actor-id



PK of video-profiles. ?

PK = (profile-id, video-id)

Ques Why Range Queries are better in BST than Hashmap.

Hashmap-

$O(1) \rightarrow$ Avg time.

Worst TC $\rightarrow O(N)$

n = Total no. of elements $\rightarrow 10^6$ (100-200)
 k = no. of elements in range query $\rightarrow 100$

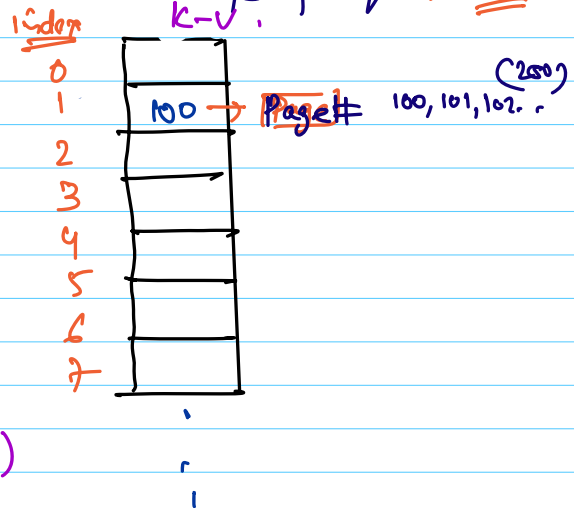
$\text{hash}(100) = 1$
 $\text{hash}(200) = 1$

$O(1) \rightarrow$ Avg

Worst $O(N)$.

HM $\rightarrow k - O(n \times k)$

HM, TC = $O(n \times k)$

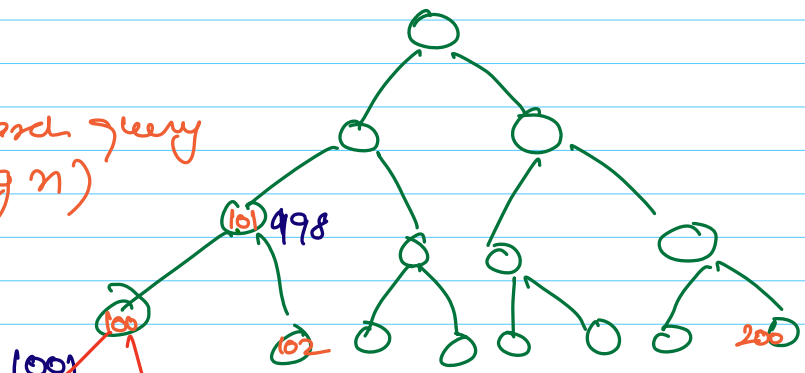


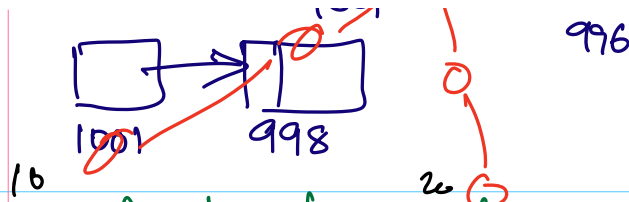
Balanced BST.

100 \in Start

Search

Start of search query
 TC = $O(\log n)$





→ in order traversal → ascending order (sorted)

TC (in order traversal)

for (range - k) → $O(k)$.

$$BBST = O(\log n) + O(k)$$

$$n \gg k.$$

$$f(n) = O(n \cdot k)$$

(left, root, right)

