

why databases?

downsides of querying data
by hand:

- error-prone
- tedious
- slow

Types of databases

relational
(SQL)

free { PostgreSQL - reddit, hipmunk
MySQL - Facebook, everybody
SQLite -
Oracle -

Google App Engine's
Database

Amazon
Dynamo

NoSQL

mongo
couch

SELECT * FROM links WHERE submitter_id = 5 OR
votes ≥ 23;

```
119 def query():
120     c = db.execute("select * from links where submitter_id = 62443 and votes > 1000")
121     link = Link(*c.fetchone())
122     return link.id
123
124 print query()
125
```

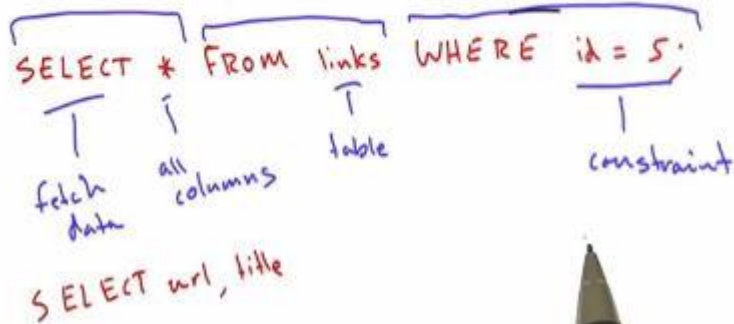
RUN

15

SQL

Structured Query Language

invented in the 1970s



Automatic Sharding and Replication

Google App Engine Datastore

Datastore is sharded
replicated

- won't have to think about scaling (too much)
- queries will be quick (because they have to be simple)
- will have to think about consistency

Automatic Sharding and Replication

Google App Engine Datastore

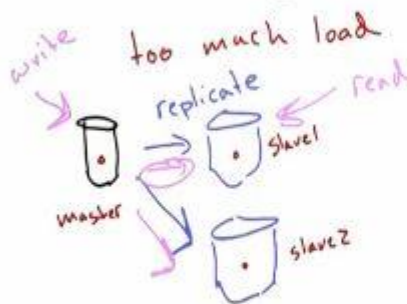
Datastore is sharded
replicated

- won't have to think about scaling (too much)
- queries will be quick (because they have to be simple)
- will have to think about consistency

```
hipmunk=# explain analyze select name from hotels where id = 51492;  
QUERY PLAN
```

```
Seq Scan on hotels (cost=0.00..52897.26 rows=1 width=23) (actual time=0.014..1  
42.342 rows=1 loops=1)  
  Filter: (id = 51492)  
  Total runtime: 142.402 ms  
(3 rows)
```

Scaling Databases



downsides

- doesn't increase write speed
- replication lag

Google
Datastore
⋮

too much data

Shard



downsides

- complex queries (range query)
- joins become difficult

Quiz

what is the query for fetching
all of the Arts from the database
sorted by creation time (most
recent first).
 \swarrow
 created

```
SELECT * FROM Art  
ORDER BY created DESC
```

Quiz

Integer Float String Date Time
 Text DateTime
Email Link Postal Address

Property

Type

title

String

art

~~String~~ Text

created

DateTime

String - < 500 chars
indexed

Text - > 500 chars
not indexed

Quiz

which is the most appropriate form element for inputting ascii art?

- ☐ `<input type="text">`
- ☒ `<textarea>`
- ☐ `<pre>`
- ☐ `<input type="password">`



Quiz

Do you understand everything there is to know about the App Engine Datastore?

- ☐ yes
- ☒ no

Quiz



which is an appropriate technique for growing a database that won't fit on one machine?

- ☐ replicate the database
- ☐ get a bigger hard disk
- ☒ shard the database
- ☐ store less data

Quiz

which is an appropriate technique for increasing the read speed from a database?

- ☐ get a faster machine
- ☒ replicate the database
- ☐ store less data
- ☐ press the turbo button

```
99 # QUIZ - Implement the function add_new_link() that both adds a link to the
100 # "links" list and updates the link_index dictionary.
101 def add_new_link(link):
102     links.append(link)
103     link_index[link.id] = link
104
105 l = Link(50, 1, 1, 1, "title", "url")
106 add_new_link(l)
107
108 print links[-1]
109 print link_by_id(50)
```

RUN

```
Link(id=50, submitter_id=1, submitted_time=1, votes=1, title='title', url='url')
Link(id=50, submitter_id=1, submitted_time=1, votes=1, title='title', url='url')
```

```
88 # QUIZ - Implement the function build_link_index() that creates a python dictionary
89 # the maps a link's ID to the link itself
90 def build_link_index():
91     index = {}
92     for l in links:
93         index[l.id] = l
94     return index
95
96 print build_link_index()
97
98
99
100 def link_by_id(link_id):
101     for l in links:
102         if l.id == link_id:
103             return l
```

RUN

```
{0: Link(id=0, submitter_id=60398, submitted_time=1334014208.0, votes=109, title='C
overtakes Java as the No. 1 programming language in the TIOBE index.',
```

```
88 # QUIZ - Implement the function link_by_id() that takes a link's ID and returns
89 # the link itself
90 def link_by_id(link_id):
91     for l in links:
92         if l.id == link_id:
93             return l
94
95 print link_by_id(24)
```

RUN

```
Link(id=24, submitter_id=48826, submitted_time=1333934004.0, votes=17, title='An R
programmer looks at Julia', url='http://www.r-bloggers.com/an-r-programmer-looks-at-
julia/')
1244
```

```
91 def query():
92     for l in links:
93         if l.id == 15:
94             return l.votes
95
96 print query() 1
```

RUN

1244

SELECT * FROM links
 WHERE votes > 10
 ORDER BY votes DESC
 (ascending by default)
 ASC - ascending
 DESC - descending

Joins

Link

ID	votes	user-id	title	url
5	206	22	1/12/12 blah	example.com

User

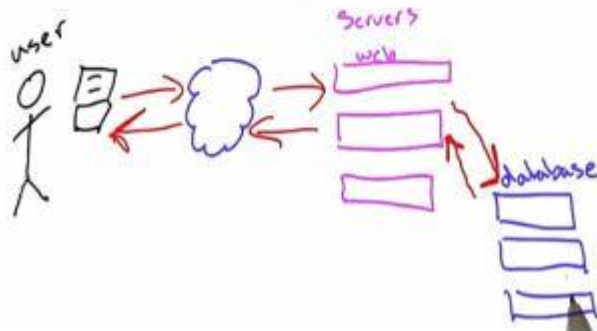
ID	name	password	date
22	Spez	hunter2	6/20/05

SELECT link.* FROM link, user WHERE link.user-id = user.id AND
 user.name = 'Spez'

Databases

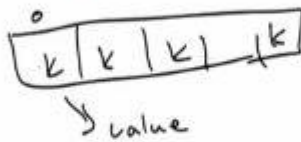
what is a database?

A program that stores and
retrieves ^{data}
large amounts
structured



Indexes for sorting

hashtable - not sorted, constant time



tree - sorted -
lookups are slower
 $\log n$

```
hipmunk=# create index hotel_id on hotels(id);  
CREATE INDEX  
hipmunk=# explain analyze select name from hotels where id = 51492;  
QUERY PLAN
```

```
Index Scan using hotel_id on hotels (cost=0.00..8.30 rows=1 width=23) (actual  
time=0.123..0.125 rows=1 loops=1)  
Index Cond: (id = 51492)  
Total runtime: 0.163 ms  
(3 rows)
```


Indexes

Sequential scans
(slow with a lot
of data)

Links = { link 1,
link 2,
link 3,
:
}

Indexes - increase the speed
of queries

index = ^{key → value} { 1: Link 1, index[2]
2: Link 2,
3: Link 3,
:
3

Datastore Types

Quiz

Integer Float String Date Time
DateTime
Email Link PostalAddress

Property

Type

title

art

created

Google App Engine Datastore

SQL \rightarrow GQL

- all queries begin with
SELECT *

- no joins

run arbitrary
queries

\rightarrow all queries must be
indexed

Google App Engine Datastore

tables \rightarrow entities

- columns are not fixed

- all have an ID

- parents / ancestors

Reddit \rightarrow Link



Reddit Hotness Algo

1. \uparrow Zombie Dogs!

2. \downarrow Ron Paul for president

3. \uparrow Look at this cat

4. \rightarrow

SELECT * FROM links
ORDER BY score DESC

Link

ID	ups	down	date	Score
	10	1	-	25

hot_idx (score)

float

+1

+ amt

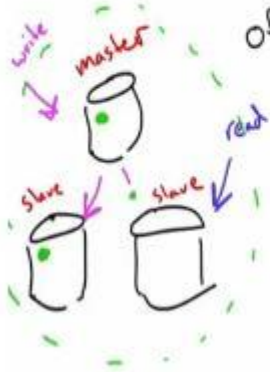
+1

+10

amt = hot(ups, downs, date)

Quiz

Replication lag is an example of the loss of which property?



- ☐ Atomicity
- ☒ Consistency
- ☐ Isolation
- ☐ Durability

ACID

Atomicity - all parts of a transaction succeed or fail together

Consistency - the database will always be consistent

Isolation - no transaction can interfere with another's

Durability - once the transaction is committed, it won't be lost.

up down
Link
score ✓

user
Karma ✓

Which statements are **true**?

☒ indexes increase the speed of database reads.

☐ indexes ~~increase~~ ^{decrease} the speed of database inserts.

which database is the
Best?

- ☐ MySQL
 - ☐ PostgreSQL
 - ☐ Google App Engine Datastore
 - ☐ Dynamo
- 

what can a database refer to?

- ☒ A program that stores and retrieves data
- ☒ the machine running that program
- ☒ a group of machines working together to store/retrieve data