



Rettwt

# Rebuilding the Twitter Platform

Daoyu Tu

Insight Data Engineering, New York

Summer 2016

# Why reinvent the wheel?

- Building a scalable & robust social media system
- Delivering real-time data to all the users



# Demo

## LaUnion PuebloEntero:

Mon Jun 27 00:16:20

RT @AmyHeartLive: Amazing bridal shower with my amazing girls! #weddingday #weddingshower #wedding #bridalshower #luckyinlove #lucky https:...

Hash tags: weddingday, weddingshower, wedding, bridalshower, luckyinlove, lucky

## jackie:

Mon Jun 27 00:08:37

pretty sure my new friend is talking to my ex and it's awkward but she doesn't know he's my ex and I have no right to ruin her happiness so

← daoyu.online 10 ⋮



# Rettiwt

## Please sign in ;^)

☐ Remember me

# Feeds Table



↑  
All my users

Mr\_Geek → Login



# A tweet comes in

**Stephen  
Curry tweets:**

Hello

userId

8:30

# A tweet comes in

Stephen  
Curry tweets:

Hello

userId

8:30



Followers: userId[]

# A tweet comes in

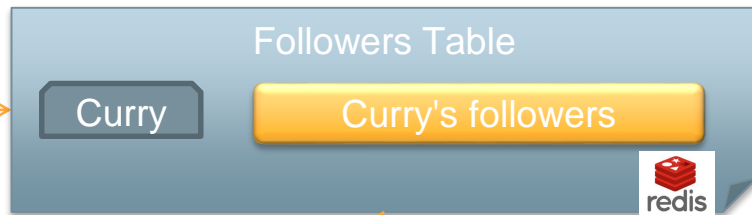
Stephen Curry tweets:

Hello

userId

8:30

All my users



Followers: userId[]

Curry's followers

Feeds Table

Follower #1

7:00

JAVA

6:00

C++

...

Follower #2

7:30

Dota2

...

Follower #n

8:00

Cute Dogs

...

Other users

8:20

Rock music

...



# A tweet comes in

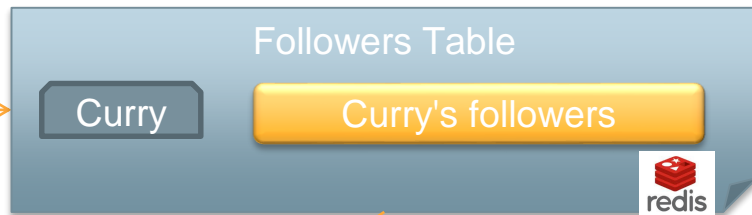
Stephen Curry tweets:

Hello

userId

8:30

All my users



Followers: userId[]

Curry's followers

Feeds Table

Follower #1

8:30

Hello

7:00

JAVA

6:00

C++

Follower #2

8:30

Hello

7:30

Dota2

...

Follower #n

8:30

Hello

8:00

Cute Dogs

...

Other users

8:20

Rock music

...



# Pipeline



# Benchmarking

- Save the tweet or tweet Id?
- Redis or Cassandra?

## Denormalized Schema

### Feeds

- uid: Long
- tweet: String
- time: timestamp

vs

## Normalized Schema

### Feeds

- uid: Long
- tid: Long
- time: timestamp



### Tweets

- tid: Long
- tweet: String

| Results | Norm               | Denorm             |
|---------|--------------------|--------------------|
| Writing | 1974 recs/s        | 691 recs/s         |
| Reading | 14ms for 20 tweets | 10ms for 20 tweets |

# Benchmarking

- Save the tweet or tweet Id?
- Redis or Cassandra?

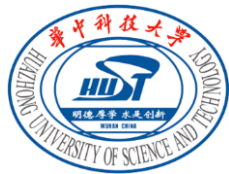


| Results | Redis         | Cassandra     |
|---------|---------------|---------------|
| Tweets  | 24.4 tweets/s | 21.9 tweets/s |

| Redis   | Cassandra  |
|---|--|
| <ul style="list-style-type: none"><li>• Fast in-memory reads</li><li>• Pure key-value storage</li></ul> | <ul style="list-style-type: none"><li>• Fault tolerance</li><li>• More storage</li></ul> |

# About Me

- Daoyu Tu
- Computer Science B.S & M.S.
- Web Developer



# Some Stats

- 2 clusters of 4 m4.xlarge nodes on AWS, one for Cassandra, one for Kafka, Spark & Redis
- 100,000 users with average degree of 94
- Writing 1974 records, processing 24.4 tweets per second