

# Big Data Infrastructure

CS 489/698 Big Data Infrastructure (Winter 2016)

Week 11: Analyzing Graphs, Redux (2/2)

March 24, 2016

Jimmy Lin

David R. Cheriton School of Computer Science

University of Waterloo

These slides are available at <http://lintool.github.io/bigdata-2016w/>

This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States  
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details



## **Theme for Today:**

How things work in the real world  
(forget everything I told you...)



**From the Ivory Tower...**





**... to building sh\*t that works**

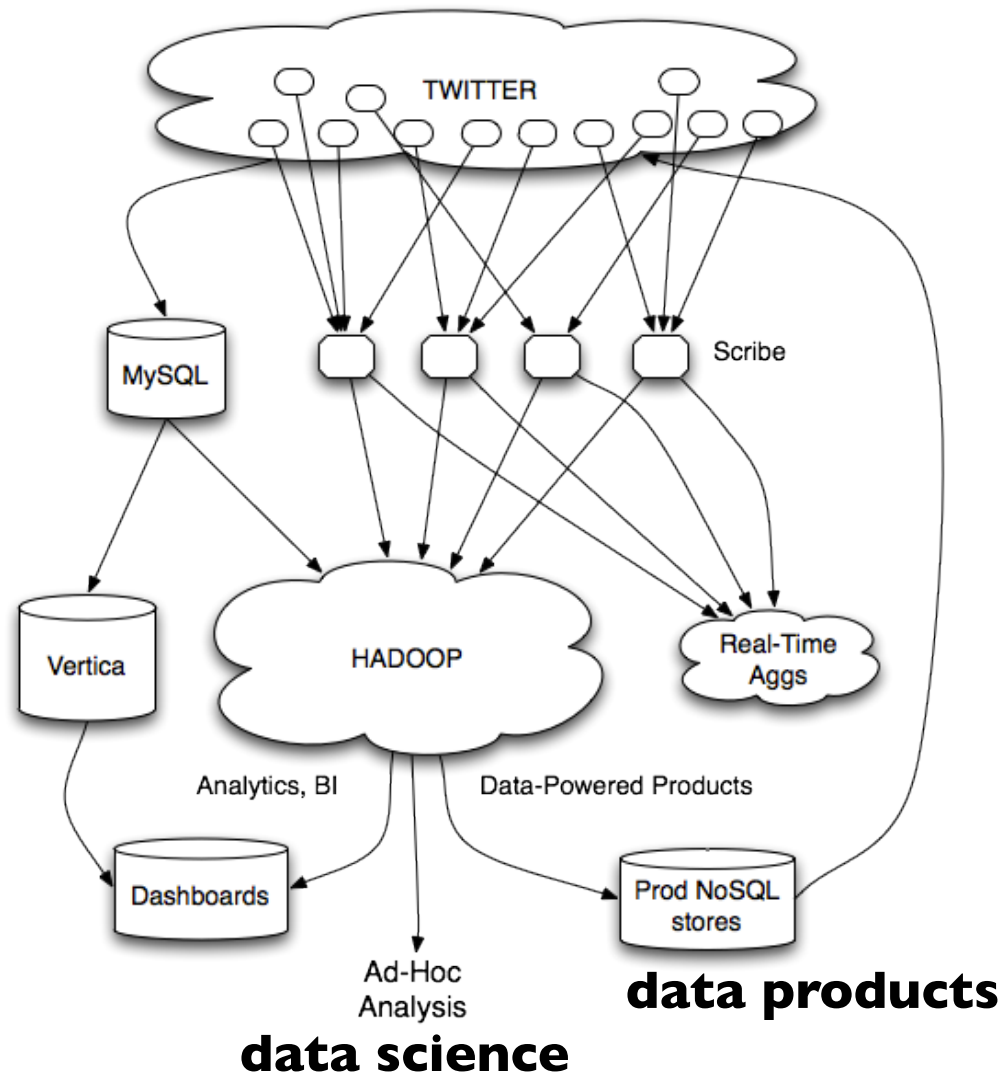




UNIVERSITY OF  
**WATERLOO**

**... and back.**





**I worked on...**

- analytics infrastructure to support data science**
- data products to surface relevant content to users**



Mishne et al. Fast Data in the Era of Big Data: Twitter's Real-Time Related Query Suggestion Architecture. SIGMOD 2013.



@lintool

Tweets

Struggling with complex data of Data Science 2/20 to rethi  
Promoted by Cloudera

TWEETS FOLLOWING FOLLOWERS

1,64

Leibert et al. Automatic Management of Partitioned, Replicated Search Services. SoCC 2011

Compose new Tweet...

Who to follow · Refresh · View all



plotly @plotlygraphs  
Follow Promoted



Brad Anderson @boorad  
Followed by Florian Leibert ...  
Follow



Sheila Morrissey @sheilaMorr  
Follow

Popular accounts · Find friends

Clinton Paquin @clintonpaquin  
Simply stated, "The only prot muscle memory" @TheChan  
View conversation

The Hill @thehill · 1h  
Republicans take debt ceiling  
View summary

Retweeted by Alex Feinberg  
Popehat @Popehat · 10h  
In a world in which few thing feed does.  
Expand

sochi| search results: #Sochi2014, #SochiProblems, Sochi, #SochiFail, Sochi 2014 @Sochi2014, Sochi Olympics 2014 @2014Sochi, Игры Сочи 2014 @sochi2014\_ru, Sochi Problems @SochiProblem, NYT Olympics @SochiNYT, Sochi Problems @SochiProblems, Search all people for sochi

Trends · Change

- #Olymp
- Ukraine
- #Conf
- Venny
- #Premi

I worked on...

– analytics infrastructure to support data science

– data products to surface relevant content to users

Gupta et al. WTF: The Who to Follow Service at Twitter. WWW 2013  
Lin and Kolcz. Large-Scale Machine Learning at Twitter. SIGMOD 2012





**circa ~2010**

~150 people total

~60 Hadoop nodes

~6 people use analytics stack daily

**circa ~2012**

~1400 people total

10s of Ks of Hadoop nodes, multiple DCs

10s of PBs total Hadoop DW capacity

~100 TB ingest daily

dozens of teams use Hadoop daily

10s of Ks of Hadoop jobs daily




# WTF

(what out of follow)

Who to follow · refresh · view all

-  **freshbooks** FreshBooks · Follow ×  
Promoted · Followed by @zappos and others.
-  **alanwarms** Alan Warms · Follow ×  
Followed by @fredwilson and others.
-  **Mozzie21** Moises Henriques · Follow ×  
can eat

Similar to @ryanhall3 · view all

-  **RunnerSpace\_com** RunnerSpace.com · Follow  
RunnerSpace.com has the latest in news and media...
-  **chrislieto** chris lieto · Follow  
Chris Lieto is a top ranked World Class Triathlete, ...
-  **runningtimes** runningtimes · Follow

Launched summer 2010





MG Siegler @parislemon · Jul 27

@kevinweil @elizabeth OMG just seeing you secured @thirdweil for baby. Most amazing handle ever. Well played. (via @amy)

Reply Retweet Favorite More

Details



Elizabeth Weil

@elizabeth



Follow

@parislemon Got it in January. Then Twitter recommended the account to my uncle. Family guessed we were pregnant. Oops. Denied it all. ;)

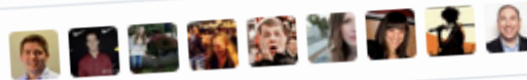
Reply Retweet Favorite Share More

RETWEETS

2

FAVORITES

20



3:39 PM - 27 Jul 2014

# #numbers

(Second half of 2012)

~175 million active users

~20 billion edges

42% edges bidirectional

Avg shortest path length: 4.05

40% as many unfollows as follows daily

WTF responsible for ~1/8 of the edges

**A talk in three episodes...**

# **PROLOGUE**



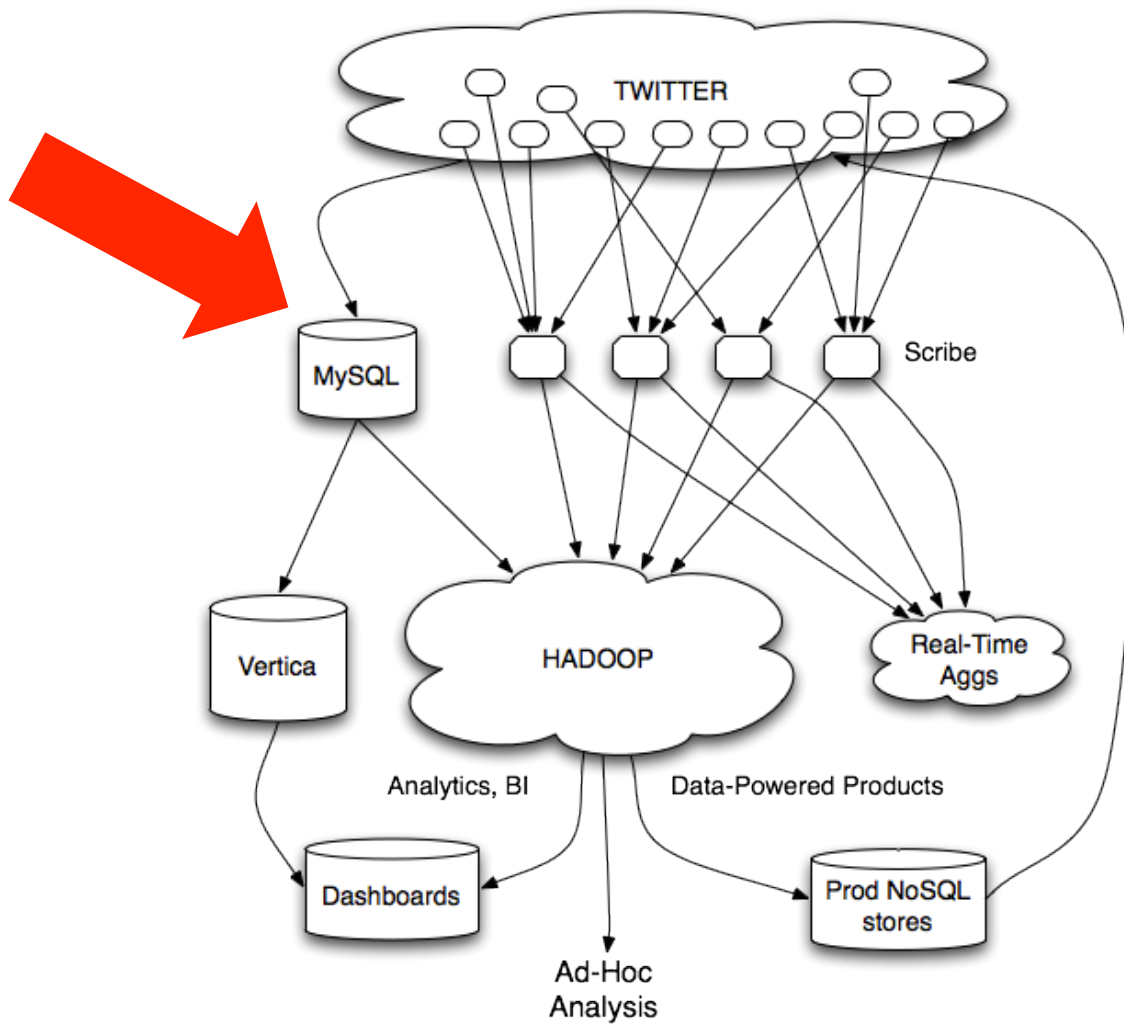


# **flockDB**

(graph database)

Simple graph operations  
Set intersection operations

**Not appropriate for graph algorithms!**





# Use Hadoop!

**MapReduce sucks for graph algorithms...**

Java verbosity

Long task startup

Stragglers

Needless graph shuffling

Frequent checkpointing

# What about?

HaLoop (VLDB 2010)

Twister (MapReduce Workshop 2010)

Pregel/Giraph (SIGMOD 2010)

Graphlab (UAI 2010)

PrIter (SoCC 2011)

Datalog on Hyracks (Tech report, 2012)

Spark/GraphX (NSDI 2012, arXiv 2014)

PowerGraph (OSDI 2012)

GRACE (CIDR 2013)

Mizan (EuroSys 2013)

...

A LONG TIME AGO IN A GALAXY FAR FAR AWAY...



CIRCA 2010



**MapReduce sucks for graph algorithms...  
Let's build our own system!**

**Key design decision:**

Keep entire graph in memory... on a single machine!

# Nuts!

## Why?

Because we can!

Graph partitioning is hard... so don't do it

Simple architecture

**Right choice at the time!**



# The runway argument





Suppose:  $10 \times 10^9$  edges  
(src, dest) pairs:  $\sim 80$  GB

$18 \times 8$  GB DIMMS = 144 GB

$18 \times 16$  GB DIMMS = 288 GB

$12 \times 16$  GB DIMMS = 192 GB

$12 \times 32$  GB DIMMS = 384 GB

# Cassovary

In-memory graph engine

Implemented in Scala

Compact in-memory representations

But no compression

Avoid JVM object overhead!

Open-source



# PageRank

“Semi-streaming” algorithm

Keep vertex state in memory, stream over edges

Each pass = one PageRank iteration

Bottlenecked by memory bandwidth

Convergence?

Don't run from scratch... use previous values

A few passes are sufficient

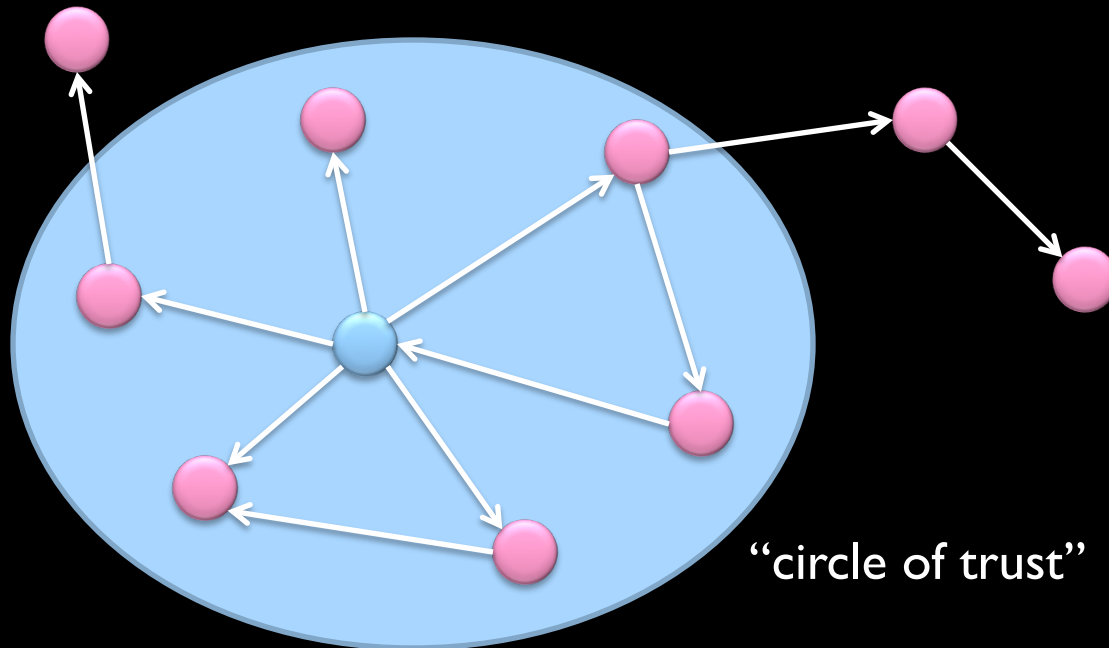


# “Circle of Trust”

Ordered set of important neighbors for a user

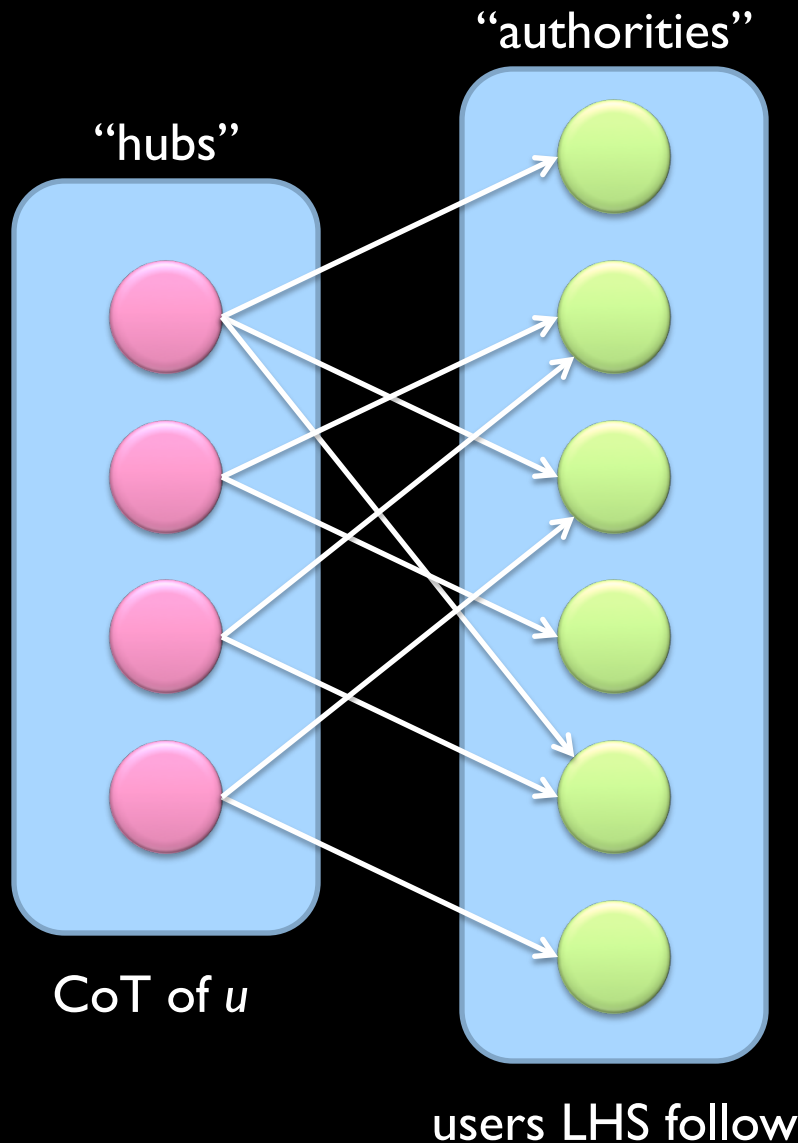
Result of egocentric random walk

Computed online based on various input parameters



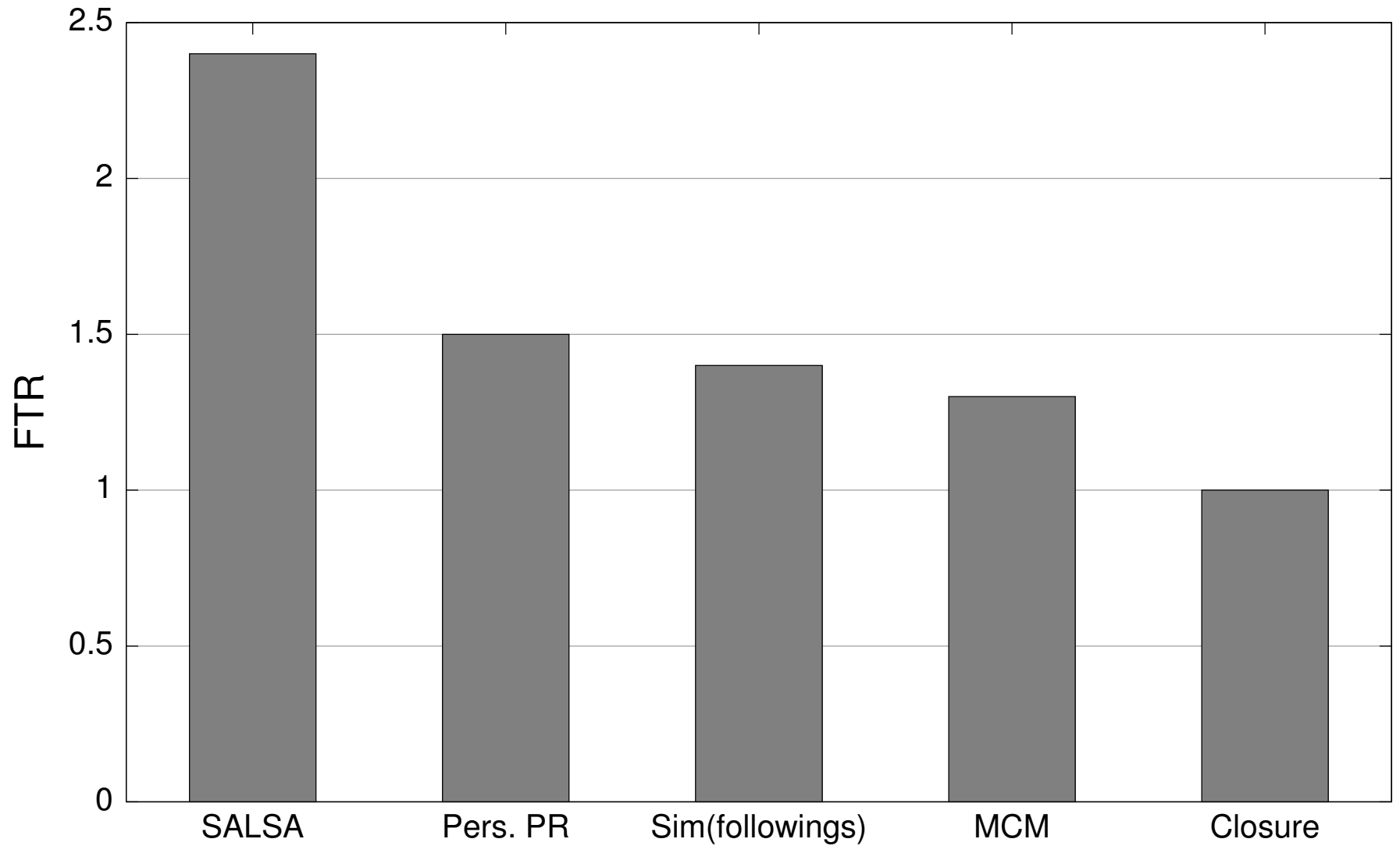
One of the features used in search

# SALSA for Recommendations

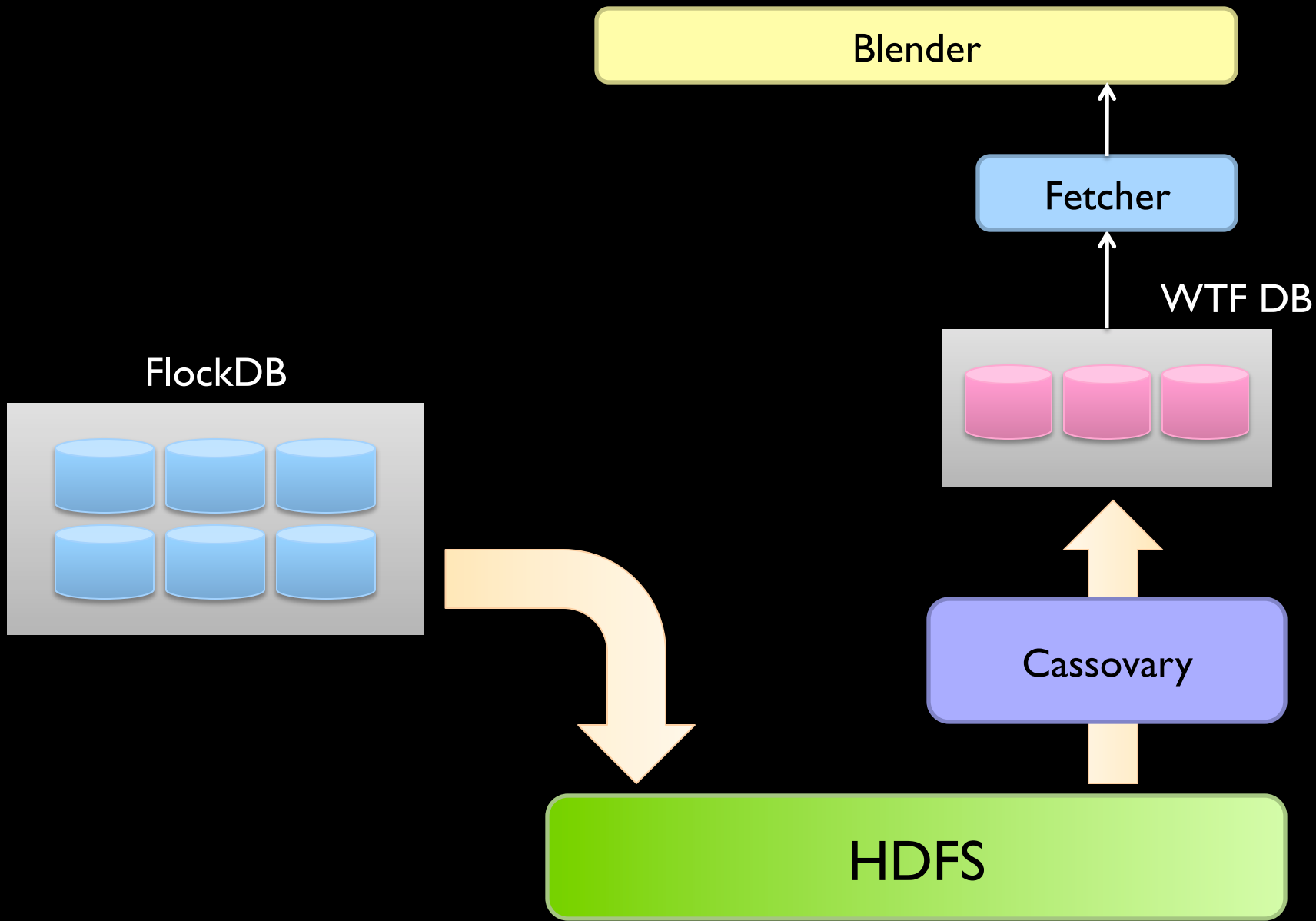


**hubs scores:**  
similarity scores to  $u$

**authority scores:**  
recommendation scores for  $u$

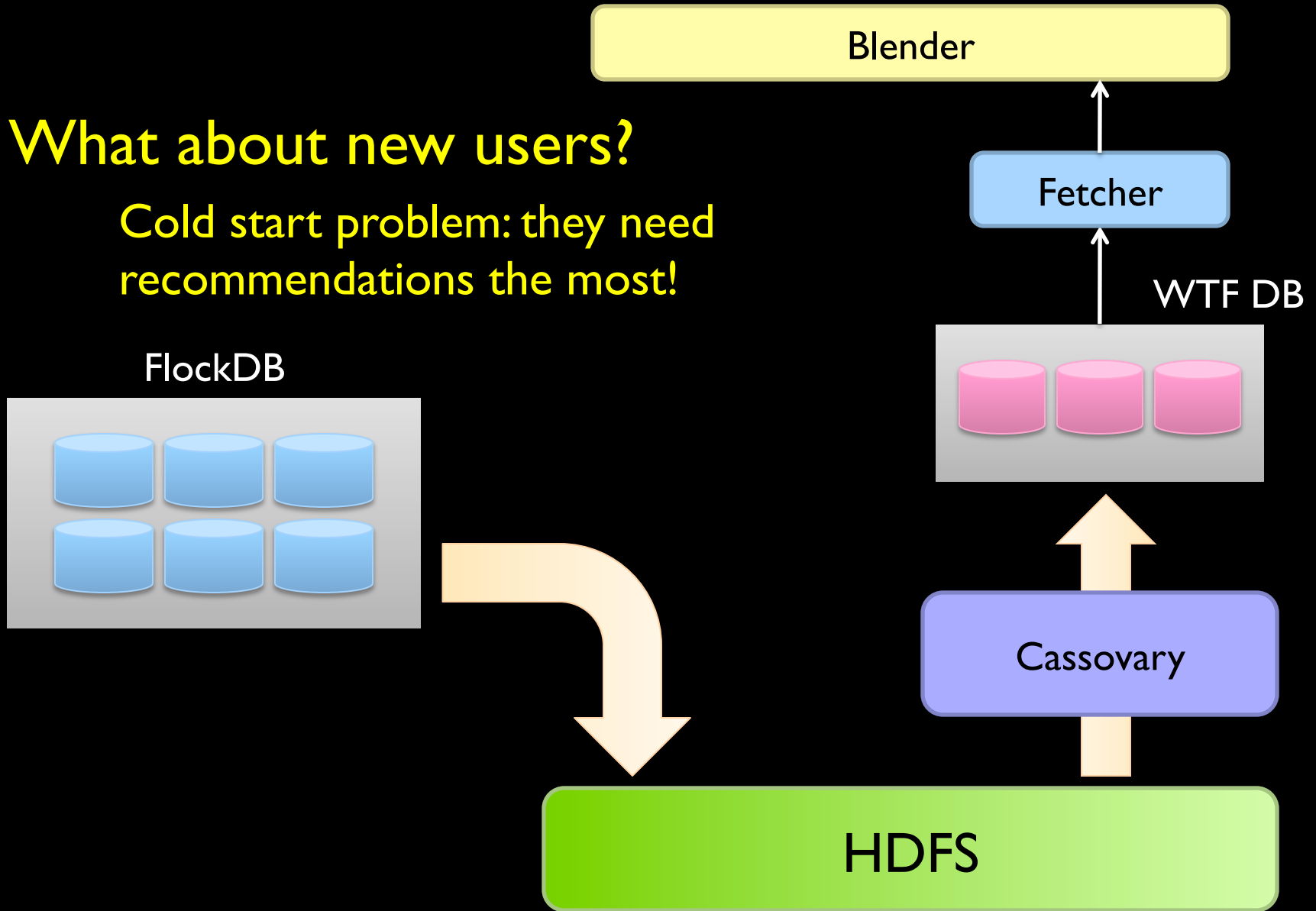


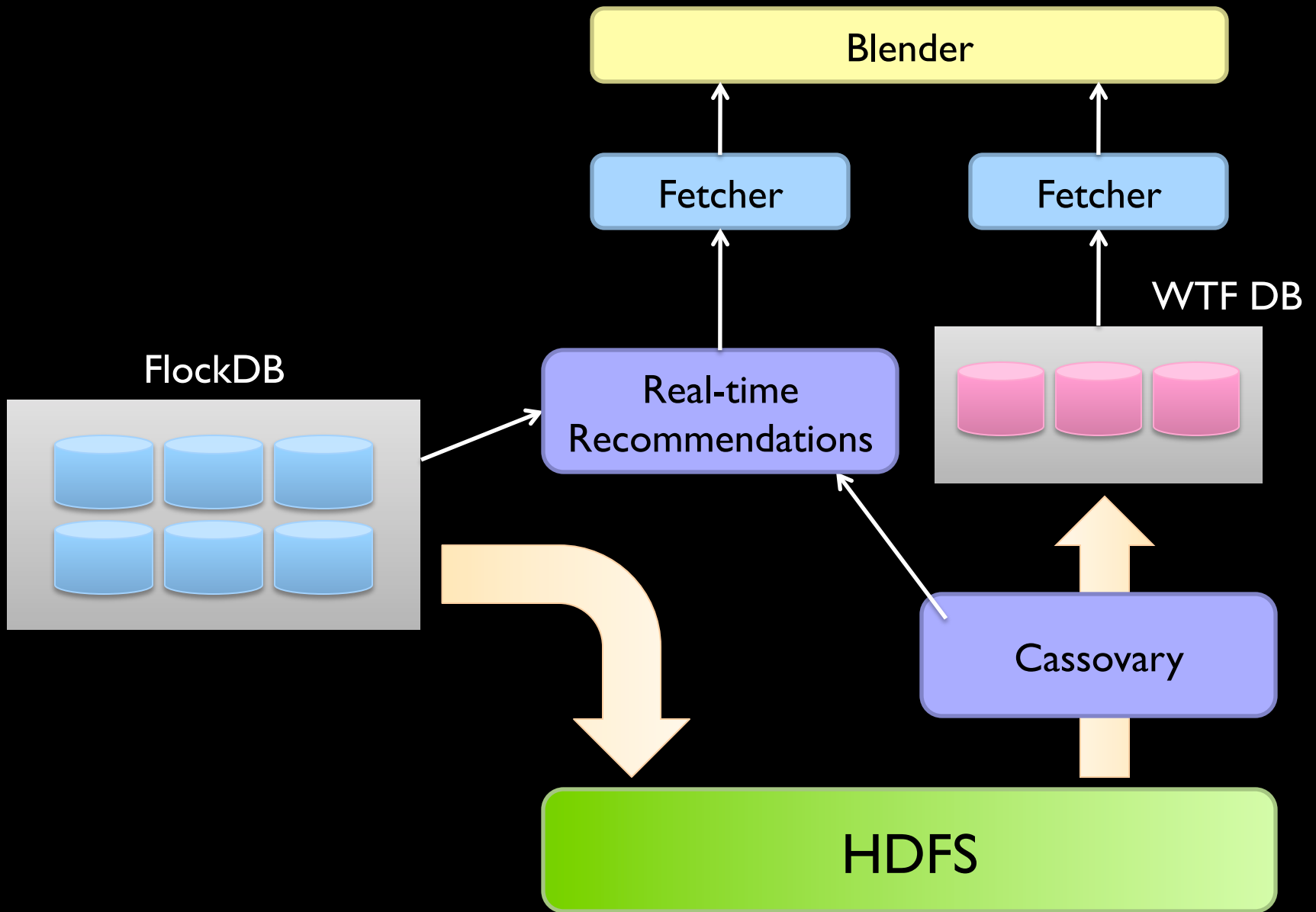




# What about new users?

Cold start problem: they need recommendations the most!







Spring 2010: no WTF  
seriously, WTF?

Summer 2010: WTF launched

THE **STAR WARS** SAGA CONTINUES...



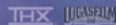
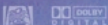
**HADOOP**  
 THE  
~~EMPIRE~~  
 STRIKES BACK

TWENTIETH CENTURY FOX PRESENTS A LUCASFILM LTD. PRODUCTION STAR WARS: EPISODE V - THE EMPIRE STRIKES BACK

STORY BY MARK HAMILL, CHARLES LASKER, AND GEORGE LUCAS. SCREENPLAY BY LEIGH BRACKETT AND LAWRENCE KASDAN. DIRECTED BY IRVYN KERSHNER

CASTING BY DAVID PROWSE. COSTUME DESIGNER KENNY BAKER. HAIR BY PETER MATHEW. MAKEUP BY FRANK OZ

EXECUTIVE PRODUCERS IRVYN KERSHNER, LEIGH BRACKETT, AND LAWRENCE KASDAN. PRODUCED BY GEORGE LUCAS AND GEORGE LUCAS. DIRECTED BY JOHN WILLIAMS



A LUCASFILM LTD. PRODUCTION  
 A TWENTIETH CENTURY RELEASE  
 WWW.STAR.WARS.COM

**CIRCA 2012**



Another “interesting” design choice:  
**We migrated from Cassovary back to Hadoop!**



**Whaaaaa?**

Cassovary was a stopgap!

**Hadoop provides:**

Richer graph structure

Simplified production infrastructure

Scaling and fault-tolerance “for free”

**Right choice at the time!**

# **Wait, didn't you say MapReduce sucks?**

What exactly is the issue?

Random walks on egocentric 2-hop neighborhood

Naïve approach: self-joins to materialize, then run algorithm

The shuffle is what kills you!

# Graph algorithms in MapReduce

Tackle the shuffling problem

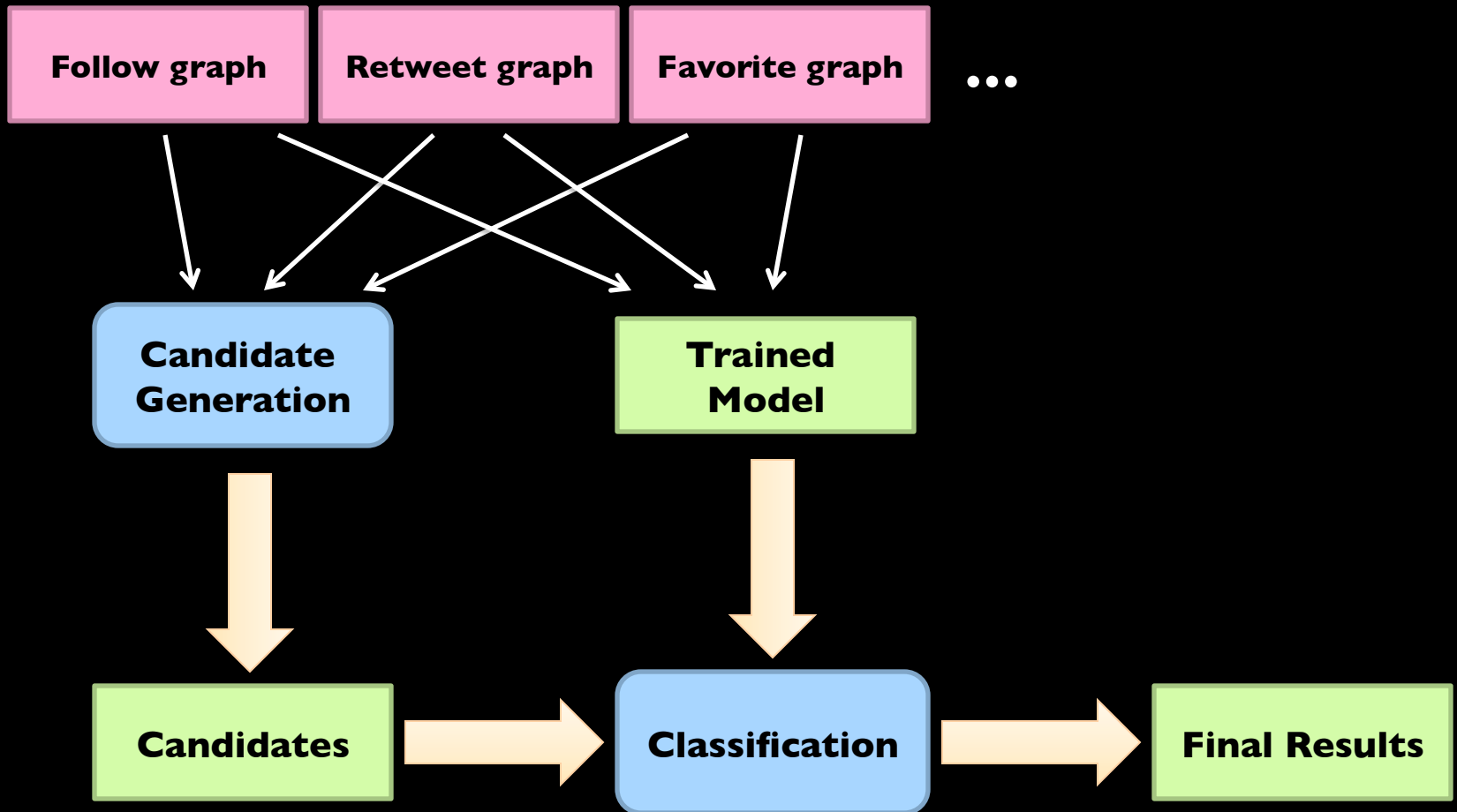
## Key insights:

Batch and “stich together” partial random walks\*

Clever sampling to avoid full materialization

\* Sarma et al. Estimating PageRank on Graph Streams. PODS 200  
Bahmani et al. Fast Personalized PageRank on MapReduce. SIGMOD 2011.

# Throw in ML while we're at it...





RETURN TO A GALAXY FAR FAR AWAY...



CUSTOM ARCHITECTURES

ALLUCASFEM LTD. PRODUCTION  
A TWENTIETH CENTURY FOX FILM  
WWW.STAR.WARS.COM

CIRCA 2013





**@dickc**  
dick costolo

Our mission: Instantly connect people everywhere to what's most meaningful to them. #mwc11

14 Feb via web ☆ Favorite ↻ Retweet ↩ Reply

**Isn't the point of Twitter real-time?**  
So why is WTF still dominated by batch processing?

## **From batch to real-time recommendations:**

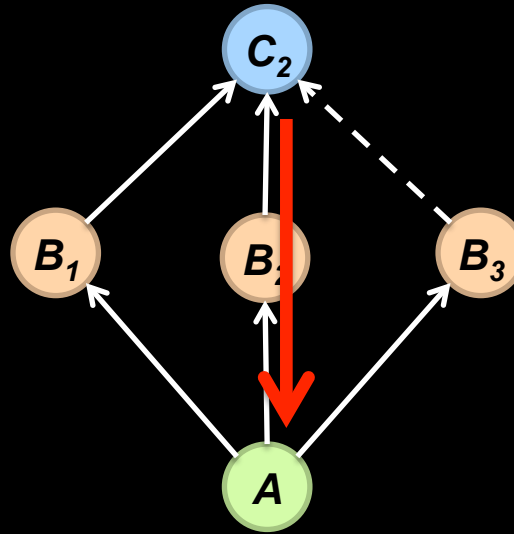
Recommendations based on recent activity

“Trending in your network”

## **Inverts the WTF problem:**

For this user, what recommendations to generate?

Given this new edge, which user to make recommendations to?



## Why does this work?

A follows B's because they're interesting  
B's following C's because "something's happening"  
(generalizes to any activity)



## Scale of the Problem

$O(10^8)$  vertices,  $O(10^{10})$  edges

Designed for  $O(10^4)$  events per second

## Naïve solutions:

Poll each vertex periodically

Materialize everyone's two-hop neighborhood, intersect

## Production solution:

**Idea #1:** Convert problem into adjacency list intersection

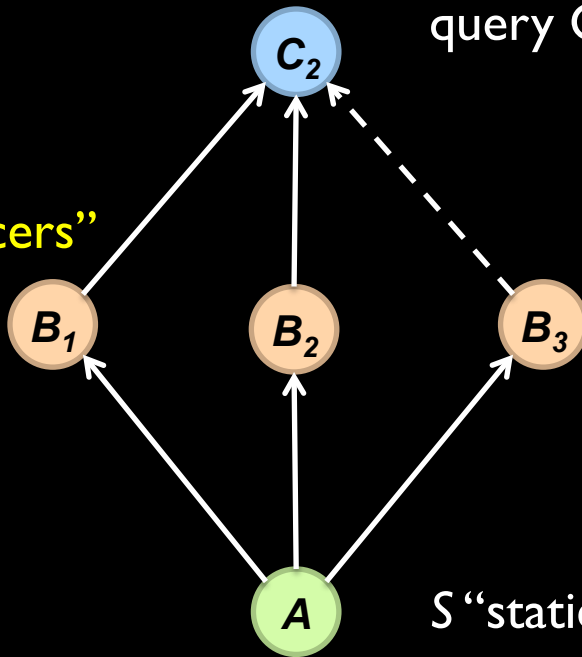
**Idea #2:** Partition graph to eliminate non-local intersections

# Single Node Solution

D “dynamic” structure:  
stores inverted adjacency lists  
query C, return all B’s that link to it

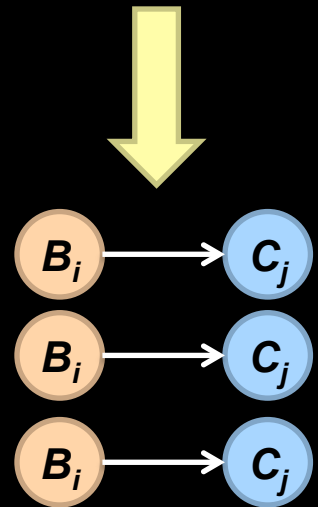
Who we’re recommending

“influencers”



Who we’re making the recommendations to

S “static” structure:  
stores inverted adjacency lists  
query B, return all A’s that link to it



# Algorithm

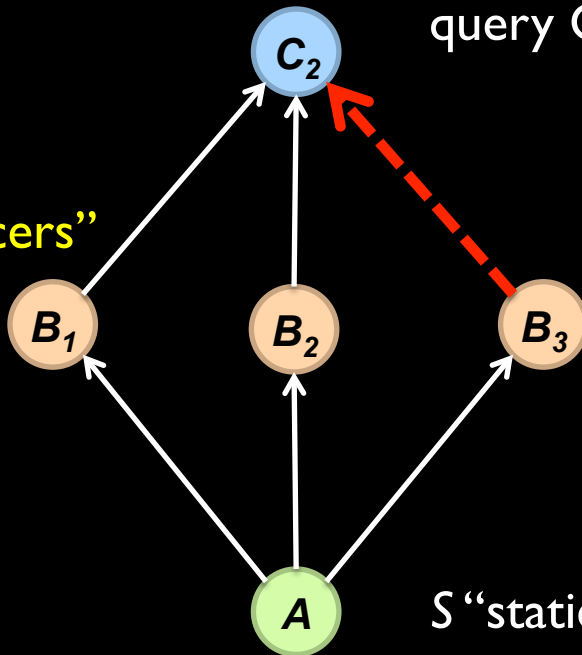
$D$  “dynamic” structure:

stores inverted adjacency lists

query  $C$ , return all  $B$ 's that link to it

Who we're recommending

“influencers”



Who we're making the recommendations to

$S$  “static” structure:

stores inverted adjacency lists

query  $B$ , return all  $A$ 's that link to it

1. Receive  $B_3$  to  $C_2$

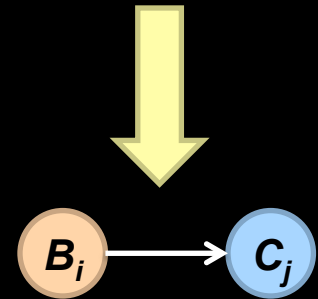
2. Query  $D$  for  $C_2$ , get  $B_1, B_2, B_3$

3. For each  $B_1, B_2, B_3$ , query  $S$

4. Intersect lists to compute  $A$ 's

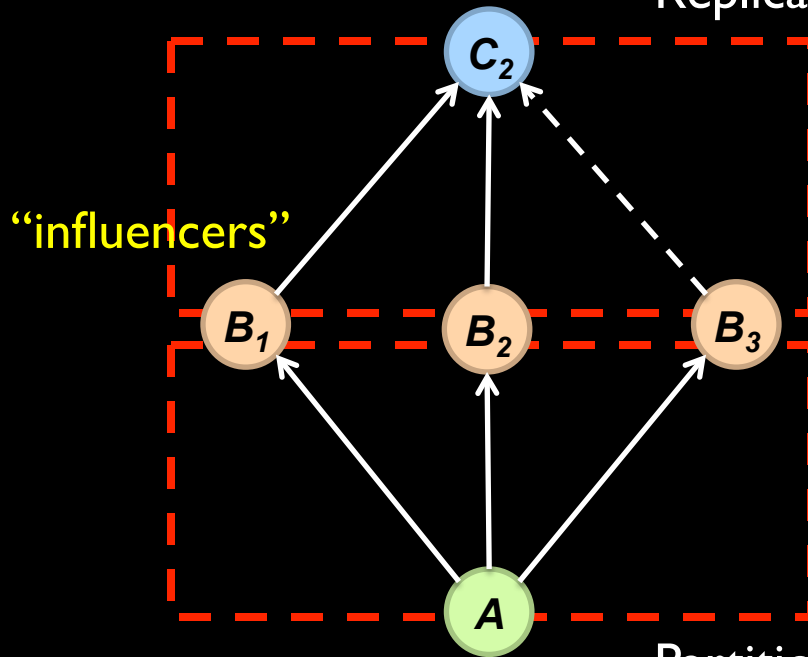
**Idea #1:** Convert problem into adjacency list intersection

# Distributed Solution

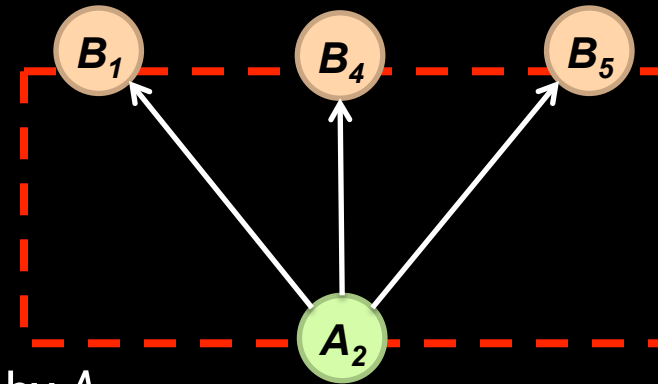


Who we're recommending

Replicate on every node



1. Fan out new edge to every node
2. Run algorithm on each partition
3. Gather results from each partition



Who we're making the recommendations to

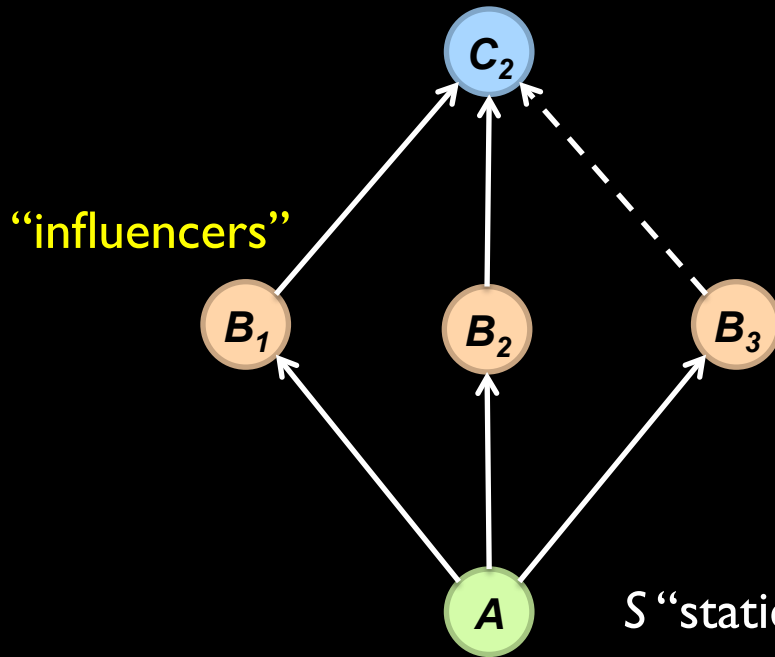
**Idea #2:** Partition graph to eliminate non-local intersections



# Notes

D “dynamic” structure:  
stores inverted adjacency lists

Who we’re recommending



Who we’re making the  
recommendations to

Memory pressure?

Why?

S “static” structure:  
stores inverted adjacency lists

# Production Status

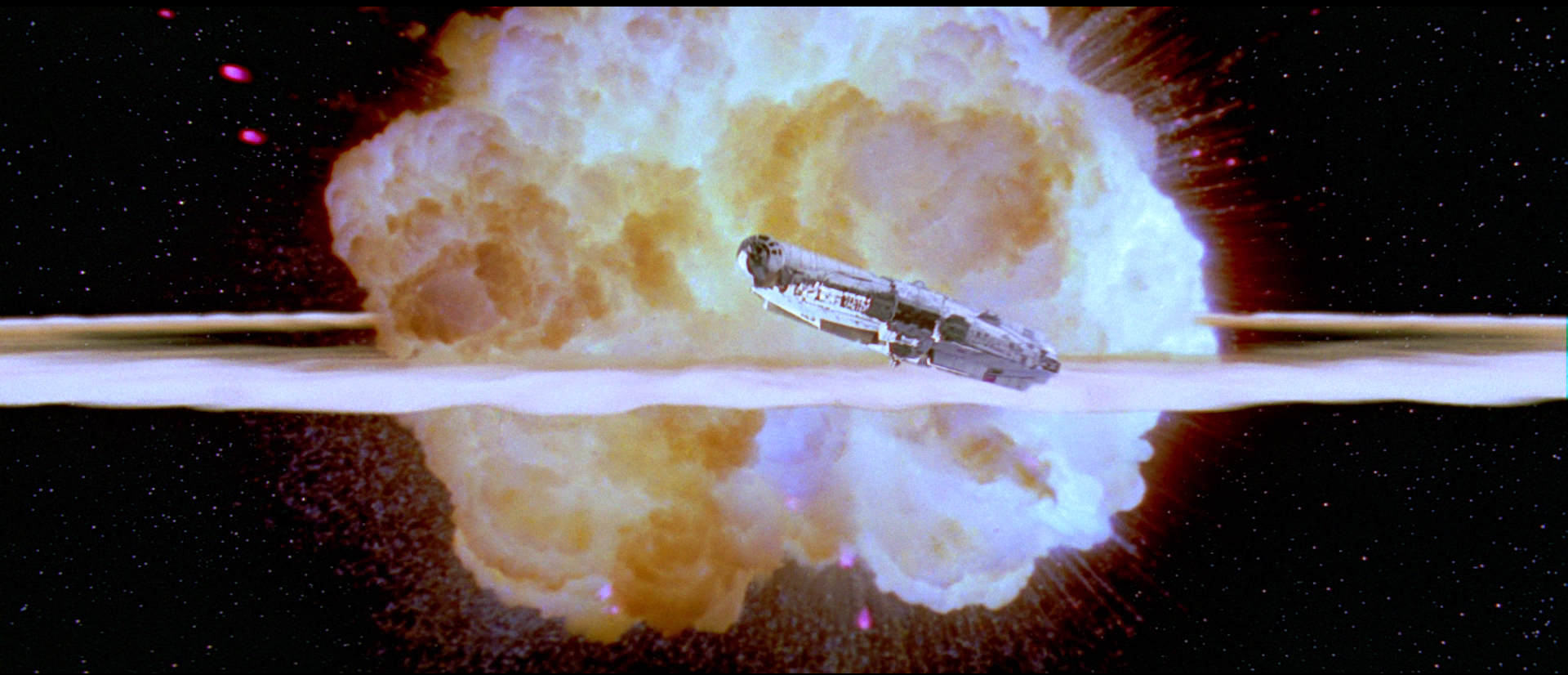
Launched September 2013

Push recommendations to Twitter mobile users

Billions of raw candidates, millions of push notifications daily

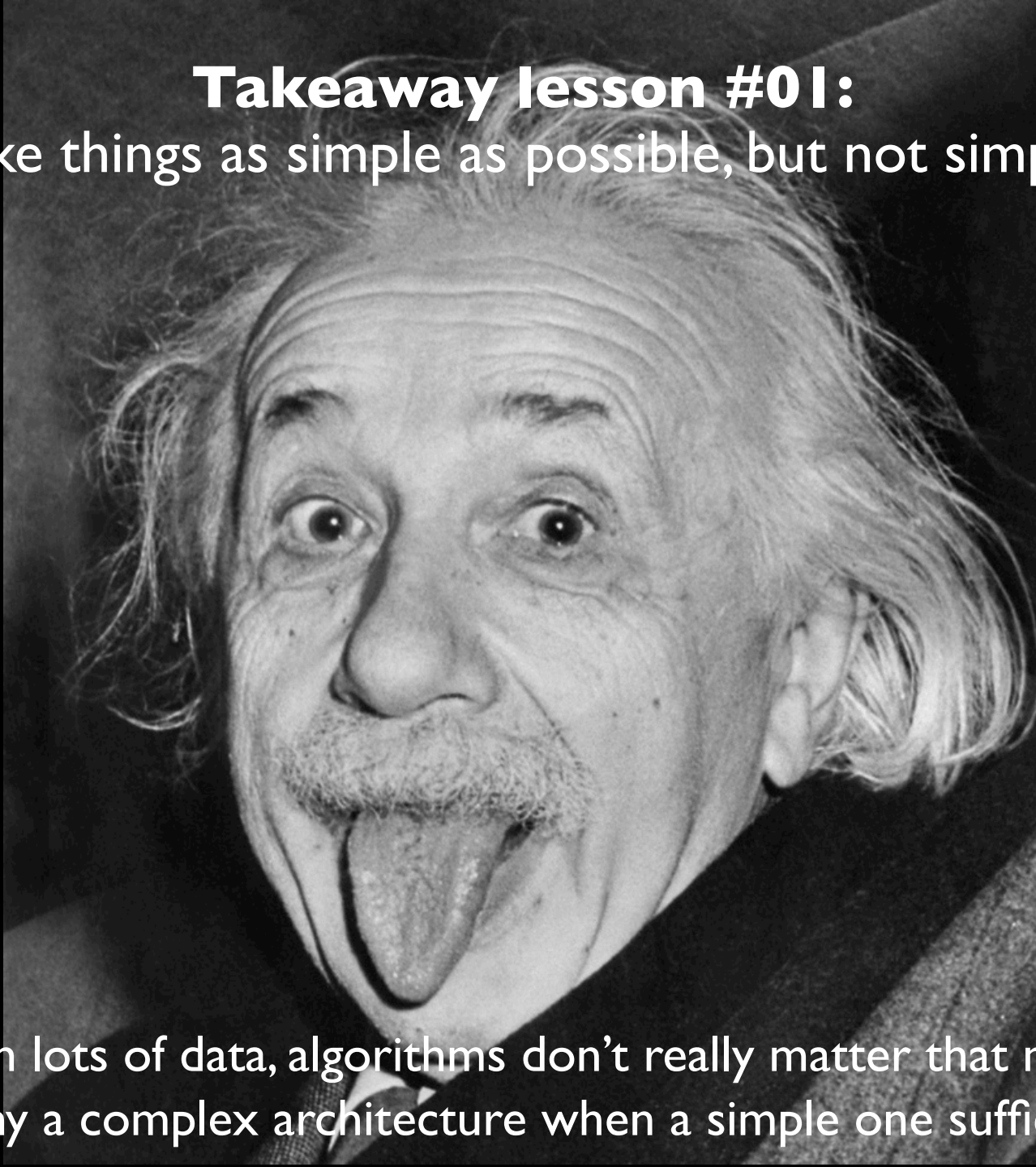
## Performance

End-to-end latency (from edge creation to delivery):  
median 7s, p99 15s



## Takeaway lesson #01:

Make things as simple as possible, but not simpler.



With lots of data, algorithms don't really matter that much  
Why a complex architecture when a simple one suffices?





**Takeaway lesson #10:**  
Constraints aren't always technical.





**Takeaway lesson #11:**  
Plumbing matters. A lot.



# Questions?

“In theory, there is no difference between theory and practice. But, in practice, there is.”

- Jan L.A. van de Snepscheut

