

Big Data Management

Data Streaming

Björn Þór Jónsson

Outline

Lecture Hours (14-16)

- Streaming Layer
 - Requirements
 - Event Queues
 - Streaming Systems
 - Revisit Kappa Architecture
- Practical Concerns
 - Cluster & Project 3
 - Course Evaluation
 - Exam

Exercise Hours (16-18)

- Do one or more of:
 - Work on Project 3
 - Look at Exam
 - Play with Flink

The Three “V”’s



Volume



Velocity



Variety



Veracity



Validity



Viability



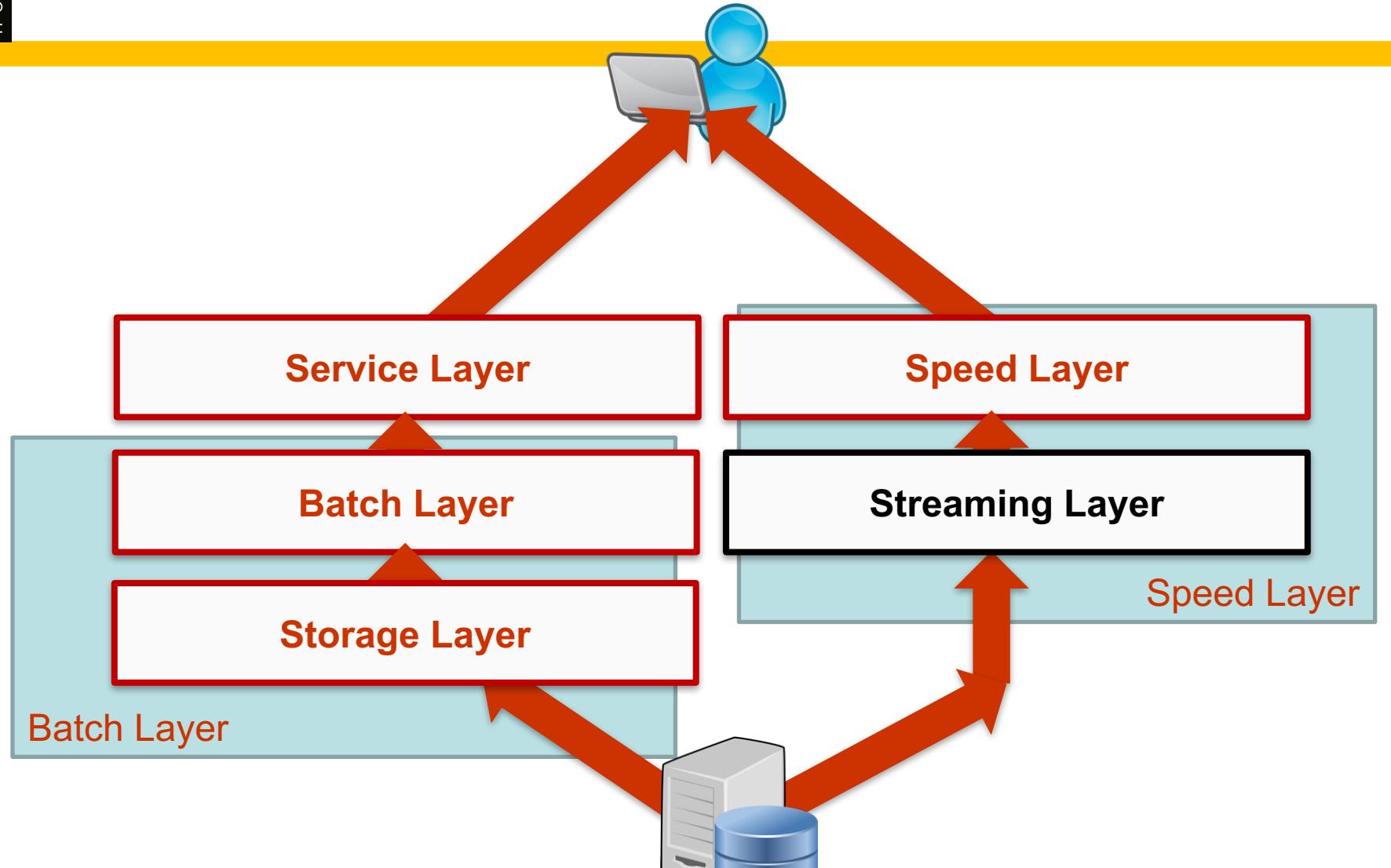
Value



Velocity ?= Real-Time

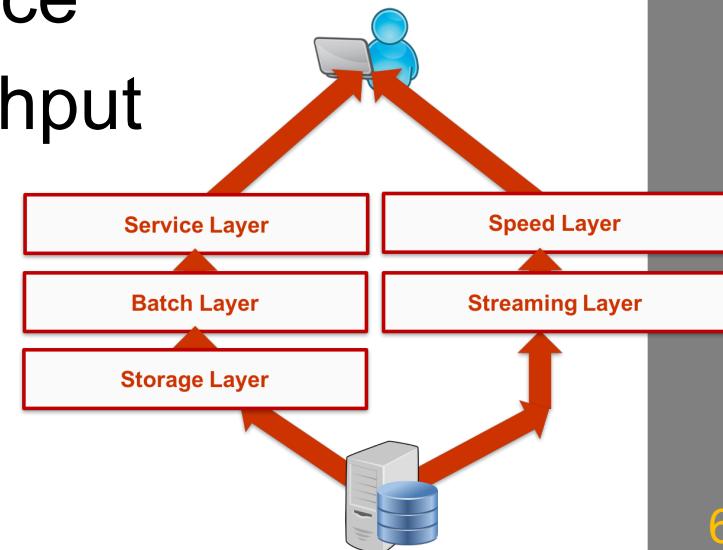
- Projects = Batch Layer + Analytics
 - What are related real-time applications?
- Other domains?

Framework Lambda Architecture



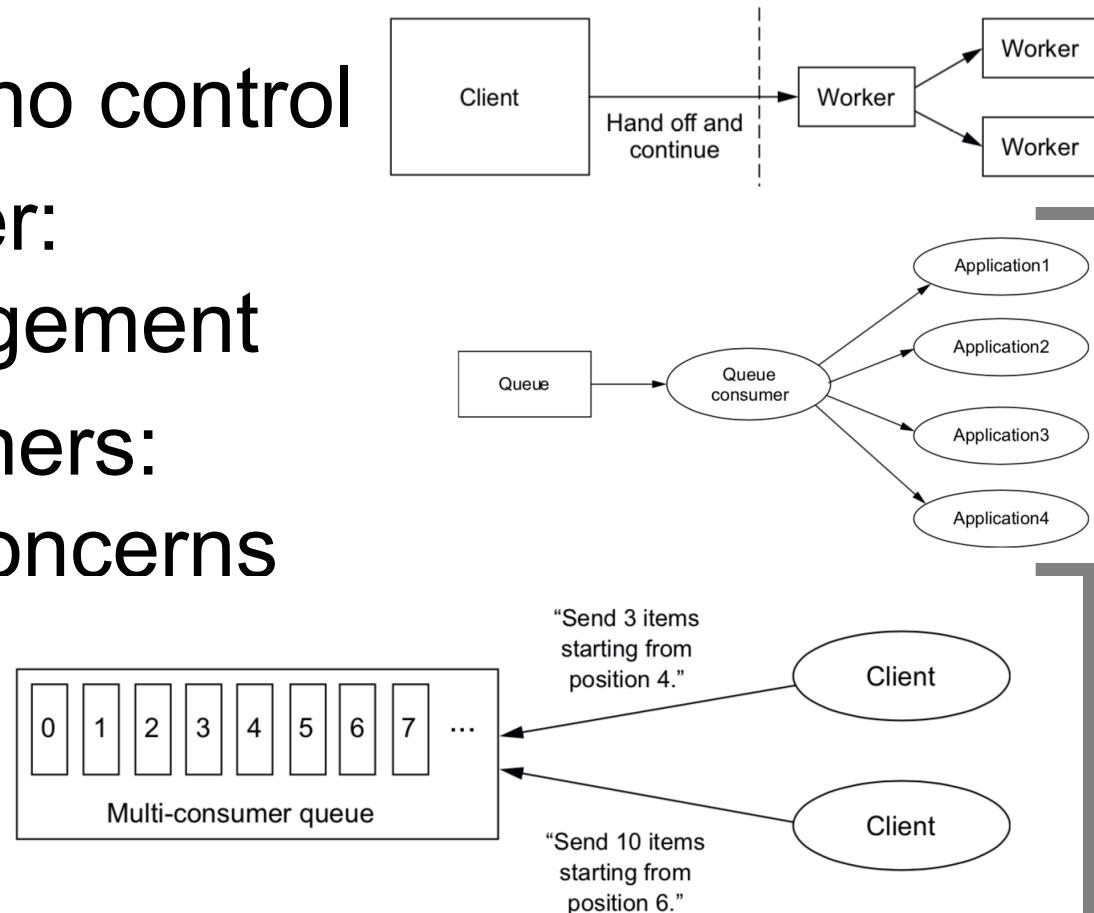
Streaming Layer

- Input: Stream of input records
- Output: Updates to Speed Layer views
- Conflicting Requirements:
 - Scaling out vs Fault tolerance
 - Low latency vs High throughput
- Most difficult part of Lambda Architecture

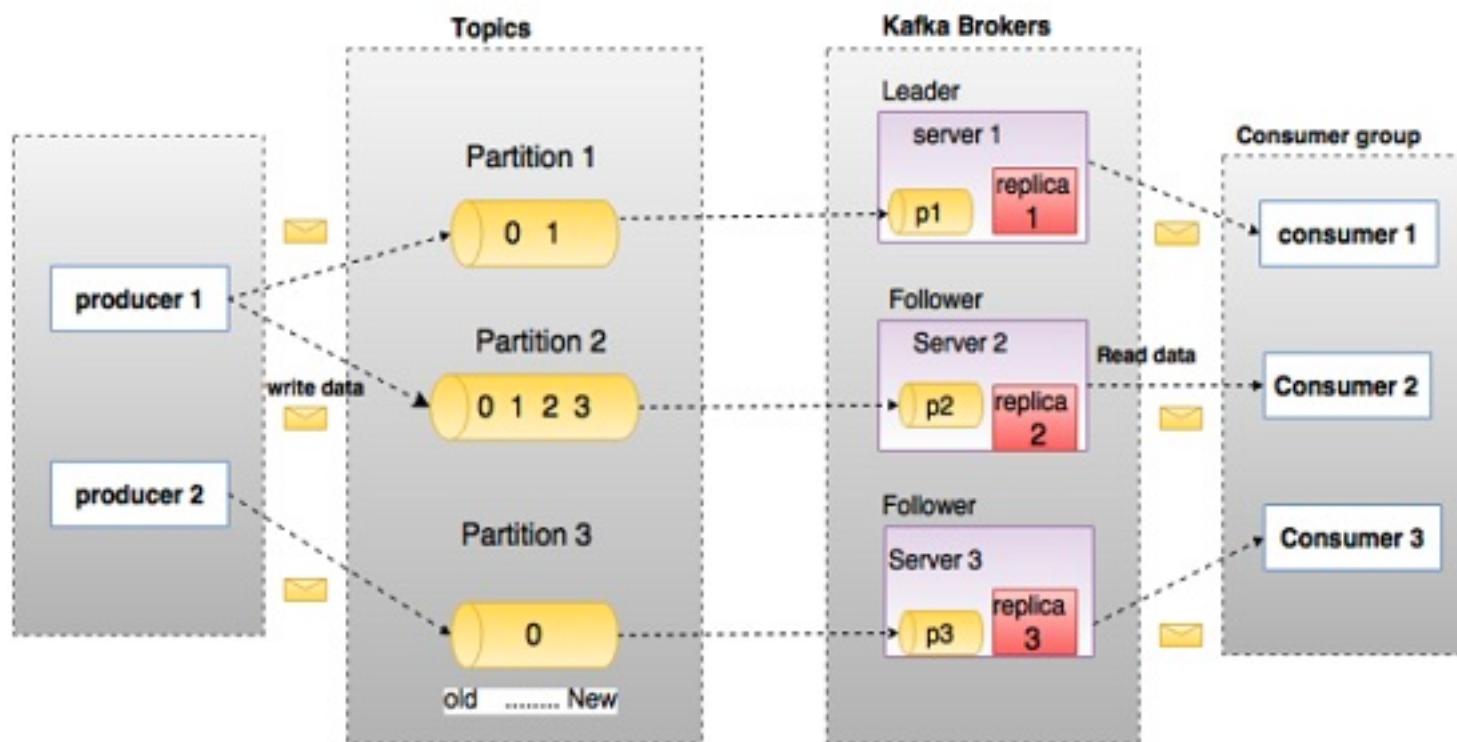


Delivering Data

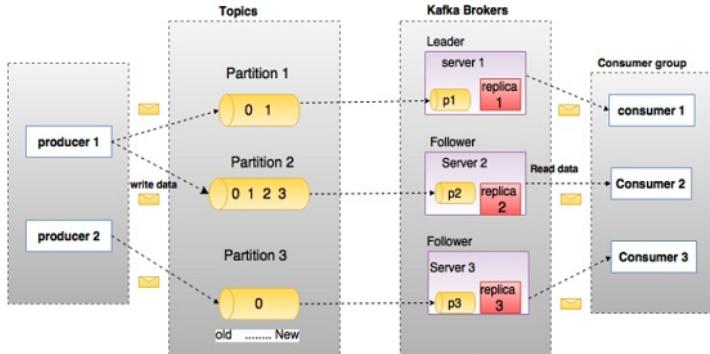
- No queue:
Direct handoff, no control
- Single consumer:
Complex management
- Multiple consumers:
Separation of concerns



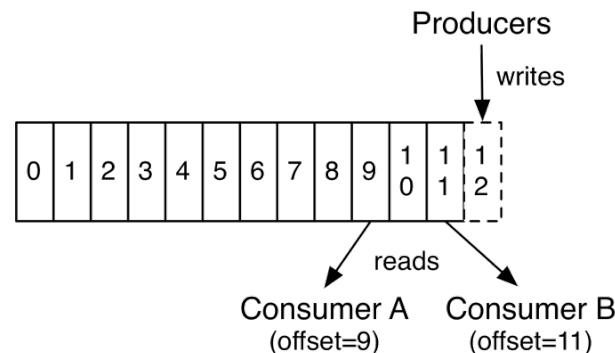
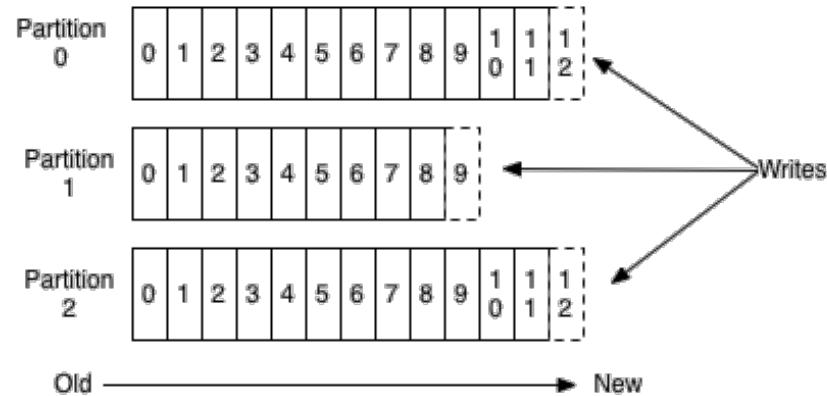
Kafka



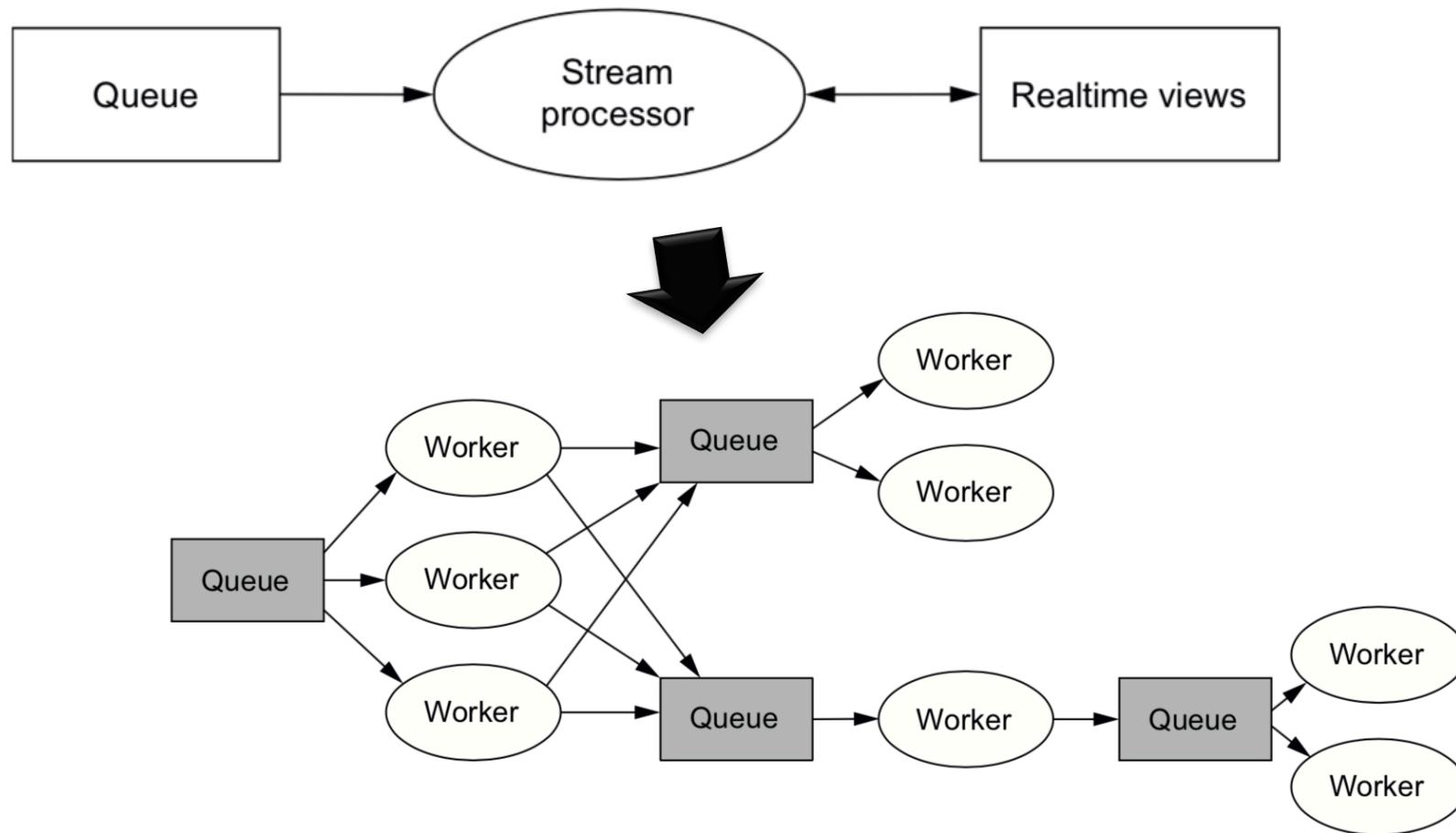
Kafka



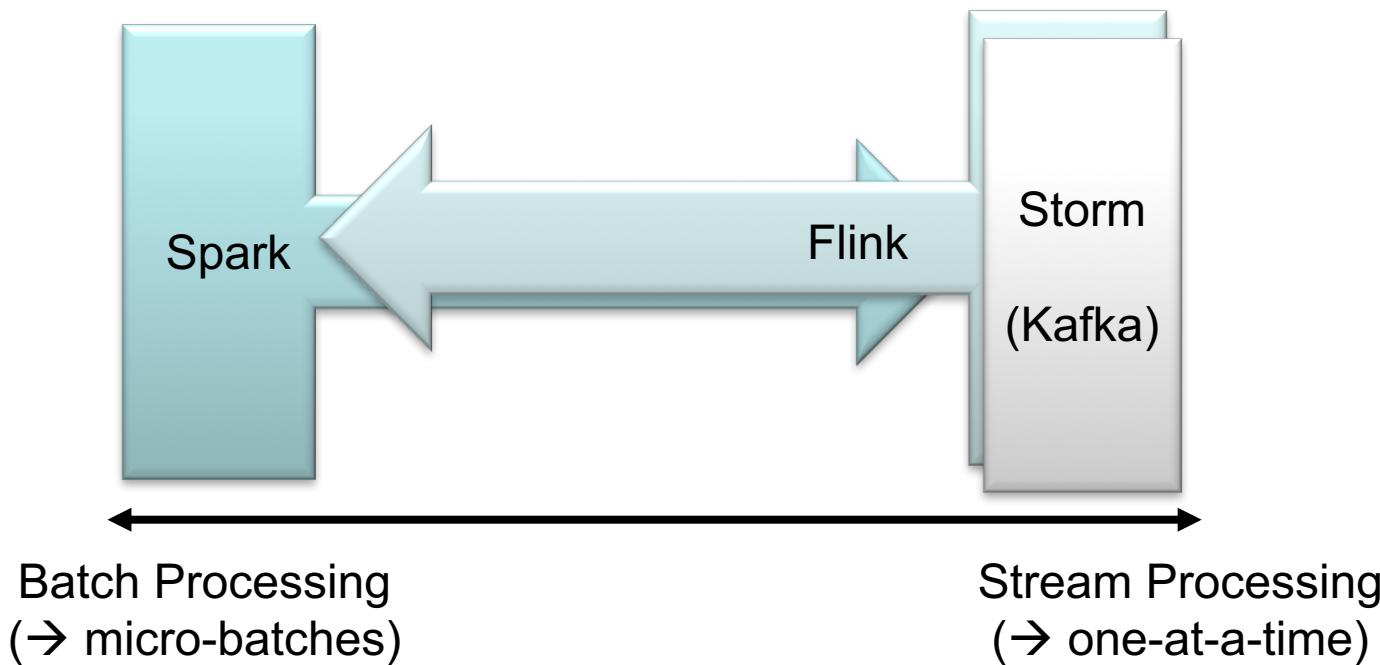
Anatomy of a Topic



Stream Processing



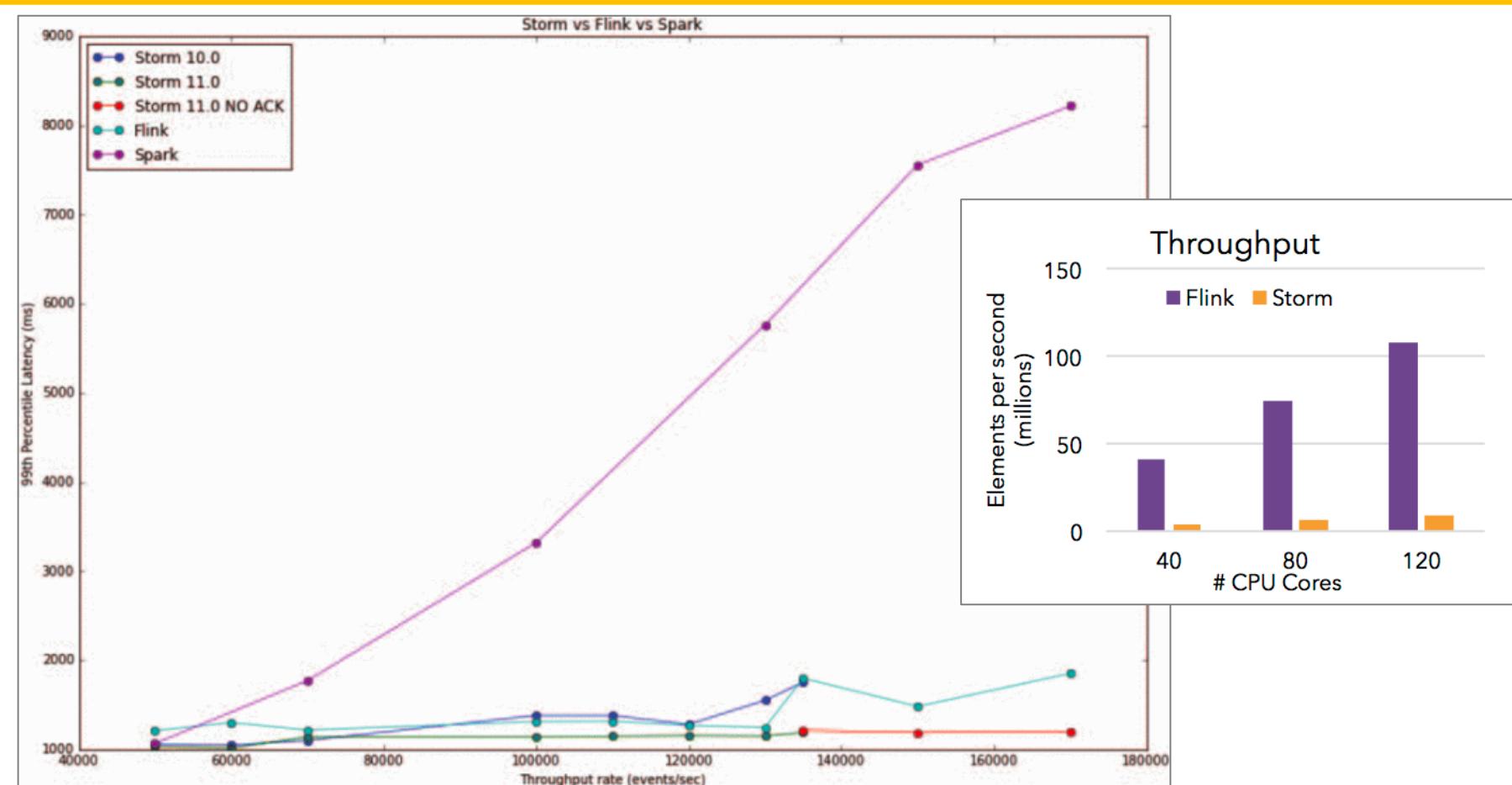
Stream Processing Paradigms



Comparison

	One-at-a-time	Micro-batched
Lower latency	✓	
Higher throughput		✓
At-least-once semantics	✓	✓
Exactly-once semantics	In some cases	✓
Simpler programming model	✓	

Latency over Throughput



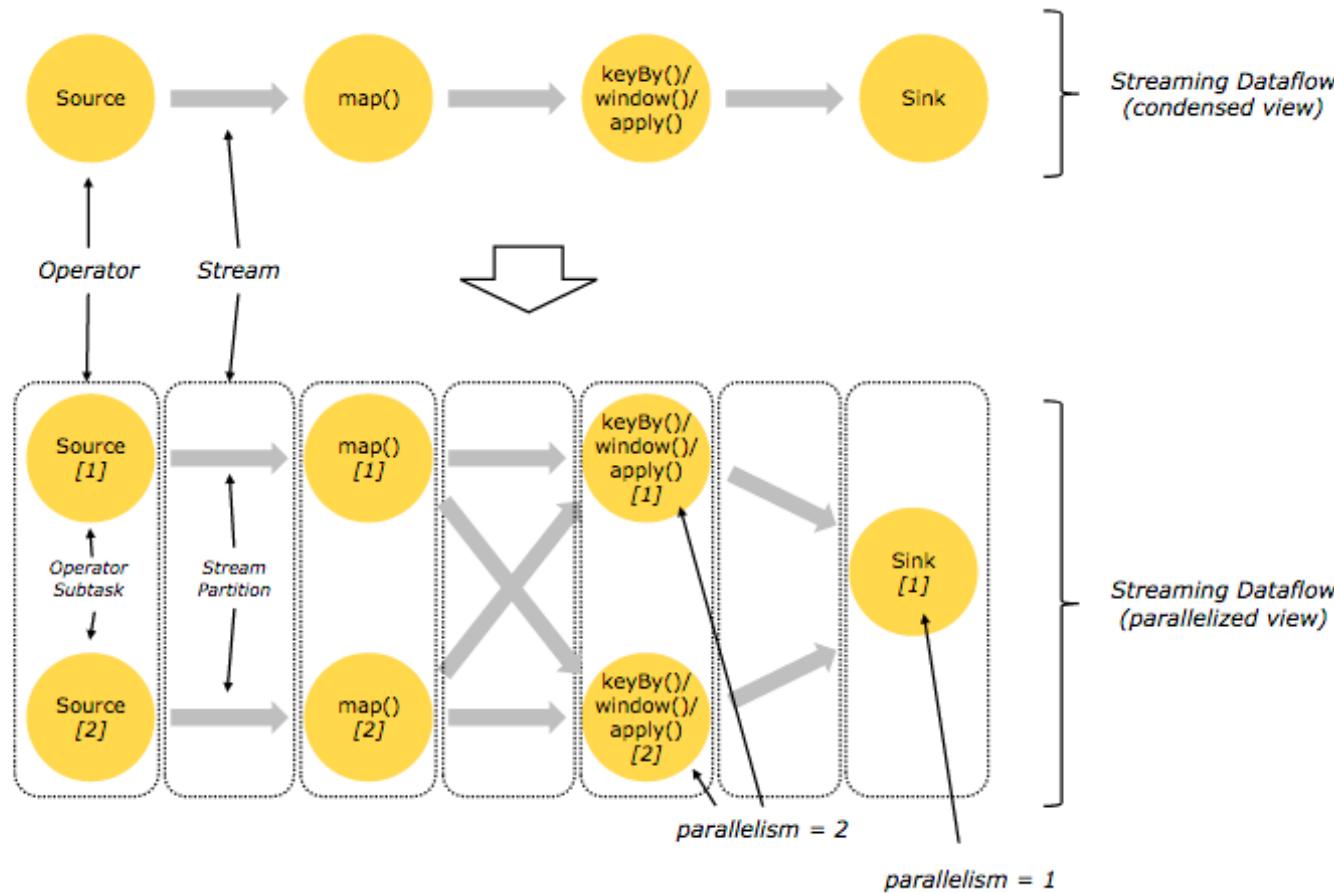
Simple pipeline

<http://ieeexplore.ieee.org/document/7530084/>

Shuffling pipeline

<https://flink.apache.org/introduction.html>

Flink



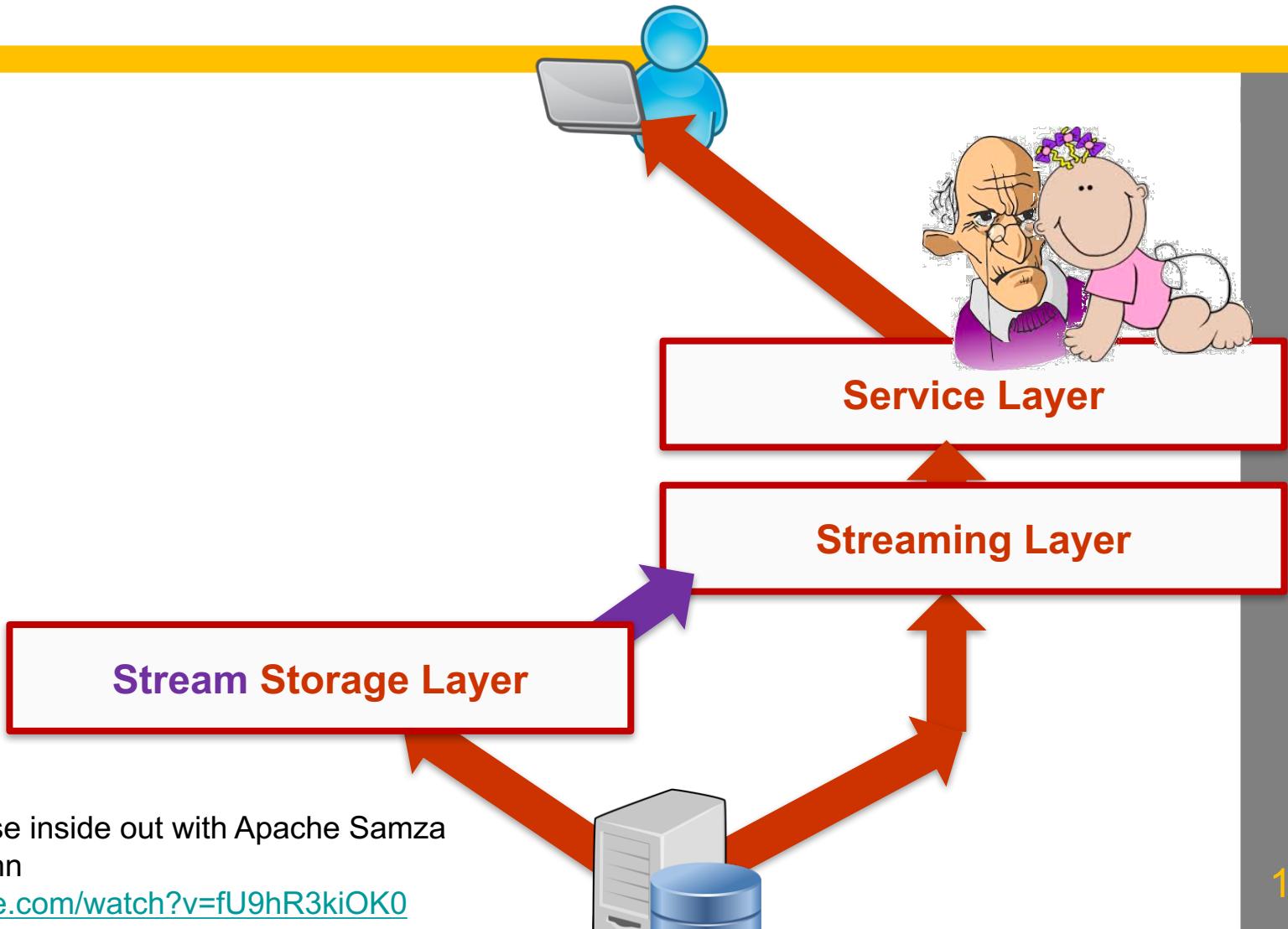
Flink Ecosystem

Deploy	Core			APIs & Libraries
Local Single JVM	Cluster Standalone, YARN	Cloud GCE, EC2	DataStream API Stream Processing	DataSet API Batch Processing
Runtime Distributed Streaming Dataflow				
			FlinkML Machine Learning	Gelly Graph Processing
			Table Relational	Table Relational
			CEP Event Processing	Table Relational

Small Demo: Streaming Word Count

- https://ci.apache.org/projects/flink/flink-docs-release-1.3/quickstart/setup_quickstart.html
 - Finding the right directory to be in was tricky
 - On my Mac (Flink installed using Homebrew):
 - cd /usr/local/Cellar/apache-flink/1.3.2/libexec

Kappa Architecture Revisited



Turning the database inside out with Apache Samza
by Martin Kleppmann
<https://www.youtube.com/watch?v=fU9hR3kiOK0>

Take Away Points

- Streaming Layer
 - Input: Stream of input records
 - Output: Updates to Speed Layer views
 - Conflicting goals:
 - Scaling out vs. Fault tolerance
 - Latency vs. Throughput
- Spark Streaming, Flink, Storm, Kafka, ...
 - Storm + Flink towards Kappa architecture

Outline

Lecture Hours (14-16)

- Streaming Layer
 - Requirements
 - Event Queues
 - Streaming Systems
- Practical concerns
 - Cluster & Project 3
 - Course Evaluation
 - Exam

Exercise Hours (16-18)

- Do one or more of:
 - Play with Flink
 - Work on Project 2
 - Look at exam...



Cluster & Project 3

- New cluster: hq.itu.dk (thanks to Sebastian!)
- Configured (thanks to Omar!)
 - Original password – ssh keys were lost!
 - Please set up again...
 - Memory for ~4 concurrent (full) pipelines
 - Hopefully sufficient due to staggered executions
 - Don't start late!
 - Project 3 data files are available in HDFS under the /sumo-data/ directory
 - `hdfs dfs -ls /sumo-data/`



Cluster & Project 3

- spark-shell / spark-submit:
 - include SPARK_MAJOR_VERSION=2
 - to use cluster add option: --master yarn
- hive:
 - alias for 'hive -hiveconf hive.execution.engine=mr'.

Course Evaluation

Count	Suggestion Category
14	Disjoint between theory and practice, exercises are unfocused / too much installation work, tools different from book
9	Problems with integration with critical course (also 4 positive mentions)
2	Overlap with SoWS
2	Practice exam not provided
2	Piazza
1	Cluster
1	Assignment 1 not useful

Exam

- Available on LearnIT “now”
- Deadline on December 20
 - Directions on LearnIT!

Exam Outline

This exam assignment is a written work based on a) the learning outcomes of the course and b) the three projects given throughout the semester. You should submit this written work individually. You can reuse work that was done as a group on these three projects, but *the text of the written work must be your own.*

This exam is composed of four questions. Your answers should be insightful, clear and concise. The document should be in a question-answer format, and must fit within (maximum) 10 normal pages. The written work must be submitted via LearnIT, where you will find all instructions related to submission. The deadline is managed by study administration, and cannot be changed in any way.

Course Outline

Calendar		Course	Reading	Lecture Topics	Exercise Topics
Week	Date	Week	Materials		
35	29.08	1	Papers	TECH: Relational model; DB systems; Big data motivation	
36	05.09	2	Papers	CRIT: Raw vs cooked data; Data valence; Personal data; Data walk	CRIT: Data walk
37	12.09	3	Papers	CRIT: Social aspects; Immutability and certainty; Perils of interpretation and bias	CRIT: Bias Project 1
38	19.09	4	Ch 1-3	TECH: Lambda Architecture: Overview; Layers and roles; Data properties	TECH: Intro to Spark Shell
39	26.09	5	M&W: Ch 4-5	Batch Layer: Storage; High volume; DHT	Project 1
40	03.10	6	Papers	Intro to data cleaning	CRIT: Andrew Clement guest lecture Project 2
41	10.10	7	M&W: Ch 6-9	Batch Layer: Batch processing; Spark	Spark; HDFS; Cluster
42					
43	24.10	8	M&W: Ch 10-13	Serving/Speed Layer: NoSQL; CAP theorem	Drill
44	31.10	9	M&W: Ch 18 + Papers: TBA	Data pipeline management; Metadata	Project 2
45	07.11	10	Papers	Data quality exercise	Project 3
46	14.11	11			
47	21.11	12	M&W: Ch 14-17	Speed Layer: Data streaming	Streaming / Project / Exam
48	28.11	13			Project 3
49	05.12	14			Portfolio Work
50	12.12	15			Portfolio Work
Exam report due 20.12					

No Class

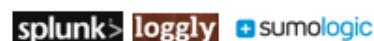
Joint Class with DIM Students

Big Data Landscape

Vertical Apps



Log Data Apps



Ad/Media Apps



Data As A Service



Business Intelligence



Microsoft Business Intelligence



Analytics and Visualization



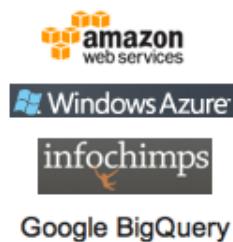
Analytics
Infrastructure



Operational Infrastructure



Infrastructure As A Service



Structured Databases



Technologies



THE BIG DATA LANDSCAPE

JANUARY 2014

Apps

Vertical

newton ellucian practicefactor
SurveyMonkey predictive policing

Operational Intelligence

splunk new relic AppDynamics
sumologic VITRIA

Data As A Service

factual FICO kaggle INRIX
GNIP apigee Placed LOGATE
DATAFIFT TOPSY LexisNexis

Consumer

Google amazon
@WalmartLabs ebay NETFLIX in

Ad/Media

METAMARKETS collective
rocketfuel FLURRY
collective DataXO
Media Science TURN inidx
bloomreach

Business Intelligence

ORACLE Hyperion
SAP Business Objects RJMetrics
Microsoft Business Intelligence
IBM JASPERREPORTS birst
openbach MicroStrategy
Autonomy bime
Chart.io GoodData
ATTIVIO Recorded Future

Analytics and Visualization

Tableau QlikView
Palantir OPERA TRIFACTA
TERADATA ASTER knoema
centrifuge SAS TIBCO
panopticon Real-Time Visual Data Analysis
Datameer IDATA plafoura
platfora ClearStory CIRRO
alteryx visual.ly uFORA
metatlayer Atigeo Alpine AVATA
SiSense

Infrastructure

Analytics

cloudera Hortonworks
MAPR Cohera HADAPT
Pivotal INFOBRIGHT
NETEZZA VERTICA
exasol kognitio

Operational

COUCHBASE mongoDB
KEROSPIKE splice
DATASTAX VoltDB
TERRACOTTA INFORMATICA
MarkLogic

As A Service

Dubale amazon
Windows Azure MORTAR
CSC Google BigQuery

Structured DB

ORACLE MySQL
SQL Server PostgreSQL
IBM DB2 SYBASE
memsql TERADATA



Technologies



BIG DATA LANDSCAPE 2017



Last updated 4/5/2017

© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark (@firstmarkcap)

mattturck.com/bigdata2017

FIRSTMARK
EARLY STAGE VENTURE CAPITAL