# Google MapReduce

*DS 5110: Big Data Systems (Spring 2023)*

Lecture 3b

Yue Cheng

UNIVERSITY *of* VIRGINIA

# Applications

| Batch | SQL | ETL | Machine learning | Emerging apps? |

# Scalable computing engines

# Scalable storage systems

Datacenter infrastructure

# The big picture (motivation)

- Datasets are <span style="color:red">too big</span> to process using a single computer

# The big picture (motivation)

- Datasets are <span style="color:red">too big</span> to process using a single computer

- Good parallel processing engines are <span style="color:red">rare (back then in the late 90s)</span>

# The big picture (motivation)

- Datasets are too big to process using a single computer

- Good parallel processing engines are rare (back then in the late 90s)

- Want a parallel processing framework that:
  - is **general** (works for many problems)
  - is **easy to use** (no locks, no need to explicitly handle communication, no race conditions)
  - can **automatically parallelize** tasks
  - can **automatically handle machine failures**

# Context (Google circa 2000)

- Starting to deal with massive datasets

- But also addicted to cheap, unreliable hardware
  - Young company, expensive hardware not practical

- Only a few expert programmers can write distributed programs to process them
  - Scale so large jobs can complete before failures

# Context (Google circa 2000)

- Starting to deal with massive datasets
- But also addicted to cheap, unreliable hardware
  - Young company, expensive hardware not practical
- Only a few expert programmers can write distributed programs to process them
  - Scale so large jobs can complete before failures

- **Key question:** how can every Google engineer be imbued with the ability to write parallel, scalable, distributed, fault-tolerant code?

- **Solution:** abstract out the redundant parts

- **Restriction:** relies on job semantics, so restricts which problems it works for

# Application: Word Count

```
cat data.txt
    | tr –s '[[:punct:][:space:]]' '\n'
    | sort | uniq -c
```

```
SELECT count(word), word FROM data
    GROUP BY word
```

# Deal with multiple files?

# Deal with multiple files?

1. Compute word counts from individual files

# Deal with multiple files?

1. Compute word counts from individual files

2. Then merge intermediate output

# Deal with multiple files?

1. Compute word counts from individual files

2. Then merge intermediate output

3. Compute word count on merged outputs

# What if the data is too big to fit in one computer?

# What if the data is too big to fit in one computer?

1. In parallel, send to worker:
   - Compute word counts from individual files
   - Collect results, wait until all finished

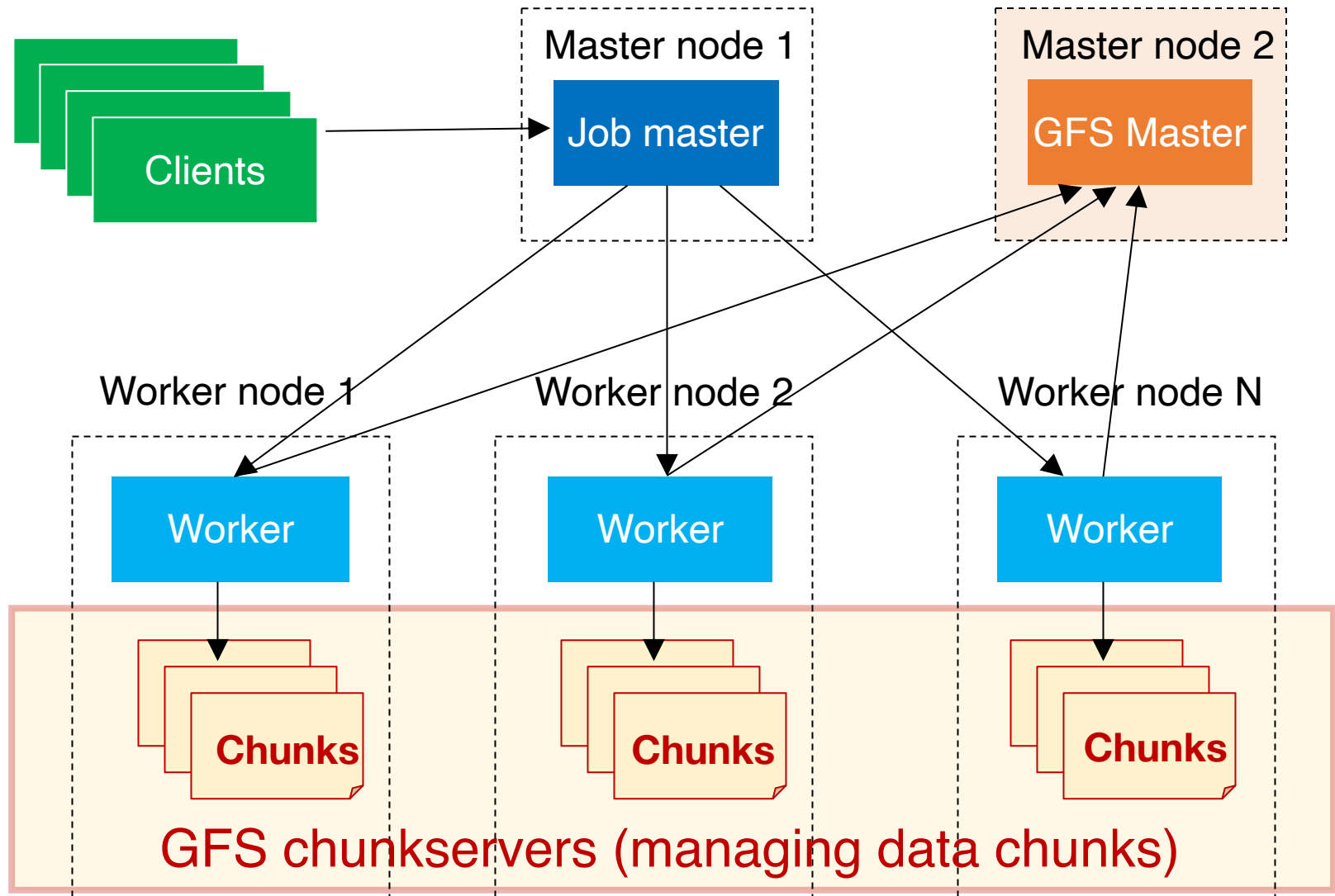# What if the data is too big to fit in one computer?

1. In parallel, send to worker:
   - Compute word counts from individual files
   - Collect results, wait until all finished

2. Then merge intermediate output

# What if the data is too big to fit in one computer?

1. In parallel, send to worker:
   - Compute word counts from individual files
   - Collect results, wait until all finished

2. Then merge intermediate output

3. Compute word count on merged intermediates

# MapReduce+GFS: Put everything together
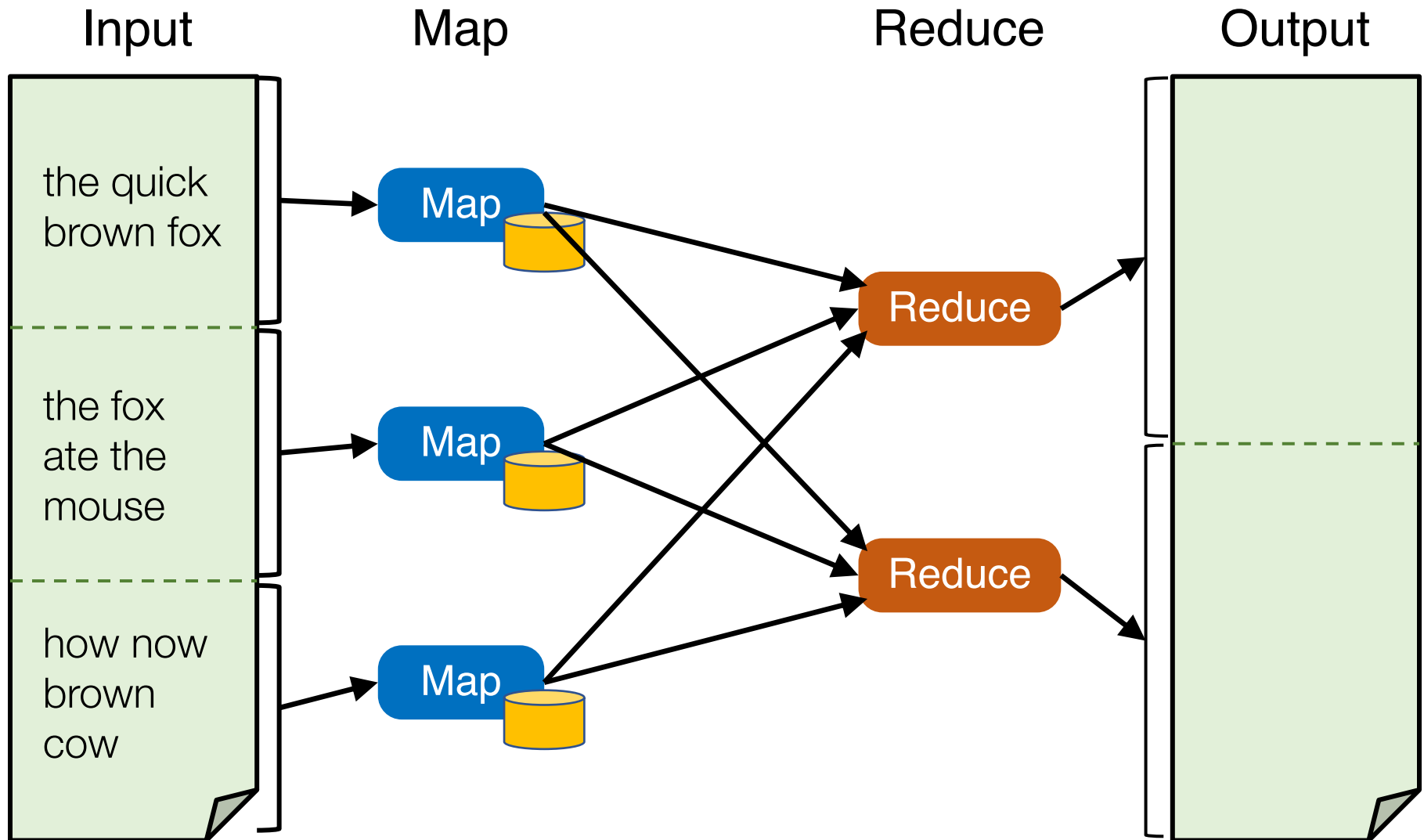
# MapReduce: Programming interface

- `map(`<span style="color:red">`k1, v1`</span>`)` → `list(`<span style="color:blue">`k2, v2`</span>`)`
  - Apply function to (`k1, v1`) pair and produce set of intermediate pairs (`k2, v2`)



- `reduce(`<span style="color:blue">`k2`</span>`, list(`<span style="color:blue">`v2`</span>`))` → `list(`<span style="color:green">`k3, v3`</span>`)`
  - Apply aggregation (reduce) function to values
  - Output results

# MapReduce: Word Count
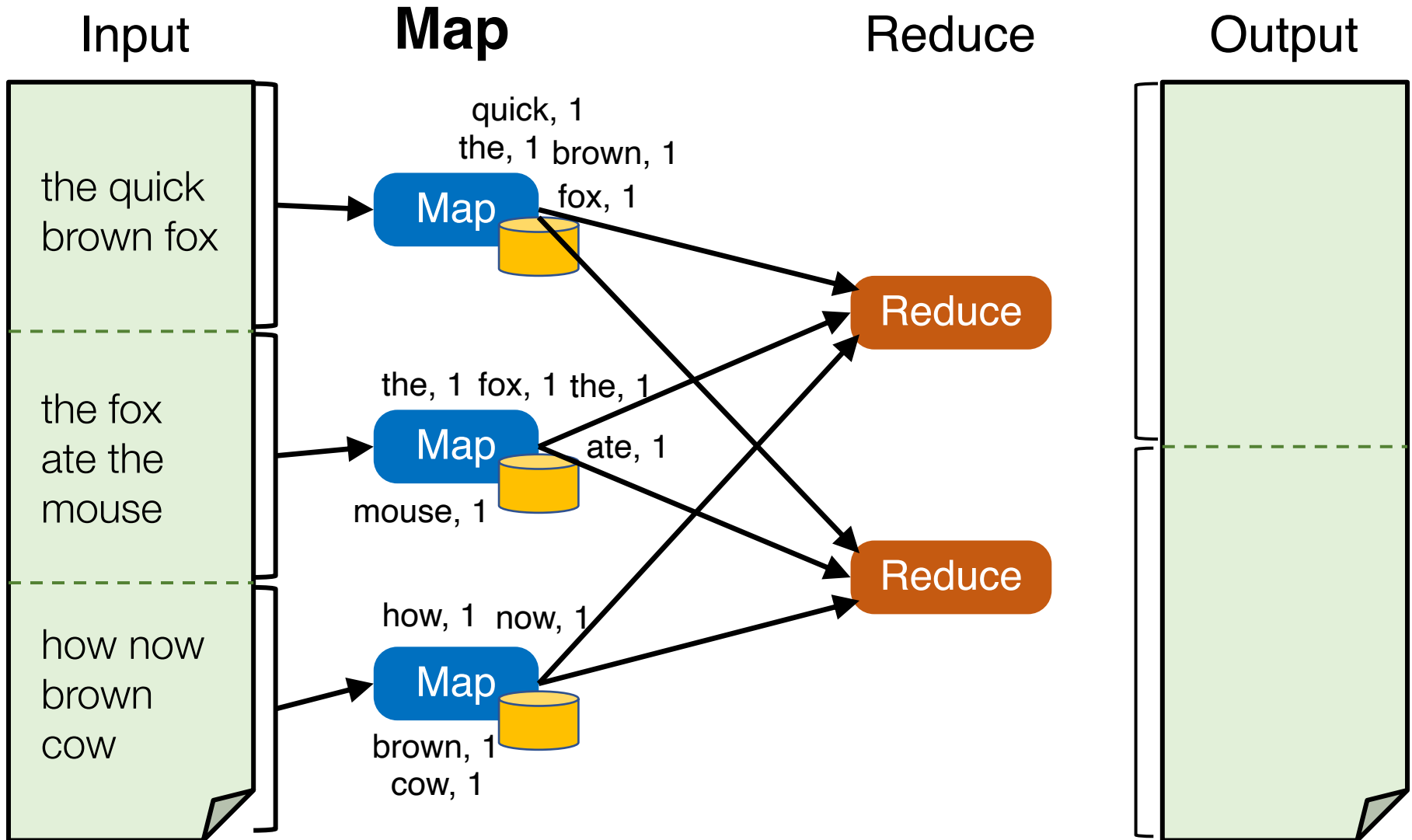
```
map(key, value):
    for each word w in value:
        EmitIntermediate(w, "1");

reduce(key, values):
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));
```

# Word Count execution

Input        Map        Reduce        Output

the quick brown fox

the fox ate the mouse

how now brown cow

Map

Map

Map

Reduce

Reduce

# Word Count execution

Input  **Map**  Reduce  Output

the quick
brown fox

the fox
ate the
mouse

how now
brown
cow

quick, 1
the, 1  brown, 1
fox, 1

Map

the, 1  fox, 1  the, 1
ate, 1

Map

mouse, 1

how, 1  now, 1

Map

brown, 1
cow, 1

Reduce

Reduce

# Word Count execution

| Input | Map | Shuffle & Sort | Reduce | Output |
|-------|-----|----------------|--------|--------|

the quick brown fox

**Map**

the fox ate the mouse

**Map**

how now brown cow

**Map**

the, 1
the, 1
brown, 1
fox, 1
how, 1
now, 1
brown, 1
the, 1
fox, 1

**Reduce**

quick, 1
ate, 1
mouse, 1
cow, 1

**Reduce**

# Word Count execution

| Input | Map | Shuffle & Sort | **Reduce** | Output |

Input

Map

Shuffle & Sort

**Reduce**

Output

the quick brown fox

the fox ate the mouse

how now brown cow

Map

Map

Map

the, 1
the, 1
brown, 1
fox, 1
how, 1
now, 1
brown, 1
the, 1    fox, 1

Reduce

quick, 1    Reduce
ate, 1
mouse, 1
cow, 1

brown, 2
fox, 2
how, 1
now, 1
the, 3

ate, 1
cow, 1
mouse, 1
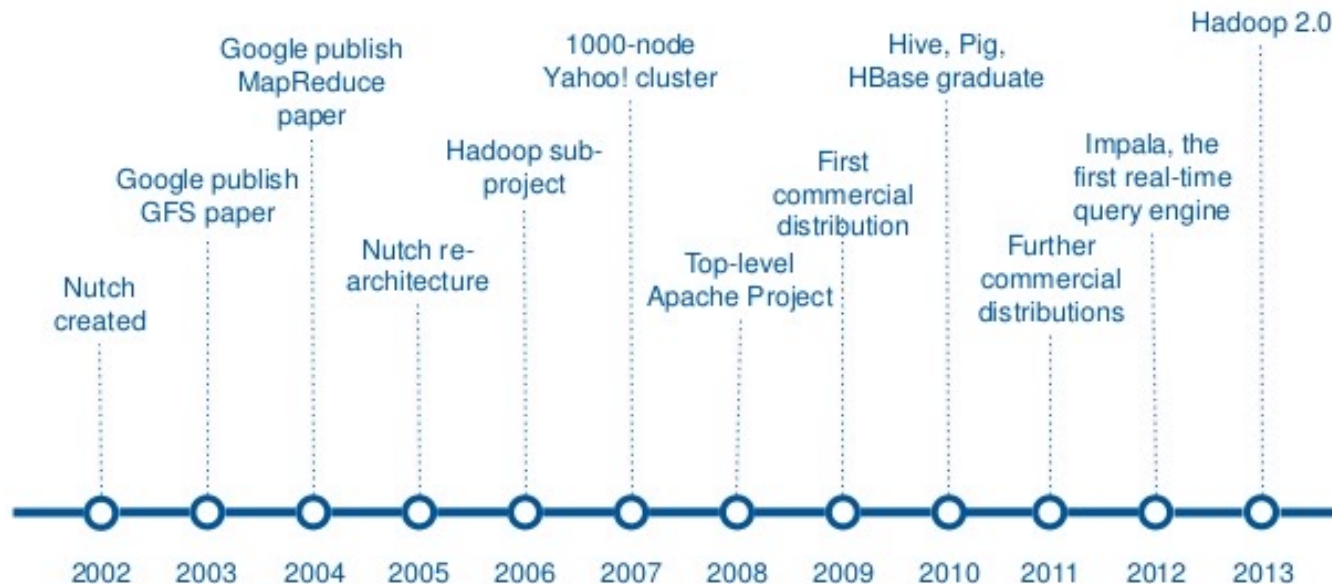quick, 1

# MapReduce data flows in paper

# How it started: Apache Hadoop

- An open-source implementation of Google's MapReduce framework
  - Hadoop MapReduce atop Hadoop Distributed File System (HDFS)

APACHE
*hadoop*

A Brief History of Hadoop

Google publish MapReduce paper

1000-node Yahoo! cluster

Hive, Pig, HBase graduate

Hadoop 2.0

Google publish GFS paper

Hadoop sub-project

First commercial distribution

Impala, the first real-time query engine

Nutch re-architecture

Top-level Apache Project

Further commercial distributions

Nutch created

2002  2003  2004  2005  2006  2007  2008  2009  2010  2011  2012  2013

4

# How it's going …

DATA & AI LANDSCAPE 2019

## INFRASTRUCTURE

**HADOOP ON-PREMISE** — cloudera, Hortonworks, MAPR, Pivotal, IBM InfoSphere, jethro

**HADOOP IN THE CLOUD** — AWS, Microsoft Azure, Google Cloud, SAP Cloud Platform, IBM InfoSphere BigInsights, arm, Qubole, CAZENA

**STREAMING / IN-MEMORY** — Amazon Kinesis, databricks, SAP Cloud Platform, ORACLE Coherence, confluent, Striim, hazelcast, GridGain, GIGASPACES, WallarooLabs, FASTDATA.io, kx

**NoSQL DATABASES** — Google Cloud, AWS, ORACLE, Microsoft Azure, mongoDB, MarkLogic, Couchbase, DATASTAX, redislabs, AEROSPIKE, ArangoDB, SCYLLA

**NewSQL DATABASES** — SAP, Clustrix, Pivotal, NuoDB, MariaDB, MEMSQL, influxdata, Cockroach Labs, VoltDB, splice, imply, paradigm4, TiDB

**GRAPH DBs** — neo4j, Amazon Neptune, IBM, ORACLE, OrientDB, Kognitio, InfiniteGraph, Objectivity

**MPP DBs** — TERADATA, VERTICA, IBM Data Warehouse Systems, Action, Exasol, dremio, Yellowbrick

**CLOUD EDW** — AWS, Google Cloud, Microsoft Azure, Pivotal, snowflake, Infoworks

**SERVERLESS** — PULSAR, nuclio, Pivotal Function Service

**DATA TRANSFORMATION** — talend, pentaho, alteryx, TRIFACTA, tamr, Paxata, StreamSets, UNIFI

**DATA INTEGRATION** — SAP Data Services, Informatica, MuleSoft, TEALIUM, snaplogic, enigma, Qlik Data Catalyst, Segment, ATTUNITY, xplenty, ZALONI, import.io, Infoworks, Fivetran, SNOWPLOW, MATILLION

**DATA GOVERNANCE** — Informatica, IBM, collibra, McAfee Skyhigh Security Cloud, Alation, IMMUTA, OKERA, unravel, Numerify, zenbase, OptsRamp, MAGNITUDE, MANTA, data.world, Waterline Data, ScienceLogic

**MGMT / MONITORING** — AWS, New Relic, actifio, rubrik, AppDynamics, dynatrace, WAVEFRONT, SignalFx, Moogsoft, pagerduty, splunk, SCALYR, VEEAM, EPIMIX

**COMPUTER VISION** — Microsoft Azure, Amazon Rekognition, clarifai, deepomatic, EVER AI, neurala, twentybn, UBIQUITY.ai, YITU, trax, BLUE VISION, Synthesia, DataGrid

**STORAGE** — AWS, Google Cloud, Microsoft Azure, IBM Storage, PURE STORAGE, ALLUXIO, wasabi, nimble storage, Qumulo, panasas, COHESITY

**CLUSTER SVCS** — IBM, AWS, Amazon ECS, Amazon EKS, Cluster Server, MESOSPHERE, packet, BlueData, Bright Computing, CYCLECLOUD

**DATA GENERATION & LABELLING** — amazon mechanical turk, upwork, appen, Verta.ai, datmo, Labelbox, scale, HIVE, Mighty AI, REVEL.IE, LIONBRIDGE, figure eight, fiddler

**AI OPS** — ALGORITHMIA, SPELL, comet, SQREAM, datatron, brytlyt, PG-Strom, BLAZINGDB, Mavidus, habana, WAVE, SambaNova, CERNAMI, PAID, Determined AI, FLOYDHUB

**GPU DBs & CLOUD** — kinetica, omni sci, Cerebras

**HARDWARE** — Google TPU, arm, intel AI, NVIDIA, IBM Power Systems, GRAPHCORE, MYTHIC

## ANALYTICS & MACHINE INTELLIGENCE

**DATA ANALYST PLATFORMS** — Microsoft, pentaho, alteryx, Digital Reasoning, guavus, AYASDI, ATTIVIO, Datameer, incorta, interana, MODE, ENDOR, ASCEND.IO, sisu, switchboard, Starburst

**DATA SCIENCE PLATFORMS** — IBM, databricks, data iku, DOMINO, rapidminer, TIBCO, ANACONDA, SAS, Altair, H2O.ai, KNIME, MathWorks

**BI PLATFORMS** — looker, einstein analytics, AWS, DOMO, ARCADIA DATA, ThoughtSpot, Qlik, SAP Lumira, ATSCALE, SSENSE, GoodData, Information Builders, birst, MicroStrategy, Keen IO

**VISUALIZATION** — tableau, Microsoft Power BI, Google Cloud, Periscope Data, plotly, zepl, GEOMDATA, ViSENZE, ELEMENT AI, deepsense.ai, CHARTIO, Toucan Toco

**MACHINE LEARNING** — Azure Machine Learning, Amazon SageMaker, Google Cloud, AutoML, Vision, H2O.ai, DataRobot, gamalon

**SEARCH** — elasticsearch, ORACLE, ENDECA, algolia, Coveo, Lucidworks, ATTIVIO, swiftype, EXALEAD, alphasense, MAANA, omni:us, SINEQUA

**LOG ANALYTICS** — splunk, sumologic, NETBASE, synthesio, tracx, kibana, TIMBER, logz.io

**SOCIAL ANALYTICS** — Hootsuite, sprinklr, NETBASE, synthesio, tracx, simplereach, bitly, SimilarWeb

**WEB / MOBILE / COMMERCE ANALYTICS** — Google Analytics, mixpanel, AMPLITUDE, Airtable, RESCI, SIGOPT, granify, custora

**HORIZONTAL AI** — IBM Watson, Cortana, Face++, vicarious, sentient, Voyager, semanticmachines, CognitiveScale, PROPHESEE, PETUUM, Cinnamon, curious AI, OSARO, Fortress

**SPEECH & NLP** — Google Cloud, twilio, amazon alexa, Amazon Translate, Mobvoi, EigenTechnologies, SoundHound Inc., PRIMER, MindMeld, NUANCE, cogito, snips, SMARTLING, Unbabel, PolyAI

## APPLICATIONS – ENTERPRISE

**SALES** — CHORUS, INSIDESALES.COM, peopleai, conversica, clari, aviso, tact.ai, TROOPS, fuse machines, Clearbit

**MARKETING - B2B** — RADIUS, App Annie, SendGrid, EVERSTRING, Lattice, MINTIGO, sense, contentsquare, TEALIUM, mparticle, Amplero, amperity, QUANTIFIND, ENGAGIO, Simon, Lytics, PERSADO, KNOTCH, mrp

**MARKETING - B2C** — zeta, bloomreach, BLUECORE, braze, ACTIONIQ

**CUSTOMER EXPERIENCE / SERVICE** — qualtrics, MEDALLIA, SurveyMonkey, User Testing, CLARABRIDGE, zendesk, Kustomer, freshdesk, INTERCOM, Drift, LIVEPERSON, Gainsight, pendo, HEAP, Amplitude, Watson Assistant, Dialogflow, DigitalGenius, ASAPP, ada, AUTOMAT, afiniti, Catflirt, netomi, clara, talla, Kasisto

**ENTERPRISE PRODUCTIVITY** — slack, ORACLE, GURU, lumiata, DIFFBOT

**HUMAN CAPITAL** — Hire Vue, pymetrics, hiQ, GIGSTER, mya, Allyo, textio, Wade & Wendy, Stella, entelo, RESTLESS BANDIT

**LEGAL** — RAVEL, Seal, Everlaw, Disco, kira, JUDICATA, BREVIA, IRONCLAD, LIONMIGHT, PREMONITION, ROSS, casetext

**REGTECH & COMPLIANCE** — text IQ, Comply Advantage, TRADESHIFT, CROSSBEAM, ONFIDO, DATA REPUBLIC, beamery

**FINANCE** — anaplan, SAP/S4 HANA, VIDADO, AppZen, mineraltree, SCALEFACTOR, botkeeper

**BACK OFFICE AUTOMATION & RPA** — UiPath, HyperScience, ANTWORKS, blueprism, WorkFusion, workato, Catalytic, KRYON, ALKYMI

**SECURITY** — TANIUM, CYLANCE, zscaler, StackPath, illumio, CODE42, CipherCloud, DARKTRACE, ANOMALI, ThreatMetrix, pindrop, exabeam, SIGNIFYD, SentinelOne, SecurityScorecard, SOCURE, Vade Secure, bitglass, BioCatch, Recorded Future, feedzai, Cyber?, BITSIGHT, AREA 1 SECURITY, sparkcognition, IronNet Cybersecurity, FORTER, riskrecon, JASK, BLUE HEXAGON, Semmle, OBSIDIAN, AXONIUS, SHIELD AI, Armorblox

## APPLICATIONS – INDUSTRY

**ADVERTISING** — AppNexus, MediaMath, Integral, criteo, X+AD, ORACLE MOAT, theTradeDesk, dstillery, TAPAD, dataxu, gumgum, Appier, tremor, yieldmo

**EDUCATION** — Liulishuo, KNEWTON, Clever, declara, PANORAMA, knowre, gradescope

**REAL ESTATE** — REDFIN, Opendoor, VTS, GEOPHY, reonomy, COMPSTAK, SPACEMAKER, SKYLINE, STREETLINE, OpenDataSoft

**GOV'T** — OPENGOV, mark43, FiscalNote, LiveStories, Passport, SmartProcure, PREMISE, PAGAYA

**INTELLIGENCE** — Palantir, Datamin?, Quantopian, Quid, PRIMER, RavenPack, cognigo

**FINANCE - INVESTING** — NUMERAI, iSENTIUM, ALGORIZ, TrueAccord, MoneyLion, aire

**FINANCE - LENDING** — ondeck, Affirm, TALA, finova, Upstart, AURA, BLEARBANC, upgrade, 100Credit, WeLab, cignifi

**INSURANCE** — Metromile, Lemonade, CYENCE, Hippo, Shift Technology, ROOT, zesty.ai, CAPE

**HEALTHCARE** — flatiron, Clover, XYRUUS, HealthTap, METABIOTA, Ginger.io, Glow, babylon, 3DMed, zebra, ovia, TEMPUS, patientslikeme, AiCure, insitro, notable, komodo health, RECURSION, prognos, notable, enlitic, BlackThorn, CITRINE, Qventus, ARTERYS, IMAGEN, Phosphorus, CLOUD MEDx, PAIGE, DATAVANT, twoXAR, deep genomics, innovaccer, DIAGNOSS PROJECT, LeanTaaS

**LIFE SCIENCES** — StandMe, color, verily, WuXiNextCODE, Clear Labs, freenome, NANOPORE, DNAnexus, NIO, DNAnexus, SOPHIA GENETICS, OWKIN

**TRANSPORTATION** — UBER, TESLA, ZOOX, CLEARPATH, cruise, WAYMO, NURO, nvitonomy, drive.ai, CAMBRIDGE MOBILE, Aurora, nauto, AIMOTIVE, G7, PILOT.AI, NIO, OPTIMUS, moovit, Ike, nexar, Kodiak, comma.ai, netradyne, Civil Maps, GAMAYA, Terravion, thinci, INRIX

**AGRICULTURE** — FARMERS, Granular, JOHN DEERE, BLUE RIVER, FarmersEdge, AgroStar, FarmLogs, TARANIS, prospera

**COMMERCE** — instacart, FAIRE, STITCH FIX, Dia & Co, HowGood, heuristics, ecrebo

**INDUSTRIAL** — AVEVA, SIEMENS, PREDIX, UPTAKE, Guardian Optical, DATAVISOR, SCORTEX, KONUX, TACHYUS

**OTHER** — eharmony, stem, Amper, ByteDance, happyco, celect, SOJERN, BOXEVER, VERDIGRIS, duetto, lakedeck, Spoke, Electric, ZINIER

## CROSS-INFRASTRUCTURE/ANALYTICS

AWS, Google Cloud, Microsoft, IBM, SAP, Hewlett Packard Enterprise, SAS, 1010DATA, vmware, TIBCO, TERADATA, ORACLE, NetApp, syncsort, MAPR, cloudera

## OPEN SOURCE

**FRAMEWORKS** — Hadoop, Spark, Flink, YARN, TEZ, MESOS, kubernetes, docker, CDAP, Red Hat, HELIX

**QUERY / DATA FLOW** — Spark SQL, HIVE, presto, APACHE DRILL, SLAMDATA, DRILL, GraphQL, Flink

**DATA ACCESS & DATABASES** — cassandra, mongoDB, redis, Zookeeper, Cockroach Labs, druid, CouchDB, SciDB, riak, HBASE, Cloud Spanner, accumulo

**ORCHESTRATION & MGMT** — talend, Apache Zookeeper, NiFi, Apache Ambari, Apache Airflow, MESOS, etcd, Kong

**STREAMING & MESSAGING** — Spark, Flink, beam, kafka, STORM, Apache RocketMQ

**STAT TOOLS & LANGUAGES** — python, Scala, Numpy, Kubeflow, mleap, DVC, SciPy, julia

**AI OPS & INFRA** — mlflow, Polyaxon, seldon

**AI / MACHINE LEARNING / DEEP LEARNING** — TensorFlow, Keras, MXNet, Caffe, Microsoft Cognitive Toolkit, OpenAI, DMTK, theano, PaddlePaddle, Apache SINGA, DIMSUM, FeatureFu, mxnet, VELES, Chainer, Michelangelo, ONNX, WEKA, YDWIG, PyTorch, neon, DSSTNE, mllib, DL4J, MAHOUT, Aerosolve, fast.ai, mlr, OpenML

**SEARCH** — elasticsearch, Solr, Lucene

**LOGGING & MONITORING** — elasticsearch, kibana, SENTRY, logstash, Prometheus, fluentbit, fluentd, Grafana, Vector

**VISUALIZATION** — matplotlib, BeakerX, TensorBoard, seaborn, jupyter, zeppelin, Bokeh

**COLLABORATION** — ANACONDA, accumulo

**SECURITY** — Apache Ranger, KNOX, Sentry, accumulo

## DATA SOURCES & APIs

**HEALTH** — Apple, VALIDIC, practice fusion, fitbit, GARMIN, HUMAN API, kinsa, MIMIC

**IOT** — GE Digital, UPTAKE, thingworx, helium, samsara, estimote

**FINANCIAL & ECONOMIC DATA** — Bloomberg, THOMSON REUTERS, DOW JONES, S&P CAPITAL IQ, CB INSIGHTS, PLAID, ENVESTNET YODLEE, PRECISIONHAWK, Descartes Labs, Eagle Alpha, THE WORLD BANK, estimize, PREMISE, Quandl, StockTwits, xignite, Thinknum, earnest, predata

**AIR / SPACE / SEA** — Orbital Insight, planet, SKYCATCH, AIRBOTICS, Spire, kespry, WINDWARD, DroneDeploy, MarineTraffic

**PEOPLE / ENTITIES** — acxiom, experian, EPSILON, InsideView, Crimson Hexagon, BASIS, Quantcast

**LOCATION INTELLIGENCE** — FOURSQUARE, mapbox, MapAnything, sense360, pitney bowes, HEXAGON, PlaceIQ, esri, factual, CARTO, Mapillary, StreetLine, cuebiq, Radar, OpenStreetMap, SAFEGRAPH

**OTHER** — DATA.GOV, IMAGENET, wiki links, wiki data, CRUX, grafiti.io

## DATA RESOURCES

**DATA SERVICES** — OPERA, LRO, GENERAL ASSEMBLY DATA SCIENCE, fractal, EXL, innoplexus

**INCUBATORS & SCHOOLS** — Six Sigma, GA, galvanize, DataCamp, DataElite, INSIGHT, The Data Incubator, PLURALSIGHT, kaggle, DataKind, METIS

**RESEARCH** — facebook research, OpenAI, MIRI, MILA, VECTOR INSTITUTE, CSAIL, ALLEN INSTITUTE for Artificial Intelligence

July 16, 2019 - FINAL 2019 VERSION

© Matt Turck (@mattturck), Lisa Xu (@lisaxu92), & FirstMark (@firstmarkcap)   mattturck.com/data2019

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

# Stragglers



# tasks

Map task completion time distribution

# Stragglers



# tasks

Map task completion time distribution

- **Tail execution time** means some workers (always) finish late

- Q: How can MR work around this?
  - Hint: its approach to **fault-tolerance** provides the right tool

# Resilience against stragglers

- If a task is going slowly (i.e., <span style="color:red">straggler</span>):
  - Launch second copy of task on another node
  - Take the output of whichever finishes first

# More design

- Master failure

- Locality

- Task granularity

# GFS usage at Google

- 200+ clusters

- Many clusters of 1000s of machines

- Pools of 1000s of clients

- 4+ PB filesystems

- 40 GB/s read/write load
  - In the presence of frequent hardware failures

* Jeff Dean, LADIS 2009

# MapReduce usage statistics over time

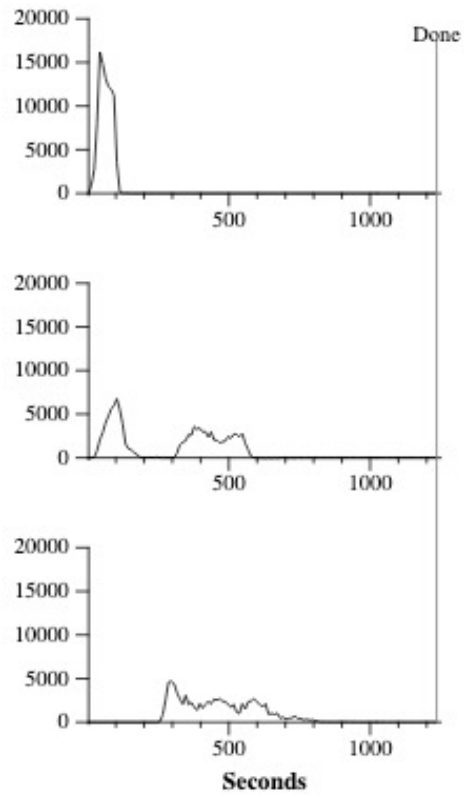|  | Aug, '04 | Mar, '06 | Sep, '07 | Sep, '09 |
|---|---|---|---|---|
| Number of jobs | 29K | 171K | 2,217K | 3,467K |
| Average completion time (secs) | 634 | 874 | 395 | 475 |
| Machine years used | 217 | 2,002 | 11,081 | 25,562 |
| Input data read (TB) | 3,288 | 52,254 | 403,152 | 544,130 |
| Intermediate data (TB) | 758 | 6,743 | 34,774 | 90,120 |
| Output data written (TB) | 193 | 2,970 | 14,018 | 57,520 |
| Average worker machines | 157 | 268 | 394 | 488 |

* Jeff Dean, LADIS 2009

# MapReduce discussion

What will likely serve as a performance bottleneck for Google's MapReduce used back in 2004 (or even earlier)? CPU? Memory? Disk? Network? Anything else?

# MapReduce discussion

What will likely serve as a performance bottleneck for Google's MapReduce used back in 2004 (or even earlier)? CPU? Memory? Disk? Network? Anything else?

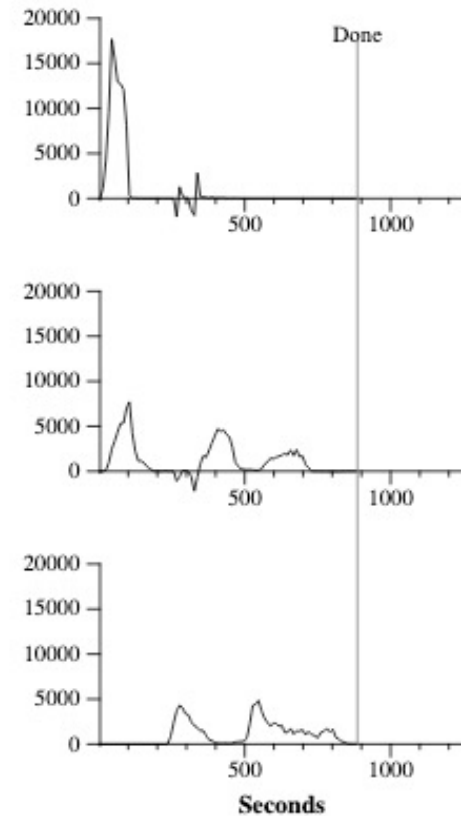How does MapReduce reduce the effect of slow network?

# MapReduce discussion



(a) Normal execution

(b) No backup tasks

(c) 200 tasks killed

# MapReduce discussion

Consider a log analytics job where you perform log-based debugging. You want to extract the timestamp info of all entries that match a keyword and then calculate the count of all matched entries:

1. Filter the entries with the keyword;

2. Calculate the count of all matched entries

What are the main shortcomings of using MapReduce to support such pipeline-like applications?

# Next step

- Look out for
  - Project suggestion doc
  - Fill the team composition form
  - Project bid and team composition due by Feb 24

- Next week: Apache Spark