# Google MapReduce

*DS 5110: Big Data Systems (Spring 2023)*
Lecture 3b

Yue Cheng

UNIVERSITY of VIRGINIA

WC.

Applications

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

Scalable computing engines

MR.

Scalable storage systems

GFS.

Datacenter infrastructure

# The big picture (motivation)

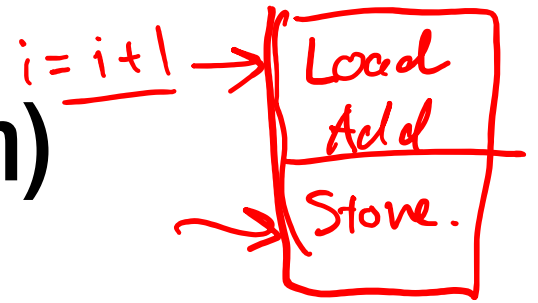- Datasets are <span style="color:red">too big</span> to process using a single computer

# 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>

MPI:
Message Passing Interface.

Decomposition.

# The big picture (motivation)

*i = i + 1* → Load / Add / Store.

- Datasets are too big to process using a single computer

Atomicity

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

expressive.

- 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 <span style="color:red">massive</span> 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

*cmd tools.*

```
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?

BSP.

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
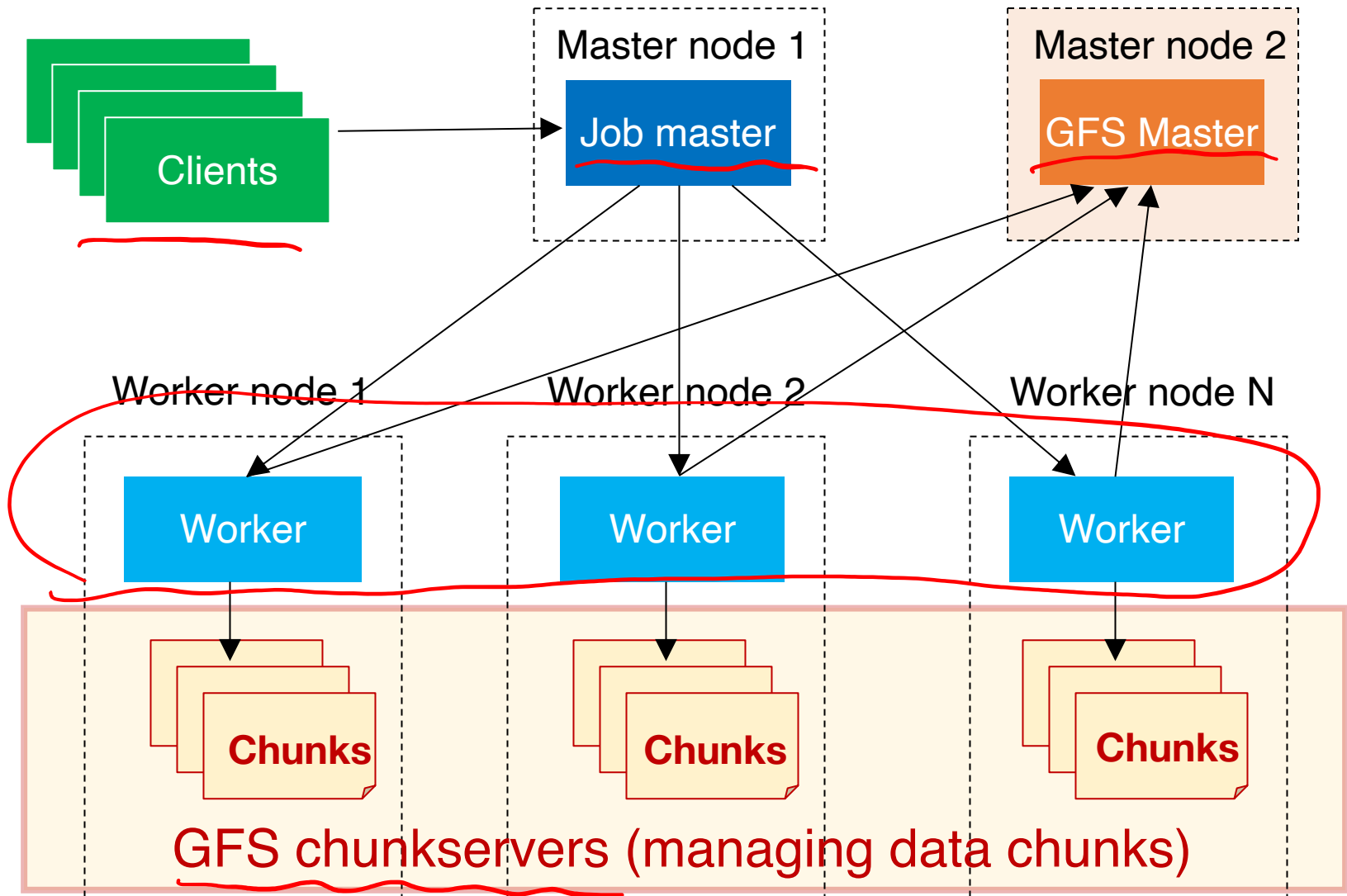
2. Then merge intermediate output

# 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`)

*Intermediate Step → Shuffle.*

- `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

*(handwritten: Ln.    Line str.)*

*(arrow)* map(key, value):
    for each word w in value:
        EmitIntermediate(w, "1");

*(handwritten: Shuffle.  w    List("1", "1", "1", ...))*

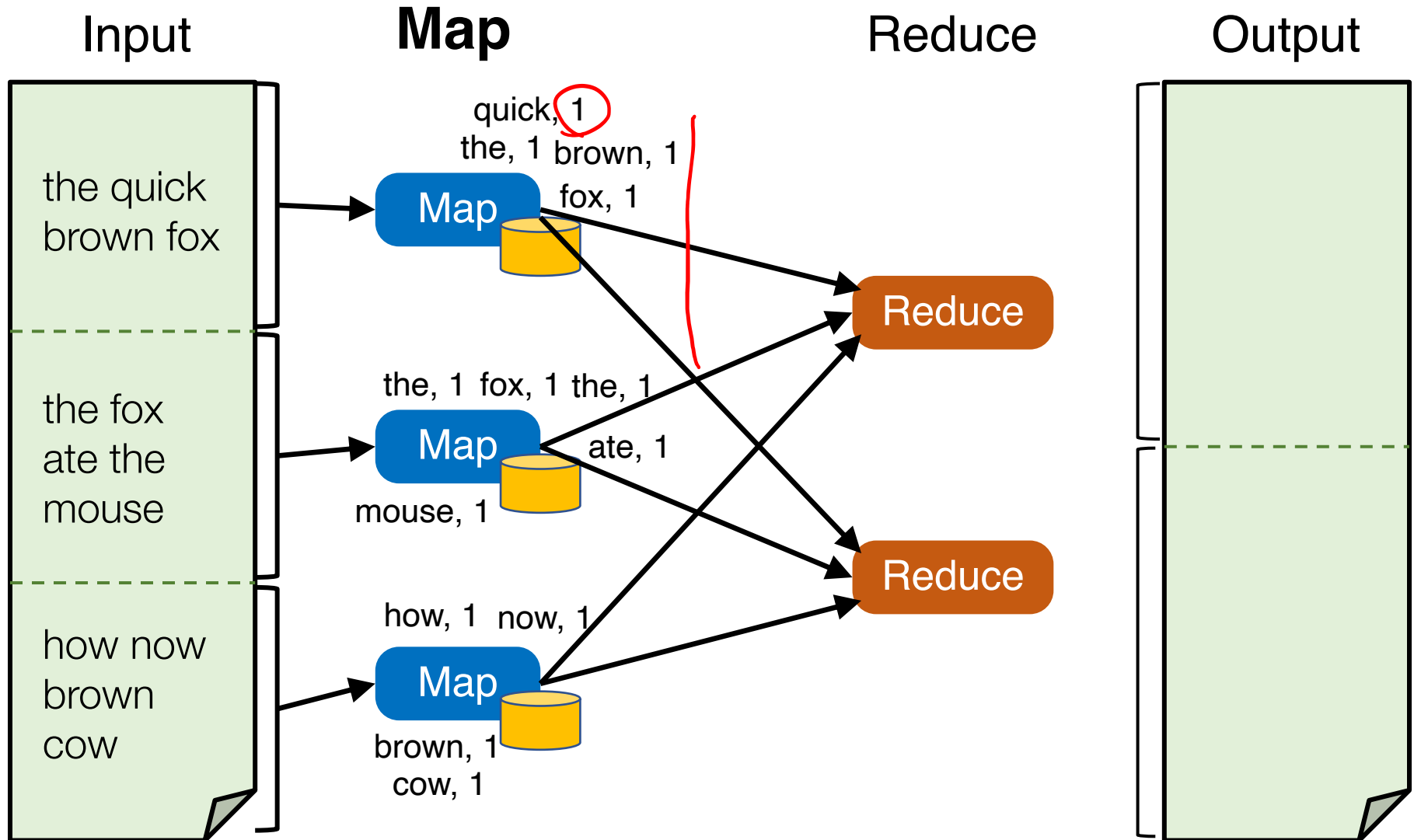*(arrow)* reduce(key, values):
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));

# Word Count execution

GFS.

| 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

Map          ate, 1

mouse, 1

how, 1  now, 1

Map

brown, 1
cow, 1

Reduce

Reduce

# Word Count execution

hash func.

All- to- All.

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

the quick brown fox

the fox ate the mouse

how now brown cow

Map

Map

Map

R

R

R

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

the quick brown fox

the fox ate the mouse

how now brown cow

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

the, 1
fox, 1

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

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

ate, 1
cow, 1
mouse, 1
quick, 1
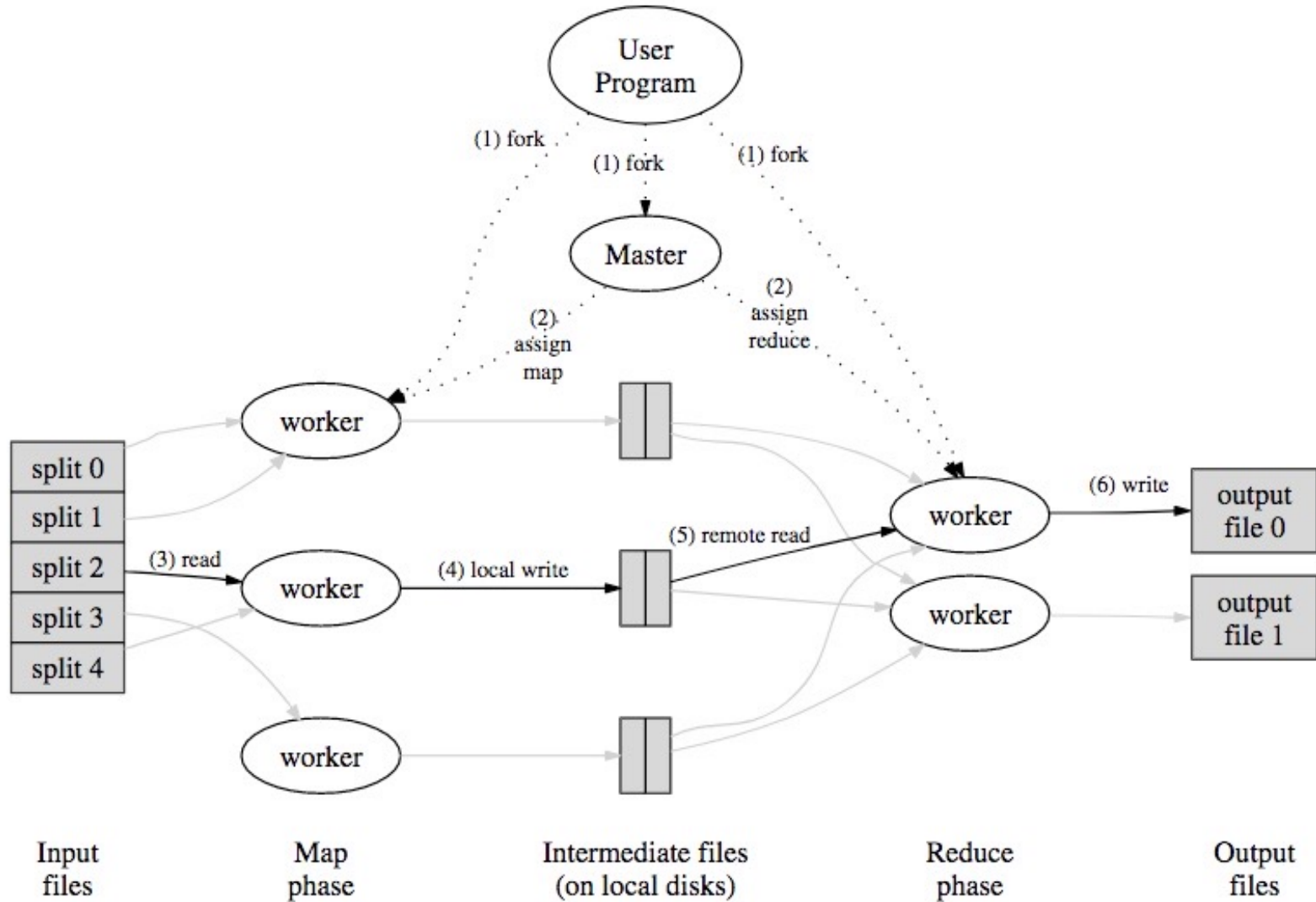
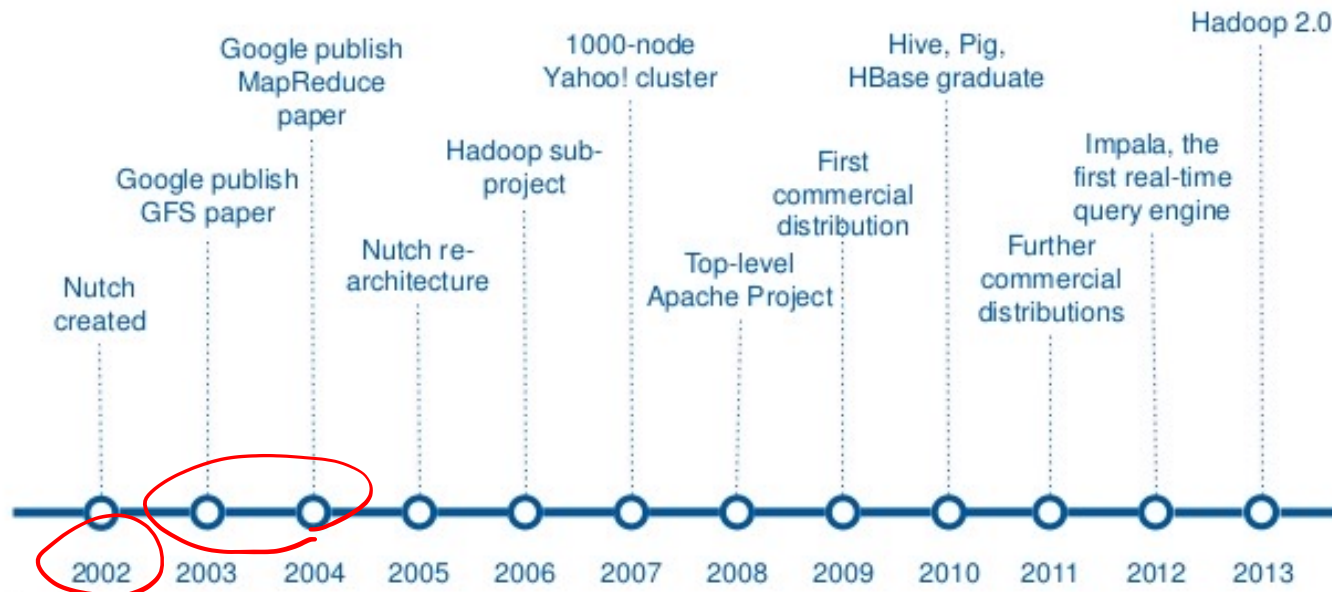# MapReduce data flows in paper

# How it started: Apache Hadoop

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



A Brief History of Hadoop

© 2013, Axidata Systems FZ-LLC

# 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 TREASURE DATA · 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 · MemSQL · MariaDB · influxdata · Cockroach Labs · VOLTDB · splice · imply · paradigm4 · TIBR

**GRAPH DBs**
neo4j · Amazon Neptune · IBM · ORACLE · OrientDB · Kognito · InfiniteGraph · Objectivity

**MPP DBs**
TERADATA · VERTICA · IBM Data Warehouse Systems · Kaction · ORACLE · 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 · ZALONI · import.io · Infoworks · Fivetran · SNOWPLOW · MATILLION

**DATA GOVERNANCE**
Informatica · IBM · SailPoint · collibra · ALEX · BigID · DATABANK · SignalFx · Alation · Waterline Data · IMMUTA · OKERA · unravel · Numerify · zenloss · OpsRamp · MANTA · data.world

**MGMT / MONITORING**
aws · New Relic · actifio · rubrik · APPDYNAMICS · dynatrace · WAVEFRONT · DATADOG · splunk · Moogsoft · pagerduty · SCALYR · VEEAM · OptScale · MAGNITUDE

**DATA TRANSFORMATION / STORAGE**
aws · Google Cloud · IBM Storage · Microsoft Azure · Amazon ECS · PURE STORAGE · ALLUXIO · wasabi · nimble storage · Qumulo · panasas · COHESITY

**CLUSTER SVCS**
IBM · ALGORITHMIA · SPELL · comet · Cnvrg · MESOSPHERE · packet · Rancher · HIVE · scale · Weight & Biases · BrightComputing · CYCLECLOUD

**DATA GENERATION & LABELLING**
amazon mechanical turk · Upwork · unity · appen · datmo · Scale · Datatron · Mighty AI · FIGURE EIGHT · LIONBRIDGE

**AI OPS**
datatron · brytlyt · PG-Strom · BLAZINGDB · Mavidius · habana · CERNAI · HALO · Paperspace

**GPU DBs & CLOUD**
kinetica · omni sci · SQREAM · Cerebras · Determined AI · FLOYDHUB

**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 · sisu · switchboard · Starburst

**DATA SCIENCE PLATFORMS**
IBM · databricks · dataiku · DOMINO · rapidminer · TIBCO · SAS · Altair · ANACONDA · KNIME · MathWorks · H2O.ai · DataRobot · gamalon

**BI PLATFORMS**
looker · einstein analytics · aws · tableau · Power BI · DOMO · ARCADIA DATA · ThoughtSpot · SAP · Qlik · ATSCALE · SSENSE · GoodData · Information Builders · birst · MicroStrategy · Keen IO

**VISUALIZATION**
tableau · Power BI · SAS · Google Cloud · Celonis · Periscope Data · zepl · GEOMDATA · plotly · CHARTIO · Toucan Toco

**MACHINE LEARNING**
Azure Machine Learning · Google Cloud AutoML, Vision · Amazon SageMaker · H2O.ai · DataRobot · gamalon · ViSENZE · ELEMENT AI · deepsense.ai

**COMPUTER VISION**
Microsoft Azure · Amazon Rekognition · clarifai · deepomatic · EVER AI · neurala · twentybn · UBIQITY.6 · YITU · trax · BLUE VISION · Numenta

**HORIZONTAL AI**
IBM Watson · Cortana · Face++ · sentient · Voyager.ai · vicarious · Affectiva · PROPHESEE · CognitiveScale · PETUUM · SI · NaraLogics · CURIOUS AI · OSARO · Pollexion · Fortress AI

**SPEECH & NLP**
Google Cloud · twilio · amazon alexa · Amazon Translate · OpenAI · semanticmachines · Mobvoi · EigenTechnologies · SoundHound Inc. · x.ai · PRIMER · volcero · cogito · snips · SMARTLING · Unbabel · PolyAI

**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

### APPLICATIONS – ENTERPRISE

**SALES**
CHORUS · INSIDESALES.COM · peopleai · conversica · Clari · aviso · tact.ai · TROOPS · fuse|machines · Clearbit

**MARKETING - B2B**
RADIUS · App Annie · EVERSTRING · Lattice · MINTIGO · 6sense · contentsquare · TEALIUM · imparticle · Amplero · amperity · QUANTIFIND · ENGAGIO · Lytics · PERSADO · KNOTCH · mrp

**MARKETING - B2C**
zeta · bloomreach · SendGrid · braze · ACTIONIQ · BLUECORE · tubular · Reflektion · Simon · remesh

**CUSTOMER EXPERIENCE / SERVICE**
qualtrics · MEDALLIA · SurveyMonkey · UserTesting · CLARABRIDGE · zendesk · Kustomer · freshdesk · INTERCOM · Drift · LIVEPERSON · Gainsight · pendo · HEAP · Amplitude · Watson Assistant · ada · AUTOMAT · afiniti · DigitalGenius · ASAPP · clara · Cee[bots · netomi · Diffbot · GURU · lumiata · talla · Kasisto

**ENTERPRISE PRODUCTIVITY**
slack · ORACLE

**HUMAN CAPITAL**
HireVue · pymetrics · hiQ · GIGSTER · mya · Allyo · textio · Wade&Wendy · Stella · Cushee · entelo · RESTLESS BANDIT · beamery

**LEGAL**
RAVEL · Seal · Everlaw · Disco · kira · JUDICATA · BREVIA · IRONCLAD · LegalVision · PREMONITION · ROSS · casetext

**REGTECH & COMPLIANCE**
Comply Advantage · Refinitiv · CROSSBEAM · DATA REPUBLIC

**FINANCE**
anaplan · Zuora · SAP/S4HANA · TRADESHIFT · VIDADO · mineraltree · SCALEFACTOR · INTOSUM · pilot

**BACK OFFICE AUTOMATION & RPA**
UiPath · HyperScience · blueprism · AppZen · WorkFusion · workato · AntWorks · Catalytic · botkeeper · KRYON · ALKYMI

**SECURITY**
TANIUM · CYLANCE · zscaler · StackPath · illumio · CODE42 · CipherCloud · DARKTRACE · ANOMALI · ThreatMetrix · SIGNIFYD · SentinelOne · SecurityScorecard · SECURE · Vade Secure · bitglass · BlueTalon · Recorded Future · feedzai · Cybex · BITSIGHT · sparkcognition · IronNet Cybersecurity · FORTER · riskrecon · JASK · AREA 1 Security · BLUEHEXAGON · Semmle · OBSIDIAN · AXONIUS · SHIELD AI · ArmorBlox

### APPLICATIONS – INDUSTRY

**ADVERTISING**
AppNexus · MediaMath · criteo · xAd · Integral · ORACLE MOAT · theTradeDesk · dstillery · TAPAD · dataxu · gumgum · Appier · yieldmo

**EDUCATION**
Liulishuo · KNEWTON · Clever · declara · kidaptive · PANORAMA · knowre · gradescope

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

**GOV'T**
OPENGOV · mark43 · FN FiscalNote · LiveStories · Passport · SmartProcure · PAGAYA

**INTELLIGENCE**
Palantir · Datamin · Quid · PRIMER · Quantopian

**FINANCE - INVESTING**
KENSHO · Quantopian · ADDEPAR · Aurora · iSENTIUM · ALGORIZ · TrueAccord · aire

**FINANCE - LENDING**
ondeck · Affirm · JIANPU.AI · TALA · even financial · Upstart · AVANT · MoneyLion · cignifi

**INSURANCE**
Metromile · Lemonade · CYENCE · Hippo · Shift Technology · ROOT · zesty.ai · TRACTABLE · CAPE

**HEALTHCARE**
flatiron · Clover · KYRUUS · HealthTap · METABIOTA · Ginger.io · Glow · babylon · 3DMed · zebra · PathAI · ovia · TEMPUS · patientslikeme · AiCure · insitro · notable · komodohealth · RECURSION · prognos · enlitic · BlackThorn · CITRINE · Qventus · ARTERYS · IMAGEN · PAIGE · DATAVANT · innovaccer

**LIFE SCIENCES**
Color · verily · WuXiNextCODE · CANOPY · DNAnexus · NIO · deep genomics · OWKIN

**TRANSPORTATION**
UBER · TESLA · WAYMO · ZOOX · CLEARPATH · cruise · NURO · nvidia · drive.ai · CAMBRIDGE · Aurora · nauto · AIMOTIVE · G7 · PILOT.AI · nexar · Kodiak · comma.ai · netradyne · Civil Maps · thinci · INRIX

**AGRICULTURE**
FARMERS · Granular · JOHN DEERE · AgroStar · FarmLogs · TARANIS · GAMAYA · Terravion · prospera

**COMMERCE**
Instacart · FAIRE · Dia & Co · RetailNext · HowGood · AgroStar

**INDUSTRIAL**
AVEVA · SIEMENS · PREDIX · UPTAKE · SCORTEX · KONUX · TACHYUS

**OTHER**
eharmony · stem · Amper · ByteDance · Clect · SOJERN · BOXEVER · VERDIGRIS · hipdeck · duetto · Electric · ZINIER · Spoke

### 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

**QUERY / DATA FLOW**
Spark SQL · HIVE · presto · Apache Drill · SLAMDATA · GraphQL · Flink

**DATA ACCESS & DATABASES**
cassandra · mongoDB · redis · 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 · Julia · SciPy

**AI OPS & INFRA**
mlflow · Kubeflow · mleap · DVC · SELDON · Polyaxon

**AI / MACHINE LEARNING / DEEP LEARNING**
TensorFlow · Keras · MLlib · Caffe · Microsoft Cognitive Toolkit · OpenAI · DMTK · theano · PaddlePaddle · Apache SINGA · DIMSUM · FeatureFu · mxnet · VELES · Chainer · Michelangelo · ONNX · WEKA · 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 · TensorBoard · seaborn · Bokeh · BeakerX · jupyter · Zeppelin

**COLLABORATION**
ANACONDA · accumulo

**SECURITY**
Apache Ranger · KNOX · Sentry · accumulo

### DATA SOURCES & APIs

**HEALTH**
VALIDIC · practicefusion · fitbit · GARMIN · HUMAN API · kinsa · MIMIC

**IOT**
GE Digital · thingworx · UPTAKE · helium · samsara · estimote

**FINANCIAL & ECONOMIC DATA**
Bloomberg · THOMSON REUTERS · DOW JONES · S&P CAPITAL IQ · CB INSIGHTS · PLAID · ENVESTNET YODLEE · PRECISIONHAWK · Descartes Labs · Quandl · Eagle Alpha · THE WORLD BANK · estimize · PREMISE · 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

**OTHER**
DATA.GOV · IMAGENET · wiki links · wiki data · LabelMe · CRUX · ili.grafiti.io

### DATA RESOURCES

**DATA SERVICES**
OPERA · LIG · Six Sigma · DATA SCIENCE · PLURALSIGHT · fractal · EXL

**INCUBATORS & SCHOOLS**
GA galvanize · DataCamp · DataElite · INSIGHT · The Data Incubator · kaggle · DataKind · innoplexus · METIS

**RESEARCH**
facebook research · OpenAI · MIRI · MILA · CSAIL · VECTOR INSTITUTE · Qi · AI2 ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

# Stragglers

hist.

# tasks

long tail.

Map task completion time distribution

# Stragglers



Map task completion time distribution

$$x \cdot 1 = x$$

$$x \cdot 0 = 0$$

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

Idempotente.

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

# Resilience against stragglers

*speculative exe.*

- If a task is going slowly (i.e., straggler):
  - Launch second copy of task on another node
  - Take the output of whichever finishes first

# More design

*Job*

- Master failure

- Locality

- Task granularity

M    R

# computers.

→ LB

# 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?
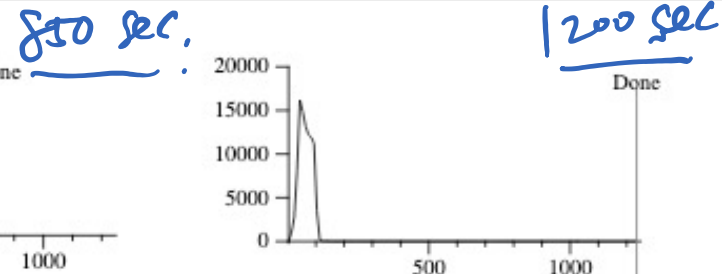
12.5 MB/s.

How does MapReduce reduce the effect of slow network?

# MapReduce discussion



(a) Normal execution  (b) No backup tasks  (c) 200 tasks killed

*Handwritten annotations:*
- Map.
- 13 GB/s.
- 850 sec.
- 1200+ sec
- Failures of 200 Workers
- 900 sec.
- ① ②
- Replication
- 700
- Reduce

# 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