

Applied Data Analysis (CS401)

Lecture 13

Scaling to
data

14 Dec 2022



model
massive

EPFL

Robert West



Announcements

- Homework H2 grades have been released
- Project milestone P3 due next week (Fri 23 Dec)
- Friday's lab session:
 - Last lab session! → Last quiz (on lecture 12)
 - Project office hour (same [sign-up protocol](#) as last week)
 - Exercises on Spark (useful for your future projects, your job, your love life)
- **Course eval is available on IS-Academia!**
 - Note: this is different from the eval from a few weeks ago...

Feedback

Give us feedback on this lecture here:

<https://go.epfl.ch/ada2022-lec13-feedback>

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more (fewer) details?
- Where is Waldo?
- ...



So far in this class...

- We made one big assumption:
 - All data fits on a single machine
 - Even more, all data fits into memory on a single machine (Pandas)
- Realistic assumption for **prototyping**, but frequently not for production code

The big-data problem

Data is growing faster than computation speed

Growing data sources

>> Web, mobile, sensors, ...

Cheap hard-disk storage

Stalling CPU speeds

RAM bottlenecks



Examples

Facebook's daily logs: 60 TB

1000 Genomes project: 200 TB

Google Web index: [100+ PB](#)

Cost of 1 TB of disk: \$50

Time to read 1 TB from disk: 3 hours (100 MB/s)



DISCLAIMER

**These numbers
(anno domini
2016) are
outdated (too
small)!**

The big-data problem

Single machine can no longer store, let alone process, all the data

Only solution is to **distribute** over a large cluster of machines

But how much data should you get?

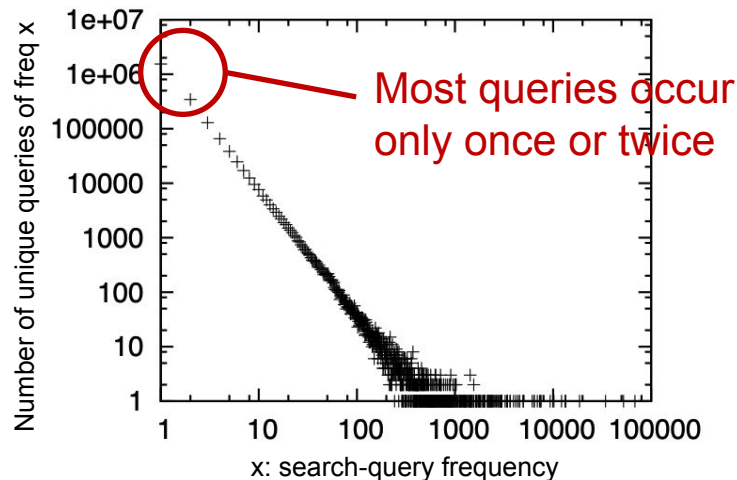
Of course, “it depends”, but for many applications the answer is:

As much as you can get

Big data about people (text, Web, social media) tends to follow heavy-tailed distributions (e.g., power laws)

Example: Web search

59% of all Web search queries are unique
17% of all queries were made only twice
8% were made three times



Hardware for big data

Budget (a.k.a. commodity) hardware
Not "gold-plated" (a.k.a. custom)

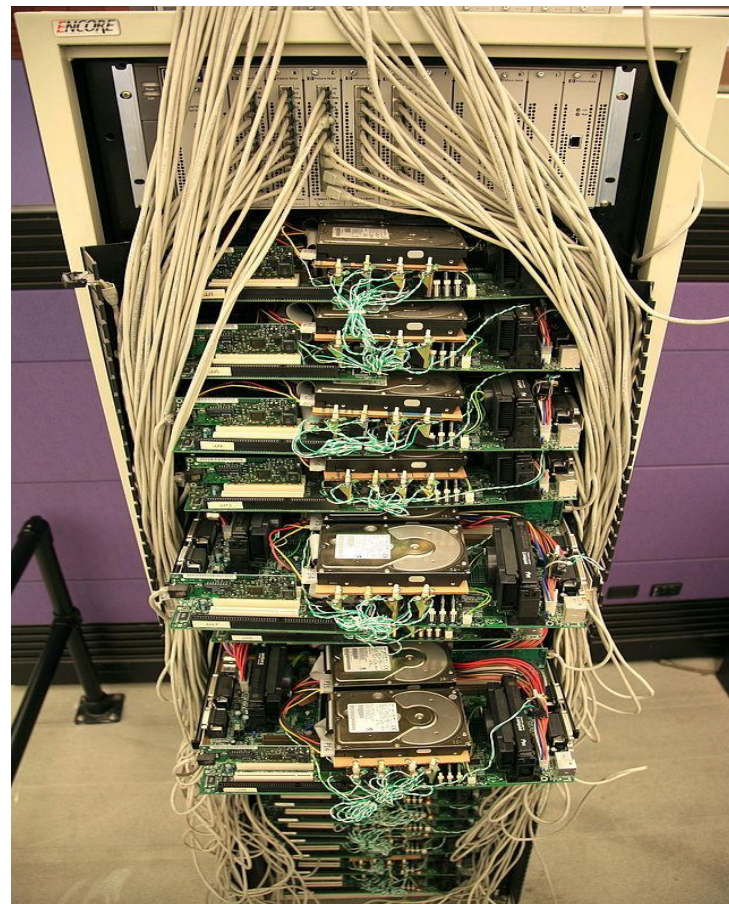
Many low-end servers

Easy to add capacity

Cheaper per CPU and per disk

Increased complexity in software:

- Fault tolerance
- Virtualization (e.g., distributed file systems)



Problems with cheap hardware

Failures, e.g. (Google numbers)

- 1-5% hard drives/year
- 0.2% DIMMs (dual in-line memory modules)/year

Commodity network (1-10 Gb/s) speeds vs. RAM

- Much more latency (100x – 100,000x)
- Lower throughput (100x-1000x)

Uneven performance

- Inconsistent hardware (e.g., old + new)
- Variable network latency
- External loads



DISCLAIMER

These numbers are constantly changing thanks to new technology!

Google datacenter

A wide-angle, low-perspective shot of a vast Google datacenter. The room is filled with rows of server racks, some of which are illuminated with blue and yellow lights. The ceiling is high and features a complex network of steel beams and pipes. The floor is made of large, light-colored tiles. The overall atmosphere is one of a highly advanced and organized technological environment.

How to program this beast?

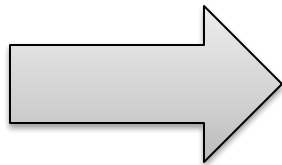
What's hard about cluster computing?

How do we split work across machines?

How do we deal with failures?

How do you count the number of occurrences of each word in a document?

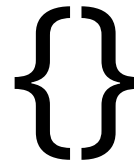
"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and
ham?"



I: 3
am: 3
Sam: 3
do: 1
you: 1
like: 1
...

A hashtable (a.k.a. dict)!

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and
ham?"



A hashtable!

"I am Sam

I am Sam

Sam I am

Do you like

Green eggs and
ham?"

{I: 1}

A hashtable!

"I am Sam

I am Sam

Sam I am

Do you like

Green eggs and
ham?"

{I: 1,
am: 1}

A hashtable!

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and
ham?"

{I: 1,
am: 1,
Sam: 1}

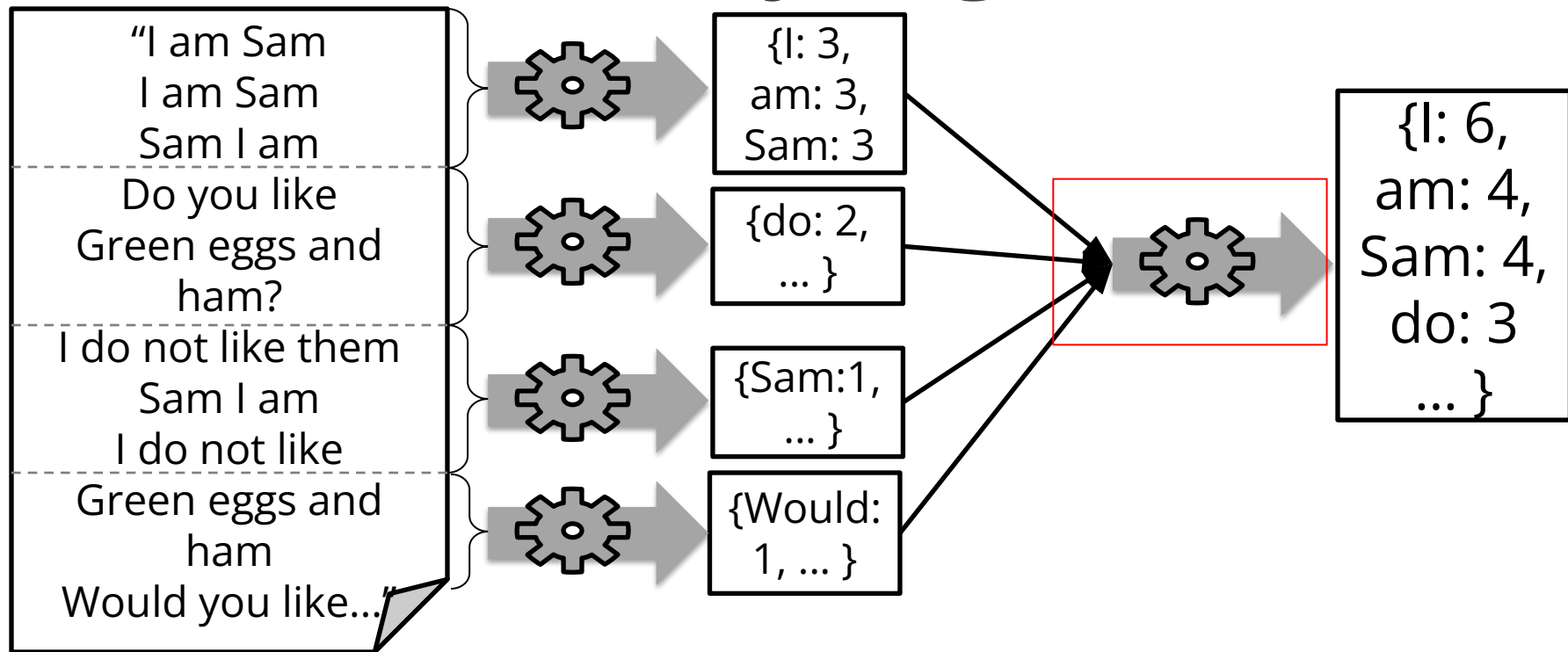
A hashtable!

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and
ham?"

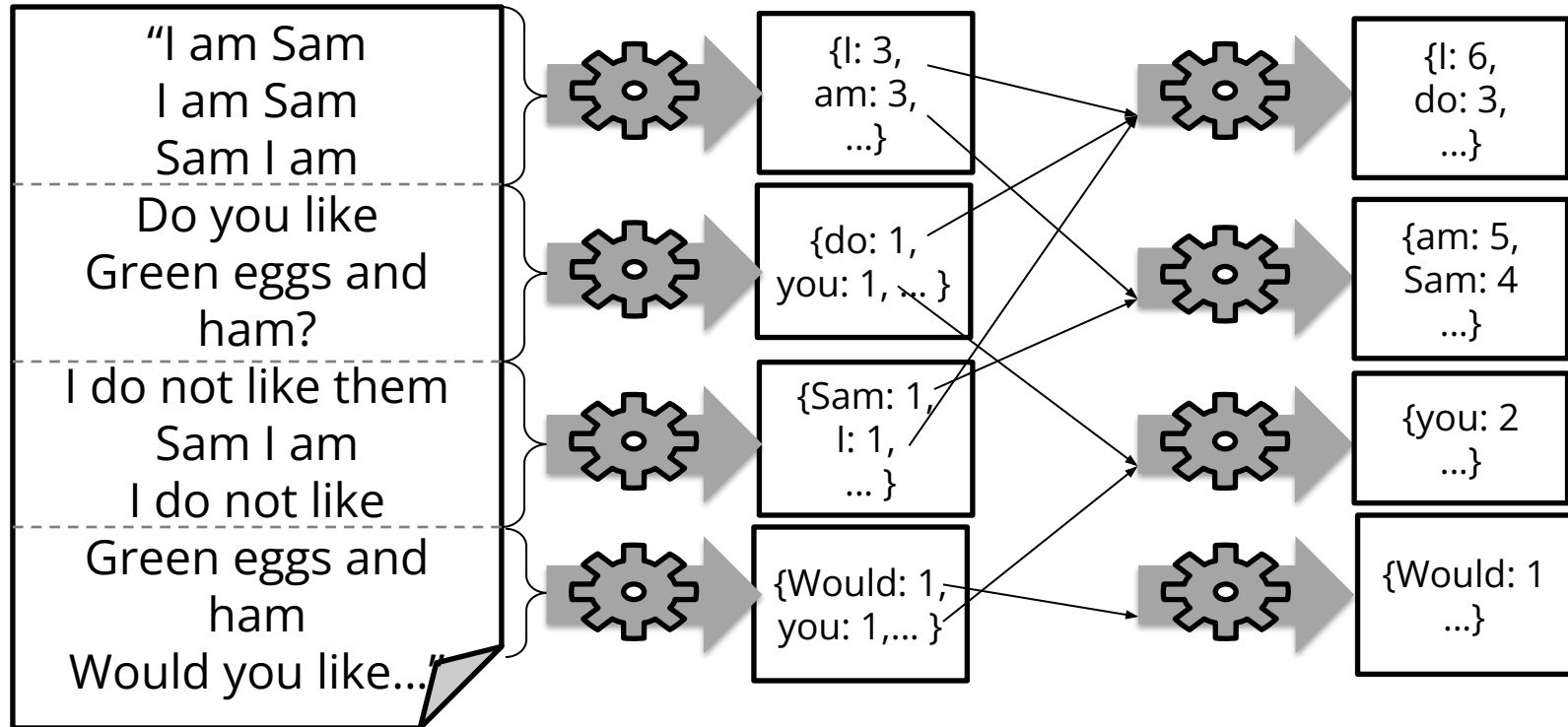
{I: 2,
am: 1,
Sam: 1}

What if the document is
really big?

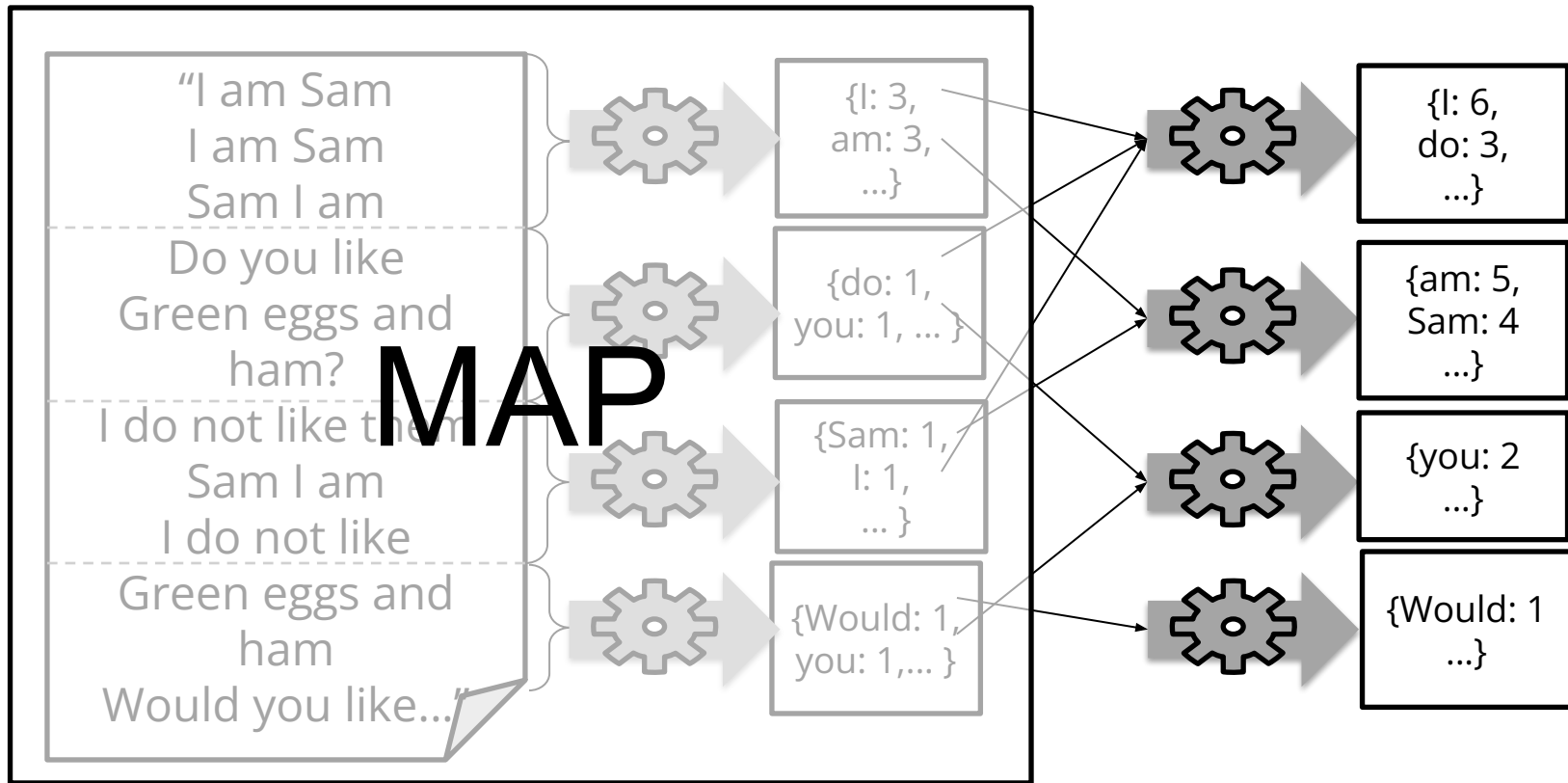
What if the document is really big?



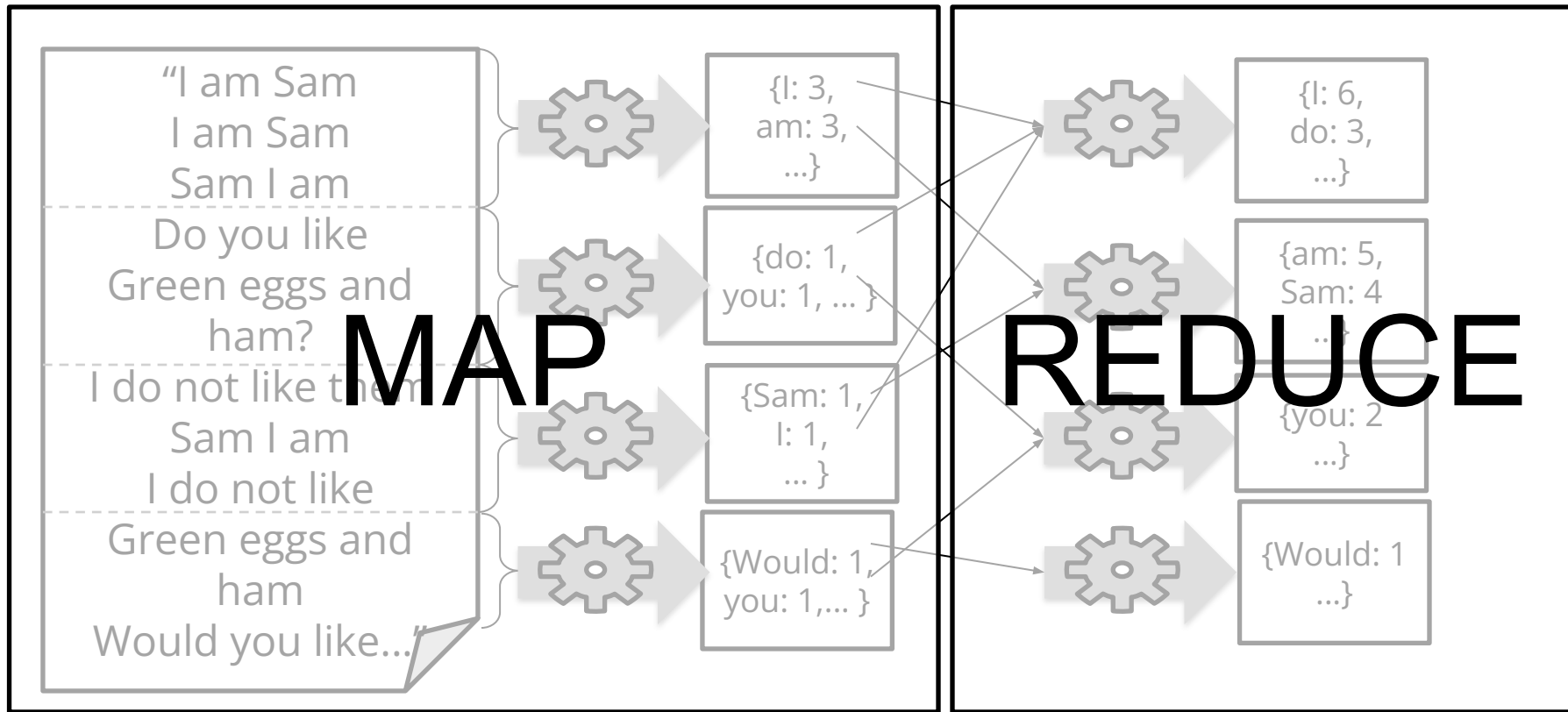
"Divide and Conquer"



"Divide and Conquer"



“Divide and Conquer”



What's hard about cluster computing?

How to divide work across machines?

- Must consider network, data locality
- Moving data may be very expensive

How to deal with failures?

- 1 server fails every 3 years => 10K servers see ~10 faults/day
- Even worse: stragglers (node not failed, but slow)

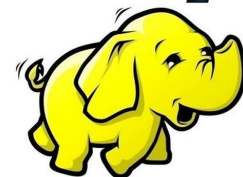
Solution: MapReduce

- Smart systems engineers have done all the work for you
 - Task scheduling
 - Virtualization of file system
 - Fault tolerance (incl. data replication)
 - Job monitoring
 - etc.
- “All” you need to do: implement Mapper and Reducer classes



Jeff Dean [\[facts\]](#)

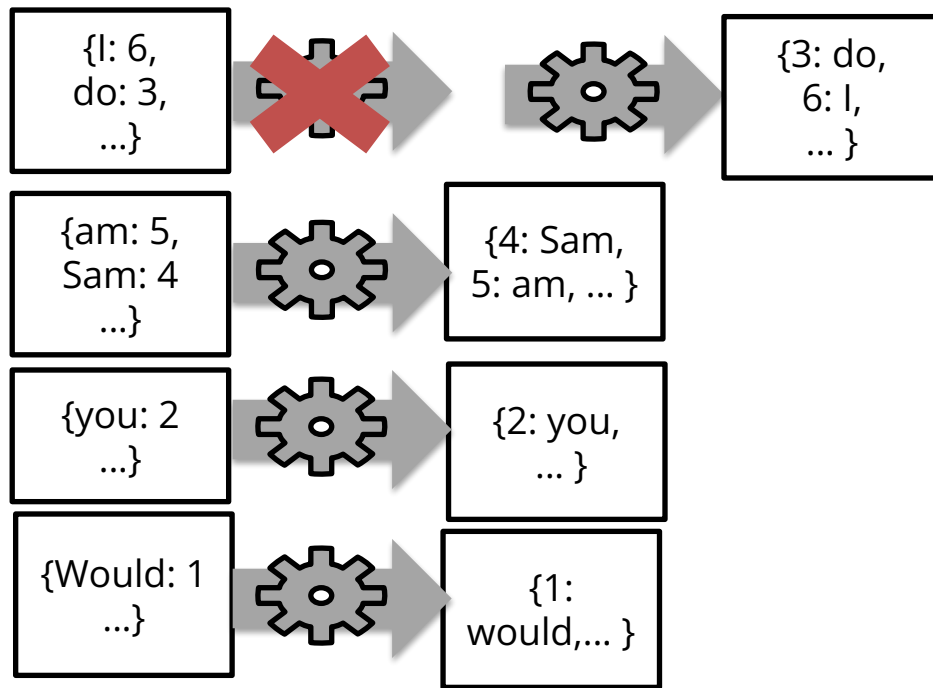
hadoop





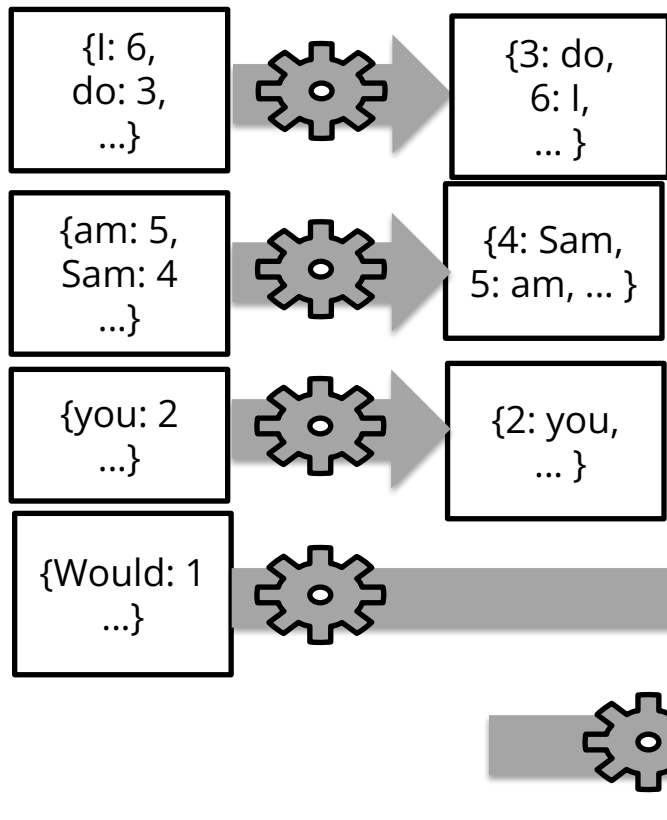
Applied Machine Learning Days '19 [\[link\]](#)

How to deal with failures?



Just launch another task!

How to deal with slow tasks?



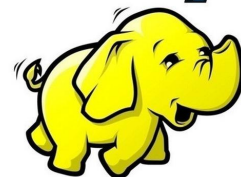
Just launch
another task!

Solution: MapReduce



Jeff Dean

hadoop



- Smart systems engineers have done all the work for you
 - Task scheduling
 - Virtualization of file system
 - Fault tolerance (incl. data replication)
 - Job monitoring
 - etc.
- **“All”** you need to do: implement Mapper and Reducer classes

Need to break more complex jobs into sequence of MapReduce jobs

Example task

Suppose you have user info in one file, website logs in another, and you need to find the top 5 pages most visited by users aged 18-25.



In MapReduce

[illegible]


Enter: *Spark*



- A high-level API for programming MapReduce-like jobs

```
sc = SparkContext()
print "I am a regular Python program, using the pyspark lib"
users = sc.textFile('users.tsv') # user <TAB> age
    .map(lambda s: tuple(s.split('\t')))
    .filter(lambda (user, age): age>=18 and age<=25)
pages = sc.textFile('pageviews.tsv') # user <TAB> url
    .map(lambda s: tuple(s.split('\t')))
counts = users.join(pages)
    .map(lambda (user, (age, url)): (url, 1))
    .reduceByKey(add)
    .takeOrdered(5)
```




- Implemented in  **Scala** (go, EPFL!)
- Additional APIs in
 - Python
 - Java
 - R



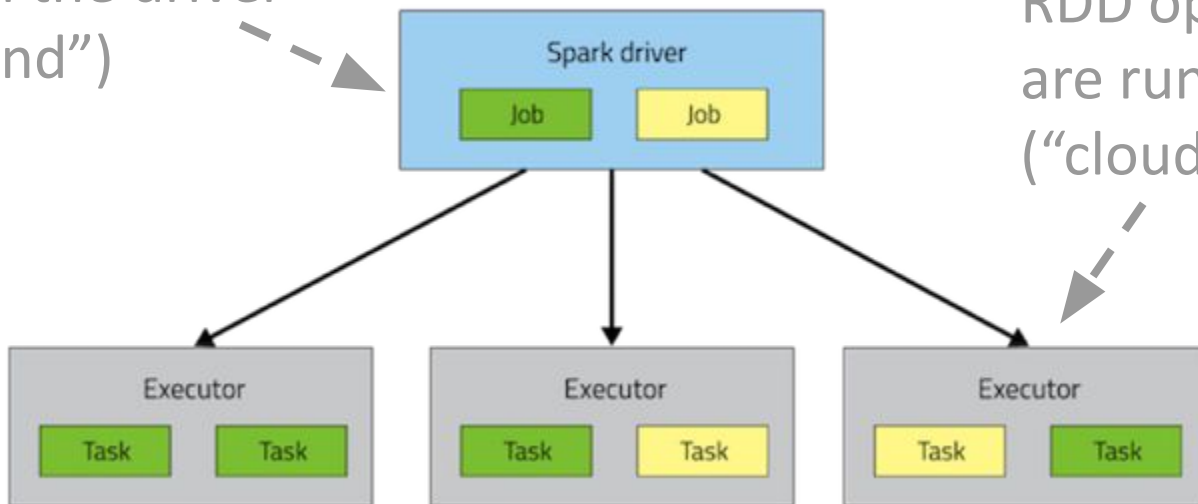
RDD: resilient distributed dataset



- To programmer: looks like one single list (each element represents a “row” of a dataset)
- Under the hood: oh boy...
 - RDDs “live in the cloud”: split over several machines, replicated, etc.
 - Can be processed in parallel
 - Can be transformed to a single, real list (if small...)
 - Typically read from the distributed file system (HDFS)
 - Can be written to the distributed file system

Spark architecture

Your Python script runs in the driver (“ground”)



RDD operations are run in executors (“cloud”)



RDD operations



- **“Transformations”**

- Input: RDD; output: another RDD
- Everything remains “in the cloud”
- Example: for every entry in the input RDD, count chars
 - `RDD:['I', 'am', 'you'] → RDD:[1, 2, 3]`

- **“Actions”**

- Input: RDD; output: a value that is returned to the driver
- Result is transferred “from cloud to ground”
- Example: take a sample of entries from RDD and print it on the driver’s shell

Lazy execution [unrelated]

- **Transformations** (i.e., RDD→RDD operations) are not executed until it's really necessary (a.k.a. “lazy execution”)
- Execution of transformations triggered by **actions**
- Why?
 - If you never look at the data, there's no point in manipulating it...
 - Smarter query processing possible:
E.g., `rdd2 = rdd1.map(f1)`
`rdd3 = rdd2.filter(f2)`
Can be done in one go -- no need to materialize rdd2



"I have good news and bad news"

RDD transformations [\[full list\]](#)

- **map**(*func*): Return a new distributed dataset formed by passing each element of the source through a function *func*
 - $\{1,2,3\}.\text{map}(\text{lambda } x: x*2) \rightarrow \{2,4,6\}$
- **filter**(*func*): Return a new dataset formed by selecting those elements of the source on which *func* returns true
 - $\{1,2,3\}.\text{filter}(\text{lambda } x: x \leq 2) \rightarrow \{1,2\}$
- **flatMap**(*func*): Similar to map, but each input item can be mapped to 0 or more output items (so *func* should return a list rather than a single item)
 - $\{1,2,3\}.\text{flatMap}(\text{lambda } x: [x,x*10]) \rightarrow \{1,10,2,20,3,30\}$

RDD transformations [\[full list\]](#)

- **sample**(*withReplacement?*, *fraction*, *seed*): Sample a fraction *fraction* of the data, with or without replacement, using a given random number generator *seed*
- **union**(*otherDataset*): Return a new dataset that contains the union of the elements in the source dataset and the argument.
- **intersection**(*otherDataset*): ...
- **distinct**(): Return a new dataset that contains the distinct elements of the source dataset.

RDD transformations [\[full list\]](#)

- **sample**(*withReplacement?*, *fraction*, *seed*): Sample a fraction *fraction* of the data, with or without replacement, using a given random number generator *seed*

Why *relative fraction*,
and not *absolute
number*?

POLLING TIME

Scan QR code or go to
<https://web.speakup.info/room/join/66626>



RDD transformations [\[full list\]](#)

- **groupByKey()**: When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable<V>) pairs.
 - $\{(1,a), (2,b), (1,c)\}.groupByKey() \rightarrow \{(1,[a,c]), (2,[b])\}$
- **reduceByKey(*func*)**: When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function *func*, which must be of type (V, V) => V.
 - $\{(1, 3.1), (2, 2.1), (1, 1.3)\}.reduceByKey(\text{lambda } (x,y): x+y) \rightarrow \{(1, 4.4), (2, 2.1)\}$

RDD transformations [\[full list\]](#)

- **sortByKey()**: When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs sorted by keys
- **join(*otherDataset*)**: When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with all pairs of elements for each key
 - $\{(1,a), (2,b)\}.join(\{(1,A), (1,X)\}) \rightarrow \{(1, (a,A)), (1, (a,X))\}$
- Analogous: **leftOuterJoin**, **rightOuterJoin**, **fullOuterJoin**
- (There are several other RDD transformations, and some of the above have additional arguments; cf. [tutorial](#))

RDD actions [\[full list\]](#)

- **collect()**: Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
- **count()**: Return the number of elements in the dataset.
- **take(*n*)**: Return an array with the “first” *n* elements of the dataset.
- **saveAsTextFile(*path*)**: Write the elements of the dataset as a text file in a given directory in the local filesystem or HDFS.
- (There are several other RDD actions; cf. [tutorial](#))

Broadcast variables

- `my_set = set(range(1e80))`
`rdd2 = rdd1.filter(lambda x: x in my_set)`
^ This is a bad idea: `my_set` needs to be shipped with every task (one task per data partition, so if `rdd1` is spread over N partitions, the above will require copying the same object N times)
- Better:
`my_set = sc.broadcast(set(range(1e80)))`
`rdd2 = rdd1.filter(lambda x: x in my_set.value)`
^ This way, `my_set` is copied to each executor only once and persists across all tasks (one per partition) on the same executor
- Broadcast variables are **read-only**

Accumulators

- `def f(x): return x*2`
`rdd2 = rdd1.map(f)`
^ How can we easily know how many rows there are in rdd1 (without running a costly reduce operation)?
- Side effects via accumulators!
`counter = sc.accumulator(0)`
`def f(x): counter.add(1); return x*2`
`rdd2 = rdd1.map(f)`
- Accumulators are **write-only** (“add-only”) for executors
- Only driver can read the value: `counter.value`

RDD persistence

```
rdd2 = rdd1.map(f1)
list1 = rdd2.filter(f2).collect()
list2 = rdd2.filter(f3).collect()
```



rdd1.map(f1)
transformation is
executed twice

```
rdd2 = rdd1.map(f1)
rdd2.persist()
list1 = rdd2.filter(f2).collect()
list2 = rdd2.filter(f3).collect()
```



Result of rdd1.map(f1)
transformation is cached
and reused (can choose
between memory and
disk for caching)

Spark DataFrames



- Bridging the gap between your experience with Pandas and the need for distributed computing
 - RDD = list of rows
 - DataFrame = table with rows and typed columns
- Important to understand what RDDs are and what they offer, but today most of the tasks can be accomplished with DataFrames (**higher level of abstraction => less code**)
- <https://databricks.com/blog/2015/02/17/introducing-dataframes-in-spark-for-large-scale-data-science.html>

Spark SQL



```
sc = SparkContext()
```

```
sqlContext = HiveContext(sc)
```

```
df = sqlContext.sql("SELECT * from table1 GROUP BY id")
```



Spark's Machine Learning Toolkit

MLlib: Algorithms [[more details](#)]

Classification

- Logistic regression, decision trees, random forests

Regression

- Linear (with L1 or L2 regularization)

Unsupervised:

- Alternating least squares
- K-means
- SVD
- Topic modeling (LDA)

Optimizers

- Optimization primitives (SGD, L-BGFS)

Example:

Logistic regression with MLlib

```
from pyspark.mllib.classification \
    import LogisticRegressionWithSGD

trainData = sc.textFile("...").map(...)
testData = sc.textFile("...").map(...)
model = \
    LogisticRegressionWithSGD.train(trainData)
predictions = model.predict(testData)
```

Remarks

- This lecture is not enough to teach you Spark!
- To use it in practice, you'll need to delve into further online material
- Also: Friday's lab session
- You can't learn it without some frustration :(
- Important skill: assess whether you'd benefit from Spark
 - E.g., >1TB: yes, you'll need Spark
 - 20GB: it depends...



Feedback

Give us feedback on this lecture here:

<https://go.epfl.ch/ada2022-lec13-feedback>

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more (fewer) details?
- Where is Waldo?
- ...

Cluster etiquette

- Develop and debug locally
 - Install Spark locally on your personal computer
 - Use a small subset of the data
- When ready, launch your script on the cluster using [spark-submit](#)
- **Never (never!) use the Spark shell (a.k.a. pyspark) -- it's hereby officially forbidden**
- Useful trench report from a dlab member:
[“What I learned from processing big data with Spark”](#)