

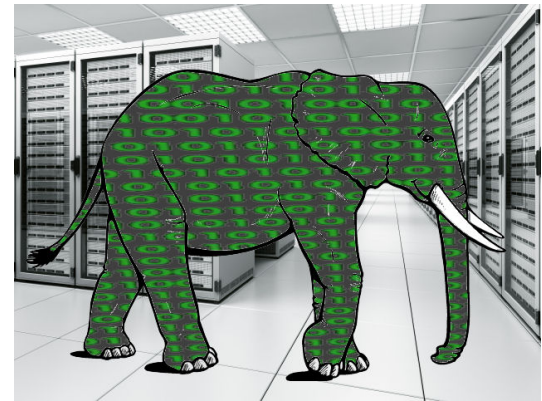
An Introduction to BIG DATA

Prof. Dr. Philippe Cudré-Mauroux

<http://exascale.info>

November 6, 2012

EPFL



Big Data & Me

- My lab @ unifr: eXascale Infolab
 - How to store, manage and query Big Data
 - Teach Big Data at Swiss Joint MSc in CS, HES Lucern and Royal Institute of Tech. (Sweden)
 - But first time *in French!*
 - Previously: M.I.T. Database Systems Lab, EPFL, U.C. Berkeley
 - Industry also (IBM Watson Research, HP, Microsoft Research Asia)



eXascale Infolab



Instant Quizz

- SQL?
- OLAP?
- 3 Vs of Big Data?
- CAP?
- Hadoop?
- Impala?

On the Menu Today

- Big Data: Context
- 3 Vs of Big Data
- Big Data & Dinosaurs
- Hadoop
 - Demo
- The future of Big Data

Exascale Data Deluge

- Science
 - Biology
 - Astronomy
 - Remote Sensing
- Web companies
 - Ebay
 - Yahoo
- Financial services, retail companies, governments, etc.



- ➔ New machines
- ➔ New data formats
- ➔ Peta & exa-scale data sets



Big Data Buzz

Between now and 2015, the firm expects big data to create some **4.4 million IT jobs** globally; of those, 1.9 million will be in the U.S. Applying an economic multiplier to that estimate, Gartner expects each new big-data-related IT job to create work for three more people outside the tech industry, for a total of almost 6 million more U.S. jobs.

Office of Science and Technology
Executive Office of the President
New Executive Office Building
Washington, DC 20502

FOR IMMEDIATE RELEASE
March 29, 2012

Contact: Rick Weiss 202 456-6037 rweiss@ostp
Lisa-Joy Zgorski 703 292-8311 lisajoy@ostp

Growth in the Asia Pacific Big Data market is expected to accelerate rapidly in two to three years time, from a mere US\$258.5 million last year to in excess of **\$1.76 billion in 2016**, with highest growth in the storage segment.

OBAMA ADMINISTRATION UNVEILS “BIG DATA” INITIATIVE: ANNOUNCES \$200 MILLION IN NEW R&D INVESTMENTS

In order to make the most of the fast-growing volume of digital data, the Obama Administration today announced a “Big Data Research and Development Initiative.” Improving our ability to extract knowledge and insights from large and complex collections of digital data, the initiative promises to help solve some of the Nation’s most pressing challenges.

To launch the initiative, six Federal departments and agencies today announced more than \$200 million in new commitments that, together, promise to greatly improve the tools and techniques needed to access, organize, and glean discoveries from huge

Big Data Everywhere!

- The Age of Big Data (NYTimes Feb. 11, 2012)
<http://www.nytimes.com/2012/02/12/sunday-review/big-datas-impact-in-the-world.html>

“GOOD with numbers? Fascinated by data? The sound you hear is opportunity knocking.”

“Welcome to the Age of Big Data. The new megarich of Silicon Valley, first at Google and now Facebook, are masters at harnessing the data of the Web — online searches, posts and messages — with Internet advertising. At the World Economic Forum last month in Davos, Switzerland, Big Data was a marquee topic. A report by the forum, “Big Data, Big Impact,” declared **data a new class of economic asset, like currency or gold.**”

10 ways big data changes everything

- Some concrete examples

- <http://gigaom.com/2012/03/11/10-ways-big-data-is-changing-everything/2/>

1. Can gigabytes predict the next Lady Gaga?
2. How big data can curb the world's energy consumption
3. Big data is now your company's virtual assistant
4. The future of Foursquare is data-fueled recommendations
5. How Twitter data-tracked cholera in Haiti
6. Revolutionizing Web publishing with big data
7. Can cell phone data cure society's ills?
8. How data can help predict and create video hits
9. The new face of data visualization
10. One hospital's embrace of big data



What can you do with the data

© Mike Franklin

- Reporting
 - Post Hoc
 - Real time
- Monitoring (fine-grained)
- Exploration
- Finding Patterns
- Root Cause Analysis
- Closed-loop Control
- Model construction
- Prediction
- ...

The 3-Vs of Big Data

- **V**olume
 - Amount of data
- **V**elocity
 - speed of data in and out
- **V**ariety
 - range of data types and sources
- [Gartner 2012] *"Big Data are high-volume, high-velocity, and/or high-variety information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization"*

More Data => Better Answers?

© Mike Jordan

- Not that easy...
- More Rows: Algorithmic complexity kicks in
- More Columns: Exponentially more hypotheses
- Another formulation of the problem:
 - Given an inferential goal and a fixed computational budget, provide a guarantee that the quality of inference will increase monotonically as data accrue (without bound)
- In other words:
 - => **Data should be a resource, not a load**

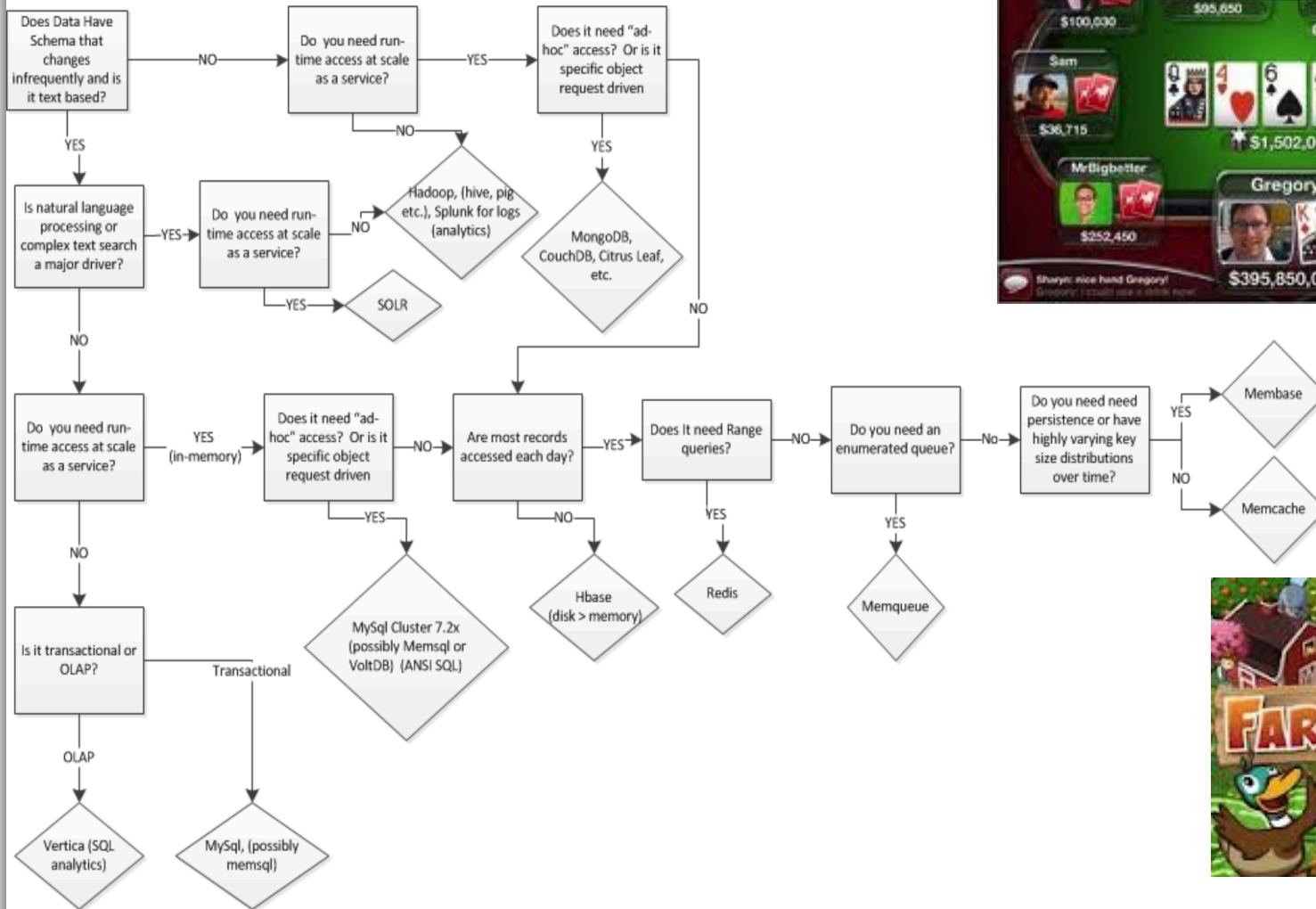
Big Data Today: A Mess

Big Data Landscape (Version 2.0)



A Concrete Example: Zynga

Dan's Scalable Database Decision Matrix at Zynga, 2012



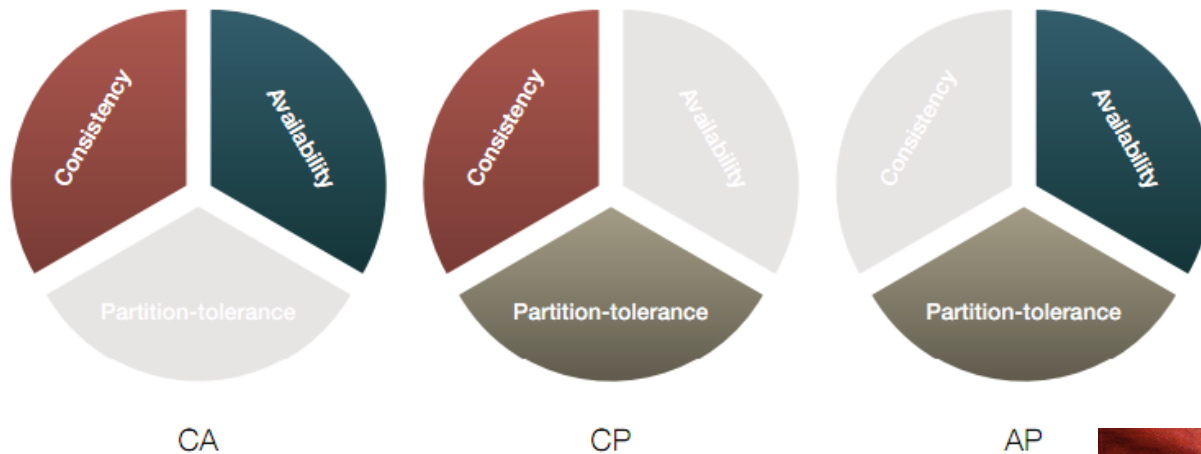
What's wrong with my old DBMS?

- Managing Big Data is hard...
 - ... extremely hard
 - Traditional DBMSs are 30 years old, were not meant for Big Data
 - One user, one CPU, one type of queries
 - Obsolete physical model (n-ary storage, B-trees, etc.)
 - Impractical logical guarantees (transactions, ACID)



What's wrong with my old DBMS?

- Managing big data is hard...
 - ... strictly-speaking, it's actually **impossible**
 - CAP theorem



➔ **Time for a serious makeover**



Leading the Pack of Wolves: Hadoop

- Google: Map/Reduce paper published 2004
- Open source variant: Hadoop
- Map-reduce = high-level programming model and implementation for large-scale parallel data processing
- Right now most overhyped system in CS

Interest over time ?

The number 100 represents the peak search volume

News headlines Forecast ?



A Few MR Numbers @ Google

	Aug. '04	Mar. '06	Sep. '07
Number of jobs (1000s)	29	171	2,217
Avg. completion time (secs)	634	874	395
Machine years used	217	2,002	11,081
map input data (TB)	3,288	52,254	403,152
map output data (TB)	758	6,743	34,774
reduce output data (TB)	193	2,970	14,018
Avg. machines per job	157	268	394
Unique implementations			
map	395	1958	4083
reduce	269	1208	2418

Data Model

- Files !
- A file = a bag of (key, value) pairs
- A map-reduce program:
 - Input: a bag of (input key, value) pairs
 - Output: a bag of (output key, value) pairs

Step 1: the MAP Phase

- User provides the MAP-function:
 - Input: one (input key, value)
 - Output: a bag of (intermediate key, value) pairs
- System applies map function in parallel to all (input key, value) pairs in the input file

Step 2: the REDUCE Phase

- User provides the REDUCE function:
 - Input: intermediate key, and bag of values
 - Output: bag of output values
- System groups all pairs with the same intermediate key, and passes the bag of values to the REDUCE function

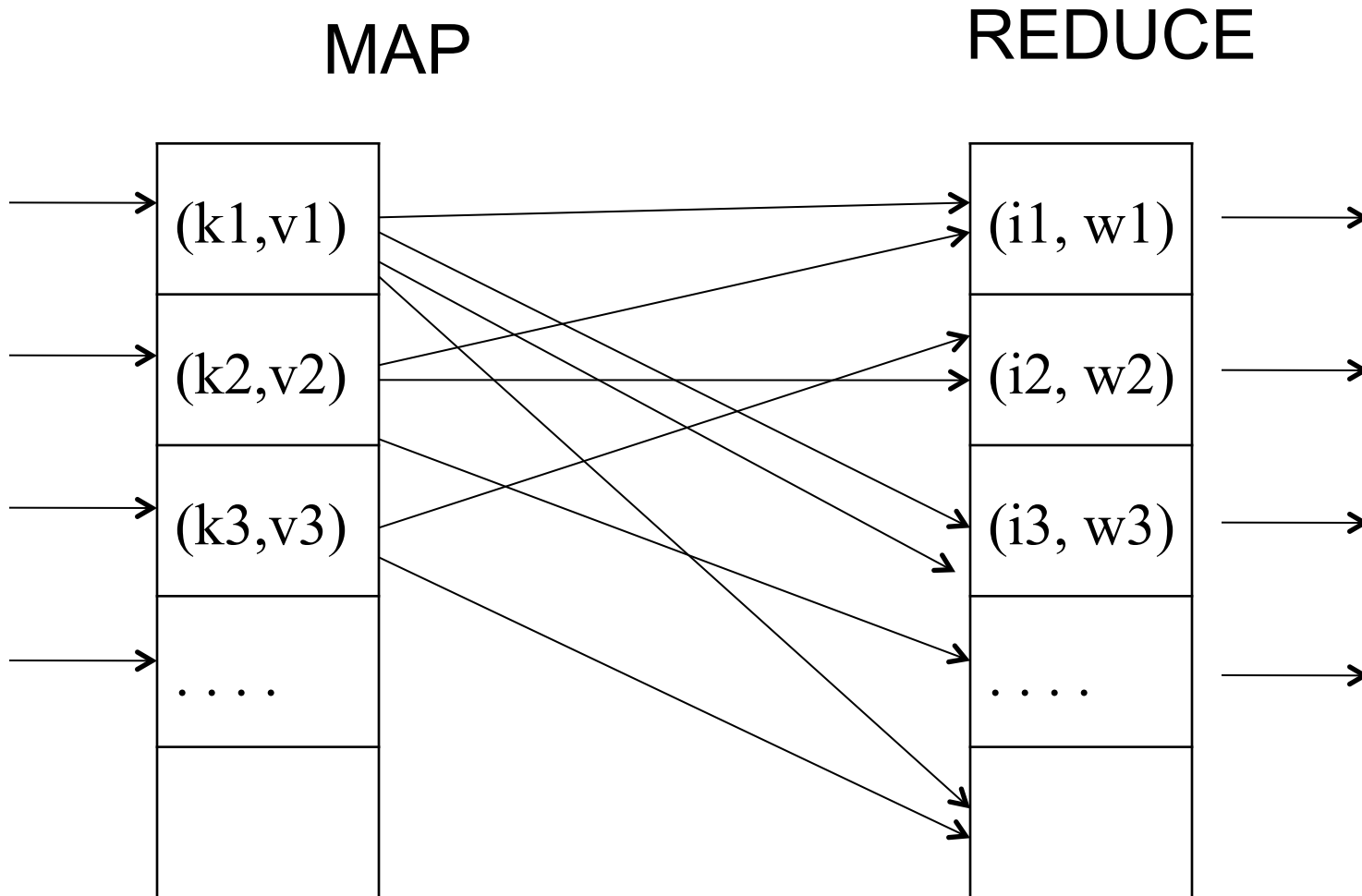
Example

- Counting the number of occurrences of each word in a large collection of documents

```
map(String key, String value):  
  // key: document name  
  // value: document contents  
  for each word w in value:  
    EmitIntermediate(w, "1");
```

```
reduce(String key, Iterator values):  
  // key: a word  
  // values: a list of counts  
  int result = 0;  
  for each v in values:  
    result += ParseInt(v);  
  Emit(AsString(result));
```

MapReduce Execution



Map = GROUP BY, Reduce = Aggregate

R(documentKey, word)

```
SELECT word, sum(1)
FROM R
GROUP BY word
```

Example: MR word length count

Abridged Declaration of Independence

A Declaration By the Representatives of the United States of America, in General Congress Assembled. When in the course of human events it becomes necessary for a people to advance from that subordination in which they have hitherto remained, and to assume among powers of the earth the equal and independent station to which the laws of nature and of nature's god entitle them, a decent respect to the opinions of mankind requires that they should declare the causes which impel them to the change.

We hold these truths to be self-evident; that all men are created equal and independent; that from that equal creation they derive rights inherent and inalienable, among which are the preservation of life, and liberty, and the pursuit of happiness; that to secure these ends, governments are instituted among men, deriving their just power from the consent of the governed; that whenever any form of government shall become destructive of these ends, it is the right of the people to alter or to abolish it, and to institute new government, laying it's foundation on such principles and organizing it's power in such form, as to them shall seem most likely to effect their safety and happiness. Prudence indeed will dictate that governments long established should not be changed for light and transient causes: and accordingly all experience hath shewn that mankind are more disposed to suffer while evils are sufferable, than to right themselves by abolishing the forms to which they are accustomed. But when a long train of abuses and usurpations, begun at a distinguished period, and pursuing invariably the same object, evinces a design to reduce them to arbitrary power, it is their right, it is their duty, to throw off such government and to provide new guards for future security. Such has been the patient sufferings of the colonies; and such is now the necessity which constrains them to expunge their former systems of government. the history of his present majesty is a history of unremitting injuries and usurpations, among which no one fact stands single or solitary to contradict the uniform tenor of the rest, all of which have in direct object the establishment of an absolute tyranny over these states. To prove this, let facts be submitted to a candid world, for the truth of which we pledge a faith yet unsullied by falsehood.

Example: MR word length count

Abridged Declaration of Independence

Map Task 1
(204 words)

A Declaration By the Representatives of the United States of America, in General Congress Assembled.
When in the course of human events it becomes necessary for a people to advance from that subordination in which they have hitherto remained, and to assume among powers of the earth the equal and independent station to which the laws of nature and of nature's god entitle them, a decent respect to the opinions of mankind requires that they should declare the causes which impel them to the change.
We hold these truths to be self-evident; that all men are created equal and independent; that from that equal creation they derive rights inherent and inalienable, among which are the preservation of life, and liberty, and the pursuit of happiness; that to secure these ends, governments are instituted among men, deriving their just power from the consent of the governed; that whenever any form of government shall become destructive of these ends, it is the right of the people to alter or to abolish it, and to institute new government, laying it's foundation on such principles and organizing it's power in such form, as to them shall seem most likely to effect their safety and happiness. Prudence indeed will

(key, value)

(yellow, 17)

(red, 77)

(blue, 107)

(pink, 3)

Yellow: 10+

Red: 5..9

Blue: 2..4

Pink: = 1

Map Task 2
(190 words)

dictate that governments long established should not be changed for light and transient causes: and accordingly all experience hath shewn that mankind are more disposed to suffer while evils are sufferable, than to right themselves by abolishing the forms to which they are accustomed. But when a long train of abuses and usurpations, begun at a distinguished period, and pursuing invariably the same object, evinces a design to reduce them to arbitrary power, it is their right, it is their duty, to throw off such government and to provide new guards for future security. Such has been the patient sufferings of the colonies; and such is now the necessity which constrains them to expunge their former systems of government. the history of his present majesty is a history of unremitting injuries and usurpations, among which no one fact stands single or solitary to contradict the uniform tenor of the rest, all of which have in direct object the establishment of an absolute tyranny over these states. To prove this, let facts be submitted to a candid world, for the truth of which we pledge a faith yet unsullied by falsehood.

(yellow, 20)

(red, 71)

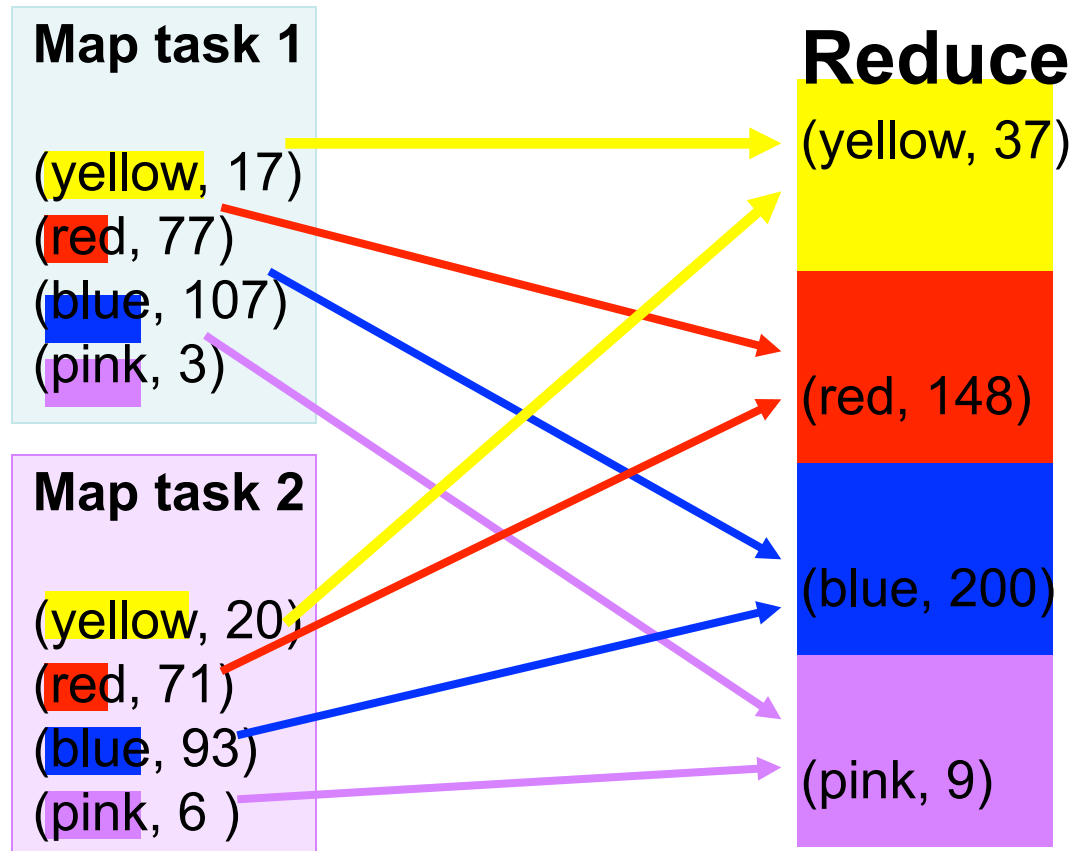
(blue, 93)

(pink, 6)

Example: MR word length count

Map is a **GROUP BY** operation

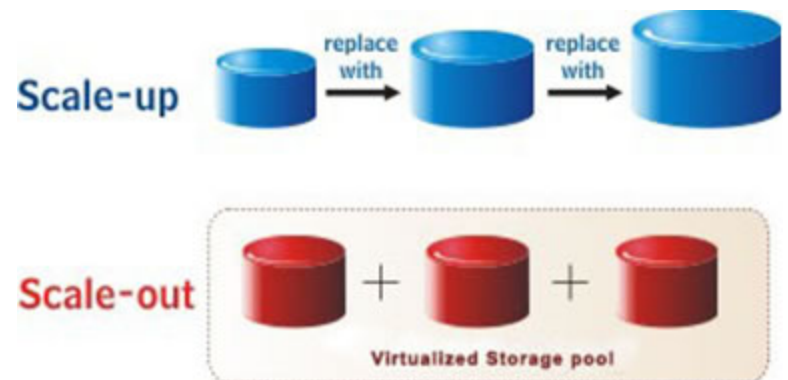
Reduce is an **AGGREGATE** operation



Map/Reduce / Hadoop Limitations

- Plenty of limitations...
 - Simplistic data model
 - Extremely impractical language (Map / Reduce)
 - No data/process affinity
 - No pipelining
 - Batch only
 - Slow
 - [...]

- ... but it works!
 - i.e., allows **scale-out**



Hadoop Demo



Services ▾

Edit ▾

Philippe Cudre-Mauroux ▾

N. Virginia ▾

Help ▾

Your Elastic MapReduce Job Flows

Create New Job Flow Terminate Debug

Show/Hide Refresh Help

Viewing: All

1 to 1 of 1 Job Flows

	Name	State	Creation Date	Elapsed Time	Normalized Instance Hours
<input checked="" type="checkbox"/>	My Job Flow	TERMINATED	2012-11-05 20:57 GMT+0	0 hours 9 minutes	3

1 Job Flow selected

Job Flow: j-11JBQQWPCPVMH

Last State Change: Terminated by user request

Description

Steps

Bootstrap Actions

Instance Groups

Monitoring

Name: My Job Flow

Creation Date: 2012-11-05 20:57 GMT+0100

Start Date: 2012-11-05 21:03 GMT+0100

End Date: 2012-11-05 21:11 GMT+0100

Higher-Level Tools

Scripting language:

Query engines:

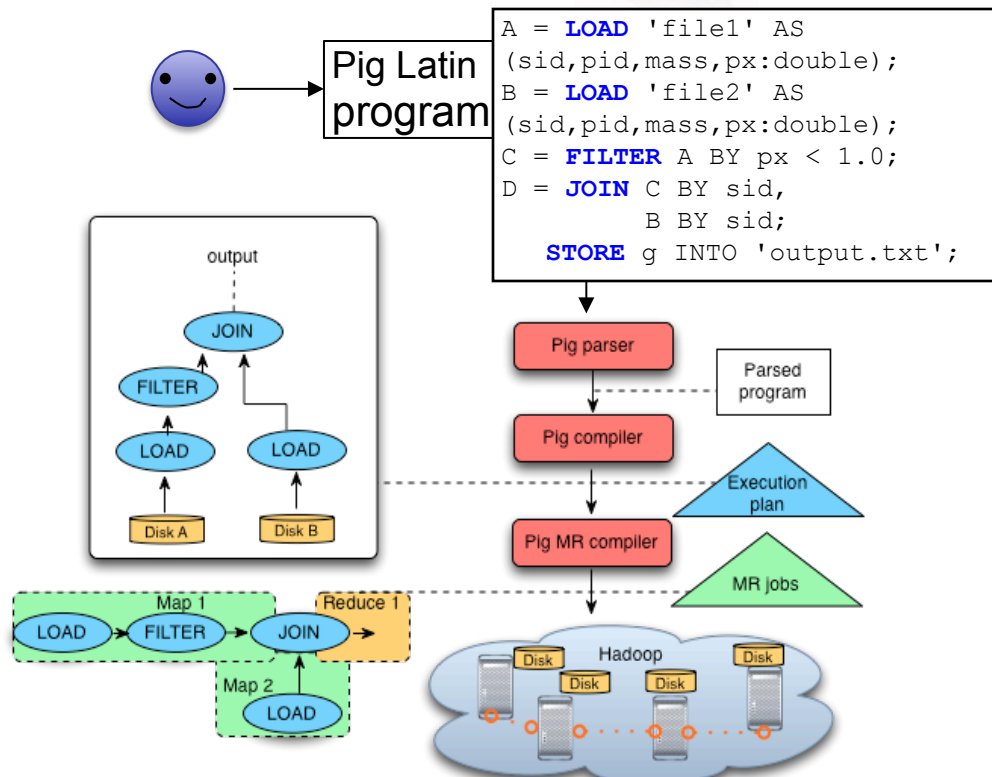
Pig Latin



pig



HIVE



The Future of Big Data?

- The end of one-size-fits-all
- Diversification of tools
 - Relational DBMSs are here to stay
 - Premium database vendors (teradata, vertica, etc.)
 - Post map/reduce solutions
 - Yarn, Impala (Dremel? Percolator?)
 - SAAS Cloud databases (Amazon, Google, Microsoft)
 - Stream data management (Truviso, Storm)
 - Data integration (Virtuoso, SAP, Oracle)
 - Key/value (Cassandra), Document (CouchDB), Graph (neo4J), Array (SciDB), [...] database systems

⇒ **Countless** ~~problems~~ **opportunities** in the coming years

What about Big Swiss Data?

- Thanks for your attention!
- Questions?

