▷ Yogesh Simmhan
▷ simmhan@iisc.ac.in

▷ Department of Computational and Data Sciences
▷ Indian Institute of Science, Bangalore

DS256 (3:1)
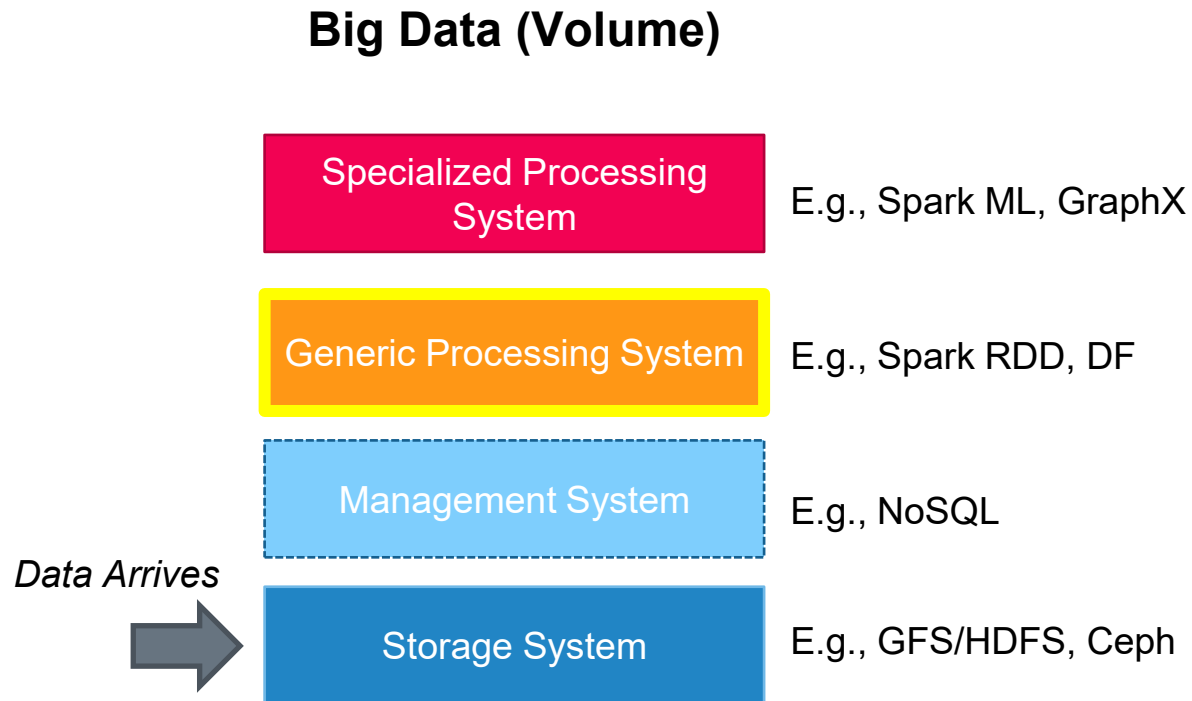
# Scalable Systems for Data Science

# Processing Large Volumes of Big Data

# Role of Big Data Processing System for Large Data Volumes

**Big Data (Volume)**

| | |
|---|---|
| Specialized Processing System | E.g., Spark ML, GraphX |
| Generic Processing System | E.g., Spark RDD, DF |
| Management System | E.g., NoSQL |
| Storage System | E.g., GFS/HDFS, Ceph |

*Data Arrives*

# Required Reading

"

- **_MapReduce: Simplified Data Processing on Large Clusters_**, *Dean and Ghemawat, USENIX OSDI, 2004*

- **_Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing_**, *Matei Zaharia, et al., USENIX NSDI, 2012*

- **_Spark SQL: Relational Data Processing in Spark_**, *Michael Armbrust, et al., ACM SIGMOD 2015*

- *Select chapters from **Learning Spark**, Holden Karau, et al., 1st Editions and **Learning Spark**, Jules S. Damji, Brooke Wenig, Tathagata Das, Denny Lee, 2nd Edition*

- **Optional: Clash of the Titans: MapReduce vs. Spark for Large Scale Data Analytics**, Juwei Shi, Yunjie Qiu, Umar Farooq Minhas, Limei Jiao Chen Wang, Berthold Reinwald, and Fatma Ozcan, *VLDB* 2015

# MapReduce

MapReduce: Simplified Data Processing on Large Clusters, Dean and Ghemawat, *USENIX OSDI*, 2004

# MapReduce

*"A **simple and powerful interface** that enables **automatic parallelization** and distribution of **large-scale computations**, combined with an implementation of this interface that achieves high performance on **large clusters** of **commodity PCs**."*

Dean and Ghermawat, "MapReduce: Simplified Data Processing on Large Clusters", OSDI, 2004

# Google's Cluster Model (2004)

▷ Machines have 2 x86 CPU, running Linux, with 2-4 GB of memory
  ○ 2020: Raspberry Pi 4B has 4-core ARM, 2-4GB RAM, Rs.5000 ☺

▷ Commodity networking – 100 Mbps or 1 Gbps per machine, low bisection bandwidth.

▷ A cluster has 100's – 1000's of machines, failures are common

▷ Storage on inexpensive IDE disks attached to servers
  ○ GFS manages the data

▷ Users submit jobs to a scheduling system
  ○ Each job is a set of tasks
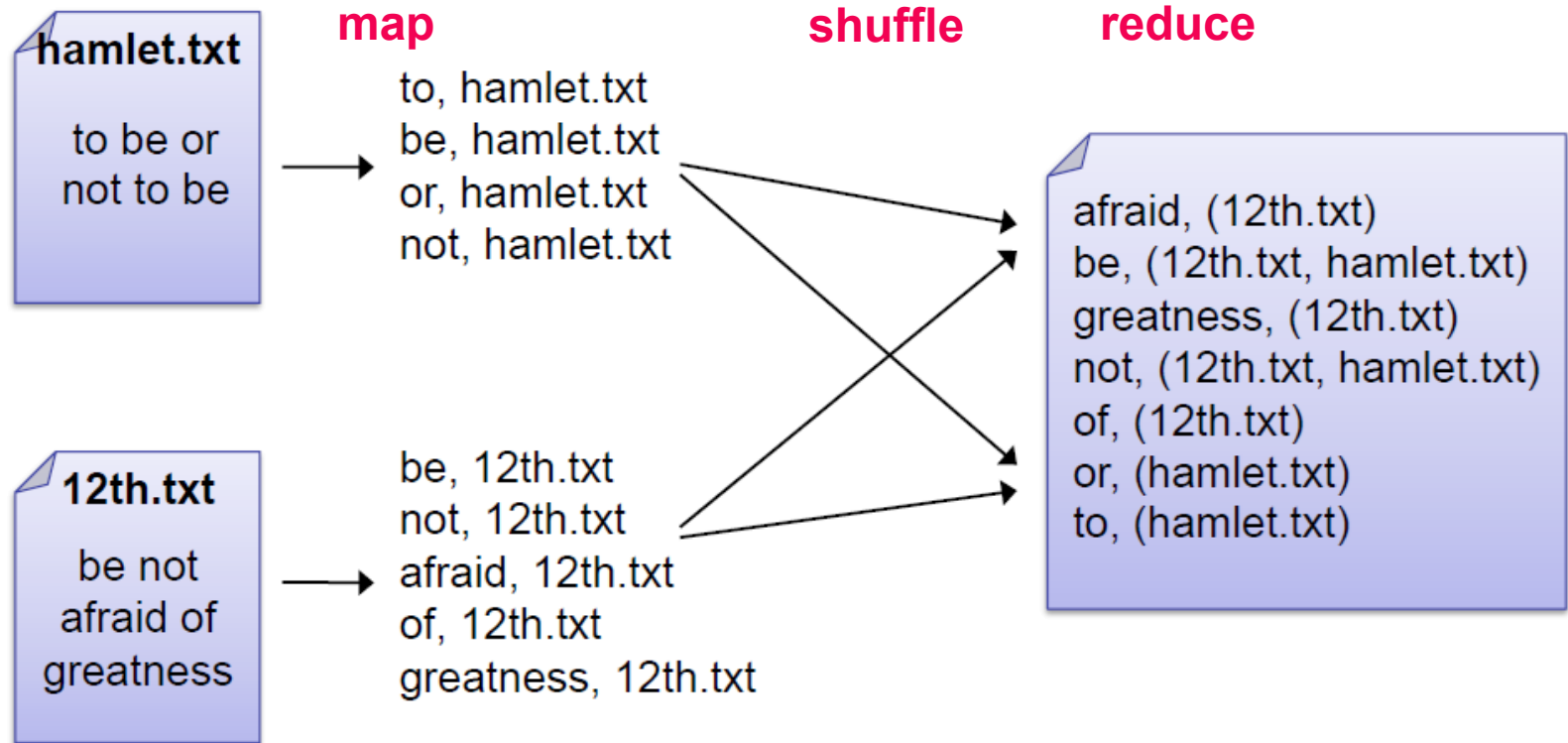  ○ Mapped by the scheduler to available machines in cluster

# MapReduce Design Pattern

▷ Programming model for distributed applications
  o Clean abstraction for programmers
  o Automatic parallelization & distribution

▷ Fault-tolerance

▷ Batch data processing system
  o Large inputs sizes

▷ Simple data-intensive applications
  o Distributed Grep: Document list ➔ Occurrence of search term
  o URL Access Frequency: URL access list ➔ <URL, freq>
  o Reverse Web-Link Graph: <target,src> ➔ <src, target[]>
  o Term-Vector per Host: <host,word[]> ➔ <host,<word,freq>[]>
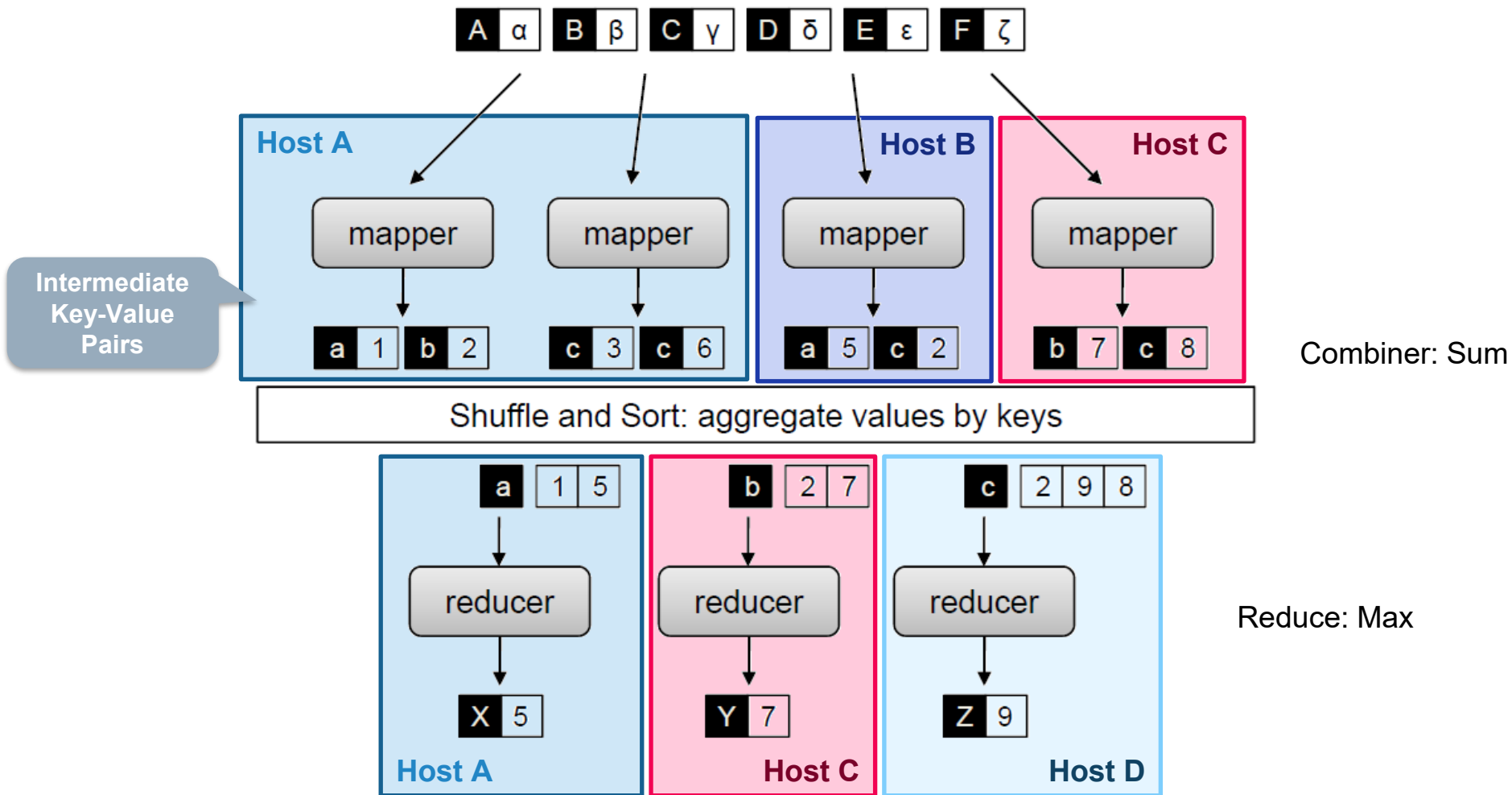
# MapReduce: Data-parallel Programming Model

▷ Process data using **map** & **reduce** user-defined functions

▷ `map(k`$_i$`, v`$_i$`)` → `List<k`$_m$`, v`$_m$`>`
   - *map* is called once on every input item
   - Emits a series of intermediate key/value pairs

▷ **`shuffle & sort phase`**
   - All map output values (**v**$_m$) with a given key (**k**$_m$) are *grouped* together, keys *sorted* within a group
   - Happens internally within the framework

▷ `reduce(k`$_m$`, List<v`$_m$`>)` → `List<k`$_r$`, v`$_r$`>`
   - *reduce* is called once on *every unique key & all its values*
   - Emits a value that is added to the output

# Inverted Index using MR



- ▷ map(url, line)
  - ○ *foreach*(**word** *in* **line.split())** *emit*(**word, url**)
- ▷ reduce(word, url[])
  - ○ *save*(**word, findDistinct(url[])**)

# **Map**-*Shuffle*-*Sort*-**Reduce**

# Histogram using MR



```
7    2    11    2
2    1    11    4
9   10     6    6
6    3     2    8
0    5     1   10
2    4     8   11
5    0     1    0
```

M   M   M   M

```
1,1  0,1  2,1  0,1
0,1  0,1  2,1  1,1
2,1  2,1  1,1  1,1
1,1  0,1  0,1  2,1
0,1  1,1  0,1  2,1
0,1  1,1  2,1  2,1
1,1  0,1  0,1  0,1
```

**Shuffle**

```
2,1  0,1 0,1  1,1
2,1  0,1 0,1  1,1
2,1  0,1 0,1  1,1
2,1  0,1 0,1  1,1
2,1  0,1 0,1  1,1
2,1  0,1 0,1  1,1
2,1           1,1
2,1           1,1
```

R   R   R

```
2,8   0,12   1,8
```

*Data transfer & shuffle* between Map & Reduce **(28 items)**
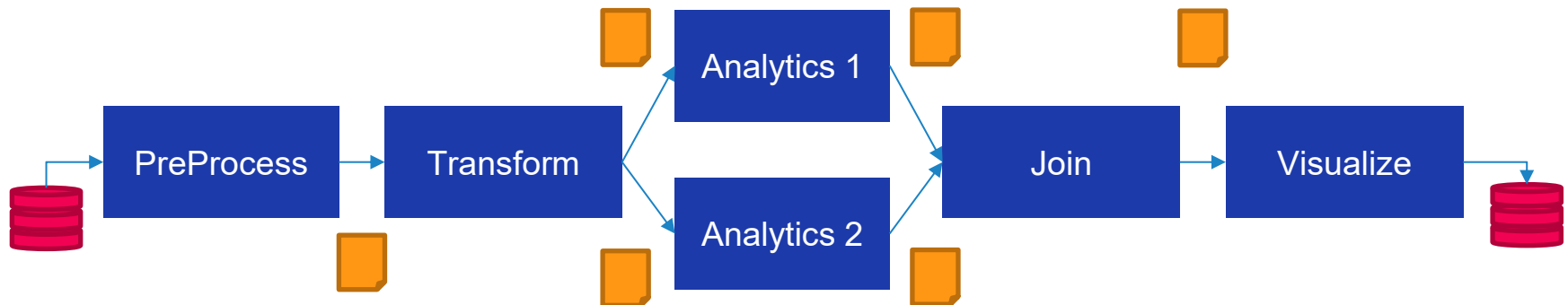
```
int bucketWidth = 4 // input

Map(k, v) {
    emit(floor(v/bucketWidth), 1)
    // <bucketID, 1>
}




// one reduce per bucketID
Reduce(k, v[]){
    sum=0;
    foreach(n in v[])  sum+=n;
    emit(k, sum)
    // <bucketID, frequency>
}
```
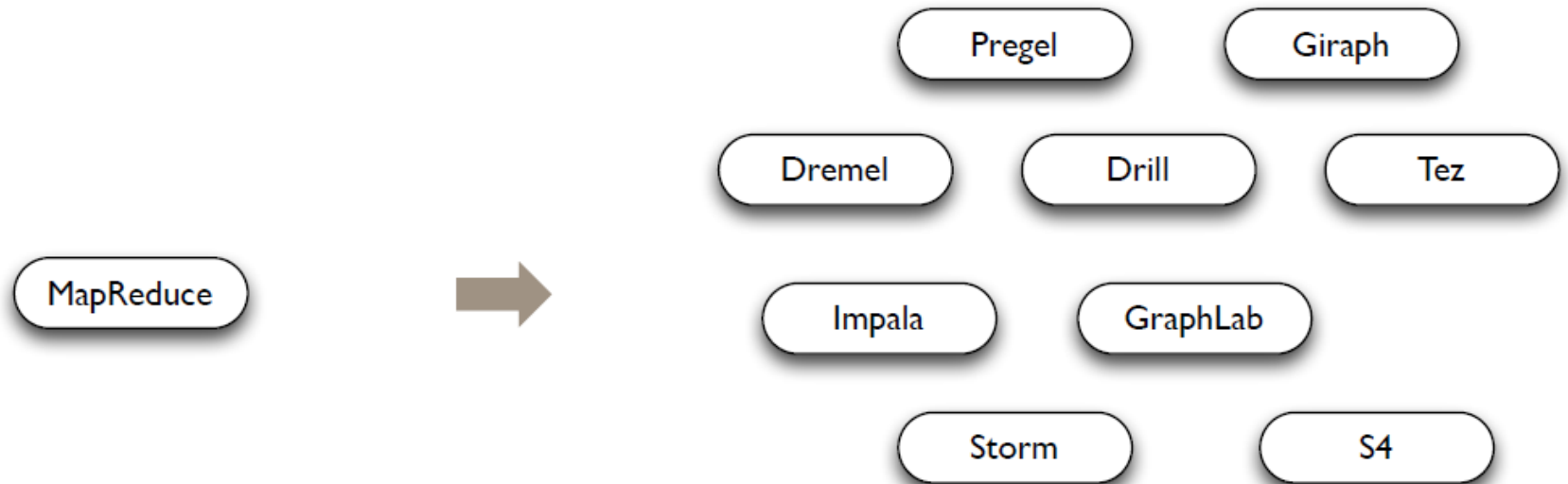
# Limitations of MapReduce

▷ Multi-stage computing not simple
  o Many different jobs

▷ Complex code for simple transformations
  o Repetitive, not *data centric*

https://stanford.edu/~rezab/sparkclass/slides/itas_workshop.pdf

# Limitations of MapReduce

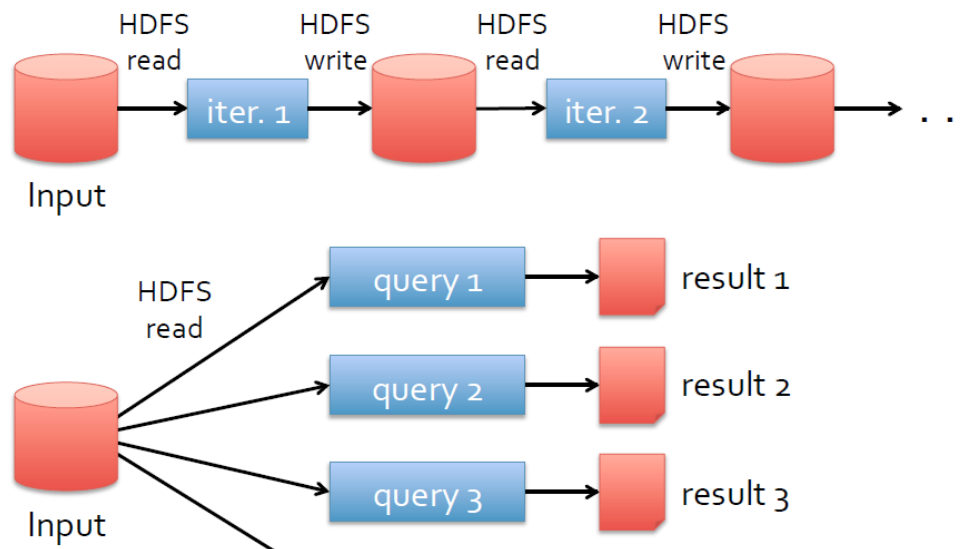▷ Limited support for non-text, Non-static data



**General Batch Processing**

**Specialized Systems:**
iterative, interactive, streaming, graph, etc.

# Limitations of MapReduce

▷ Poor performance for:
  - o Complex, multi--stage applications (e.g. iterative machine learning & graph processing)
  - o Interactive *ad hoc* queries

© Matei Zaharia

# Latency & Bandwidth

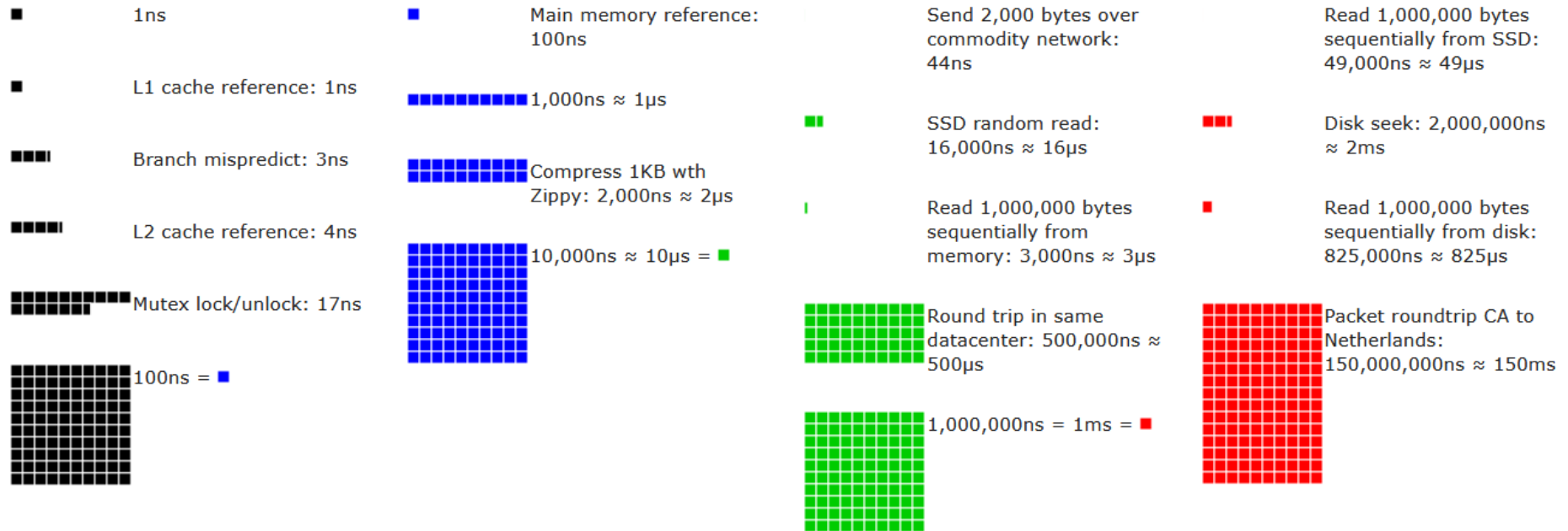- ▷ L1 cache reference
- ▷ L2 cache reference
- ▷ **Main memory reference**
- ▷ **Send 1K bytes over 1 Gbps network**
- ▷ **Read 4K randomly from SSD***
- ▷ **Read 1MB sequentially from memory**
- ▷ **Round trip within same datacenter**
- ▷ **Read 1MB sequentially from SSD***
- ▷ **Send 1MB over 1 Gbps network**
- ▷ **Disk seek**
- ▷ **Read 1MB sequentially from disk**
- ▷ **Send packet CA->NL->CA**

# Latency & Bandwidth

| | | | | |
|---|---|---|---|---|
| ▷ | L1 cache reference | 0.5 ns | | |
| ▷ | L2 cache reference | 7 ns | | |
| ▷ | **Main memory reference** | **100 ns** | | |
| ▷ | **Send 1K bytes over 1 Gbps network** | **10,000 ns** | **10 μs** | |
| ▷ | **Read 4K randomly from SSD*** | **150,000 ns** | **150 μs** | |
| ▷ | **Read 1MB sequentially from memory** | **250,000 ns** | **250 μs** | |
| ▷ | **Round trip within same datacenter** | **500,000 ns** | **500 μs** | |
| ▷ | **Read 1MB sequentially from SSD*** | **1,000,000 ns** | **1,000 μs** | **1 ms** |
| ▷ | **Send 1MB over 1 Gbps network** | **8,250,000 ns** | **8,250 μs** | **8 ms** |
| ▷ | **Disk seek** | **10,000,000 ns** | **10,000 μs** | **10 ms** |
| ▷ | **Read 1MB sequentially from disk** | **20,000,000 ns** | **20,000 μs** | **20 ms** |
| ▷ | **Send packet CA->NL->CA** | **150,000,000 ns** | **150,000 μs** | **150 ms** |

https://gist.github.com/jboner/2841832

**Latency Numbers Every Programmer Should Know**

2020

■ 1ns

■ L1 cache reference: 1ns

■■■■ Branch mispredict: 3ns

■■■■■ L2 cache reference: 4ns

■■■■■■■■■■■ Mutex lock/unlock: 17ns

100ns = ■

■ Main memory reference: 100ns

■■■■■■■■■■ 1,000ns ≈ 1μs

■■■■■■■■■■ Compress 1KB wth Zippy: 2,000ns ≈ 2μs

10,000ns ≈ 10μs = ■

Send 2,000 bytes over commodity network: 44ns

■■ SSD random read: 16,000ns ≈ 16μs

I Read 1,000,000 bytes sequentially from memory: 3,000ns ≈ 3μs

Round trip in same datacenter: 500,000ns ≈ 500μs

1,000,000ns = 1ms = ■

Read 1,000,000 bytes sequentially from SSD: 49,000ns ≈ 49μs

■■■■ Disk seek: 2,000,000ns ≈ 2ms

■ Read 1,000,000 bytes sequentially from disk: 825,000ns ≈ 825μs

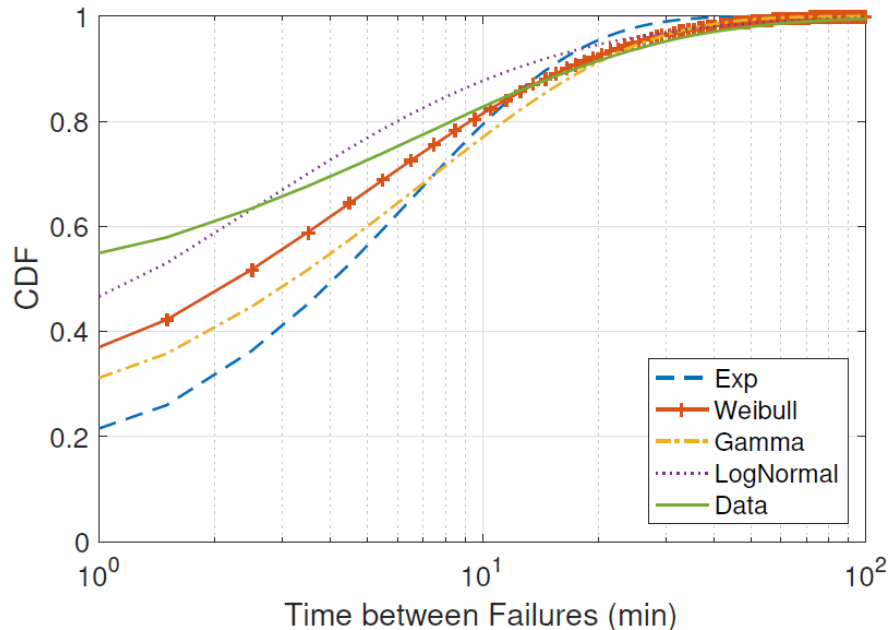Packet roundtrip CA to Netherlands: 150,000,000ns ≈ 150ms

"

*Bandwidth of Memory ≫ Network or Disk*

# MTTF in Data Center

*"The MTBF (mean time between failures) across all data centers we investigate (with hundreds of thousands of servers) is only 6.8 minutes, while the MTBF in different data centers varies between 32 minutes and 390 minutes."*

→ **MTBF with 1000 servers is 680mins**
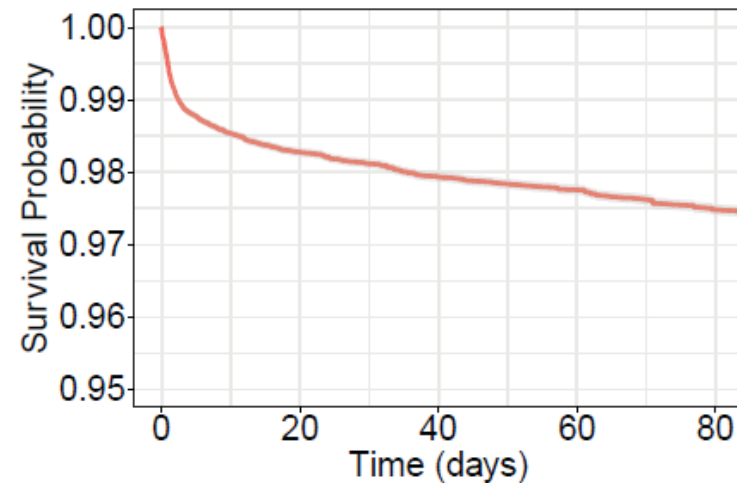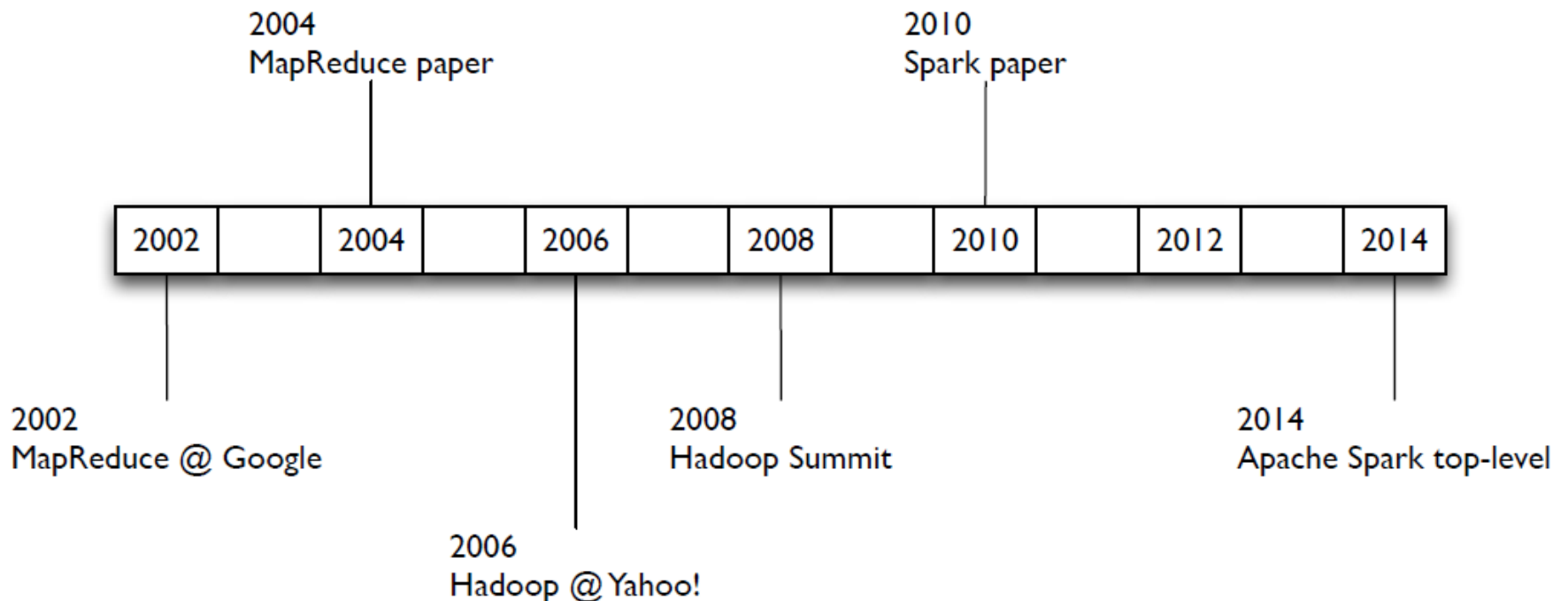→ **MTBF with 100 servers is 6800mins (4.7 days)**





**Figure 4: Kaplan Meier survival estimate of datacenter switches, shaded region shows 95% confidence intervals.**

Surviving switch failures in cloud datacenters, ACM SIGCOMM Computer Communication Review, Volume 51Issue 2April 2021
What Can We Learn from Four Years of Data Center Hardware Failures?, DSN, 2017

20

"

*Failures may be infrequent during the lifetime of an application execution*

# From MapReduce to Spark

▷ ## Google's MapReduce
  o Programming Model
  o Apache Hadoop runtime environment



2004
MapReduce paper

2010
Spark paper

| 2002 | | 2004 | | 2006 | | 2008 | | 2010 | | 2012 | | 2014 |

2002
MapReduce @ Google

2008
Hadoop Summit

2014
Apache Spark top-level

2006
Hadoop @ Yahoo!

https://stanford.edu/~rezab/sparkclass/slides/itas_workshop.pdf

# Apache Spark

**Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing,** Zaharia, Chowdhury, Das, Dave, Ma, McCauly, Franklin, Shenker and Stoica, *USENIX NSDI*, 2012

**Learning Spark**
Holden Karau, Andy Konwinski, Patrick Wendell & Matei Zaharia, O'Reilly, First and Second Editions

https://databricks.com/p/ebook/learning-spark-from-oreilly

# The Spark Ecosystem

▷ **Core Spark Engine**
  o RDDs, Transformations, Actions, batch processing

▷ **Higher level abstractions**
  o Data frames, SQL-like queries
  o Discretized streams, semi-realtime data
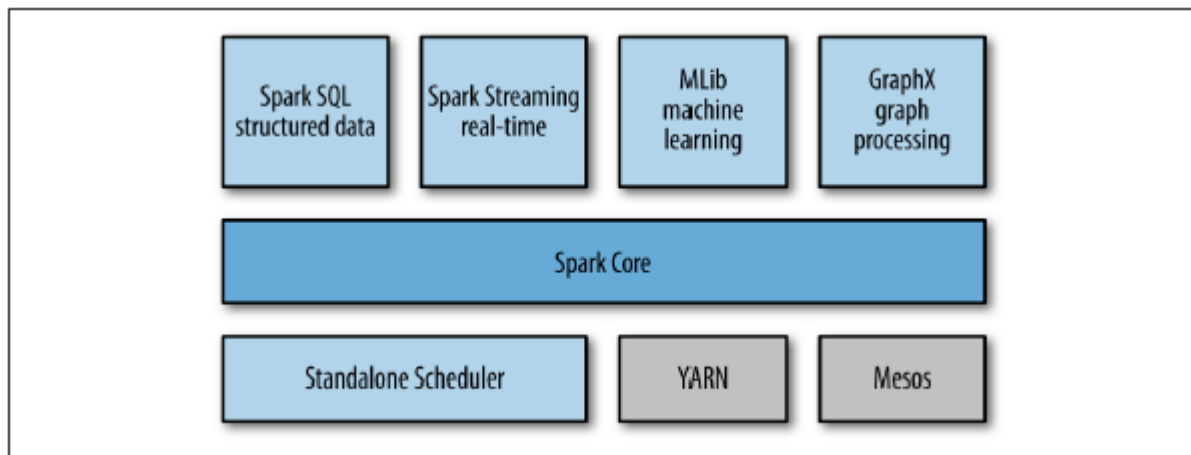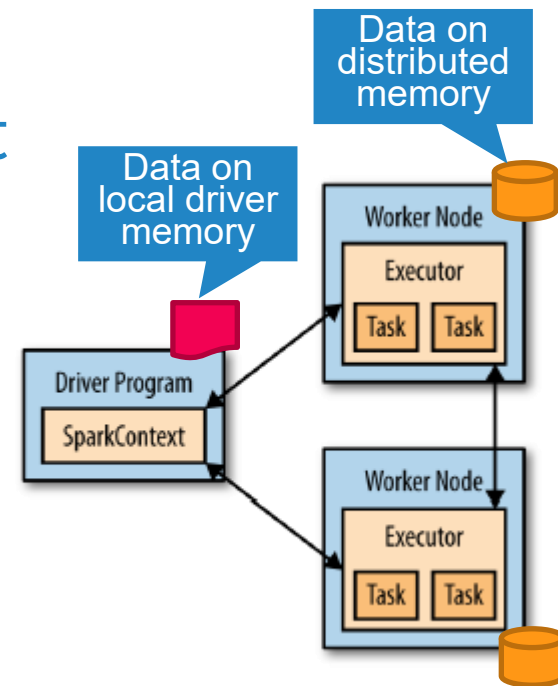  o Machine learning libraries, **Mllib**
  o Linked data analytics, **GraphX**



*Figure 1-1. The Spark stack*
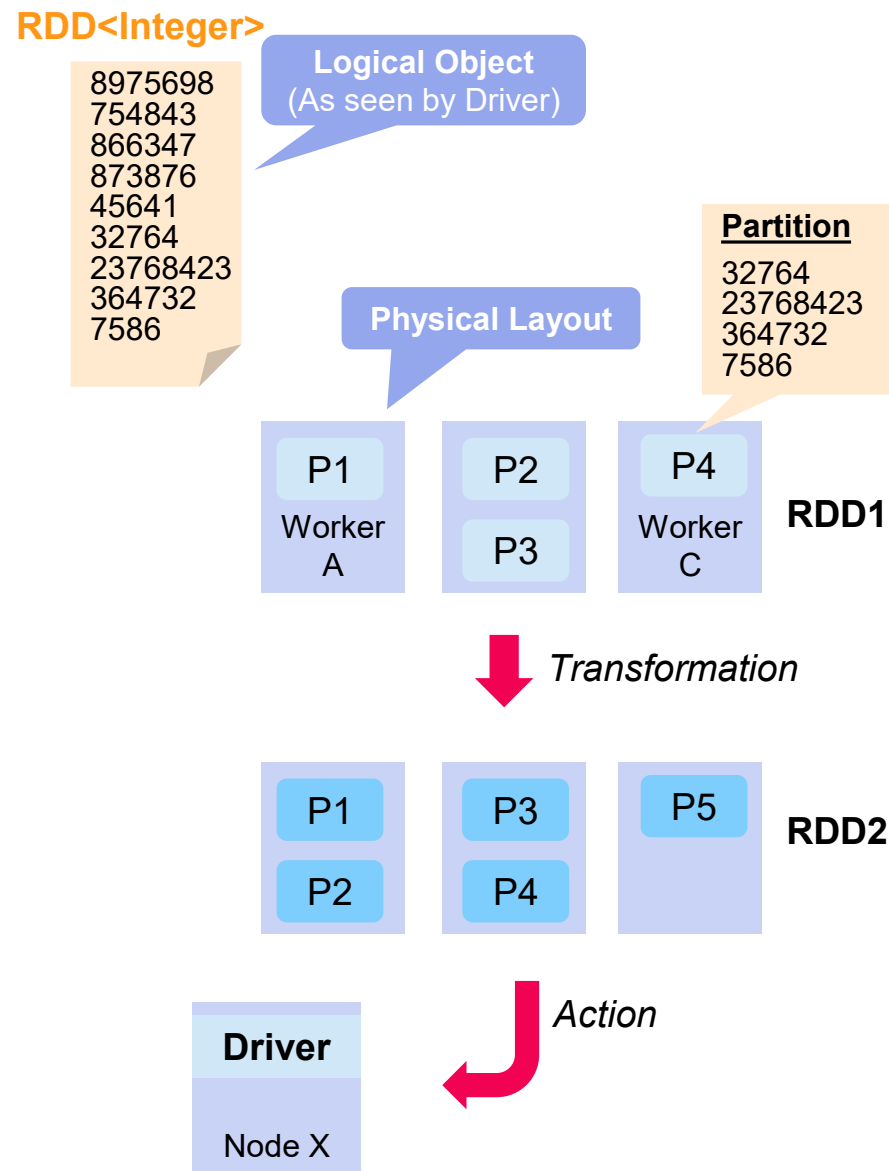
# Spark: A Distributed Execution Engine

▷ **Driver**: User program for application, uses Spark Context, local variables

▷ **Spark Context**: Gives access to distributed computing environment

▷ **Worker**: Machines on which actual heavy-lift happens

▷ **Executor**: Spark execution environment in a worker, Process, exclusive to an application

▷ **Task**: Single operation on data, thread



Data on distributed memory

Data on local driver memory

Worker Node
Executor
Task    Task

Driver Program
SparkContext

Worker Node
Executor
Task    Task

# Spark RDD
## Resilient Distributed Dataset

▷ **Collection of homo-geneous objects**
  - ○ Order is not preserved*

▷ **Distributed** on workers
  - ○ 1 or more **Partitions**

▷ **Read-only**, immutable

▷ Can be **rebuilt**

▷ Can be **cached**

▷ MR like data-parallel operations
  - ○ **Execute** on workers

**RDD<Integer>**

8975698
754843
866347
873876
45641
32764
23768423
364732
7586

**Logical Object**
(As seen by Driver)

**Partition**

32764
23768423
364732
7586

**Physical Layout**

| P1 | P2 | P4 |
| Worker A | P3 | Worker C |

**RDD1**

*Transformation*

| P1 | P3 | P5 |
| P2 | P4 | |

**RDD2**

*Action*

**Driver**

Node X

# Creating and Operating on an RDD

▷ Users can provide driver code in multiple languages: *Scala, Java, **Python (PySpark*)**
  - o Spark offers equivalent transformations and actions in each language

▷ Actual Spark execution environment is in Scala

▷ Create RDD by loading data
  - o HDFS, local vars, filesystem, NoSQL DB, etc.
  - o Data is loaded on partitions on different workers

▷ RDD Object offers a logical view of the dataset
  - o Can perform operations on the object

*Example 3-1. Creating an RDD of strings with textFile() in Python*

```
>>> lines = sc.textFile("README.md")
```

# Creating and Operating on an RDD

▷ Users can provide driver code in multiple languages: *Scala, Java, **Python (PySpark**)*
- o Spark offers equivalent transformations and actions in each language

▷ Actual Spark execution environment is in Scala

▷ Create RDD by loading data
- o HDFS, local vars, filesystem, NoSQL DB, etc.
- o Data is loaded on partitions on different workers

▷ RDD Object offers a logical view of the dataset
- o Can perform operations on the object

*Example 3-1. Creating an RDD of strings with textFile() in Python*

```
>>> lines = sc.textFile("README.md")
```

# Language Bindings

▷ Users can provide driver code in multiple languages
  - Scala, Java, Python

▷ Spark offers equivalent transformations and actions in each language

▷ Logic within transformations and actions can also be in these languages

▷ Actual Spark execution environment is in Scala
  - Standard data structure mapping
  - Python code is pickled (de/serialize) and shipped remotely

https://cwiki.apache.org/confluence/display/SPARK/PySpark+Internals

# Passing Functions to Spark Operations

▷ Lambda syntax
  - o **Functions** are *input parameters* to other functions
  - o Pass short functions concisely, inline

```python
pythonLines = lines.filter(lambda line: "Python" in line)
```

```python
def hasPython(line):
    return "Python" in line

pythonLines = lines.filter(hasPython)
```

# Programming with RDDs

**Learning Spark**

Holden Karau, Andy Konwinski, Patrick Wendell & Matei Zaharia,
O'Reilly, First Edition

**Chapter 3**

# Basics of Transformations

▷ Returns a new RDD, computed lazily

▷ Transforms tend to be **element-wise** operations
  o **Iterate** through each item, apply the operation, e.g.
    *Filter*

▷ **Filter** on *inputRDD* does not affect *inputRDD*
  o Returns a new RDD, *warningsRDD*

▷ **Union** operates on two RDDs
  o One of them is an input parameter

```python
inputRDD = sc.textFile("log.txt")
errorsRDD = inputRDD.filter(lambda x: "error" in x)
warningsRDD = inputRDD.filter(lambda x: "warning" in x)
badLinesRDD = errorsRDD.union(warningsRDD)
```

# Basics of Actions

▷ *Actually* triggers operations, returns a final result to driver
  o Force any required transformations to be executed
  o Count, Take, Collect

▷ Result of action must fit in memory of driver
  o Else, can write RDD to HDFS, saveAsTextFile

▷ RDDs are computed from scratch when actions are called...See *persist/cache*

```
print "Input had " + badLinesRDD.count() + " concerning lines"
print "Here are 10 examples:"
for line in badLinesRDD.take(10):
    print line
```

# Lazy Evaluation

▷ Transformations are lazily evaluated
  - o Calling a transform does NOT immediately execute it
▷ *Action* **triggers** execution of *dependent transformations*
▷ E.g., load().map().count()
  - o Load & Map do not execute till we see Count
▷ Allows Spark to reduce the number of passes through the data
  - o Materializes RDD only when required
  - o Reused RDDs that have been materialized earlier
  - o Immutability!

# Lineage Graph

▷ Keeps track of operations used to derive an RDD

▷ Helps *lazily materiali[ze]* RDD

▷ Helps *recover* RDD or their partitions that are lost