# Scalding:
## Big Data Programming with Scala

**Taewook Eom**

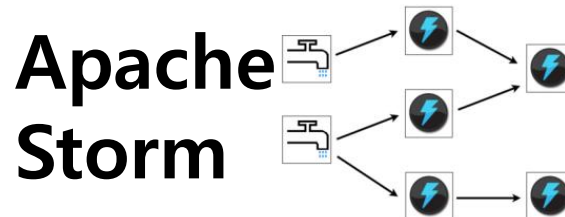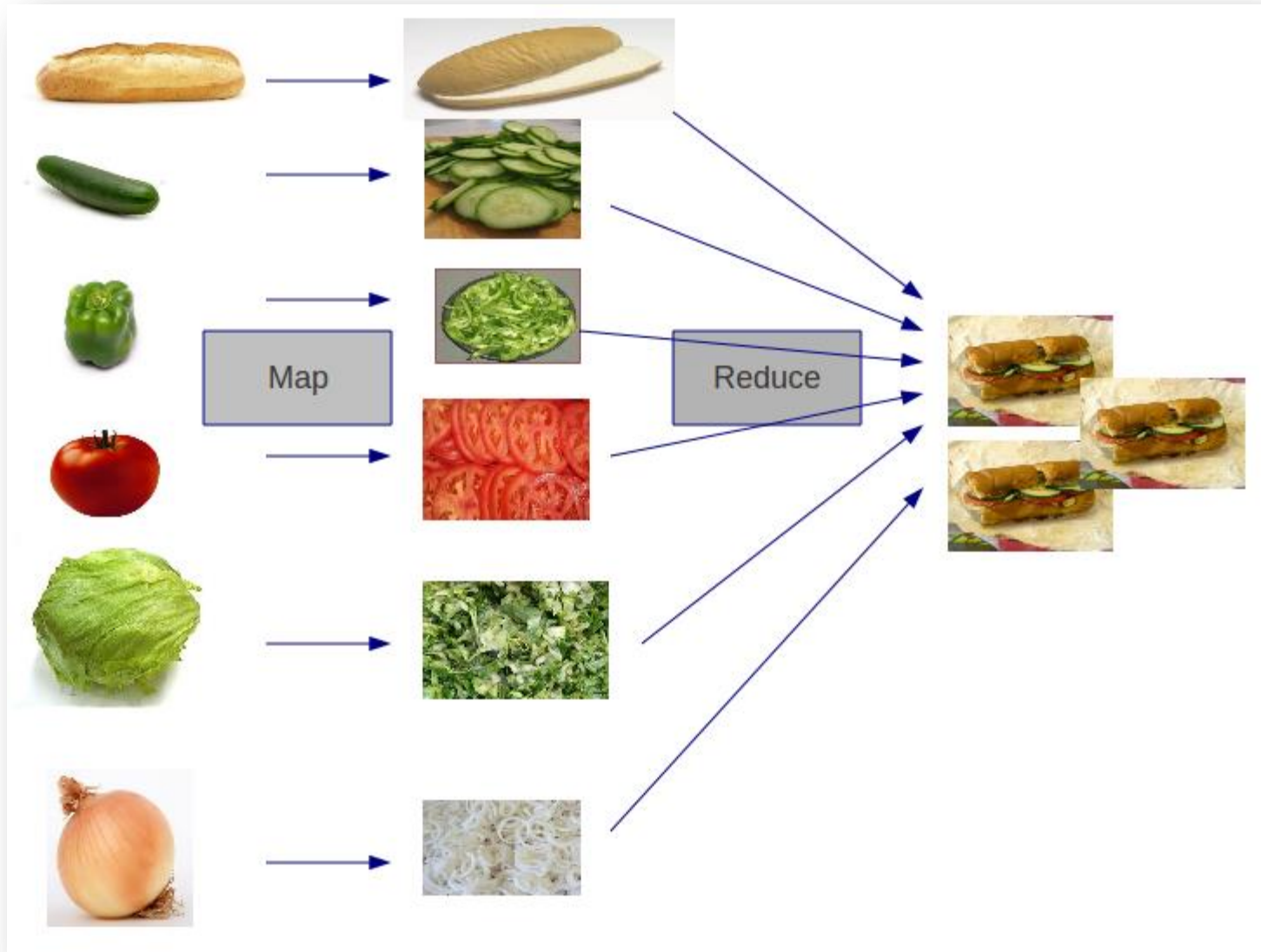Data Infrastructure Team
SK planet

taewook@sk.com

2015-01-29

# Big Data Processing
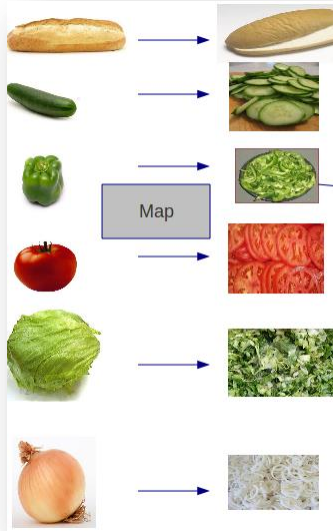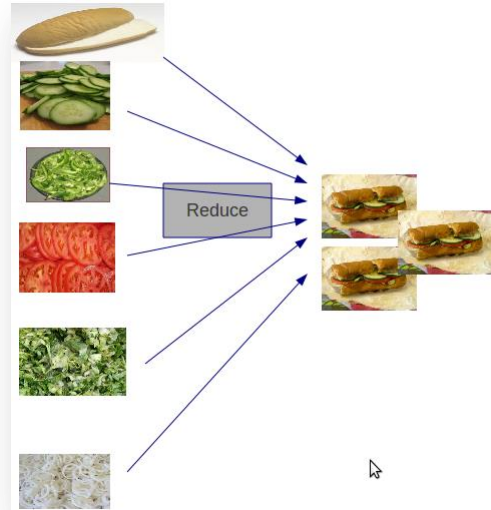
# MapReduce, MR, Map-Reduce

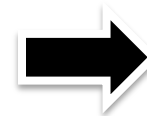# Data Processing Pattern with MR



select
function
where(filter)

group by
Join
order by
windowing function
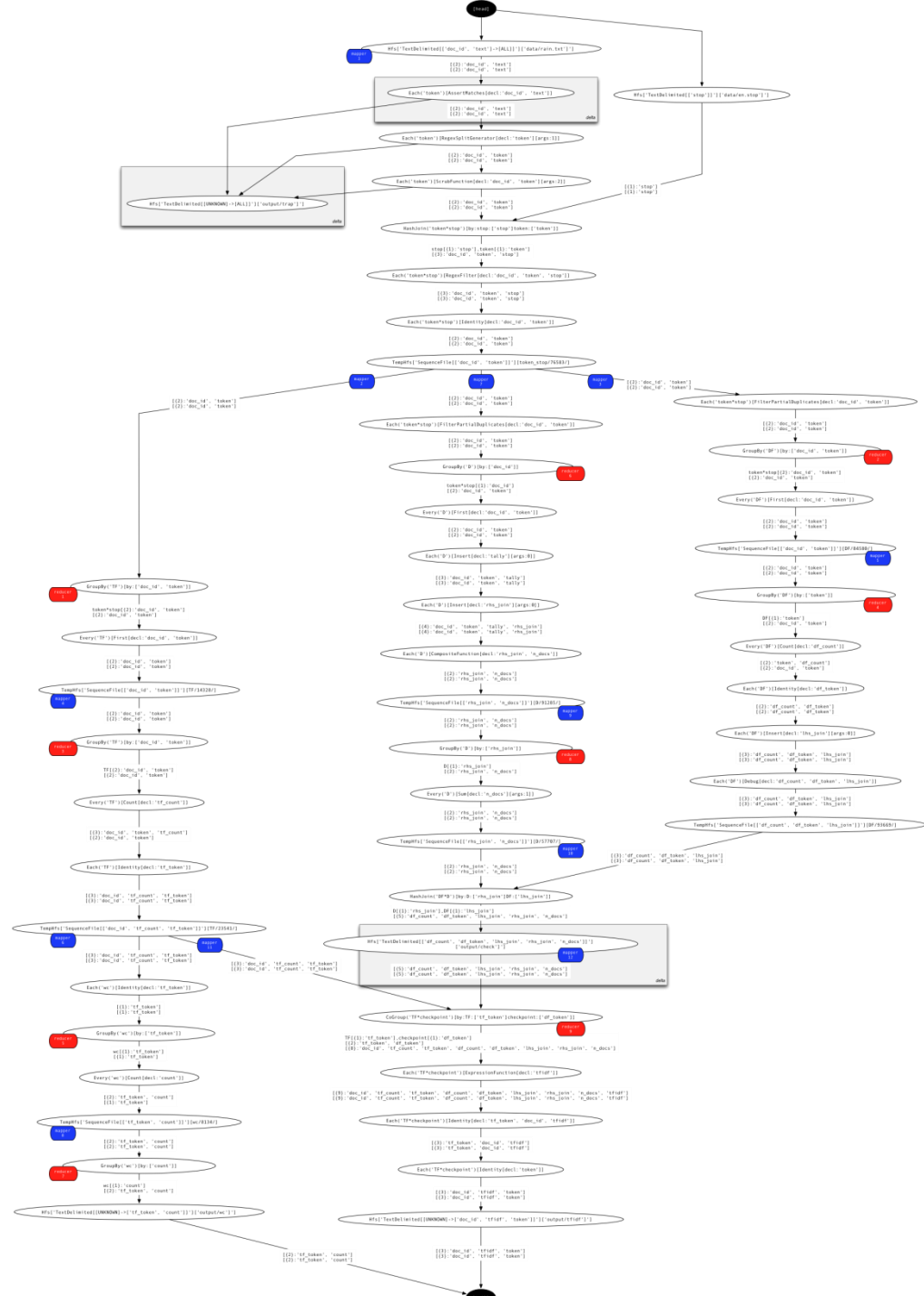analytics function

$$M = M+$$
$$MR = M+RM*$$
$$MRMR... = (M+RM*)+$$
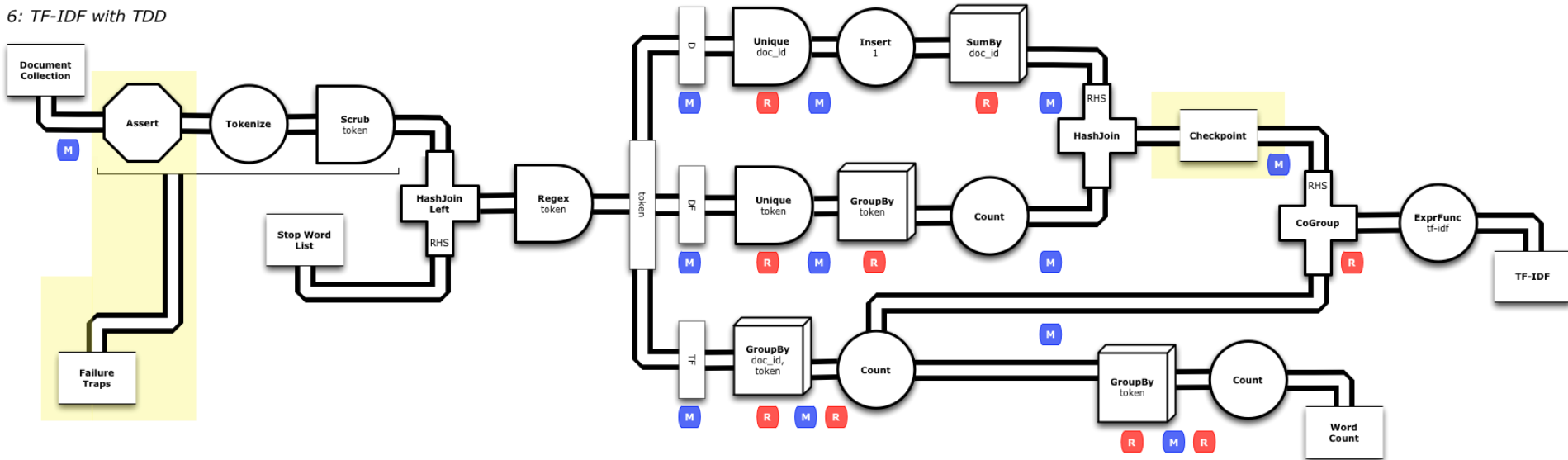
➡ **Workflow** management

# Data Workflow = DAG

## (Directed Acyclic Graph)

# Cascading

6: TF-IDF with TDD

- **Pipe** abstraction = Plumbing
- Operators like **SQL**
- **DAG** based **workflow** management

# http://www.slideshare.net/taewook/programming-cascading

# Object-Oriented vs. Functional

**OOP** focuses on the **differences in the data**

  **Data** and the **operations** upon it are tightly **coupled**
  The central model for abstraction is the **data** itself

**FP** concentrates on **consistent data structures**

  **Data** is only loosely coupled to functions
  The central model for abstraction is the **function**, not the data structure

**FP describe what** they want done, **not how** to do it
**OOP** uses mostly **imperative** techniques

```
1  var sumOfSquares = pipe(map(square), reduce(add, 0));
2
3  console.log(sumOfSquares([2, 3, 5]));
```

```
1  var sumOfSquares = function(list) {
2      var result = 0;
3      for (var i = 0; i < list.length; i++) {
4          result += square(list[i]);
5      }
6      return result;
7  };
8
9  console.log(sumOfSquares([2, 3, 5]));
```

# Data Processing, Functional Programming

**SQL**    uses a **consistent data structure** (table: rows x cols)
uses **functions** that can be **combined**
is **declarative** not imperative

## Data is Immutable

## ➜ Transformable

## by Composable Functions

# Why



http://www.scala-lang.org/

**Sca**lable **La**nguage
　　　　➔ **Big Data**

Seamless **Java Interop**
　➔ **Hadoop** runs on the **JVM**

**Functional**
　　　　➔ **Data Processing**

**REPL**(Read-Evaluate-Print Loop)
　➔ **Interactive** data analysis

# Scalding

**Scala DSL** for Cascading

**Simple** and **concise** syntax

maintained by Twitter

```java
public class Main {
    public static void main(String[] args) {
        String docPath = args[0];
        String wcPath = args[1];
        String stopPath = args[2];

        Properties properties = new Properties();
        AppProps.setApplicationJarClass(properties, Main.class);
        FlowConnector flowConnector = new Hadoop2MR1FlowConnector(properties);

        Tap docTap = new Hfs(new TextDelimited(true, "\t"), docPath);
        Tap wcTap = new Hfs(new TextDelimited(true, "\t"), wcPath);

        Fields stop = new Fields("stop");
        Tap stopTap = new Hfs(new TextDelimited(stop, true, "\t"), stopPath);

        Fields token = new Fields("token");
        Fields text = new Fields("text");
        RegexSplitGenerator splitter = new RegexSplitGenerator(token, "[ \\[\\]\\(\\),.]");
        Fields fieldSelector = new Fields("doc_id", "token");
        Pipe docPipe = new Each("token", text, splitter, fieldSelector);

        Fields scrubArguments = new Fields("doc_id", "token");
        docPipe = new Each(docPipe, scrubArguments, new ScrubFunction(scrubArguments), Fields.RESULTS);

        Pipe stopPipe = new Pipe("stop");
        Pipe tokenPipe = new HashJoin(docPipe, token, stopPipe, stop, new LeftJoin());
        tokenPipe = new Each(tokenPipe, stop, new RegexFilter("^$"));

        Pipe wcPipe = new Pipe("wc", tokenPipe);
        wcPipe = new Retain(wcPipe, token);
        wcPipe = new GroupBy(wcPipe, token);
        wcPipe = new Every(wcPipe, Fields.ALL, new Count(), Fields.ALL);

        FlowDef flowDef = FlowDef.flowDef().setName("wc")
                .addSource(docPipe, docTap).addSource(stopPipe, stopTap)
                .addTailSink(wcPipe, wcTap);

        Flow wcFlow = flowConnector.connect(flowDef);
        wcFlow.writeDOT("dot/wc.dot");
        wcFlow.complete();
    }
}
```
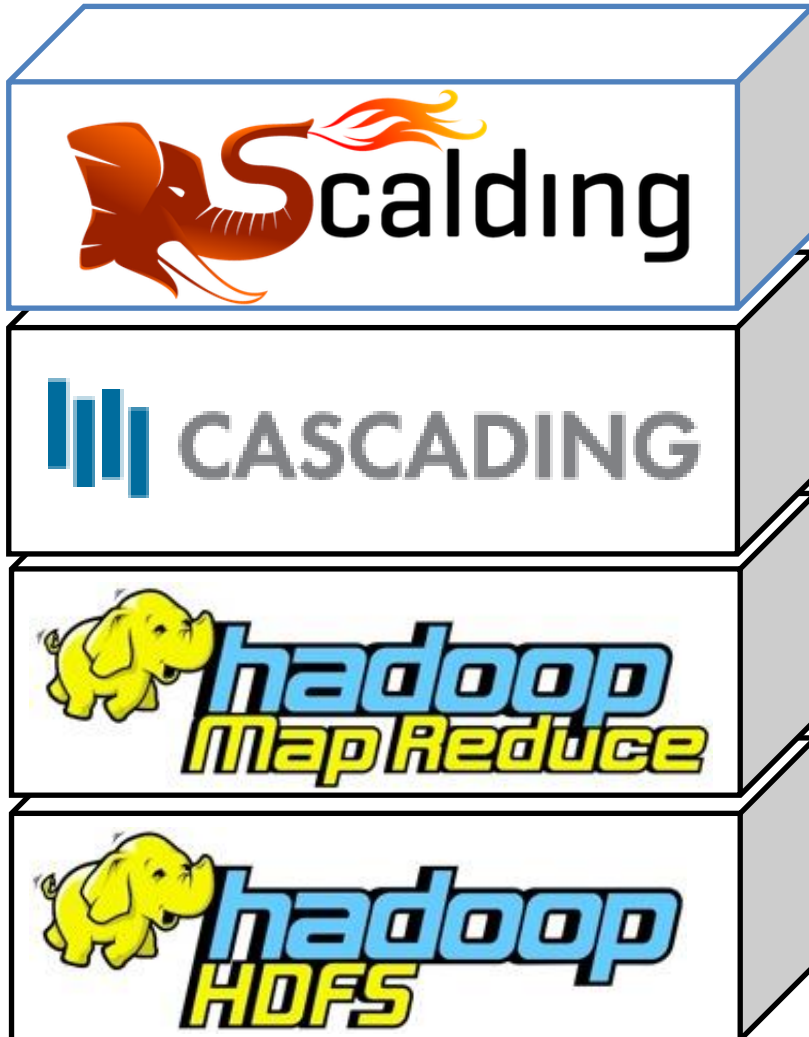
```scala
import ...

object Part4 {
  def main(args: Array[String]) {
    new Part4(Args(List("--local", "", "--input", "data/rain.txt",
        "--output", "data/output.txt", "--stop", "data/en.stop"))).run
  }
}

class Part4(args: Args) extends Job(args) {

  def scrub(text: String): String = {
    text.trim.toLowerCase.replaceAll( """[\[\]\(\),-]""", " ")
  }

  val input = Tsv(args("input"), ('docId, 'text))
  val output = Tsv(args("output"))
  val stop = Tsv(args("stop"), 'stopword).read

  input.read
    .mapTo('text -> 'stext) { text: String => scrub(text)}
    .flatMap('stext -> 'word) { stext: String => stext.split( """\s+""")}
    .project('word)
    .joinWithSmaller('word -> 'stopword, stop, joiner = new LeftJoin)
    .filter('stopword) { stopword: String => stopword == null || stopword.isEmpty}
    .groupBy('word) { group => group.size}
    .write(output)
}
```

```java
public class ScrubFunction extends BaseOperation implements Function {
    public ScrubFunction(Fields fieldDeclaration) {
        super(2, fieldDeclaration);
    }

    public void operate(FlowProcess flowProcess, FunctionCall functionCall) {
        TupleEntry argument = functionCall.getArguments();
        String doc_id = argument.getString(0);
        String token = scrubText(argument.getString(1));

        if (token.length() > 0) {
            Tuple result = new Tuple();
            result.add(doc_id);
            result.add(token);
            functionCall.getOutputCollector().add(result);
        }
    }

    public String scrubText(String text) {
        return text.trim().toLowerCase();
    }
}
```

# **UDF**(User-defined Function)

# "If you need to write UDF's all the time, something is wrong with you."

- Various authors of non-scalding frameworks who happened to be completely WRONG

# Etsy's Data-Driven Culture

At **Etsy**, it's not just **engineers** who **write** and **deploy code** – our **designers** and **product managers** regularly do too.

## Data is for Everyone

- Every person in product is a data producer
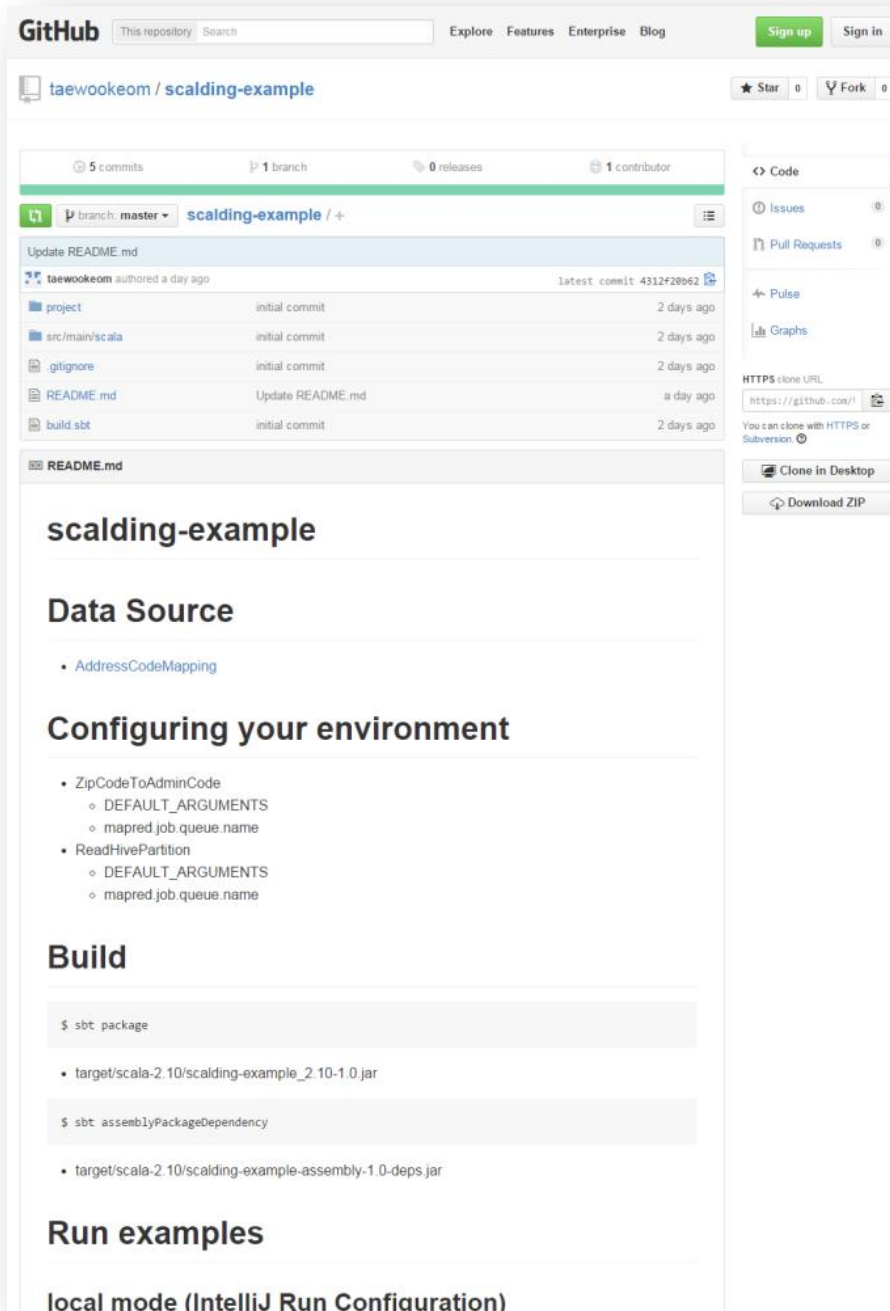- Every person in the company CAN BE a data consumer

## Learned to get data

- Wrote a scalding job to get the data
- Looked at a full month of data to check for consistency

## Why does this matter?

- Supports a more inclusive culture, welcoming people from all over the company
- If you can answer your own questions, you are more free to ask questions than if you rely on others
- Empowers product managers, developers, and designers, marketers, merchandisers, etc to be data-driven

# https://github.com/taewookeom/scalding-example



## SBT Build script
- build.sbt, project/plugins.sbt
- libraryDependencies
- Main-Class in META-INF/MANIFEST.MF

## Splitting **project** and **deps JARs**

## Run **command** and **arguments**

# Next Try



[https://spark.apache.org/](https://spark.apache.org/)

**Apache Spark™** is a fast and general engine for large-scale data processing.



Prabhakar Gopalan
@PGopalan
Follow

@strataconf summary so far after the series of keynotes: Hadoop disappears, Spark is everywhere

10:56 PM - 16 Oct 2014 · Manhattan, NY, United States

2 FAVORITES



Andre Luckow
@drelu
Follow

Full house with people standing during Spark talk. Next year it can be called Spark-World. #hadoopworld #strataconf pic.twitter.com/cBjJ934OYC

2:52 AM - 18 Oct 2014 · Manhattan, NY, United States

3 RETWEETS  9 FAVORITES

# Questions?

Questions.foreach( answer(_) )
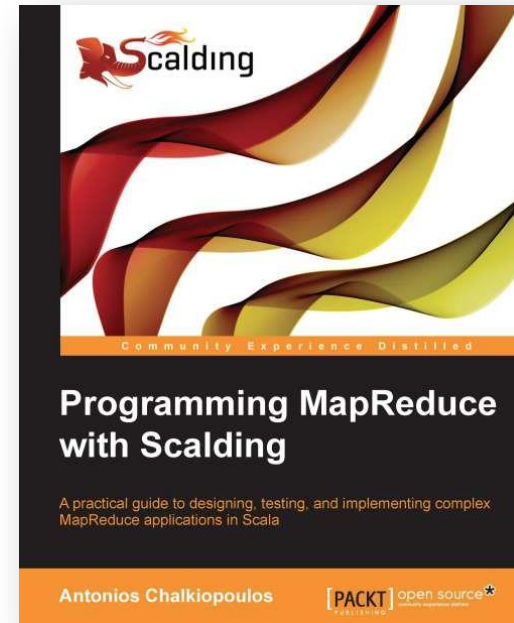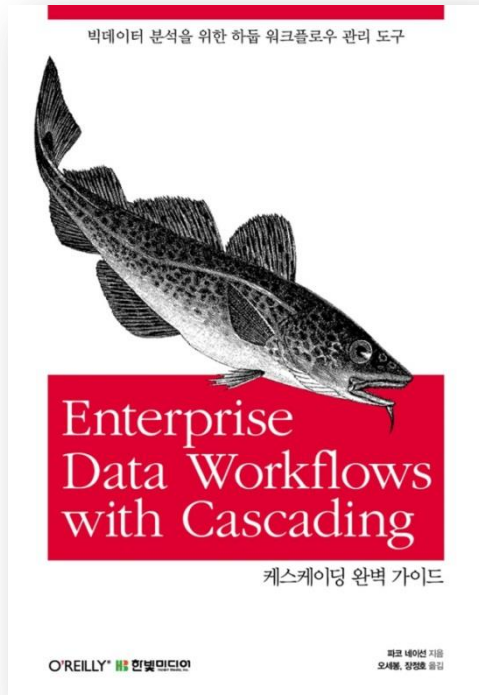
# Learning Scala





http://www.slideshare.net/deview/a4de-view2012-scalamichinisougu Scala, 미지와의 조우
http://www.slideshare.net/kthcorp/scala-15041890 꽃보다 Scala http://goo.gl/O382Fh
https://twitter.github.io/scala_school/ko/index.html 스칼라 학교!
http://refcardz.dzone.com/refcardz/scala Refcardz: Getting Started with Scala
http://wrobstory.gitbooks.io/python-to-scala/ Python To Scala
http://mbonaci.github.io/scala/ Java developer's Scala cheatsheet

# Learning Scalding





http://docs.cascading.org/tutorials/scalding-data-processing/
https://github.com/twitter/scalding/wiki/Getting-Started
https://github.com/twitter/scalding/wiki/Fields-based-API-Reference
https://github.com/twitter/scalding/tree/master/tutorial
https://github.com/scalding-io/ProgrammingWithScalding
http://sujitpal.blogspot.kr/2012/08/scalding-for-impatient.html
https://github.com/snowplow/scalding-example-project

Michael Feathers
@mfeathers

OO makes code understandable by encapsulating moving parts. FP makes code understandable by minimizing moving parts.

12:27 AM - 4 Nov 2010

212 RETWEETS  96 FAVORITES

데이터 프로그래머는 사고 방식이 프로그램에 대한 흐름제어 중심에서 데이터에 대한 함수 중심으로 바뀌어야 한다. 그래서 Scala 같은 함수형 언어에 대한 개념이 필요하다.

8:07 AM - 13 Jan 2015