DataFlow Model and Stream Processing

Saurabh Agarwal

Map Reduce - Word Count



Map Reduce - Word Count



Spark - Motivation

- Most Real Applications require multiple MR Steps
 - Indexing Pipeline 21 Steps
 - Analytics Query 5-11 Steps

Spark - Motivation

- Most Real Applications require multiple MR Steps
 - Indexing Pipeline 21 Steps
 - Analytics Query 5-11 Steps



Spark - Motivation

- Most Real Applications require multiple MR Steps
 - Indexing Pipeline 21 Steps
 - Analytics Query 5-11 Steps

Write to file System, expensive to reuse data (pagerank, interactive data)

execute ()

- Clean programmable API

- Fault tolerance and in-memory processing

Spark - Example

- Read from a Log and filter errors

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
messages = errors.map(_.split('\t')(2))
messages.cache()
messages.filter(_.contains("foo")).count
messages.filter(_.contains("bar")).count
```



Stream Processing

The Dataflow Model: A Practical Approach to Balancing Correctness, Latency, and Cost in Massive-Scale, Unbounded, Out-of-Order Data Processing

many users

Motivation

Consider an initial example: a streaming video provider wants to monetize their content by displaying video ads and billing advertisers for the amount of advertising watched. The platform supports online and offline views for content and ads. The video provider wants to know how much to bill each advertiser each day, as well as aggregate statistics about the videos and ads. In addition, they want to efficiently run offline experiments over large swaths of historical data.

Advertisers/content providers want to know how often and for how long their videos are being watched, with which content/ads, and by which demographic groups. They also want to know how much they are being charged/paid. They want all of this information as quickly as possible, so that they can adjust budgets and bids, change targeting, tweak campaigns, and plan future directions in as close to real time as possible. Since money is involved, correctness is paramount.

watched avide how many finished vide how long prohow many people in last 1 hr t Desirable adtributes of stream processing ß

Solution

Separate user API from Execution

Decompose Queries into

🧫 (1) Speed Scalability Query modification E Fault tolerant (4) Minimize gramework Ś over head Correctness What to compute -> _ Time Stamp (6) Where in **Event Time** they are being computed event occured

- When in **Processing Time** they are being materialized
- the data How earlier results are related a ctually brocks _

Streaming vs Batch











Windowing

- Required for some operations, unnecessary for others

- Windowing in Batch data ??



Data Flow Model API

- ParDo:

- GroupByKey:

- Windowing
 - AssignWindow
 - MergeWindow

Data Flow Model API



Fixed Interval

```
PCollection<KV<String, Integer>> output = input
.apply(Window.trigger(Repeat(AtPeriod(1, MINUTE)))
.accumulating())
.apply(Sum.integersPerKey());
```



Fixed Data Count

PCollection<KV<String, Integer>> output = input
.apply(Window.trigger(Repeat(AtCount(2)))
.discarding())
.apply(Sum.integersPerKey());



Micro-Batch Processing



