Sock Apache-Spark Ecosystem and Key Innovations Linxiao Bai Nilesh Patil

Spark Ecosystem

Introduction

Spark was designed for massive parallel computation on large-scale distributed cluster. It is currently the most popular open-source project of Apache Software Foundation, and the first choice of "big-data engine" for clustercomputation.

Although Spark Core learns from Hadoop MapReduce framework at lower level, with extensive optimizations and design choices, it quickly outperformed its predecessor, providing faster, more resilient, more convenient clustercomputation experience. Also, with extensive modules and libraries like Spark SQL, Spark Streaming, GraphX, and MLlib, it is a powerful platform to cater different computation needs, and different environments.

These powerful capabilities of Spark, understanding its key features and advantages is important for anyone who works in data-driven fields.

Extended Modules

Spark Core

Objective

Our goal is to provide an overview of Apache Spark ecosystem and several key Spark modules. Including lowerlevel interfaces to machine and higher-level interfaces to users.

Also, we will compare Spark with its predecessor Hadoop, and review its advantages. The intention is to provide Spark intermediate learners with a better understanding of this powerful platform and offer a different perspective.

Resource Management

Distributed File System

Modules and Innovations

User

Other Extended Modules

API for graphs and large scale graph-parallel computation. GraphX



A Spark scalable machine learning library.

Spark

API for scalable

The resilience of Spark relies on DAGScheduler, a process that record the lineage of all RDD partitions.

Compared to Hadoop:

Intermediate result of computations will not be serialized and saved to disk. RDD stays in memory as a reusable object.

Compared to Hadoop:

Spache SQL

The most popular Spark extension. Powerful relational data processing capability

Implement DataFrame like in

Python and R

Uses schema at each data





streaming applications

Notebook like interactive kernel



- One executor fails, all fail.
- Computation start-over.

hedoop

- Relies on saving result to HDFS.
- ➢ Huge disk I/O, waste of
 - communication.

Spark Core

- > One executor fails, others move on
- > Only the faulty partition needs re-
- compute
- Relies on lineage to trace-back for
- computing the faulty partition.
- Minimum I/O, Maximum efficiency
- Allocate help to struggling executors.



- Result of map-reduce save to disk, and make multiple copies.
- Serialization/deserialization repetitively between computations
- Huge disk I/O, waste of communication.
- Lazy execution, execute tasks batch-wisely

Spache

- Result saved as reusable object in memory for next computation.
 - \succ No serialization/deserialization.
 - Minimum disk I/O. Lightning fast !

partition to keep track of the relational transformation.

Uses "Catalyst" to optimize

the query logic.

Upper-level



Encapsulation

Resilience

Spark encapsulates parallelism to friendly APIs. User programs without worrying about too much details.

Compared to Hadoop:

Spark offers great generality for application development and cluster deploying.

In-Memory

Generality

Compared to Hadoop:



Resource Manager

Wide-raged support includes:



Spark-Standalone Spark-Local





Friendly interfaces, users

only worry about job itself.

- \succ Minimum declarations, Users are responsible for details usually only **SparkContext**. like exception handling.
- Massive declaration before start.
- Redundant coding, similar
 - modules appears every time.

of code

- "WordCount" takes 120 lines
- Parallelisms are wellencapsulated. "WordCount" only 5 lines of code.
- hedoop
- Needs to be deployed on Hadoop cluster. Heartbeat-based communication, Slow.
- Non-universal API.
- Java indigenous.



- Support platforms of different kinds and sizes. From local PC to large-scale cluster.
- Multiple language kernels
- Universal API. Akka fast
 - communication.
- Powerful extensions
- and libraries.



Lower-level

File System

Wide-raged support includes:







Machine