**Universität Bielefeld**

# Hands on analysis of NGS data with Sparkhit

de.NBI summer school 2017

(30.06.2017)

**Liren Huang**

**Jan Krüger**

**Supervisor: Alexander Sczyrba**

Bielefeld University

**The future is already here -**
      **it is just not very evenly distributed**
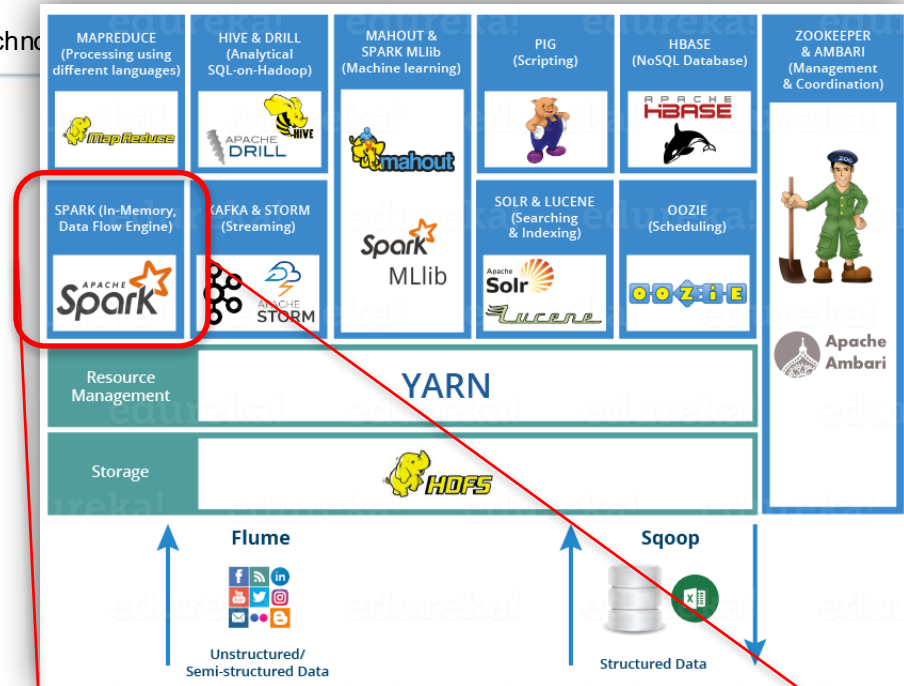
**William Ford Gibson**

# Overview

- **Introduction**
    - Apache Spark and RDD
    - Sparkhit, a toolkit for NGS data analysis

- **Hands on section**
    - Spark shell programming with RDD`s interface
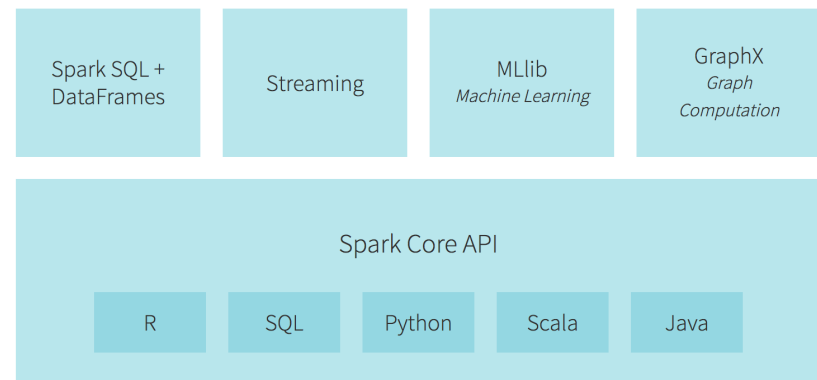    - Analyzing NGS data with Sparkhit

# Apache Spark

is a fast and general engine for large-scale data processing

- An Extended Map-Reduce model

- A distributed programming engine that can interact with most tools in Hadoop eco-system

- Its core is a distributed data abstraction called RDD (resilient distributed dataset)
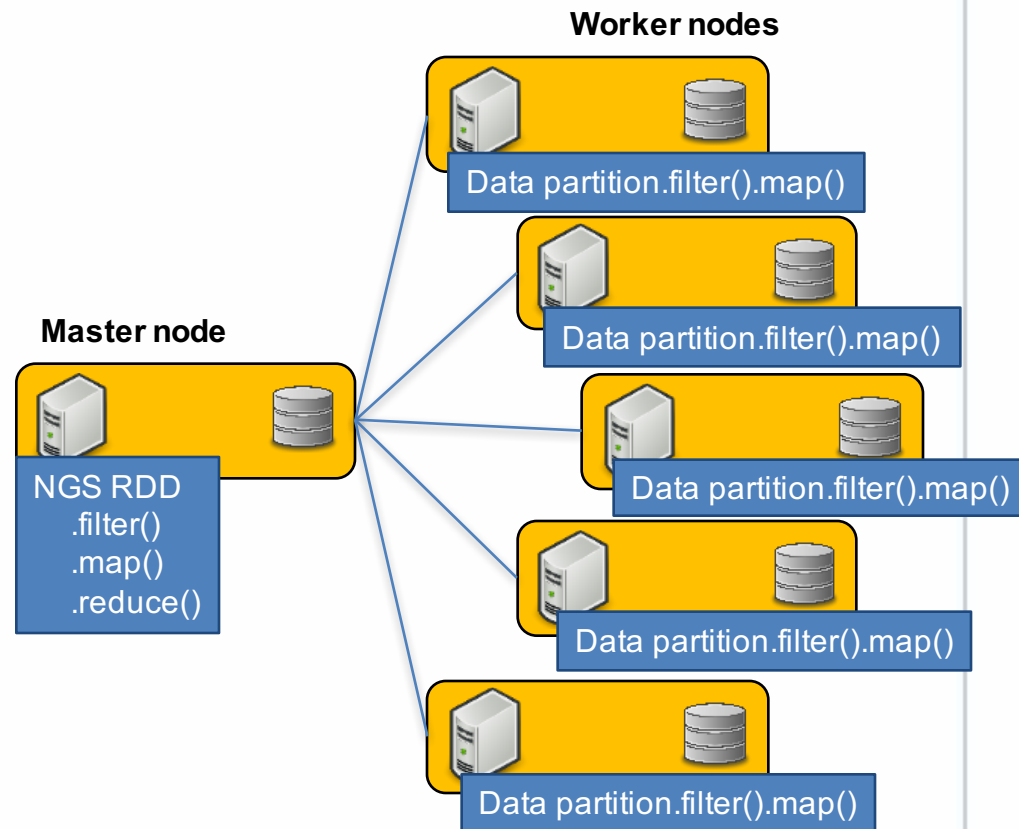


**Edureka Inc.**



**Databrick Inc.**

4

Universität Bielefeld

# RDD parallelization

- A RDD is an object.

- A RDD consists of several Data partitions across the cluster.

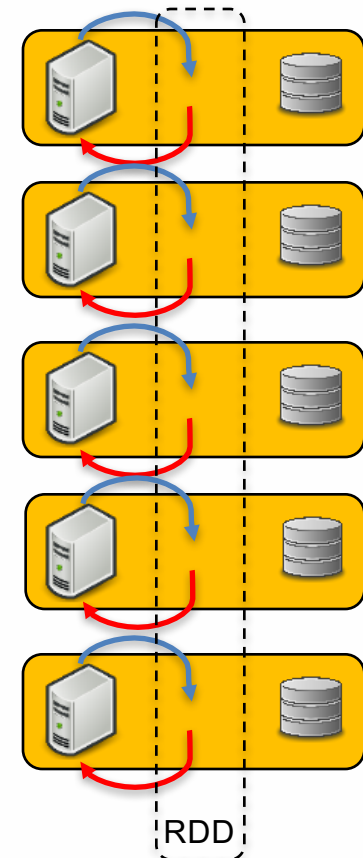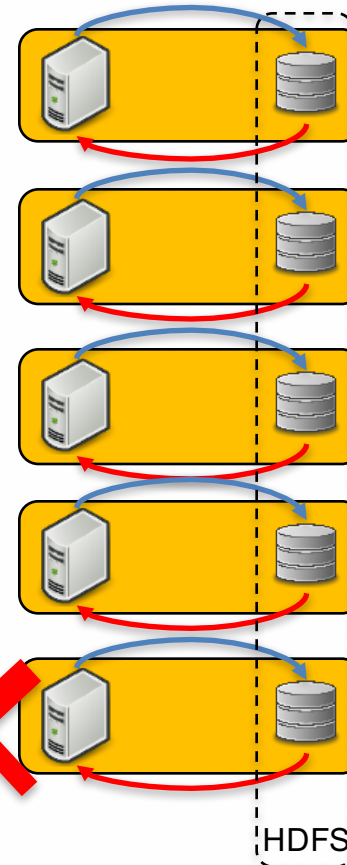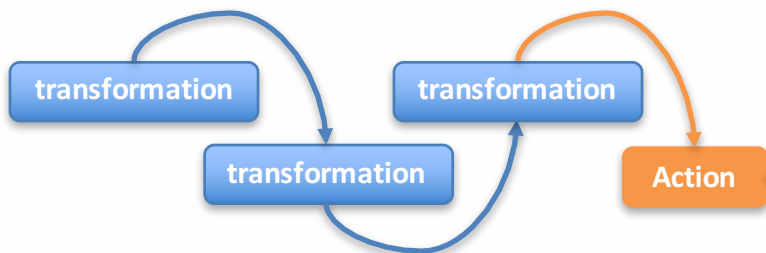- An operation to RDD is parallelized to each partition

The benefit is you focus on your algorithm while Spark distributed the workload for you.

**Worker nodes**

Data partition.filter().map()

Data partition.filter().map()

**Master node**

Data partition.filter().map()

NGS RDD
.filter()
.map()
.reduce()

Data partition.filter().map()

Data partition.filter().map()

# RDD cache

(Resilient distributed dataset)

- Distributed in memory computation for faster iterative algorithms

- Two types of operations
  - Transformation
  - Action

- Lazy feature (will see later), related to fault tolerance mechanism.
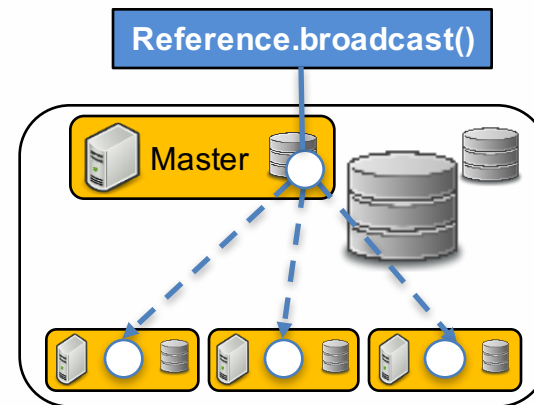
# Sparkhit-mapper

Sparkhit is a bioinformatics toolkit build on the Apache Spark platform

Here we describe a fragment recruitment application (short read mapping) call Sparkhit-mapper
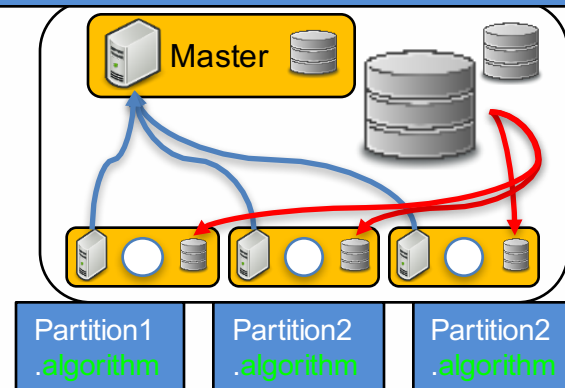
1, build reference index

2, broadcast to each worker nodes

3, each worker applies implemented alignment-algorithm for recruiting the fragments.



**https://rhinempi.github.io/sparkhit/**

www.uni-bielefeld.de

# Acknowledgement

- Dr. Alexander Sczyrba, de.NBI Bielefeld

- Jan Krüger, de.NBI Bielefeld

- Dr. Burkhard Linke, de.NBI Giessen

**Online tutorial:**

**https://rhinempi.github.io/sparkhit/usecase.html**

→ www.uni-bielefeld.de