Search                                  search                          blog   about   tags
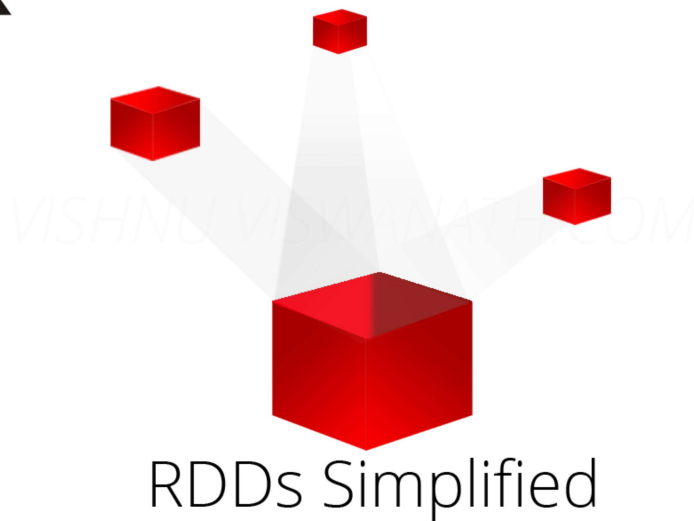
# Spark RDDs Simplified

February 4, 2016. Estimated read time: 3 minutes



Spark RDDs are very simple at the same time very important concept in Apache Spark.
Most of you might be knowing the full form of RDD, it is **Resilient Distributed Datasets**.
*Resilient* because RDDs are immutable*(can't be modified once created)* and fault
tolerant, *Distributed* because it is distributed across cluster and *Dataset* because it holds
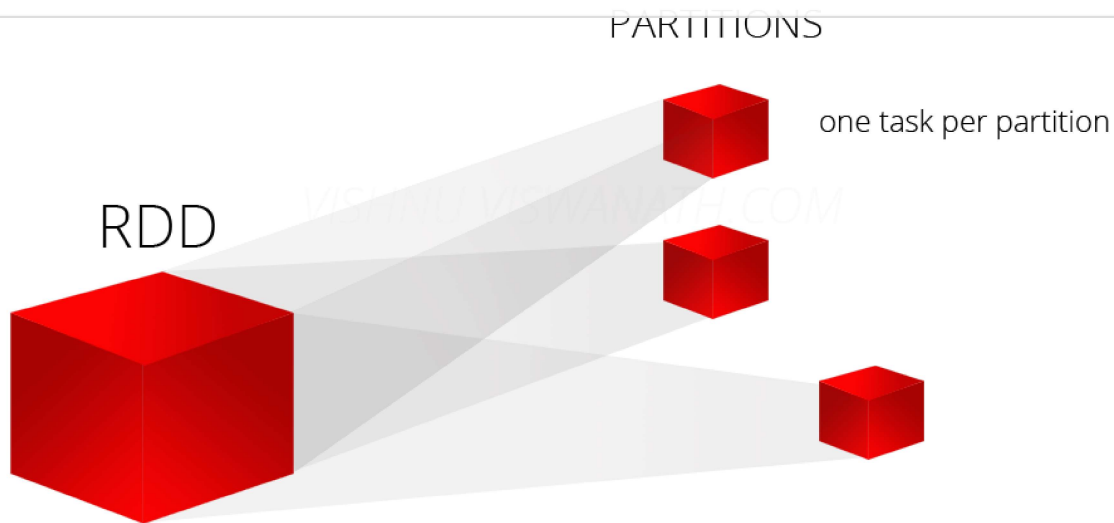data.

So why RDD? Apache Spark lets you treat your input files almost like any other variable,
which you cannot do in Hadoop MapReduce. RDDs are automatically distributed across
the network by means of Partitions.

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞   ✖

Email address                                          Subscribe         *Copyright © 2018 Vishnu Viswanath*

PARTITIONS

one task per partition

RDD

RDDs are divided into smaller chunks called Partitions, and when you execute some action, a task is launched per partition. So it means, the more the number of partitions, the more the parallelism. Spark automatically decides the number of partitions that an RDD has to be divided into but you can also specify the number of partitions when creating an RDD. These partitions of an RDD is distributed across all the nodes in the network.

# Creating an RDD

Creating an RDD is easy, it can be created either from an external file or by parallelizing collections in your driver. For example,

```
val rdd = sc.textFile("/some_file",3)
val lines = sc.parallelize(List("this is","an example"))
```

◄                                                              ►

The first line creates an RDD from an external file, and the second line creates an RDD from a list of Strings. *Note that the argument '3' in the method call sc.textFile() specifies*

66 Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? 99   ✖

*partitions, then you can simply call sc.textFile("some_file").*

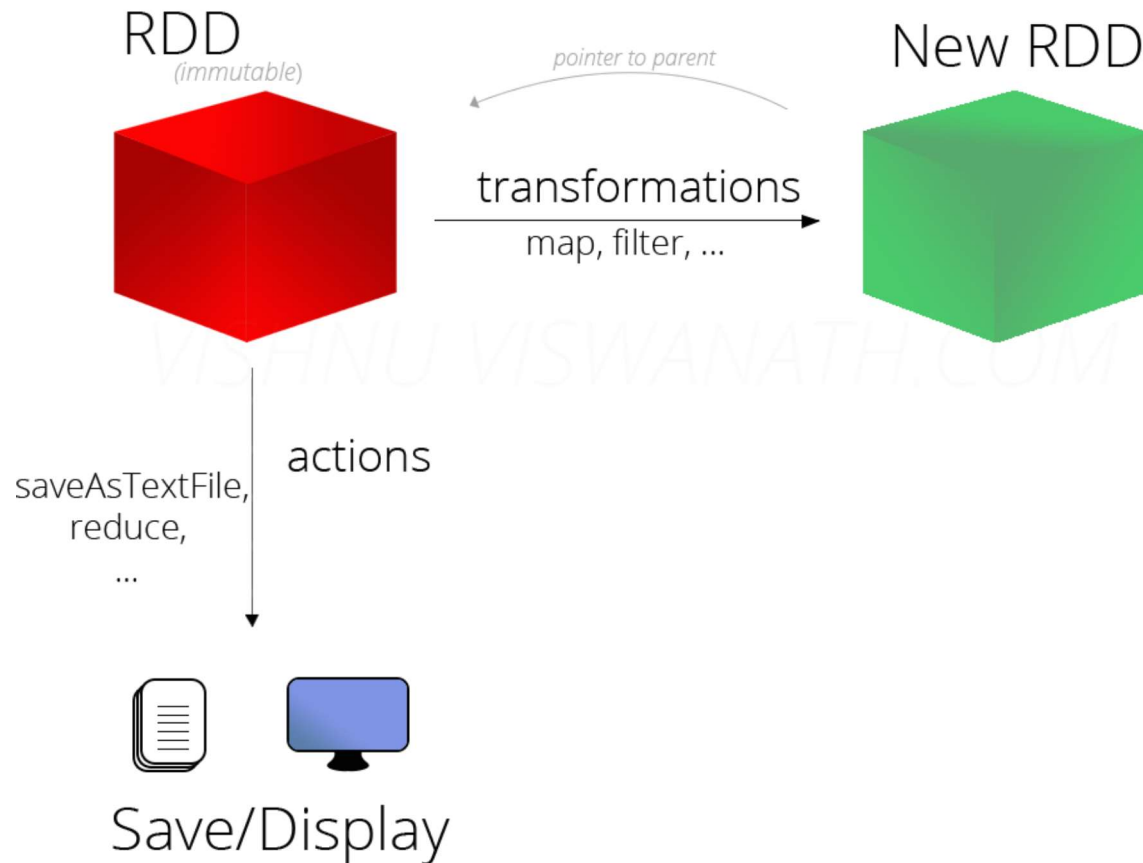Email address                                      Subscribe          *Copyright © 2018 Vishnu Viswanath*

Search                                    search                          blog    about    tags

*Actions*. **Transformation** applies some function on a RDD and creates a new RDD, it does not modify the RDD that you apply the function on.*(Remember that RDDs are resilient/immutable).* Also, the new RDD keeps a pointer to it's parent RDD.



When you call a transformation, Spark does not execute it immediately, instead it creates a **lineage**. A lineage keeps track of what all transformations has to be applied on that RDD, including from where it has to read the data. For example, consider the below example

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞     ✖

| Email address | | Subscribe |
|---|---|---|

*Copyright © 2018 Vishnu Viswanath*

rdd = sc.textFile("spam.txt")    filtered = rdd.filter()         filtered.count()

```scala
val rdd = sc.textFile("spam.txt")
val filtered = rdd.filter(line => line.contains("money"))
filtered.count()
```

*sc.textFile() and rdd.filter()* do not get executed immediately, it will only get executed once you call an *Action* on the RDD - here filtered.count(). An **Action** is used to either save result to some location or to display it. You can also print the RDD lineage information by using the command `filtered.toDebugString` *(filtered is the RDD here).*

> *RDDs can also be thought of as a set of instructions that has to be executed, first instruction being the load instruction.*
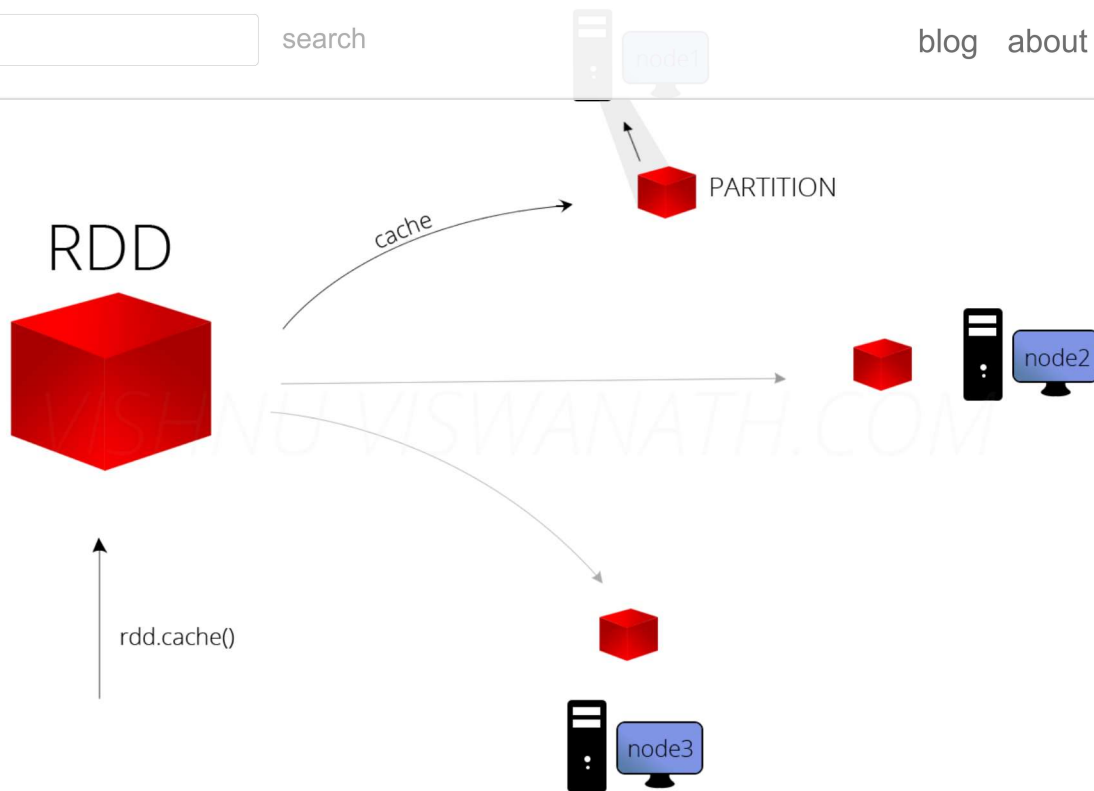
# Caching

You can cache an RDD in memory by calling `rdd.cache()`. When you cache an RDD, it's Partitions are loaded into memory of the nodes that hold it.

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞   ✖

Email address                          Subscribe              *Copyright © 2018 Vishnu Viswanath*

Caching can improve the performance of your application to a great extent. In the previous section you saw that when an action is performed on a RDD, it executes it's entire lineage. Now imagine you are going to perform an action multiple times on the same RDD which has a long lineage, this will cause an increase in execution time. Caching stores the computed result of the RDD in the memory thereby eliminating the need to recompute it every time. You can think of caching as if it is breaking the lineage, but it does remember the lineage so that it can be recomputed in case of a node failure.

This concludes the basics of RDD. If you would like to read more, Part2 talks about Persistence, Broadcast variables and Accumulators. Thanks for reading!

Continue reading

Tags: ApacheSpark Scala BigData Hadoop

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞   ✖

| Email address | Subscribe |

*Copyright © 2018 Vishnu Viswanath*

**Realtime Processing using Storm-Kafka- …**

6 years ago • 1 comment

This is part two of the series Realtime Processing using Storm and Kafka. In this …

**Deep Learning - ANN, RNN, LSTM networks**

4 years ago • 2 comments

Long Short Term Memory(LSTM) model is a type supervised Deep …

**Query Apacl**

5 years

This is Querya Flink. I

---

**25 Comments**       **vishnuviswanath.com**       🔓              1  **Login**

♡ **Favorite**  5              🐦 **Tweet**       f  **Share**              Sort by Best

| Join the discussion… |

LOG IN WITH

OR SIGN UP WITH DISQUS ?

| Name |

**Ye Joo Park** • 4 years ago
Very nice post and amazing illustrations, thanks!

1 ∧ | ∨ • Reply • Share ›

> **Vishnu Viswanath**  Mod ➜ Ye Joo Park • 4 years ago
> Thank you!
>
> ∧ | ∨ • Reply • Share ›

**Raj** • 4 years ago
Nice post, got a solid reinforcement to my learnings..

I recently worked on sparklyr where I had to do

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞   ✖

increases, the work increases, so, we use cache to make it

---

| Email address |       Subscribe       *Copyright © 2018 Vishnu Viswanath*

Search        search with cache? Can I have multiple        blog   about   tags
caches in the same lineage and access a particular cache to do
another analysis/transformations?

1 ∧ | ∨ • Reply • Share ›

**Vishnu Viswanath**  Mod  → Raj • 4 years ago
Thank you and good question. Caching and
checkpointing have some similarities but are not the
same. Caching tells spark to store the computed value in
memory/external storage. Caching does not break the
lineage, i.e., the RDD still have the lineage info.
RDD.checkpoint() stores the computed value to the
storage and it breaks the lineage. While applying
performing an action, spark first checks if the RDD is
cached, if not it checks if the RDD is checkpointed, if not
it computes the values from the source by applying all
the transformations before that point.
Yes, you can have multiple caches in the same linage
and use it for doing another transformation.

∧ | ∨ • Reply • Share ›

**Lyubcho Dimov** • 4 years ago
Amazingly well explained!! Kudos! : )

1 ∧ | ∨ • Reply • Share ›

**Vishnu Viswanath**  Mod  → Lyubcho Dimov • 4 years ago
Thank you :)

∧ | ∨ • Reply • Share ›

**Rajesh** • 5 years ago
Remember that RDD.cache() is also lazy. Let's see following
code -

L1 val rdd = sc.textFile("spam.txt")
L2 val filteredRDD = rdd.filter(line => line.contains("money"))
L3 filteredRDD.cache() // This is lazy as well
L4 filteredRDD.count()
L5 some code
L6 some code
L7 filteredRDD.count()

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞  ✖

One might think line L3 has loaded the file and cached filtered

Email address                              Subscribe

                                                    *Copyright © 2018 Vishnu Viswanath*

Line L4 is an action (1st time the filteredRDD is executed), when this is executed - filteredRDD is loaded, cached, and counted.

When line L7 is executed (2nd time the filteredRDD is executed) - the operation will take the data from the cache and count the lines.

1 ∧ | ∨  •  Reply  •  Share ›

**Vishnu Viswanath**  Mod  → Rajesh • 5 years ago

good point Rajesh.

∧ | ∨  •  Reply  •  Share ›

**Prasiddh Nandola** • 6 years ago

I have one doubt: we have only rdd.persist() and rdd.cache() methods. But, i do not want to persist entire rdd, instead i want to persist only some partitions of a particular rdd. how can i do that?

1 ∧ | ∨  •  Reply  •  Share ›

**Tapan Behera** → Prasiddh Nandola • 4 years ago

create one more RDD based on filteration and cached it. So it will cahced the new RDD not the old one.

1 ∧ | ∨  •  Reply  •  Share ›

**Vishnu Viswanath**  Mod  → Prasiddh Nandola
• 6 years ago • edited

I am not aware of such an option. what is the use case you are having ?. may be you can split your RDD into multiple RDDs based on some condition and choose to cache only the RDD you want.

∧ | ∨  •  Reply  •  Share ›

**Prasiddh Nandola** → Vishnu Viswanath
• 6 years ago

I want only some part of the rdd to persist because, Each partition has either high/low incremental

> Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ✖

updates (resulting in updates on partition 1) on

| Email address | | Subscribe |

*Copyright © 2018 Vishnu Viswanath*

materialize partition2 in memory not partition 1.

| Search | search | blog  about  tags |

**Prasiddh Nandola** • 6 years ago

nice one

1 ∧ | ∨ • Reply • Share ›

> **Vishnu Viswanath**  Mod → Prasiddh Nandola
> • 6 years ago
>
> Thank you
>
> ∧ | ∨ • Reply • Share ›

**Matt Gardner** • 6 years ago

Great blog post to explain the basic concepts - I shared this with my intro Data Science class I am teaching. Thank you.

1 ∧ | ∨ • Reply • Share ›

> **Vishnu Viswanath**  Mod → Matt Gardner • 6 years ago
>
> Glad you liked it, and hopefully it will help the students that you are teaching. Thanks!
>
> ∧ | ∨ • Reply • Share ›

**likitha prakash** • 3 years ago

Clear and simple explanation!!! Saved lot of my time :) Appreciate it!!!

∧ | ∨ • Reply • Share ›

**amit shah** • 5 years ago

Nice post. I like the simplicity with which the concept is explained. I have a question about RDD's - At what instance in time would the rdd be in-memory? From the above explanation it seems that the RDD would be discarded from the RAM after the lineage execution is completed. To keep it in-memory, we need to call rdd.persist() or rdd.cache(). Is my understanding right?

∧ | ∨ • Reply • Share ›

> **Vishnu Viswanath**  Mod → amit shah • 5 years ago
>
> Hi Amit, Thank you. Glad you liked the post.
> You are right, an RDD is kept in memory if you call

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞  ✖

recomputed every time. This is usually done when

| Email address | | Subscribe |

the same RDD. But note that your cluster should have
~~enough~~ ~~hold~~ the RDD in memory, otherwise
~~you will have to do~~ some other form of persistence (may
be memory + disk ). If you don't do any caching, during
the processing of the RDD, each partition of the RDD
will be brought into memory and will be processed.

 ∧ | ∨ • Reply • Share ›

**amit shah** ➔ Vishnu Viswanath • 5 years ago

Got it. Thanks for the detailed reply. While
reading more about RDD's I have some new
questions which could be worth discussing

1. What is the underlying physical representation
of a Spark RDD?
2. What is a pairRDD? While reading about
RDD's I understand that spark provides different
types of RDD's meant for specific needs. I am
looking for more details.
3. How can we execute SQL queries on a RDD?
I understand that the dataframe or the dataset
api is generally used for this but I want to know if
it's possible to execute SQL queries on RDD's. I
read a bit about SchemaRDD's but wasn't able to
grasp it completely.

I would appreciate your thoughts on this.

 ∧ | ∨ • Reply • Share ›

**Vishnu Viswanath** **Mod** ➔ amit shah
• 5 years ago

1. As far as I know, Spark by default uses
Java Serialization, you can also use Kryo
serializer. So the physical representation
would depend on the type of Serializer
used. (https://spark.apache.org/do...

2. A pairRDD is an RDD with key and
value. This is useful in cases where you
want to perform operations based on the
key. E.g., aggregateByKey, countByKey,
reduceByKey etc.
~~(http://spark.apache.org/doc~~

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞ ✖

RDD. You need to convert that RDD into

Email address                          Subscribe

*Copyright © 2018 Vishnu Viswanath*

Search

search

blog  about  tags

❝ Hello, Would you like to subscribe so that I can keep you posted on my new aritcles? ❞  ✖

Email address

search

Subscribe

*Copyright © 2018 Vishnu Viswanath*