

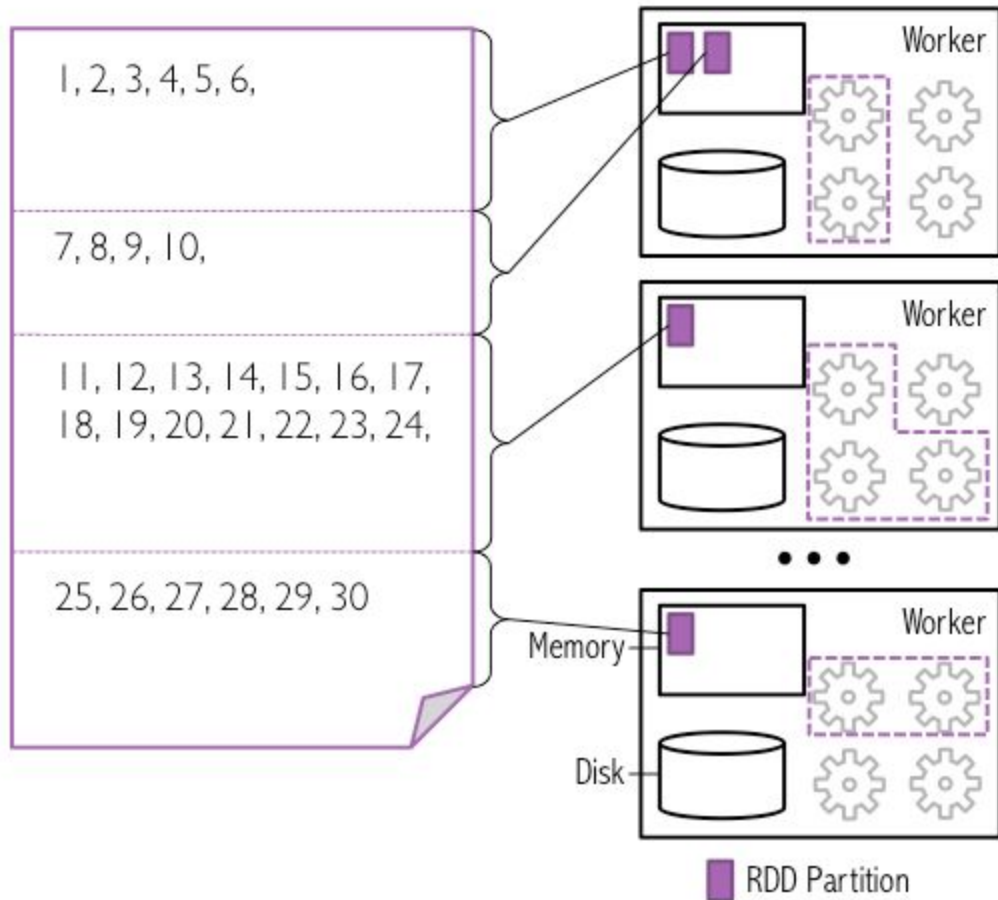
RDD Operations

...

Transformation and Action Zoo

Dataset is broken into
partitions

Partitions are each stored
in a worker's memory



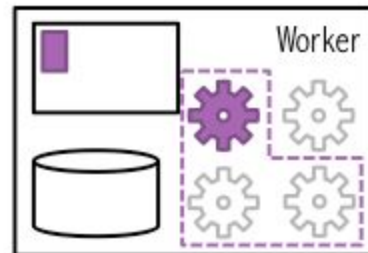
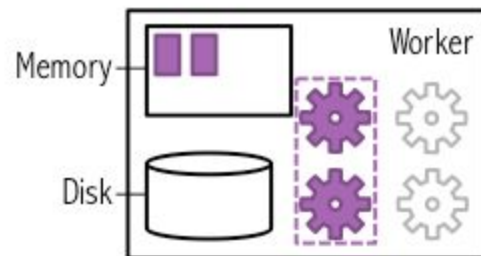
One task is launched
for each partition



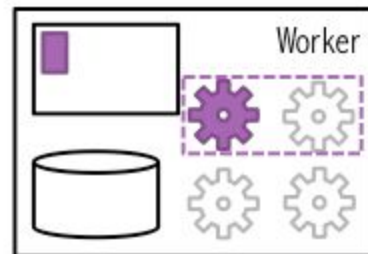
Running task



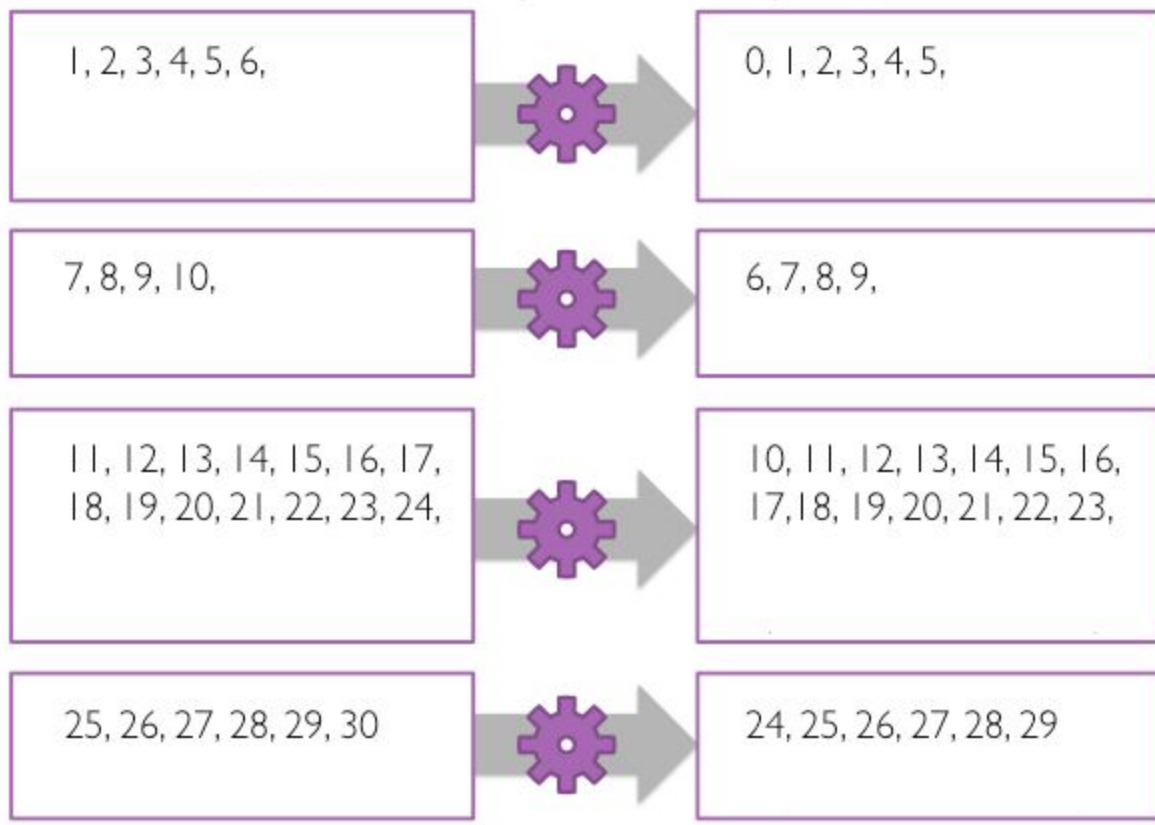
Unused core



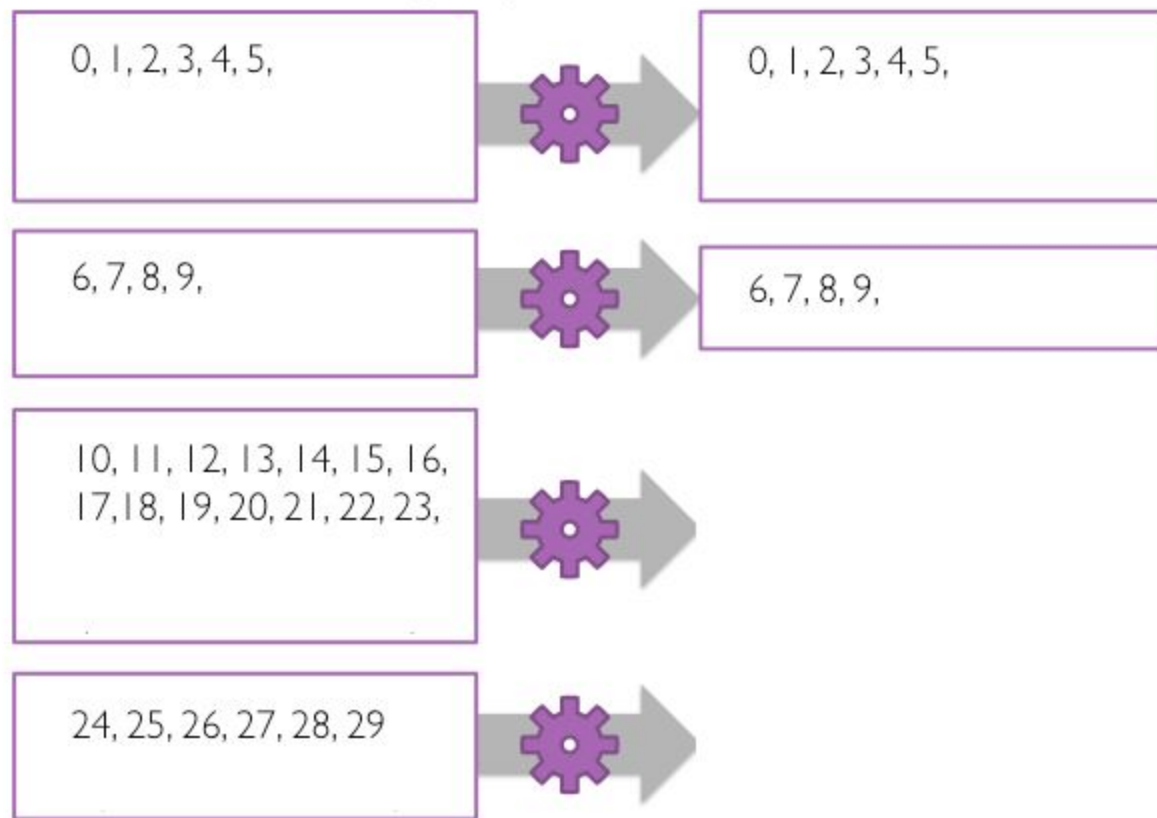
...



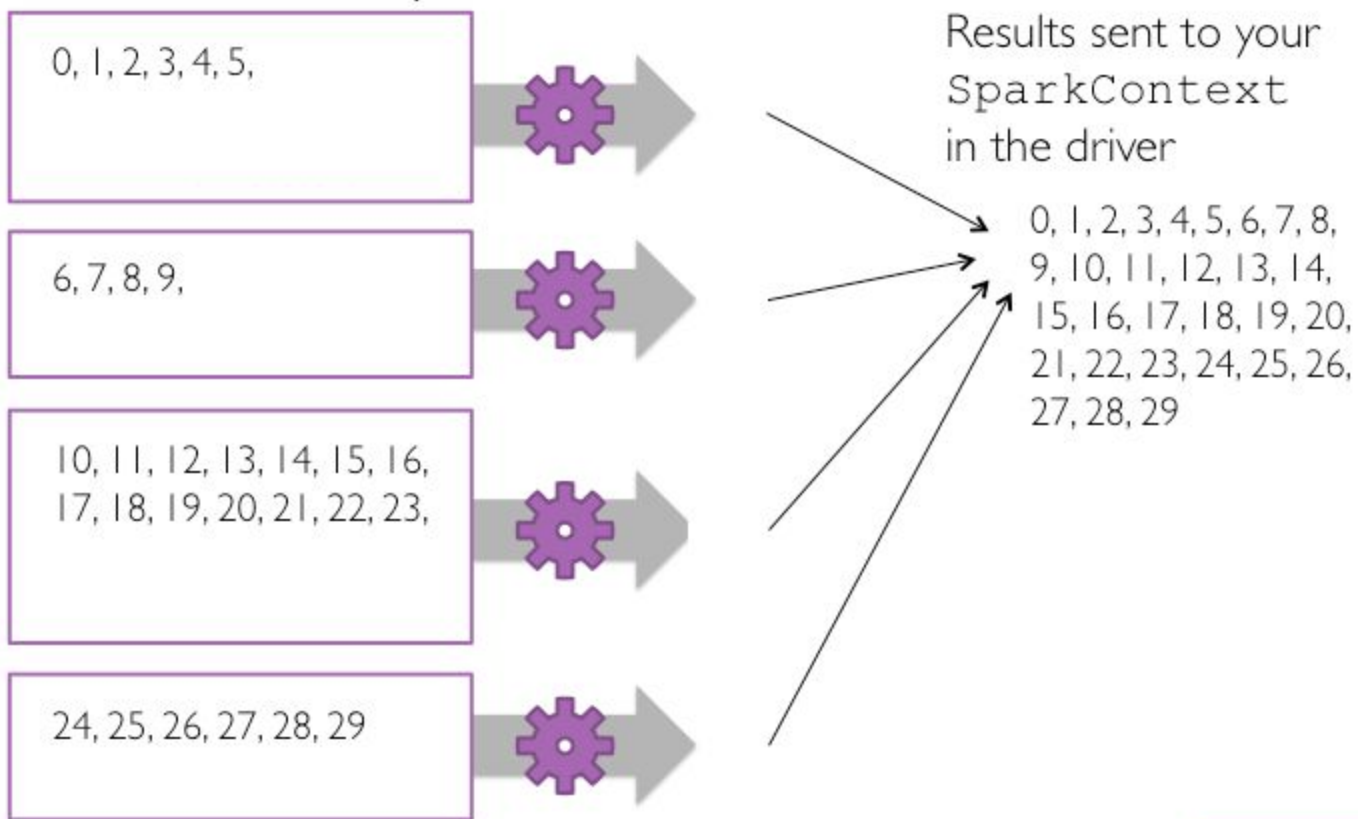
`map (f) :` Each task makes a new partition by calling `f (e)` on each entry `e` in the original partition



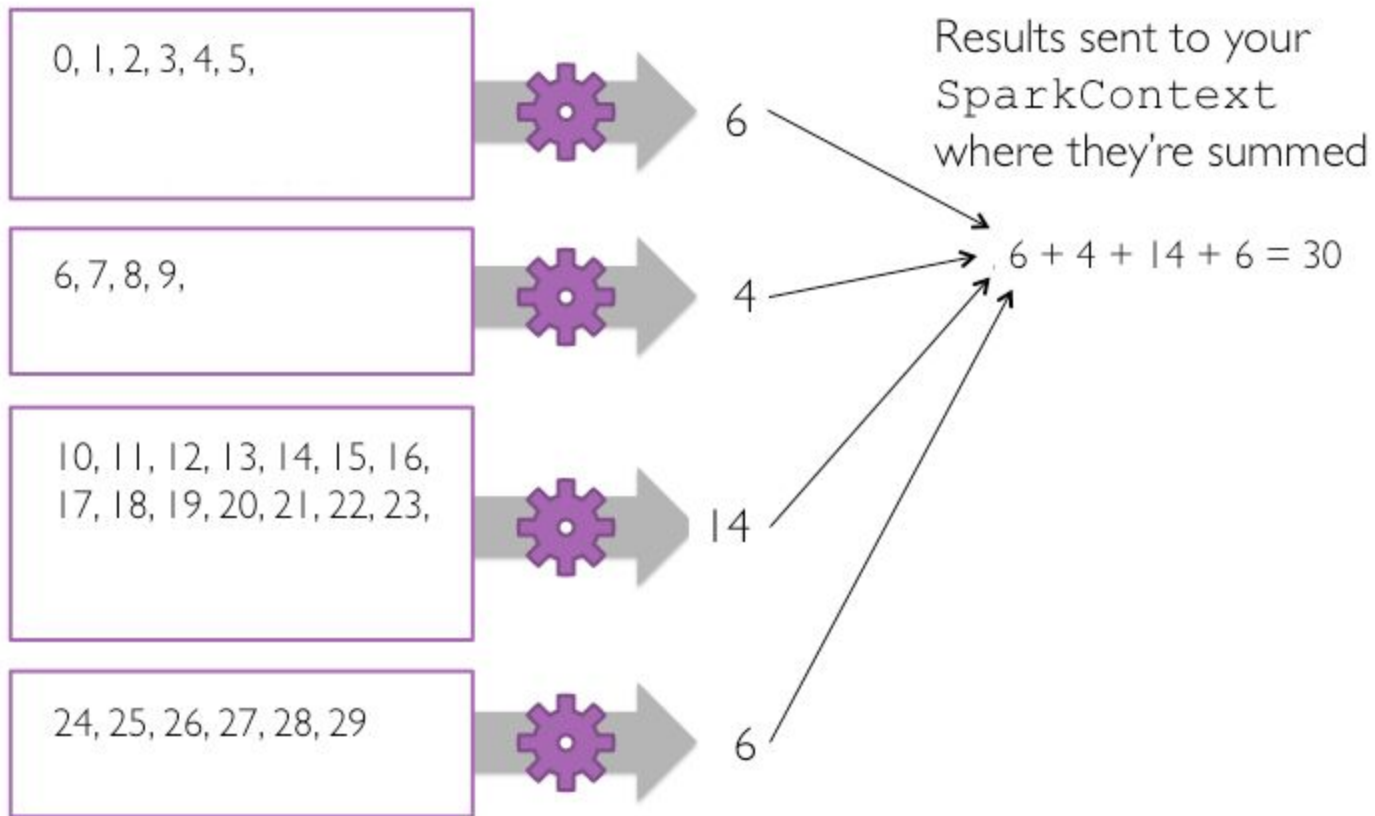
`filter()` : Each task makes a new partition with the entries in the original partition that have a value less than 10



`collect()` : Gathers the entries from all partitions into the driver



count () : Each task counts the entries in one partition



Operations on RDD

Transformations and Actions

Once created, RDDs offer two types of operations:

- Transformations (akin to map)
- Actions (akin to reduce)

Laziness

Transformation

Transformations construct a new RDD from a previous one

Examples

- Applying functions on each element of an RDD
- Filtering data that matches a predicate

Action

Actions compute a result based on an RDD

- Result is either
 - returned to the driver program, or
 - saved to an external storage system

Examples

- return the first element of an RDD
 - return the RDD as a Python object
-

Laziness

Why is it useful?

```
lines = sc.textFile("large_textfile")  
lines.first()
```

Don't need to read the whole file to
extract the first line.

Exploring RDDs

Hands-on Activity

Spark_Activities_02_Transformations
_and_Actions.ipynb

Activity 2

Transformations

List of common Transformations

<http://spark.apache.org/docs/latest/programming-guide.html#transformations>

Spark_Activities_02_Transformations_and_Actions.ipynb

Activity 3

map()

Elementwise mapping

Can be performed on all RDDs

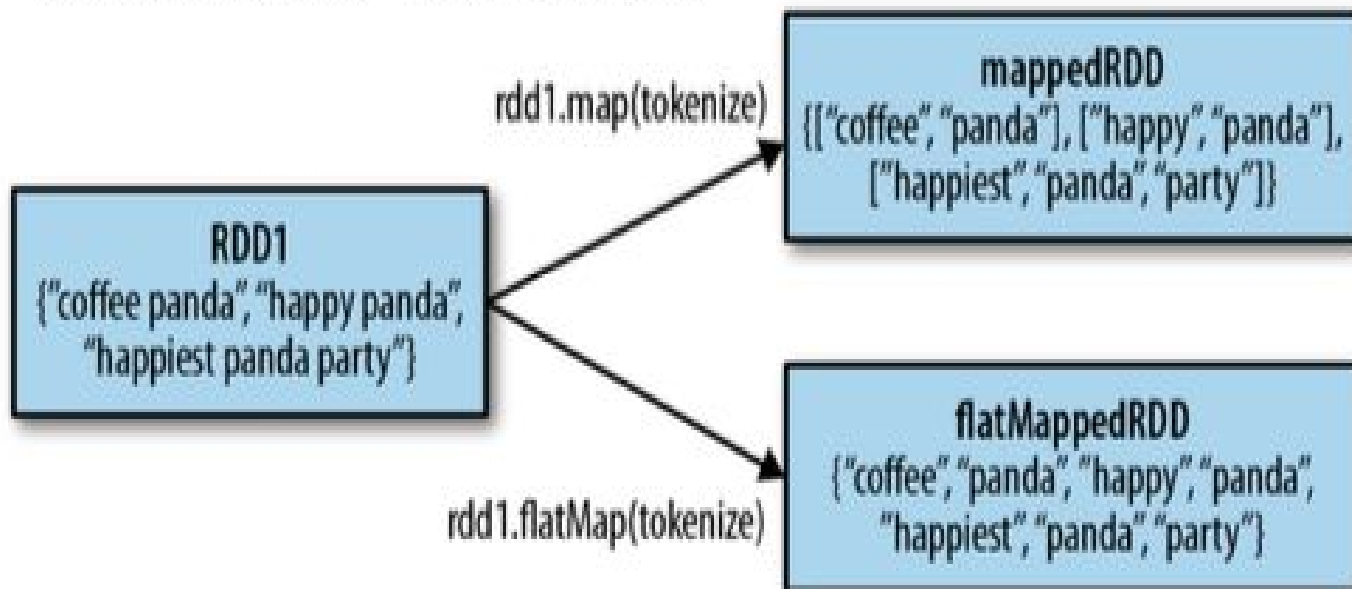
- Applies given function to each element of an RDD (maps existing value to new values) and returns the new RDD
- Return type doesn't have to be the same as input type

flatMap()

multiple output elements for each
input element

- Instead of returning a single element, it returns an iterator with our return values
- Rather than producing an RDD of iterators, it returns an RDD that consists of the elements from all of the iterators

`tokenize("coffee panda") = List("coffee", "panda")`



Set Operations

- Both RDD have to be of the same type
- The RDDs don't have to be sets themselves

- `distinct()`
- `union()`
- `intersection()`
- `subtract()`

Let's draw quick Venn Diagrams

RDD1
{coffee, coffee, panda,
monkey, tea}

RDD2
{coffee, money, kitty}

RDD1.distinct()
{coffee, panda,
monkey, tea}

RDD1.union(RDD2)
{coffee, coffee, coffee,
panda, monkey,
monkey, tea, kitty}

RDD1.intersection(RDD2)
{coffee, monkey}

RDD1.subtract(RDD2)
{panda, tea}

filter()

Filter out unwanted elements

Can be performed on all RDDs

- Returns an RDD that only has elements that pass the filter() function (the argument)
- Are input RDD type and output RDD type will always be the same?

Function name	Purpose	Example	Result
<code>map()</code>	Apply a function to each element in the RDD and return an RDD of the result.	<code>rdd.map(x => x + 1)</code>	{2, 3, 4, 4}
<code>flatMap()</code>	Apply a function to each element in the RDD and return an RDD of the contents of the iterators returned. Often used to extract words.	<code>rdd.flatMap(x => x.to(3))</code>	{1, 2, 3, 2, 3, 3, 3}
<code>filter()</code>	Return an RDD consisting of only elements that pass the condition passed to <code>filter()</code> .	<code>rdd.filter(x => x != 1)</code>	{2, 3, 3}
<code>distinct()</code>	Remove duplicates.	<code>rdd.distinct()</code>	{1, 2, 3}
<code>sample(withReplacement, fraction, [seed])</code>	Sample an RDD, with or without replacement.	<code>rdd.sample(false, 0.5)</code>	Nondeterministic

Function name	Purpose	Example	Result
<code>union()</code>	Produce an RDD containing elements from both RDDs.	<code>rdd.union(other)</code>	{1, 2, 3, 3, 4, 5}
<code>intersection()</code>	RDD containing only elements found in both RDDs.	<code>rdd.intersection(other)</code>	{3}
<code>subtract()</code>	Remove the contents of one RDD (e.g., remove training data).	<code>rdd.subtract(other)</code>	{1, 2}
<code>cartesian()</code>	Cartesian product with the other RDD.	<code>rdd.cartesian(other)</code>	{(1, 3), (1, 4), ... (3,5)}

Actions

List of Common Actions

<http://spark.apache.org/docs/latest/programming-guide.html#actions>

reduce()

takes a function that operates on two elements of the type in input RDD and returns a new element of the same type

Example

- Add, Multiply etc
- Paste (strings)
- Cumulative Sum

collect()

- The dataset must fit into a single machine

first()

- Returns the first element

take()

- Order may not be preserved

More Operation

Numeric RDD Operations

- `stats()`
- `max()`
- `mean()`
- `min()`
- `stdev()`

Spark_Activities_02_Transformations
_and_Actions.ipynb

Activity 5

Exercises

Exercises

Spark_Activities_02_Transformations
_and_Actions.ipynb

Activity 6
