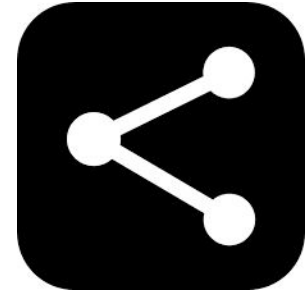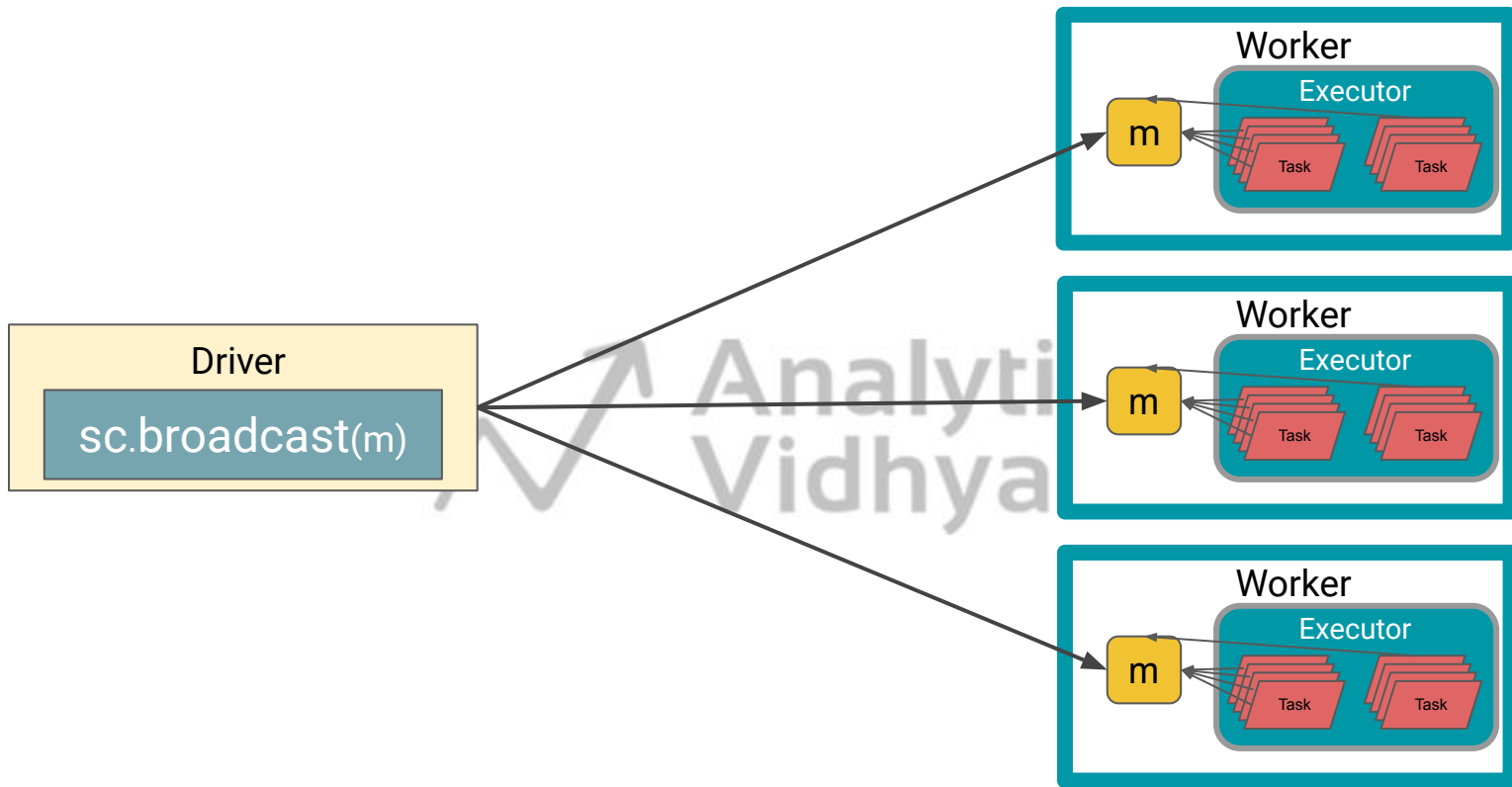# Advanced Programming in Spark

Shared Variables

# Shared Variables

- Broadcast Variables
- Accumulators

# Broadcast Variable

- Read-only variable cached on each node rather than shipping a copy of it with tasks

- They give every node a copy of a large input dataset in an efficient manner

- Attempts to distribute broadcast variables using efficient broadcast algorithms to reduce communication cost

# How to create Broadcast Variables?

Variable to be broadcasted

```
# create a list
my_list = ["India", "Australia", "Italy", "Sri-Lanka", "Singapore"]

# create the braodcast variable
broadcast_variable = sc.broadcast(my_list)
```

**broadcast** keyword with Spark Context

# When to create Broadcast Variables?

- When tasks across multiple stages need the same data.

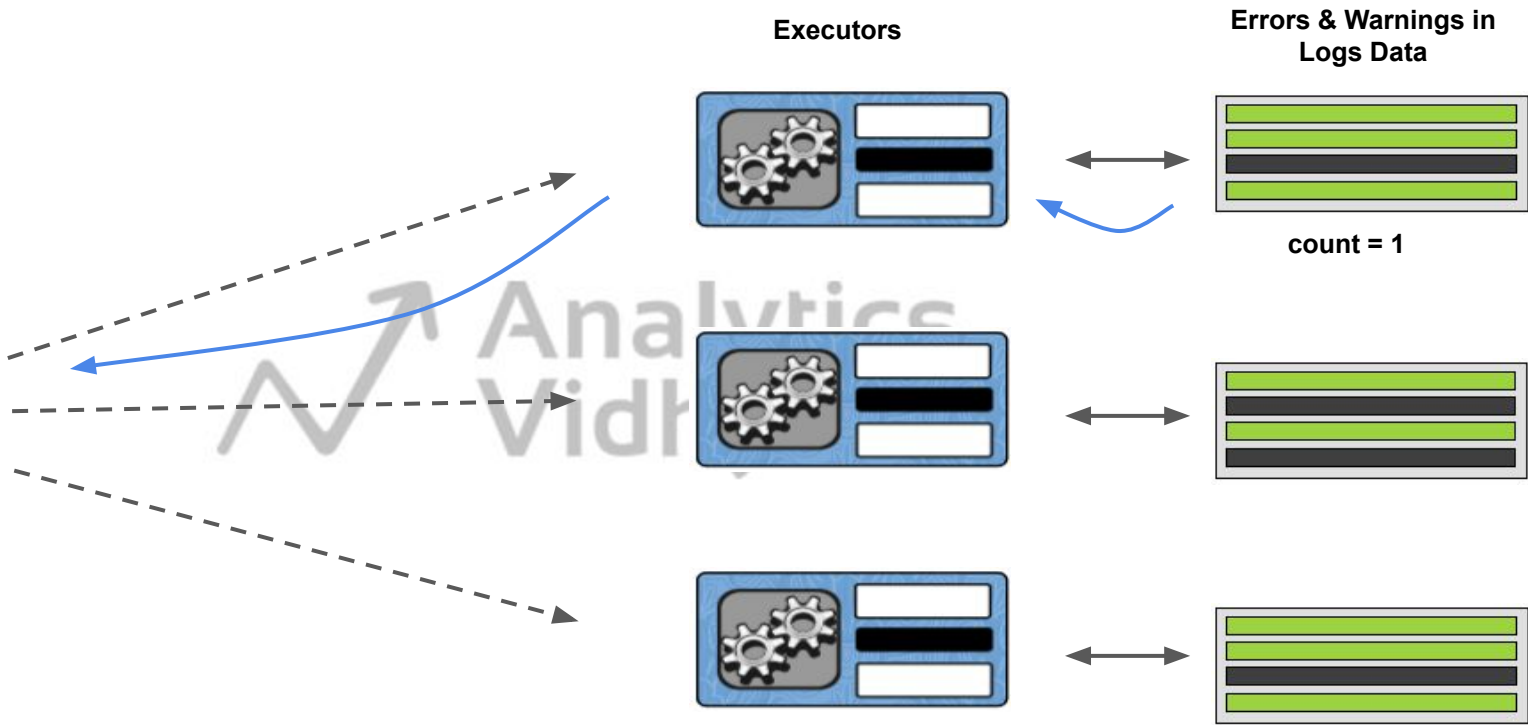- When caching the data in deserialized form is important.

# Accumulators

- Are only created through an associative or a commutative operation.

- Can be used to implement counters or sums.

- Spark natively supports accumulators of numeric types, and programmers can
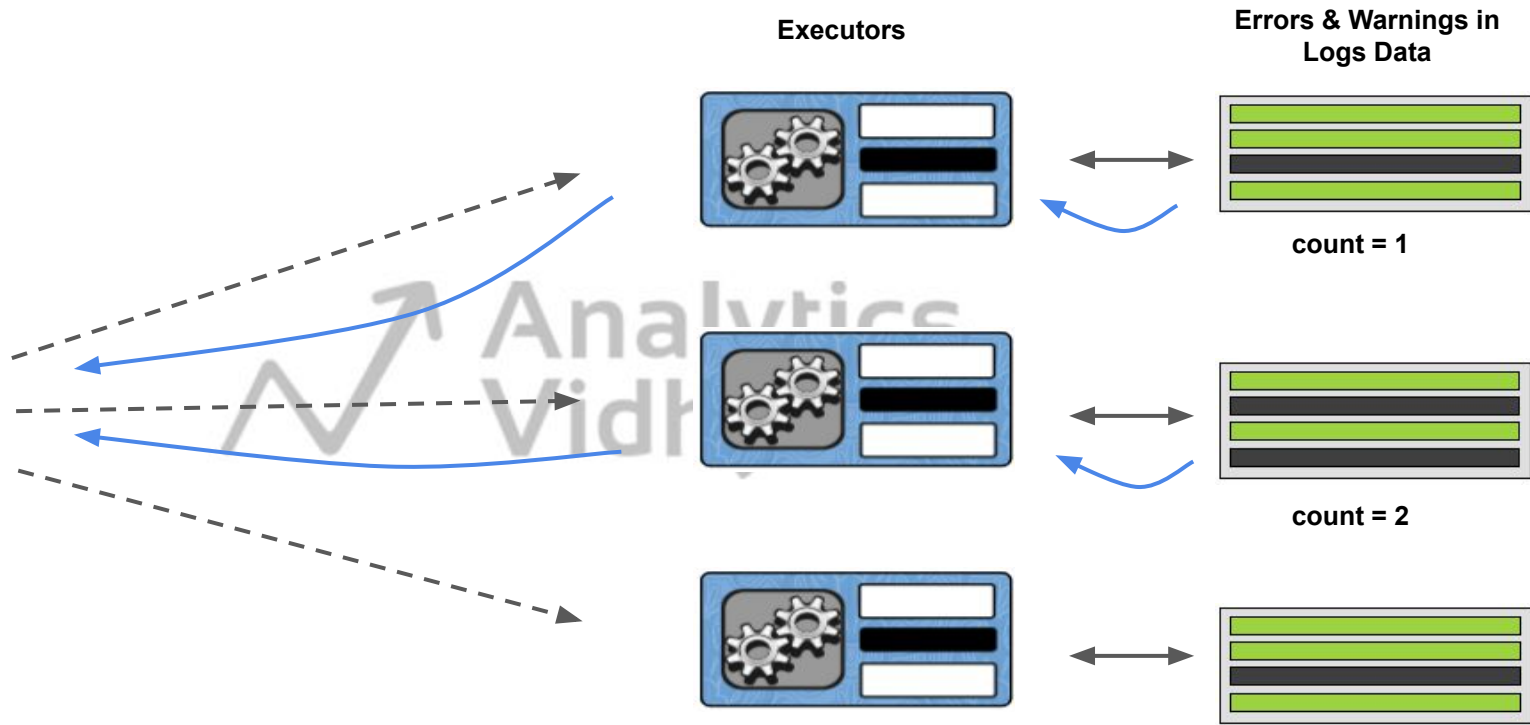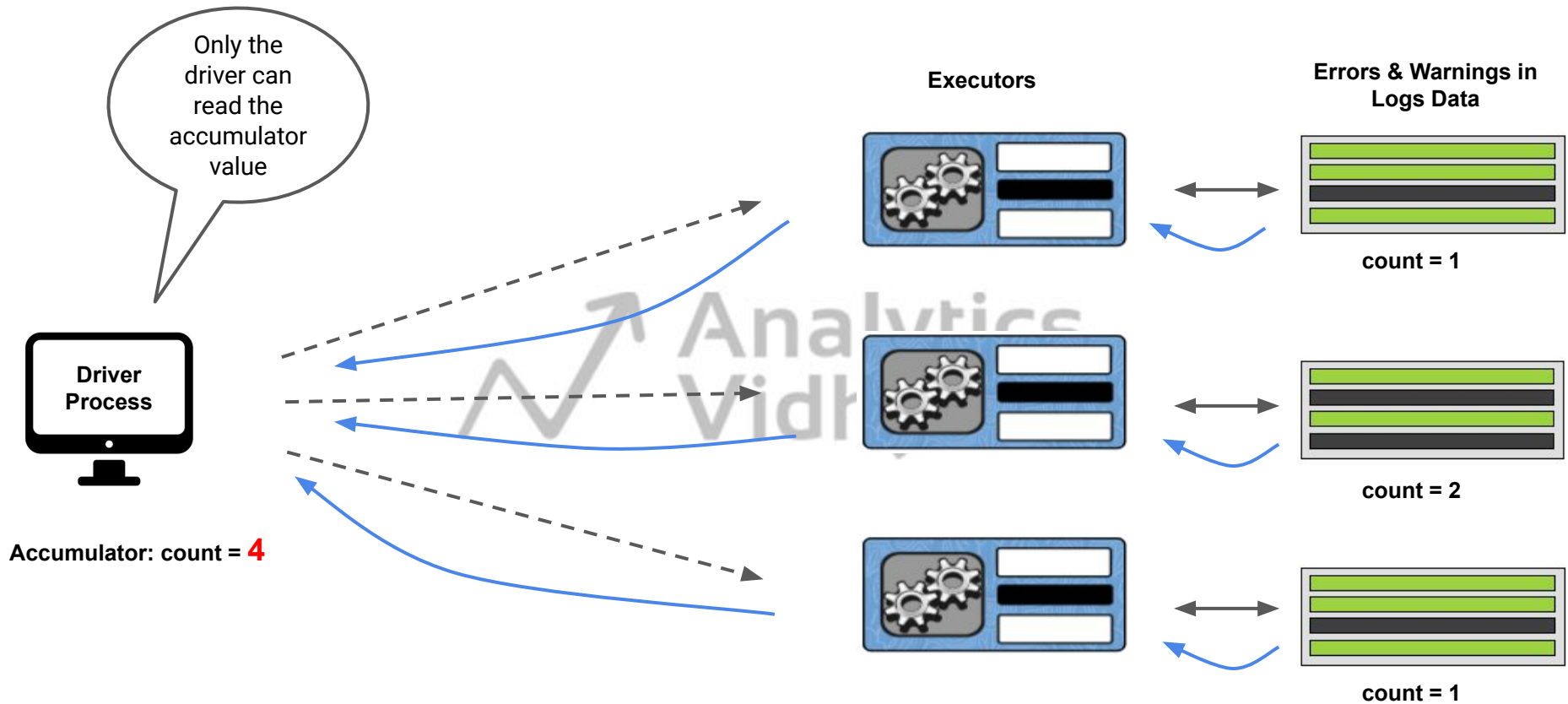
  add support for new types

**Executors**

**Errors & Warnings in Logs Data**

count = 1

**Driver Process**

Accumulator: count = **0**

# How to create Accumulators?

**Accumulator** keyword with Spark Context

**Function that will execute on each executor.**

**foreach function**

**driver program can read the accumulator value.**

```
num = sc.accumulator(10)

def f(x):
    global num
    num+=x

rdd = sc.parallelize([20,30,40,50])
rdd.foreach(f)

final = num.value
```

Analytics Vidhya

Thank You!!