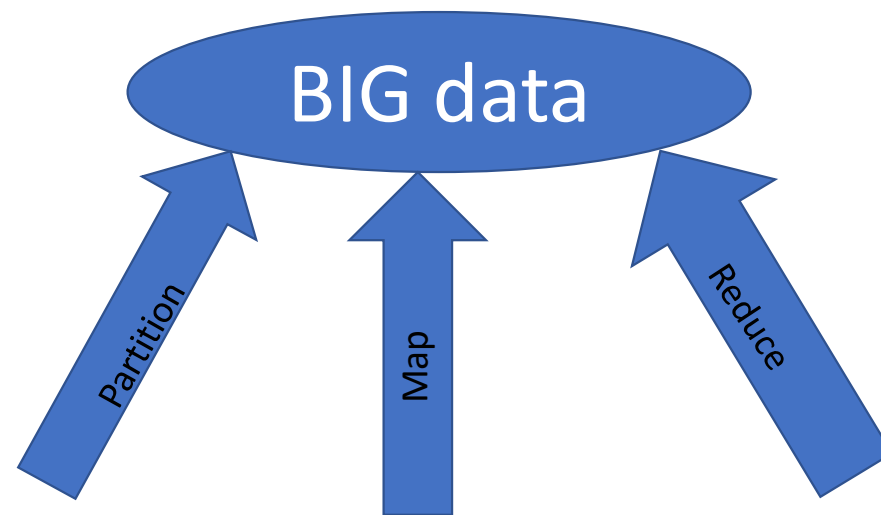


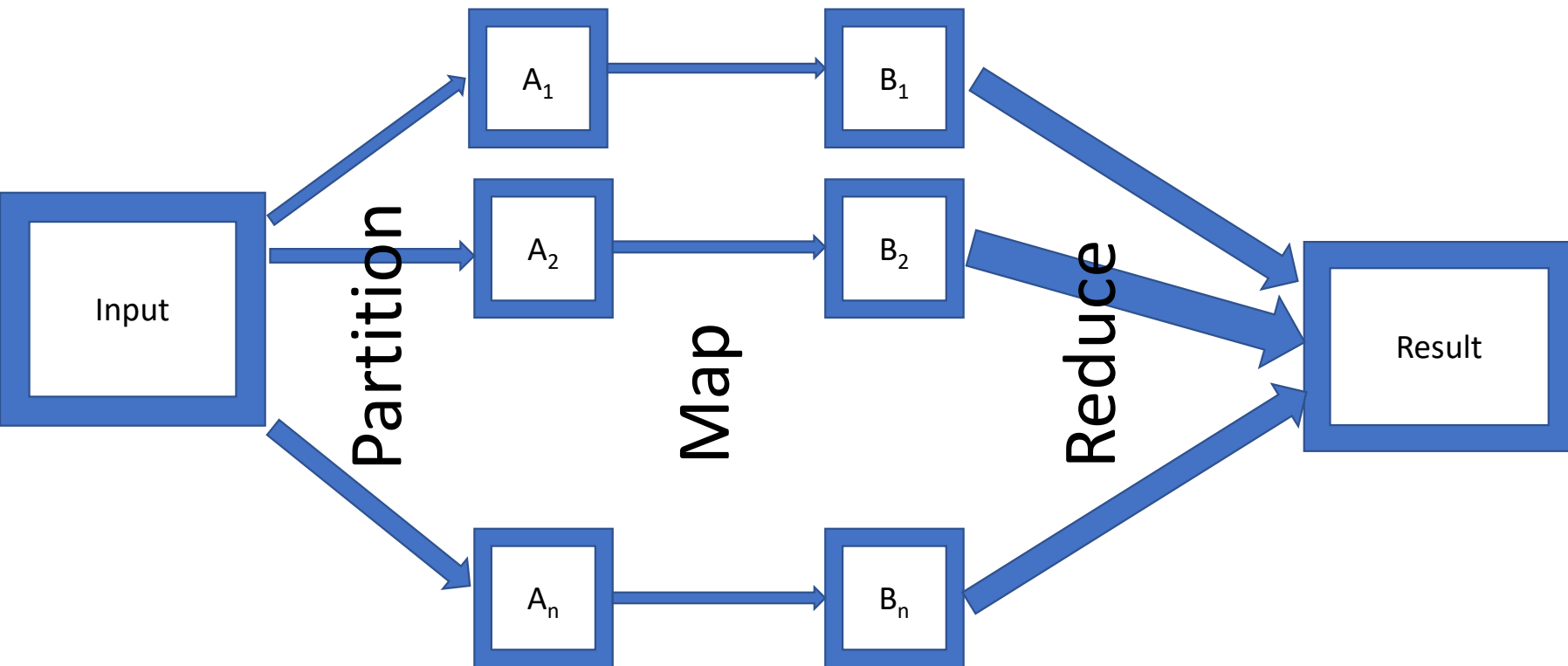
What is BIG data?

- Not solvable by “normal” DB (Oracle etc)
- Not solvable by “distributed” DB (mongoDB, Cassandra)
- “Lax” notions of transactions and availability.
- Fits a certain model (map-reduce) of problems



Three legs of the BIG data stool

Partition, Map (Parallel), Reduce (Concurrent)



Partition

Partitions the “work” into pieces

Map (Math function, especially for Spark)

$f(x) \Rightarrow y$

Reduce (Math function, especially for Spark)

Reduces a collection to a single value

Associate operation (like addition, multiplication, min, max)

$f(y1, y2) \Rightarrow z$

$f(y2, y1) \Rightarrow z$

Simplest MapReduce (Do Nothing)

Map: $f(x) \Rightarrow x$

Reduce: $f(y_1, y_2) \Rightarrow [y_1, y_2]$

Count # of lines

Map: $f(x) \Rightarrow 1$

Reduce: $f(y_1, y_2) \Rightarrow y_1 + y_2$

Count lines with Colorado

Map: $f(x) \Rightarrow 1$ if(line contains Colorado)
 $\Rightarrow 0$ otherwise

Reduce: $f(y1, y2) \Rightarrow y1 + y2$

Count lines for each State

Map: $f(x) \Rightarrow (\text{Colorado}, 1)$

$\Rightarrow (\text{Alaska}, 1)$

.

.

.

if(line contains Colorado)

if(line contains Alaska)

And so on

Reduce**ByKey**: $f(y1, y2) \Rightarrow \text{Sum by Key (State), (Repartition is implied)}$

Produces Grouping behavior. Can implement org based grouping

Hadoop (first generation)

Maximum **ONE** of

Partition

Map

Reduce

Spark (2nd generation)

Arbitrary number of

Partition

Map

Reduce

Chained in a pipeline.

Allows for partition optimization (reduces network traffic between workers)

Very hard problem. At the level of writing an OS.

Spark Streaming

Spark behavior applied to a “units of work” in time.