

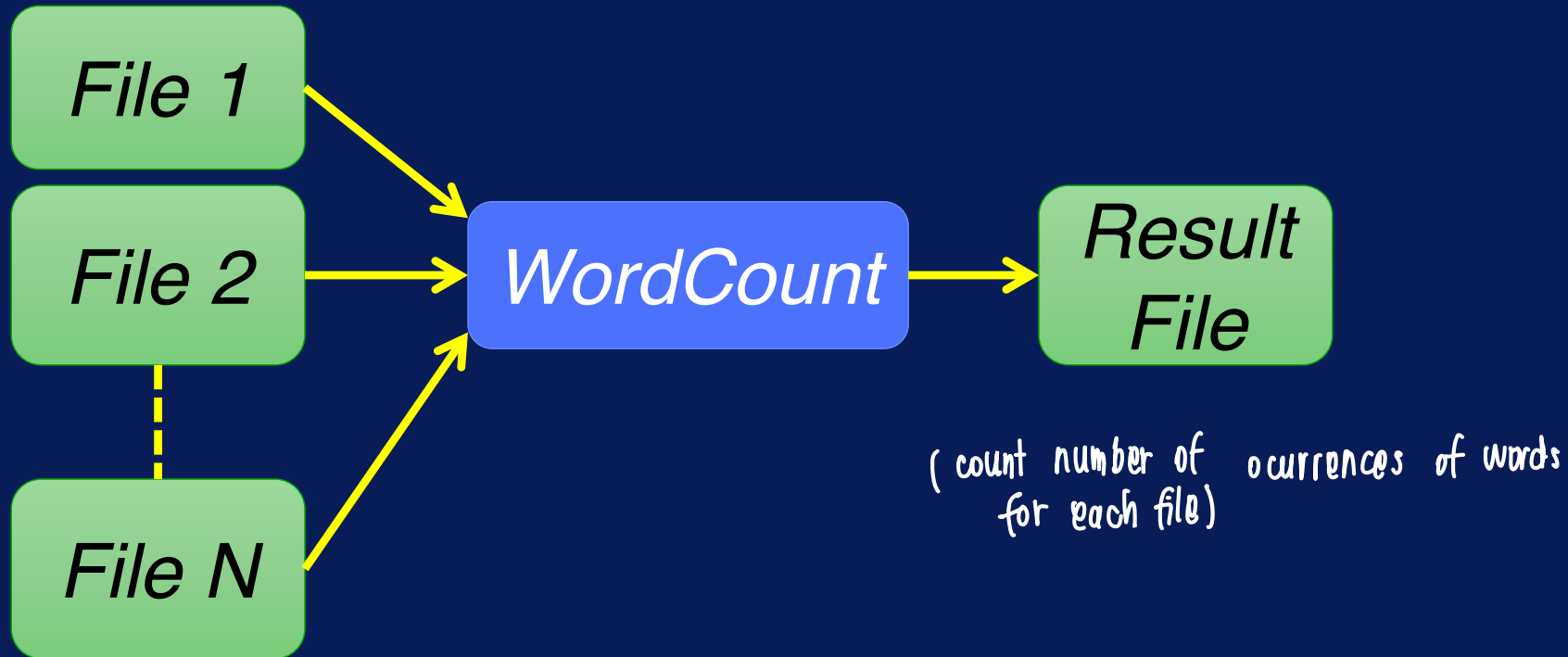
Big Data Processing Pipelines:

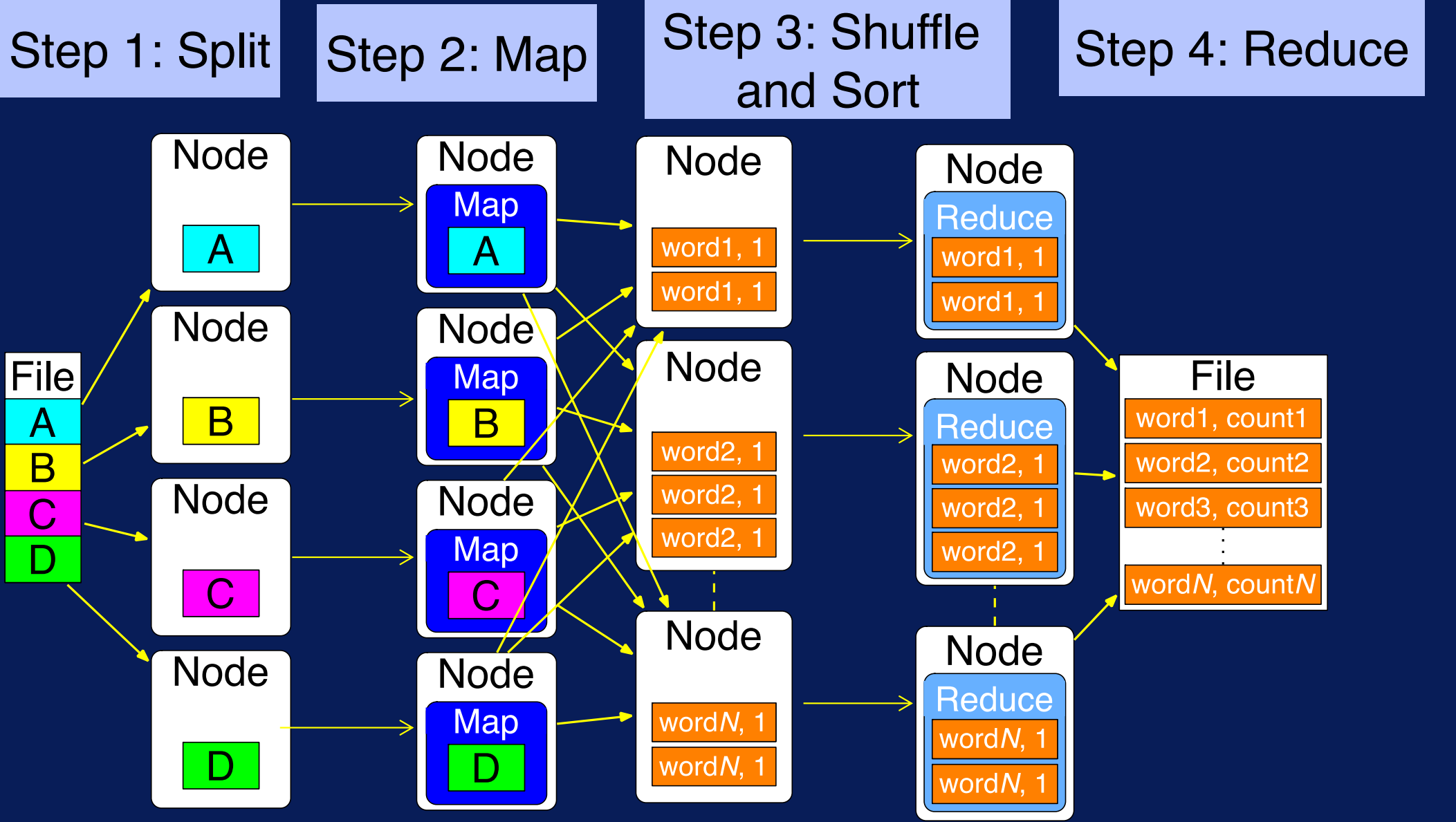
A Dataflow Approach

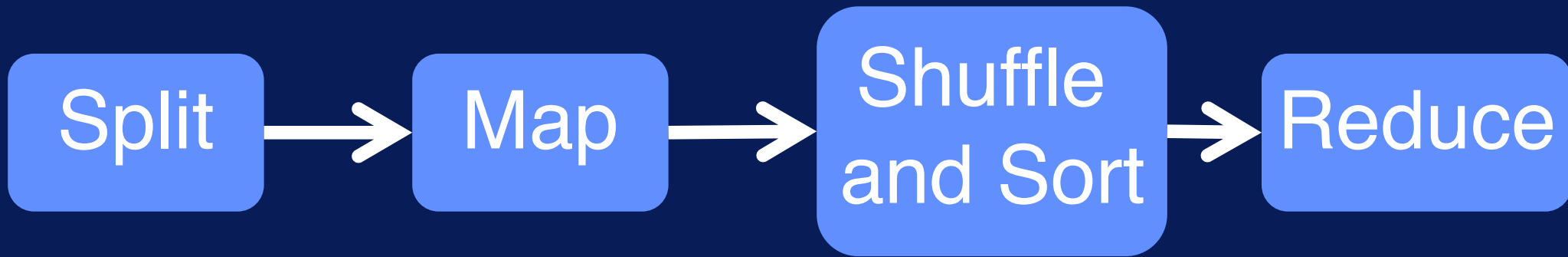
After this video you will be able to..

- Summarize what dataflow means and its role in data science
- Explain “split->do->merge” big data pipeline with examples
- Define the terms data parallel

Example MapReduce Application: WordCount










Represents a large
number of applications.

Big Data Pipelines



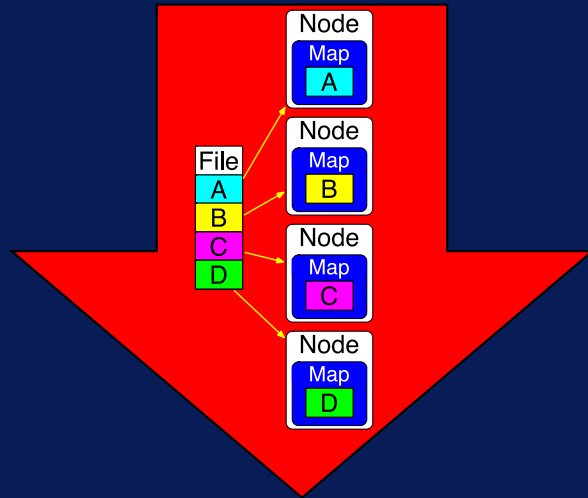
cat  sort

A UNIX pipe provides one-way communication
between two processes on the same computer

parallelism: running the same functions
simultaneously for the elements
or partitions of a dataset on
multiple cores

* occurs in every step





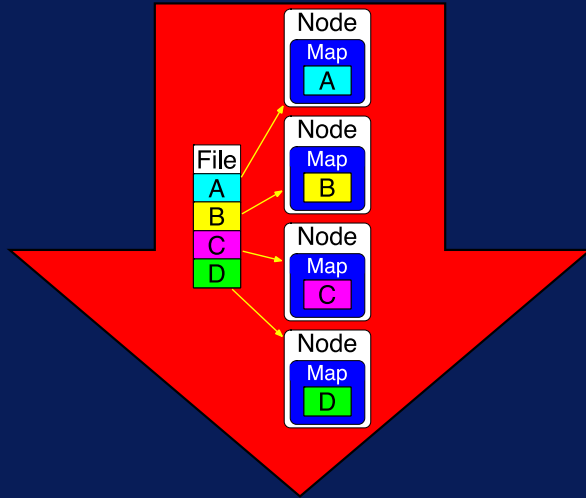
Parallelization
over the input

* decide on data
granularity of each
parallel computation

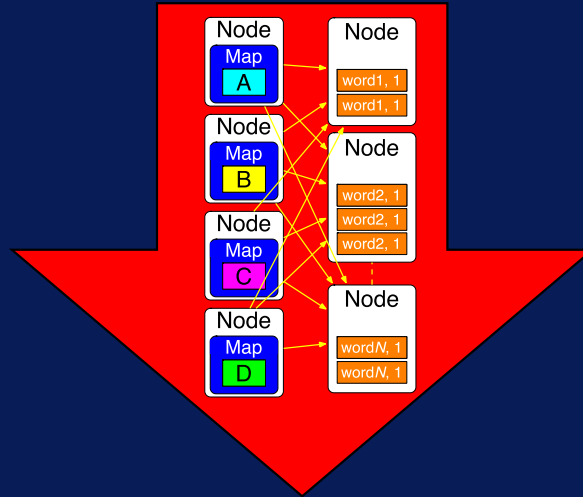
Map

Shuffle
and Sort

Reduce



Parallelization
over the input



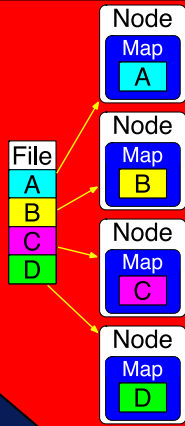
parallel grouping

Parallelization
data sorting

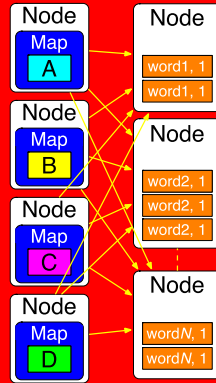
Map

Shuffle and Sort

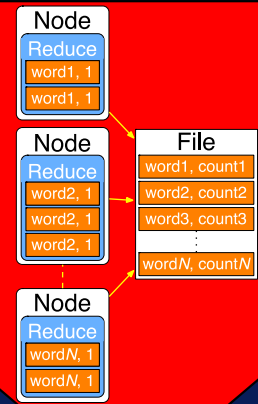
Reduce



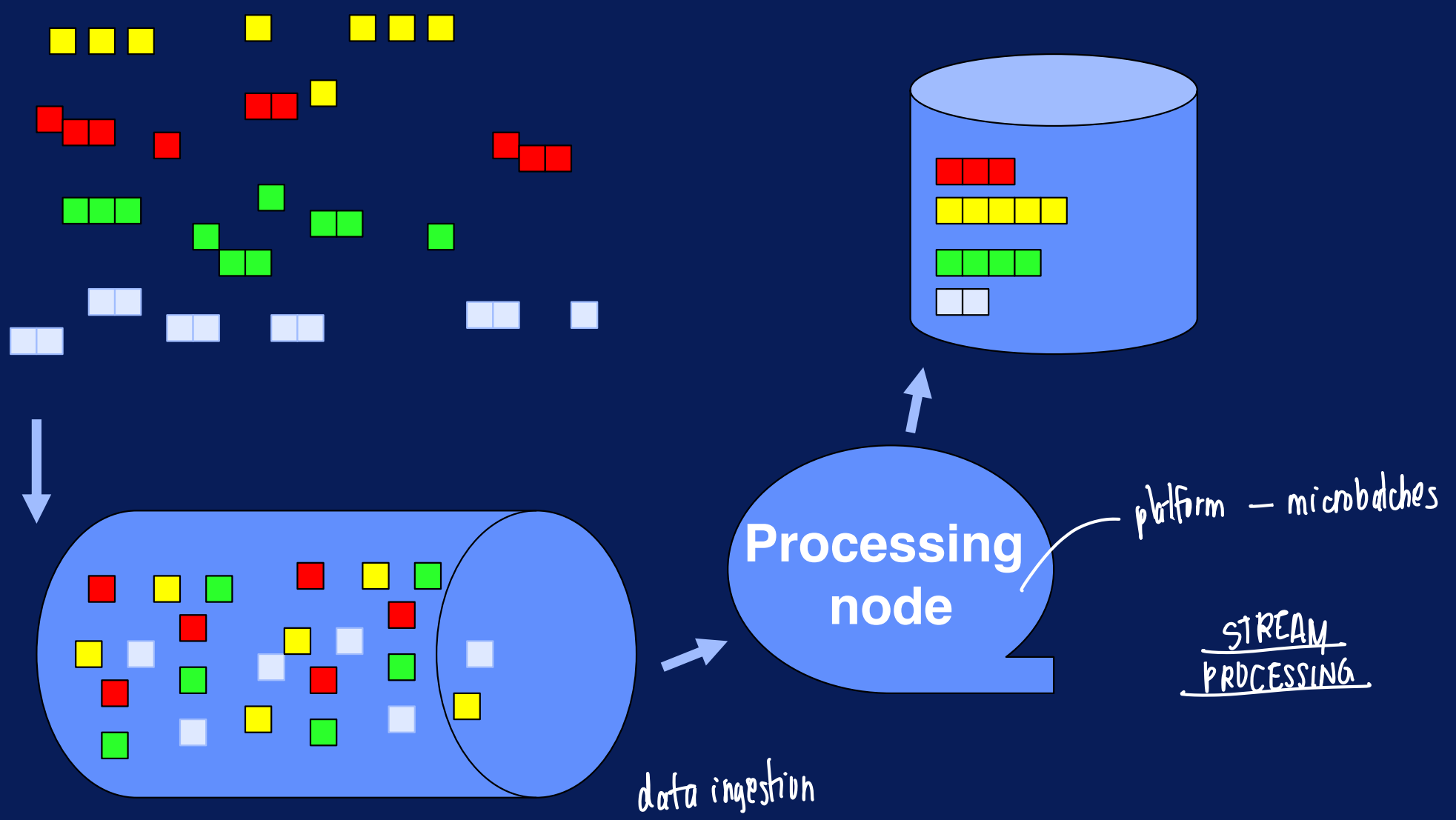
Parallelization
over the input

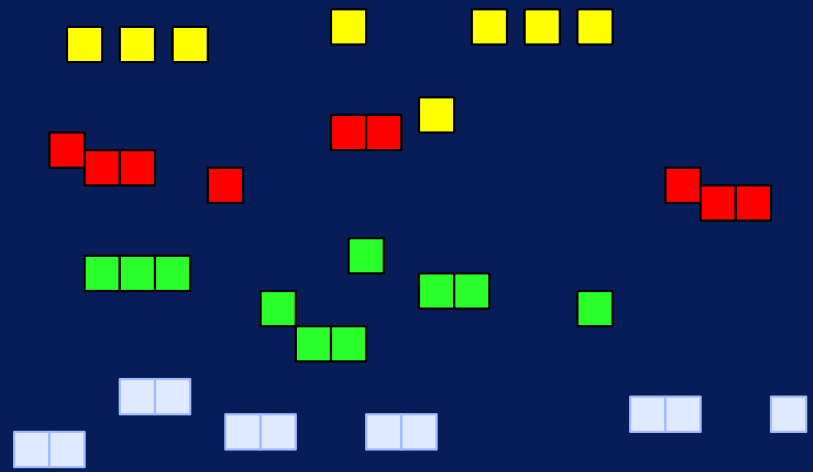


Parallelization over
intermediate data

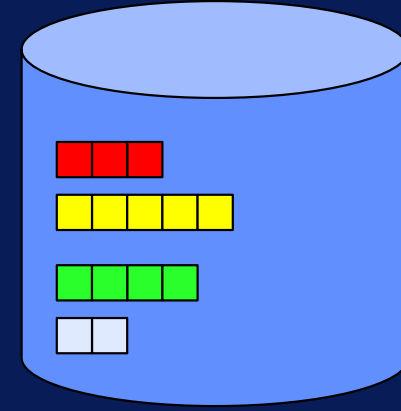


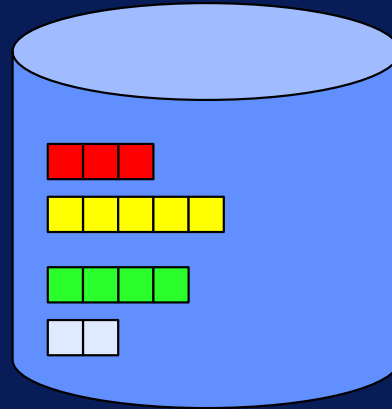
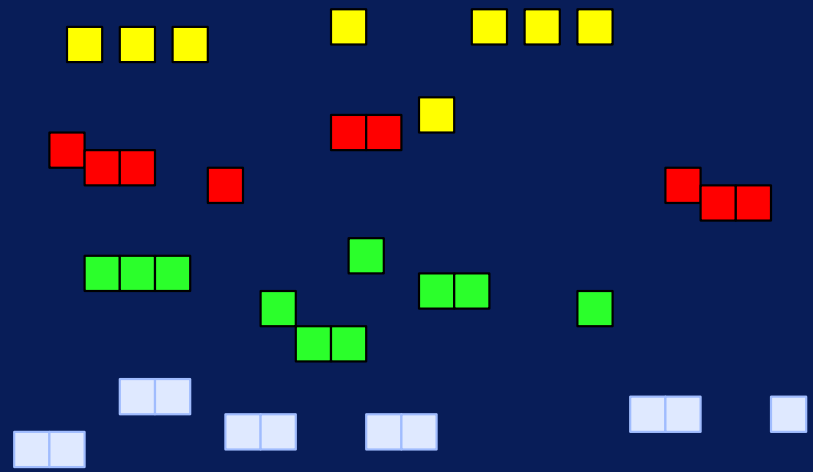
Parallelization
over data groups

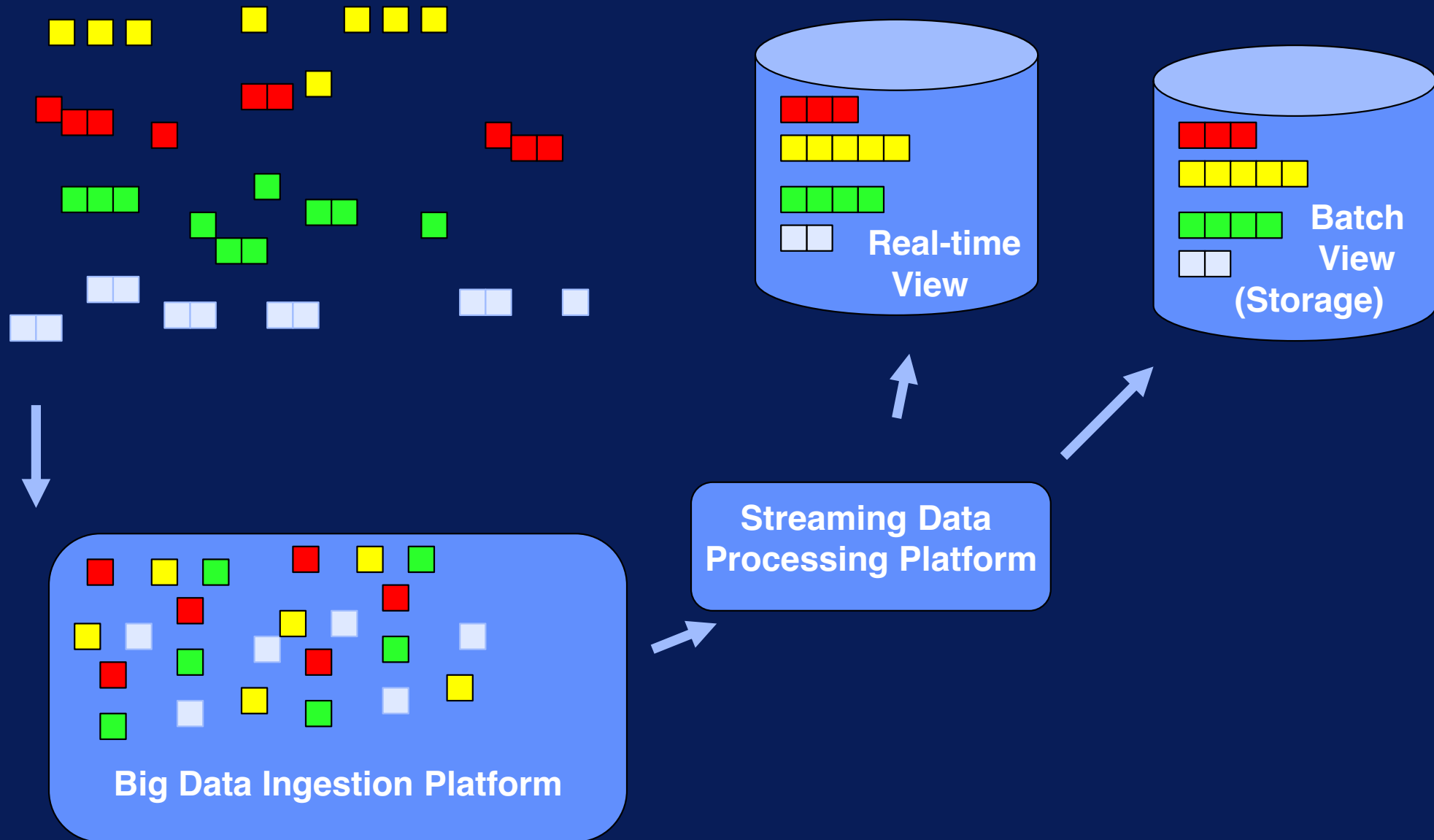




Processing
node









big data pipelines get created to process data
through an aggregated set of steps that
can be represented with the split → do → merge
with data parallel scalability