

SPARK FOR DATA ENGINEERING

Background

DATA LIFECYCLE

Raw Data (TB-PB)

Integrate
Explore
Aggregate

Use Cases
Deploy
Monitor

Insight (KB-GB)

Collect

Analyze

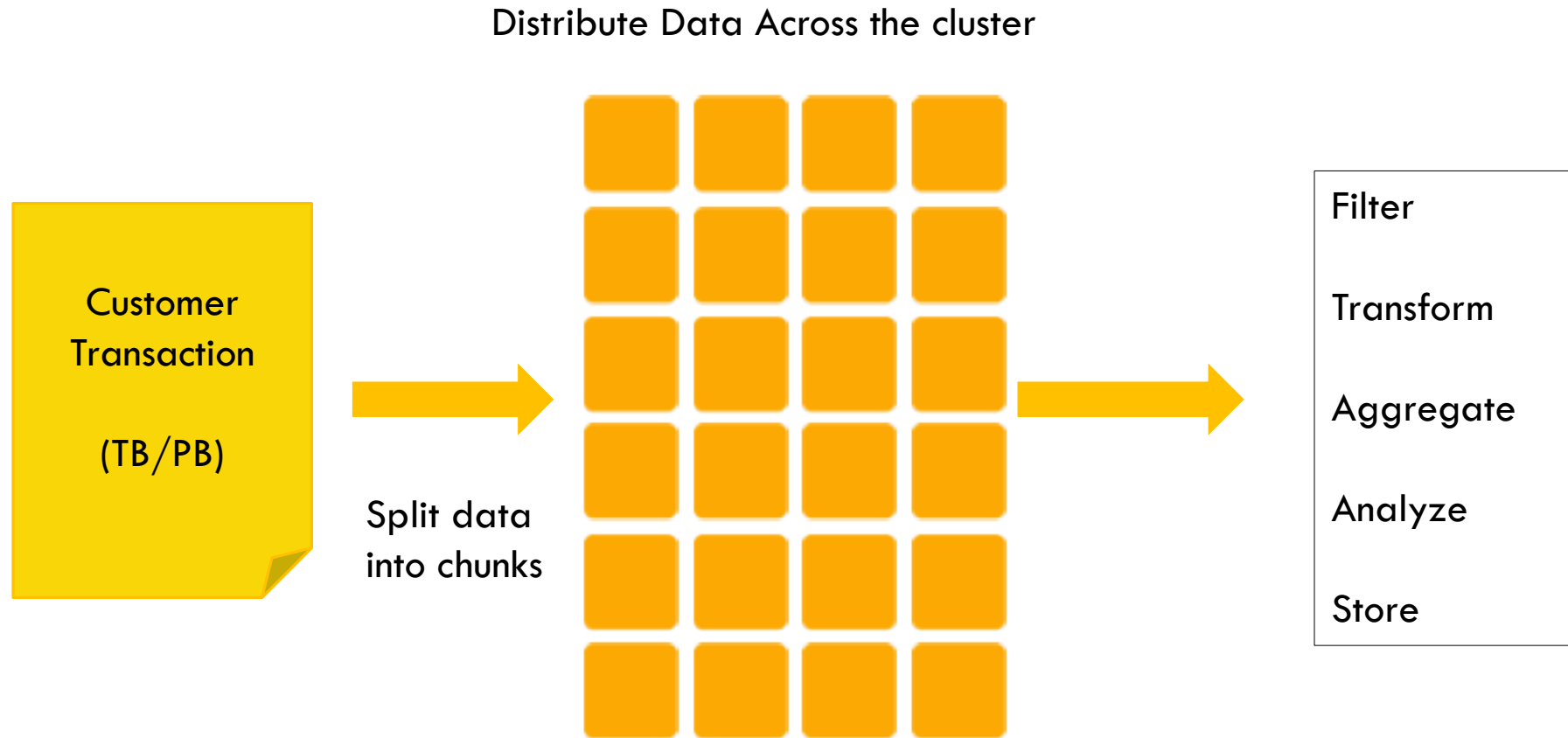
Clean

Organize

Transform

Insight

HOW SPARK WORKS?



WHY SPARK?

Data Collection



Batch



Real Time

Data Processing and Transformation



Storage and Analysis



WHY SPARK?

Amazon S3, Google
GCS, Azure, HDFS



Kafka, Kinesis, pub
sub



Baremetal

Cloud Managed Services

Cloud Kubernetes Services

WHY SPARK?

Baremetal

Yarn

Standalone

Apache Mesos

Kubernetes

Cloud Managed Services



 databricks

Kubernetes Services



Azure Kubernetes Service (AKS)



Google Kubernetes Engine



Amazon
EKS

SPARK ARCHITECTURE

