# The Stratosphere platform for big data analytics

Presentation by Daniyal Warsi, Timo Kraus & Sebastian Hofmann

# Table of Contents
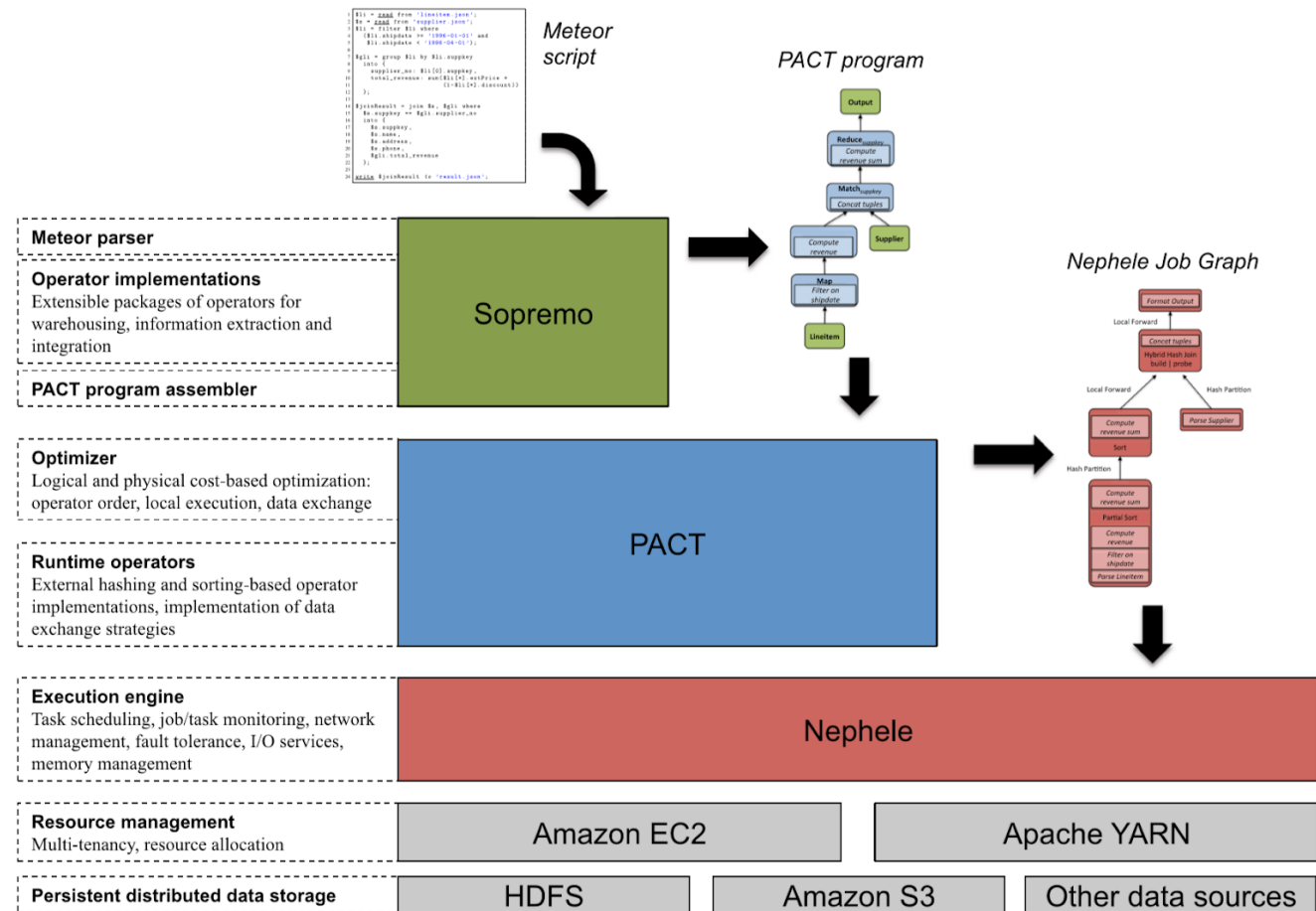
- Goal

- Stratosphere's Architecture
  - Sopremo
  - PACT
  - Nephele

- Fog Readiness

- IoT Use Case

# Goal

The Stratosphere platform …

- provides a Big Data Analytics Platform that brings together high & low Level Programming.

- enhances traditional RDBMS with the Ability to cope with heterogenous Datasets.

- enables direct Connection to external Datasets.

- extends the MapReduce Operators for efficient Execution.

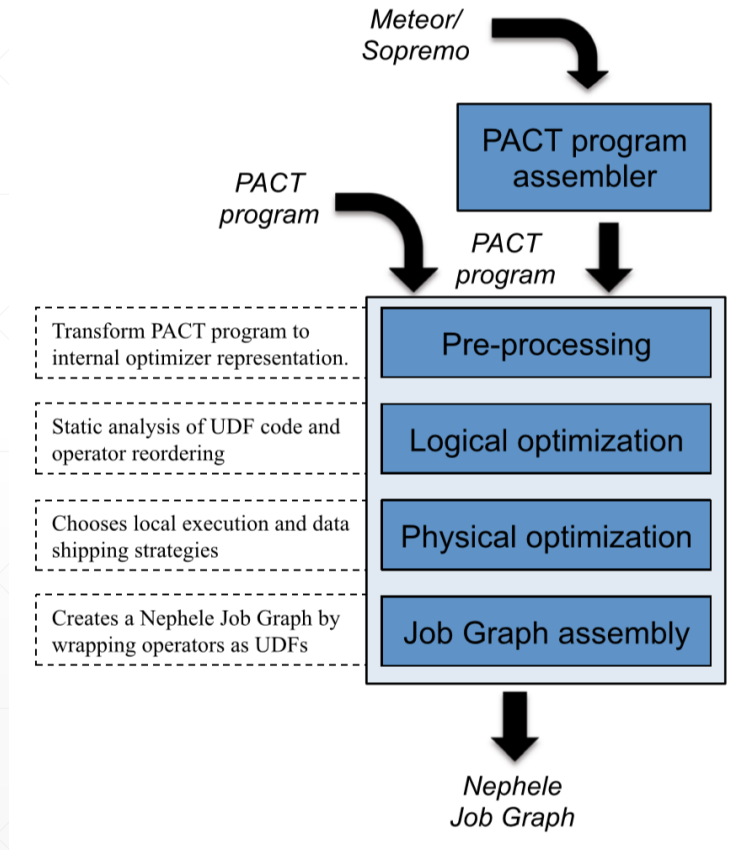# Stratosphere's Architecture



(Alexandrov et al., 2014, p.942)

# Sopremo

- Declarative Query Language

- Meteor as the textual Interface for the Sopremo Layer

- Diverse set of composed Operators (Select, Project, etc.)

- Outputs a Logical Query Plan
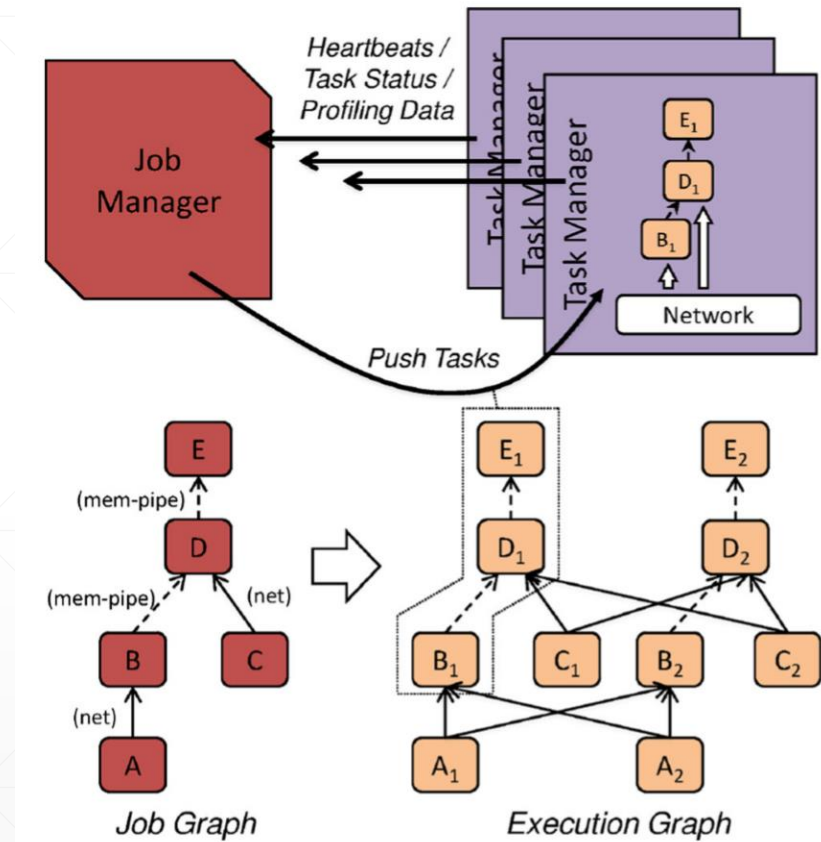
- May be enriched with Metadata for lower levels

# PACT



- Generalization of MapReduce

- Allows the Implementation of more complex Operators

- Includes a Query Optimizer for performant Execution

- Decomposes and reorganizes Operators

- Allows incremental Iterations for Pipelining

(Alexandrov et al., 2014, p.951)

# Nephele



- Layer containing the Parallel Execution Engine

- Distributed Master/Worker Execution

- Includes Memory & I/O Services

- Collects Statistics for PACT Optimization

(Alexandrov et al., 2014, p.954)

# Fog Readiness

## Pros

- Nephele Master/Worker Distribution

- Bandwidth via Data Compression on Nodes

- Heterogenous Data Handling in IOT Use Cases

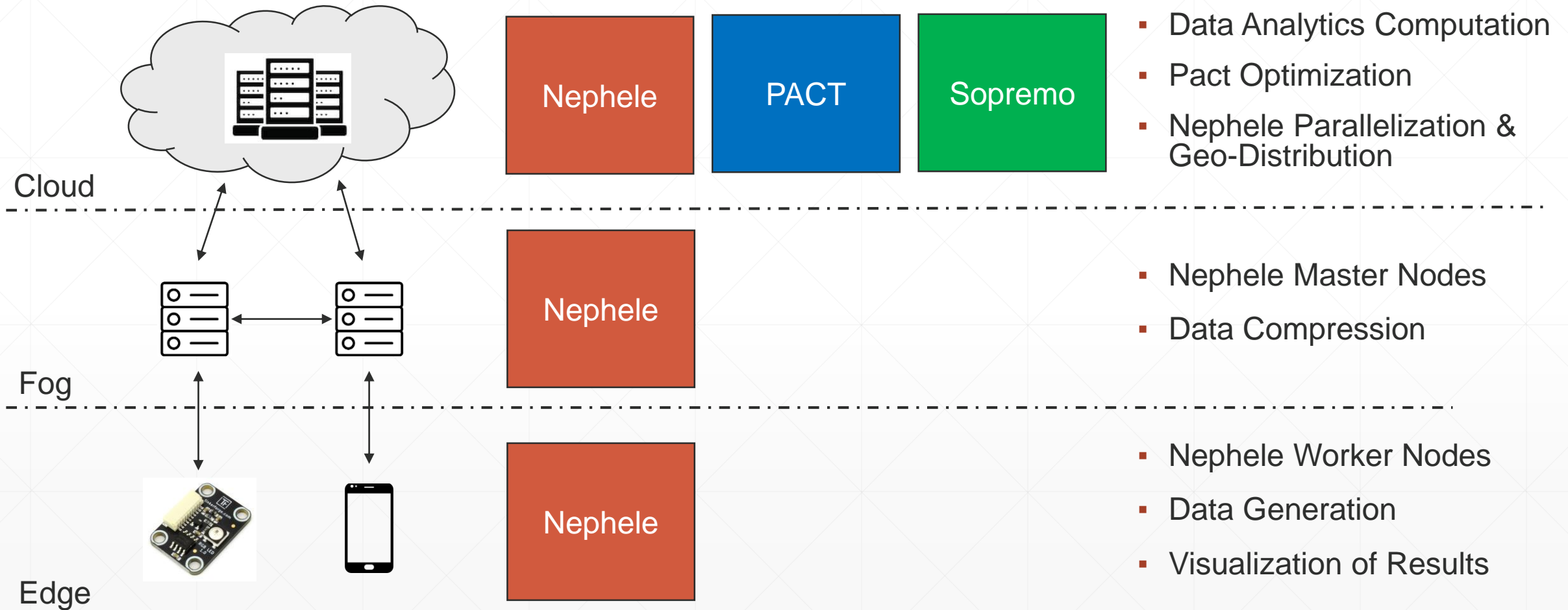- Fault Tolerance via Buffering and Checkpointing

## Cons

- Job Manager = Single Point of Failure

- Complexity in Query Optimization

- Testing on multiple Layers adds Complexity

## Neutral

- Nephele Parallelization supports Custom Events → Implementation for Geo-Distribution

- Security Aspects not discussed

# Stratosphere in an IoT Use Case



Cloud

- Data Analytics Computation
- Pact Optimization
- Nephele Parallelization & Geo-Distribution

Fog

- Nephele Master Nodes
- Data Compression

Edge

- Nephele Worker Nodes
- Data Generation
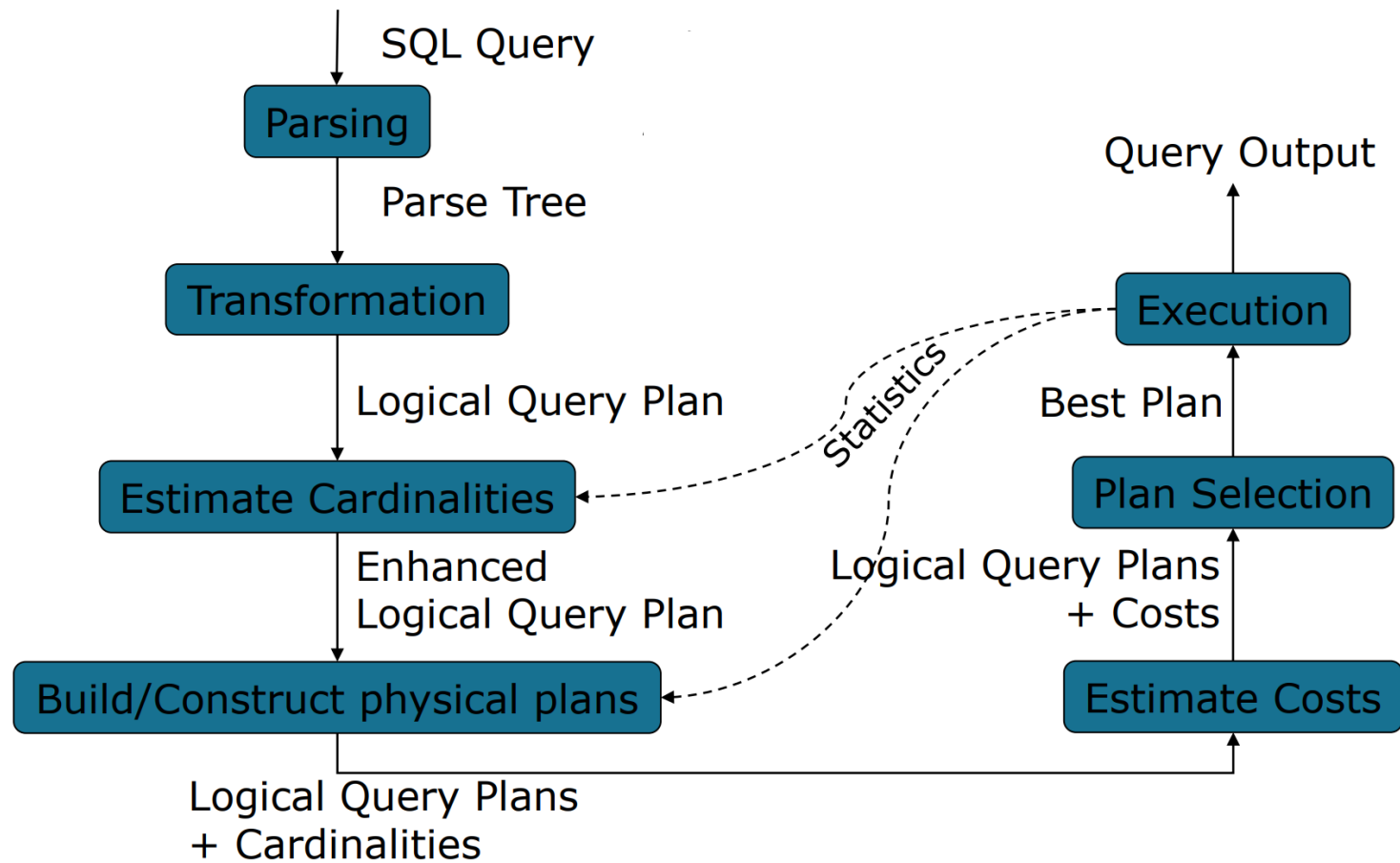- Visualization of Results
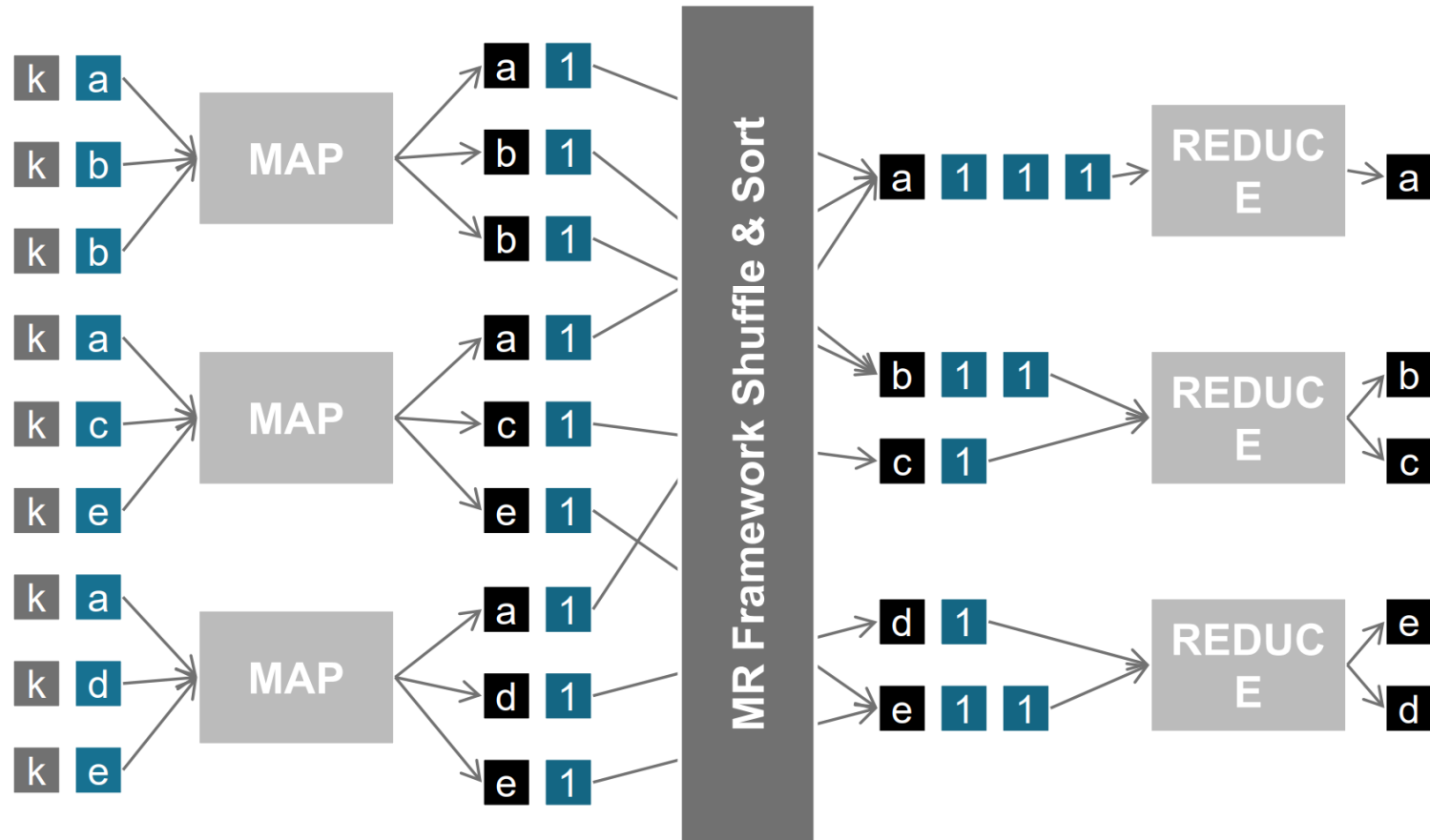
# Sources & Related Work

- **Original Paper:** Alexandrov et al. (2013): The Stratosphere platform for big data analytics

- Dean, Ghemawat (2004): MapReduce: Simplified Data Processing on Large Clusters

- Ang, Seng (2016): Big Sensor Data Applications in Urban Environments

- Shvachko et al. (2010): The Hadoop Distributed File System

- Rabl (2018): Internals of Database Systems: Web-Scale Data Management - Analytical Processing

# Annex

# Query Processing

# MapReduce



(Rabl, 2018, p.18)

# Bulk vs. incremental Iteration



(a)

(b)