

# FlowPROPHET: Generic and Accurate Traffic Prediction for Data-parallel Cluster Computing

**Hao Wang**<sup>1,2</sup>, Li Chen<sup>2</sup>, Kai Chen<sup>2</sup>, Ziyang Li<sup>2,3</sup>, Yiming Zhang<sup>3</sup>,  
Haibing Guan<sup>1</sup>, Zhengwei Qi<sup>1</sup>, Dongsheng Li<sup>3</sup>, Yanhui Geng<sup>4</sup>

<sup>1</sup> Shanghai Jiao Tong University

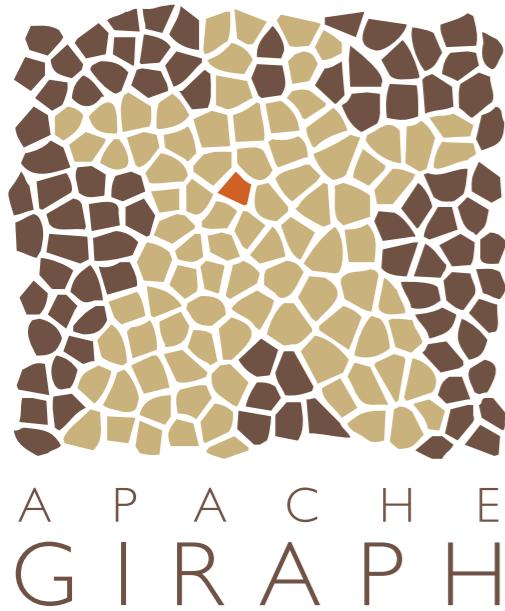
<sup>2</sup> Hong Kong University of Science and Technology

<sup>3</sup> National University of Defense Technology

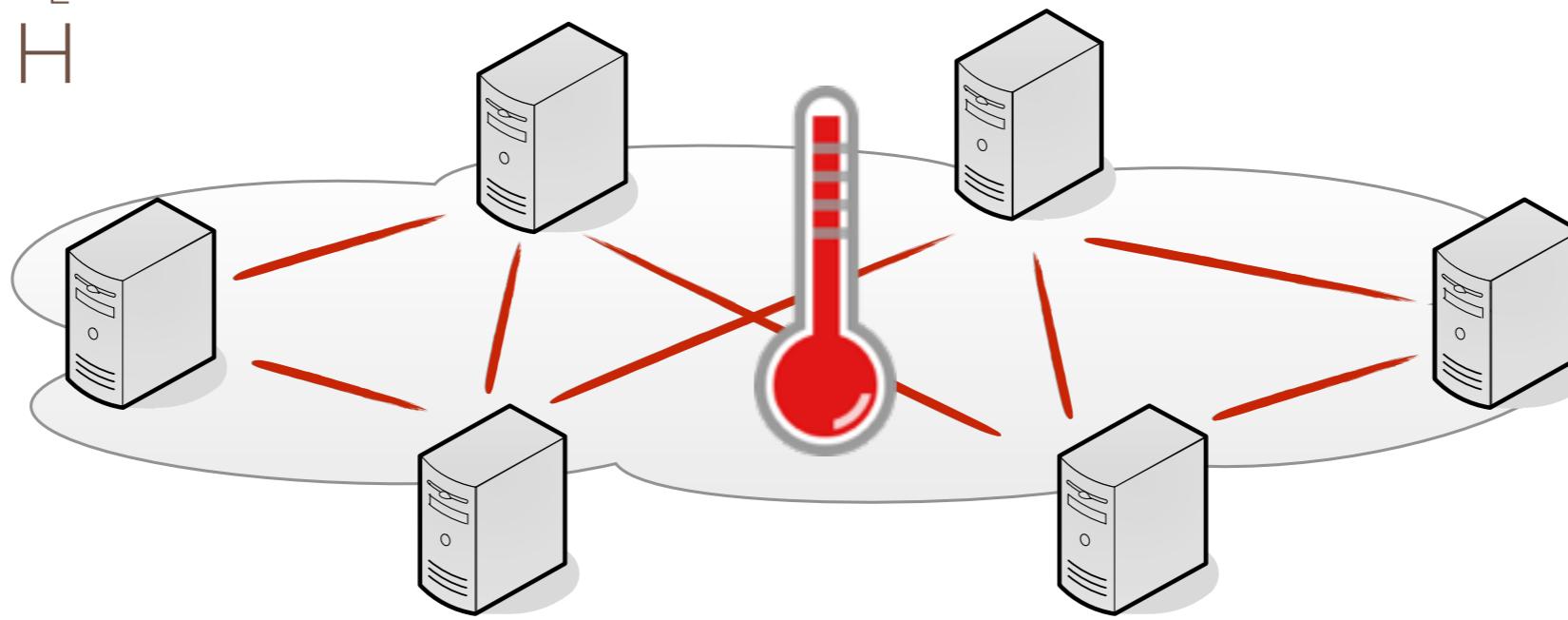
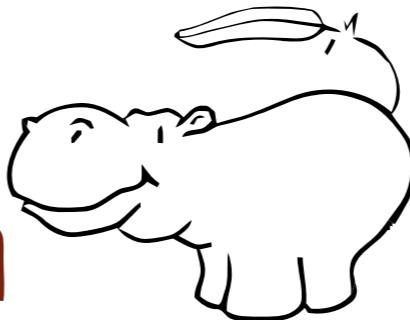
<sup>4</sup> Huawei Technologies Co. Ltd.



The word cloud is composed of various terms related to machine learning, data analysis, and algorithms. Key words include 'machine', 'algorithms', 'data', 'statistical', 'learning', 'graphs', 'clustering', 'analysis', 'guarantees', 'extract', 'use', 'amounts', 'classification', 'able', 'similarity', 'analyze', 'becomes', 'principled', 'discipline', 'arise', 'best', 'across', 'properties', 'compounds', 'understood', 'communities', 'Machine', and 'urgent'. The words are in different sizes and colors (brown, purple, blue, orange) and are interconnected by thin lines.



Apache Hama



Apache Spark logo: The word 'Spark' in a bold, black, sans-serif font, with an orange star icon above the letter 'k'.  
Dryad logo: The word 'Dryad' in a large, dark gray, sans-serif font.  
Microsoft logo: The Microsoft logo, which consists of four colored squares (red, green, blue, yellow) arranged in a 2x2 grid.

## **Flow-based optimization mechanisms:**

- PDQ [Sigcomm'12], pFabric [Sigcomm'13],  
PASE [Sigcomm'14], Varys [Sigcomm'14],  
Baraat [Sigcomm'14]

## **Architectural bandwidth provisioning:**

- c-Through [Sigcomm'10], Helios [Sigcomm'11],  
Mordia [Sigcomm'13], OSA [NSDI'12]

## **Traffic engineering:**

- Hedera [NSDI'10], MicroTE [CoNEXT'11],  
D<sup>3</sup> [Sigcomm'11]

# Knowing the Flow Information

Flow-based optimization mechanisms:

## Ahead of Time

- PDQ [Sigcomm'12], probNIC [Sigcomm'13],  
PASE [Sigcomm'14], Varys [Sigcomm'14],  
Baraat [Sigcomm'14]



Architectural bandwidth provisioning:

- c-Through [Sigcomm'10], Helios [Sigcomm'11],  
Mordia [Sigcomm'13], OSA [NSDI'12]

Traffic engineering:



- Hedera [NSDI'10], MicroTE [CoNEXT'11],  
D<sup>3</sup> [Sigcomm'11]

# FlowPROPHET

- Generic for DCFs
- Accurate and fined-grained
- Ahead-of-time
- Scalable and low-overhead

# Toy Example: Word Count

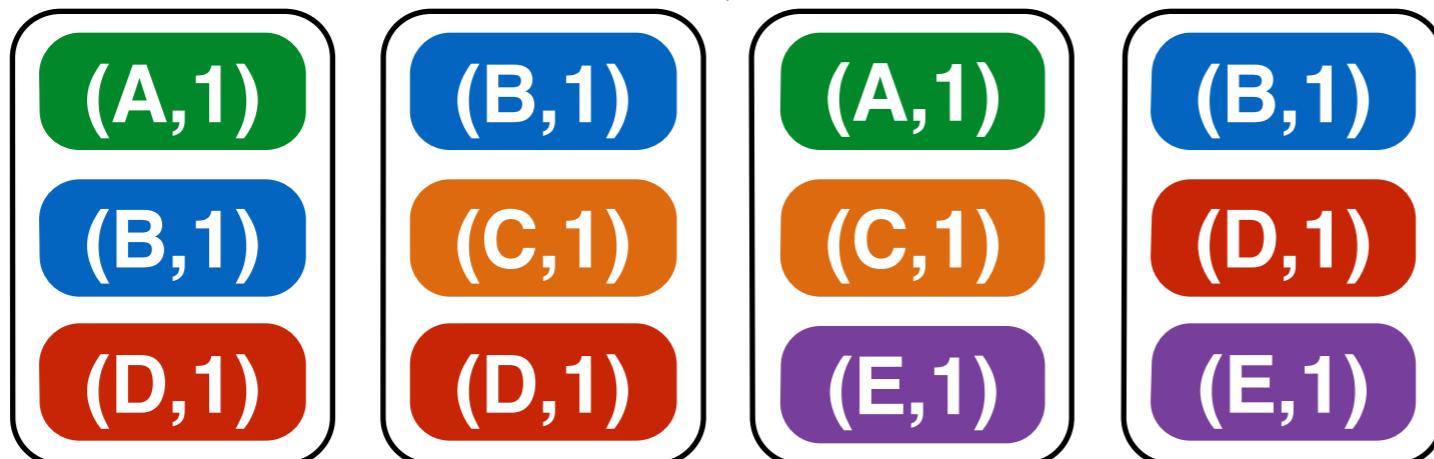
# Logical View



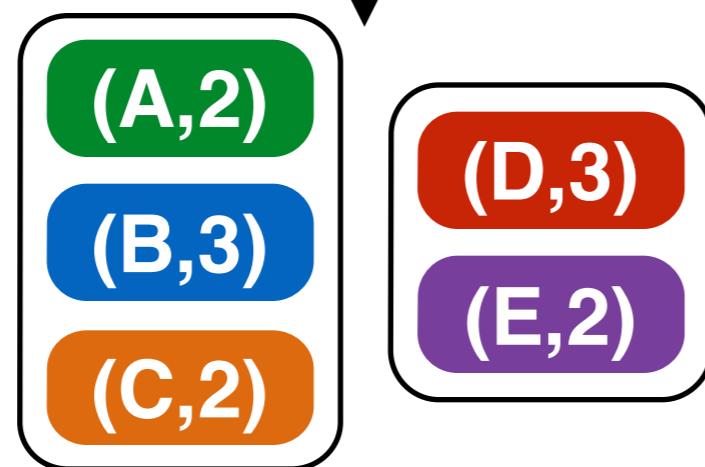
A...**D**...**C**...**A**  
**B** **D**...**E** **C**...**A**...**A** **D**...**E** **B**...



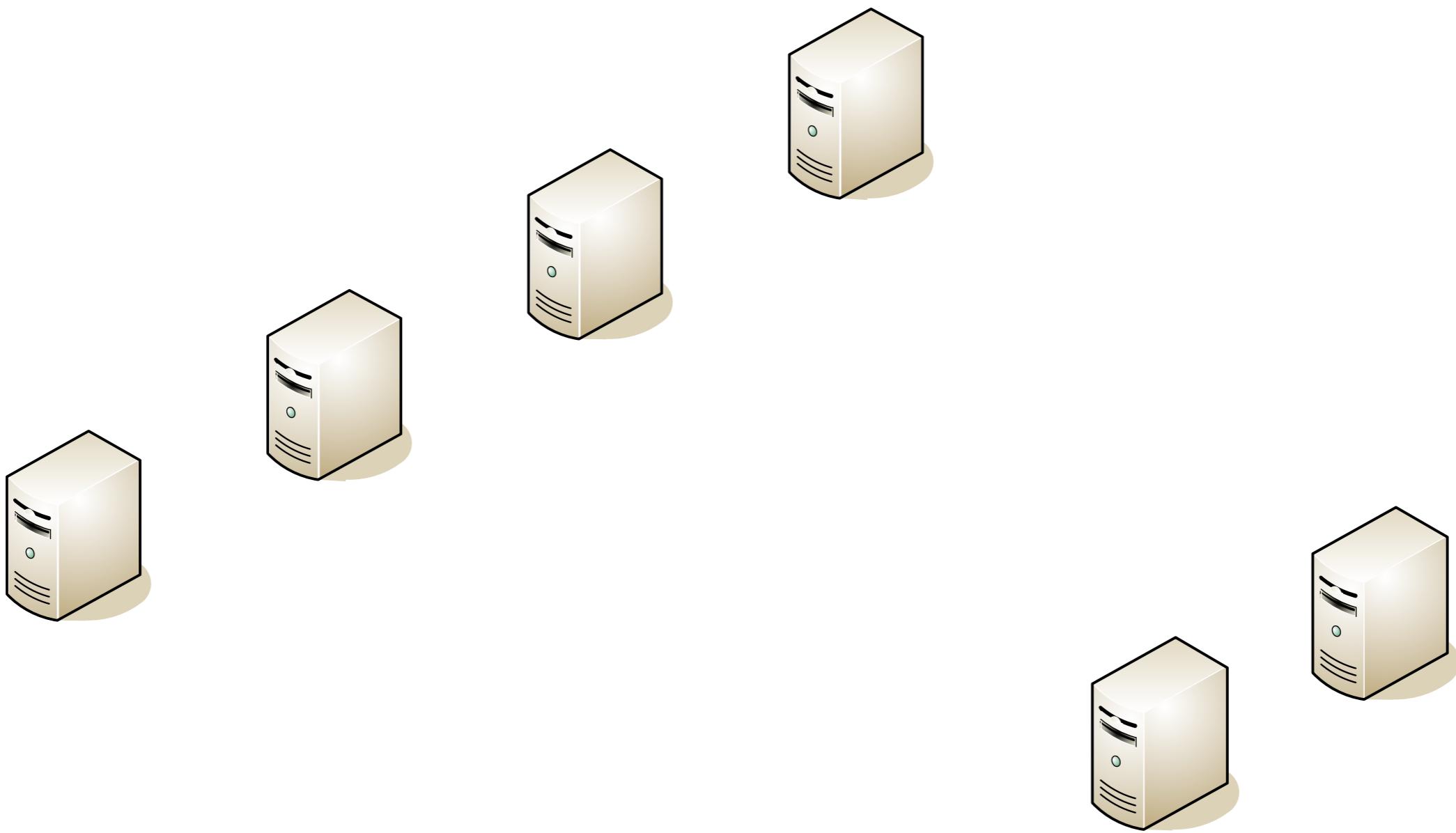
..... map()



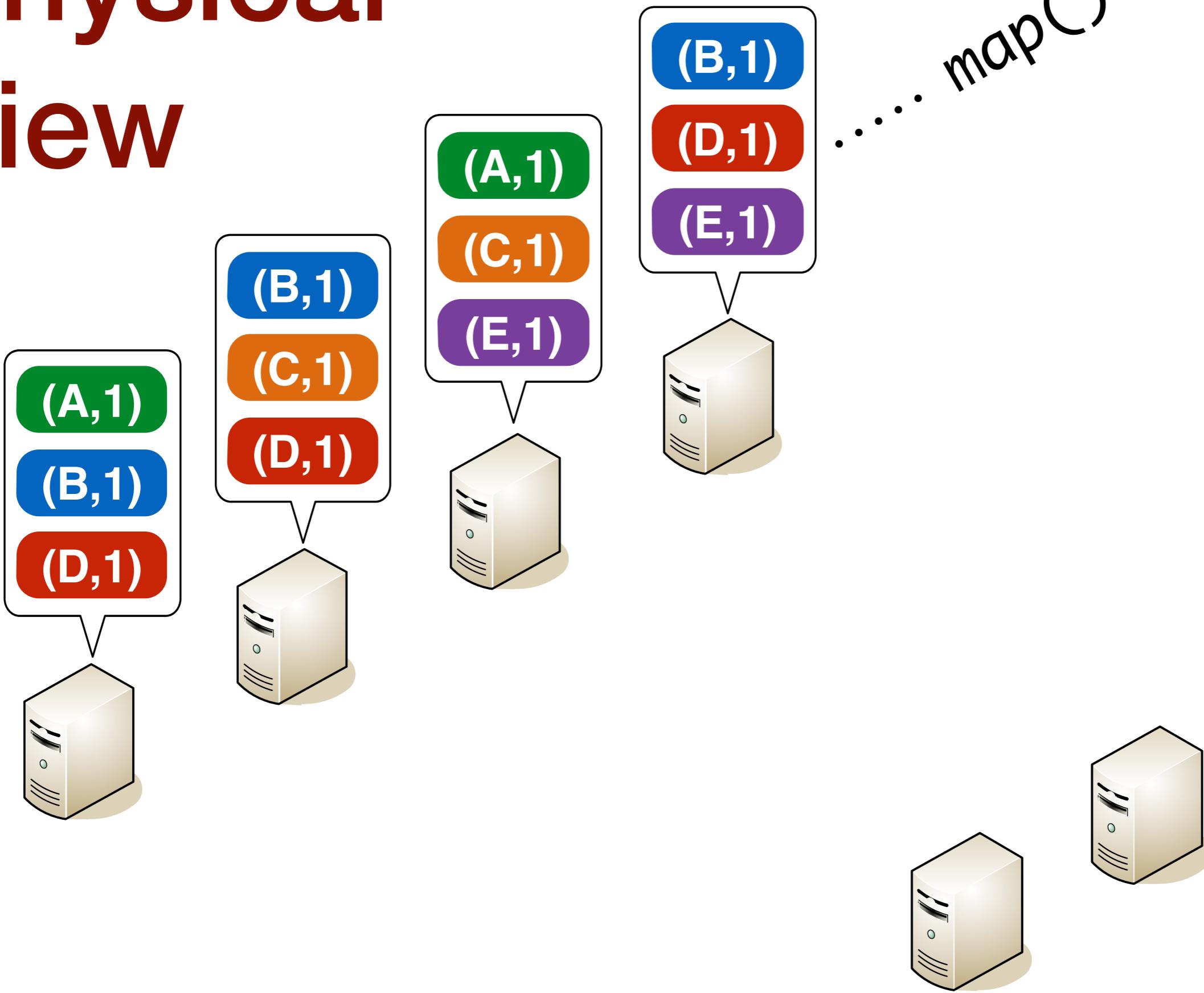
..... reduce()



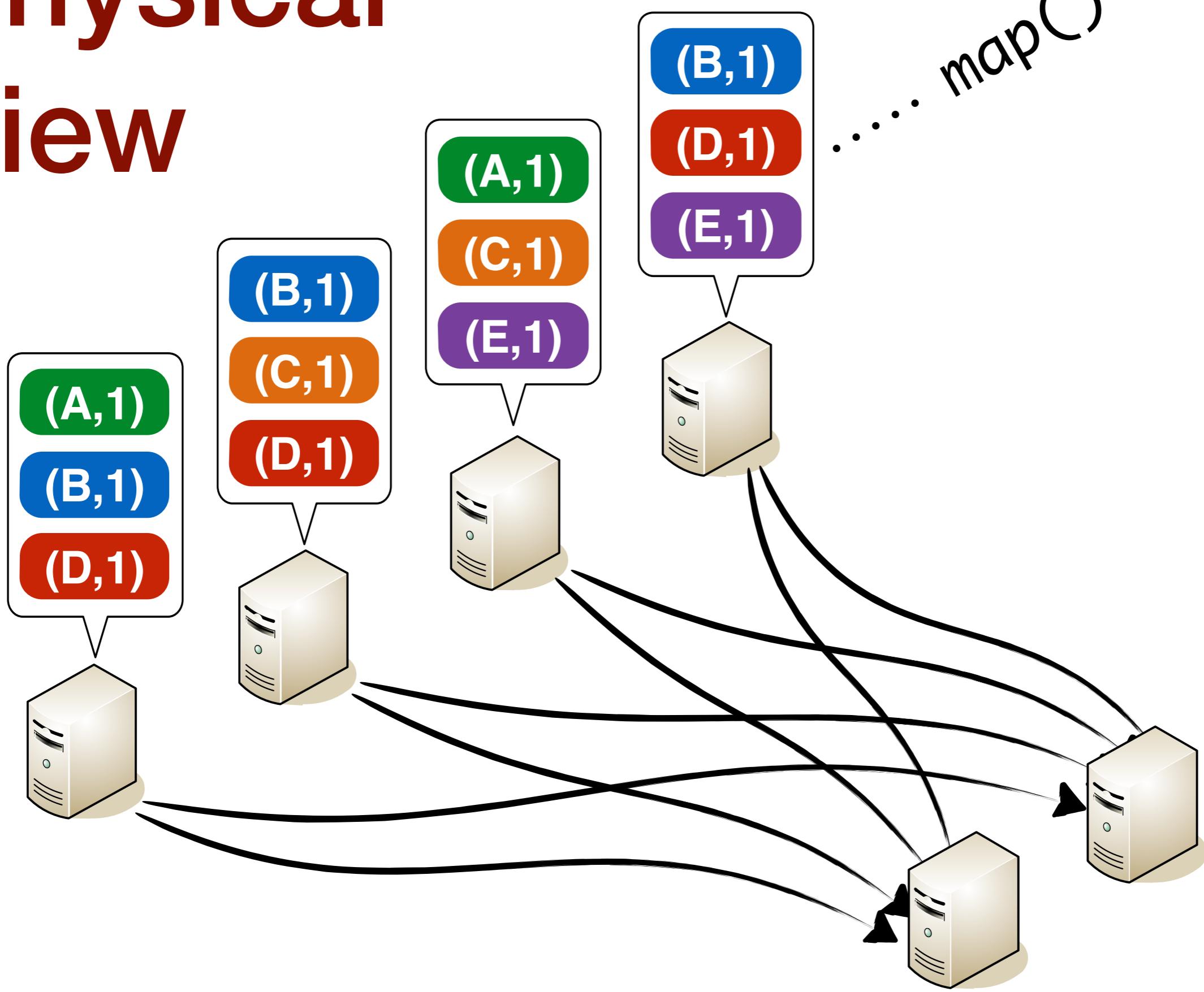
# Physical View



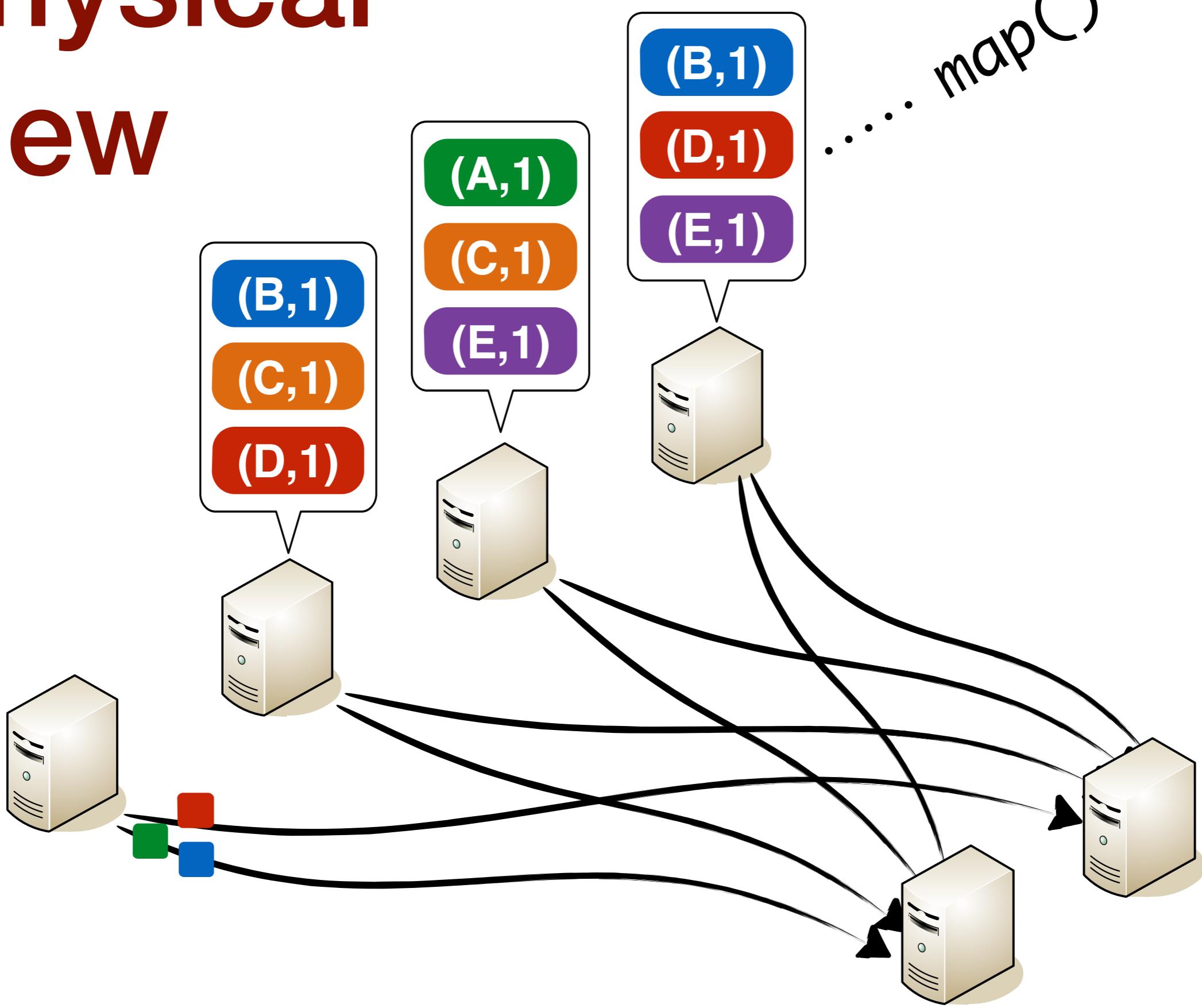
# Physical View



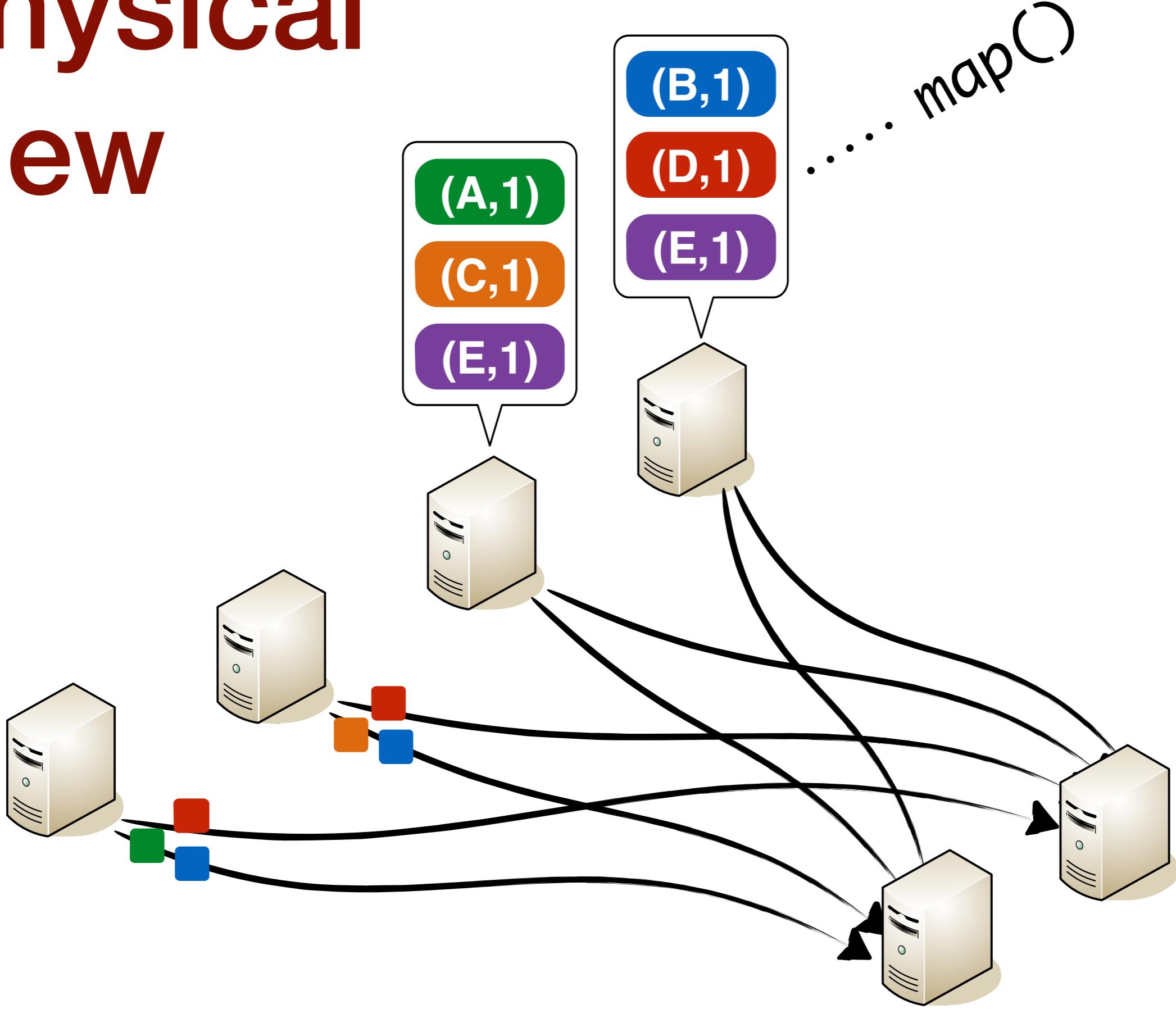
# Physical View



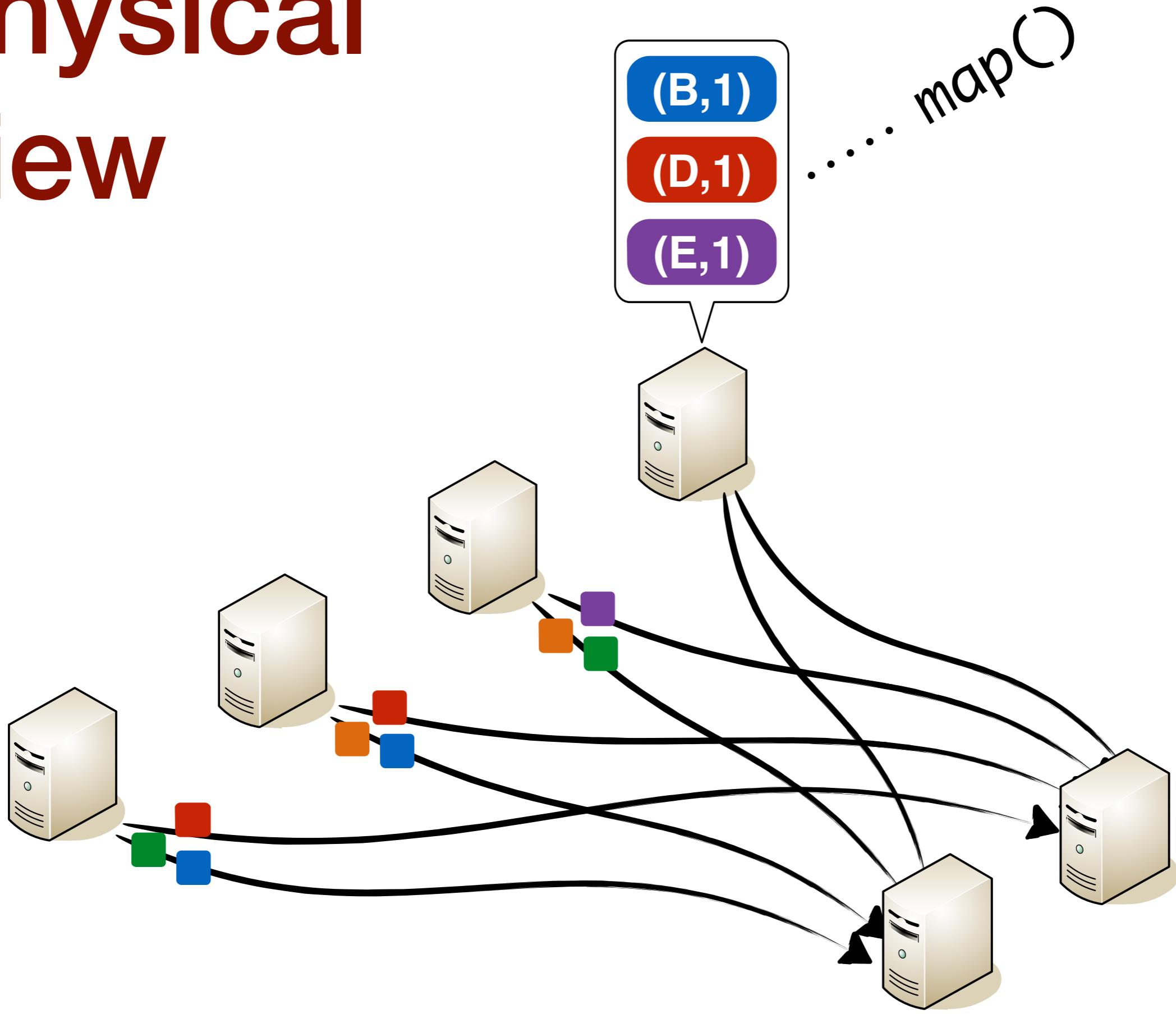
# Physical View



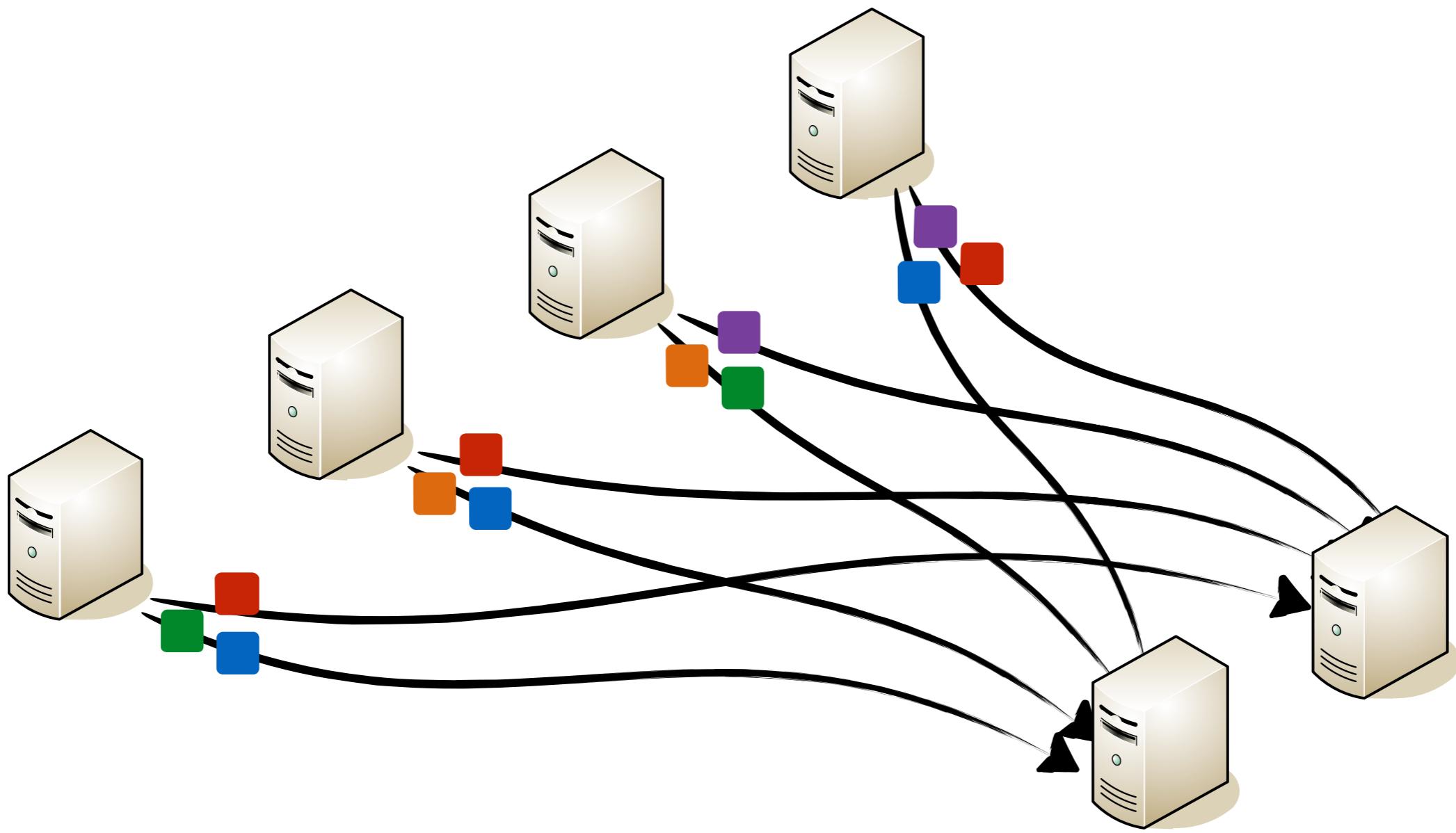
# Physical View



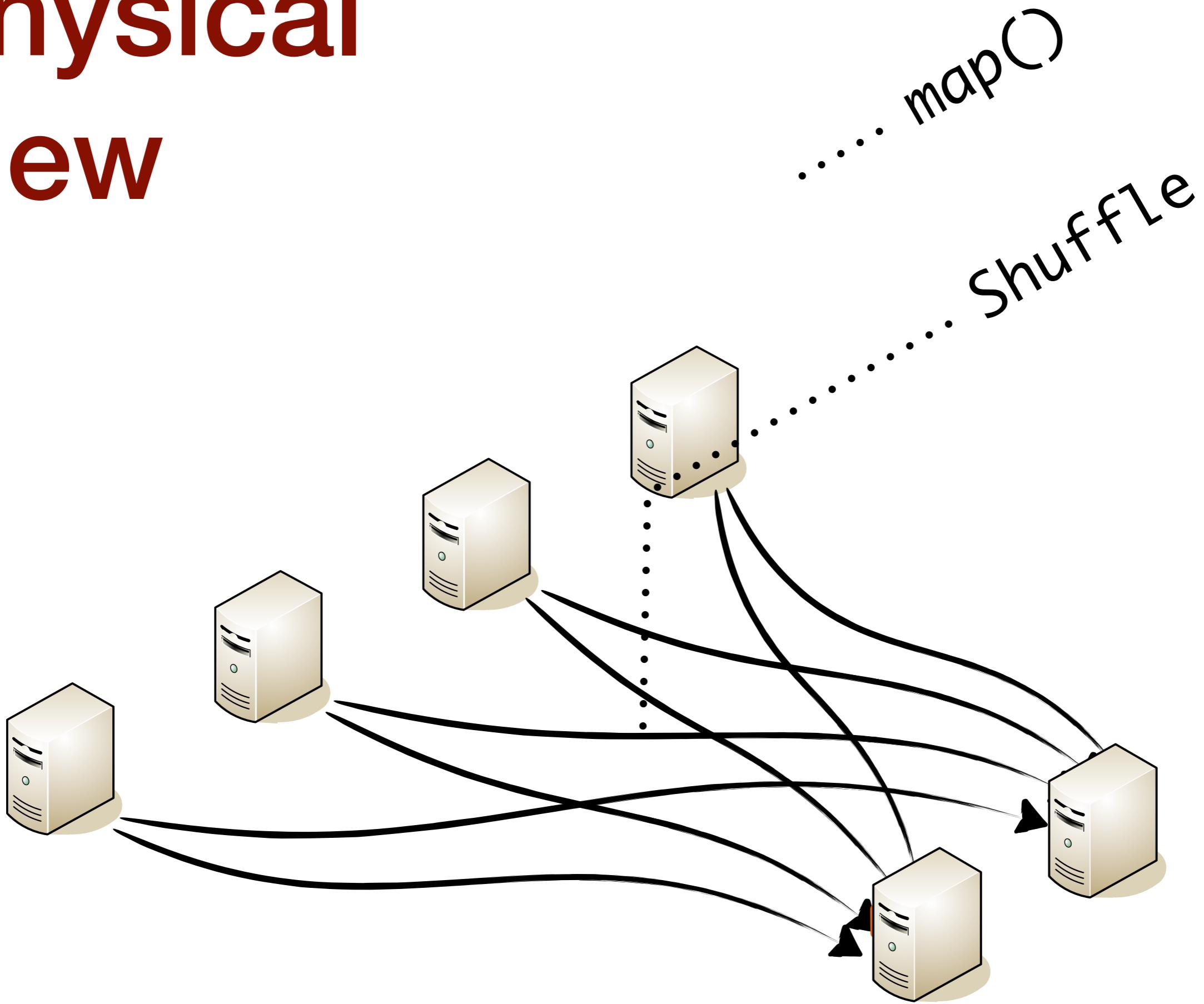
# Physical View



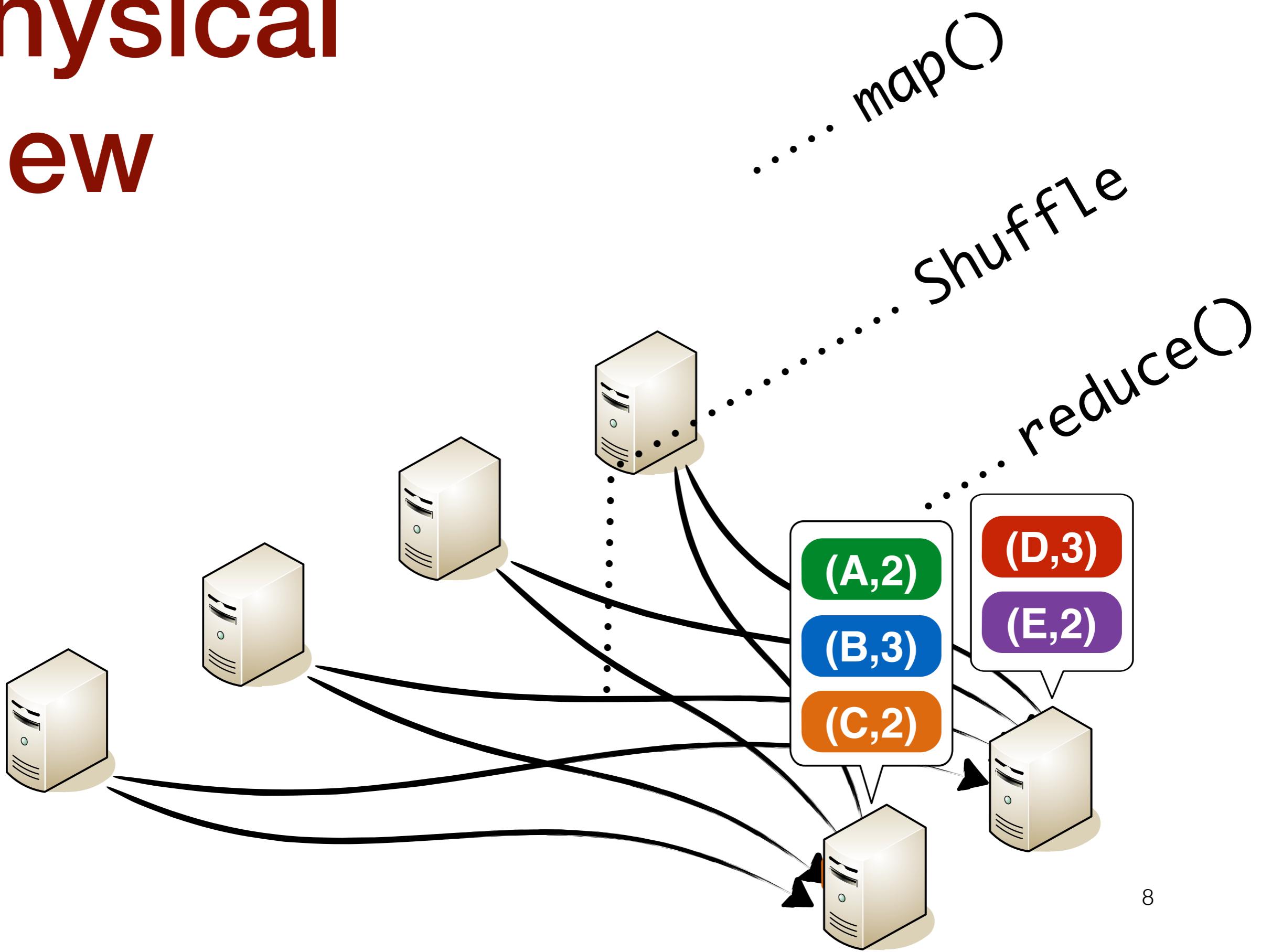
# Physical View



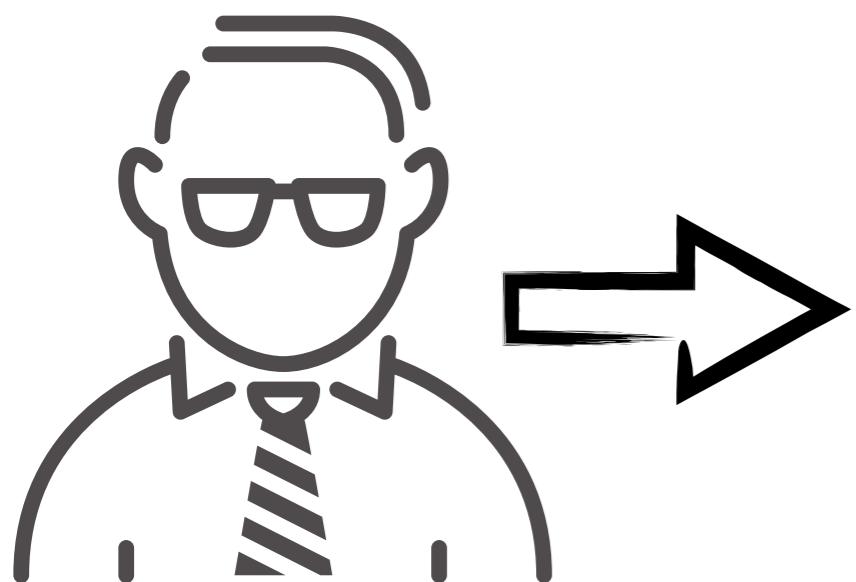
# Physical View



# Physical View



# User

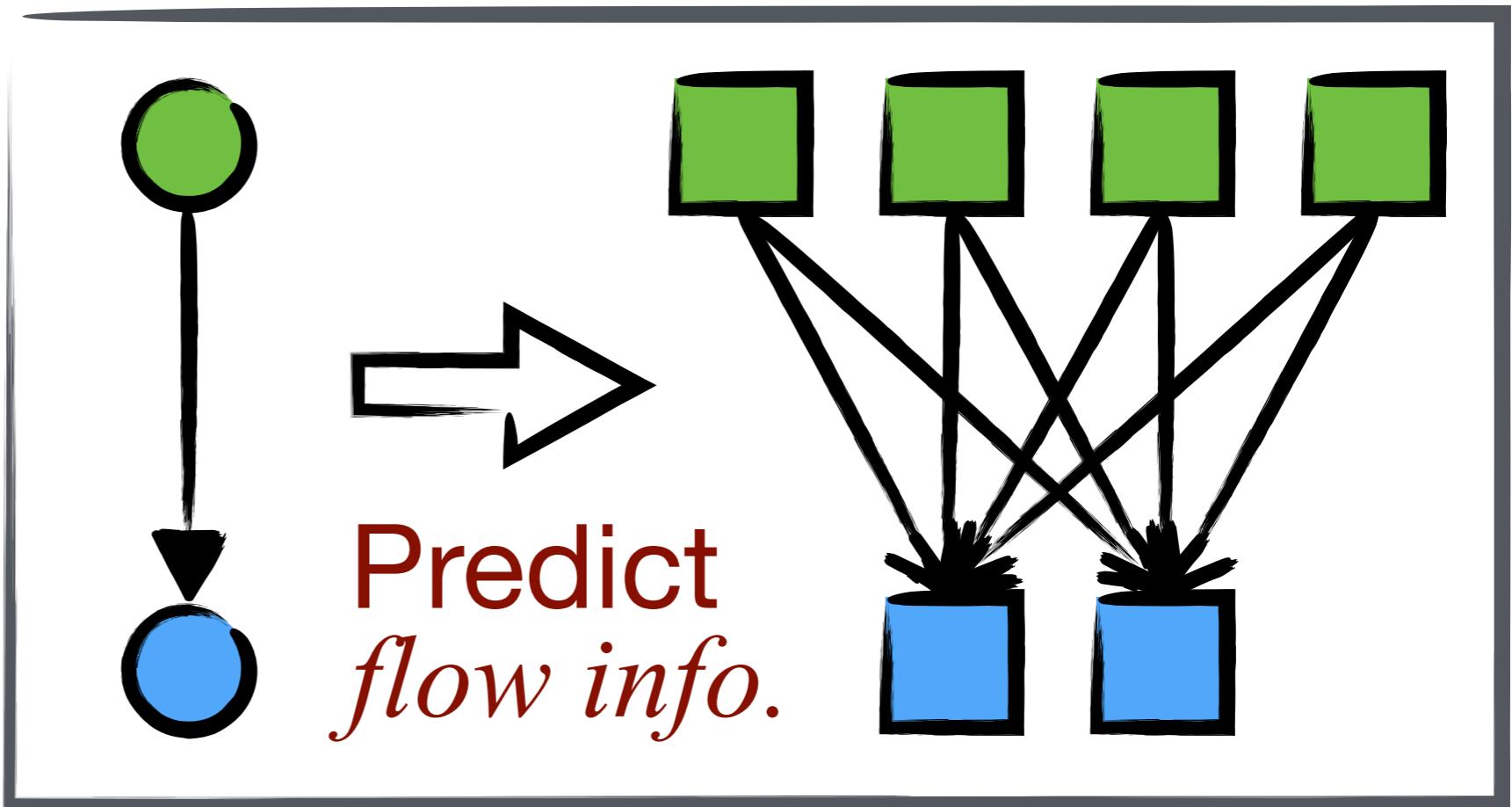
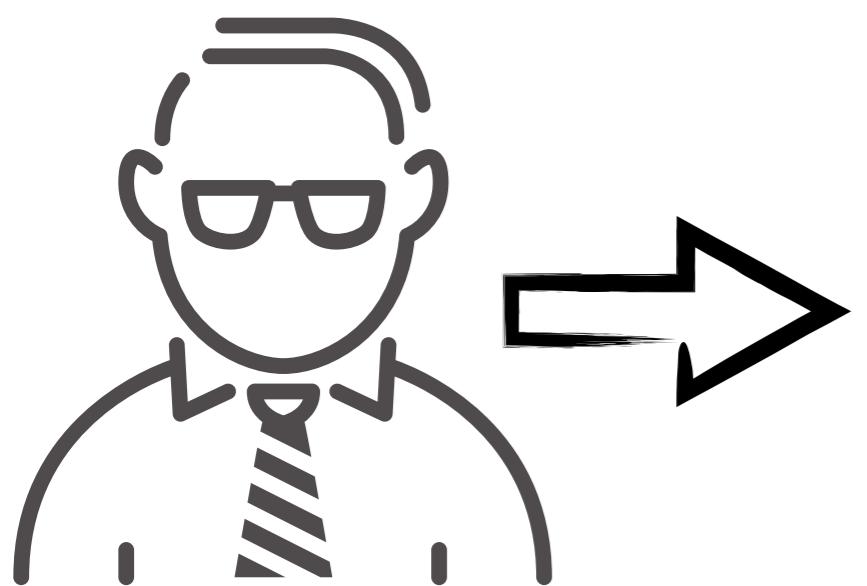


Distributed Computing  
Frameworks

# User

# Logical View

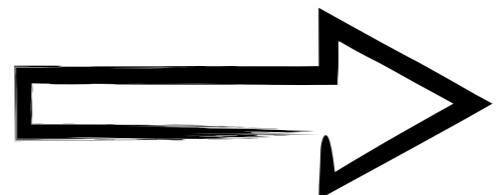
# Physical View



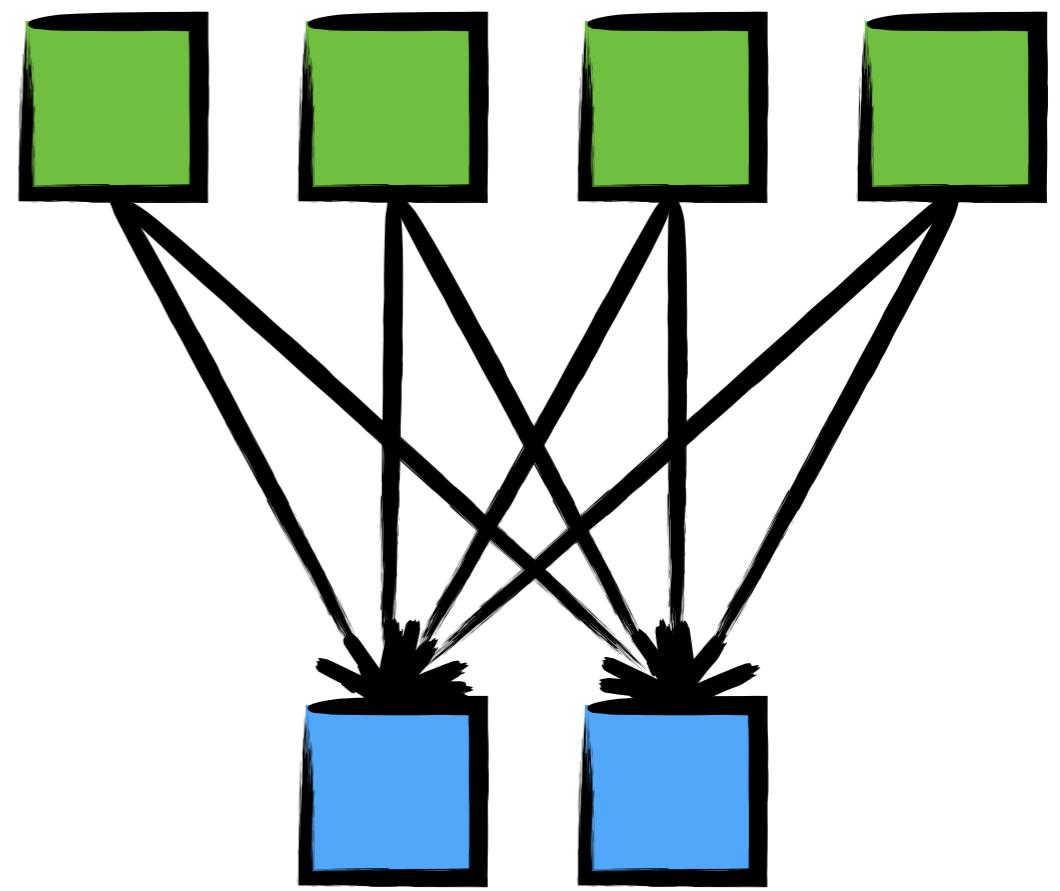
# Logical View



Predict



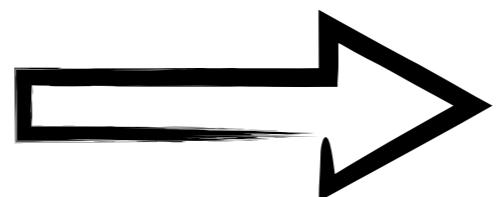
# Physical View



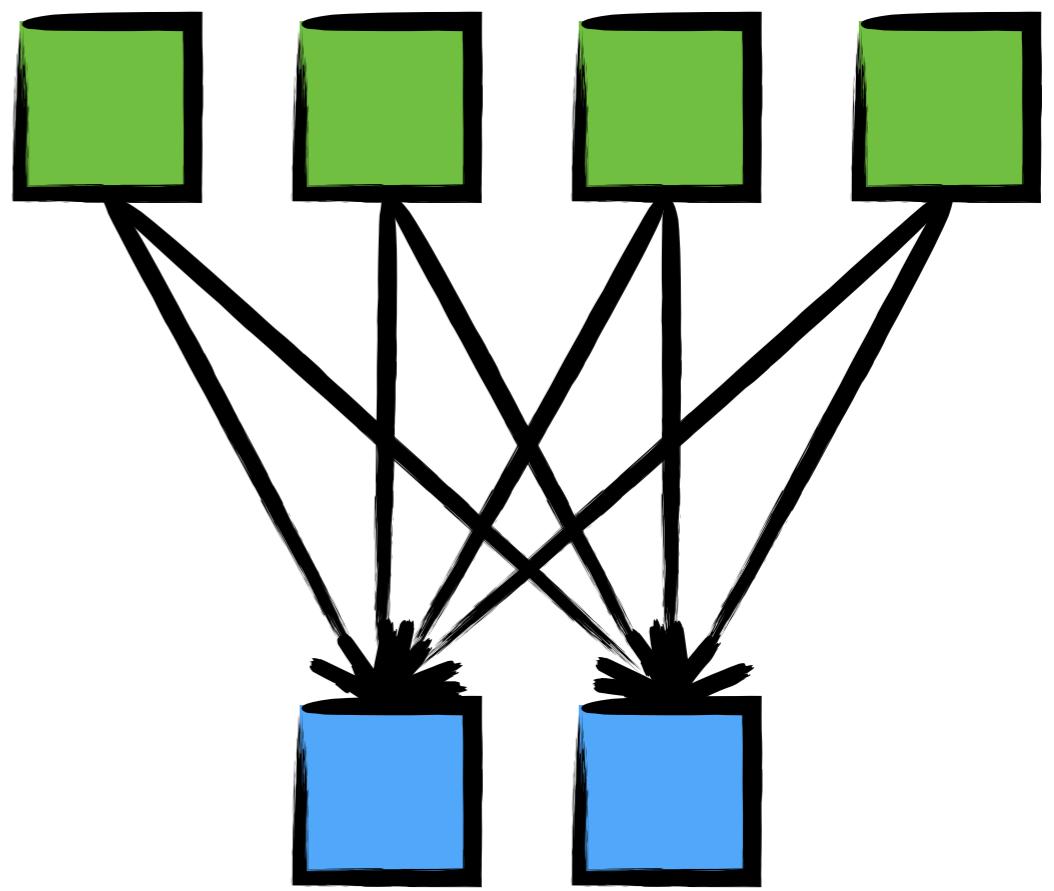
# Logical View



Predict

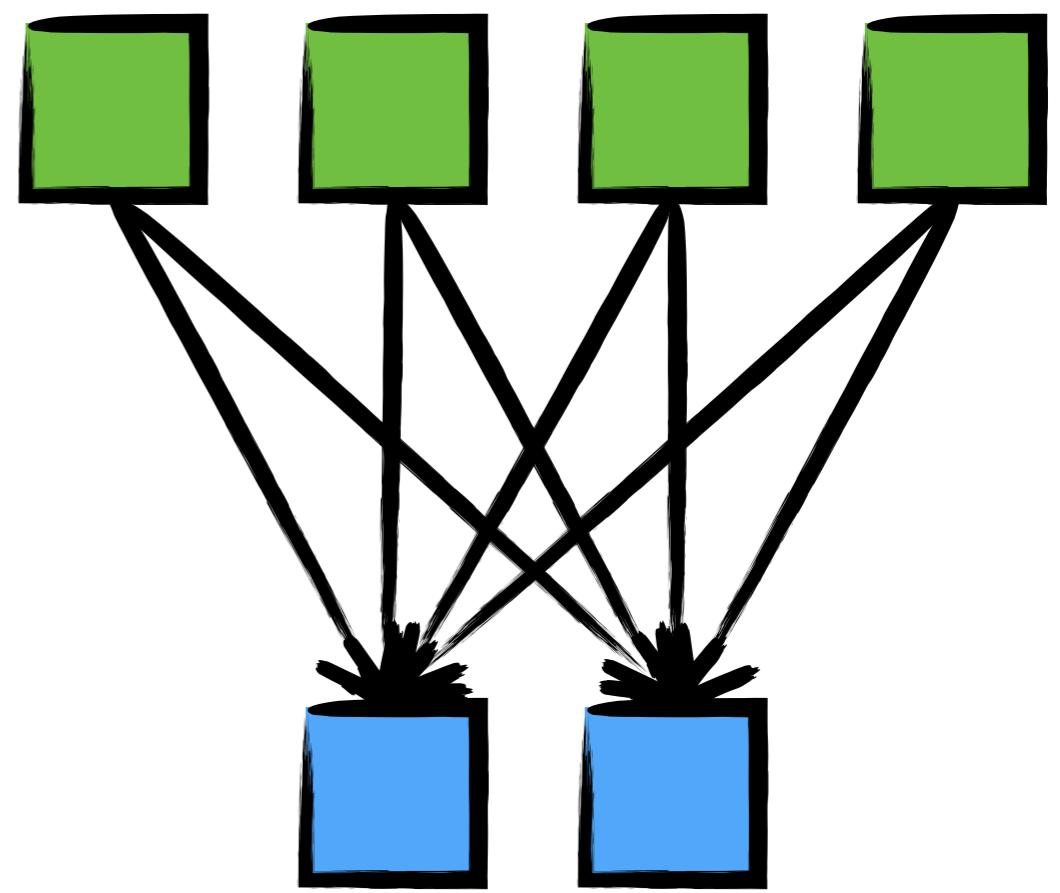
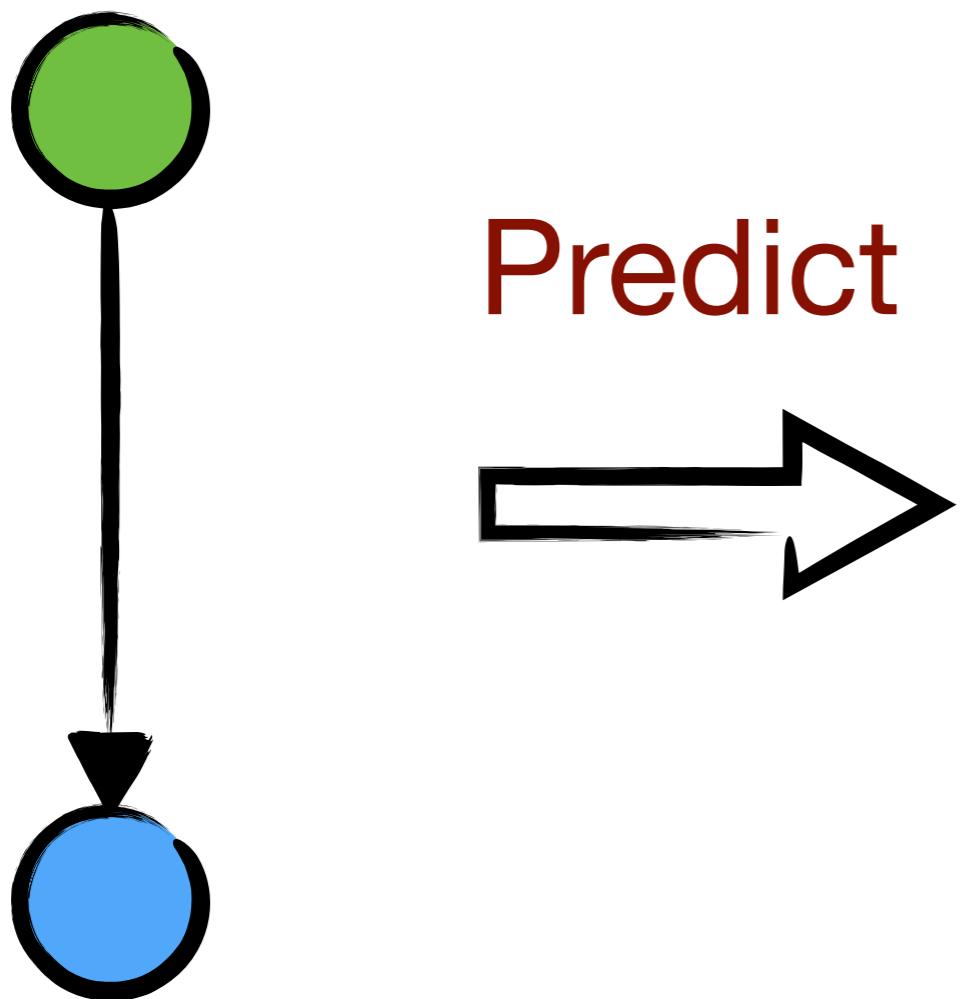


*Flow info.*

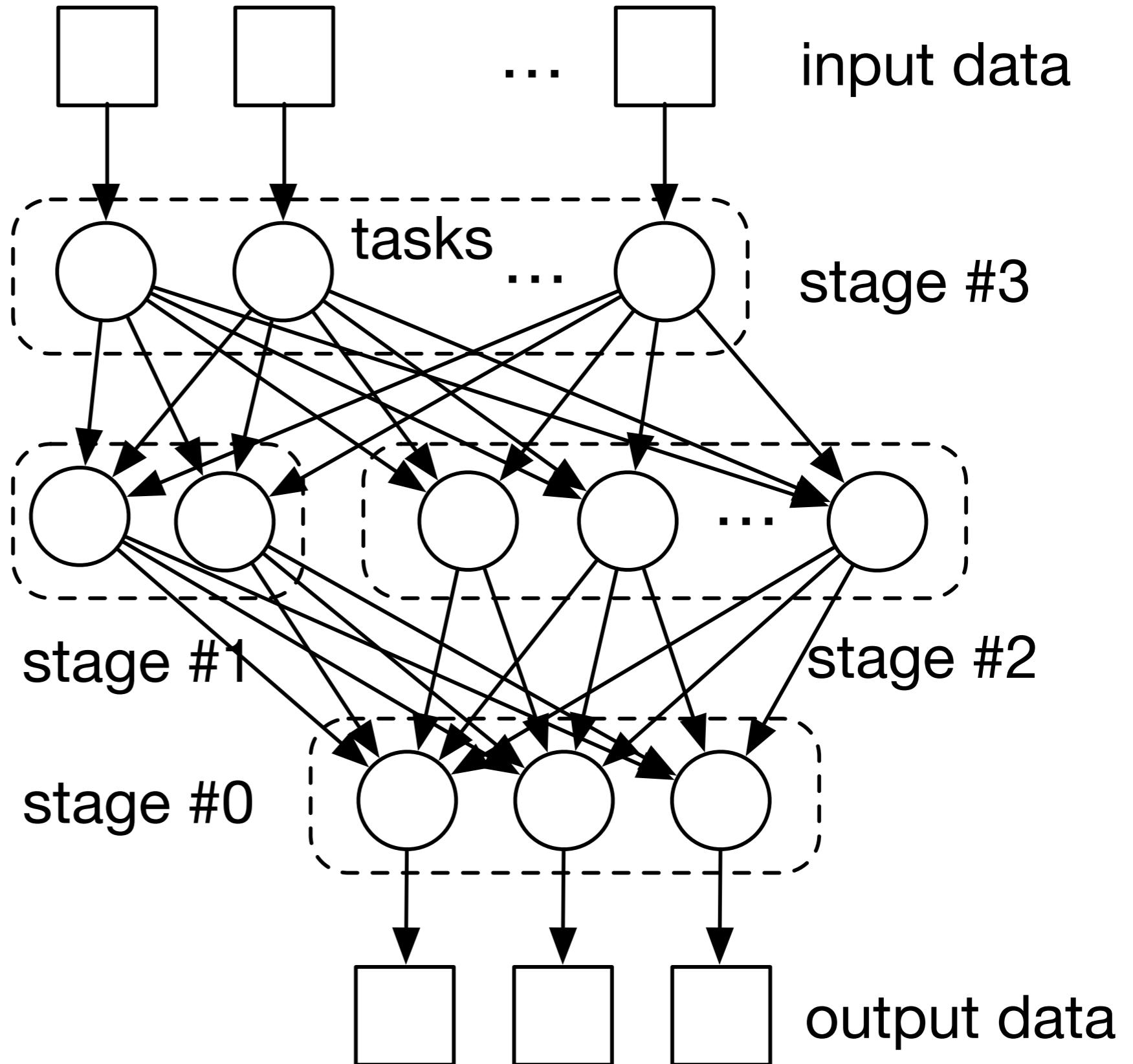


# DAG

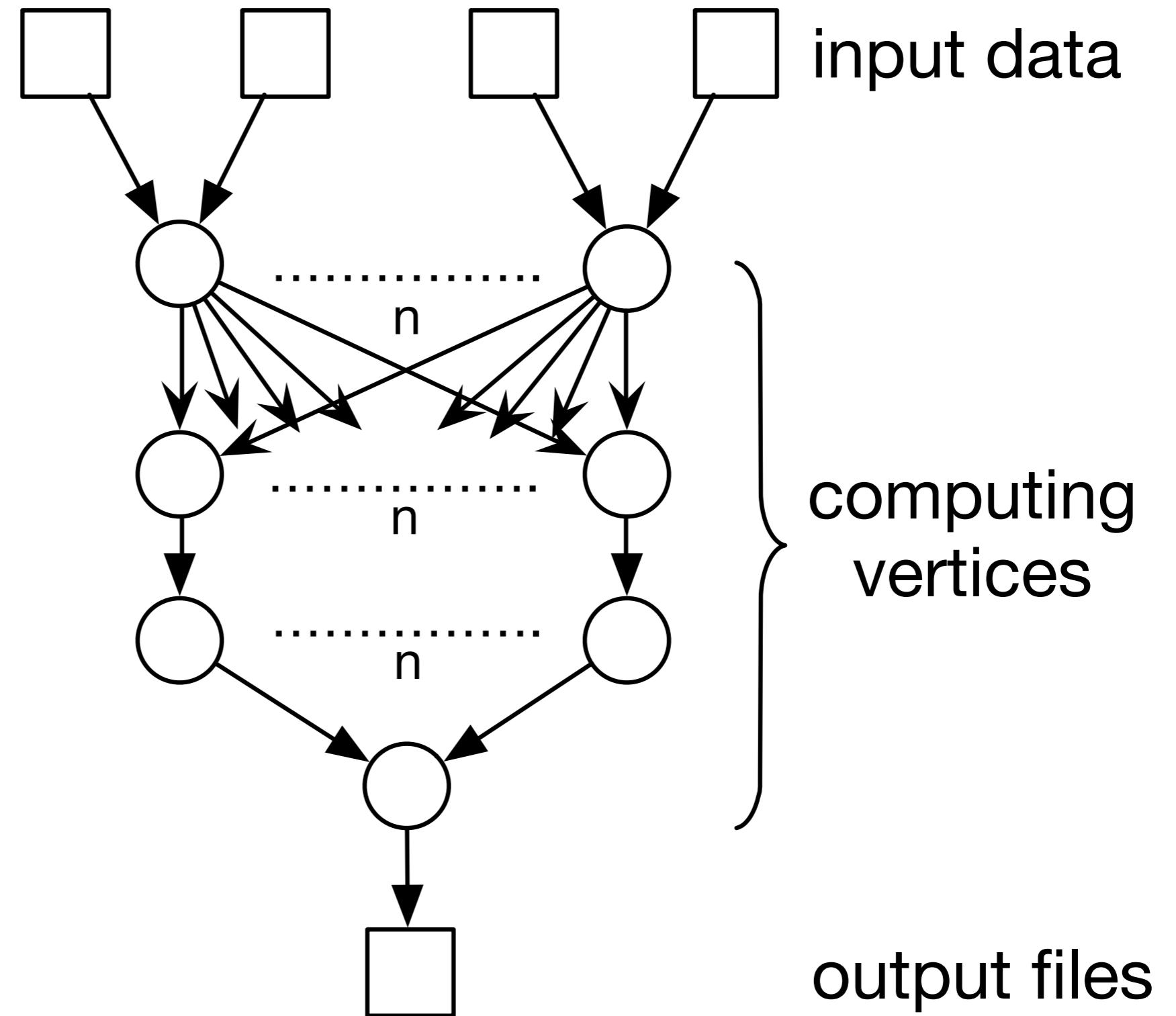
*Flow info.*

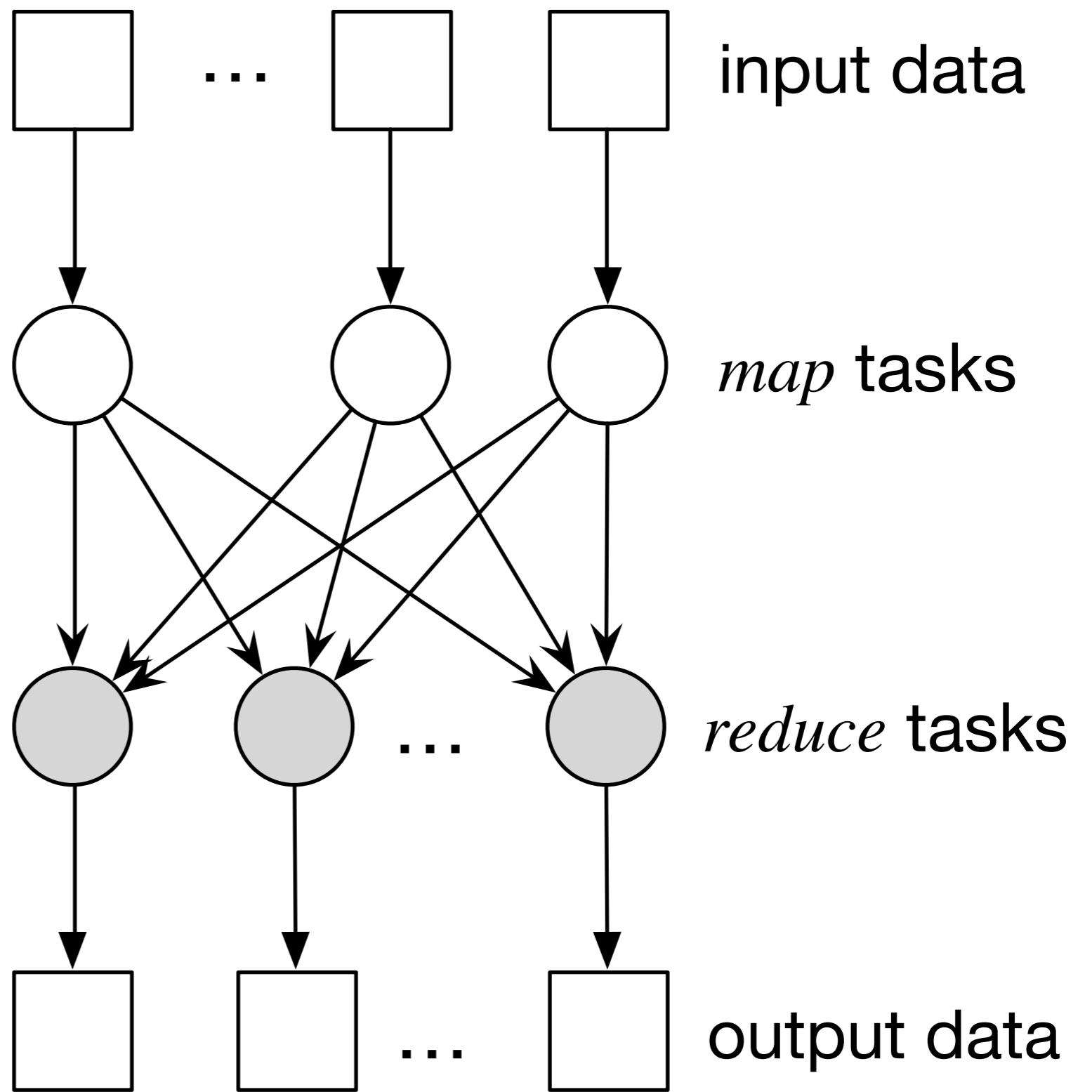


# **Directed Acyclic Graph (DAG)**

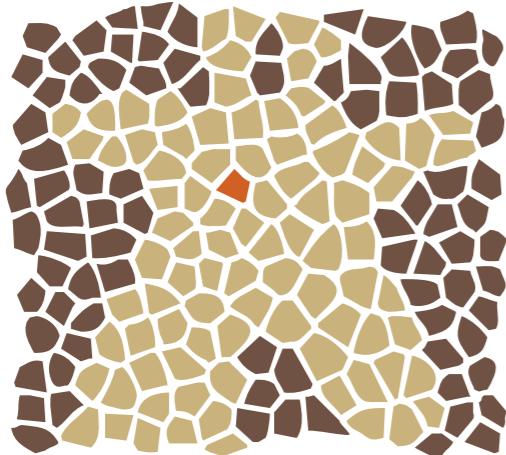


Dryad  
Microsoft

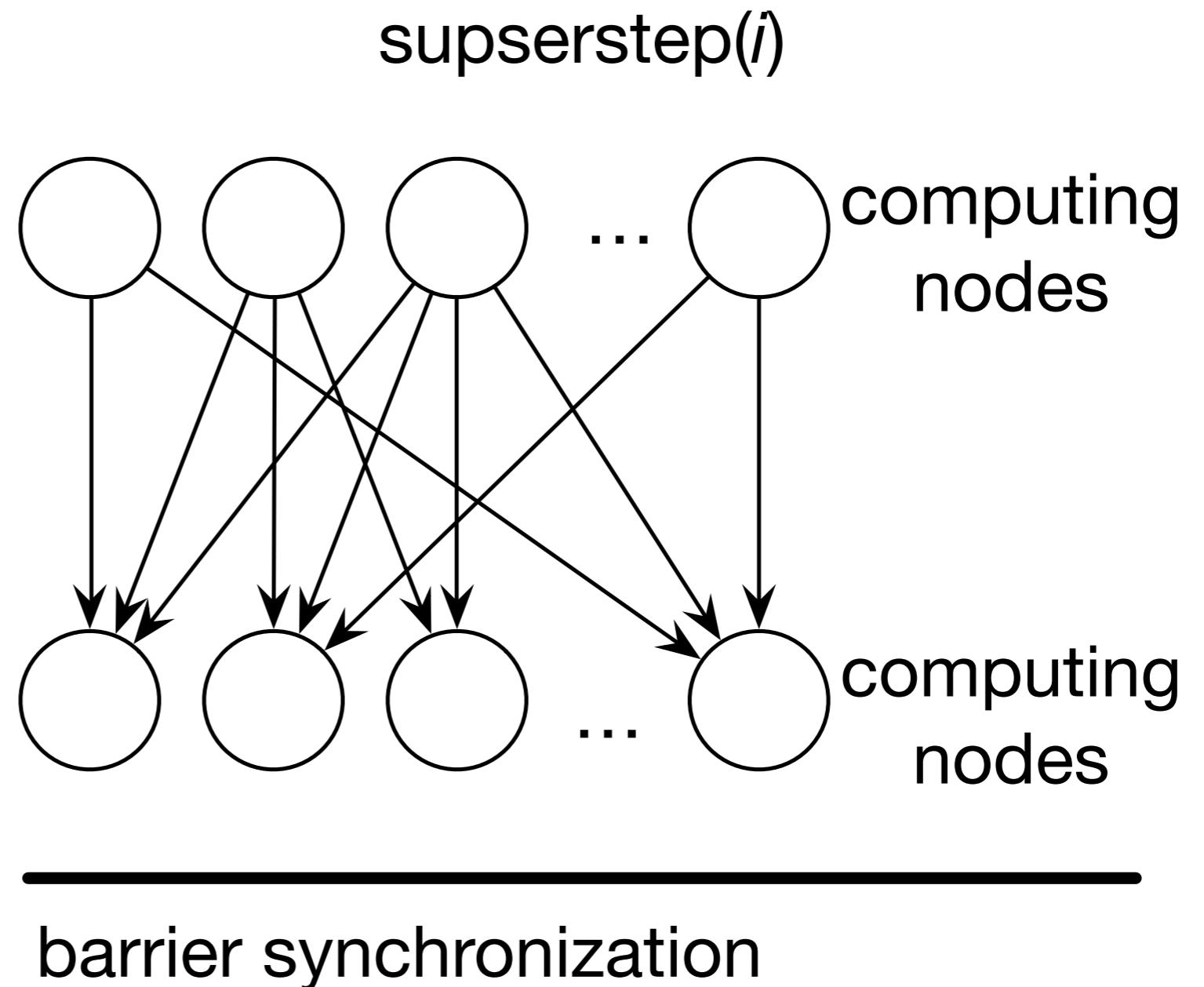


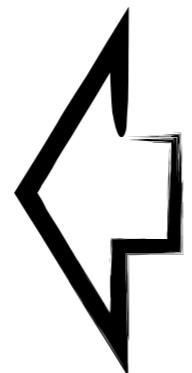
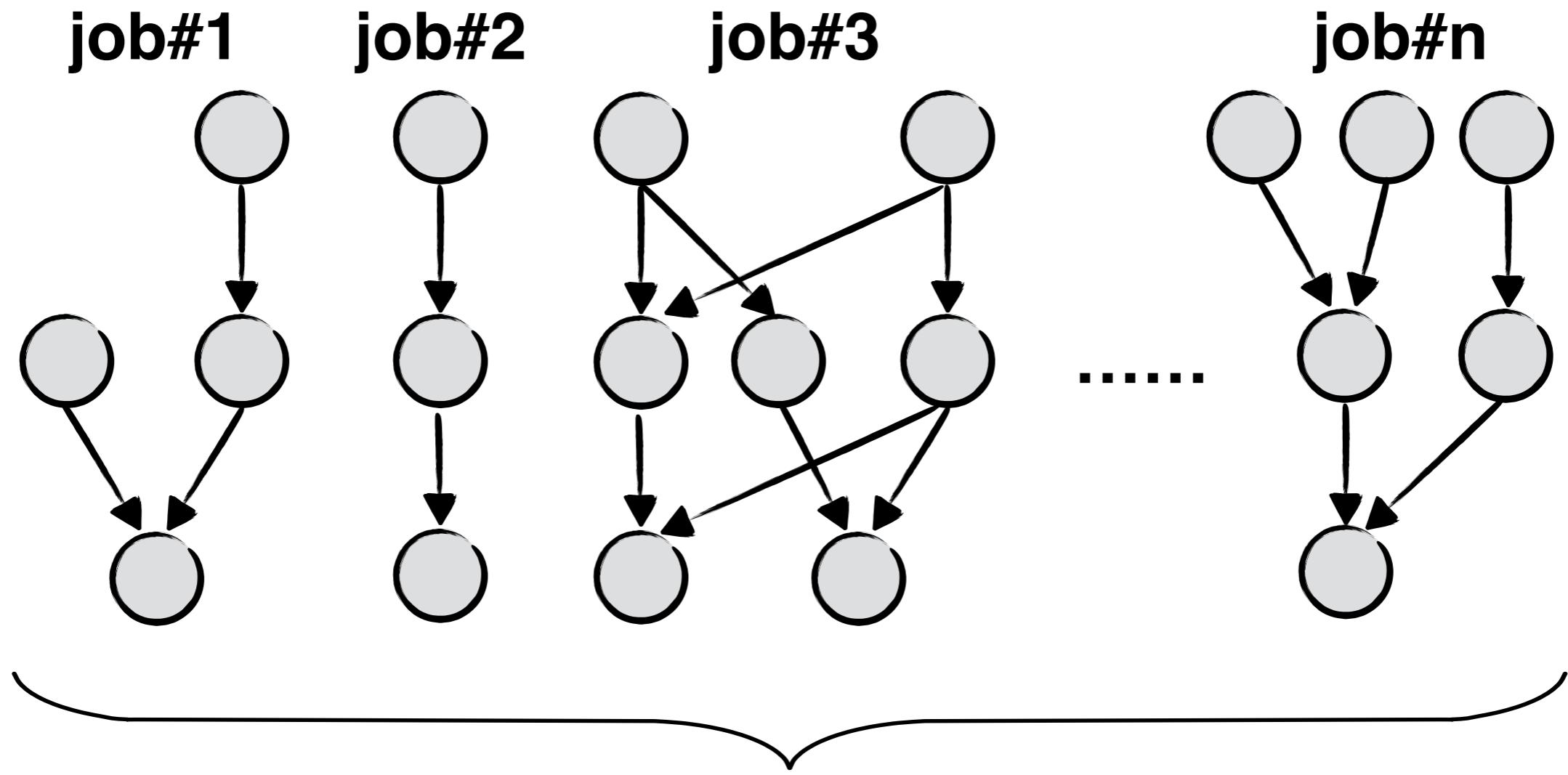


(BSP Model)

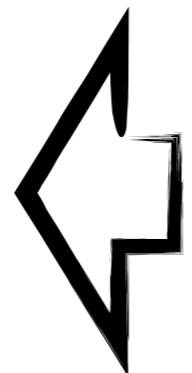
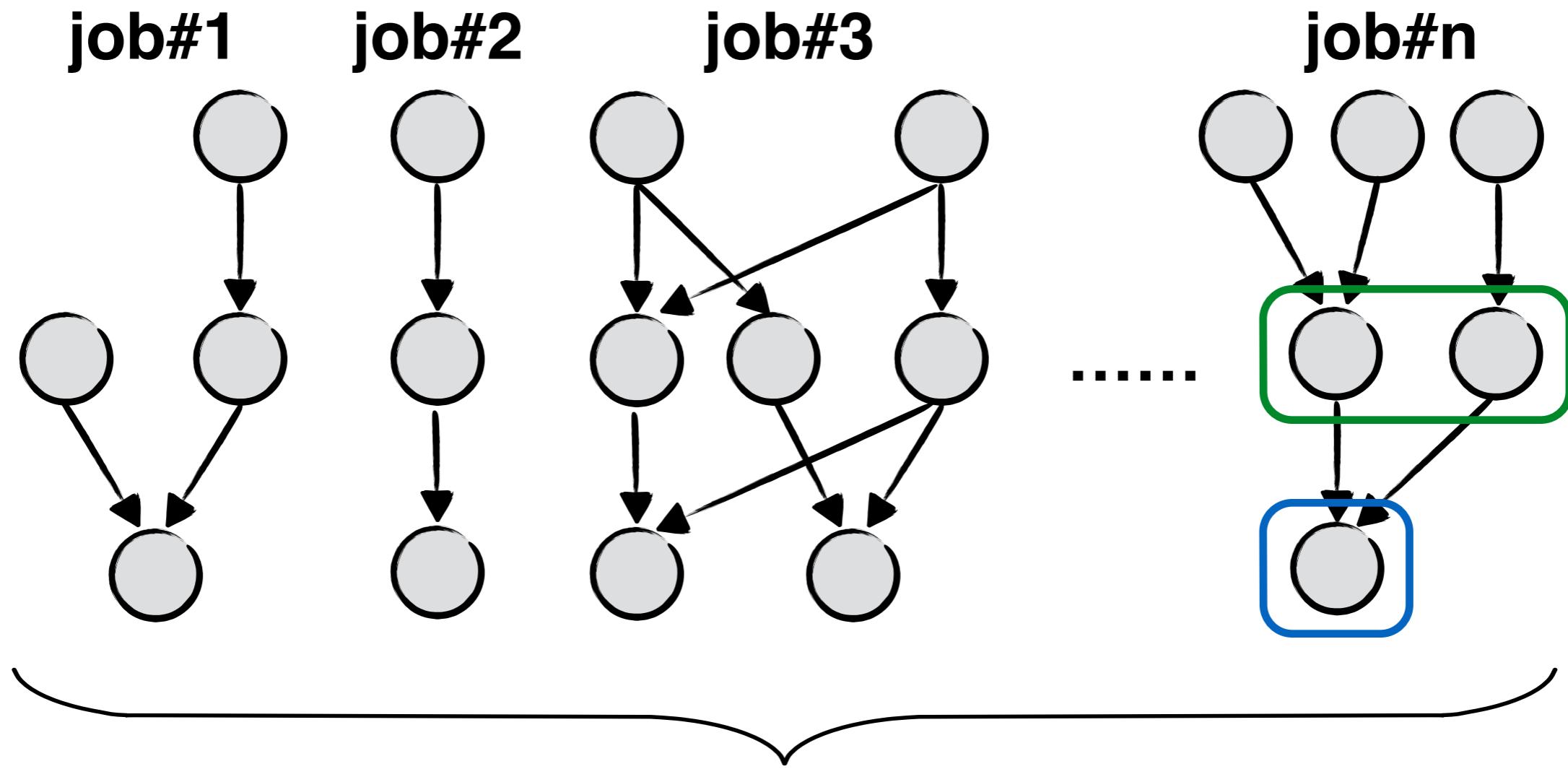


A P A C H E  
G I R A P H

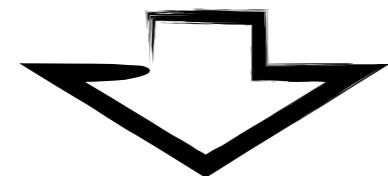
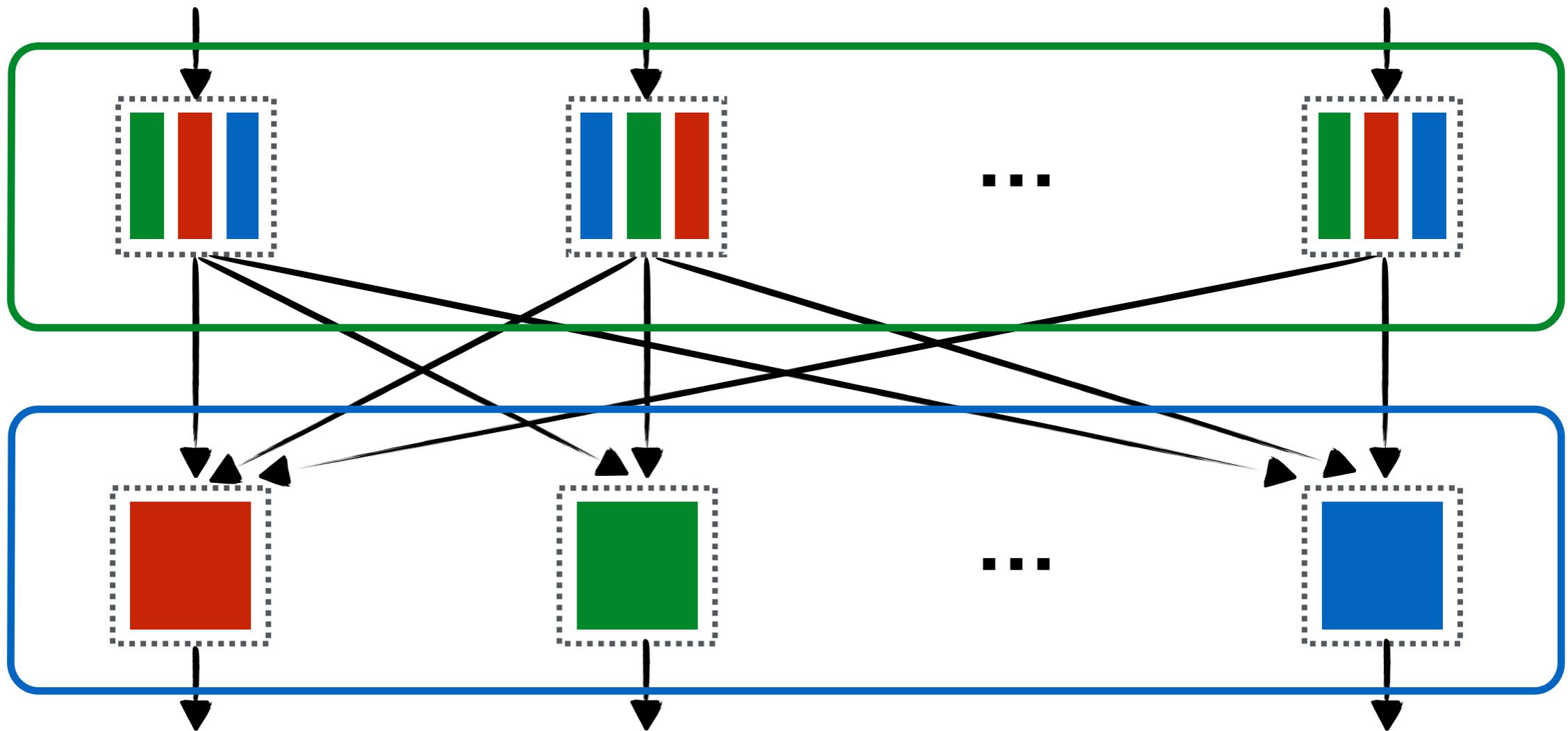




**Application  
Submit**



**Application  
Submit**



# Task Assignment

Worker#1



Worker#2

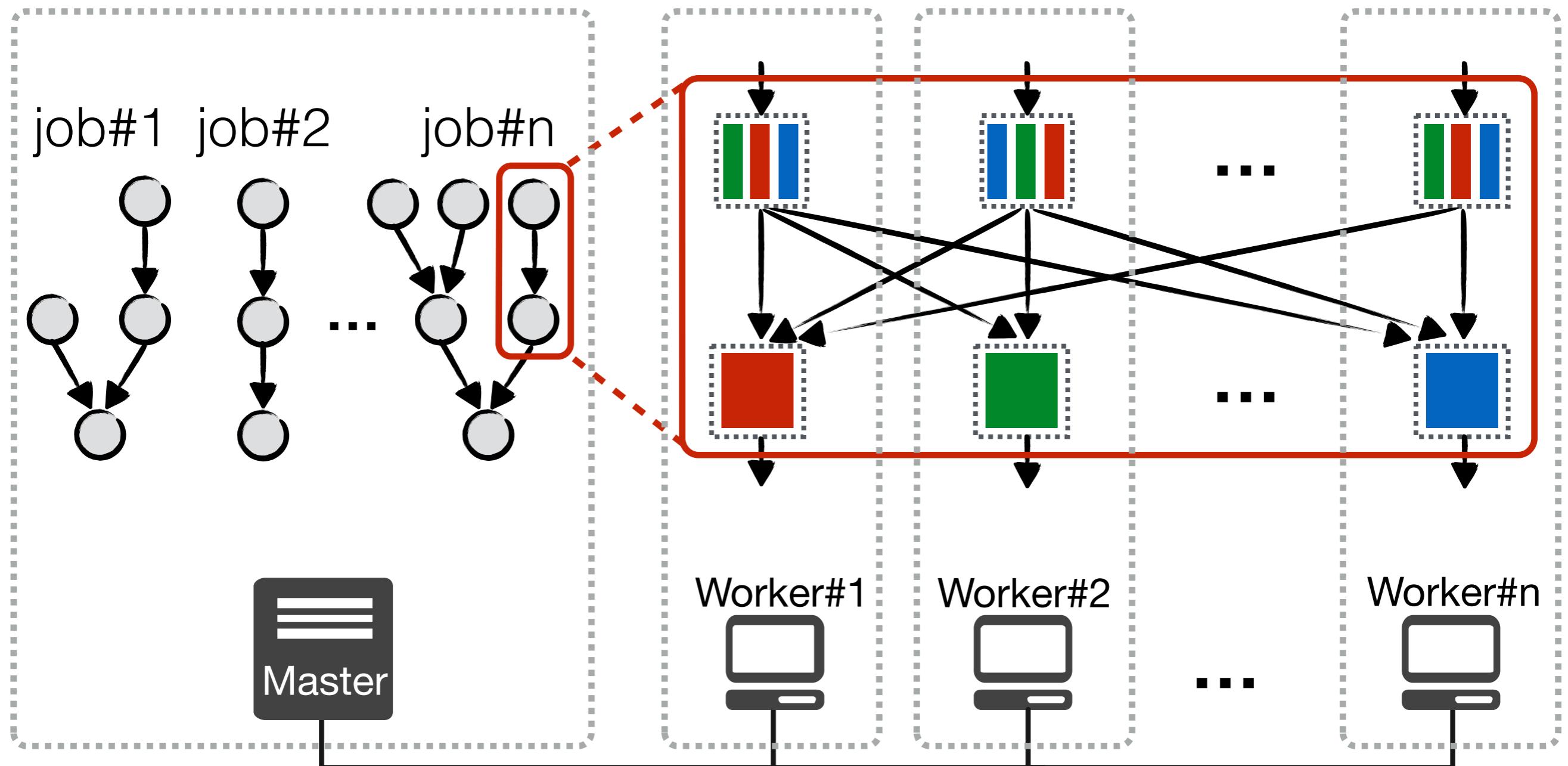


...

Worker#n



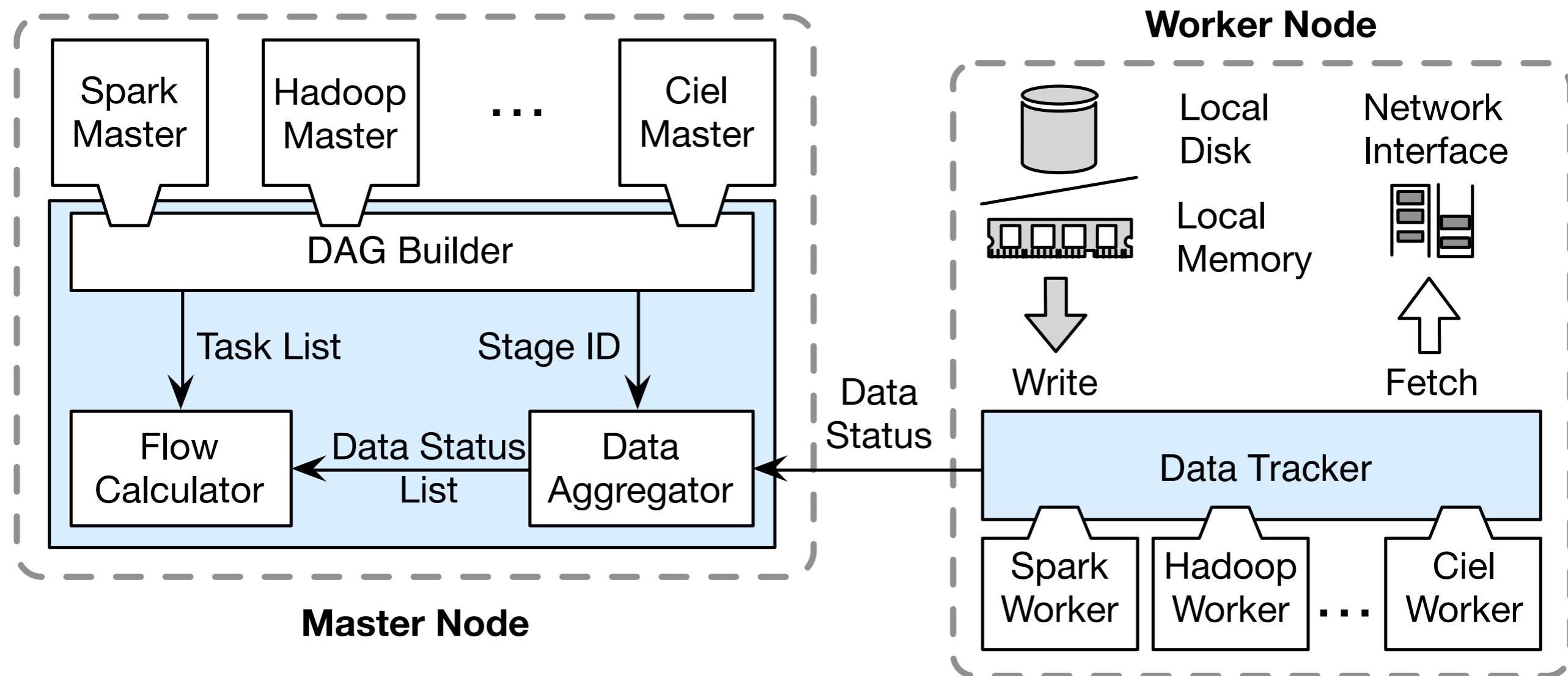
# LIFE CYCLE



# OBSERVATION

— *DAG contains necessary time, data, and flow dependencies for accurate flow prediction.*

# ARCHITECTURE



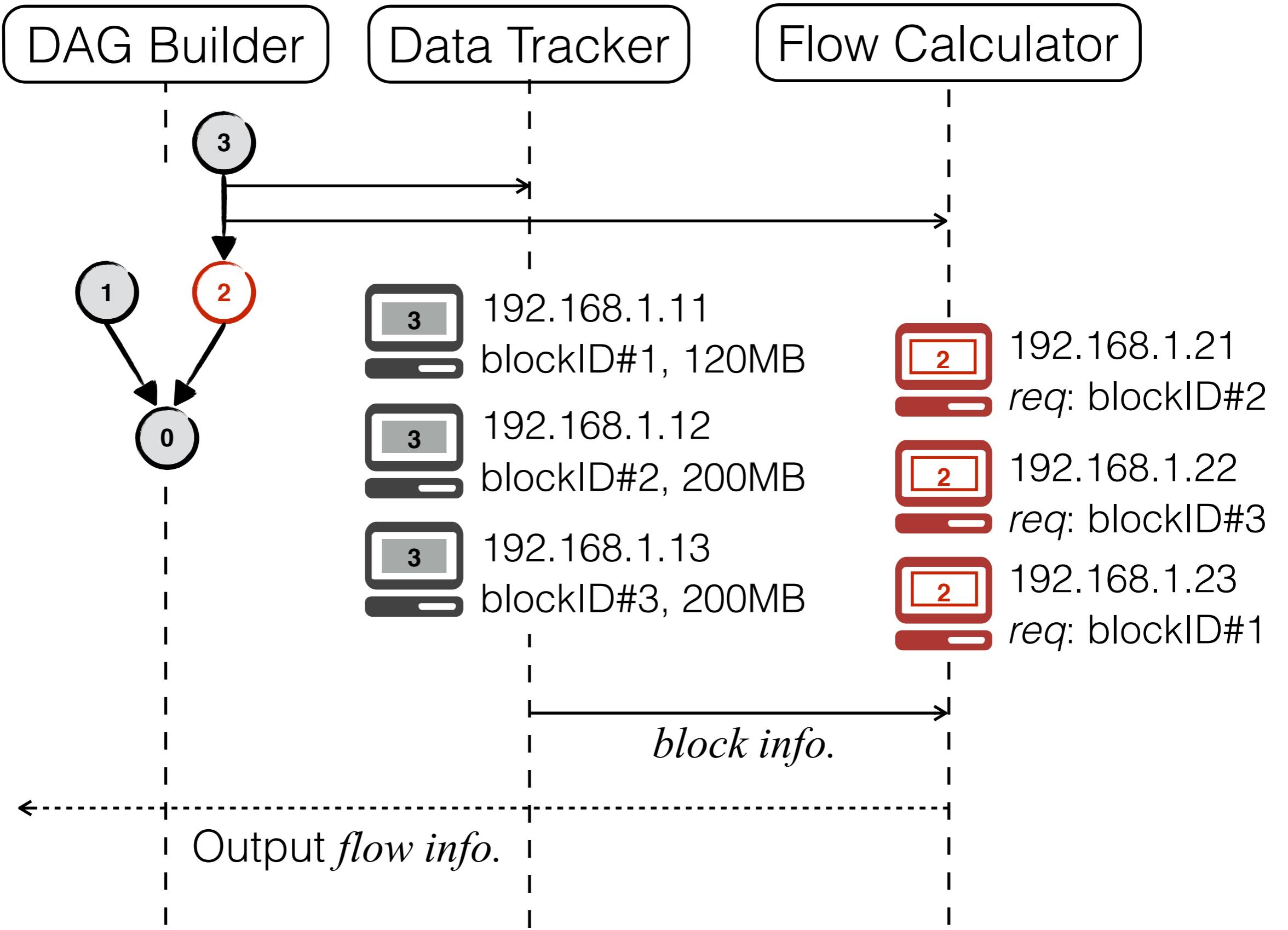
# API EXAMPLES

- Required APIs for DCF master

Event Definition	Trigger Condition
newStageEvent(stageID, childStageID)	a new stage is created
stageStartEvent(List[task], stageID)	a stage is beginning
stageFinishedEvent(stageID)	a stage is finished

- The DAG Builder event handlers

Event Definition
newStageHandler(newStageEvent) ⇒ (currentStage, childStage)
stageStartHandler(stageStartEvent) ⇒ Event(List[task], List[stageID])
stageFinishedHandler(stageFinishedEvent) ⇒ Event(stageID)



# FlowPROPHET

- Generic
- Accurate and fined-grained
- Ahead-of-time
- Scalable and low-overhead

# TESTBED

- Dell PowerEdge R320 x 37
- Intel Xeons E5-1410 2.8GHz CPU
- 24GB 1600MHz DDR3
- Broadcom Gigabit Ethernet NIC
- Pronto-3295 Gigabit Ethernet Switch

# BENCHMARKS

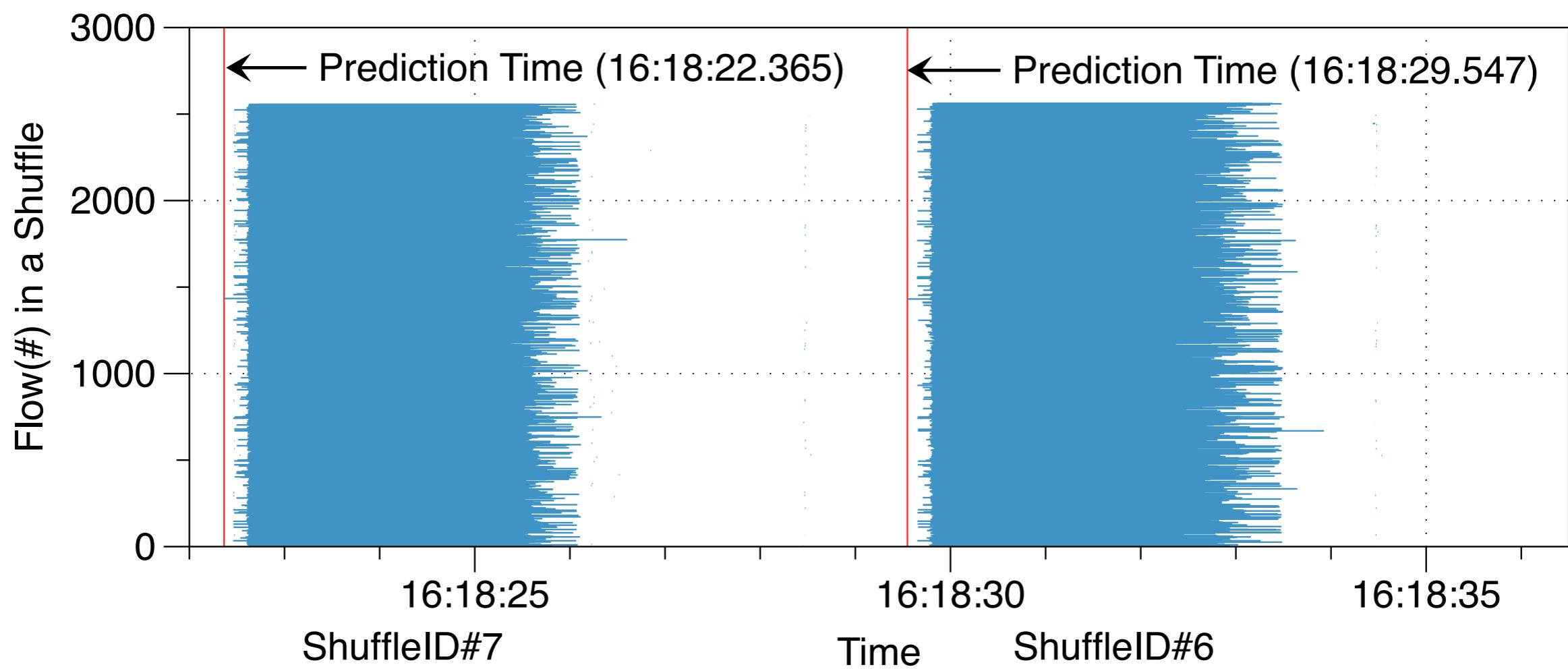
- WikiPageRank
- SparkPageRank
- Spark K-means
- Hadoop TeraSort
- $\pi$  (Pi)
- WordCount

# METRICS

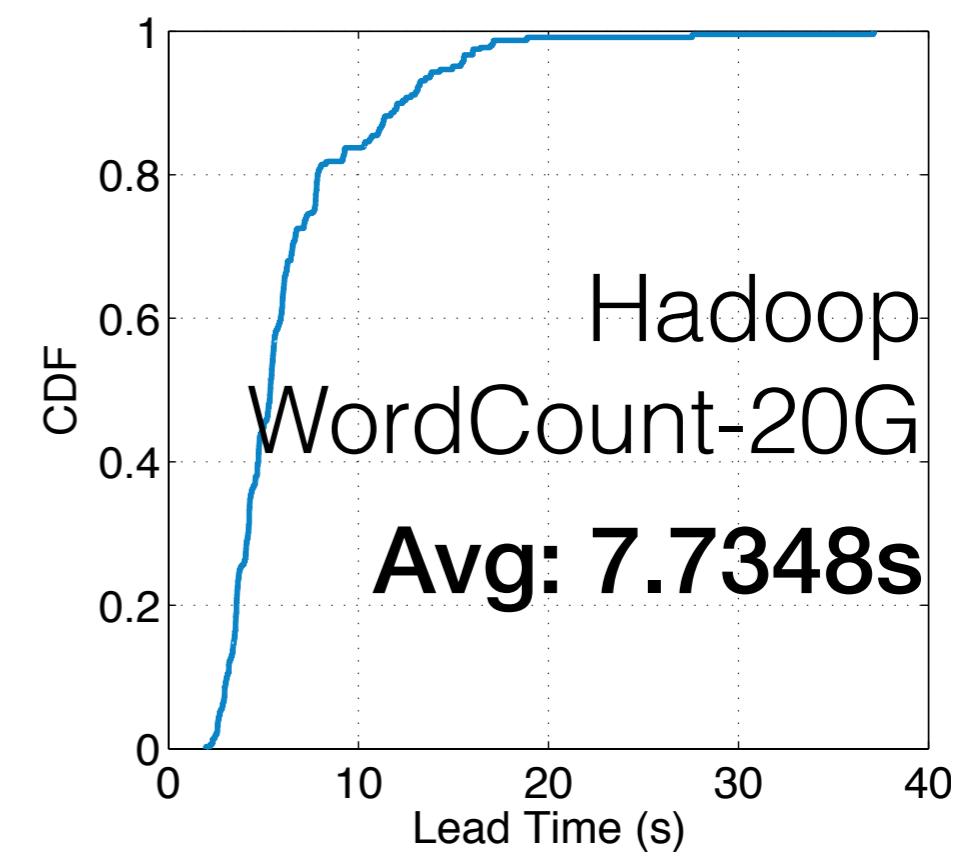
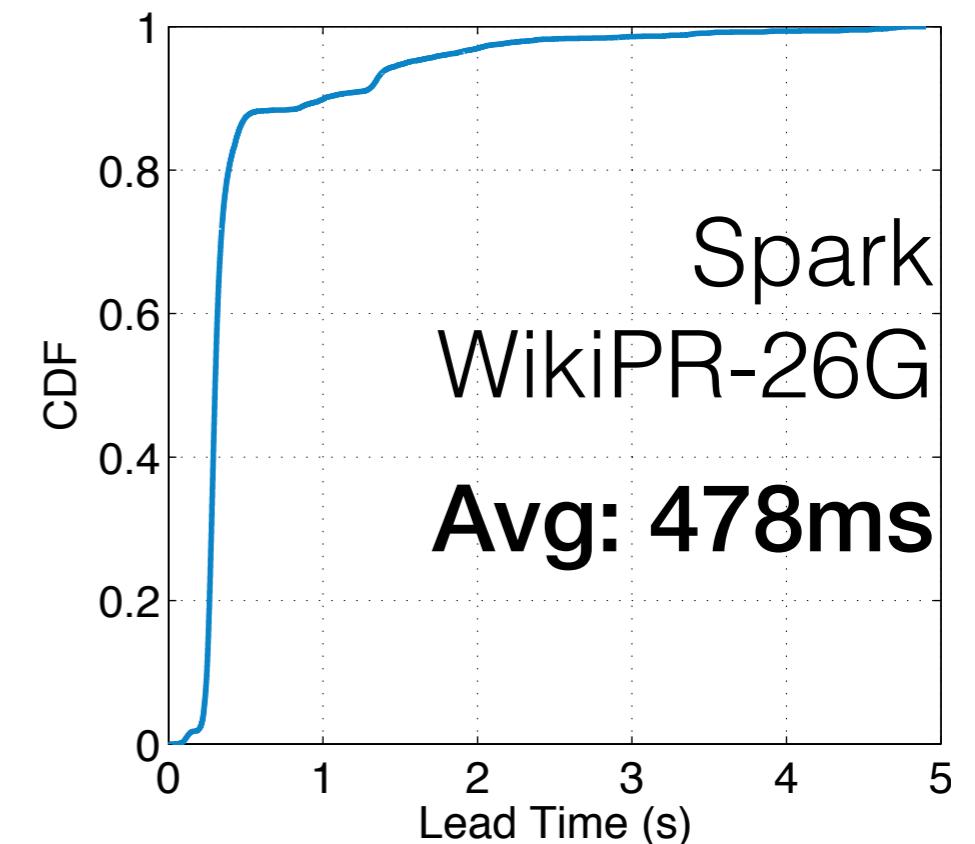
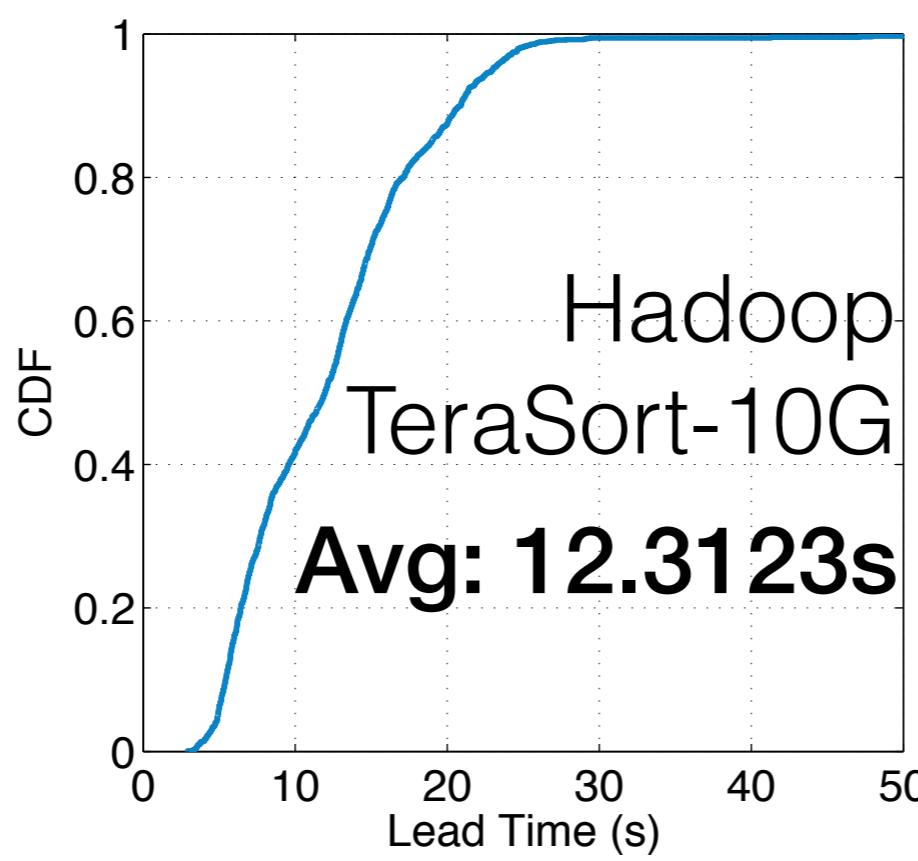
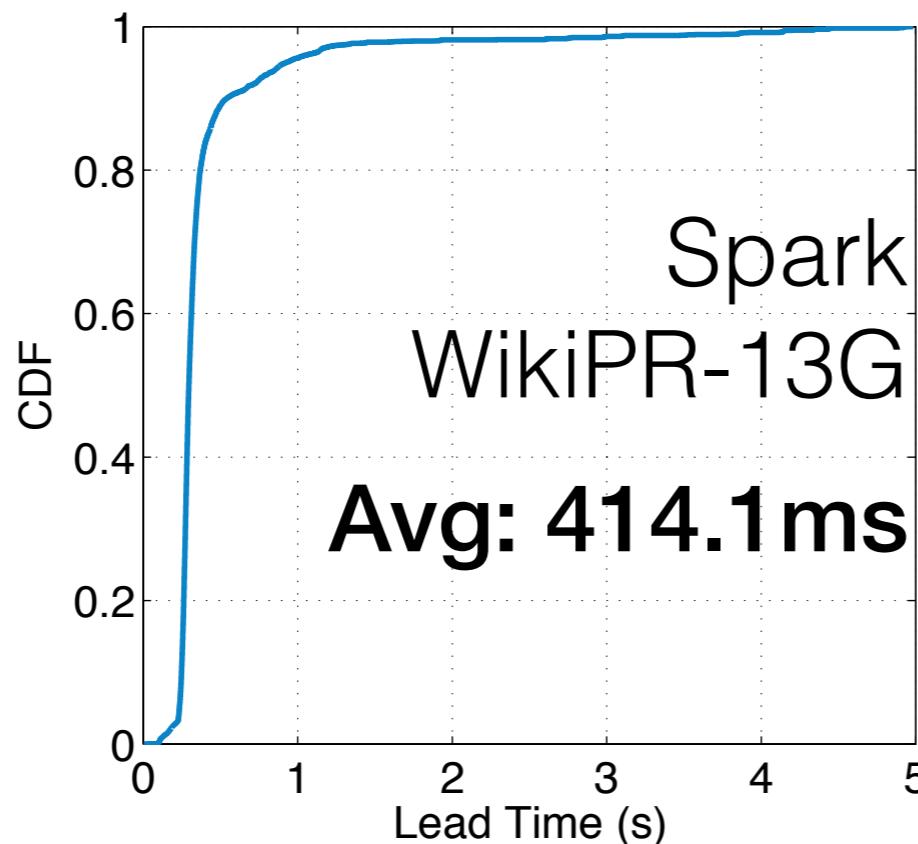
- Time advance
- Prediction accuracy
- Overhead
- Scalability
- Benefits

# TIME ADVANCE

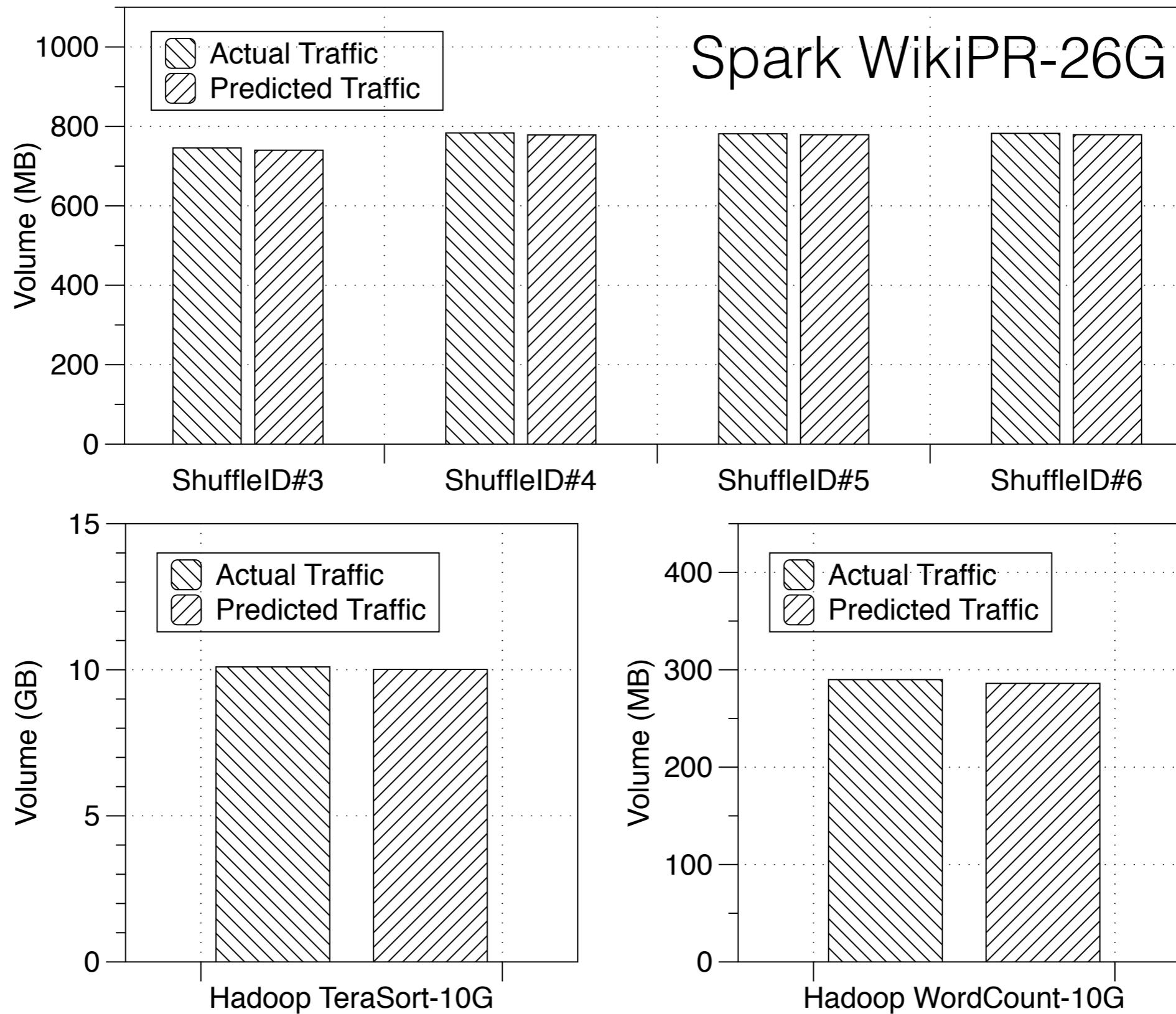
- WikipediaPageRank-13G (Spark)



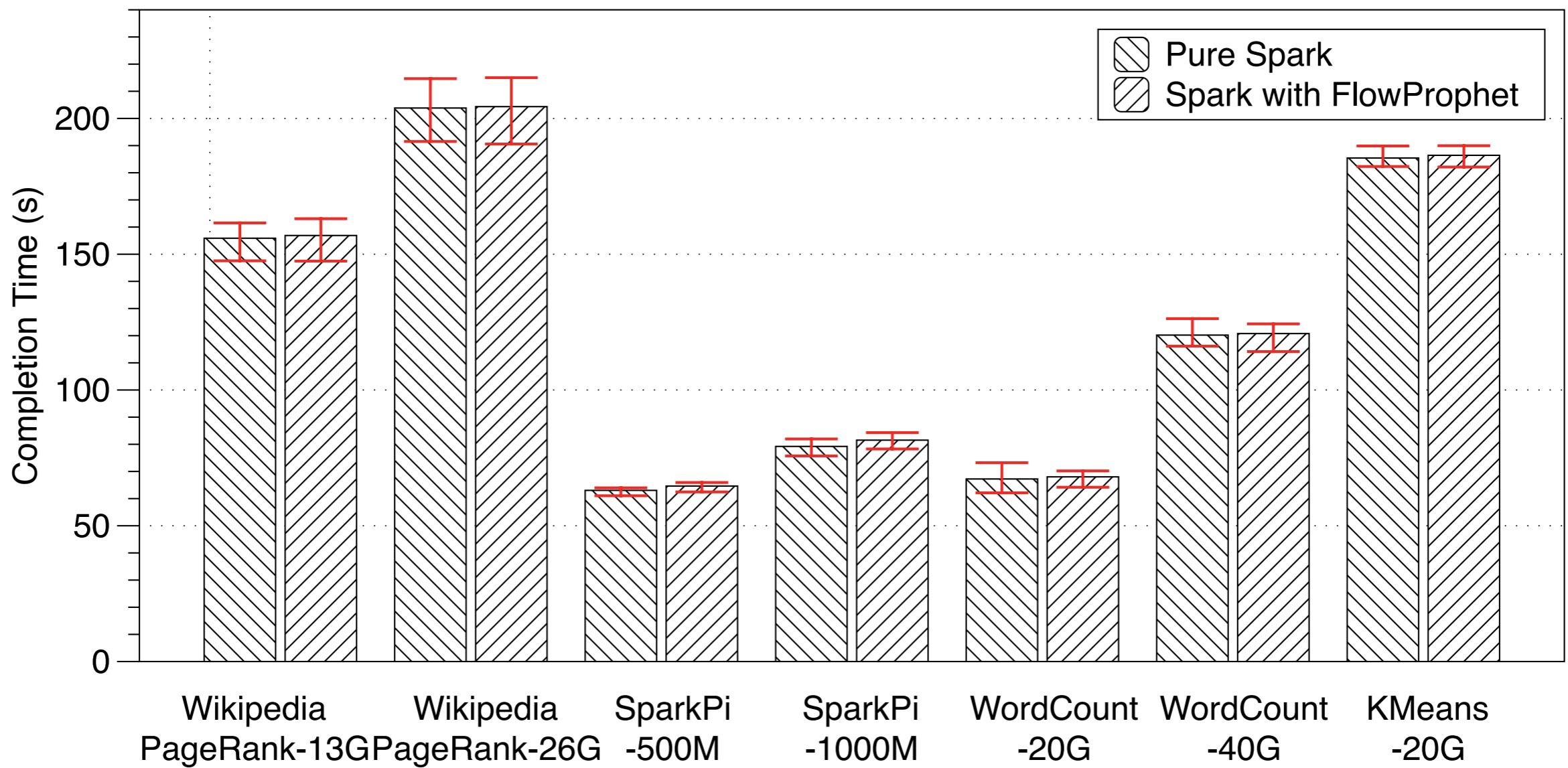
# CDF OF LEAD TIME



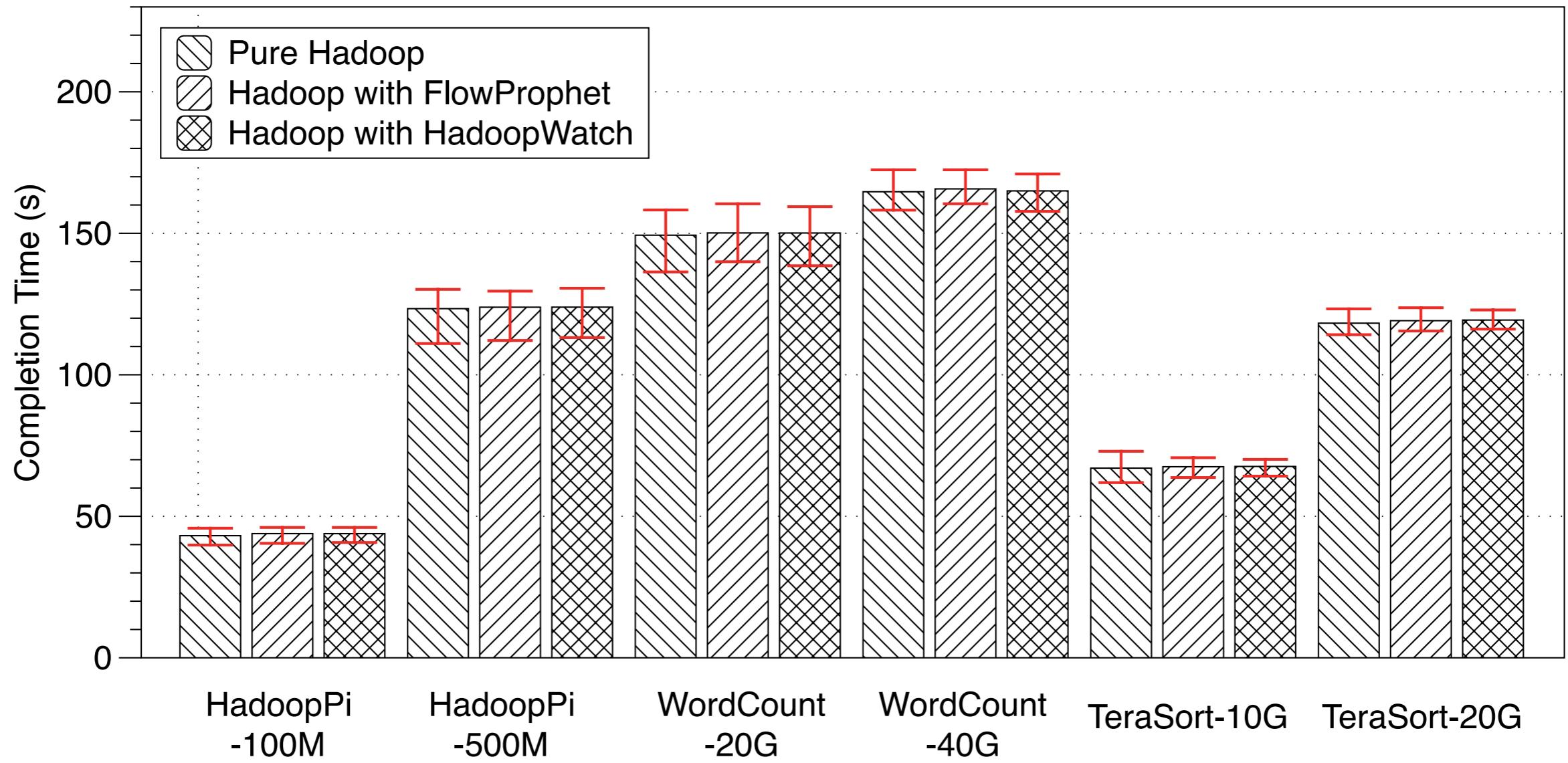
# PREDICTION ACCURACY



# OVERHEAD

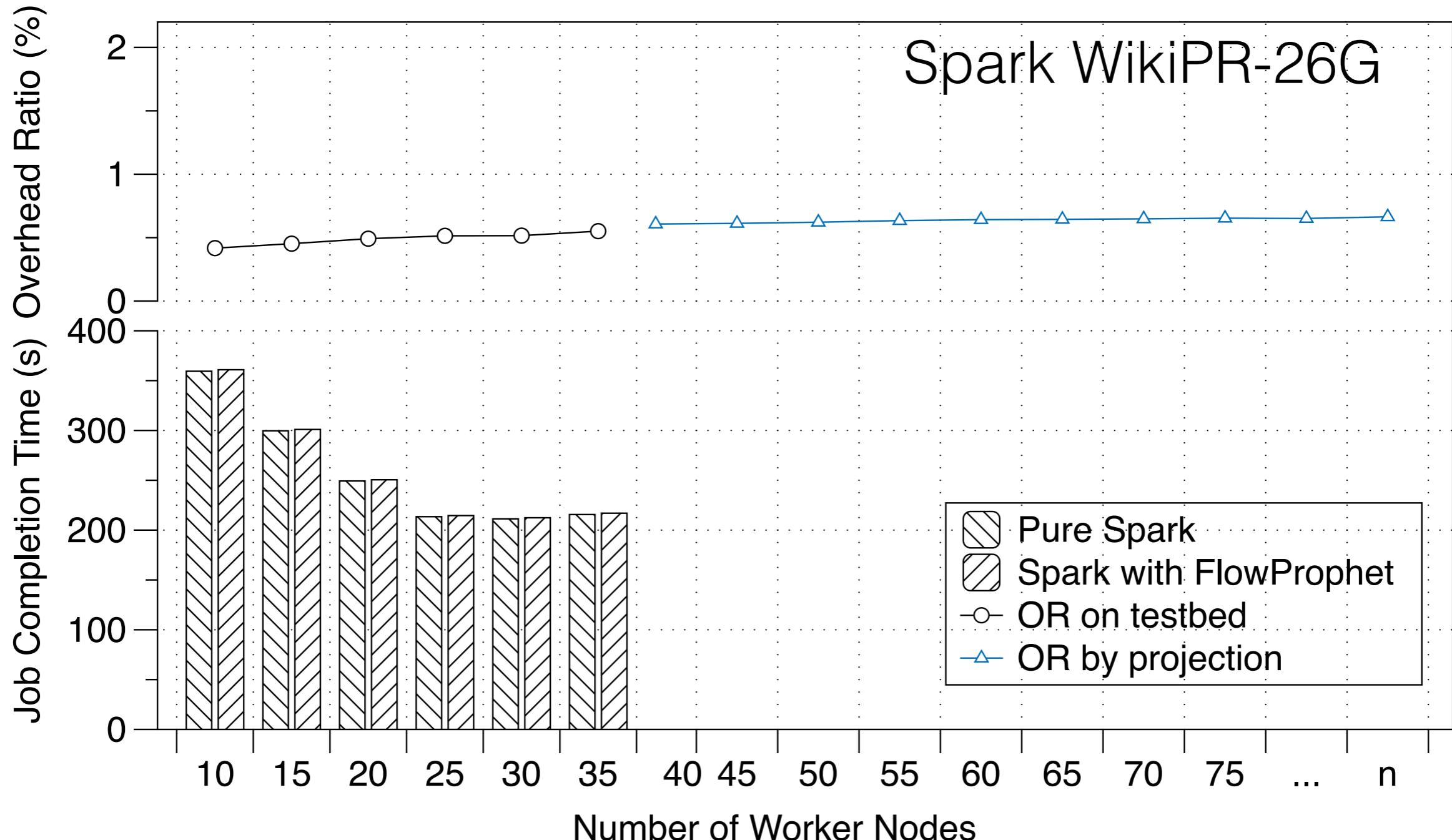


# OVERHEAD

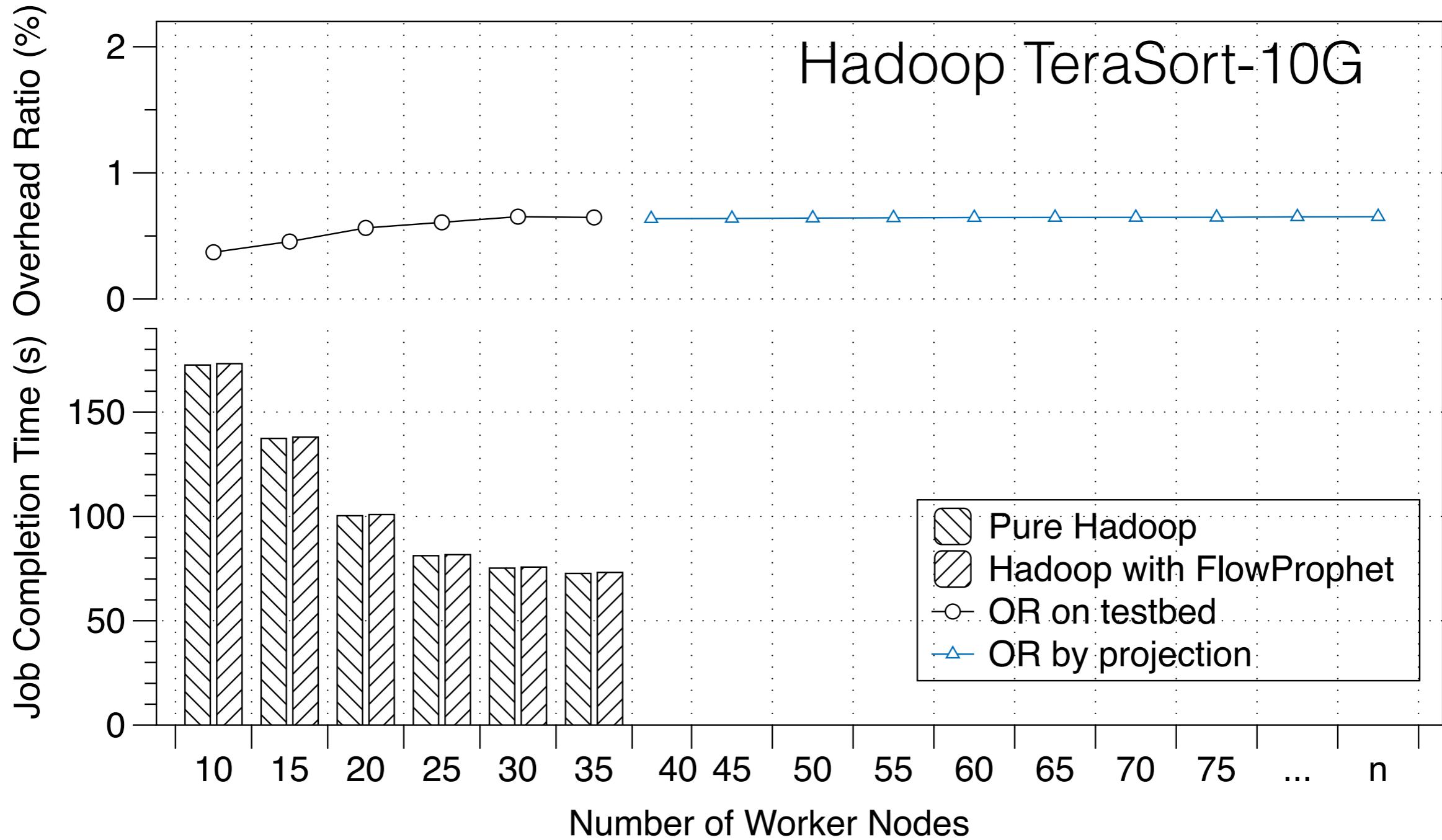


# SCALABILITY

- Overhead Ratio (OR) :  $OR = \frac{t_{enabled} - t_{disabled}}{t_{disabled}}$

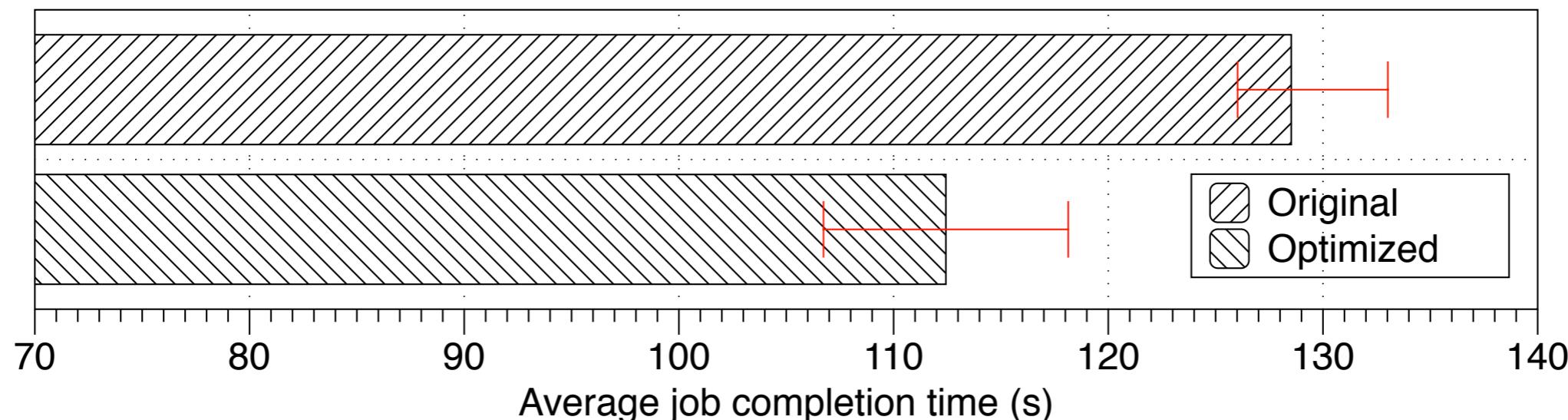
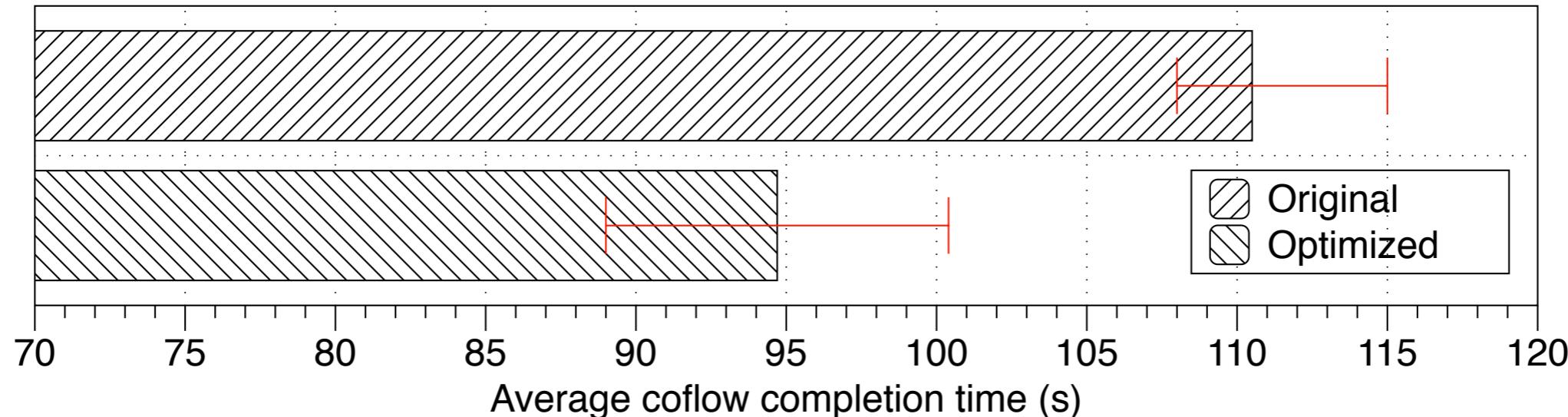


# SCALABILITY



# BENEFITS

- Hadoop TeraSort-25G
- 12.52% JCT reduction by a simple network scheduler



# RELATED WORK

- Analyze past statistics
  - Traffic Engineering with Estimated Traffic Matrices
- Monitor buffers or counters in switches
  - c-Through, Hedera, Helios
- Tracing and profiling toolkits
  - X-Trace
- File system monitoring
  - HadoopWatch

# SUMMARY

- DCF execution pattern
- DAG for predicting flows
- Design and implementation
- Evaluation on testbed



Thank | Q&A