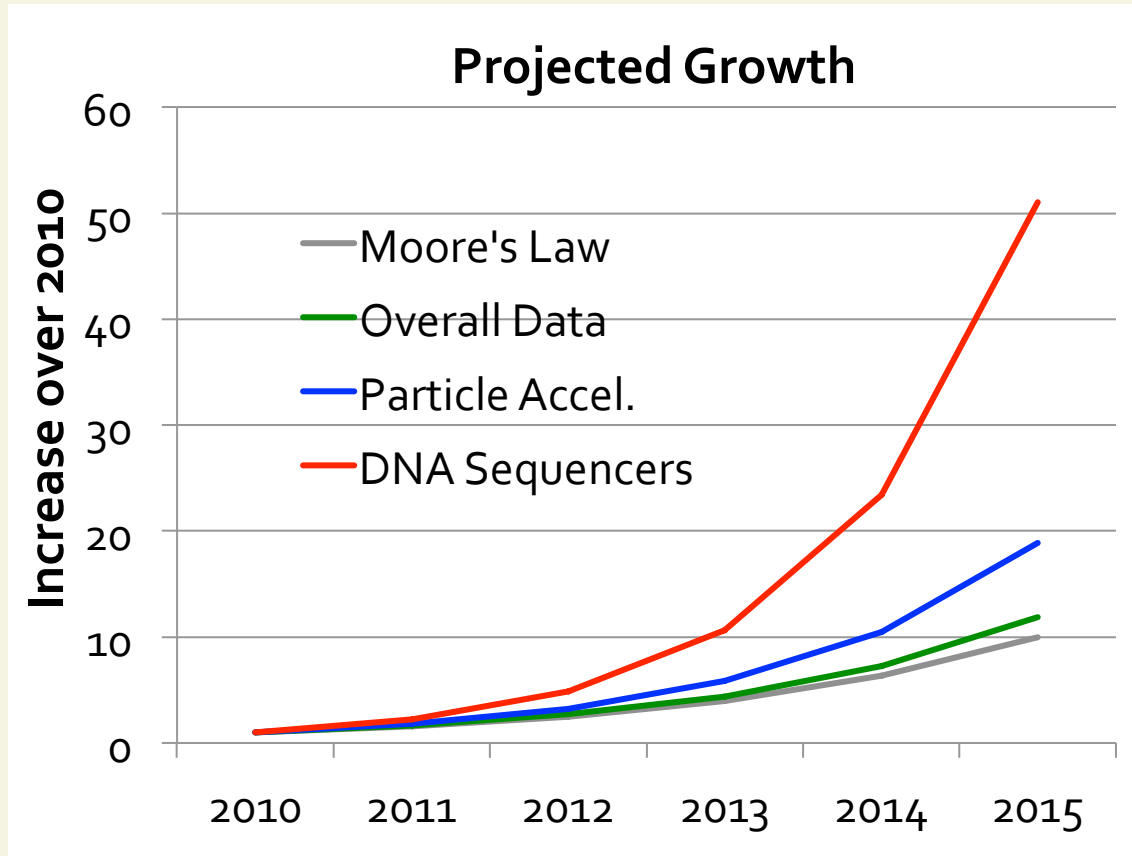# The Power of Choice in Data-Aware Cluster Scheduling

Shivaram Venkataraman[1] , Aurojit Panda[1]
Ganesh Ananthanarayanan[2], Michael Franklin[1], Ion Stoica[1]
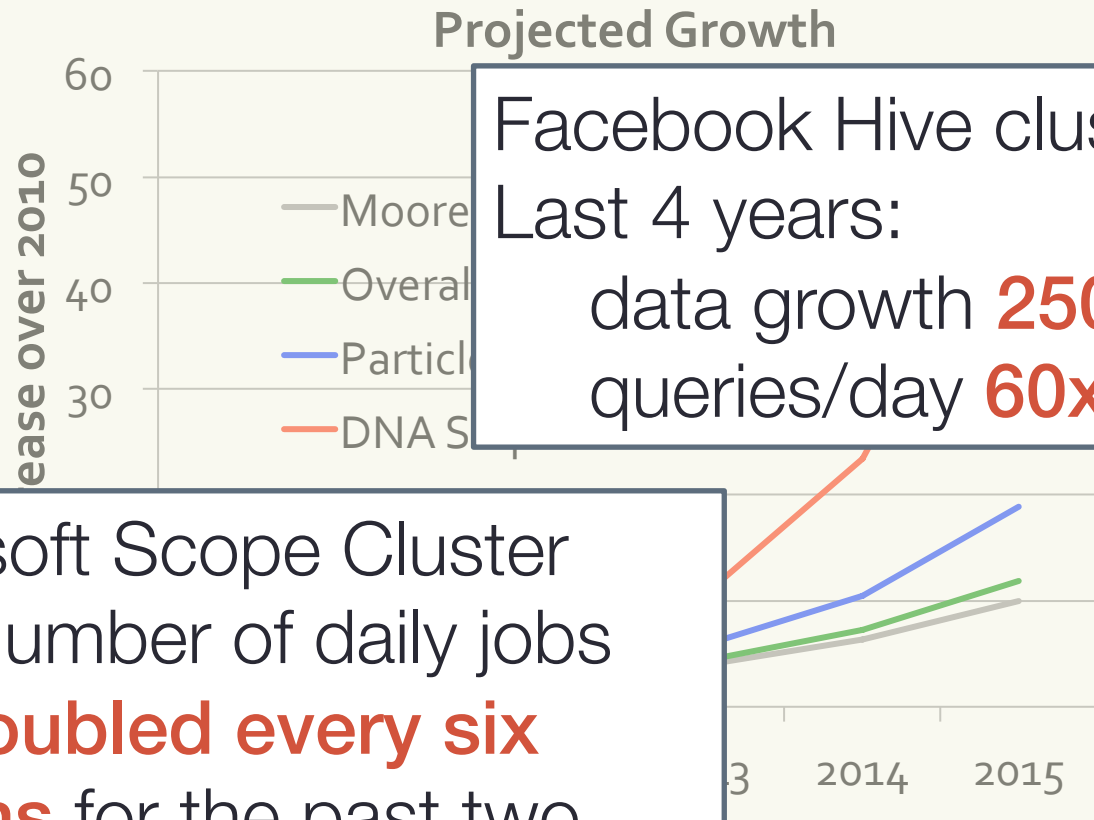
[1]UC Berkeley, [2]Microsoft Research

amplab

# Trends: Big Data

## Projected Growth

Increase over 2010

- Moore's Law
- Overall Data
- Particle Accel.
- DNA Sequencers

2010  2011  2012  2013  2014  2015

Data grows faster than Moore's Law

# Trends: Big Data

**Projected Growth**



(y-axis label) ease over 2010

60
50 — Moore...
40 — Overal...
30 — Particl...
— DNA S...
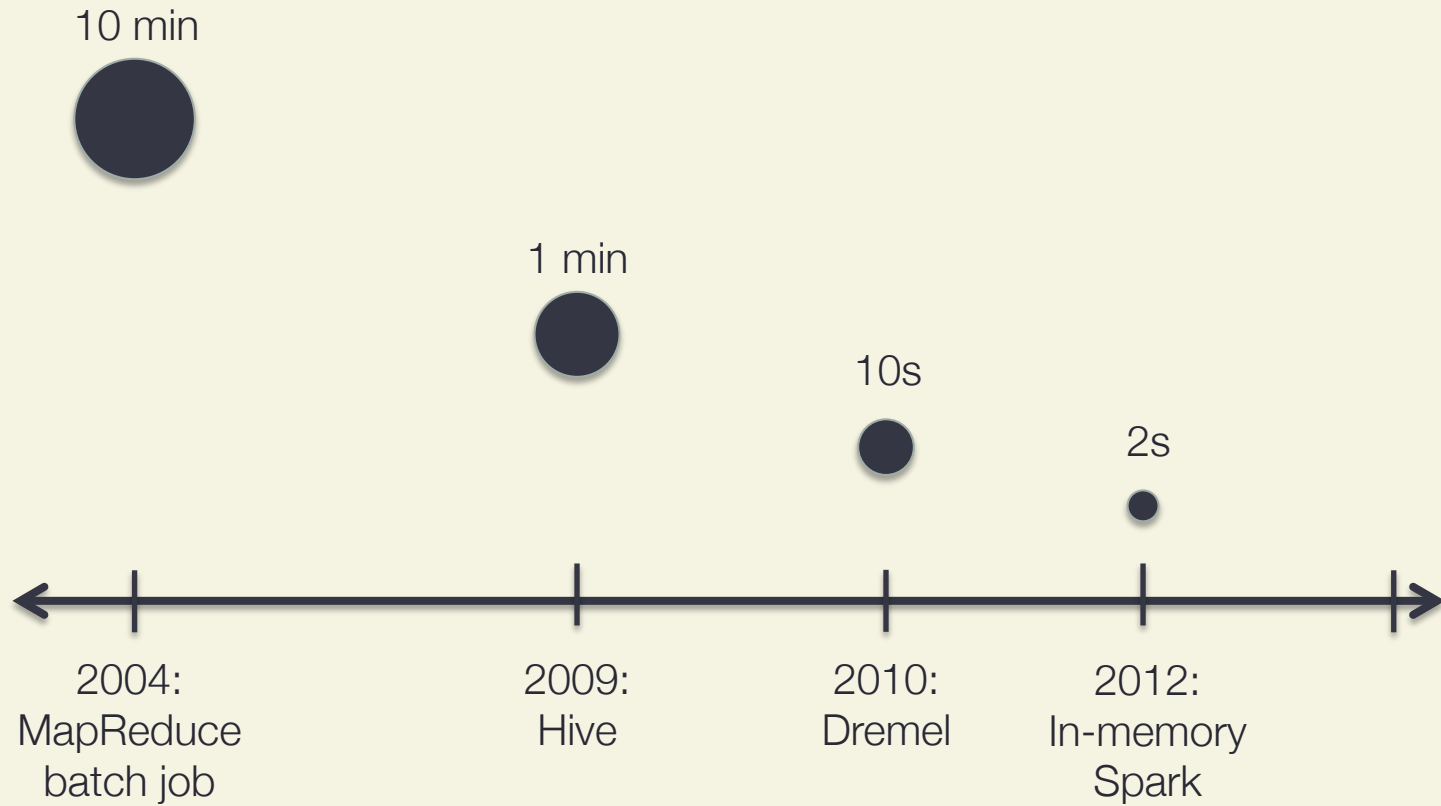
Facebook Hive cluster
Last 4 years:
   data growth **2500x**
   queries/day **60x**

Microsoft Scope Cluster
"The number of daily jobs has **doubled every six months** for the past two years."

...3   2014   2015

# Trends: Low Latency



10 min

1 min

10s

2s

2004:
MapReduce
batch job

2009:
Hive

2010:
Dremel

2012:
In-memory
Spark

# Big Data or Low Latency ?

SQL Query : 2.5 TB on 100 machines

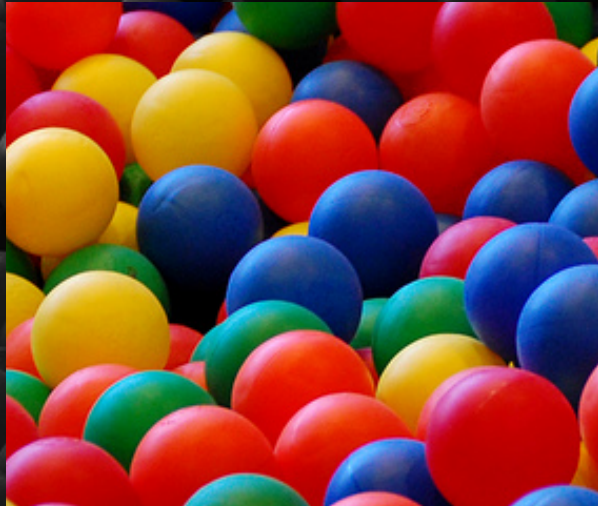> 15 minutes          1 - 5 Minutes          < 10s

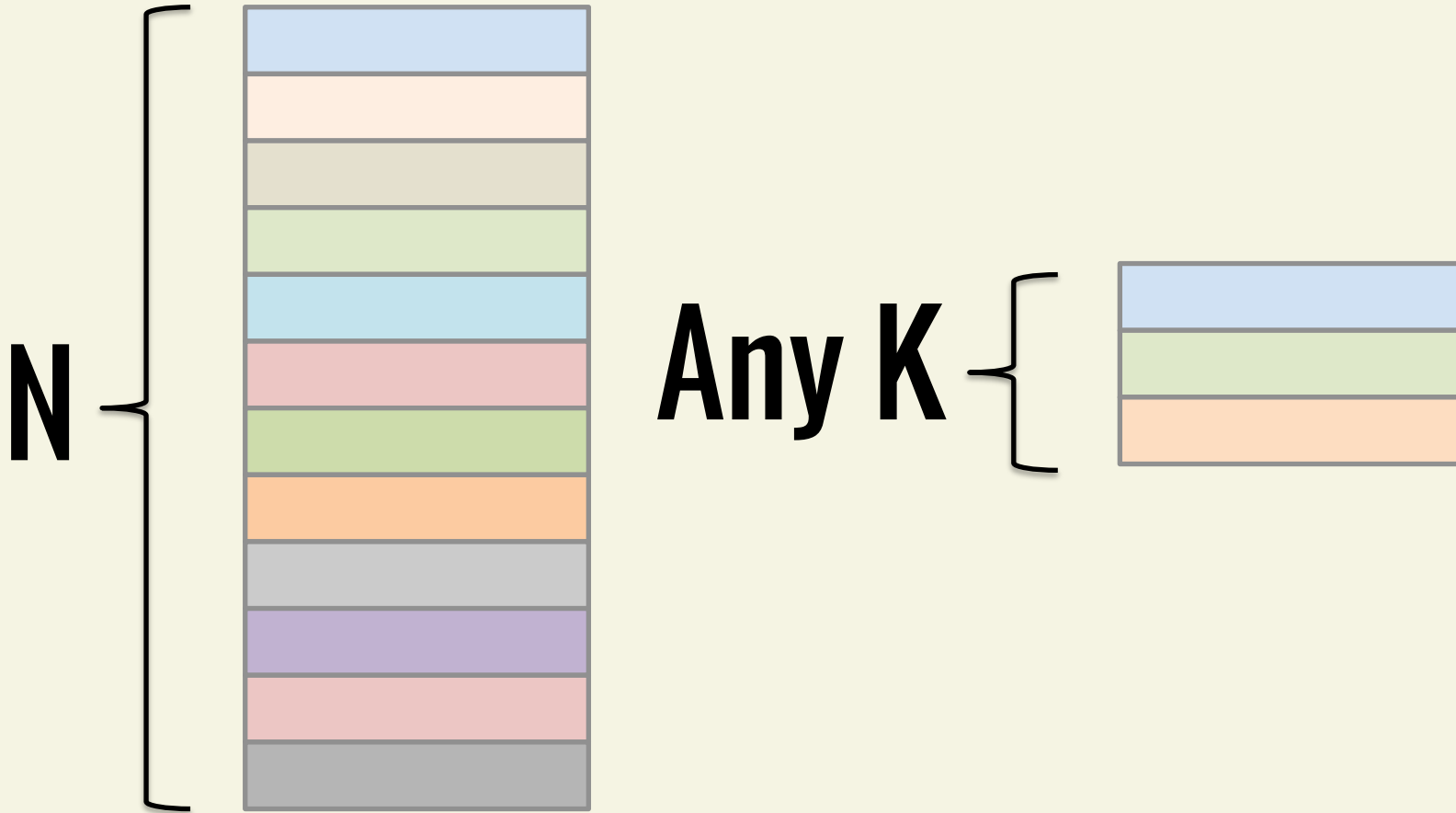# Sampling

# **Applications**

Approximate Query Processing
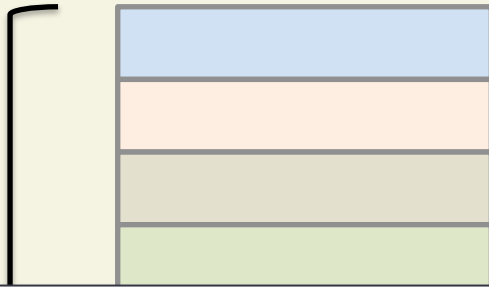    blinkdb, presto, minitable

Machine learning algorithms
    stochastic gradient, coordinate descent

# Choices
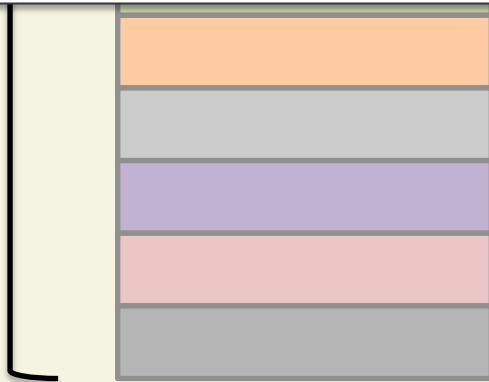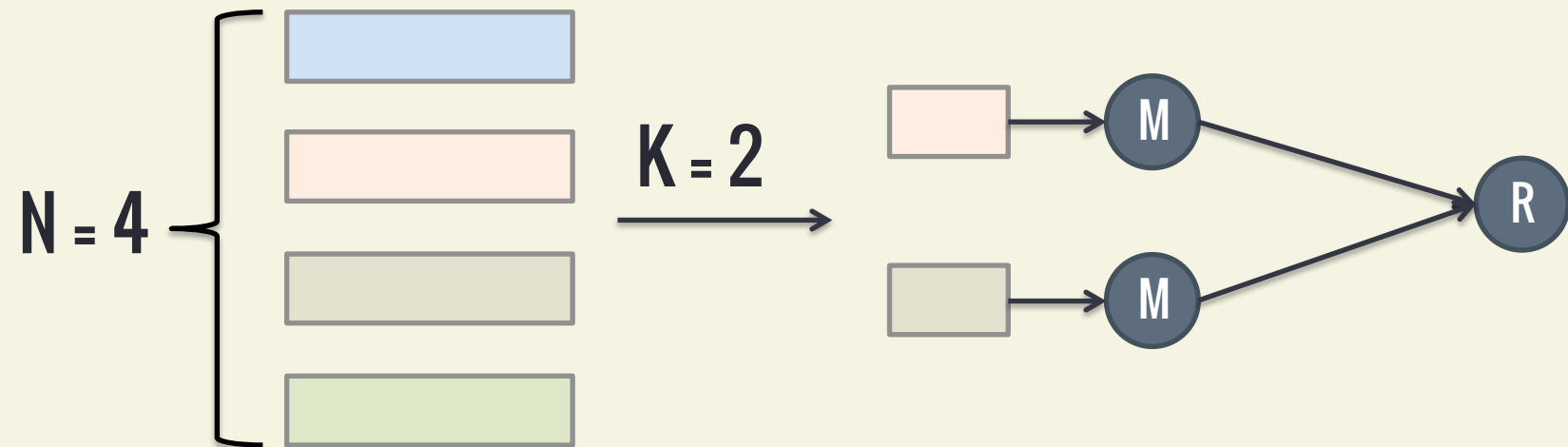
N

Any K

Sampling → Smaller Inputs + Choice

# Example

N = 4

K = 2

M

M

R

Available (N) = 2

Required (K) = 2

Rack

Time

Available Data

Running

Unavailable Data

Busy

# Choice-Aware

# Choice-Aware

**Available (N) = 4**  **Launched (M) = 3**  **Required (K) = 2**



Rack

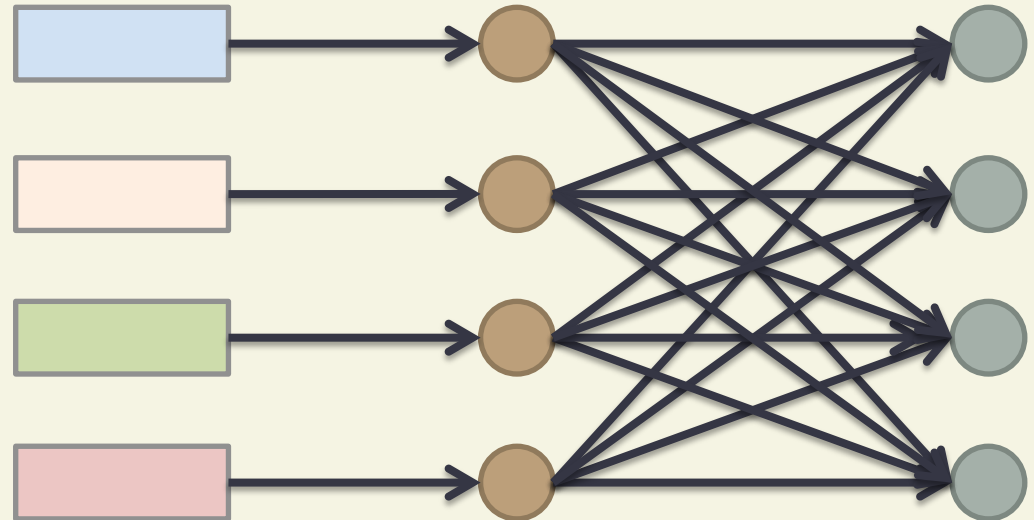Time

Available Data  Running

Busy

# KMN Scheduler

- How much can KMN improve locality
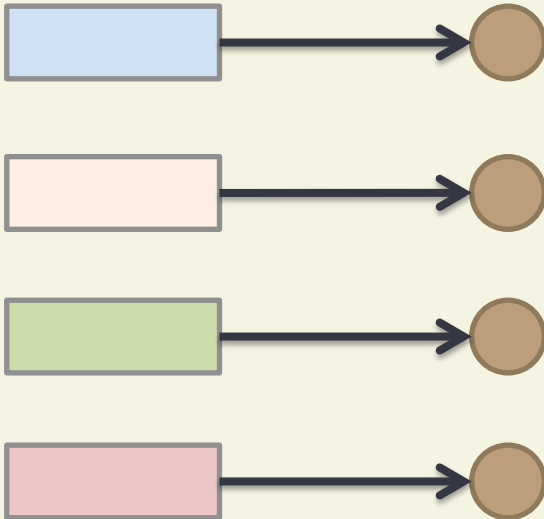- Propagate benefits across stages
- Handling stragglers

**KMN Scheduler**

**One-to-One**          **Many-to-One**
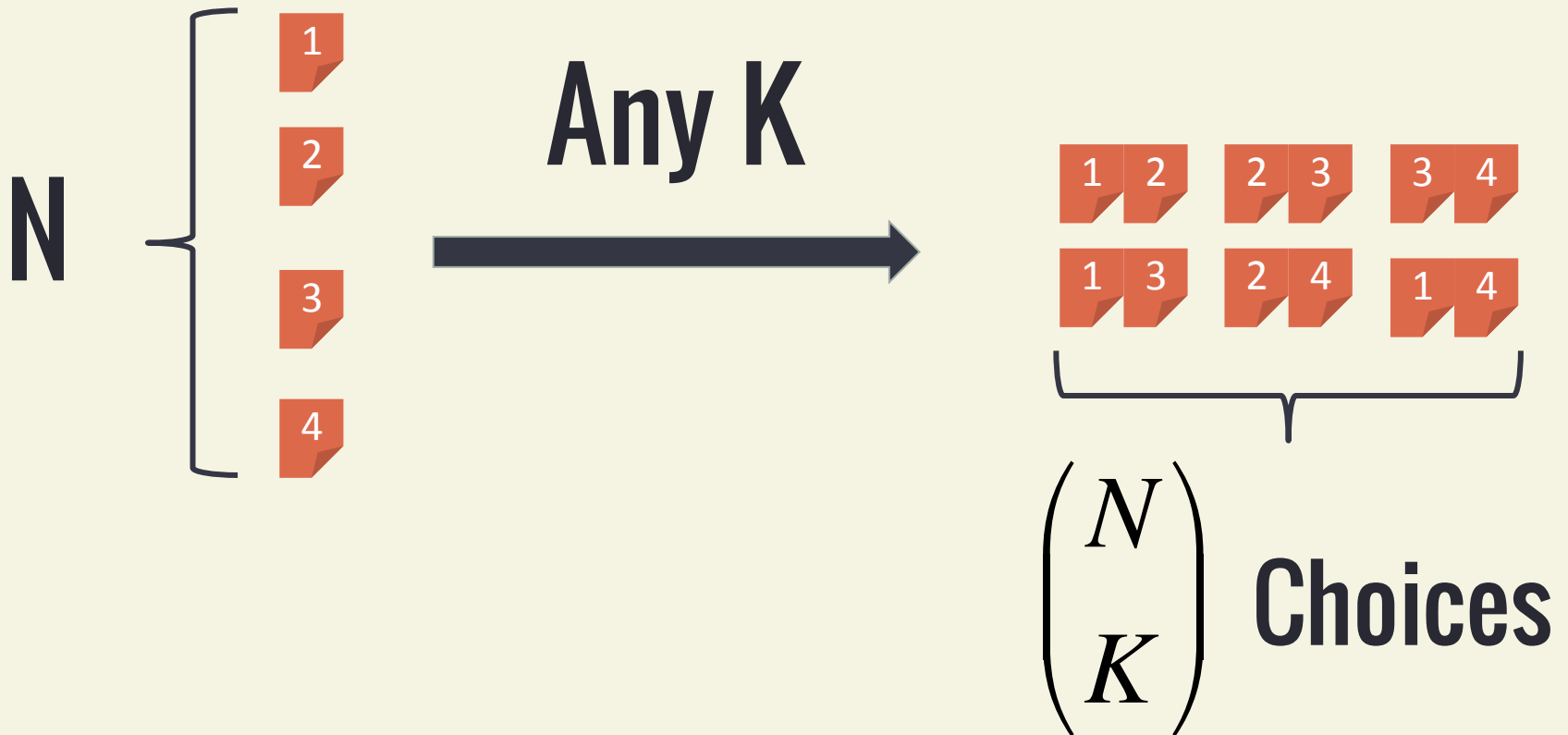
# One-to-One Stages

**Locality**

Disk ~ 100MB/s

Network ~ 10 Gbps (~1GB/s)

Memory ~ 50GB/s

# KMN Locality

N

1
2
3
4

**Any K**

→

1 2    2 3    3 4
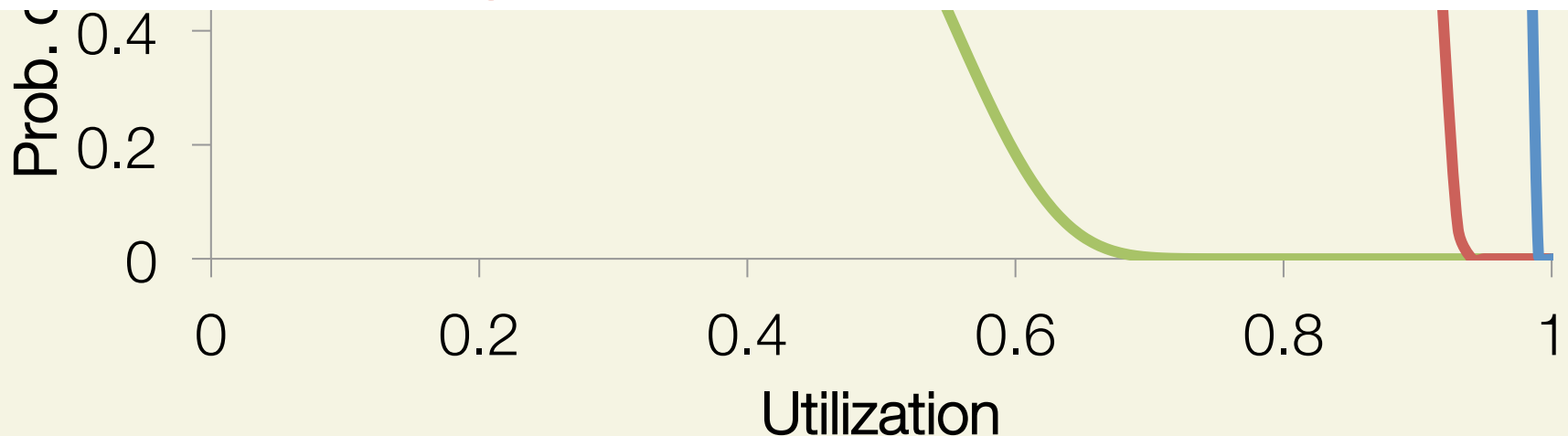1 3    2 4    1 4

$$\binom{N}{K}$$ **Choices**

# Locality, K=100

K – Number of blocks chosen
N – Number of blocks available

K/N=1.0  K/N=0.5  K/N=0.1
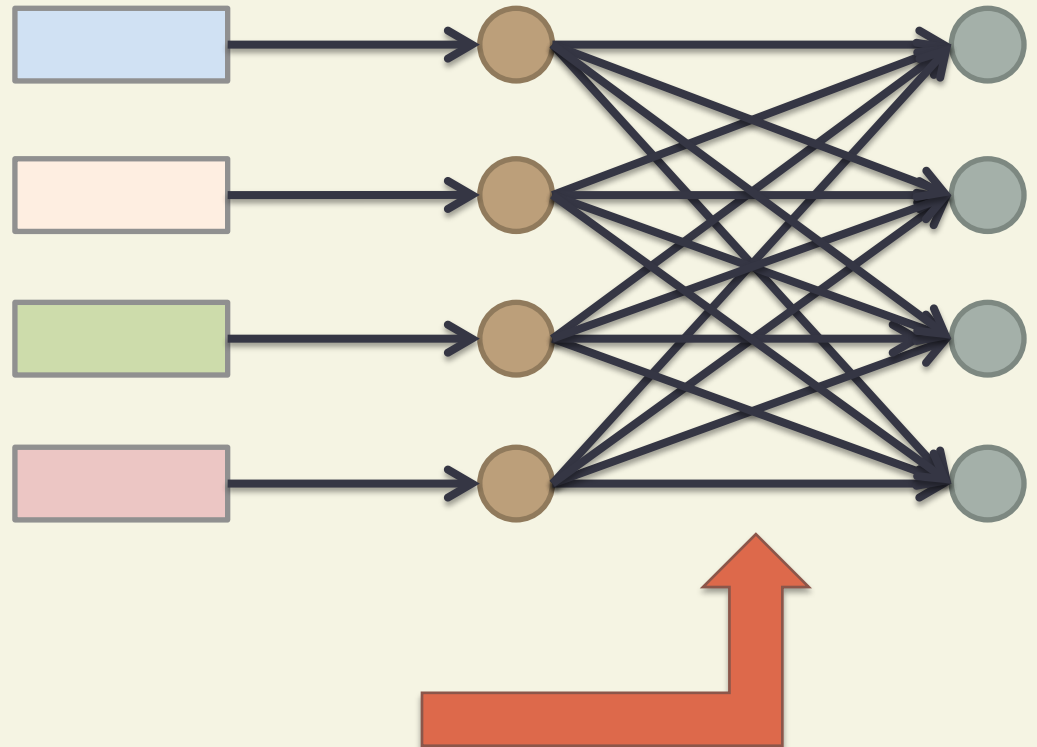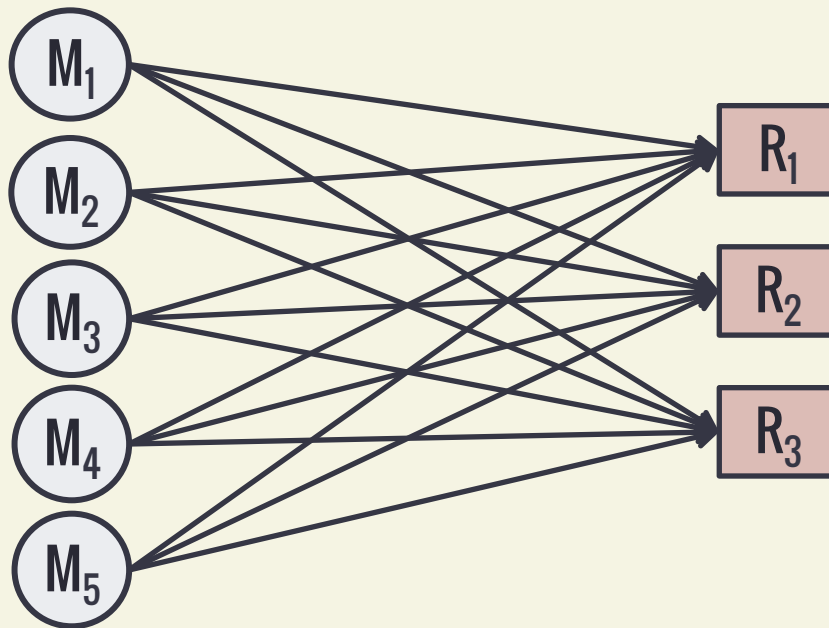


KMN significantly improves locality
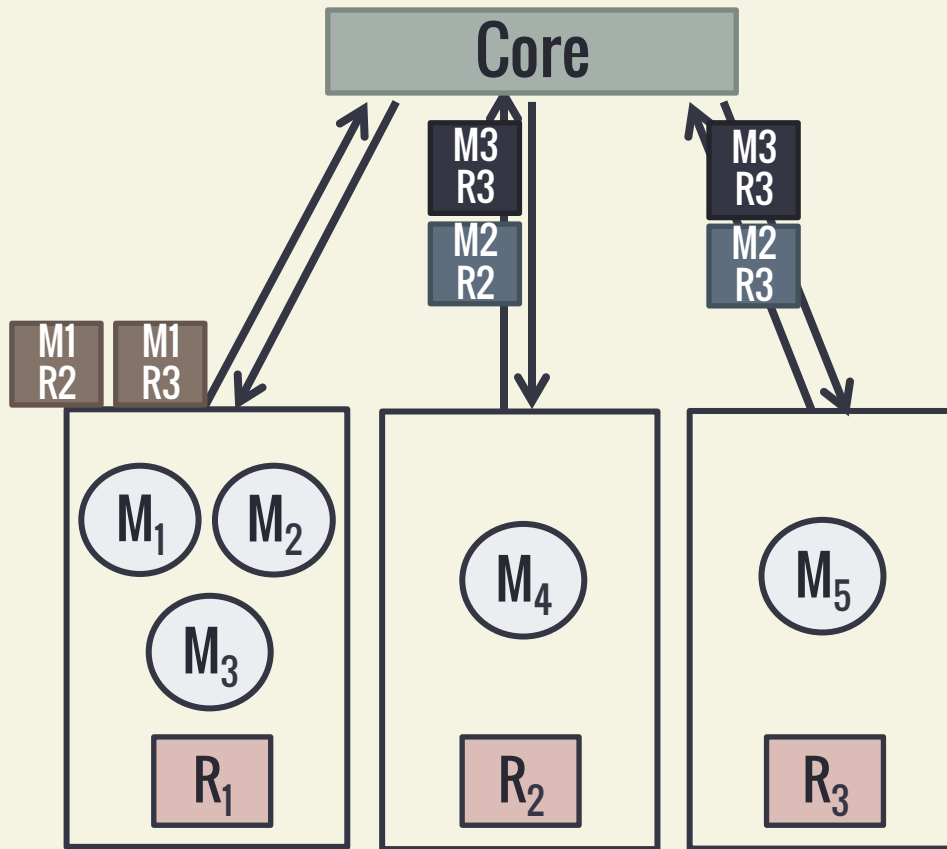
# Many-to-One Stages
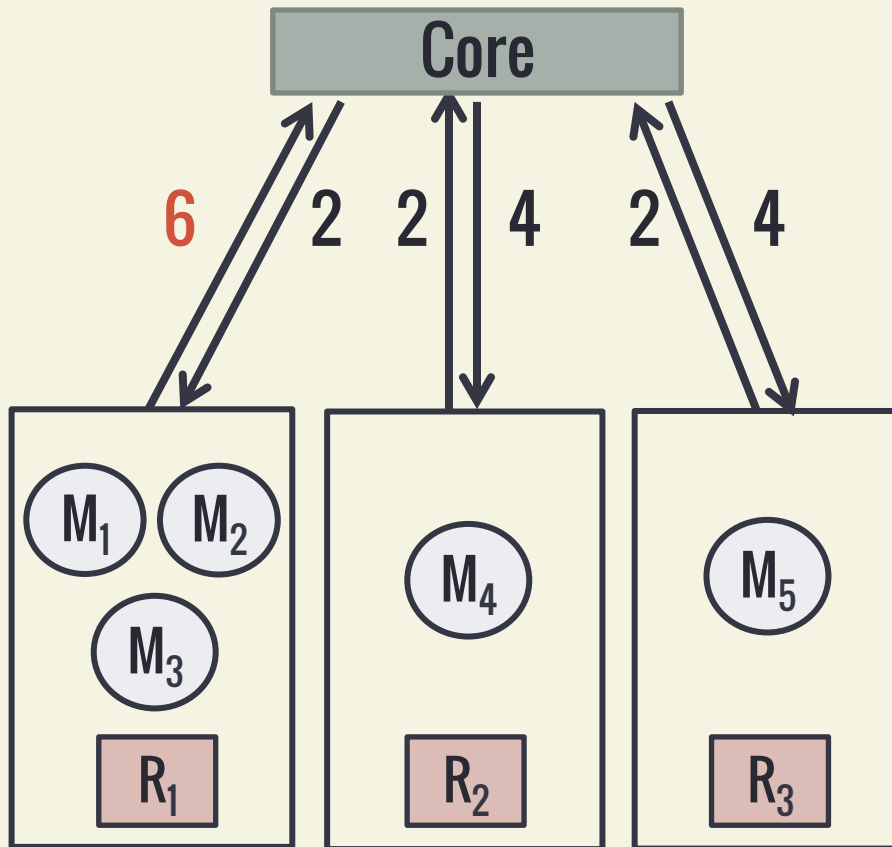
# Many-to-One Stage



**15 transfers**

# Many-To-One Transfers
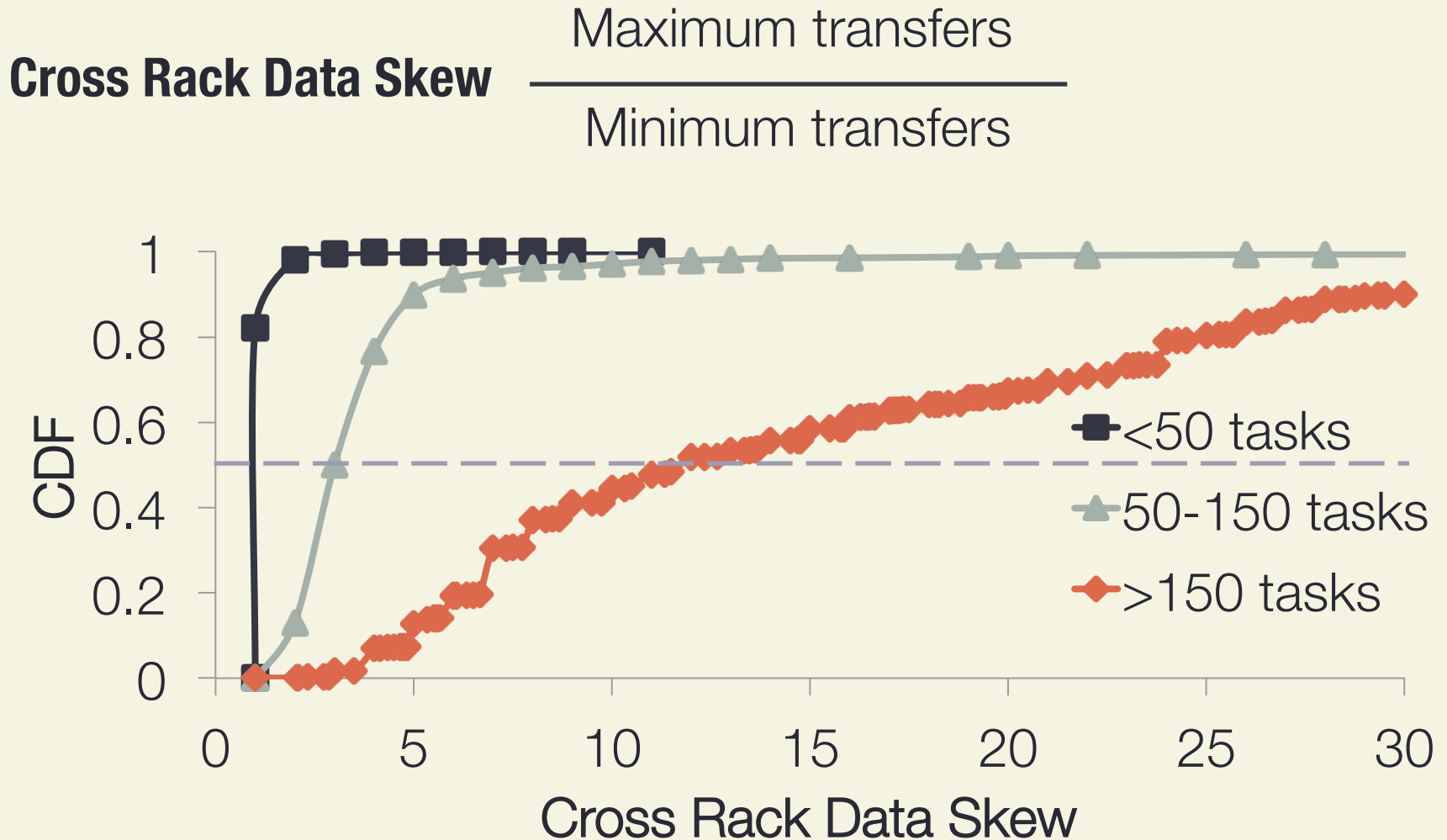
# Bottleneck Link



**Bottleneck Link**

Link with Max. transfers

**Cross Rack Data Skew**

$$= \frac{\text{Maximum transfers}}{\text{Minimum transfers}}$$

$$= \frac{6}{2} = 3$$

# Facebook Trace

**Cross Rack Data Skew** $\dfrac{\text{Maximum transfers}}{\text{Minimum transfers}}$



CDF vs Cross Rack Data Skew

- ■ <50 tasks
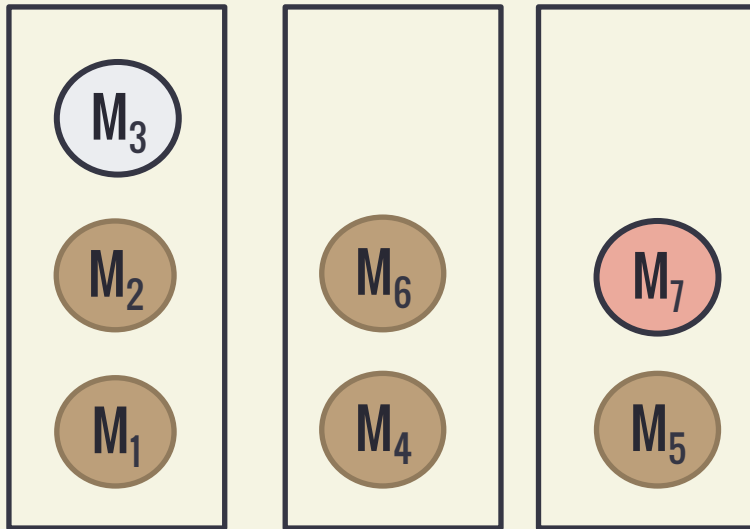- ▲ 50-150 tasks
- ◆ >150 tasks

# Power of Choice



**Cross Rack Data Skew = 3**

Load balancing: balls and bins

Insight: Run extra tasks (M > K)

# Power of Choice



Technique:

Spread out choice of K tasks to reduce skew

**M = 7, K = 5**
**Cross Rack Data Skew = 2**

# Handling Stragglers



$M_1$
$M_2$
$M_3$

Rack

$M_4$
$M_6$

Rack

$M_5$
$M_7$

Time

Stragglers

vs.

Cross-Rack
Data Skew

# Using KMN

```
// Create Spark RDD
file = sc.textFile("tpc-h.data")

// Select a 10% sample using KMN
sample = file.blockSample(0.1)

// RDD operations
sample.map { li =>
  (li.linestatus, li.quantity)
}.collect()
```

# **Also in the paper**

User-defined sampling functions

Placing reduce tasks

Killing extra tasks

# Evaluation

Facebook traces replay

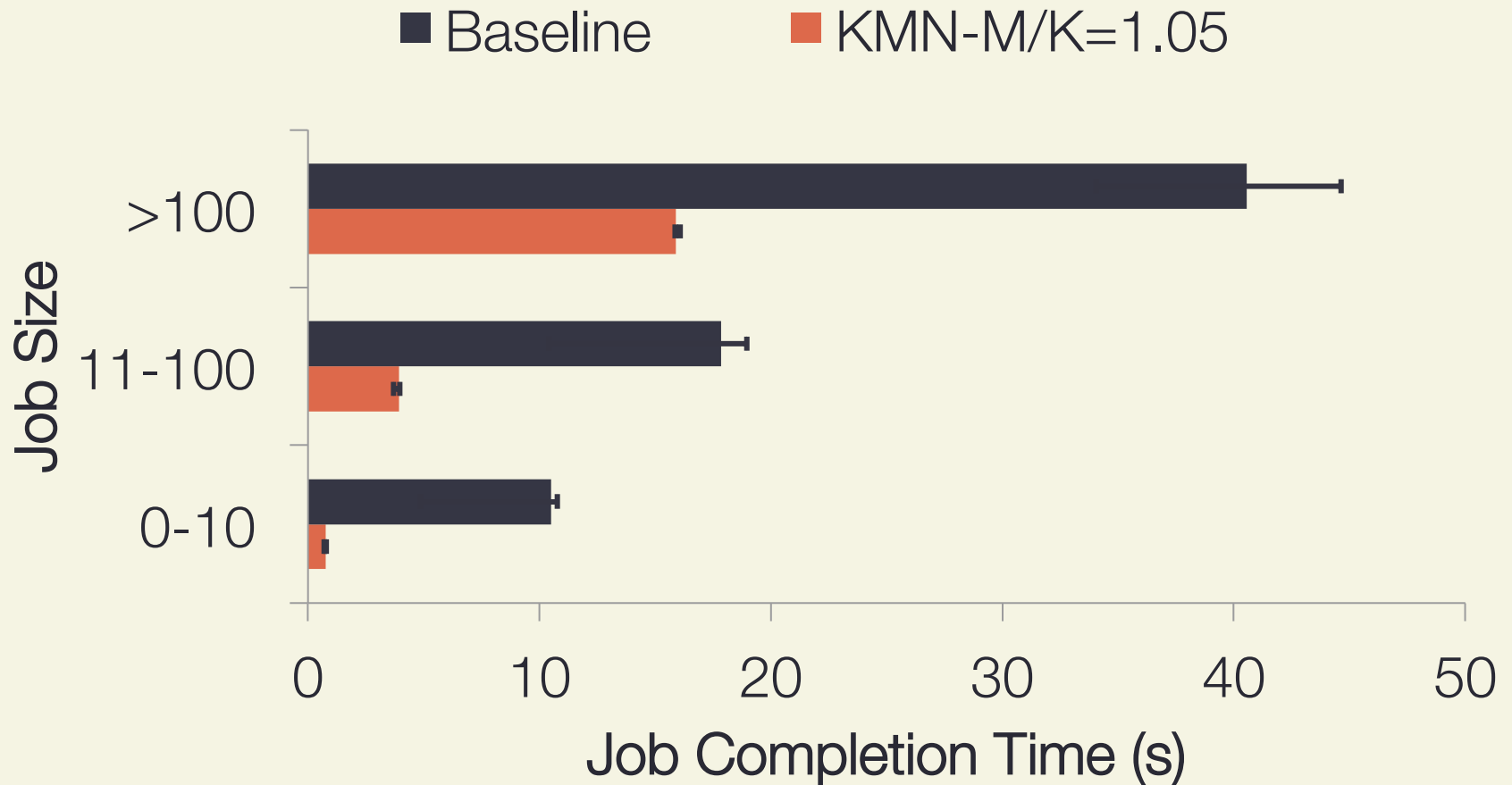Long DAGs (Stochastic Gradient Descent)

SQL queries from Conviva

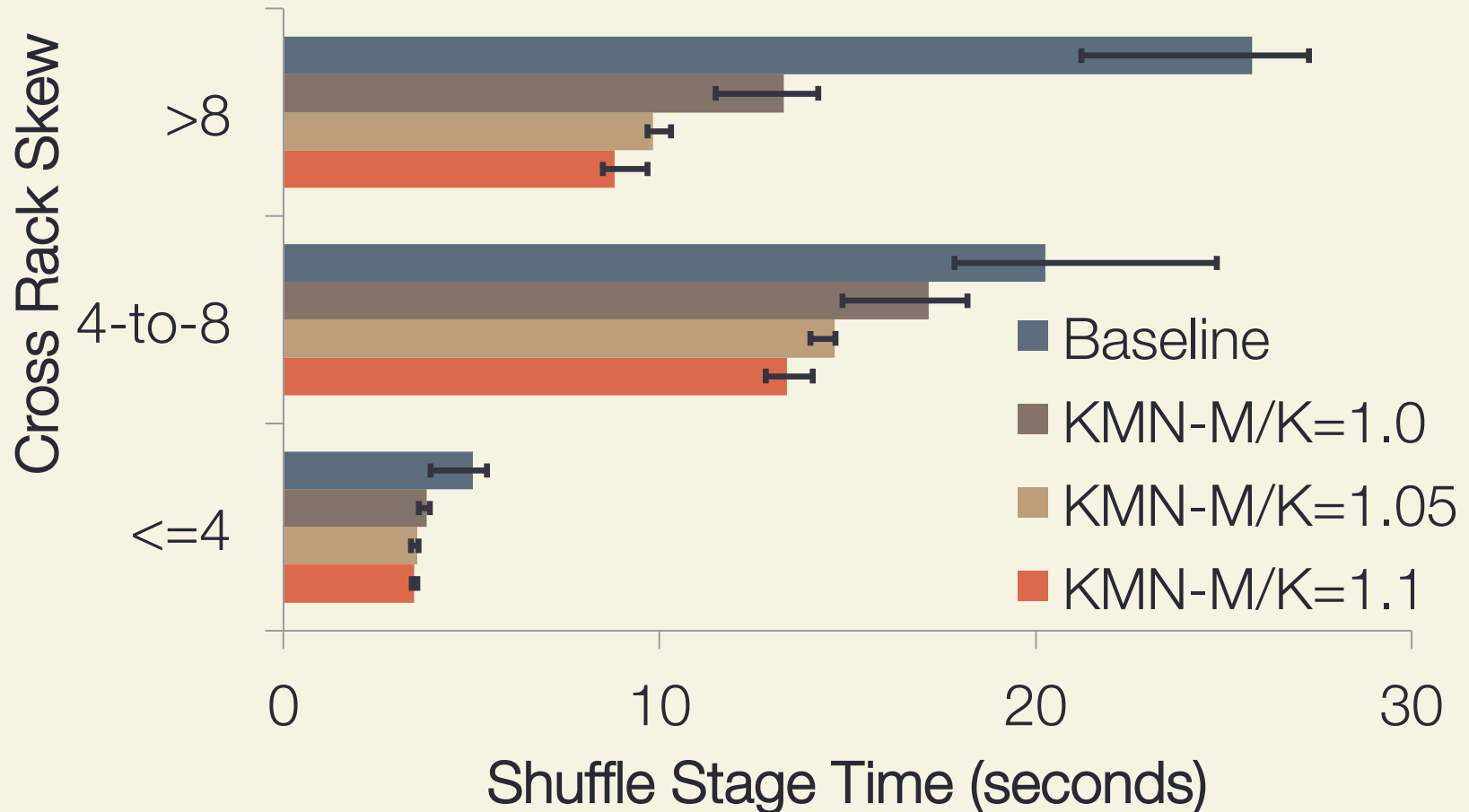Reducer placement

Varying Utilization

Baseline: Use a pre-selected random sample
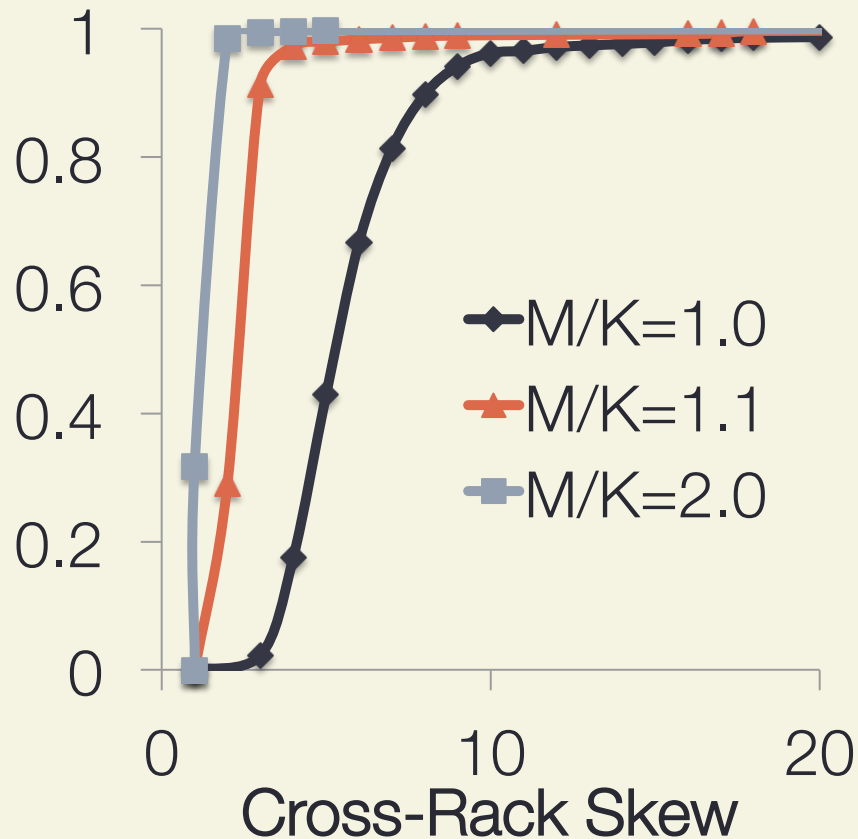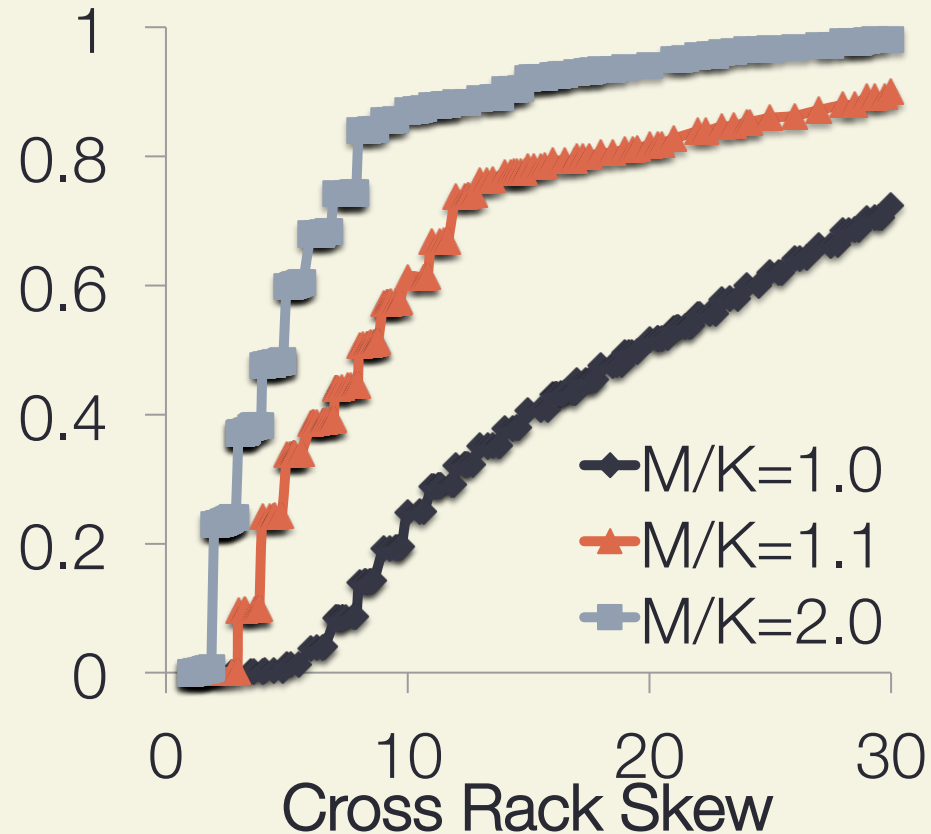Setup: 100 m2.4xlarge EC2 machines, 60GB RAM/mc

# Facebook Overall

Legend: Baseline, KMN-M/K=1.05

Chart — Job Completion Time (s) vs Job Size (0-10, 11-100, >100)

# Cross Rack Skew



Chart: Shuffle Stage Time (seconds) by Cross Rack Skew

Legend:
- Baseline
- KMN-M/K=1.0
- KMN-M/K=1.05
- KMN-M/K=1.1

Y-axis (Cross Rack Skew): >8, 4-to-8, <=4

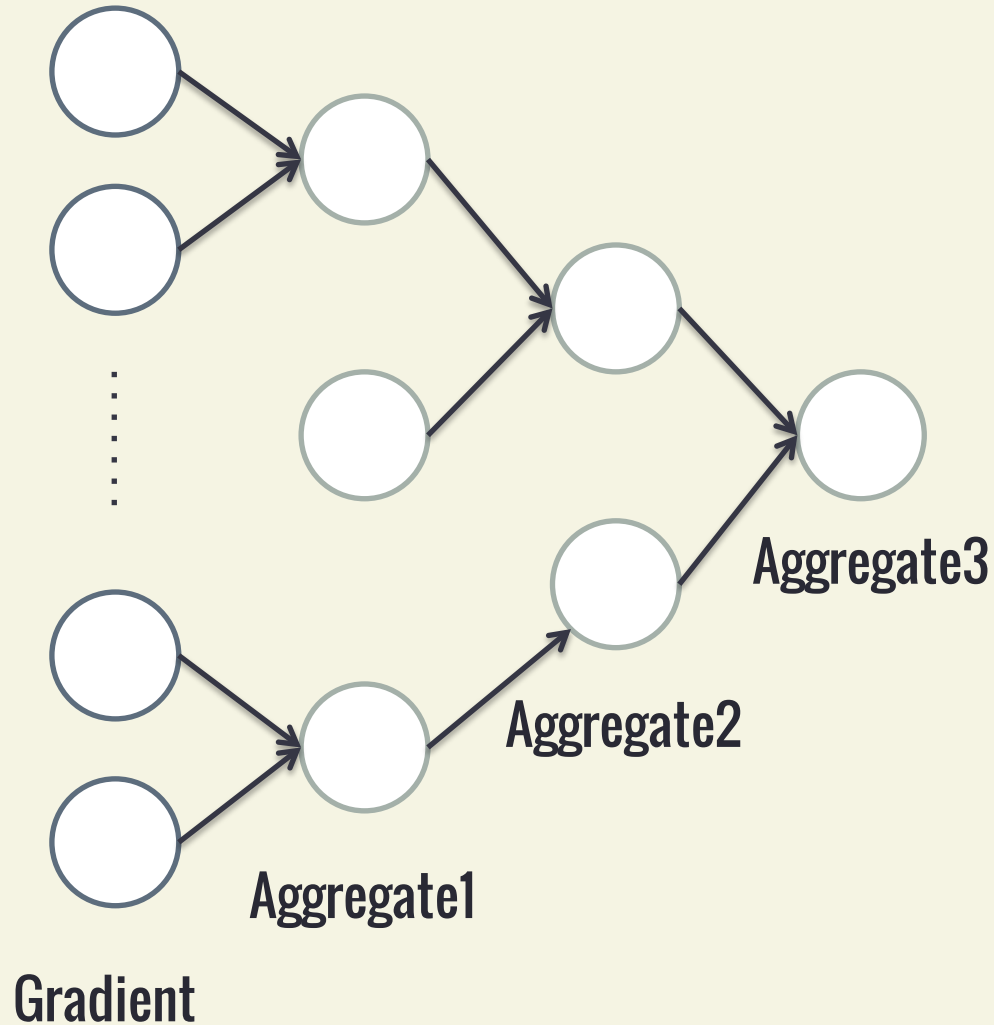X-axis: Shuffle Stage Time (seconds): 0, 10, 20, 30

# How many extra tasks ?
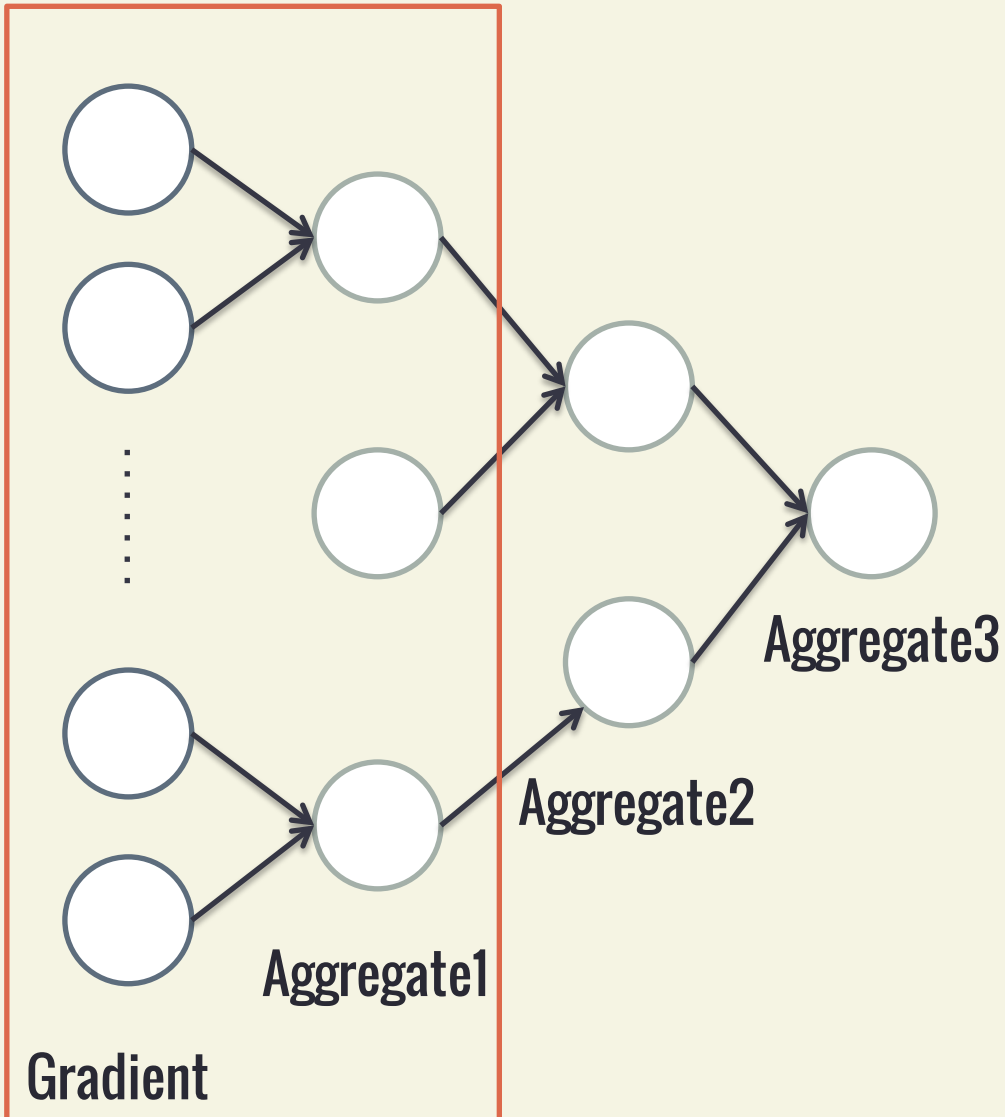


50 - 150 tasks

>150 tasks

# KMN: How many stages ?



Stochastic
Gradient
Descent

Aggregate3

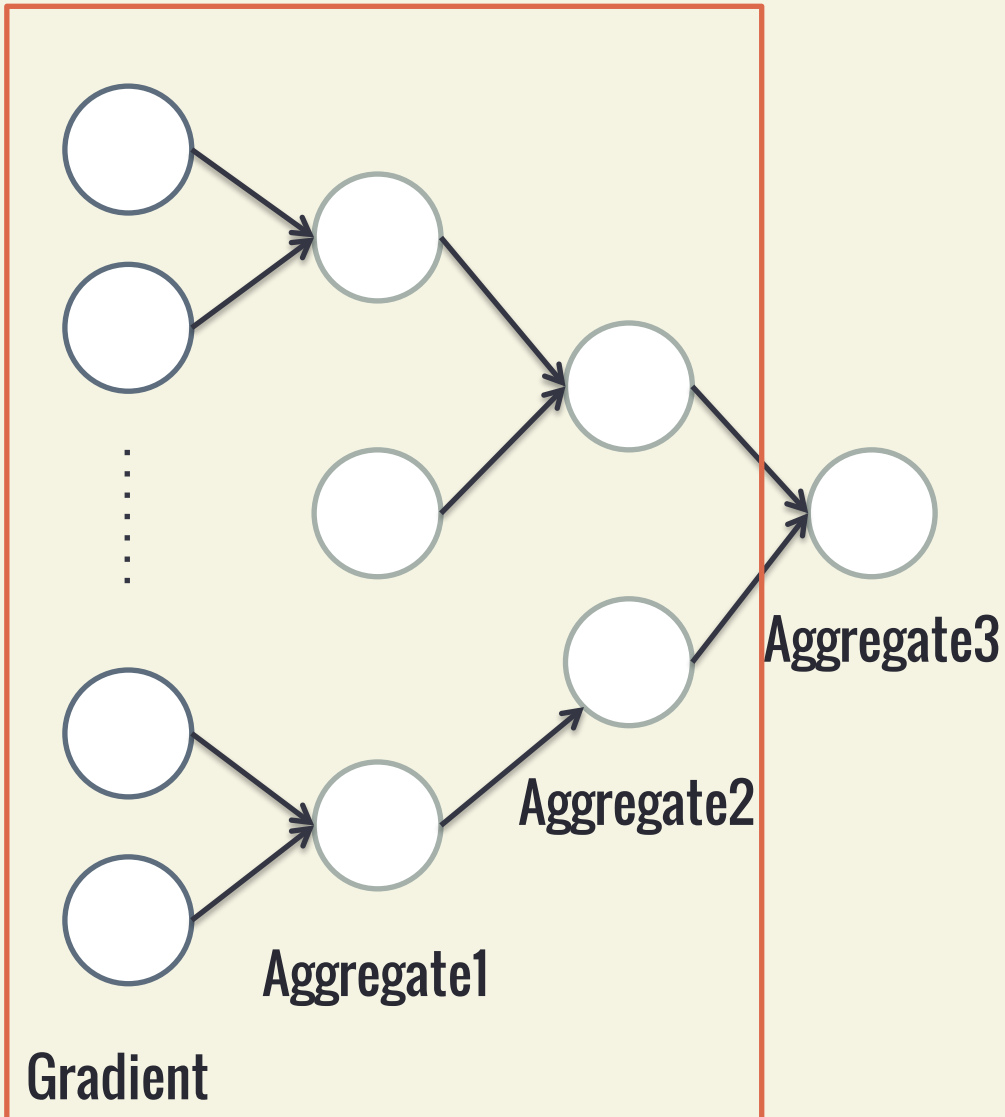Aggregate2

Aggregate1

Gradient

# KMN: How many stages ?



| KMN Stages | Time (s) |
|------------|----------|
| Gradient | 15.27 |

# KMN: How many stages ?



| KMN Stages | Time (s) |
|---|---|
| Gradient | 15.27 |
| Gradient + Agg1 | 12.72 |

# KMN: How many stages ?



| KMN Stages | Time (s) |
|---|---|
| Gradient | 15.27 |
| Gradient + Agg1 | 12.72 |
| Gradient + Agg2 | 11.79 |

# KMN: How many stages ?



Diagram showing Gradient nodes feeding into Aggregate1, Aggregate2, and Aggregate3 stages.

| KMN Stages | Time (s) |
|---|---|
| Gradient | 15.27 |
| Gradient + Agg1 | 12.72 |
| Gradient + Agg2 | 11.79 |
| Gradient + Agg3 | 12.09 |

# **Related Work**

Power of Choice

    Power-of-Two choices [TPDS'01]

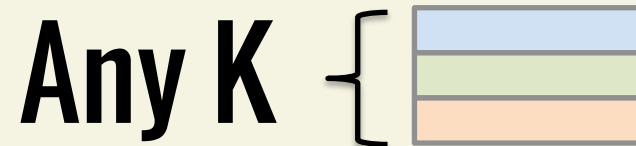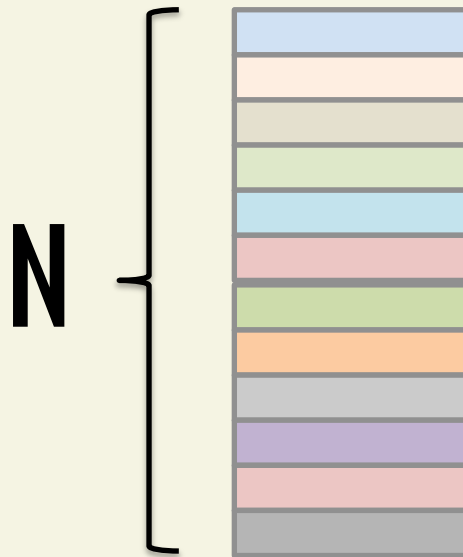    Sparrow [SOSP'13]

Improving Cluster Scheduling

    Quincy [SOSP'09]

    alsched [SOCC'12]

    Dolly [NSDI'13]

# KMN Scheduler

N

Any K

Emerging applications: ML algorithms, AQP
Improves locality, Balances network transfers