

On the Latency-Accuracy Tradeoff in Approximate MapReduce Jobs

Juan F. Pérez

Universidad del Rosario, Colombia

Robert Birke and Lydia Chen

IBM Research Zurich, Switzerland

IEEE INFOCOM – May, 2017



Universidad del
Rosario



MACC
Applied Mathematics
and Computer Science



Outline

- 1 Introduction
- 2 System Operation
- 3 The Model
- 4 Experimental assessment
- 5 Wrap-up

What?

MapReduce Jobs

- Large data sets
- Exploit large number of parallel resources
- Open source implementations (Hadoop)
- Well studied for batch jobs
- **Online** job execution
- Deadlines: constraints on response time

MapReduce Jobs

- Large data sets
- Exploit large number of parallel resources
- **Online** job execution
- Deadlines: constraints on response time

Online MapReduce Jobs

- Deadlines: constraints on response time
- Process **all** data → *exact* answer
- **Approximate** Computing
- Process **some** data → *approximate* answer
- Save computing time and resources
- Comply with deadlines

Why?

Latency-Accuracy Tradeoff

- Drop more data → Reduce latency → Reduce accuracy
- Drop less data → Increase latency → Increase accuracy
- Solve the **Latency-Accuracy Tradeoff**

Goal?

Solve Latency-Accuracy Tradeoff

- Given a *Latency* objective → Determine *Achievable Accuracy*
- Given an Accuracy objective → Determine Achievable Latency
- Additional considerations:
 - Jobs arrive in a random fashion
 - Queueing times waiting for jobs in front to be processed
 - How to drop data for approximate computations?

Solve Latency-Accuracy Tradeoff

- Given a *Latency* objective → Determine *Achievable Accuracy*
- Given an Accuracy objective → Determine Achievable Latency
- Additional considerations:
 - Jobs arrive in a random fashion
 - Queueing times waiting for jobs in front to be processed
 - How to drop data for approximate computations?

How?

Our Approach to the Latency-Accuracy Tradeoff

- Stochastic models to predict **Job Latency**
 - Capture Online MapReduce operation
 - Capture Approximate Computing
 - Data Drop options

Our Approach to the Latency-Accuracy Tradeoff

- Stochastic models to predict **Job Latency**
- Capture Online MapReduce operation
- Capture Approximate Computing
- Data Drop options

Outline

1 Introduction

2 System Operation

3 The Model

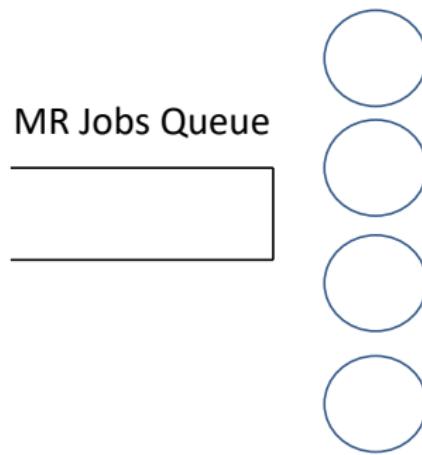
4 Experimental assessment

5 Wrap-up

Basic Operation

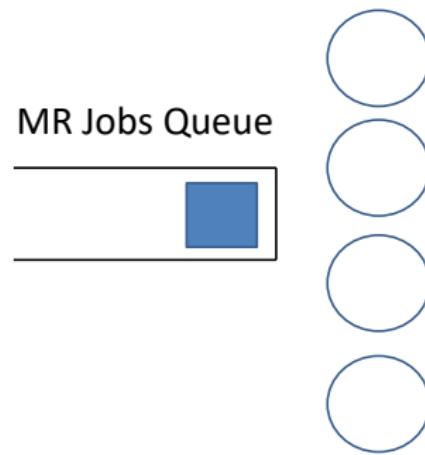
Basic Operation

C computing slots



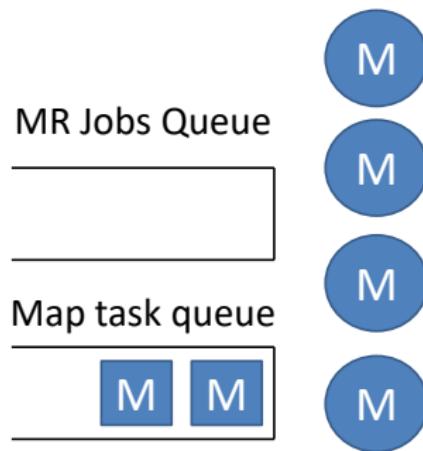
Basic Operation

C computing slots



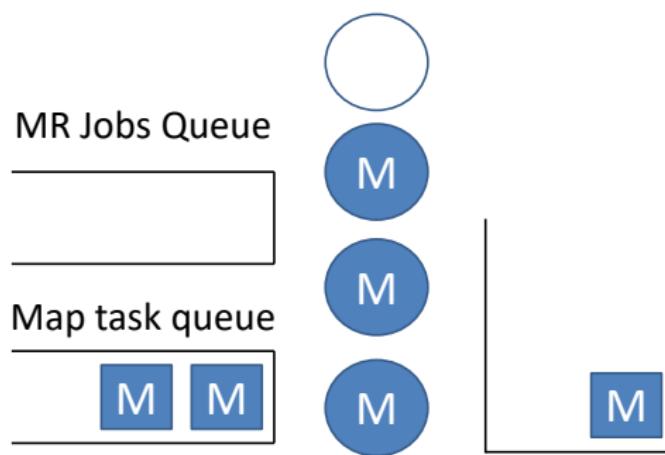
Basic Operation

C computing slots



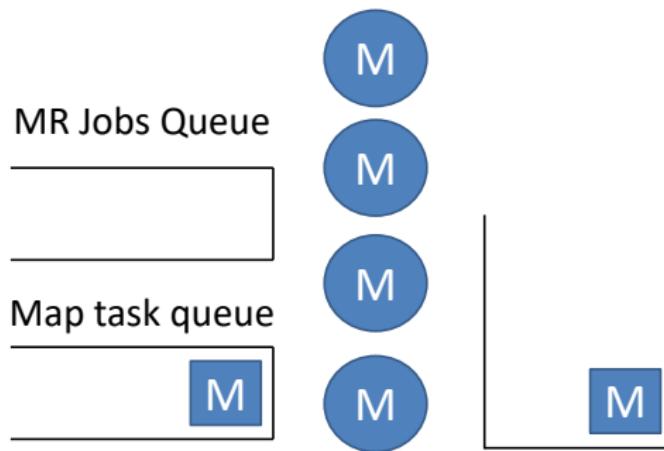
Basic Operation

C computing slots



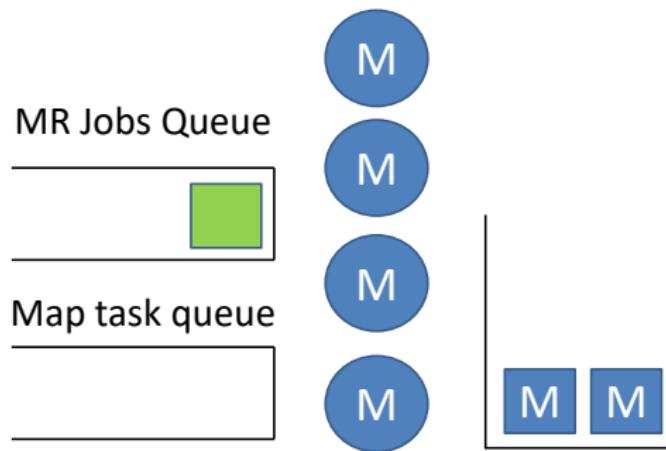
Basic Operation

C computing slots



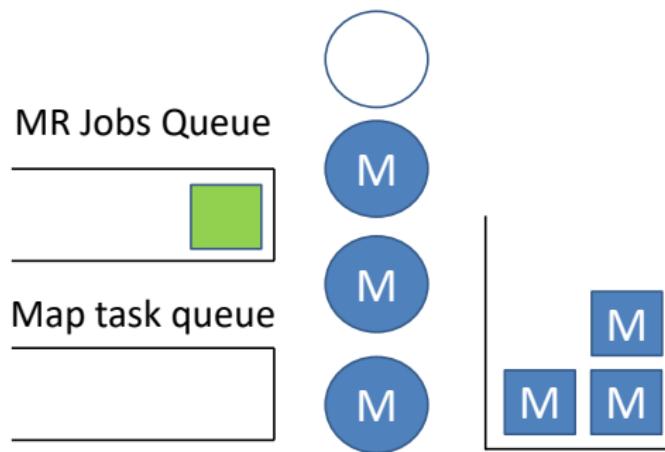
Basic Operation

C computing slots



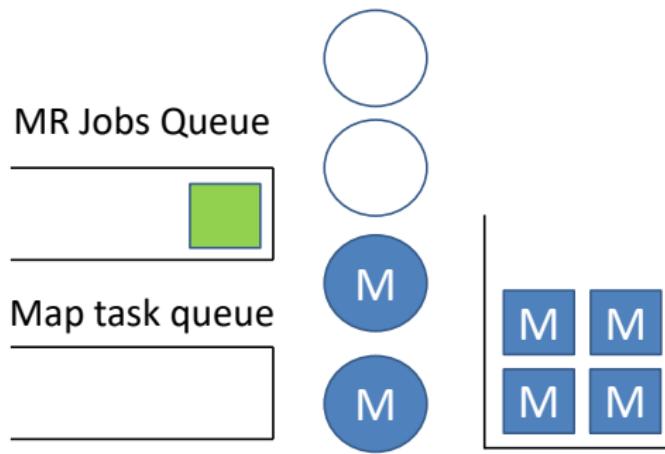
Basic Operation

C computing slots



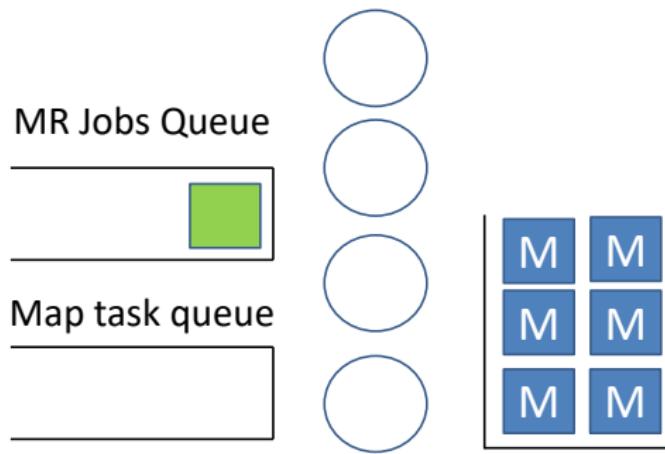
Basic Operation

C computing slots



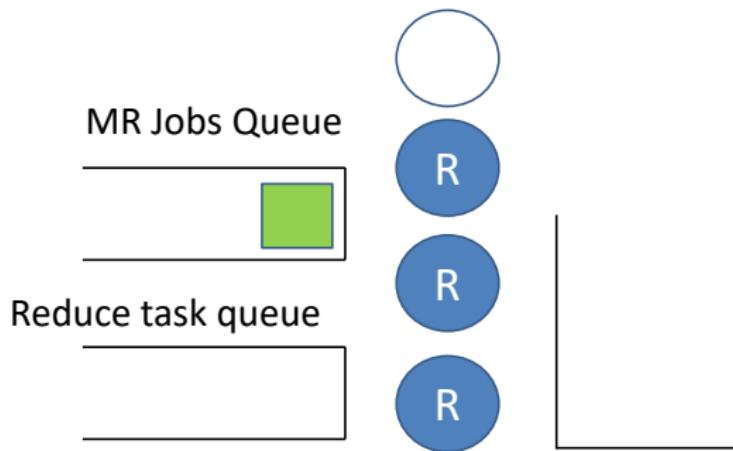
Basic Operation

C computing slots



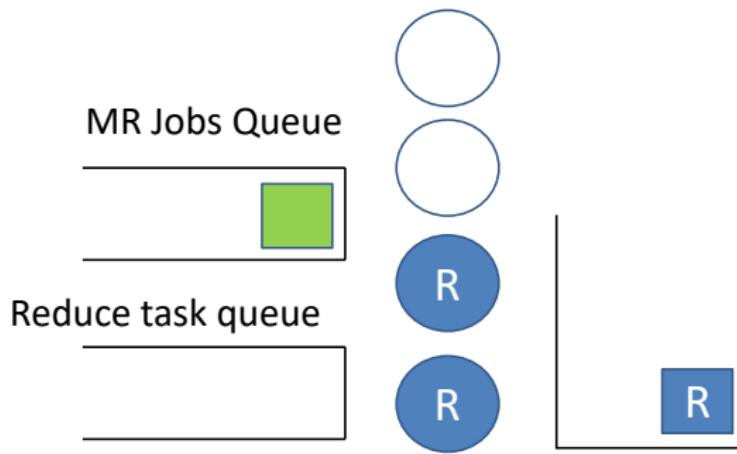
Basic Operation

C computing slots



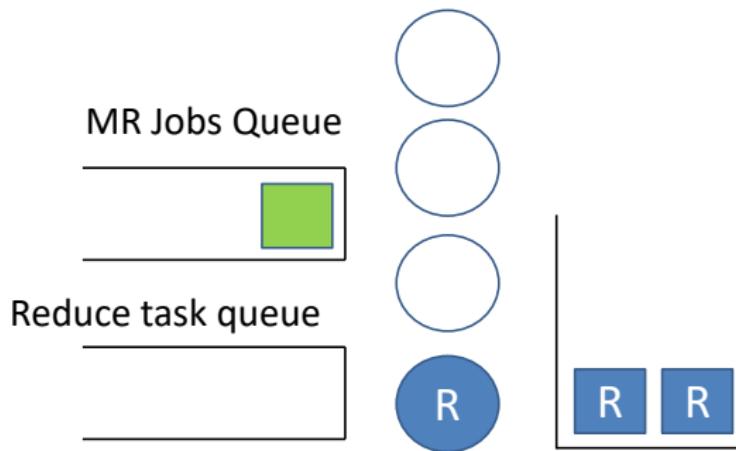
Basic Operation

C computing slots



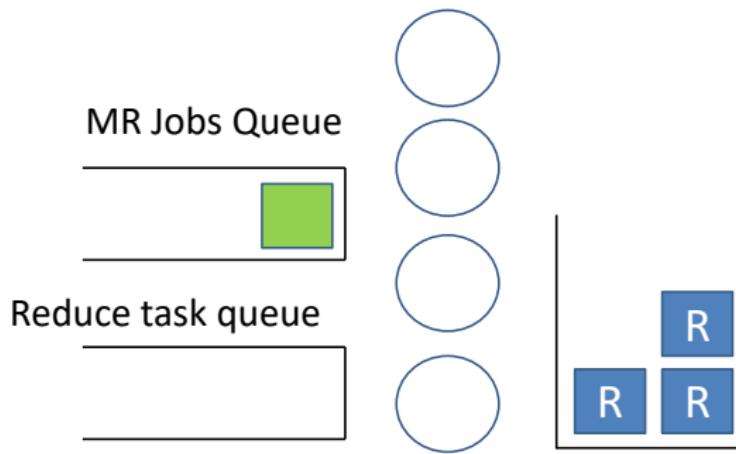
Basic Operation

C computing slots



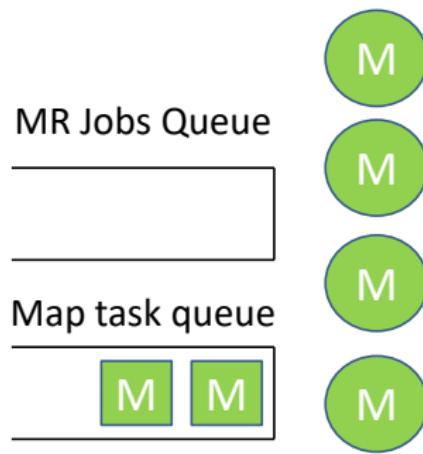
Basic Operation

C computing slots



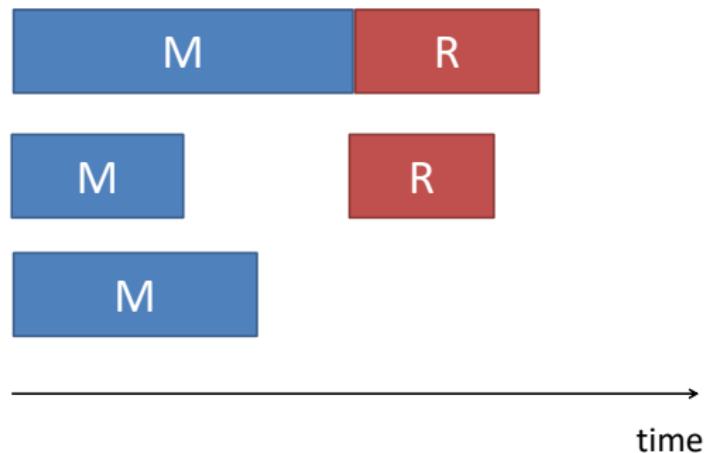
Basic Operation

C computing slots



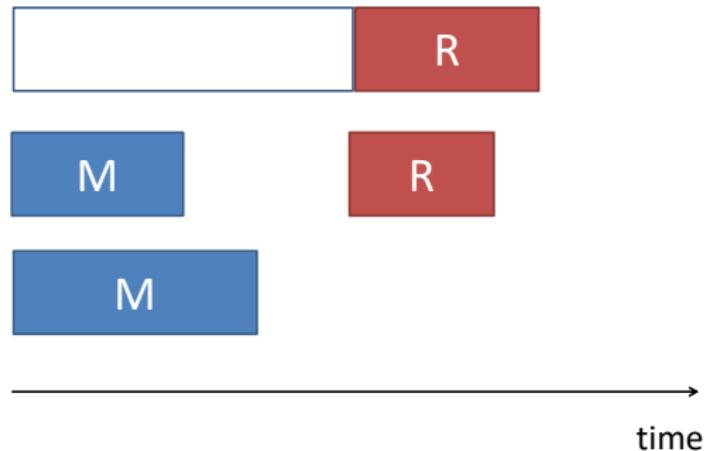
Approximate Computing

Data Drop

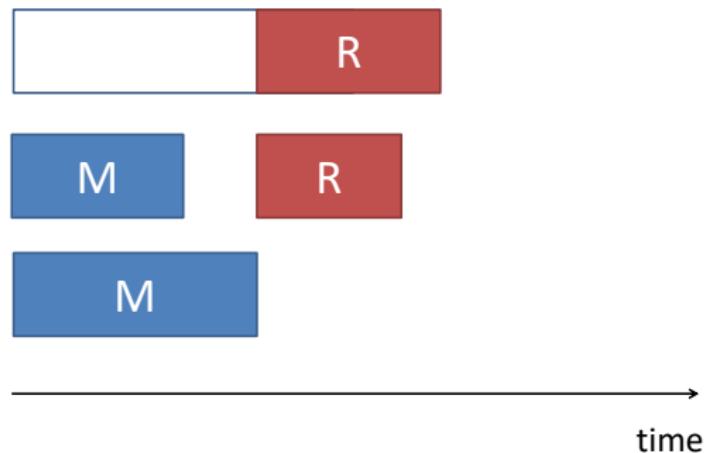


Early Dropping

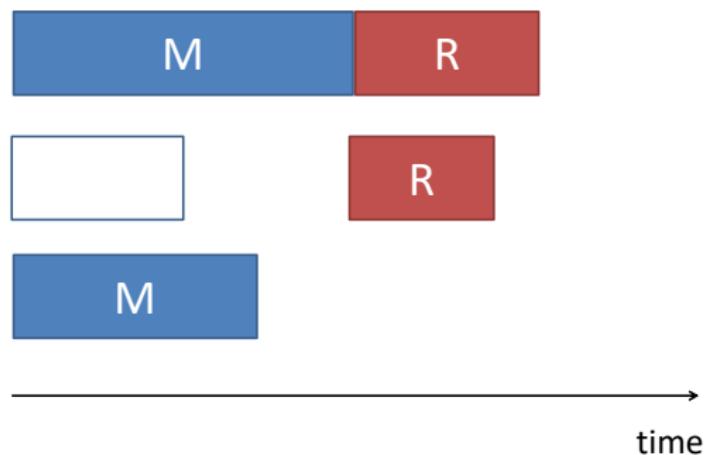
Data Drop: Task drop - Early



Data Drop: Task drop - Early



Data Drop: Task drop - Early

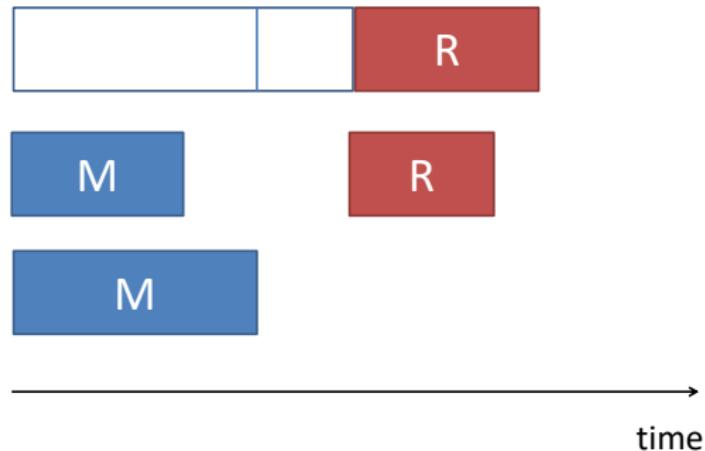


Straggler Dropping

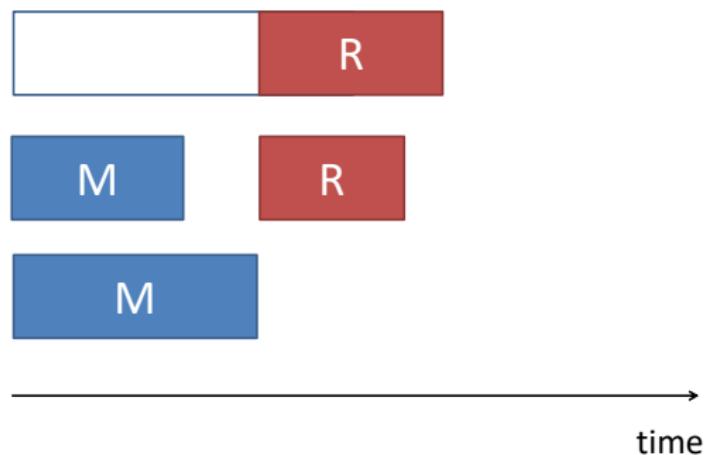
Data Drop: Task drop - Straggler



Data Drop: Task drop - Straggler

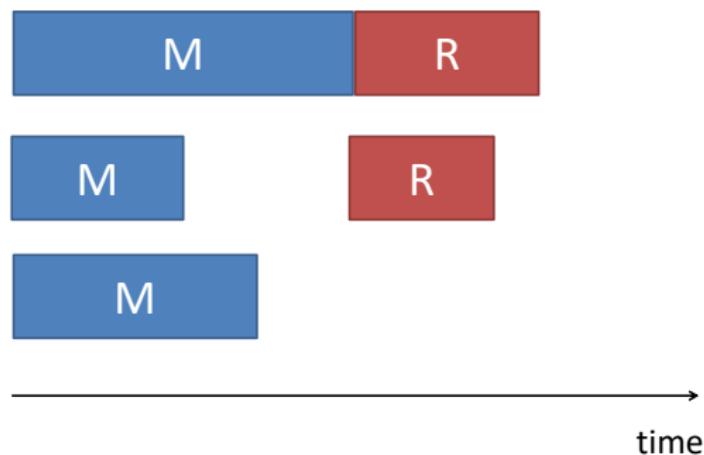


Data Drop: Task drop - Straggler

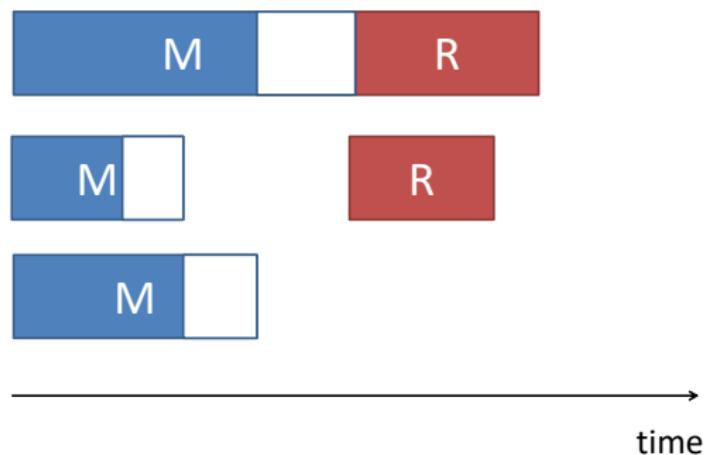


Input Sampling

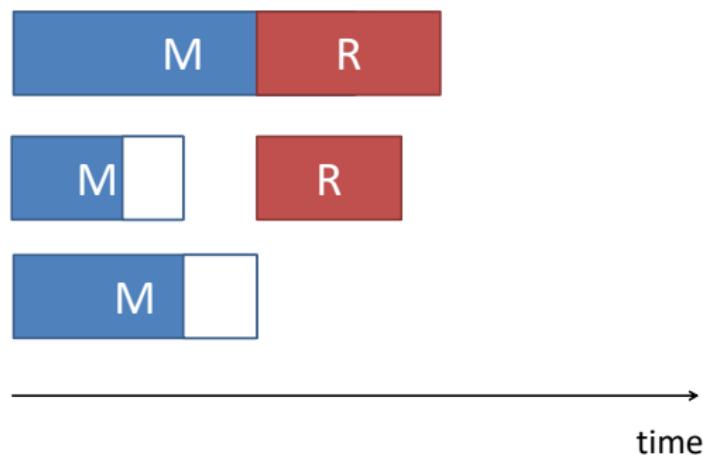
Data Drop: Input Sampling



Data Drop: Input Sampling



Data Drop: Input Sampling



Error

$$t_{\bar{N}_m - 1, 1 - \alpha/2} \sqrt{N_m \left(\left(\frac{1}{\theta_m} - 1 \right) s_u^2 + \left(\frac{1}{\eta_m} - 1 \right) \bar{s}_i^2 \right)},$$

Outline

1 Introduction

2 System Operation

3 The Model

4 Experimental assessment

5 Wrap-up

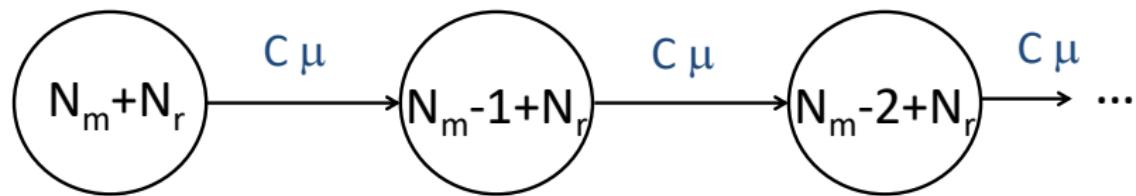
Basic Operation - Main parameters

- C computing slots
- N_m : number of map tasks
- N_r : number of reduce tasks
- $1/\mu_m$: mean map task execution time
- $1/\mu_r$: mean reduce task execution time

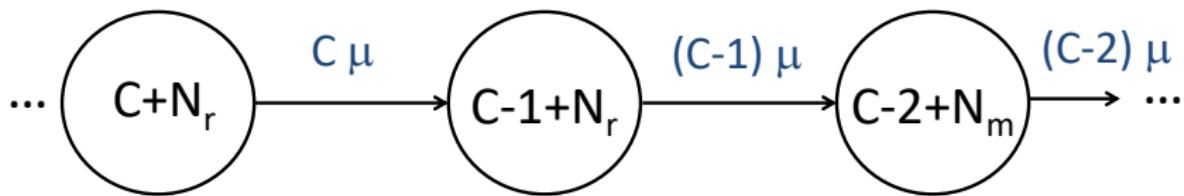
Basic Operation - Job Service Time

- Job Service Time: **Phase-type** distribution
- Phase: number of tasks to complete

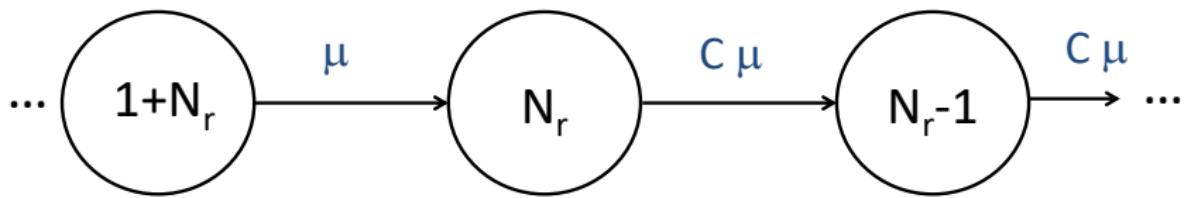
Basic Operation - Job Service Time



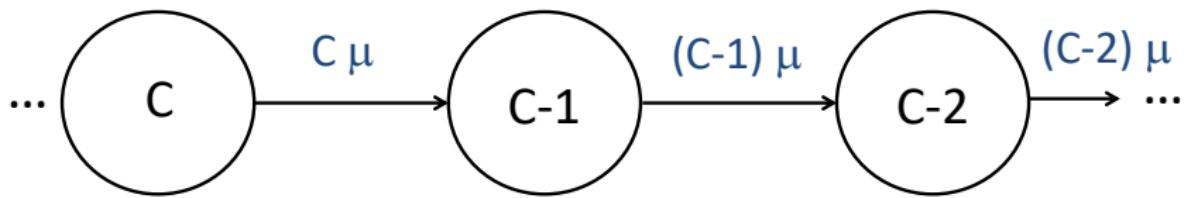
Basic Operation - Job Service Time



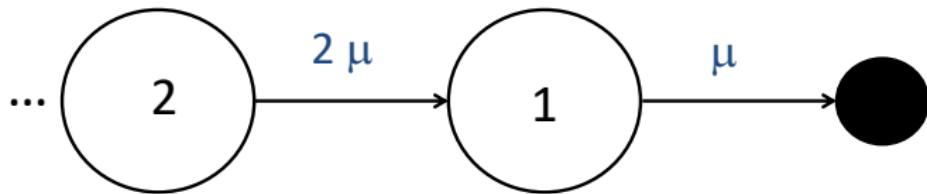
Basic Operation - Job Service Time



Basic Operation - Job Service Time



Basic Operation - Job Service Time



Modeling Approximate Computation

Introducing Early Task Dropping

- Execute fraction θ_m map tasks only
- Drop fraction $1 - \theta_m$ map tasks
- Effective number of map tasks: $\bar{N}_m = \lceil N_m \theta_m \rceil$
- Replace N_m by \bar{N}_m

Introducing Input Sampling

- Process fraction η_m of each map task data only
- Ignore fraction $1 - \eta_m$ of each map task data
- Effective mean map task execution time: $\frac{1}{\bar{\mu}_m} = \frac{\eta_m}{\mu_m}$
- Replace μ_m by $\bar{\mu}_m$

Early Dropping + Input Sampling

Transition Rates:

$$f(n) = \begin{cases} C\bar{\mu}_m, & \bar{N}_r + C < n \leq \bar{N}_r + \bar{N}_m, \\ (n - \bar{N}_r)\bar{\mu}_m, & \bar{N}_r < n \leq \bar{N}_r + C, \\ C\bar{\mu}_r, & C < n \leq \bar{N}_r, \\ n\bar{\mu}_r, & 0 < n \leq C. \end{cases}$$

Straggler Dropping

- Start with $N_m + N_r$ tasks
- Decrease one by one until $N_m - (\bar{N}_m - 1) + N_r$ are left
- Jump down to N_r
- Effective number of map tasks executed: \bar{N}_m

Straggler Dropping + Input Sampling

Transition Rates:

$$f(n) = \begin{cases} \min\{n - N_r, C\}\bar{\mu}_m, & N_r + N_m - \bar{N}_m + 1 \leq n \\ & \leq N_r + N_m, \\ \min\{n, C\}\bar{\mu}_r, & N_r - \bar{N}_r + 1 \leq n \leq N_r. \end{cases}$$

Analysis

Analysis

- FCFS scheduling
- Job service time: phase-type (PH) - (α, \mathbf{G})
- Job arrival process: Markovian arrival process (MAP)
- MAP/PH/1 queue
- Sengupta's paper: age-based analysis
- $(X(t), A(t), S(t))$: age of job in service, arrival process phase, service phase

Analysis

Solve a matrix-integral equation

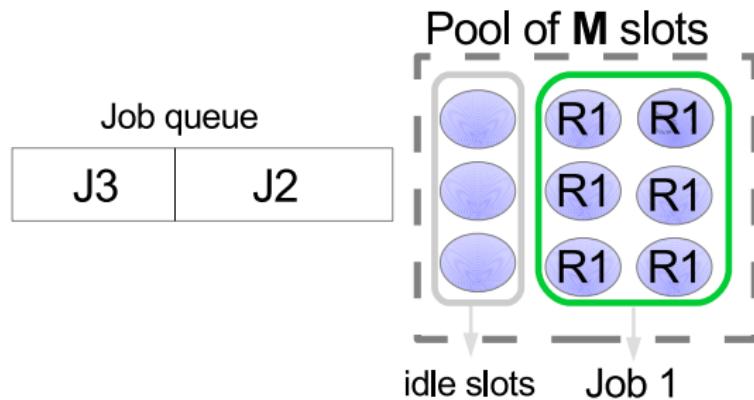
$$\mathbf{T} = \mathbf{I}_d \otimes \mathbf{G} + \int_0^{\infty} \exp(\mathbf{T}x) \left(\exp(\mathbf{D}_0x) \mathbf{D}_1 \otimes \tilde{\mathbf{G}} \right) dx,$$

- Obtain waiting time distribution (β_w, \mathbf{B}_w)
- Known service-time distribution (α, \mathbf{G})
- Obtain response time distribution (β_r, \mathbf{B}_r) :

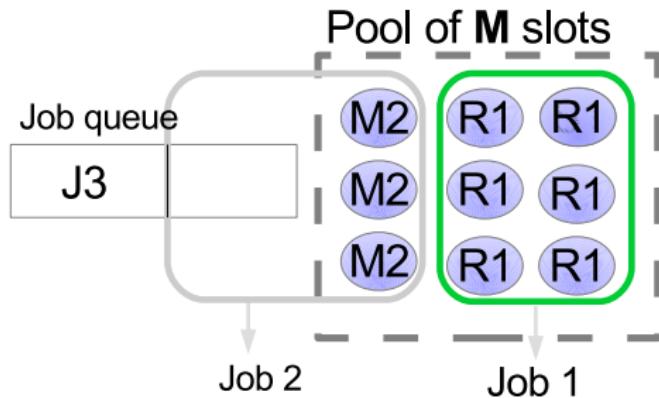
$$\beta_r = [\phi \beta_w \quad (1 - \phi)\alpha], \quad \mathbf{B}_r = \begin{bmatrix} \mathbf{B}_w & (-\mathbf{B}_w \mathbf{1})\alpha \\ \mathbf{0} & \mathbf{G} \end{bmatrix},$$

Overlapping Execution

FCFS Scheduling - Non-overlapping



Overlapping scheduling



- Allow jobs to start Map phase as soon as slots become available.
- At most 2 jobs in service: one in Map phase, one in Reduce phase.

Overlapping Scheduling

- No longer a MAP/PH/1 queue
- Service times depend on system state
- High load: high chance back-to-back jobs
- Low load: similar to non-overlapping

Overlapping Scheduling

- Generalize Sengupta's analysis
- $(X(t), A(t), M(t), R(t), Y(t))$
- $M(t)$: number of remaining map tasks of oldest job
- $R(t)$: number of remaining reduce tasks of oldest job
- $Y(t)$: number of remaining map tasks of youngest job

Overlapping Scheduling

Three phases:

M one job, map phase

R one job, reduce phase

R/M two jobs, oldest in reduce phase, youngest in map phase

Overlapping Scheduling

(Steady-state) Service time distribution: $\text{PH}(\beta_s, B_s)$

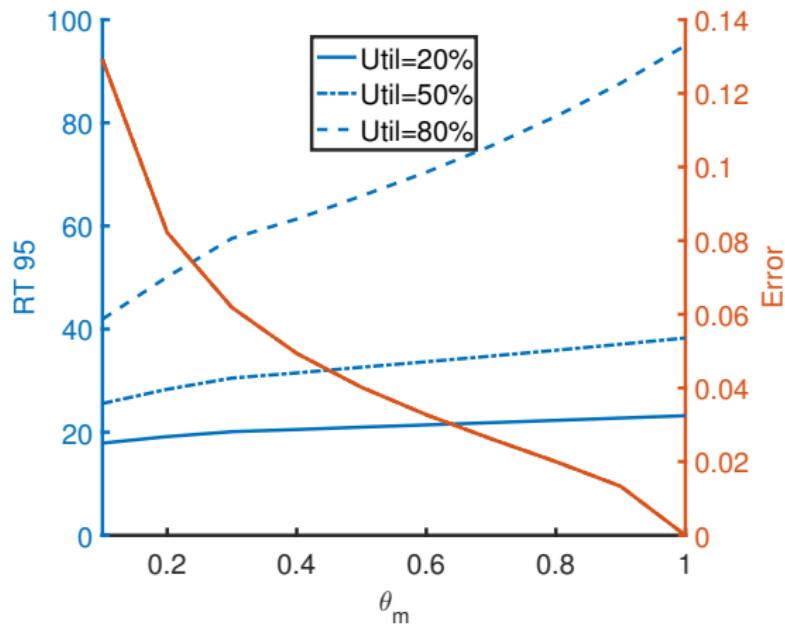
$$B_s = \begin{bmatrix} Q^{R/M, R/M} & Q^{R/M, M} & Q^{R/M, R^*} \\ 0 & Q^{M, M} & Q^{M, R^*} \\ 0 & 0 & B^{R^*, R^*} \end{bmatrix}, \quad (1)$$

$$\beta_s = [\beta_s^{R/M} \quad \beta_s^M \quad \beta_s^{R^*}]. \quad (2)$$

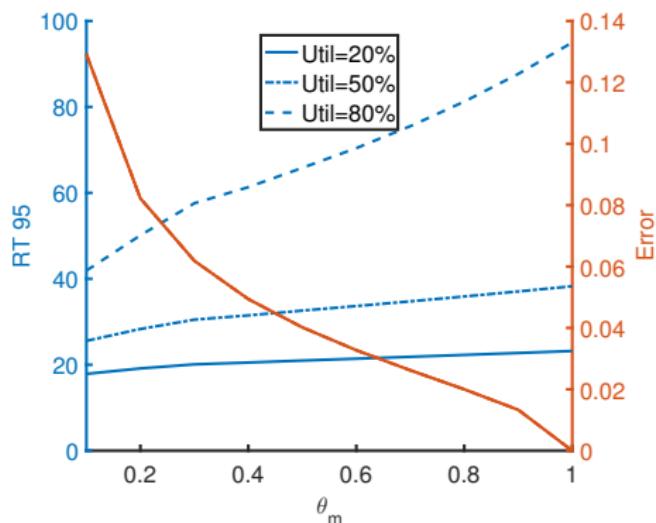
Outline

- 1 Introduction
- 2 System Operation
- 3 The Model
- 4 Experimental assessment
- 5 Wrap-up

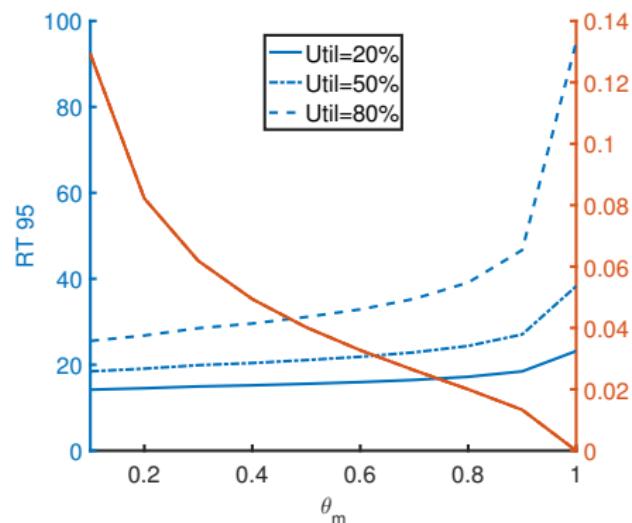
FCFS - Early Dropping



FCFS - Early vs Straggler Dropping

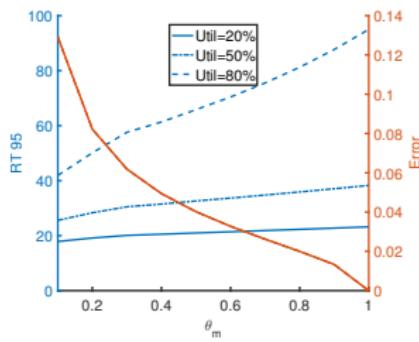


(a) Early dropping

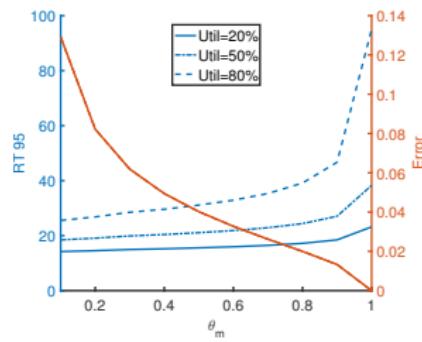


(b) Straggler dropping

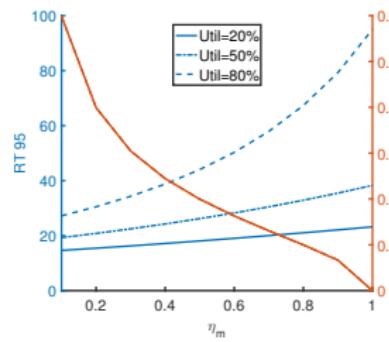
FCFS - Early vs Straggler Dropping vs Input Sampling



(a) Early dropping



(b) Straggler dropping



(c) Input sampling

FCFS - Early vs Straggler Dropping vs Input Sampling

Observation 1

Best Mechanisms are:

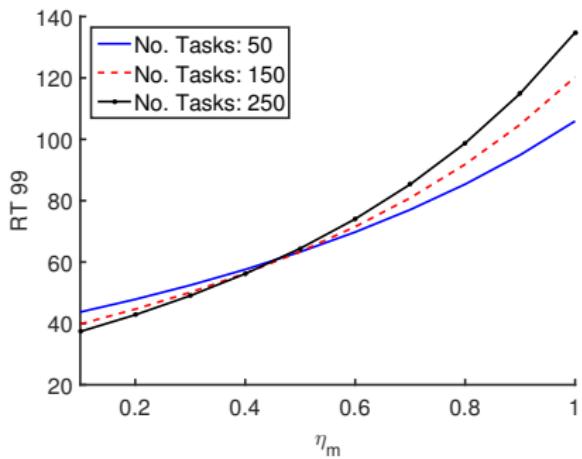
- ① Straggler dropping
- ② Input sampling
- ③ Early dropping

Observation 2

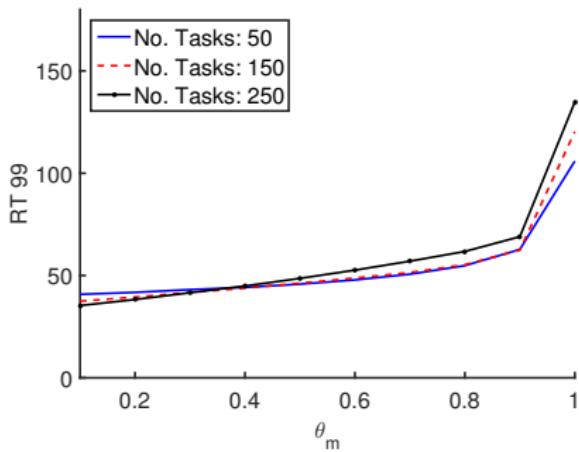
An $x\%$ data/task dropping removes

- ① $x\%$ **longest** tasks (straggler dropping)
- ② $x\%$ **all** tasks, including longest (input sampling)
- ③ $x\%$ tasks (early dropping)

Varying number of tasks

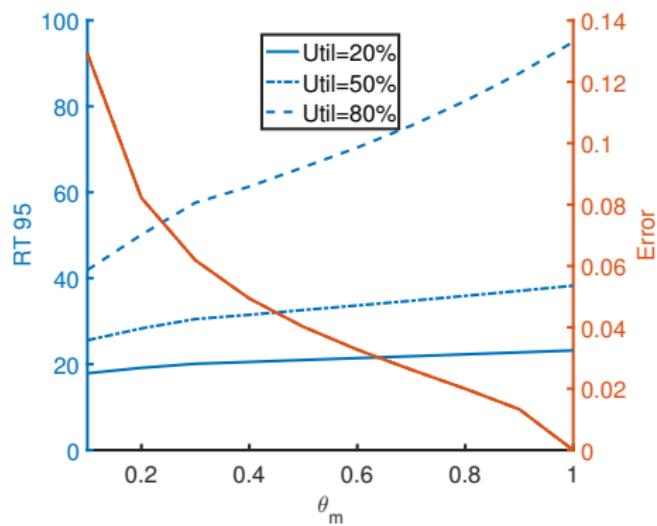


(a) Input sampling

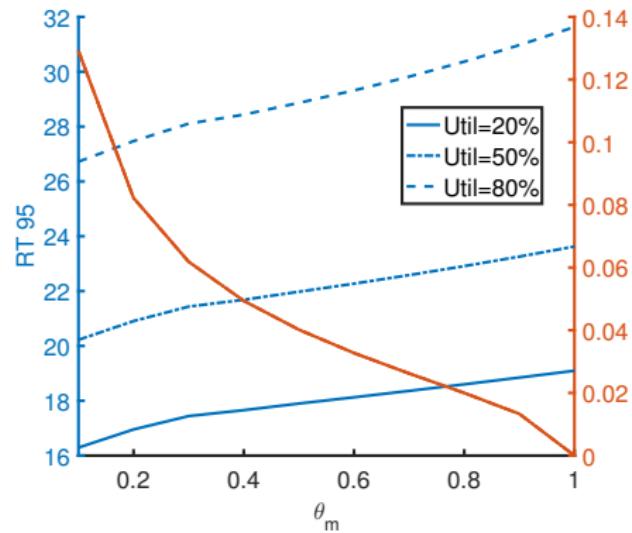


(b) Straggler dropping

FCFS vs. Overlapping (Early dropping)

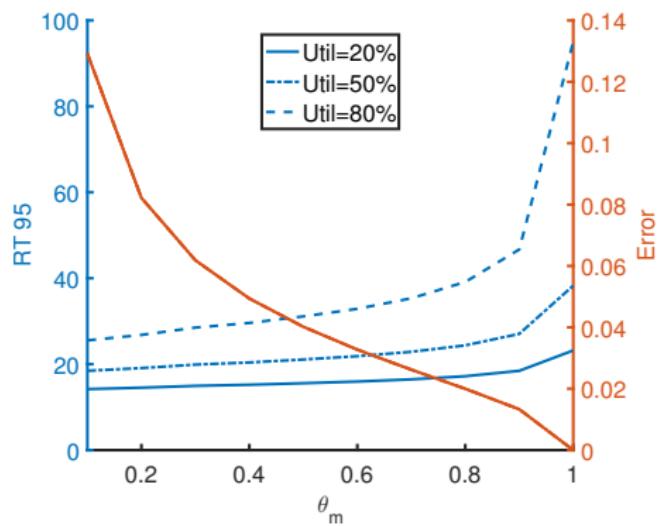


(a) FCFS

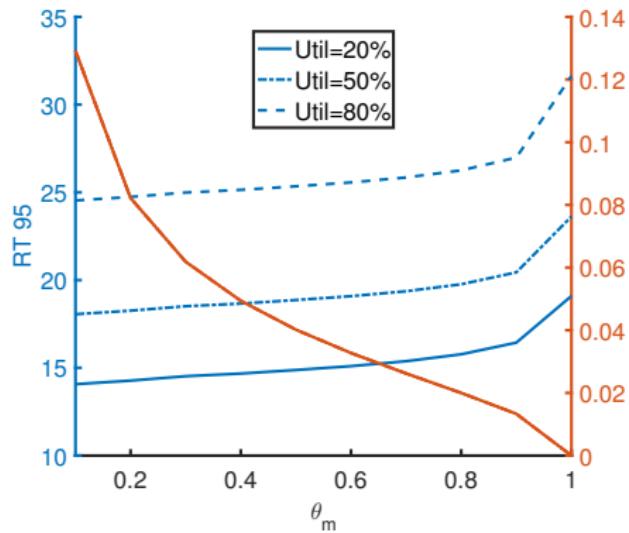


(b) Overlapping

FCFS vs. Overlapping (Straggler dropping)

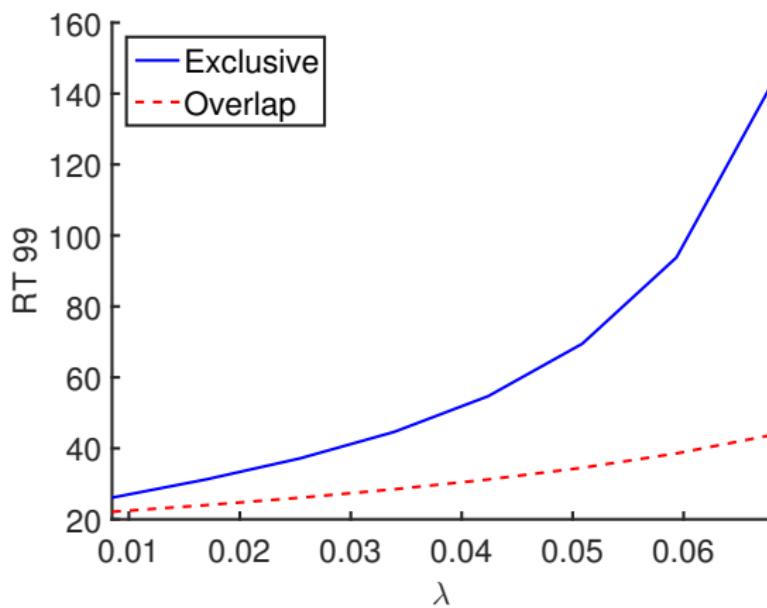


(a) FCFS



(b) Overlapping

FCFS vs. Overlapping (No approximation)



Outline

1 Introduction

2 System Operation

3 The Model

4 Experimental assessment

5 Wrap-up

Wrap-up

Model captures

- Straggler Dropping > Input Sampling > Early Dropping
- Overlapping > FCFS Scheduling
- Latency - Accuracy Tradeoff

Given a:

- **Latency Budget:** what accuracy can be expected/offered?

Wrap-up

Model captures

- Straggler Dropping > Input Sampling > Early Dropping
- Overlapping > FCFS Scheduling
- Latency - Accuracy Tradeoff

Given a:

- **Latency Budget:** what accuracy can be expected/offered?
- **Accuracy Budget:** what latency can be expected/offered?

Wrap-up

Model captures

- Straggler Dropping > Input Sampling > Early Dropping
- Overlapping > FCFS Scheduling
- Latency - Accuracy Tradeoff

Given a:

- **Latency Budget:** what accuracy can be expected/offered?
- **Accuracy Budget:** what latency can be expected/offered?

Wrap-up

Model captures

- Straggler Dropping > Input Sampling > Early Dropping
- Overlapping > FCFS Scheduling
- Latency - Accuracy Tradeoff

Given a:

- **Latency Budget:** what accuracy can be expected/offered?
- **Accuracy Budget:** what latency can be expected/offered?

Thank you!