# USC University of Southern California

# Scheduling Jobs Across Geo-distributed Datacenters
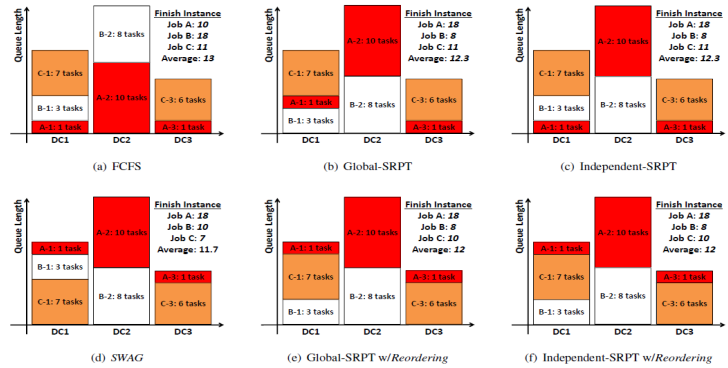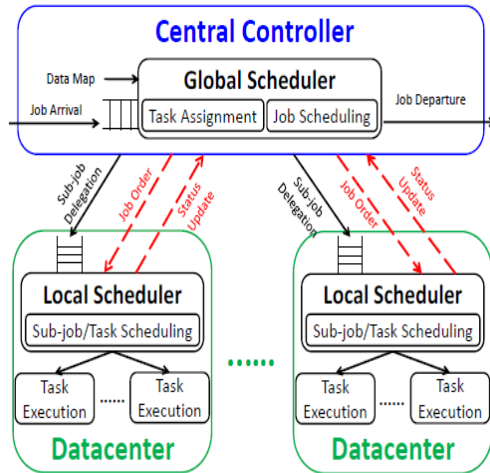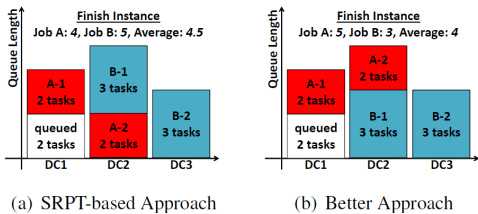## Chien-Chun Hung, Leana Golubchik, Minlan Yu

## Motivation and System Architecture

- Tasks of the jobs are distributed across the datacenters for data locality to save bandwidth and completion time.
- The imbalance in tasks distribution and the workloads at each datacenters necessitate new scheduling techniques.



(a) FCFS  (b) Global-SRPT  (c) Independent-SRPT
(d) SWAG  (e) Global-SRPT w/Reordering  (f) Independent-SRPT w/Reordering

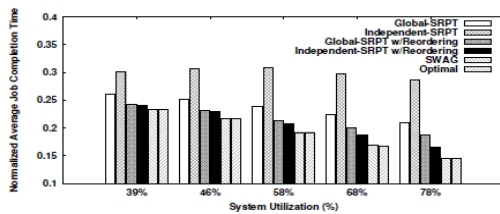| Job ID | Arrival Sequence | #Remaining Tasks in DC1 | #Remaining Tasks in DC2 | #Remaining Tasks in DC3 | Total #Remaining Tasks |
|---|---|---|---|---|---|
| A | 1 | 1 | 10 | 1 | 12 |
| B | 2 | 3 | 8 | 0 | 11 |
| C | 3 | 7 | 0 | 6 | 13 |



(a) SRPT-based Approach  (b) Better Approach

### *Reordering*-based Scheduling Approach
✓ Serve as a post-processing adjustment to improve any scheduling results.
✓ Yield the resources to other tasks if not hurting its job's overall completion time.
✓ Provably do no harm to the average job completion time for any job scheduling.

### *Workload-Aware Greedy Scheduling* (SWAG)
✓ A generic scheduling solution that computes the job order for all the jobs.
✓ Prioritize the jobs based on estimated finish times along with current workload.
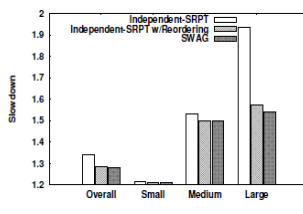✓ Greedily schedule the job that can finish quickly across all the datacenters.

## Performance Improvements



## Fairness



## Performance Sensitivity



★ SWAG and Reordering result in a significant performance improvement, up to **50%** and **30%** respectively, over SRPT-based scheduling.

★ SWAG and Reordering improve average job completion time while maintaining reasonable fairness, even for the large jobs, compared to SRPT-based scheduling.
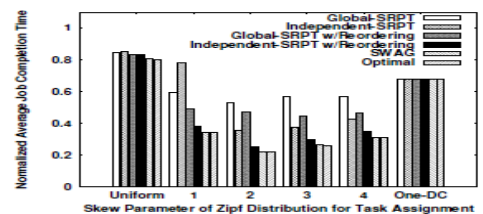
★ The biggest improvements are observed when the system is highly-loaded or there exists a high skew in workload, either in job sizes or in task assignments.

★ Without workload skew or in lightly-loaded systems, SWAG and Reordering exhibit similar performance compared to SRPT-based scheduling.

## Summary and Extensions

- SWAG vs. Reordering
  - SWAG provides greater improvements with reasonable overhead.
  - Reordering is light-weight and easily added to any scheduling approach.
- Heterogeneous datacenter capacity (#slots)
- Scheduling jobs with DAG of tasks
- Flow scheduling for intermediate data shuffling