# Exploiting Bounded Staleness to Speed Up Big Data Analytics
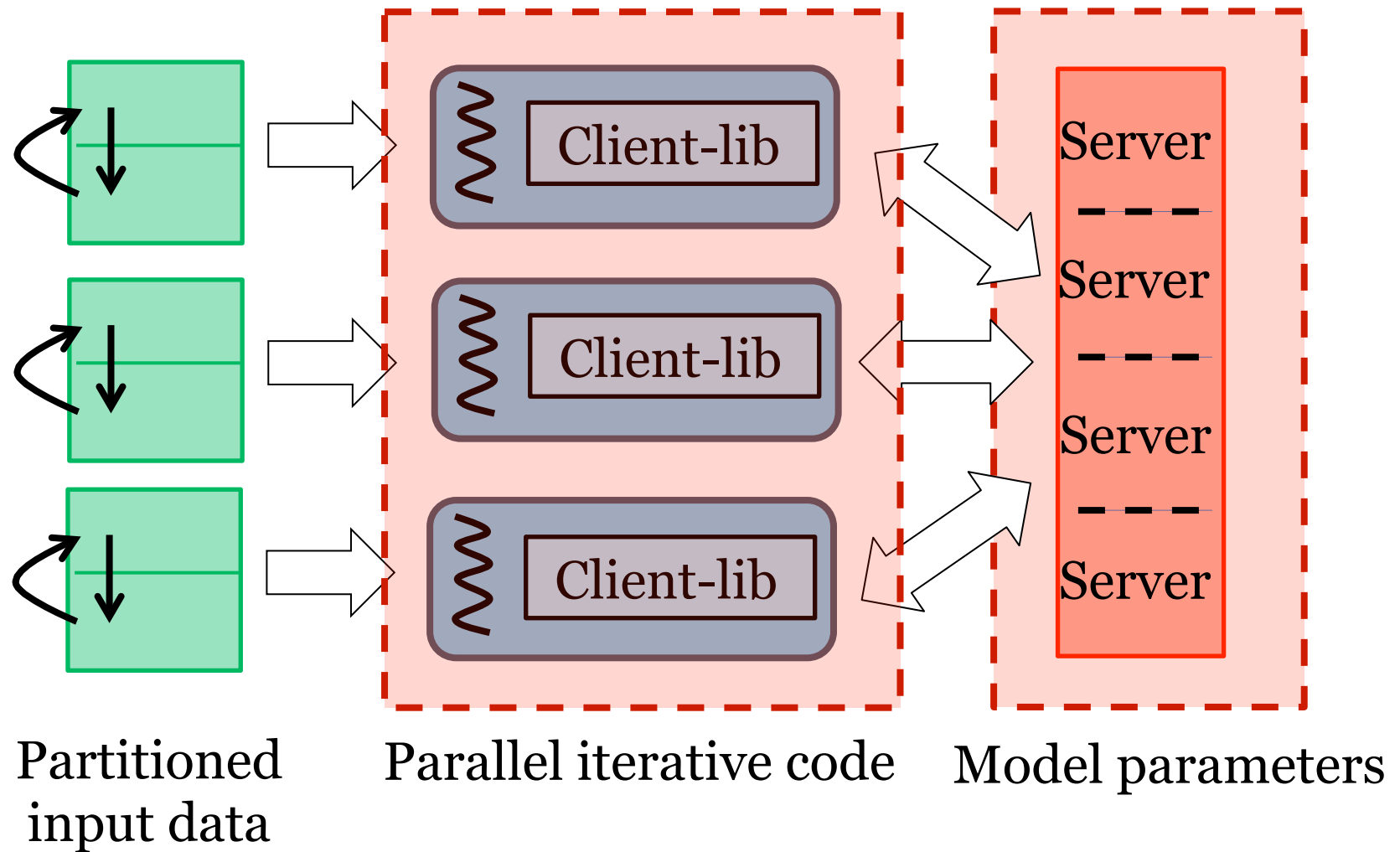
Garth Gibson & the BigLearning team
Carnegie Mellon University

http://www.istc-cc.cmu.edu/

Intel Science & Technology
Center for Cloud Computing

# Parallel ML Systems Architecture



Partitioned input data
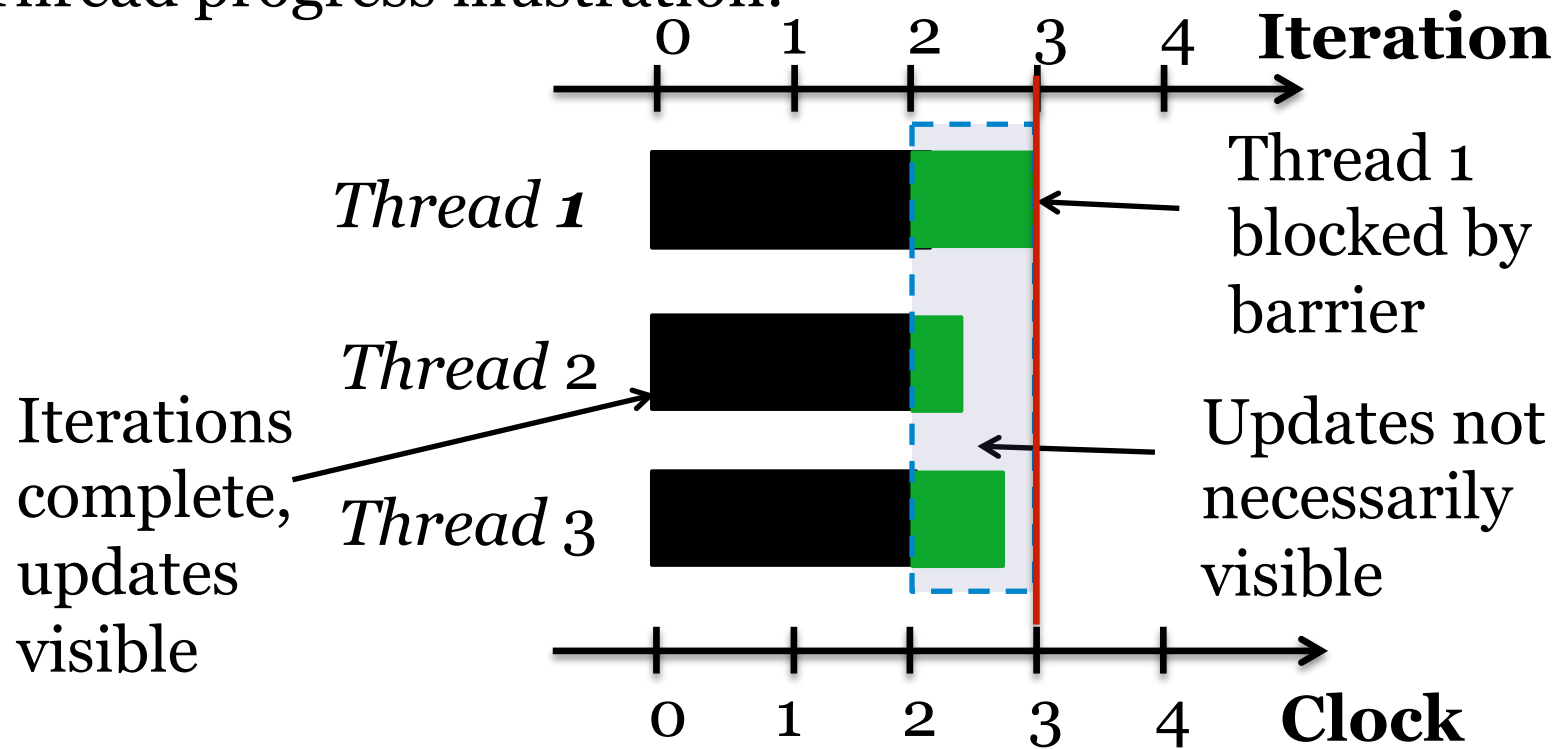
Parallel iterative code

Model parameters

# Agenda: Bound Staleness Project Suite

- Compare Bounded Async Bulk Synch Parallel (A-BSP) vs Stale Synch Parallel (SSP)
- Repetition-exploiting optimizations (to BSP)
- Managed (extra) Bandwidth SSP (MBSSP)
- Convergence-guided Scheduling (STRADS)
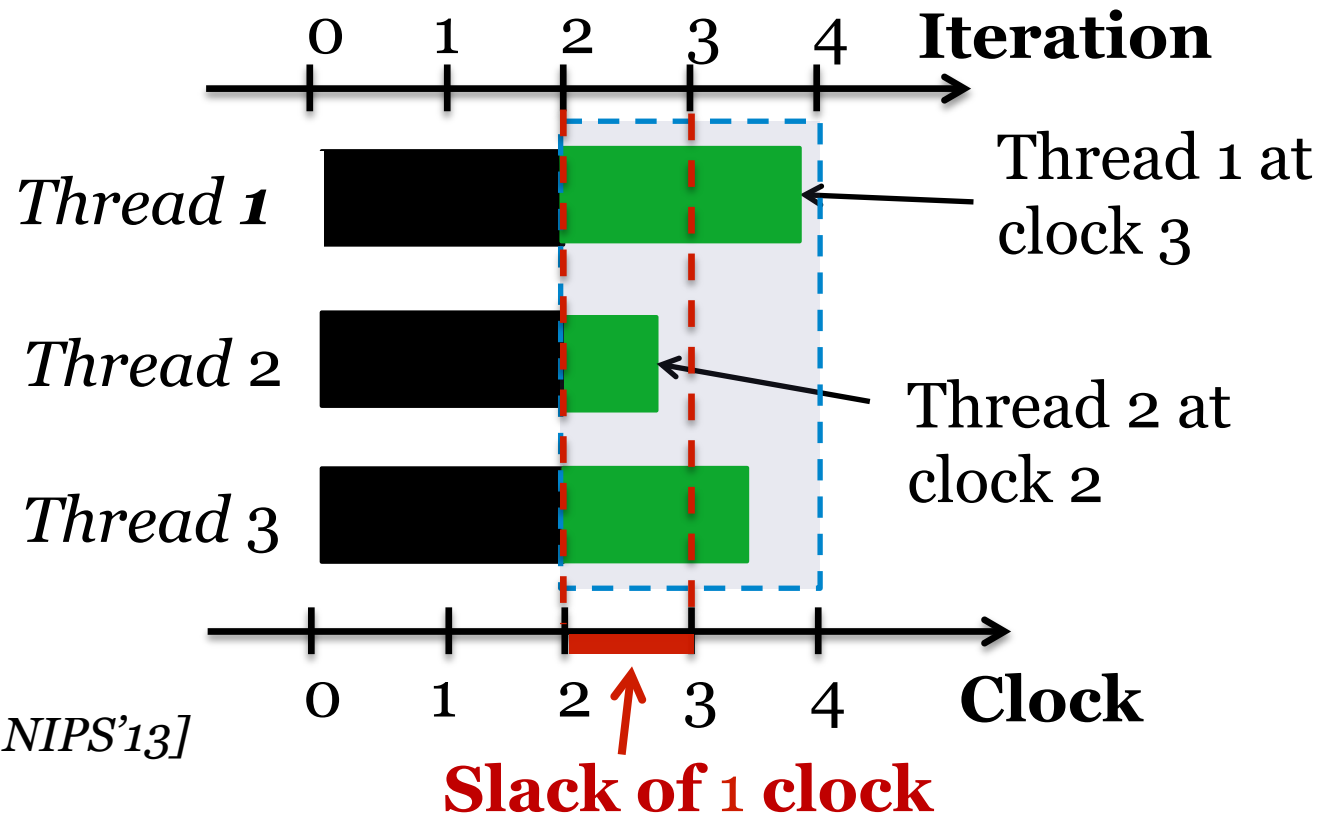
# Bulk Synchronous Parallel

- A barrier every (logical) clock
  - chunk of work, often 1 iteration on all input data

Thread progress illustration:



Thread 1 blocked by barrier

Iterations complete, updates visible

Updates not necessarily visible
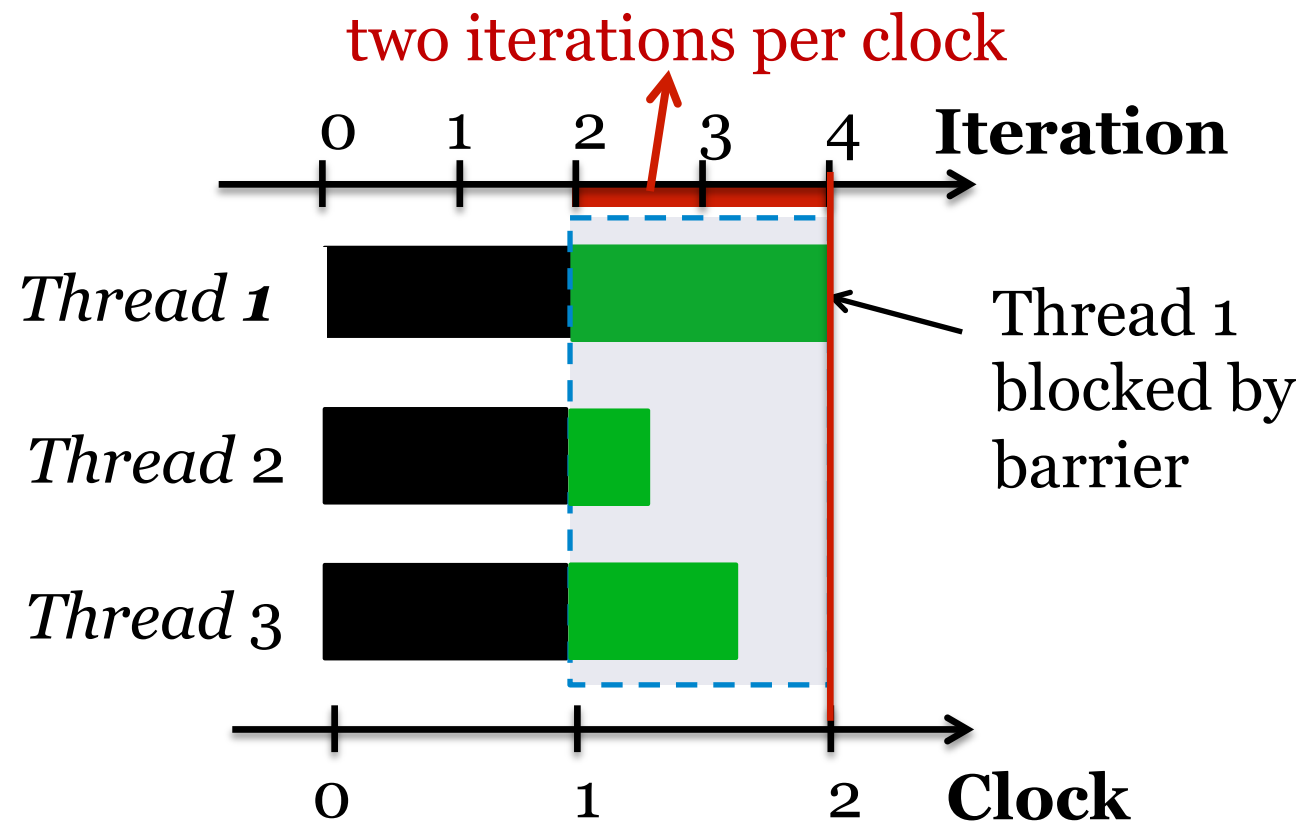
# Stale Synchronous Parallel (SSP)

- Threads allowed to be slack clocks ahead of slowest thread, possibly reading stale data



*[HotOS'13, NIPS'13]*

# Arbitrarily-sized BSP (A-BSP)

▫ Work in each clock can be more than one iteration
  - Less synchronization overhead (bounded asynch)



two iterations per clock

Iteration
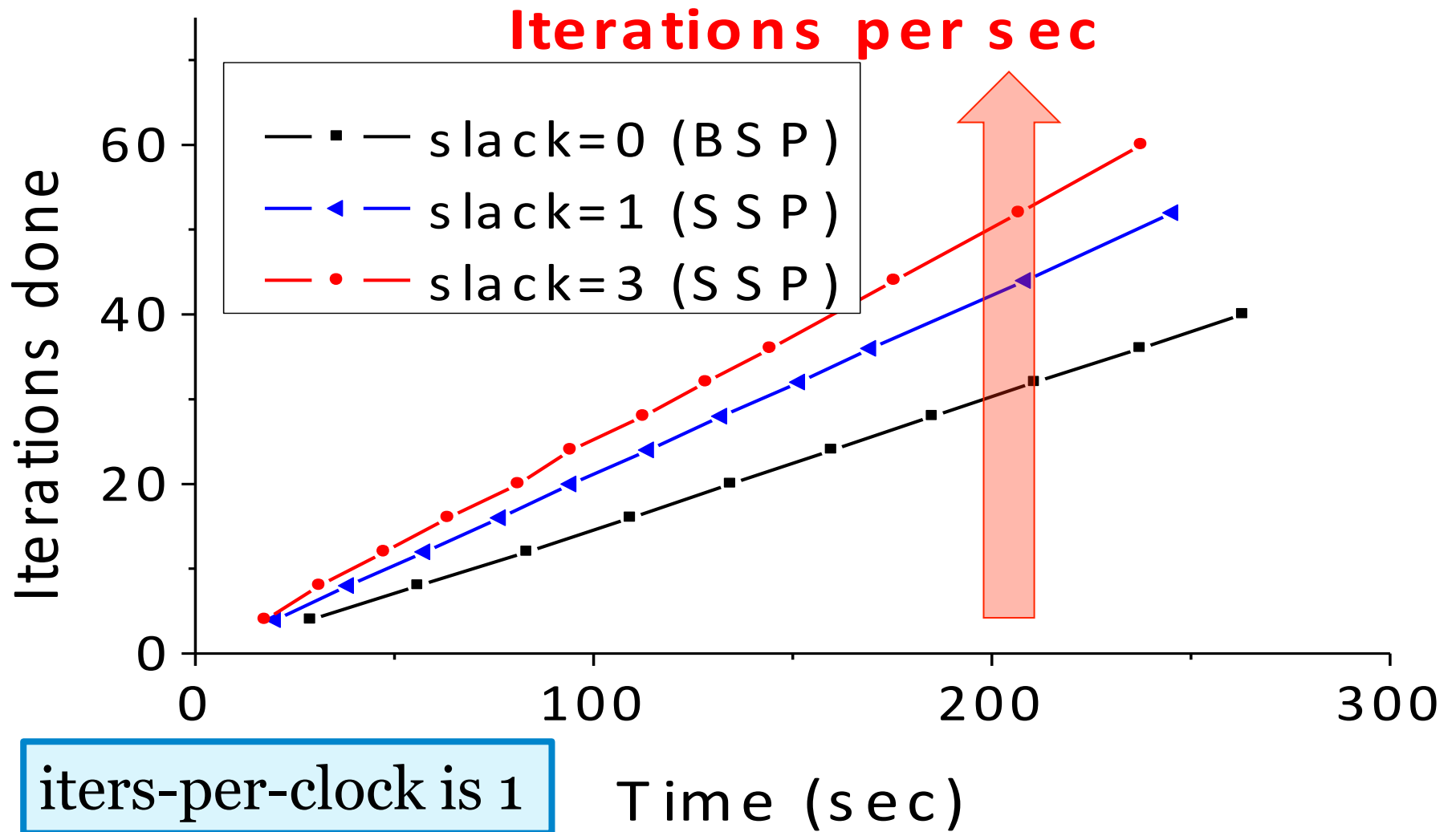
Thread 1 blocked by barrier

Clock

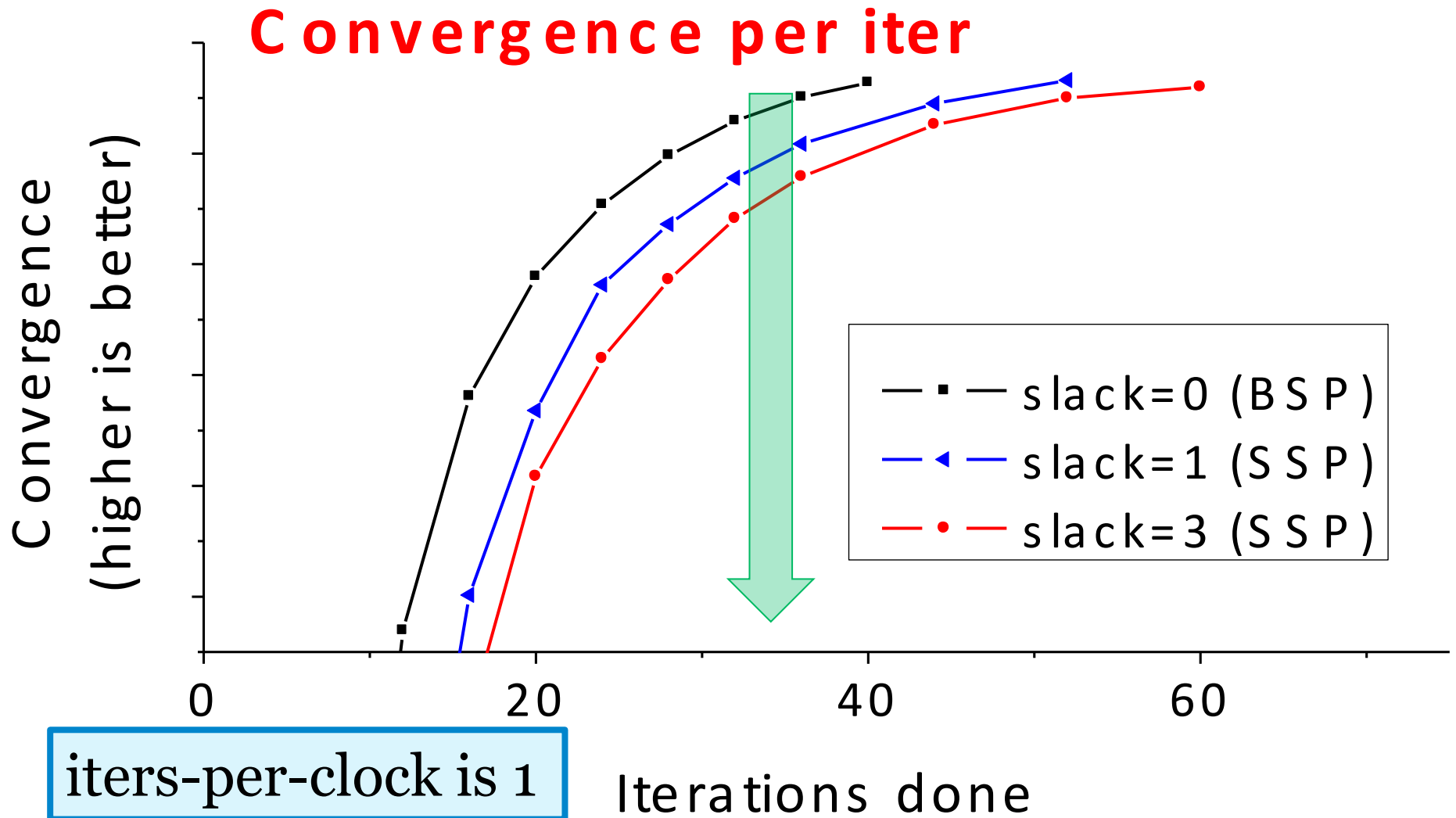  - A-BSP is SSP with a slack of zero

# Application Benchmark Example

- Topic Modeling
  - Algorithm: Gibbs Sampling on LDA
  - Input: *NY Times* dataset
    - 300k docs, 100m words, 100k vocabulary
  - Solution quality criterion: Loglikelihood
    - How likely the model generates observed data
    - Becomes higher as the algorithm converges
    - A larger value indicates better quality

- Hardware information
  - 8 machines, each with 64 cores & 128GB RAM
- Basic configuration
  - One client & tablet server per machine
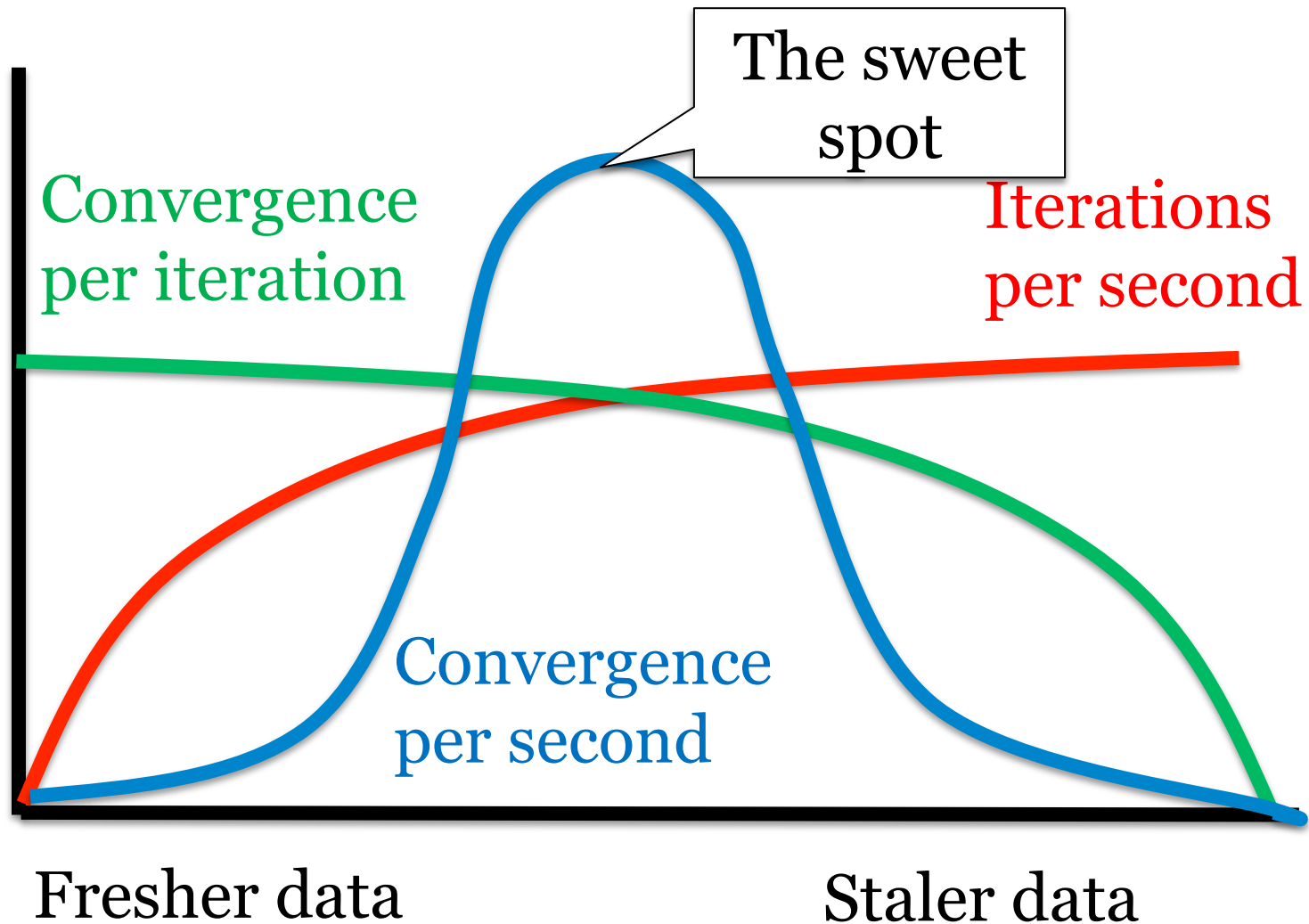  - One computation thread per core

Convergence per iter

Convergence (higher is better)

slack=0 (BSP)
slack=1 (SSP)
slack=3 (SSP)

Iterations done

iters-per-clock is 1

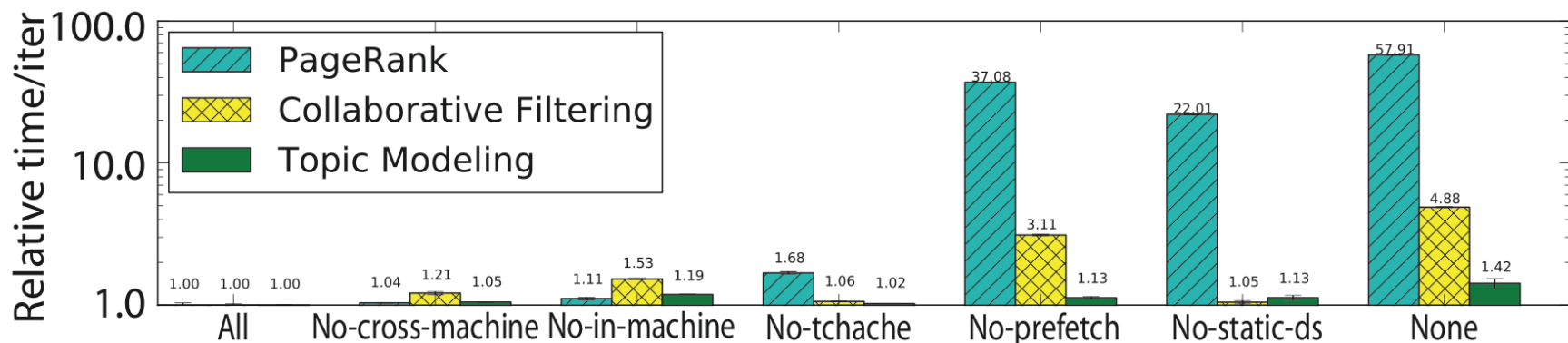[ATC'14]

- Iterative code often very repetitive – exploit!
  - Virtual iteration
- Affinity allocation, static & precomputed policies, multiple levels of cache, update prefetching
- 

*[under submission]*

Lead: Henggang Cui



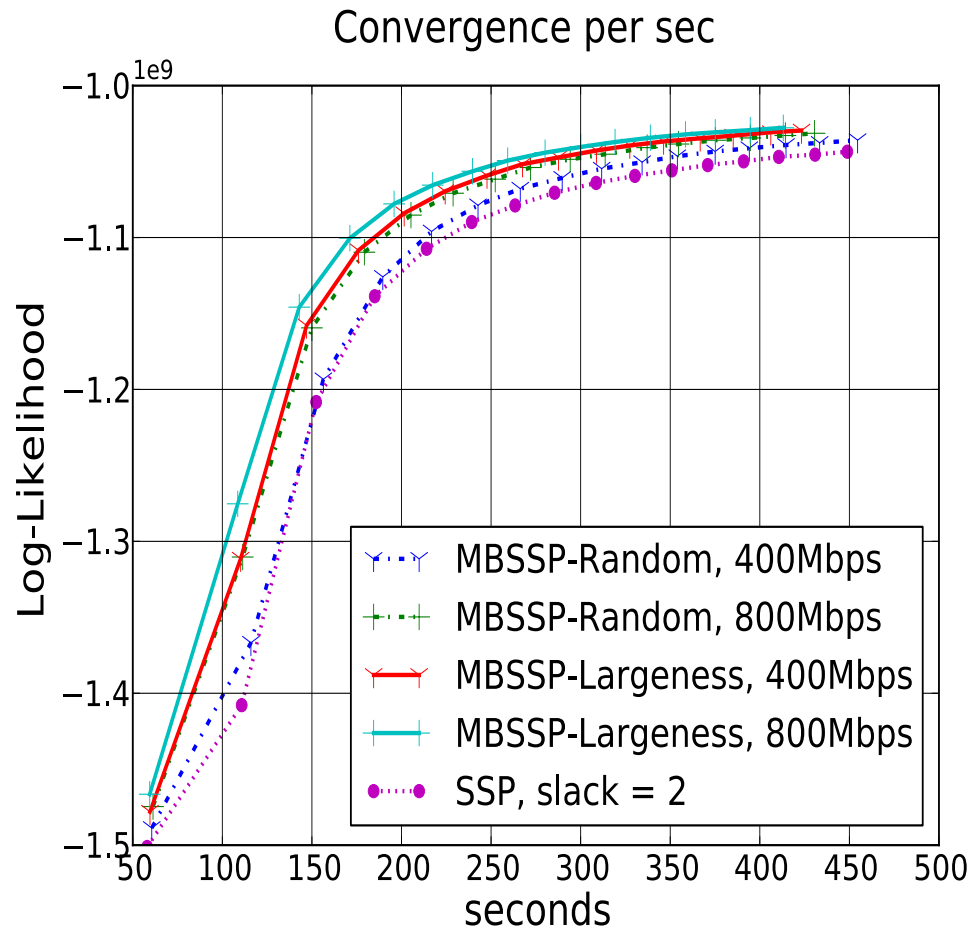- **Optimization effectiveness break-down:**

# Managed Bandwidth SSP (MBSSP)

- In SSP, communication and computation are overlapped, but every update is treated equally
- But not every update is equally important to convergence (e.g. small vs. large deltas)
- MBSSP exploits network bandwidth not fully utilized to transmit pending updates sooner
- Early transmissions may speed convergence
  - And may allow greater staleness (latency hiding)
- What to send early?  Random vs delta ordered

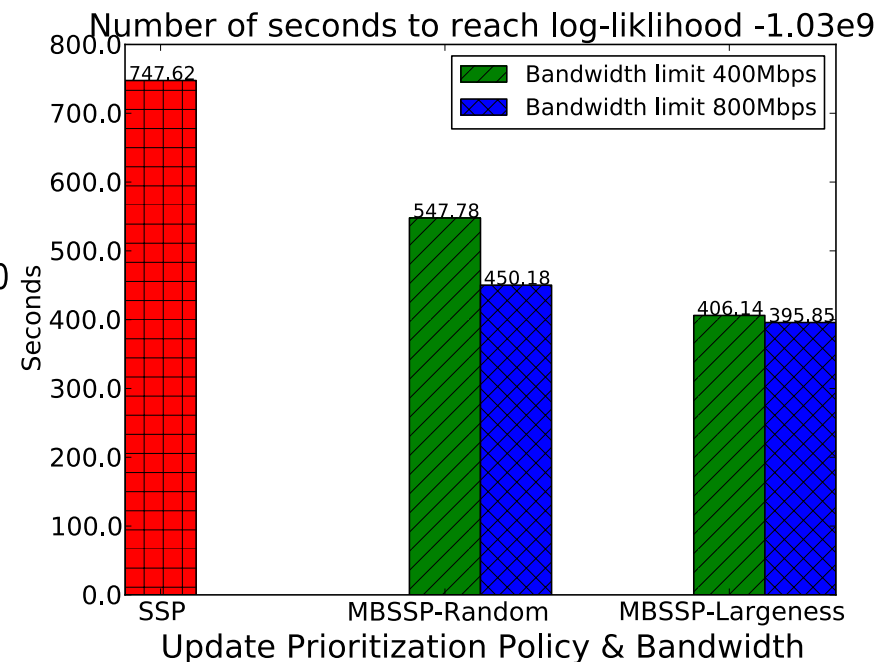- Leads: Jinliang Wei, Wei Dai

*[under submission]*

# Absolute Convergence Improved 40%

## Convergence per sec



LDA (Gibbs Sampling)
NYT Dataset
8x64 core nodes
1GE network
Fits in memory

Early transmission reduces
time needed to converge

Delta-importance-ordered
achieves as much benefit as
random early send with half
the extra bandwidth

Number of seconds to reach log-liklihood -1.03e9

# MBSSP Vision

- It is beneficial to send out early model refinements even with bounded bandwidth.

- Early communication improves convergence enabling much larger staleness (latency hiding).

- Application-specific policies for preferring model refinements can make a big difference.
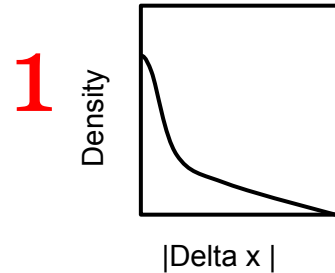
# STRADS: Up Stack to ML Scheduling

- **Uniform parameter update is not optimal**
  - Use deeper knowledge of ML algorithms to update parameters at different rates for best convergence speed (like MBSSP)

- **Random parameter selection for parallel update risks divergence (e.g. Shotgun Lasso )**
  - Control errors when selecting parameters to update in parallel

- **Leads: Jin Kyu Kim, Seunghak Lee**

Weight of model parameters
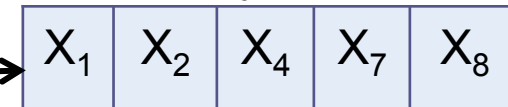|Delta of x|

| | |
|---|---|
| $\Delta x_1$ | 0.11 |
| $\Delta x_2$ | 0.01 |
| $\Delta x_3$ | 0.003 |
| $\Delta x_4$ | 0.15 |
| $\Delta x_5$ | 0.0001 |
| $\Delta x_6$ | 0.001 |
| $\Delta x_7$ | 0.07 |
| $\Delta x_8$ | 0.0003 |
| $\Delta x_9$ | 0 |

Sampling based
on delta distribution

**1**

Density

|Delta x |

One update set
in ready queue

| $X_1$ | $X_2$ | $X_4$ | $X_7$ | $X_8$ |
|---|---|---|---|---|

Dependency
Checker & Filter

**2**

.......

Workers in remote machines

New delta info

*[arXiv(1406.4580)'14]*

16

# Benefits of Two Scheduling Policies



w/o scheduling ·······
w/t STRADS ——
70 cores ·······
120 cores ·······

w/o scheduling diverges beyond 80 cores

STRADS still converge with 120 cores
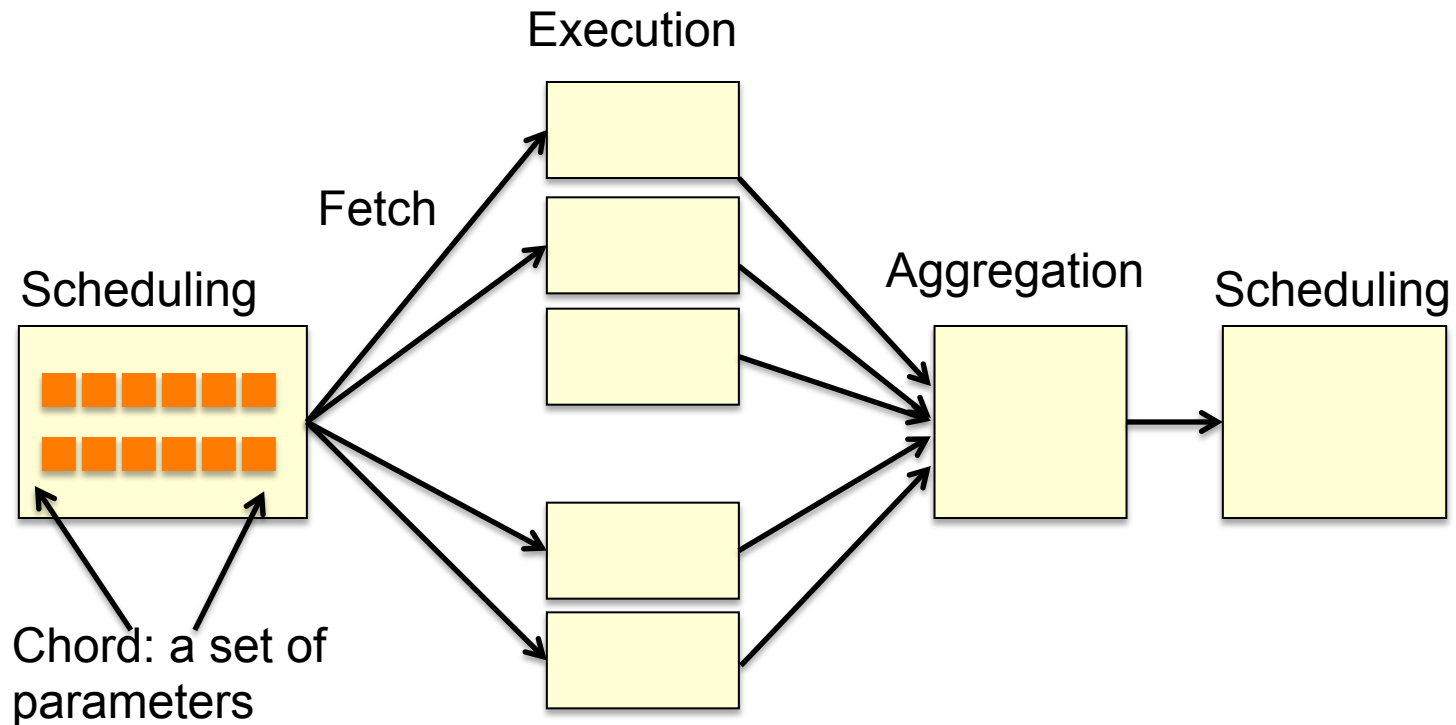
Application: Lasso

Synthetic data: 450 by 100K
- Parameters are highly correlated.

w/o scheduling:
Limit the degree of parallelism to 70 cores to avoid divergence

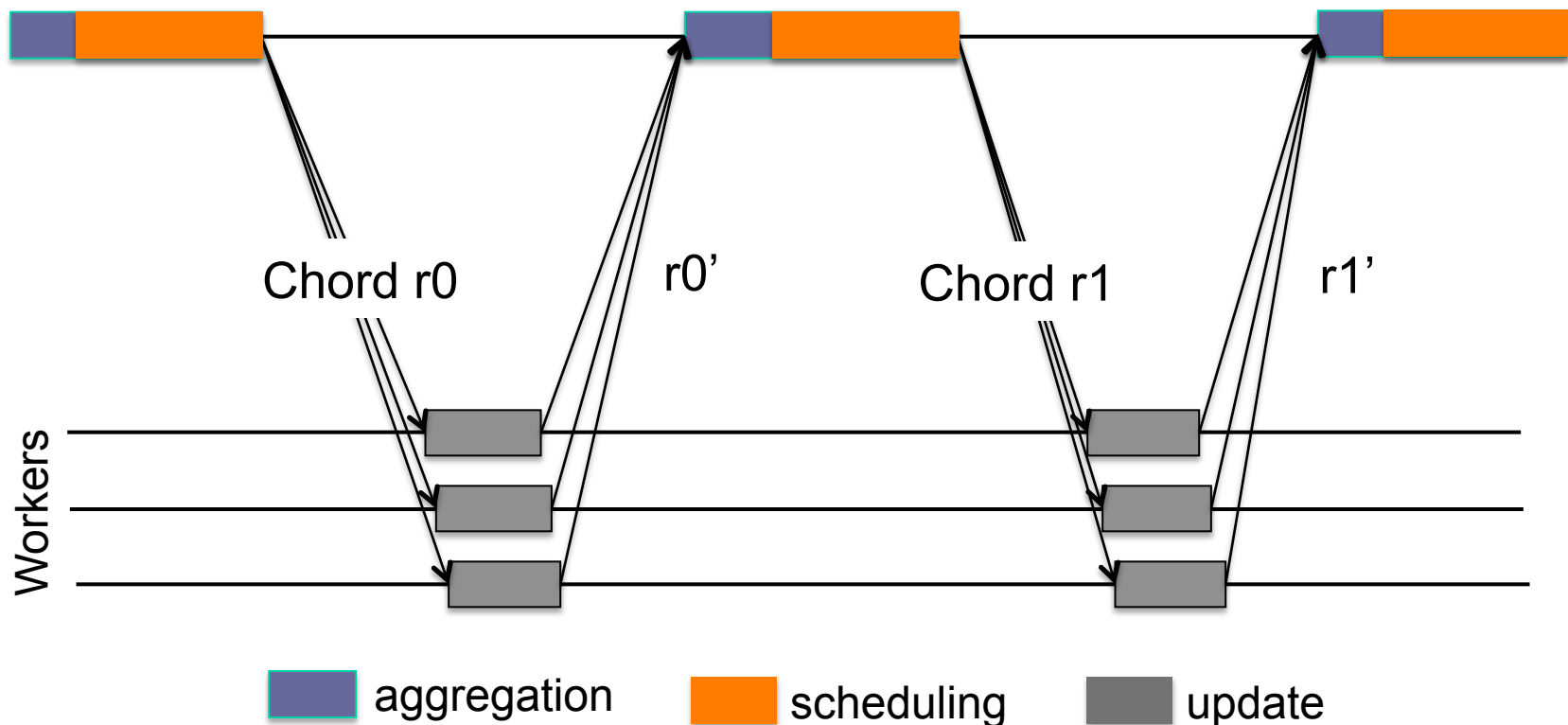## Scheduling/Fetch/Execution/Aggregation model



- Scheduling selects a chord to minimize aggregate errors of parallel update
- Parameters of a chord are selected to be approximately independent

# System Issue: Pipeline Scheduling

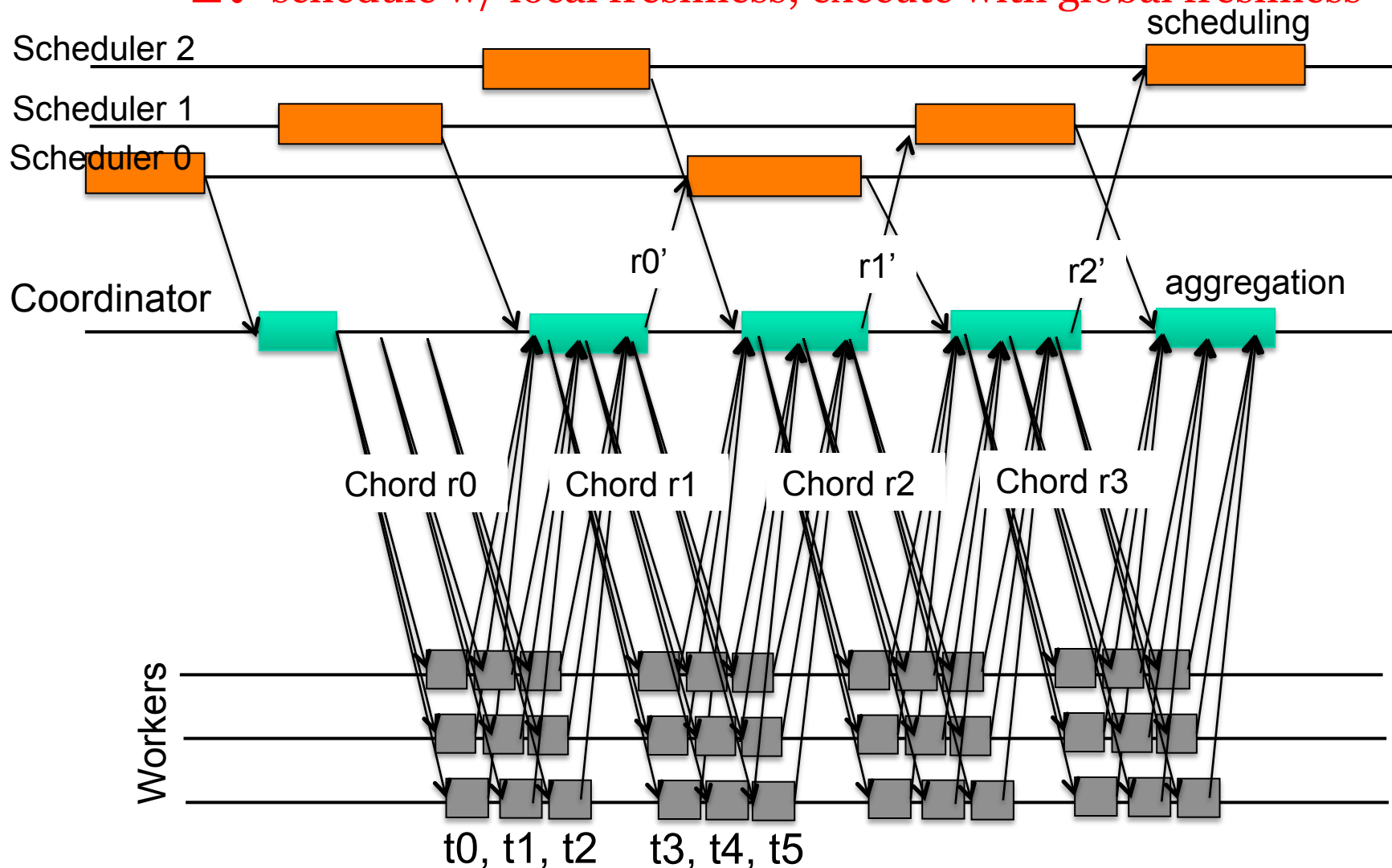Serial execution of chords is a performance bottleneck



Approach: Make scheduling decisions with latest data only for the scheduler's partition of the (big) model parameters
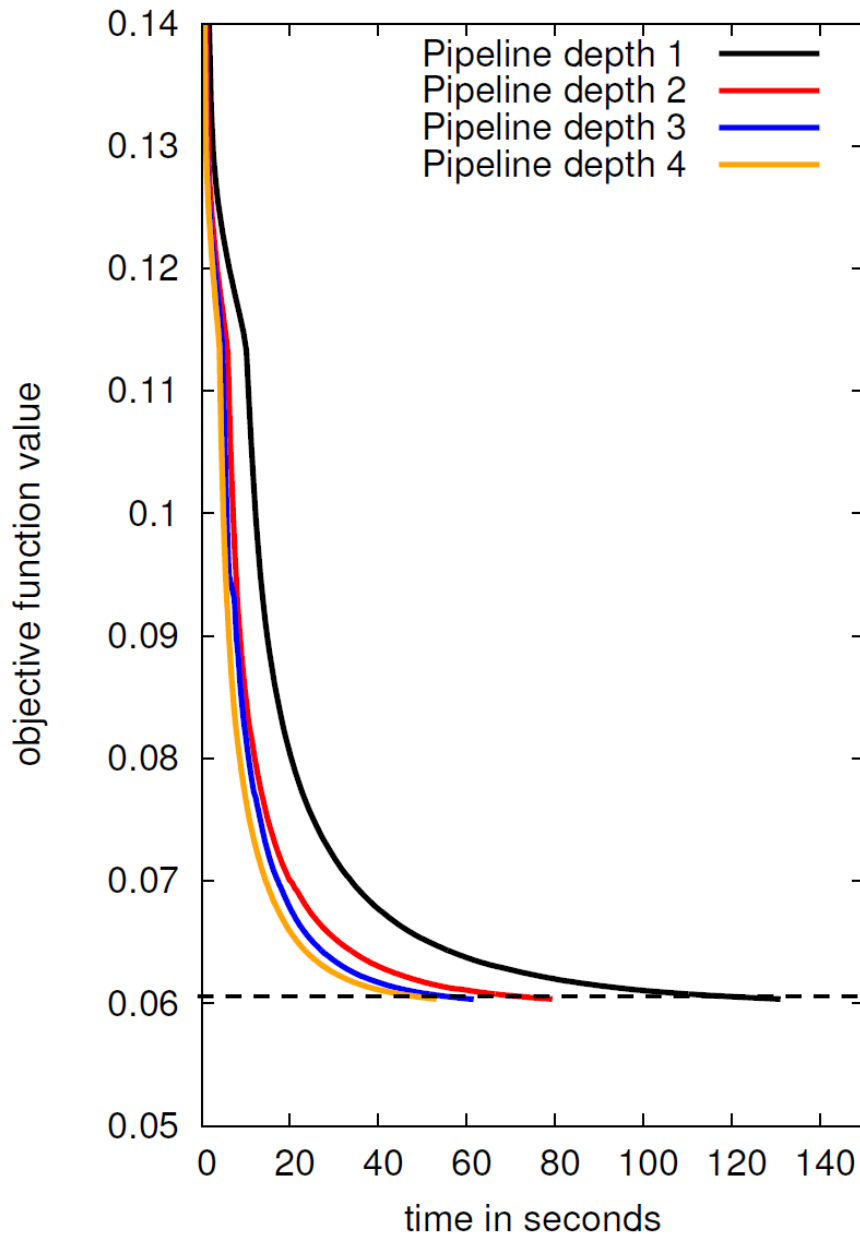
# One pipeline is not enough

**1:** schedule w/ local freshness; execute with global freshness



**2:** relax freshness of least important updates (relative to next Chord)

# STRADS Dual Pipeline Convergence



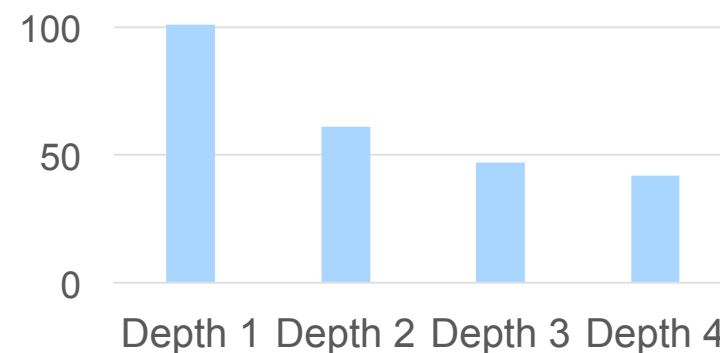Pipeline Experiment with 1M parameters

Application: Lasso
Data: Synthetic data
50K samples, 1M dimension

Depth refers to second pipeline

Time to objective value 0.061

# STRADS Vision

- STRADS' scheduling policies show order of magnitude faster convergence speed compared to parallel ML apps w/o scheduling
- ML Apps (esp. with divergence risks) benefit from significant scheduling and bounded staleness to fully utilize parallelism
- Concept of "iteration" is lost when importance guides update frequency (don't just delay communication, delay computation too)
  - Staleness can still bound minimal update frequency
- Fully utilizing hardware when scheduling is non-trivial adds additional reasons for exploiting staleness induced error tolerance

- Three canonical ML applications (Lasso, Logistic Regression, SVM) implemented of STRADS framework so far.

- Compare Bounded Async Bulk Synch Parallel (A-BSP) vs Stale Synch Parallel (SSP)
  - Similar best case speedups
  - SSP tolerates (transient) stragglers (see paper)
- Repetition-exploiting optimizations (to BSP)
- Managed extra Bandwidth SSP (MBSSP)
  - Smart early-notify speeds convergence
- Convergence-guided Scheduling (STRADS)
  - Up the ML stack to control update order too
  - Escape straightjacket of "the iteration"
  - Tackle divergence head on & use staleness for latency hiding to better utilize hardware

# Contributing Students

- Carnegie Mellon Univ students on this project
  - Henggang Cui
  - Qirong Ho
  - Jinliang Wei
  - Wei Dai
  - Jin Kyu Kim
  - Seunghak Lee
  - Abhimanu Kumar
  - James Cipar
  - Alexey Tumanov
  - Lianghong Xu
  - Jesse Haber-Kucharsky