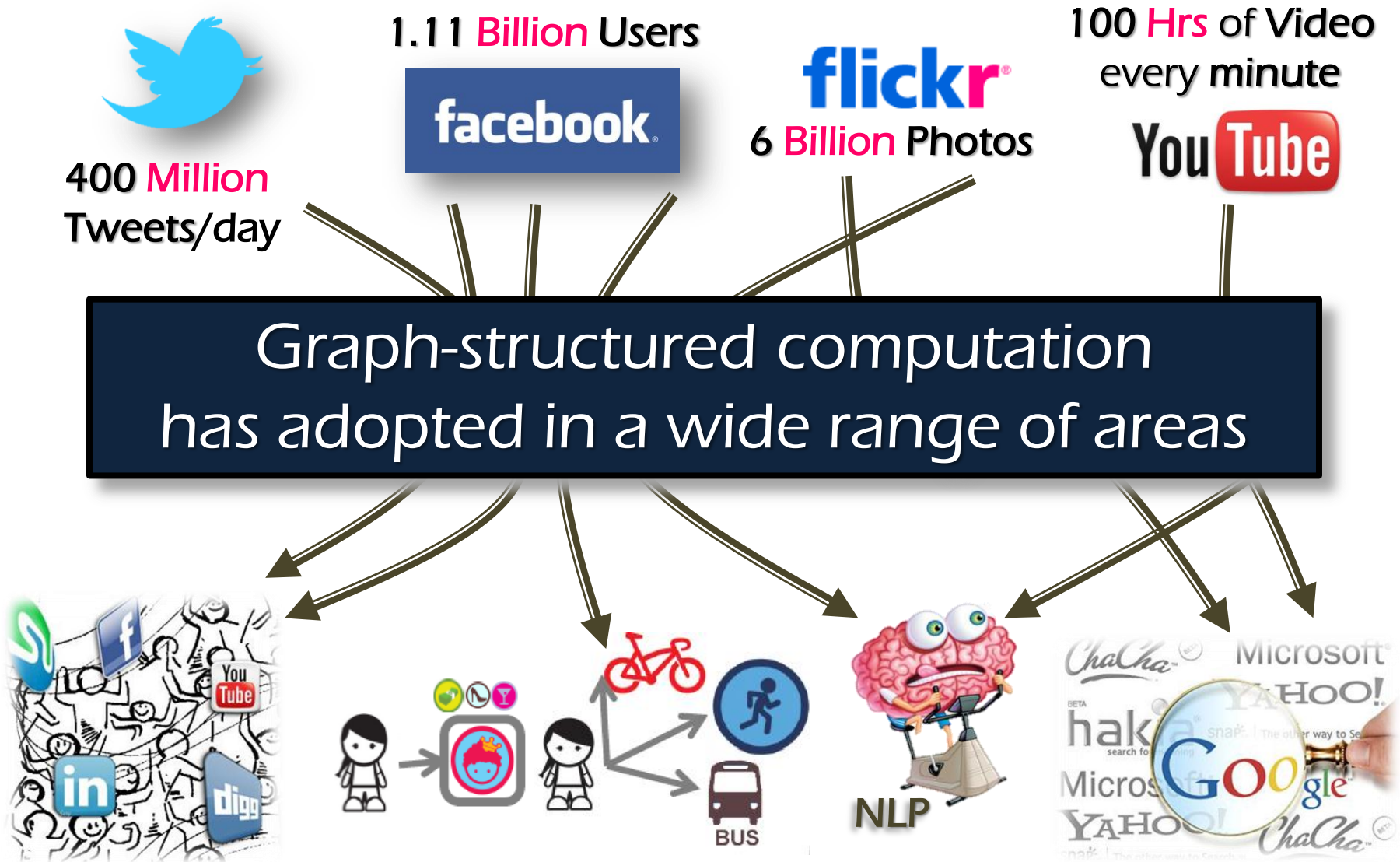# **SYNC** or **ASYNC** ?
# Time to Fuse for Distributed Graph-Parallel Computation

Chenning Xie[+], Rong Chen[+], Haibing Guan[*],
Binyu Zang[+] and Haibo Chen[+]

Institute of Parallel and Distributed Systems [+]
Department of Computer Science [*]
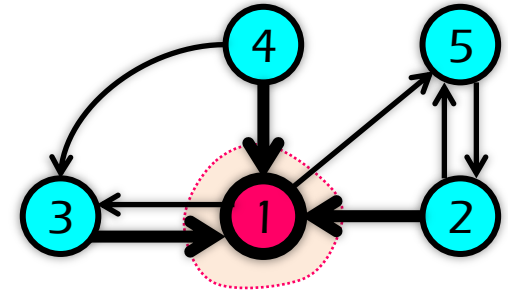Shanghai Jiao Tong University

# Big Data → Graph Computation

**400 Million** Tweets/day

**1.11 Billion** Users

facebook

**100 Hrs** of **Video** every **minute**

flickr

**6 Billion** Photos

You Tube

Graph-structured computation has adopted in a wide range of areas

NLP

# Graph-parallel Computation

"_Think as Vertex_",   e.g. **PageRank** :   $R_i = \alpha + (1-\alpha) \sum_{i,j \in E} W_{i,j} R_j$

## Characteristics

□   Linked set  →  data dependence

□   Rank of who links it  →  predictable accesses

□   Convergence  →  iterative computation

# <u>Distributed</u> Graph Computation

- ☐ Larger Graph
- ☐ Complicated Computation
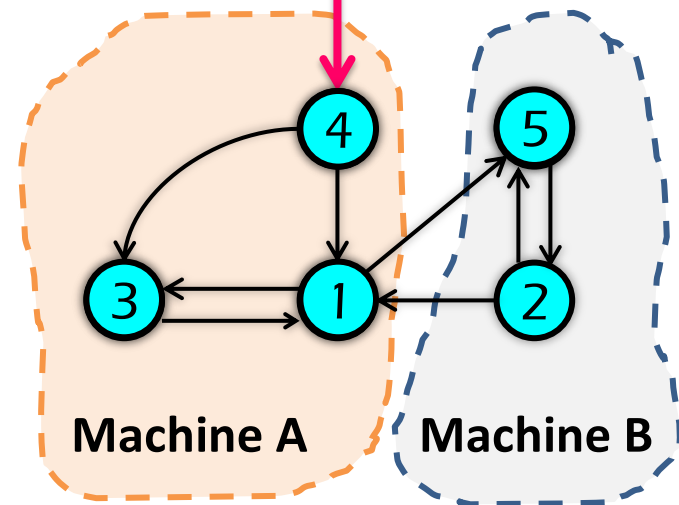- ☐ Storage support

# Distributed Graph Computation

User define the logic (e.g. Pagerank) :

☐ Input: $R_j$ (Data of neighbor j )

☐ Compute():

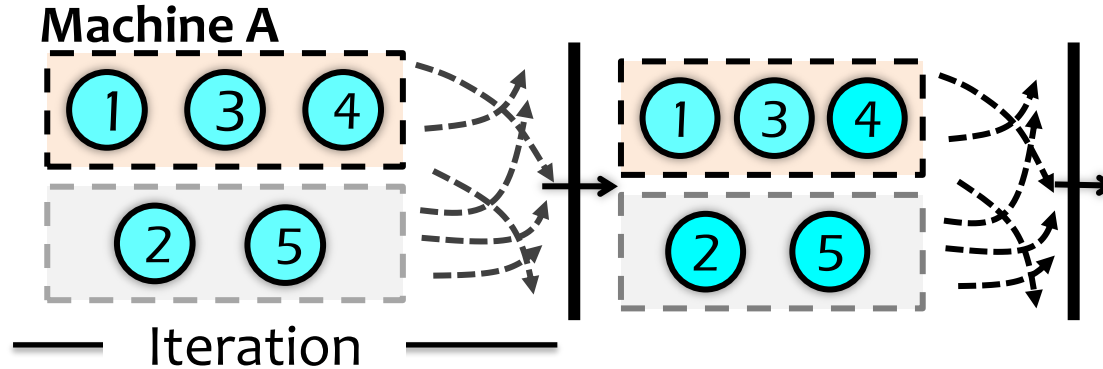$$Vertex\_Data = \alpha + (1-\alpha) \sum W_j R_j$$

## Framework:

☐ Load & partition over **cluster**

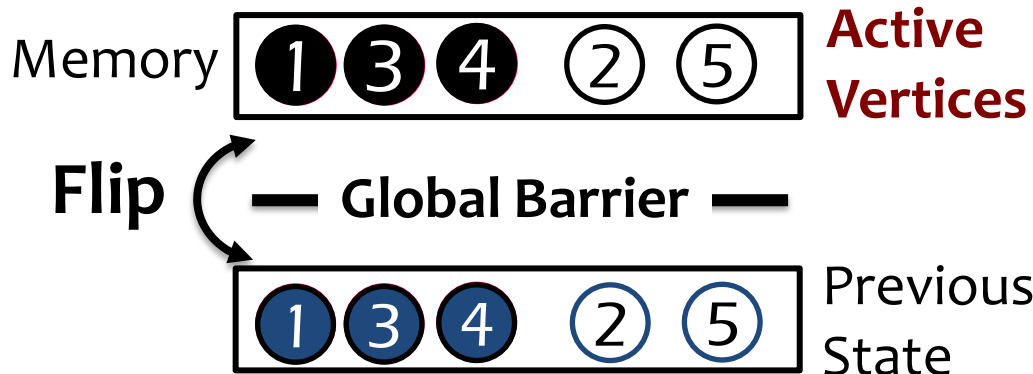☐ **Schedule** compute() **Repeatedly** until get convergence



Machine A          Machine B

# Existing Scheduling Modes - <u>**Synchronous**</u>

( Sync Mode )

## Scheduling:

**Machine A**



— Iteration —

## Internal state (e.g. Machine A):

Memory



**Active Vertices**
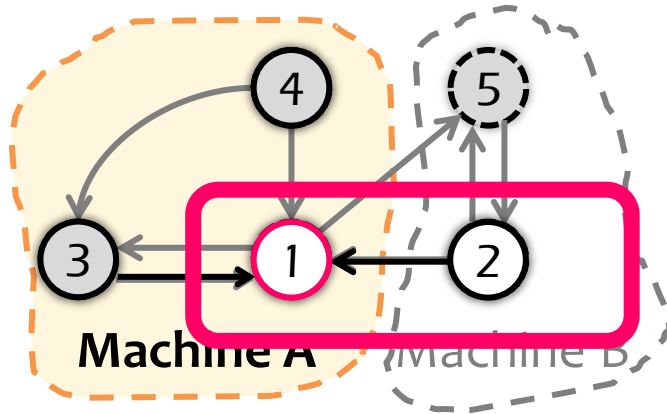
**Flip** — **Global Barrier** —

Previous State

## Pseudocode

**while** (*iteration ≤ max*) *do*

if $V_a == \emptyset$ then break

$V'_a \leftarrow \emptyset$

**foreach** $v \in V_a$ do

$A \leftarrow$ **compute**($v$)

$V'_a \leftarrow V'_a \cup A$

**barrier to update**
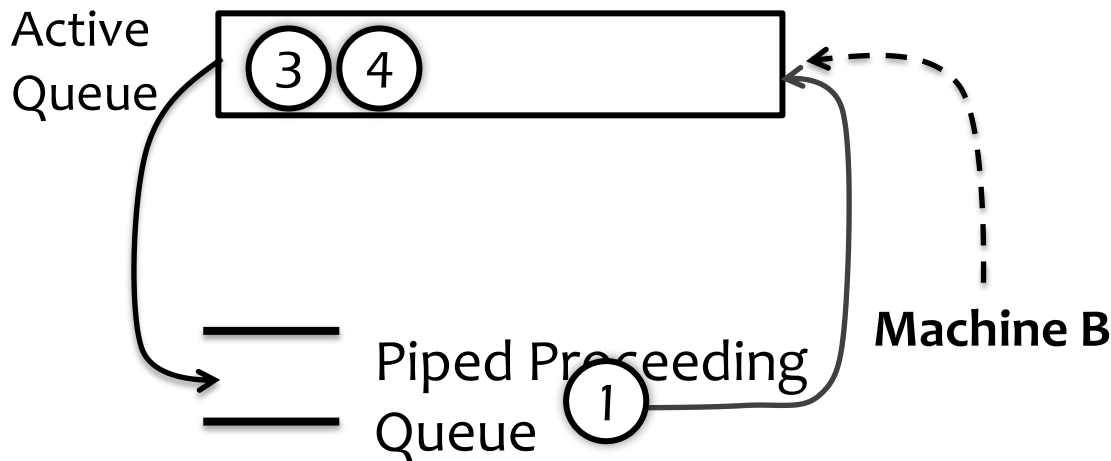
$V_a \leftarrow V'_a$

iteration ++

# Existing Scheduling Modes - __Asynchronous__

( Async Mode )

## Scheduling:



**Machine A**    Machine B

## Internal State (e.g. Machine A):

Active Queue
Piped Proceeding Queue
Machine B

## Pseudocode

**while** $(V_a \neq \emptyset)$ *do*

   $v = $ **dequeue**$(V_a)$

   $A \leftarrow$ **compute**$(v)$

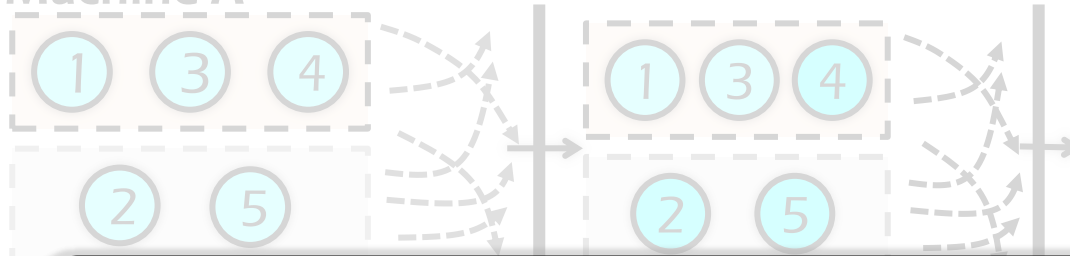   $V'_a \leftarrow V'_a \cup A$

   **signal across machines**

# Propagate ASAP,
to converge faster

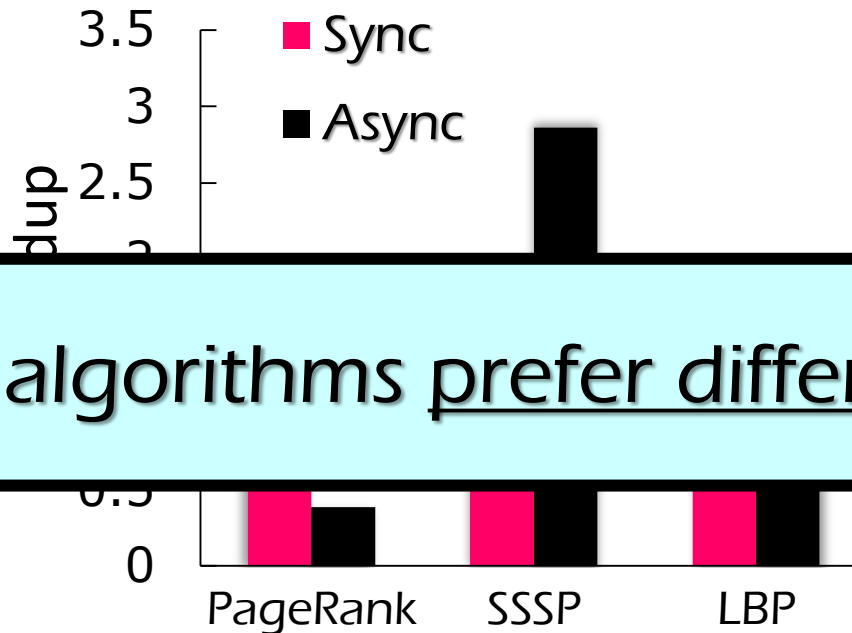# Existing Scheduling Modes

Synchronous

Machine A

1 3 4

2 5

## Which could get a **<u>better</u>** performance?

4 5

3 1 2

Machine A    Machine B

Asynchronous

# Algorithms: Sync vs. Async

- Same Configuration + Different Algorithms ?



Different algorithms prefer different modes

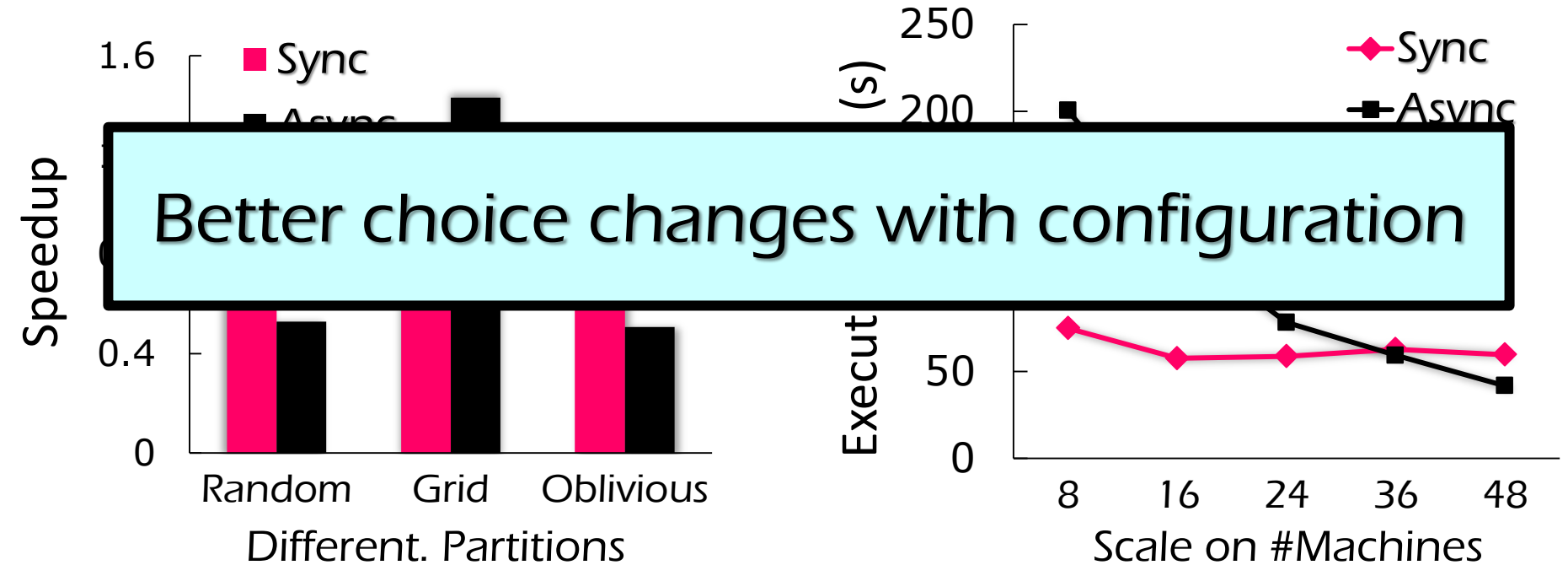Legend: Sync, Async

Chart axis values: 3.5, 3, 2.5, 2, 0.5, 0
Categories: PageRank, SSSP, LBP

**Large active vertex set** with collecting all data from neighbors

Require **fast broadcast** of shortest path value

Belief Propagation Algorithm

# Configuration: Sync vs. Async

- Same Configuration + Different Algorithms:    Uncertain
- **Different Configuration** + Same Algorithms (LBP) ?



Better choice changes with configuration

Speedup — Different. Partitions — Random, Grid, Oblivious — Sync, Async

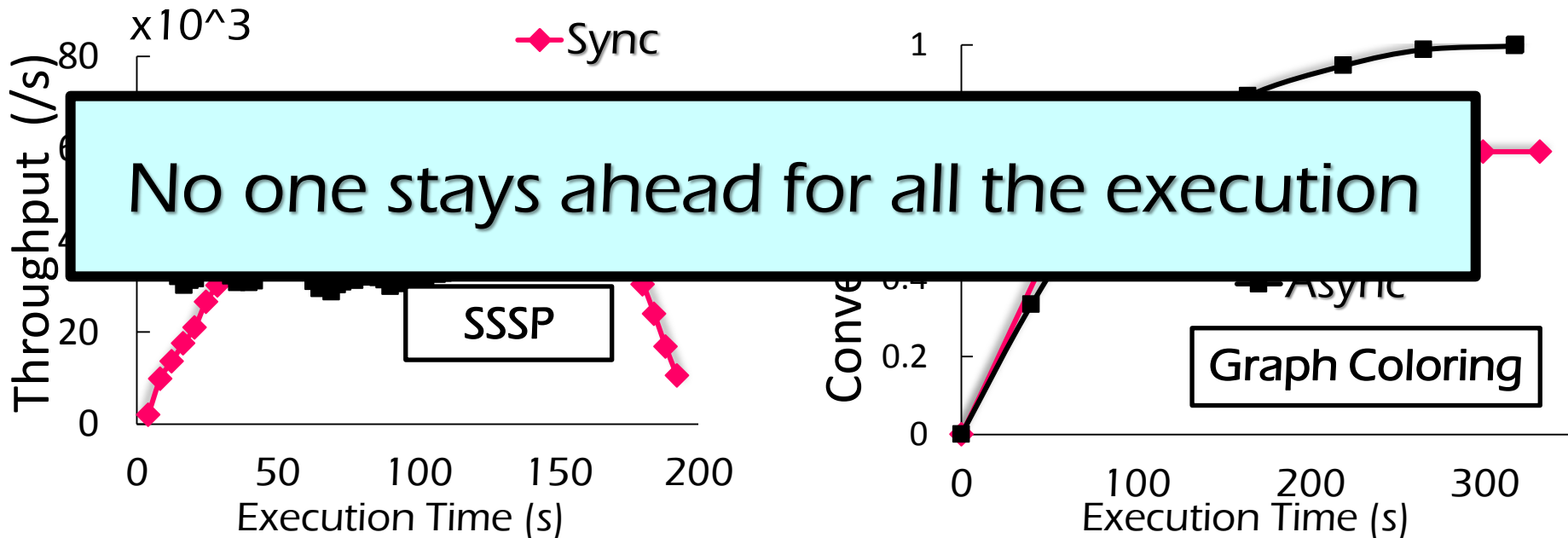Execut (s) — Scale on #Machines — 8, 16, 24, 36, 48 — Sync, Async — 250, 200, 50, 0

Partition methods affect
load balance & communication

Sync mode batches heavy load,
Async mode scales better.

# Stages: Sync vs. Async

- Same Configuration + Different Algorithms:     Uncertain
- Different Configuration + Same Algorithms:     Uncertain
- <u>Same</u> Configuration + <u>Same</u> Algorithm  ?



No one stays ahead for all the execution

SSSP

Sync

Async

Graph Coloring

Throughput (/s)   x10^3   80 ... 20 ... 0

Execution Time (s)   0   50   100   150   200

Conv...   1 ... 0.4 ... 0.2 ... 0

Execution Time (s)   0   100   200   300

Async mode starts **_faster,_**
Sync mode grows with **_a peak_**.

Sync **_faster_** but not converge,
Async _slower_ but  **converged**.

# Summery: Sync vs. Async

| Properties | SYNC | vs. | ASYNC |
|---|---|---|---|
| → Communication | Regular | | Irregular |
| → Convergence | Slow | | Fast |

**Better choice is Unintuitive**

**Single mode alone may be still Suboptimal**

| | | | |
|---|---|---|---|
| → Workload | Heavyweight | | Lightweight |
| → Scalability | \| Graph \| | | \| Machines \| |

# Contributions

First comprehensive study on Sync & Async modes

**PowerSwitch** – adaptive, fast & seamless switches

Hybrid Execution Mode ( **Hsync** Mode ):

- □ Dynamically and **transparently** support the correct mode switches

Switch Timing Model:

- □ Determine the more efficient mode combined with online sampling, offline profiling and heuristics

# Agenda

## How to Switch - the Hsync mode

- ☐ Internal state conversion
- ☐ Consistency & correctness

## When to Switch – the timing model

- ☐ Performance Metrics
- ☐ Current mode prediction
- ☐ The other mode estimation

## Implementation

## Evaluation

# Challenges of switches

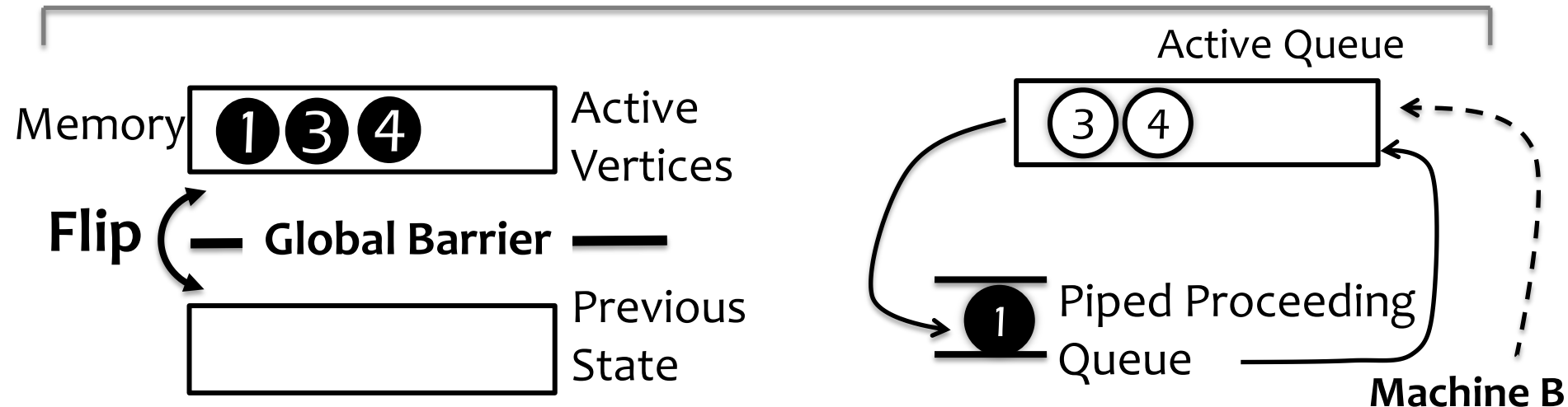> Convert state at Consistent switch points

## Sync mode

- ☐ Vertex update: **unordered**
- ☐ Flip in **global barrier**

## Async mode

- ☐ **Priority/FIFO queue**
- ☐ Dequeue and enqueue

Internal state of one machine

Memory — Active Vertices: **1** **3** **4**

**Flip** — Global Barrier — Previous State

Active Queue: ③ ④

Piped Proceeding Queue: **1**
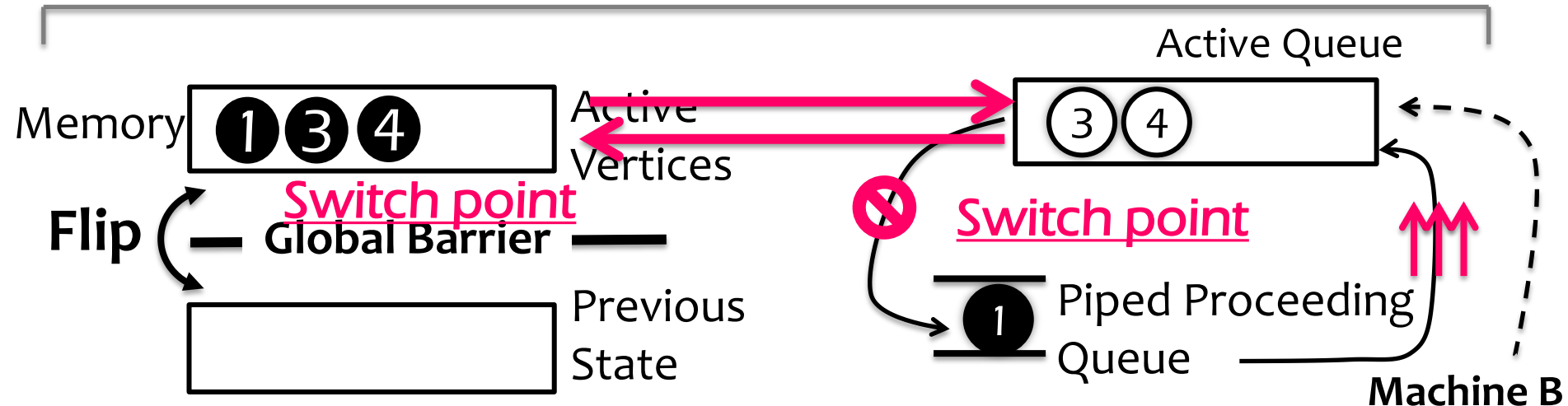
**Machine B**

# Challenges of switches   Hsync mode

Consistent switch points :

- Sync -> Async:          global barrier
- Async -> Sync:          suspend & wait

State transfer:          active vertex set

Internal state of one machine

Active Queue

Memory  ❶ ❸ ❹   Active Vertices   ③ ④

**Flip**   Switch point
**Global Barrier**

Switch point

Previous State

❶ Piped Proceeding Queue

**Machine B**

# Agenda

## How to Switch  -  the Hsync mode

- ☐ Internal state conversion
- ☐ Consistency & correctness

## When to Switch – the timing model

- ☐ Performance Metrics
- ☐ Current mode prediction
- ☐ The other mode estimation

## Implementation

## Evaluation

# Switch timing - affected by lots of factors

Challenges:

☐ How to quantify the real-time performance?

☐ How to obtain the metrics?

Performance Metrics

☐ Throughput $= \dfrac{|V_{compute}|}{T_{interval}} * \boldsymbol{\mu}$

**Convergence ratio** $\boldsymbol{\mu} = \dfrac{|NumTask_{async}|}{|NumTask_{sync}|}$

by sampling specific input pattern,

e.g. power-law, large diameter, high density...

# Predict <u>Throughput</u> for <u>Current</u> mode

Sync

☐ Iteration as interval

☐ Throughput

$$\frac{|V_{next}|}{}$$

Async

☐ Constant interval

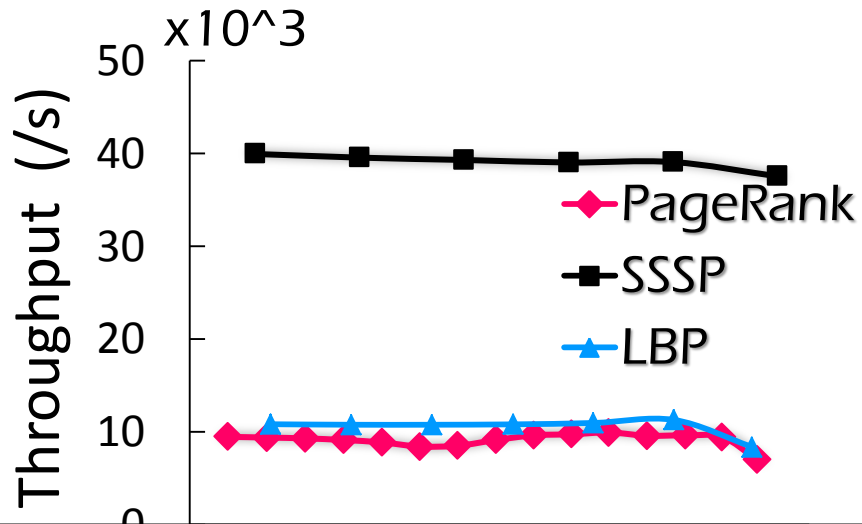☐ Throughput

$$\frac{1}{}$$

Calculate the <u>next interval</u> based on:

<u>Current</u>  +  <u>History accumulation</u>

$$T_{vert_{next}} = \alpha \cdot T_{vert_{current}} + (1 - \alpha) \cdot T_{vert_{history}}$$

# Predict for <u>Other offline</u> mode

No more execution information

Predict **<u>Async</u>** when in sync mode:



## Solution

**Online sampling** : on subset of input in Async before start

**Offline profiling** : build Neural Network model, refer to paper

# Predict for <u>Other offline</u> mode

Predict **<u>Sync</u>** when in async mode:

☐ Hard to predict exactly

☐ Heuristic:  Sync makes <u>high utilization</u> of resource.

$Thro_{Sync} > Thro_{Async}$ , if <u>workload is enough</u>

Condition:

1. Number of active vertices increases

2. Workload :  $\dfrac{|V_{new}|}{T} > Thro_{Async}$

Async -> Sync
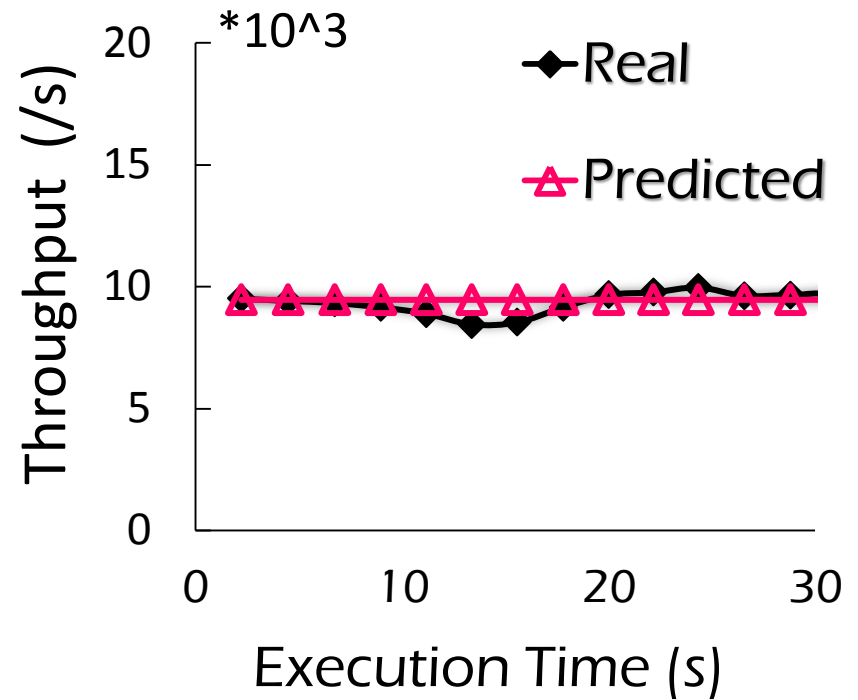
# Prediction Accuracy

PageRank:  Predicted throughput vs. Real sampled
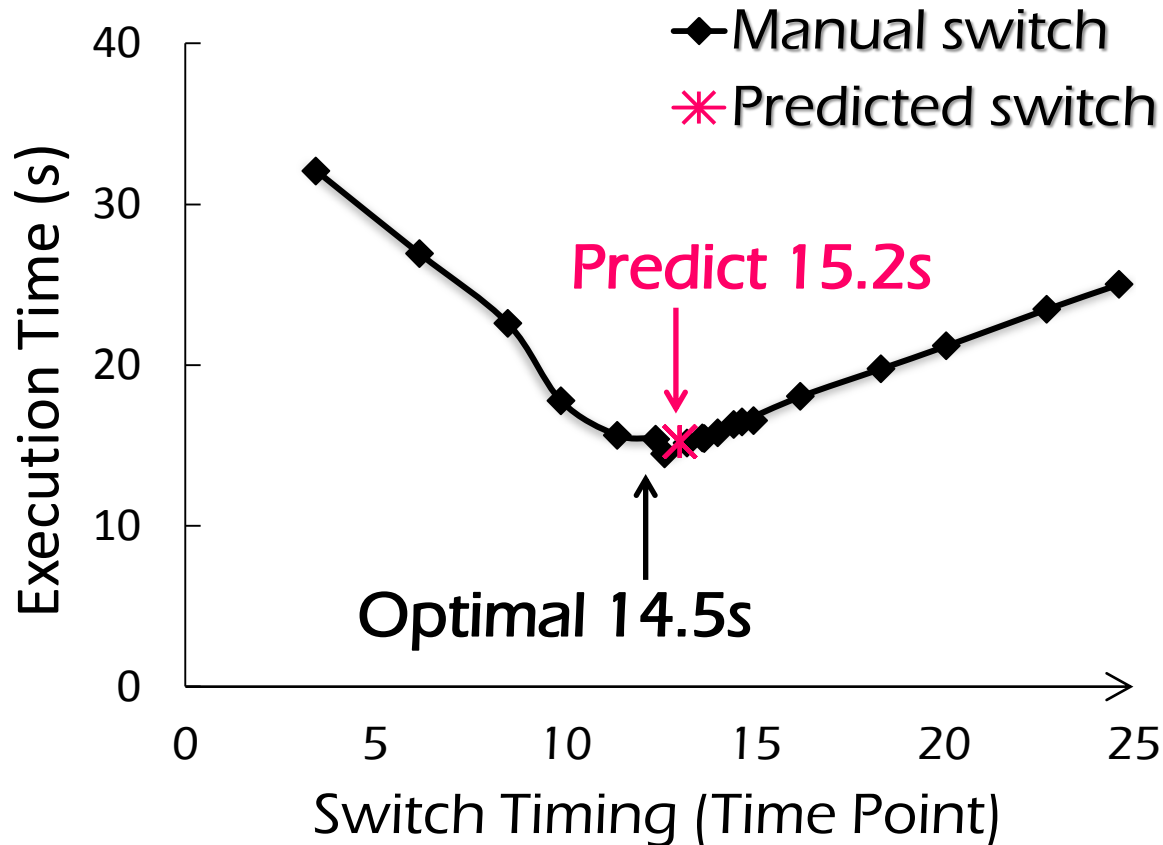


Sync mode

Async mode

# Prediction Accuracy

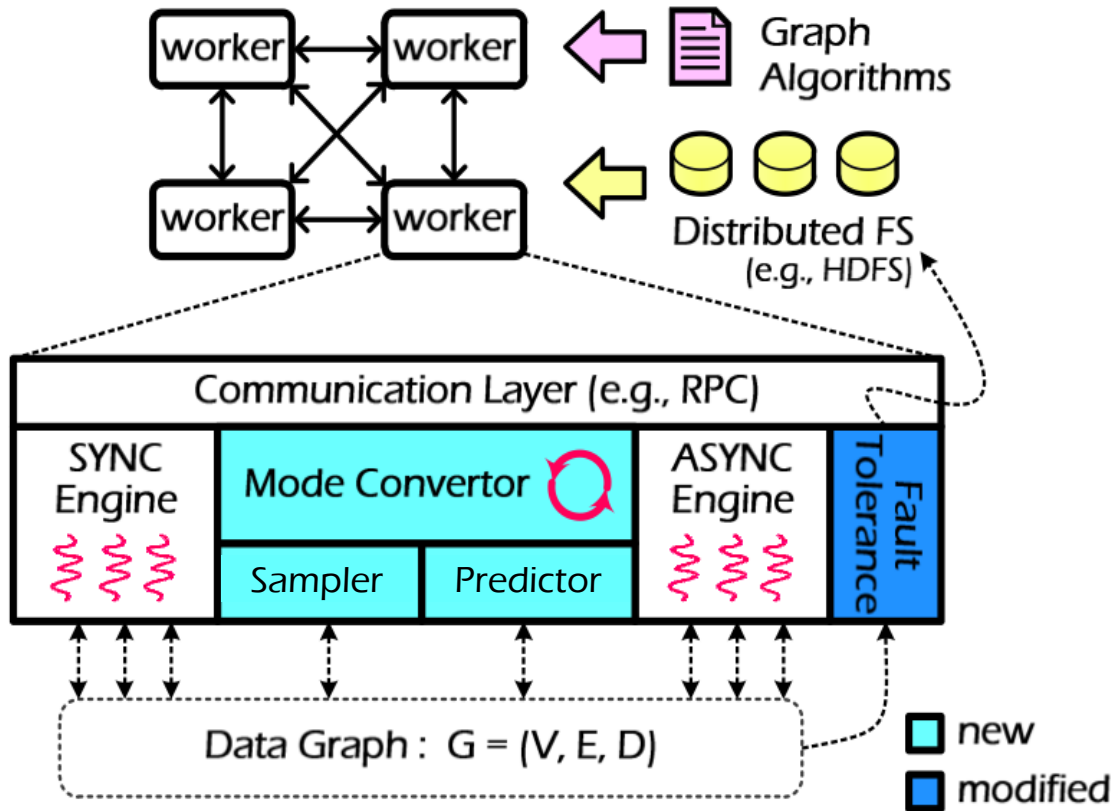PageRank:   Predicted switch timing vs. Optimal

# Implementation

## PowerSwitch:

- Based on latest GraphLab (PowerGraph) v2.2 with both Sync & Async modes.

- Provide the same graph abstraction **transparent** & **compatible** to all apps of GraphLab

## Open Source

http://ipads.se.sjtu.edu.cn/projects/powerswitch.html

# Implementation - Architecture



New

- Mode switcher
- Sampler
- Predictor

Extension

- Fault tolerance

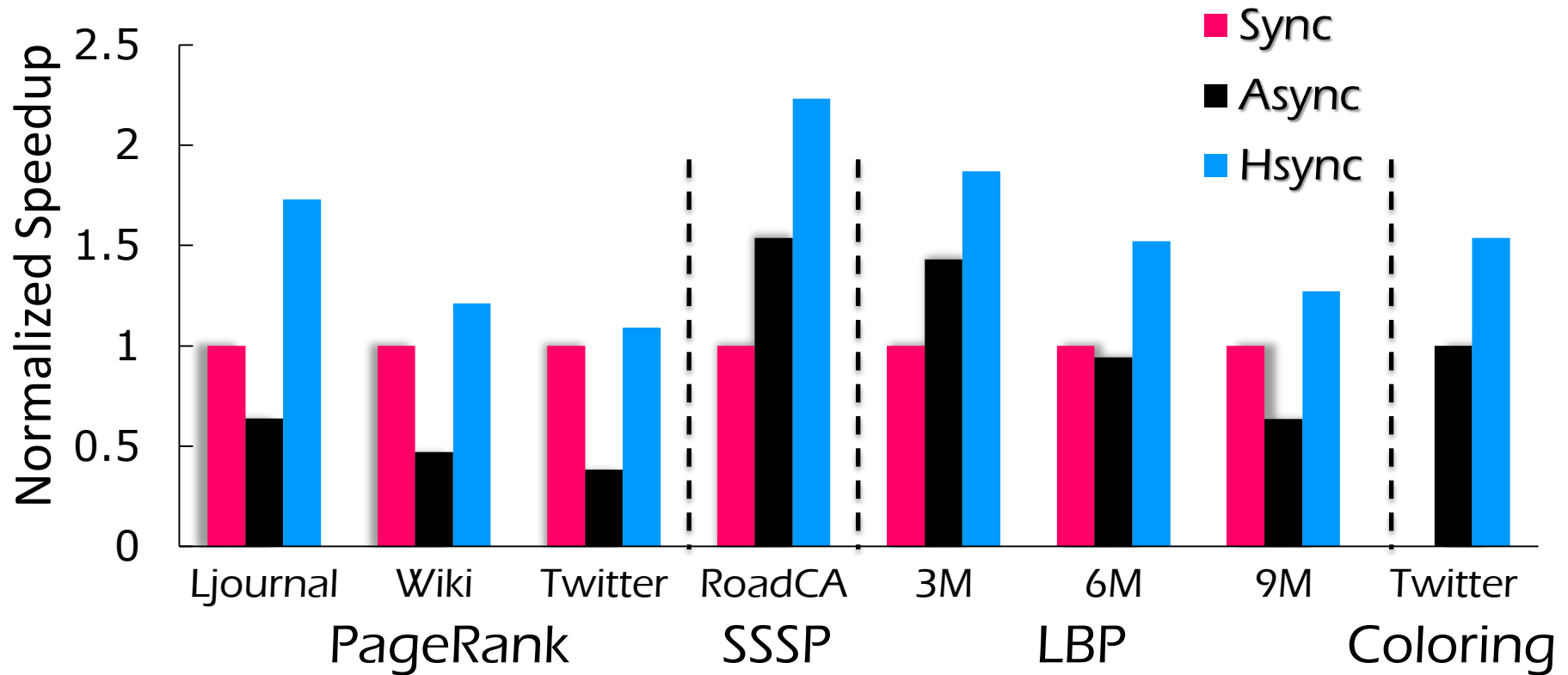# Evaluation

**Baseline**: original SYNC & ASYNC mode

## Configuration

- 48-node EC2-like cluster (VM based).
- Each node has 4 AMD Opteron cores, 12GB of RAM, connected with 1 GigE network.

## Algorithms and Data Set

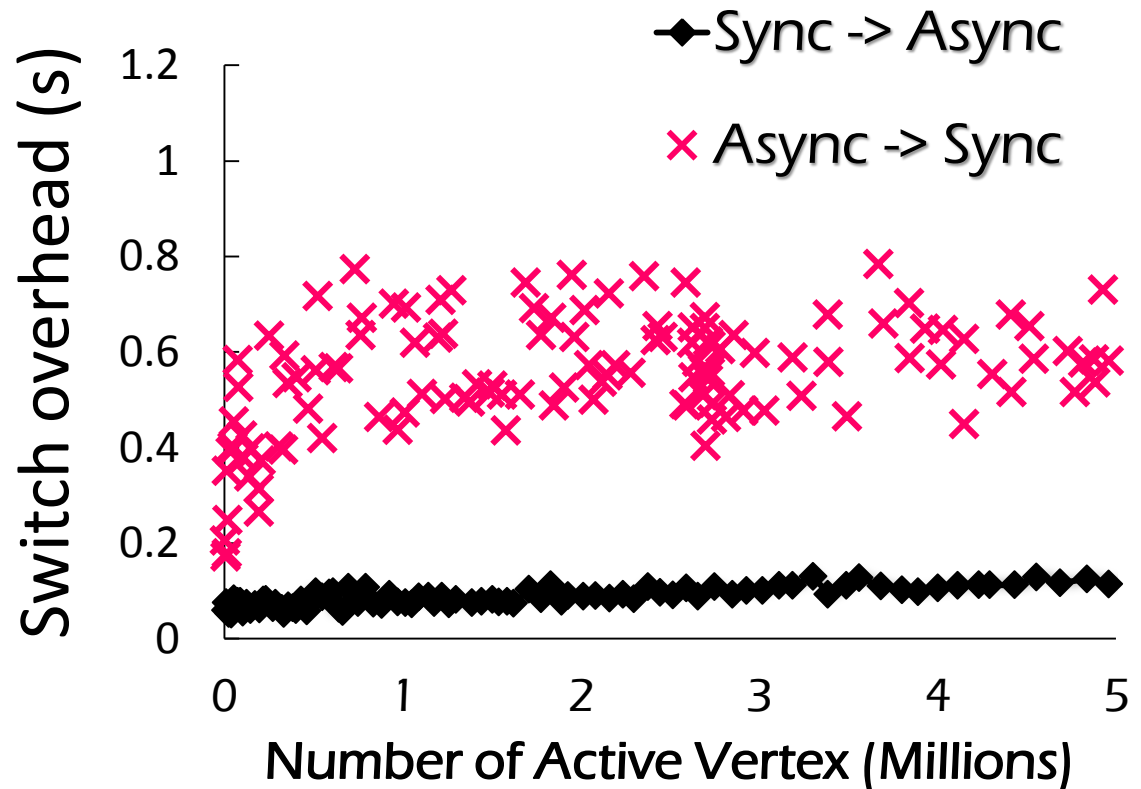| Algorithm | Graph | |V| | |E| |
|-----------|-------|-----|-----|
| PageRank | LJournal | 5.4M | 79M |
| | Wiki | 5.7M | 130M |
| | Twitter | 42M | 1.47B |
| LBP | SYN-ImageData | 1-12M | 2-24M |
| SSSP | RoadCA | 1.9M | 5.5M |
| Coloring | Twitter | 42M | 1.47B |

# Performance Overview



Outperform the baseline with <u>best mode from 9% to 73%</u> for all algorithms and dataset
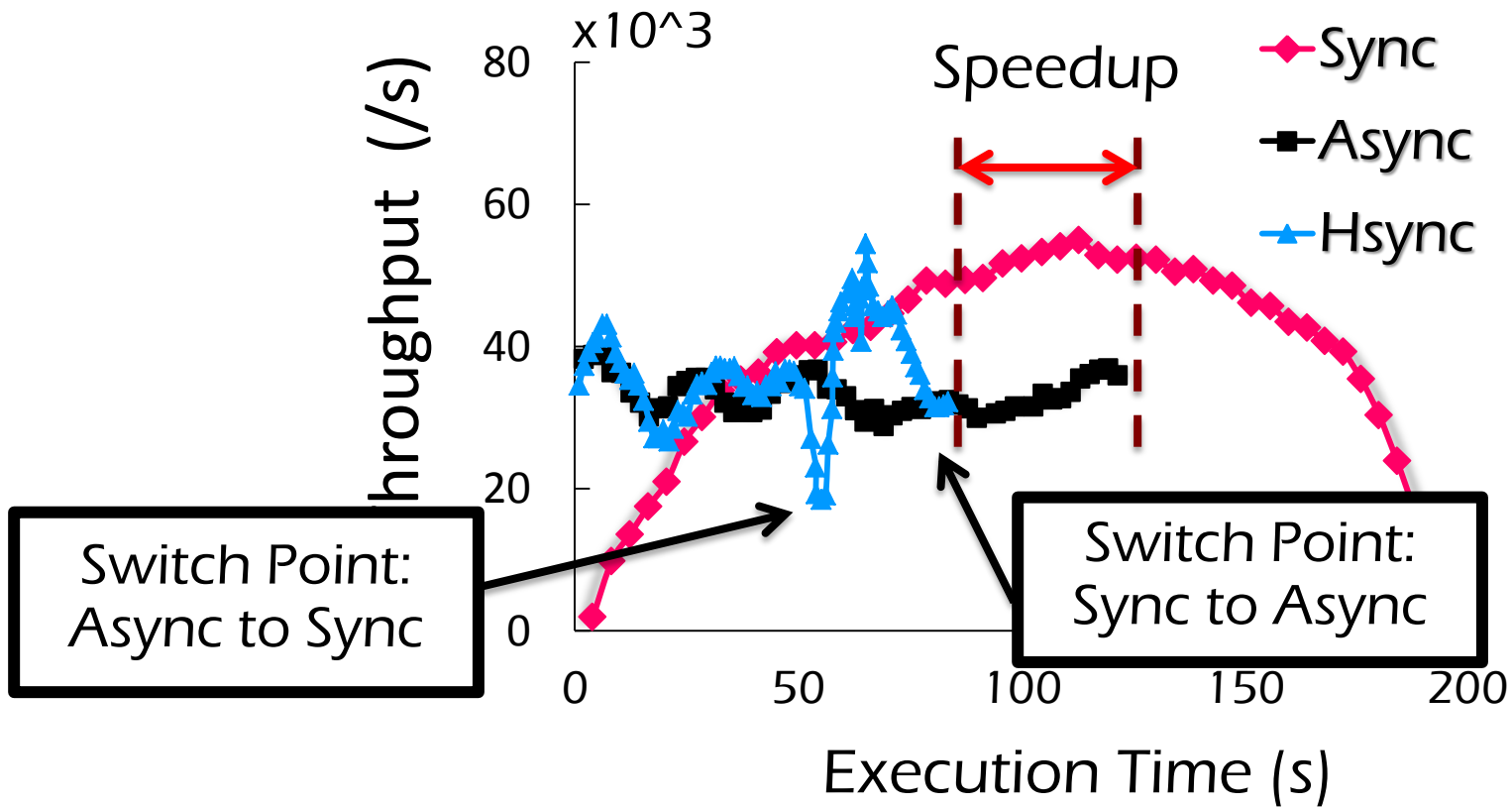
# Switch Overhead

Sync->Async : 0.1s

Async->Sync : 0.6s



**Overhead grows slightly**
with active vertex number increasing.

# Case: Single Source Shortest Path (SSSP)



Execution Mode:  Async -> Sync-> Async

# Conclusion

## PowerSwitch

- A comprehensive analysis to the performance of Sync and Async modes for different algorithms, configuration and stages

- A Hsync mode that dynamically switch modes between Sync & Async to pursue optimal performance

- An effective switch timing model to predict suitable mode with sampling & profiling

- Outperforms GraphLab with best mode from 9% to 73% for various algorithms and dataset

# Thanks



PowerSwitch

http://ipads.se.sjtu.edu.cn/
projects/powerswitch.html

Institute of Parallel and
Distributed Systems



Questions