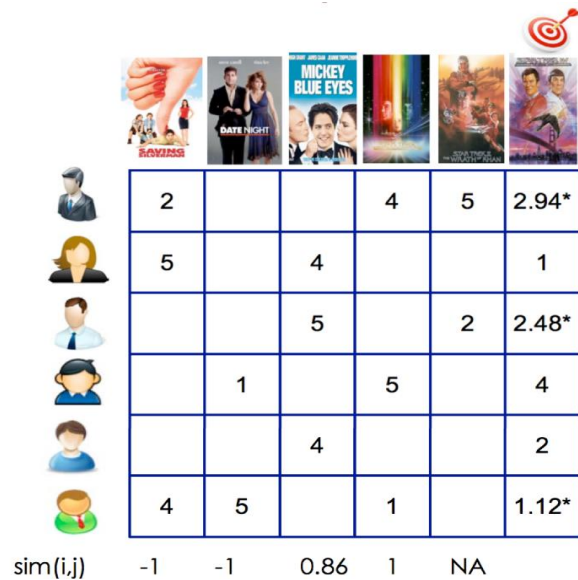


# TuX<sup>2</sup>: Distributed Graph Computation for Machine Learning

Wencong Xiao<sup>\*◇</sup>, Jilong Xue<sup>◇†</sup>, Youshan Miao<sup>◇</sup>, Zhen Li<sup>\*◇</sup>,  
Cheng Chen<sup>◇</sup>, Ming Wu<sup>◇</sup>, Wei Li<sup>\*</sup>, Lidong Zhou<sup>◇</sup>

*<sup>\*</sup>SKLSDE Lab, Beihang University; <sup>◇</sup>Microsoft Research; <sup>†</sup>Peking University*

# Machine Learning(ML) in real world



## Recommendation

**Other Movies You Might Enjoy**

**Amelie**

Add

★★★★☆

Not Interested

**Y Tu Mama Tambien**

Add

★★★★☆

Not Interested

**Guys and Dolls**

Add

★★★★☆

Not Interested

**Mostly Martha**

Add

★★★★☆

Not Interested

**Only Human**

Add

★★★★☆

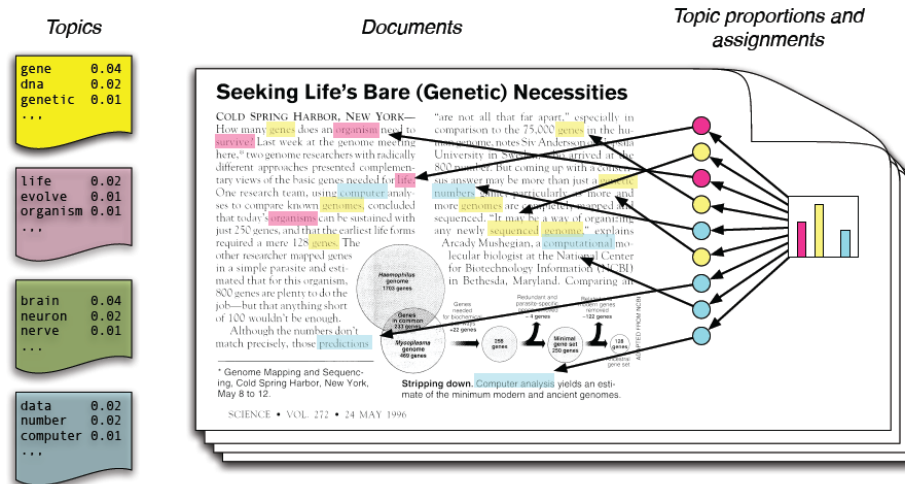
Not Interested

**Russian Dolls**

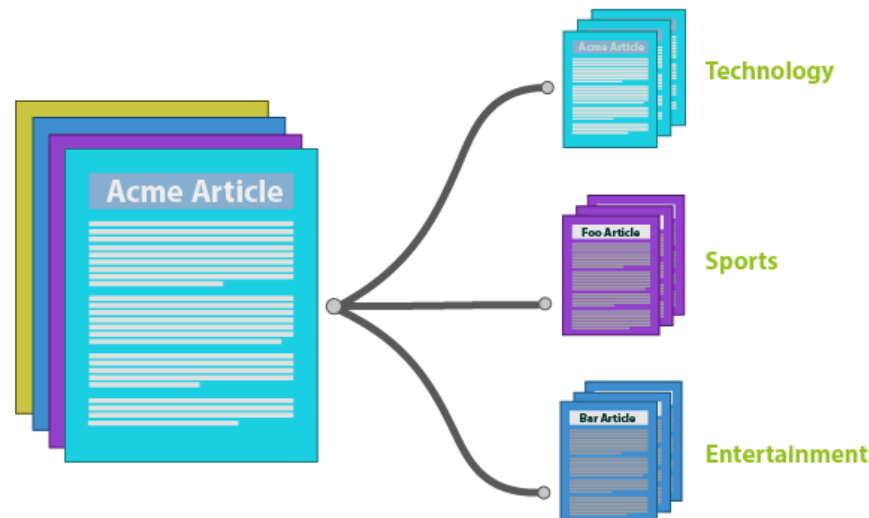
Add

★★★★☆

Not Interested



## Topic Model



## Click Prediction

**Women's Shoes**

2,060,000 RESULTS Any time Near Seattle, Washington Change

**Discount Women's Shoes - Save Up To 75% Off On Top Brands.**

Women's Shoes: Boots, Pumps, Flats, Sandals, Wedges | DSW

Women's Shoes: Boots, Pumps, Flats, Sandals, Wedges | Shoes.com

Women's Shoes: Overstock.com Buy Heels, Sandals, & Boots

Related searches for Women's Shoes

DSW Shoes Easy Sock Shoes Pumps Shoes Women's Dress Shoes Nine West Shoes JCPenney Women Shoes

**Women's Shoes at JCPenney**

Shop For Women's Shoes at JCPenney. Exciting Deals. Save on Top Brands.

**Official Payless Website**

Find Women's Shoes in Large Sizes. Shop Payless & Start Saving Today!

**TOMS Women Shoes**

Wear TOMS in Your Wedding and Help Children in Need.

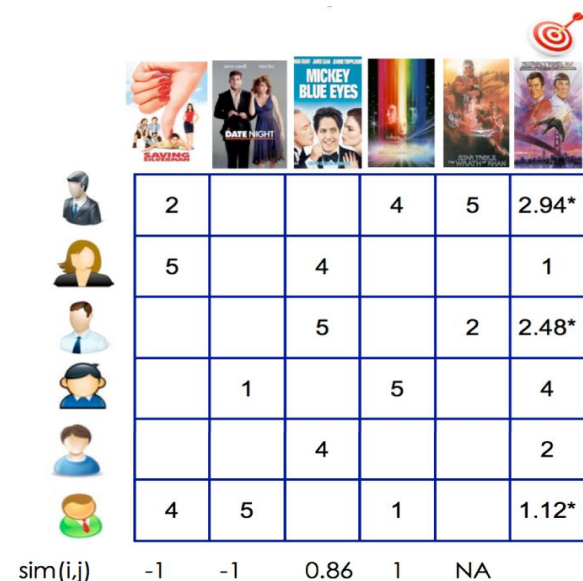
**Naturalizer Shoes**

Save Up To 70% Off At Naturalizer.com. Buy Now & Save!

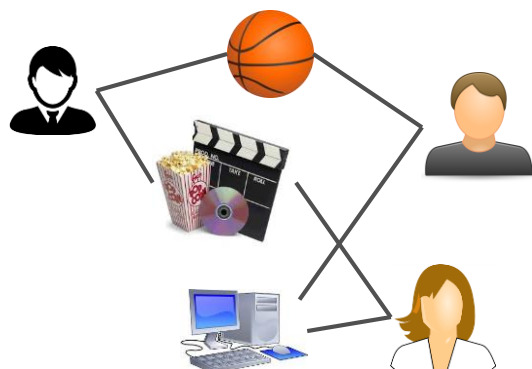
**Ann Taylor Official Site**

The Perfect Pair Of Shoes For Every Outfit! Shop Flats, Heels & Wedges.

# Graph Structures in Machine Learning



Recommendation



Topics

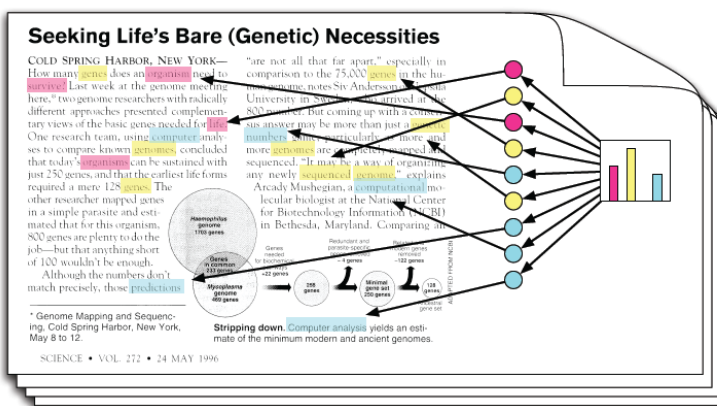
gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

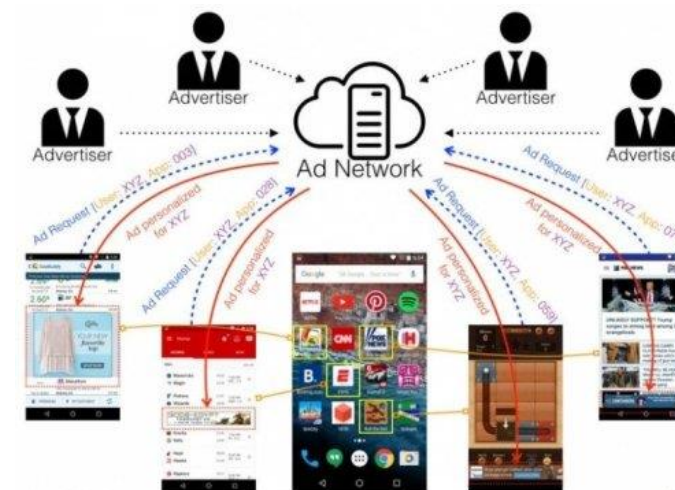
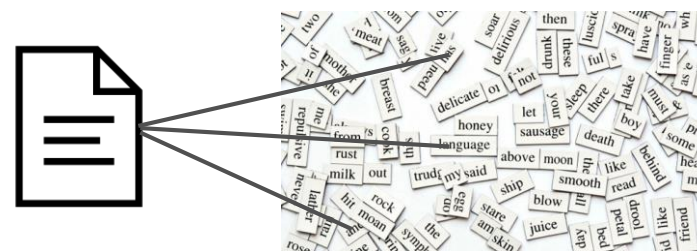
data 0.02  
number 0.02  
computer 0.01  
...

Documents

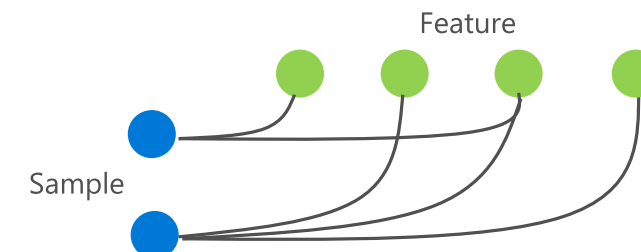


Topic proportions and assignments

Topic Model



Click Prediction



# Advantages of Graph Engine

## Simple programming model (e.g. GAS)

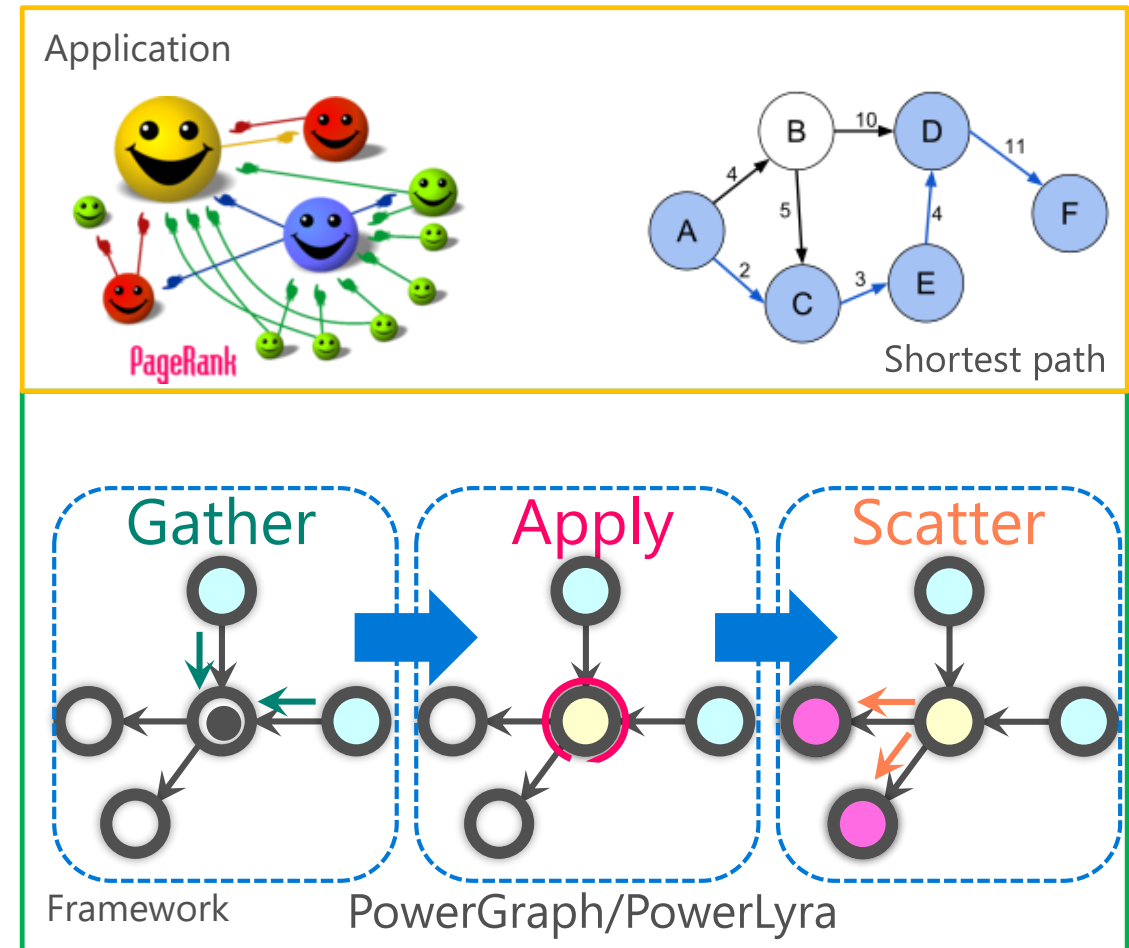
- PageRank, Shortest path, etc.

## Graph-aware optimization

- Data layout [Grace(ATC'12), Naiad(SOSP'13)]
- Partitioning [PowerLyra(EuroSys'15)]

## Scalability to trillion-edge

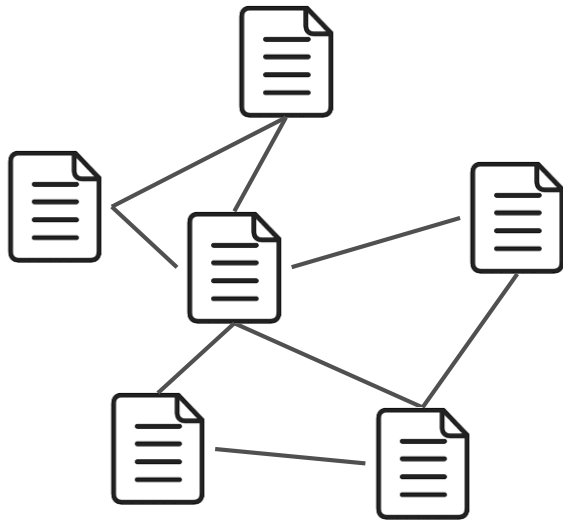
- GraM (SoCC'15)
- Chaos (SOSP'15)
- One Trillion Edges (VLDB'15)



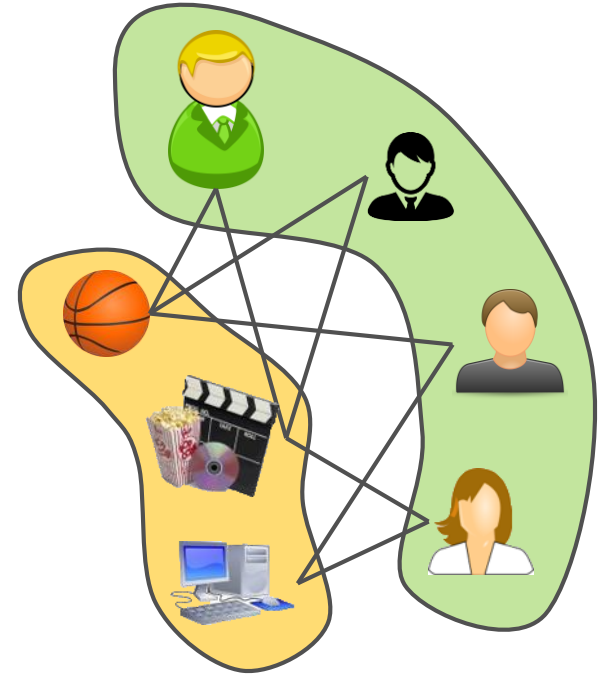
\*GAS figure from PowerLyra slides

# Gaps for ML on Graph Engine

## 1. Heterogeneous vertices



PageRank  
for WebPage Ranking

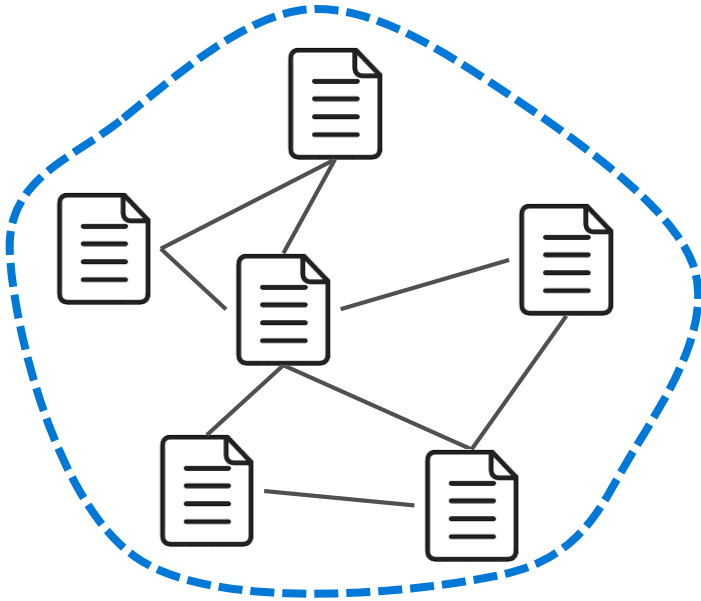


Matrix Factorization(MF)  
for Recommendation

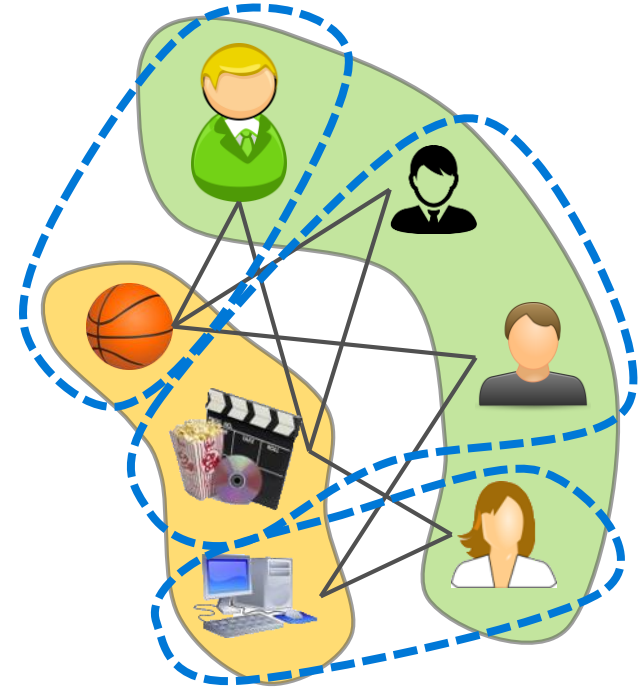


# Gaps for ML on Graph Engine

## 2. Mini-Batch



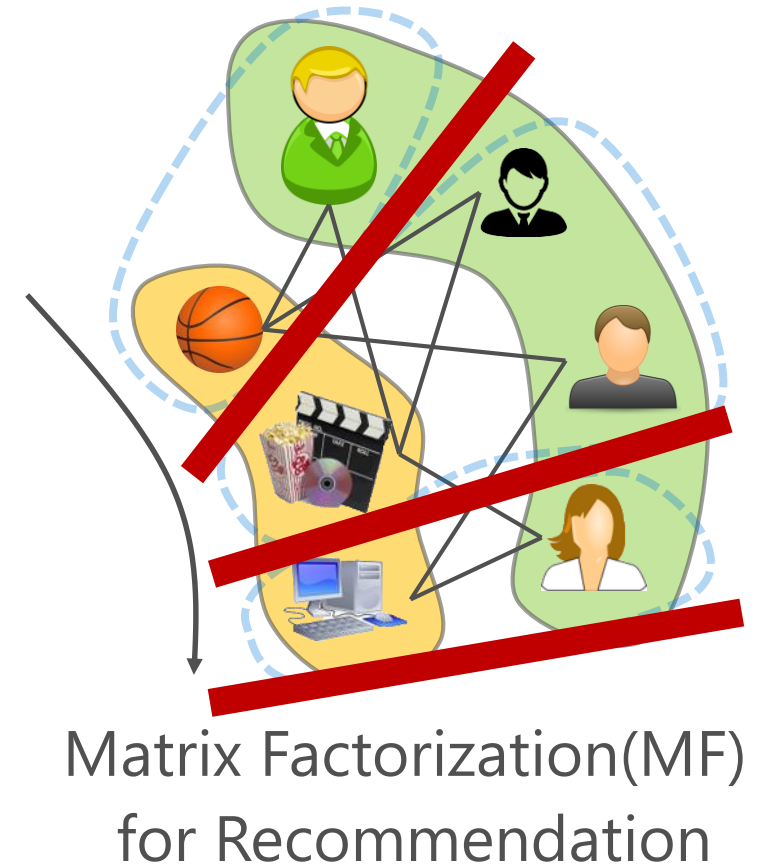
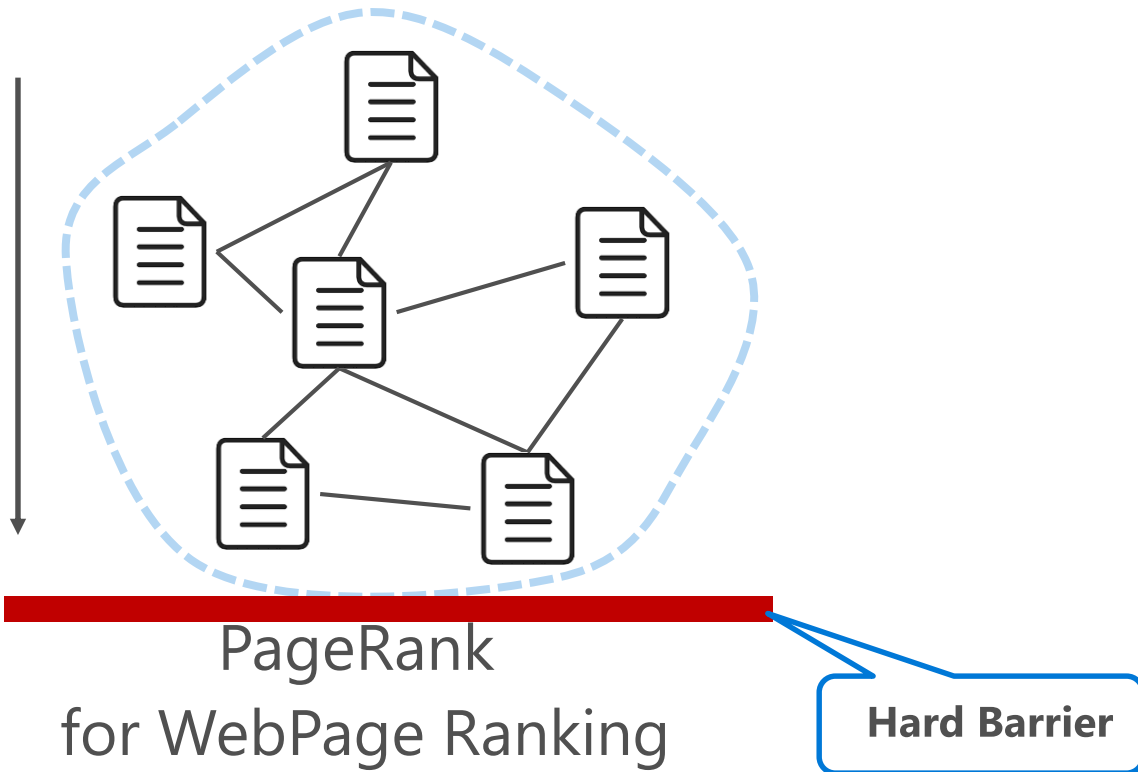
PageRank  
for WebPage Ranking



Matrix Factorization(MF)  
for Recommendation

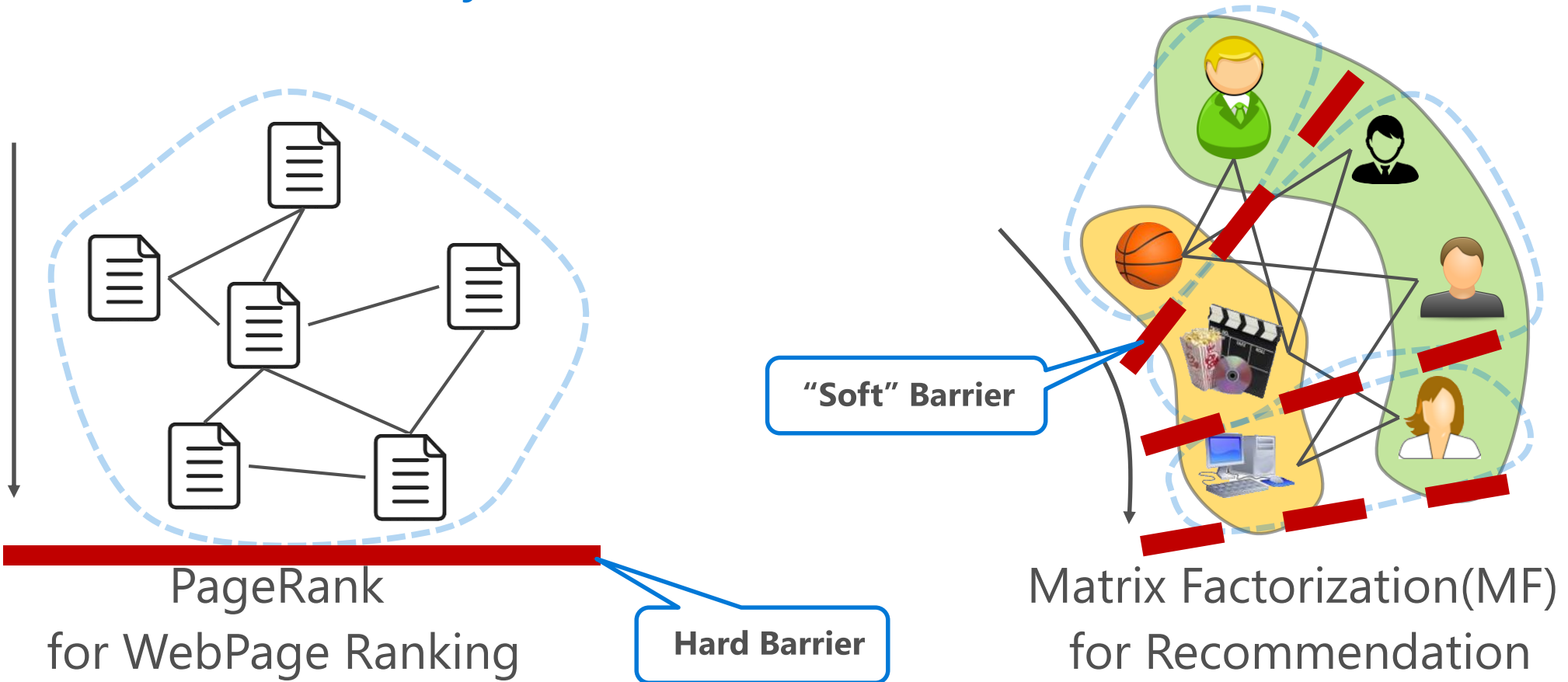
# Gaps for ML on Graph Engine

## 3. Flexible consistency



# Gaps for ML on Graph Engine

## 3. Flexible consistency







# We propose: TuX<sup>2</sup>

## Bridge Graph and ML research in one system

## Extend for distributed machine learning

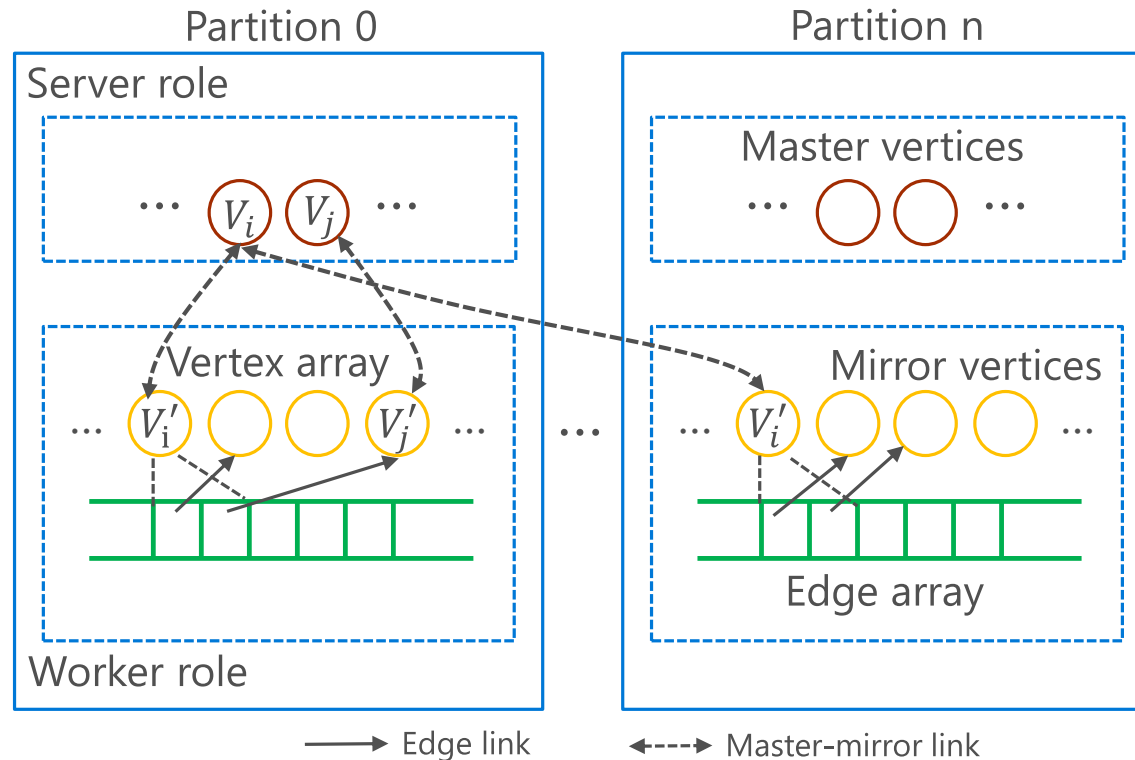
- Scheduling: Stale Synchronous Parallel (SSP) based scheduling
- DataModel: Heterogeneous data model
- Programming: MEGA (Mini-batch, Exchange, GlobalSync, and Apply) graph model

## Outperform both Graph and ML systems on ML algorithms

- **10x**  vs. PowerGraph/PowerLyra
  - Mainly due to MEGA model and heterogeneity optimization
- **48%**  vs. Petuum/Parameter-Server(P-S)
  - Mainly due to graph-based optimization

# System Architecture

## Overview



## Vertex-cut approach

- Effective for power-law graph
- Naturally fits P-S model
  - Master vertices as the global state
  - Mirror vertices as the local cache

# Key designs

Scheduling: Stale Synchronous Parallel (SSP) based scheduling

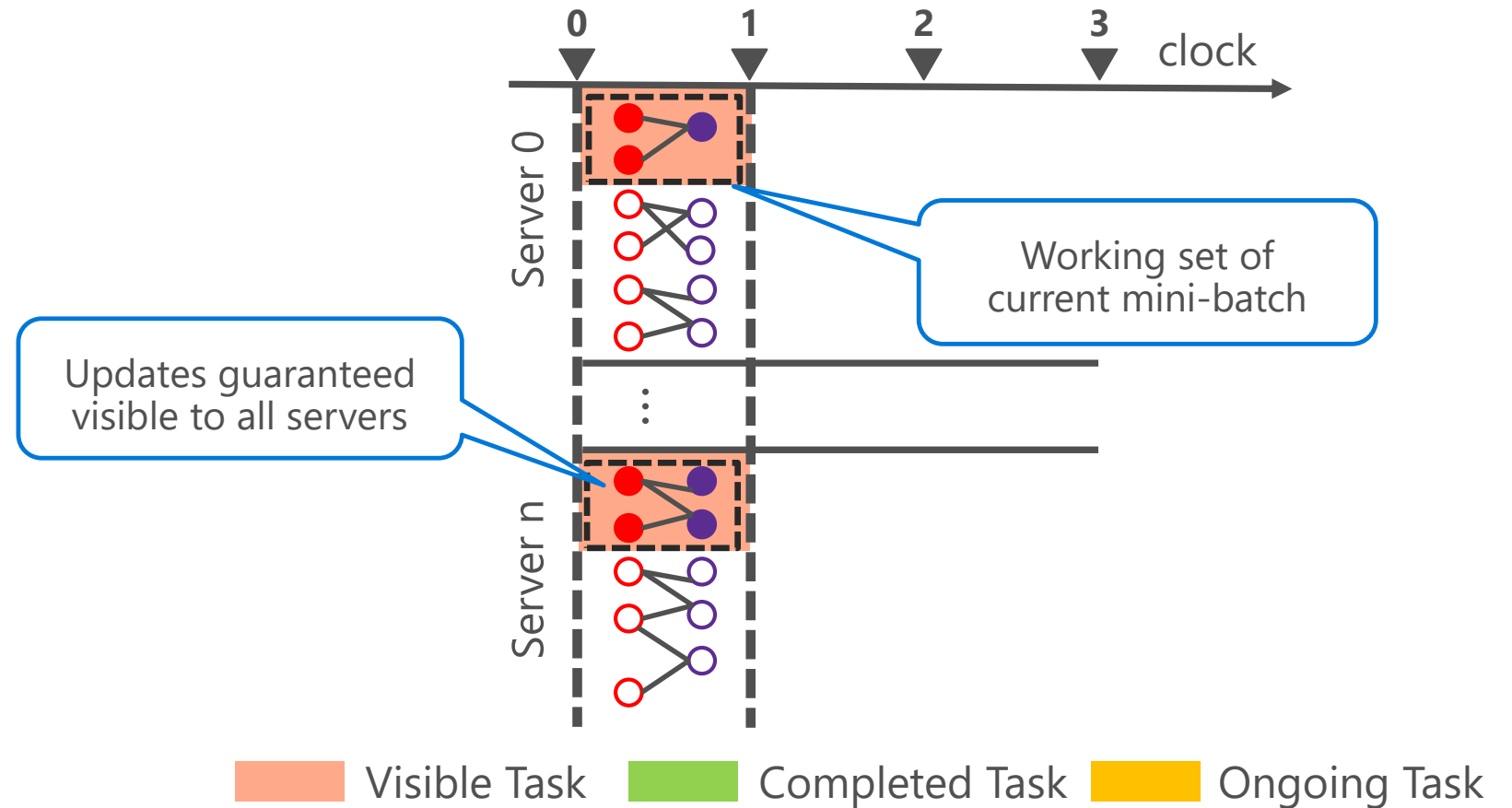
DataModel: Heterogeneous data model

Programming: MEGA graph model

# Stale Synchronous Parallel in TuX<sup>2</sup>

## Slack of 1 clock as an example

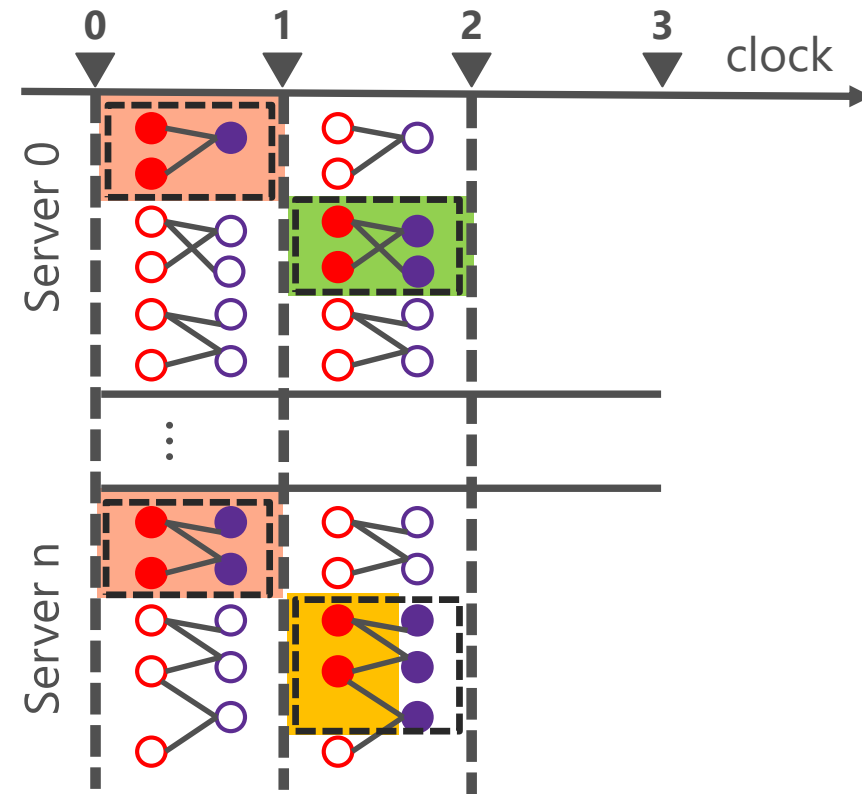
- All servers finish clock1



# Stale Synchronous Parallel in TuX<sup>2</sup>

## Slack of 1 clock as an example

- Slowest server (n) is in clock2
- Fastest server (0) finishes clock2

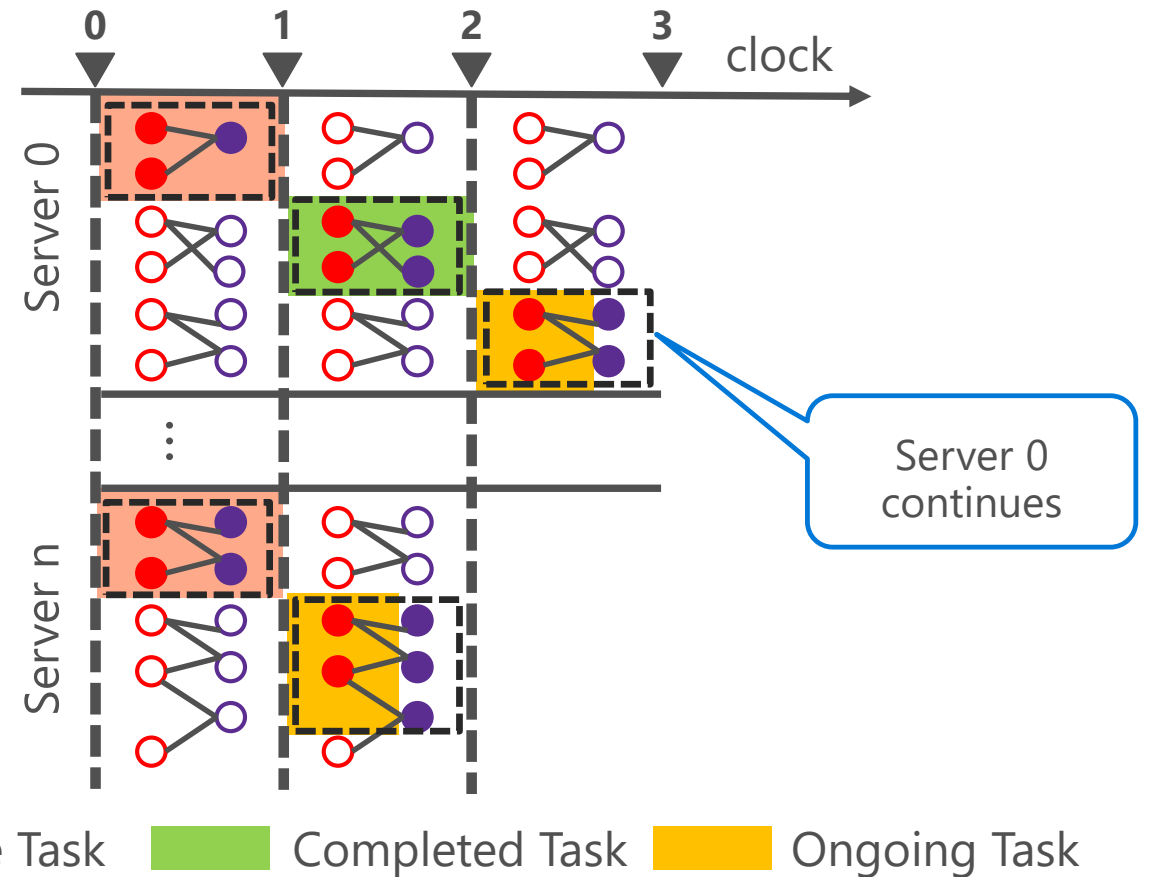


Visible Task   Completed Task   Ongoing Task

# Stale Synchronous Parallel in TuX<sup>2</sup>

## Slack of 1 clock as an example

- Slowest server (n) is in clock2
- Fastest server (0) finishes clock2
  - within the staleness bound
  - continue

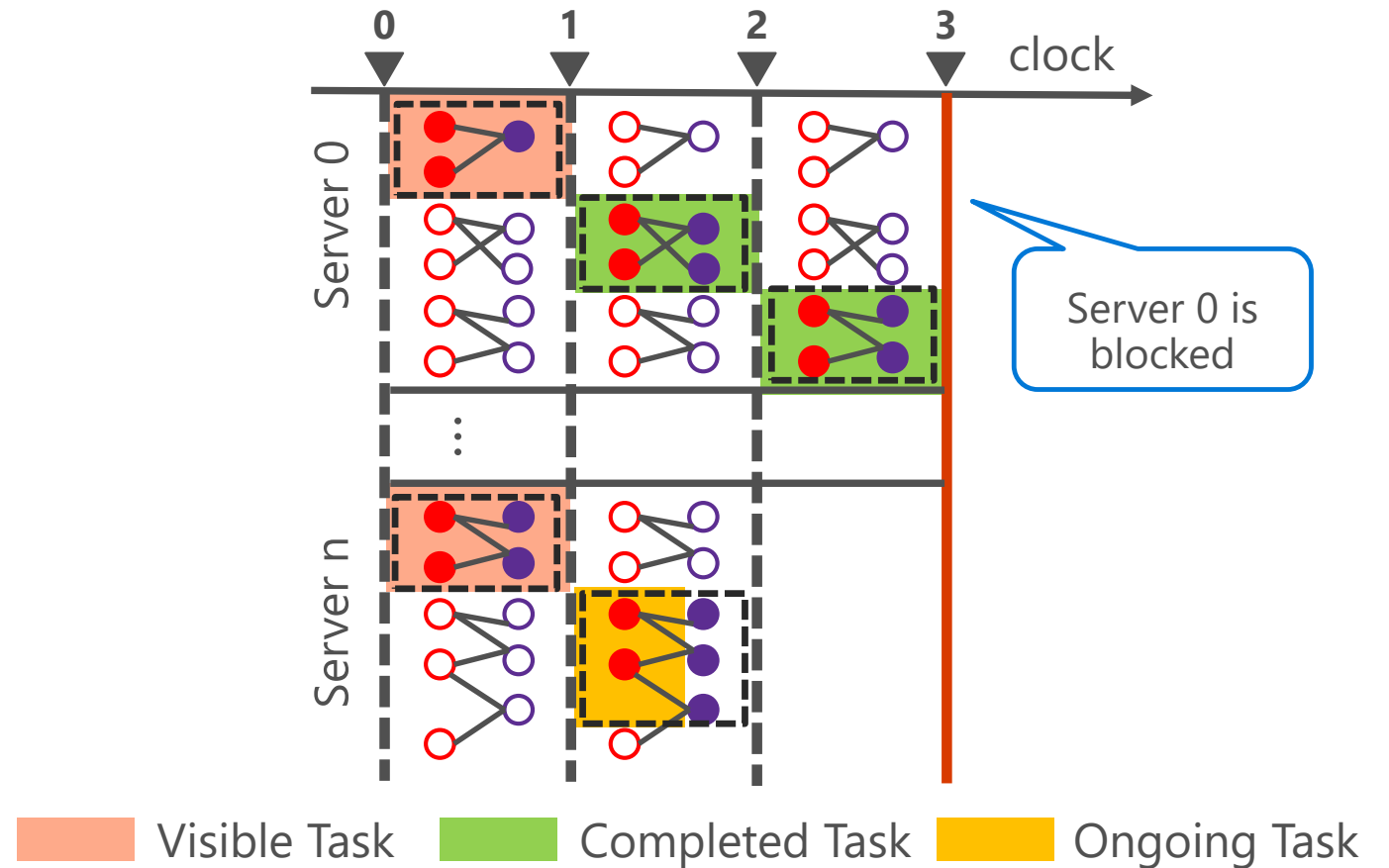




# Stale Synchronous Parallel in TuX<sup>2</sup>

## Slack of 1 clock as an example

- Slowest server (n) is in clock2
- Fastest server (0) finishes clock3
  - reaching the max slack bound
  - blocked



# Key designs

Scheduling: Stale Synchronous Parallel (SSP) based scheduling

DataModel: Heterogeneous data model

Programming: MEGA graph model

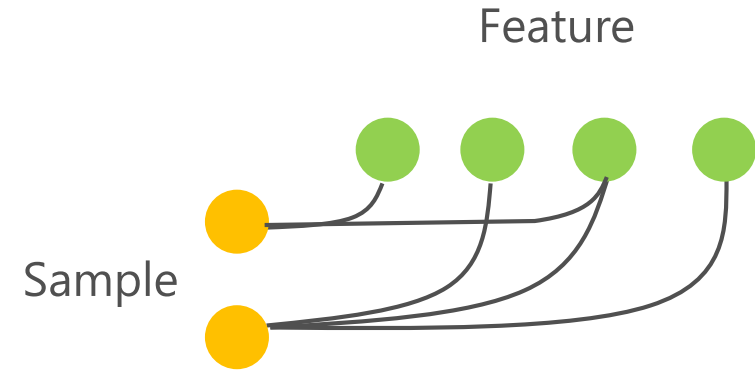
# Heterogeneity in ML

## Heterogeneous Vertices

- Different properties
  - E.g. Logistic Regression
    - Sample: Label; Feature: Weight, Gradient

## Benefit

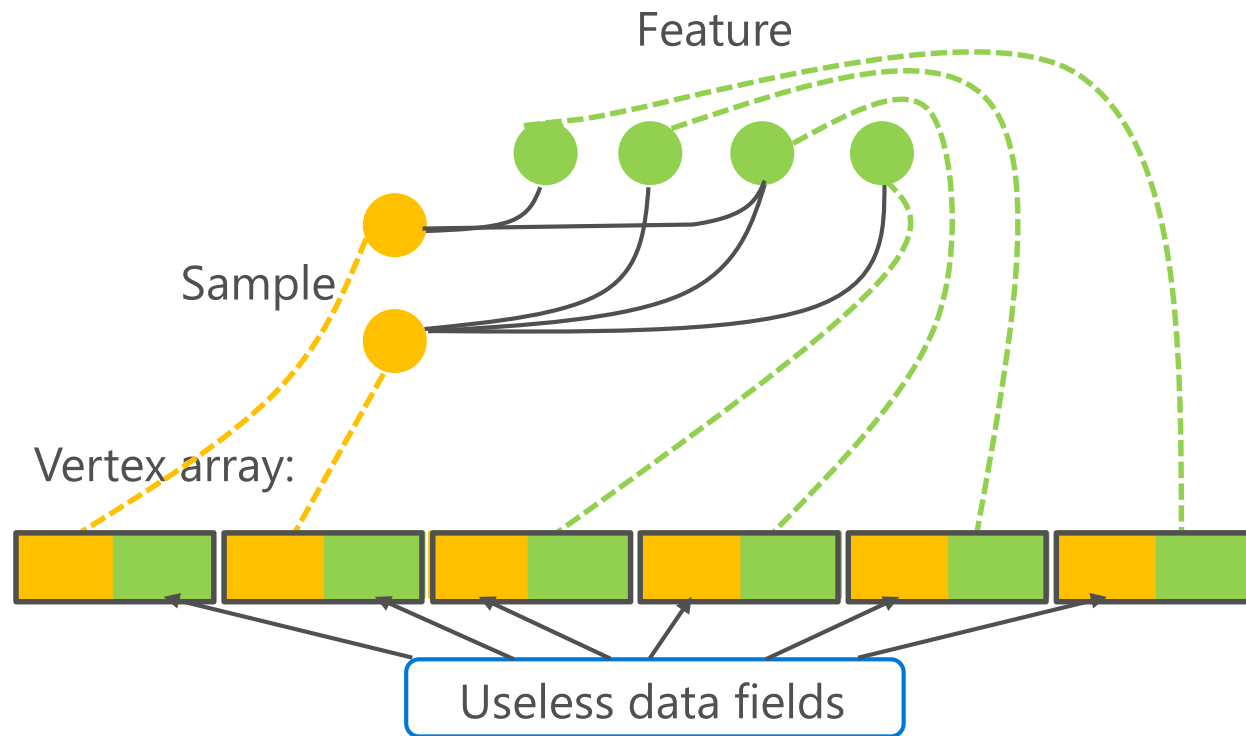
- Heterogeneity for compact data structure
  - Heterogeneity for efficient execution
- Heterogeneity for less network traffic



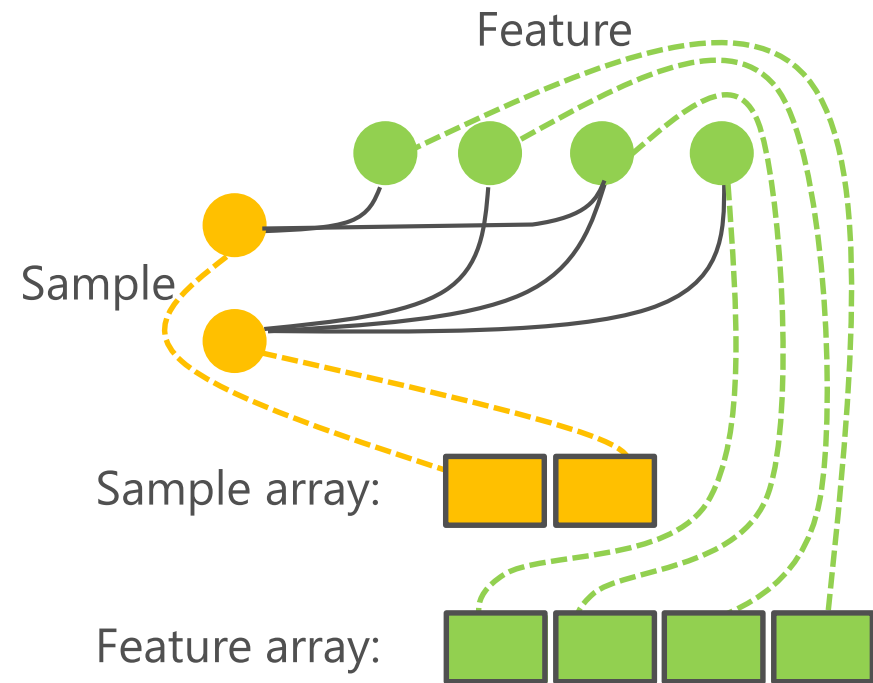
# Heterogeneity for compact data structure

## E.g. Logistic Regression

- Sample: Label; Feature: Weight, Gradient



Homogeneous Vertex Data Structure

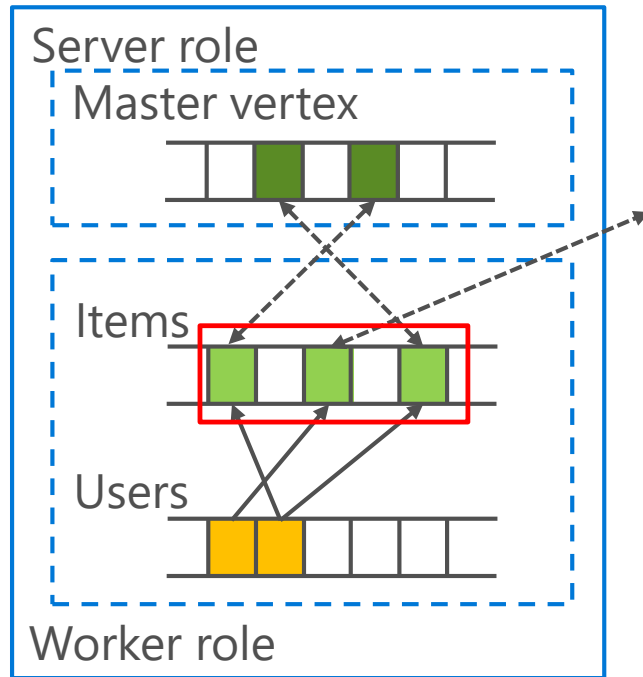


Heterogeneous Vertex Data Structure

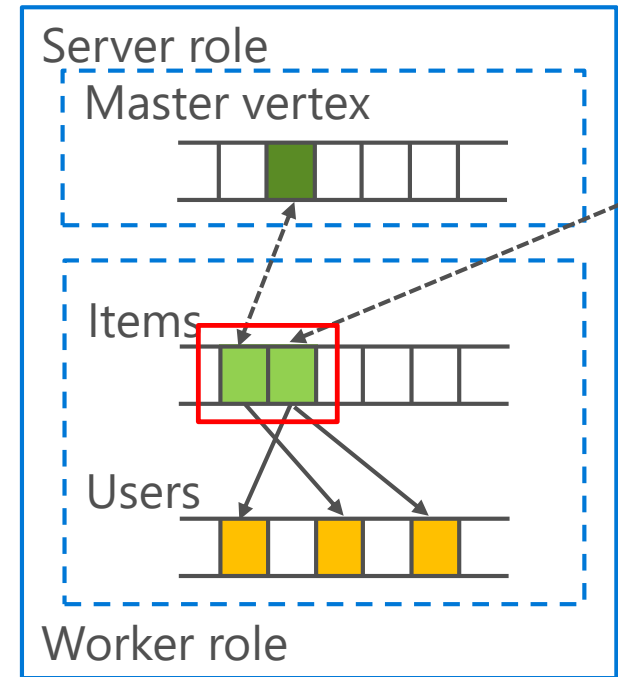
# Heterogeneity for efficient execution

## E.g. Mini-Batch MF for recommendation

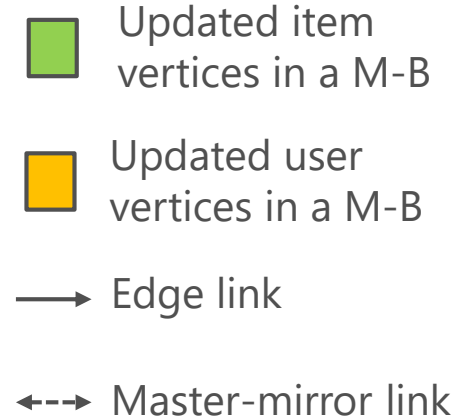
- Benefits of scanning items
  - Sequential access for locality when syncing
  - Less overhead tracing the updated vertices



Scan user vertices



Scan item vertices



# Key designs

Scheduling: Stale Synchronous Parallel (SSP) based scheduling

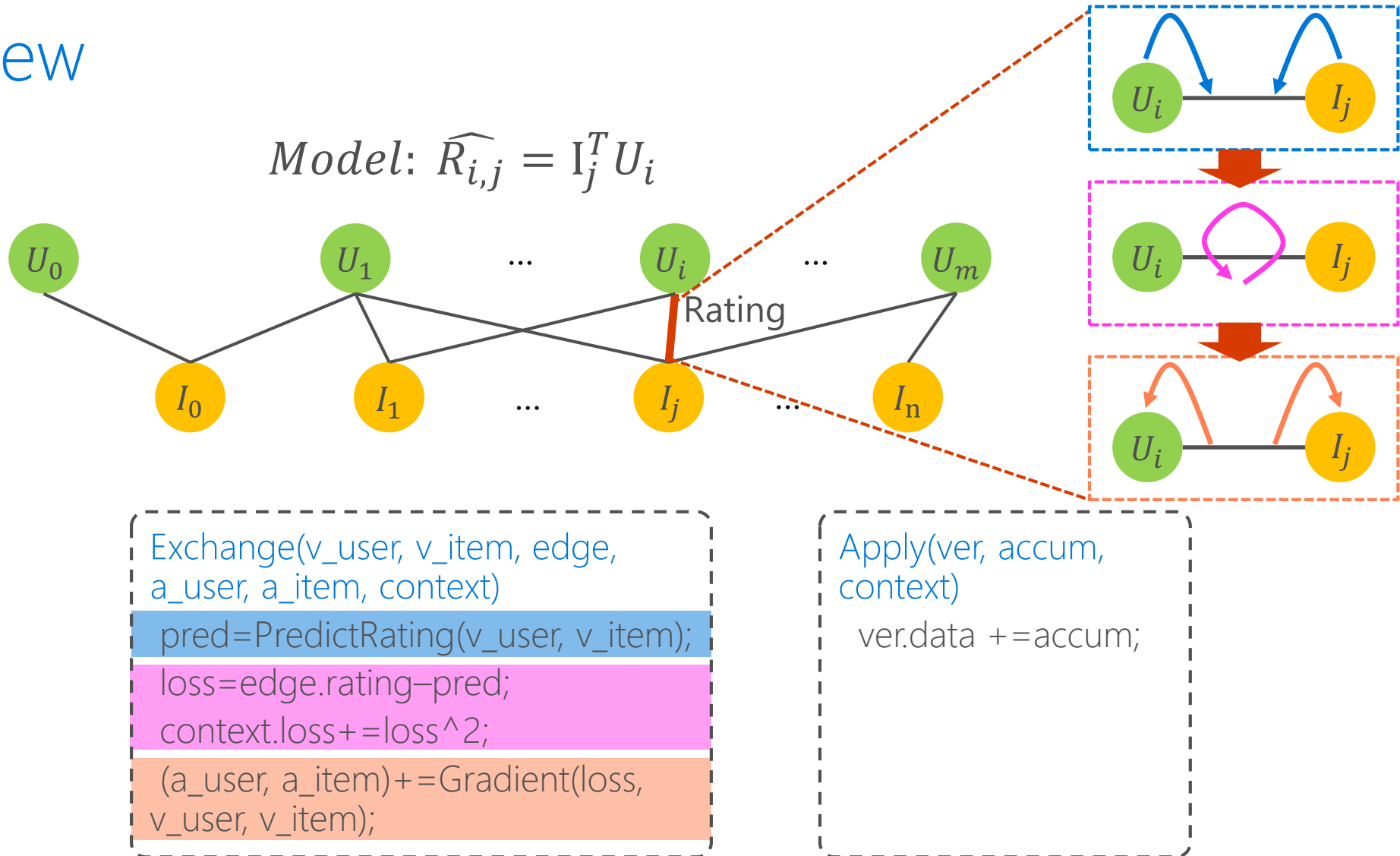
DataModel: Heterogeneous data model

Programming: MEGA graph model



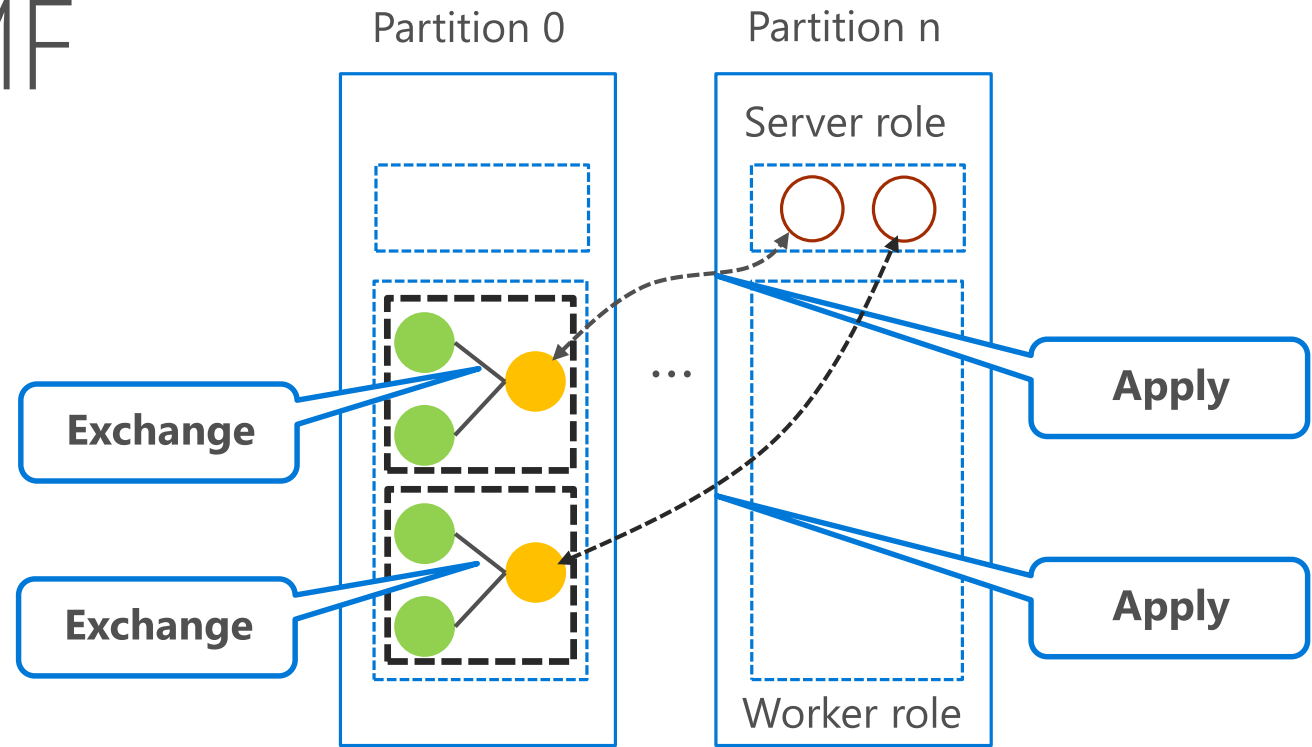
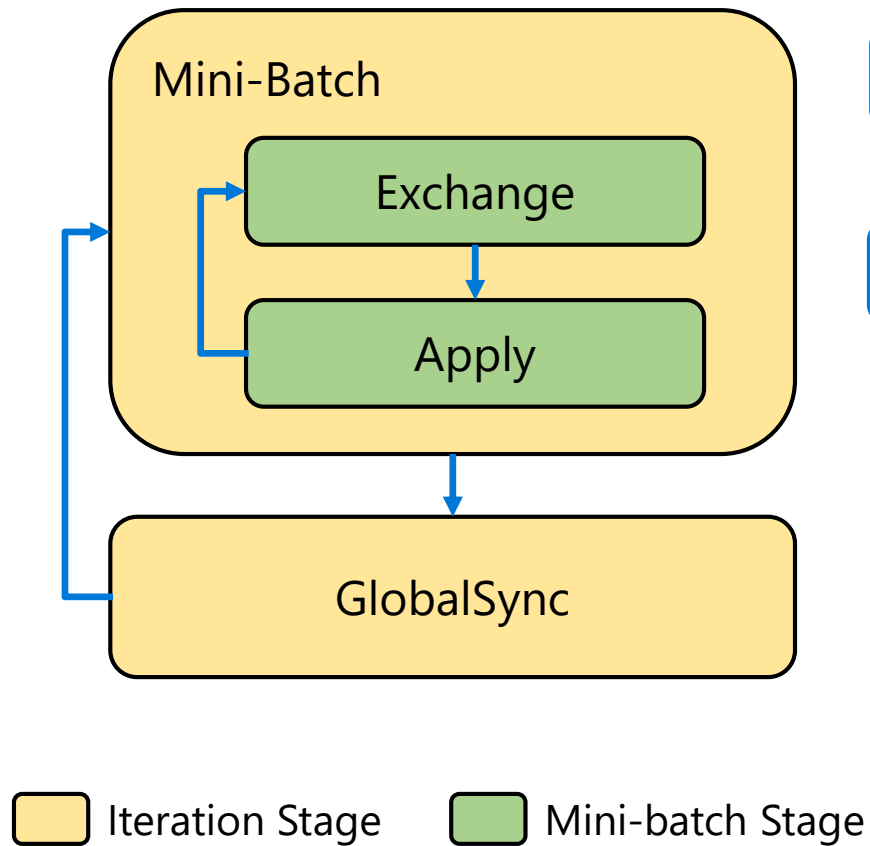
# MEGA: e.g. Mini-batch MF for recommendation

## Graph View



# Example: Mini-batch MF

## Compose stage



```
StageSequenceBuilder(ExecStages)
mbStage = new MiniBatchStage();
mbStage.SetBatchSize(100, asEdge);
mbStage.Add(ExchangeStage);
mbStage.Add(ApplyStage);
ExecStages.Add(mbStage);
ExecStages.Add(GlobalSyncStage);
```

# Experiment setup

## Machine information

- 16 CPU cores, 256GB memory, 54Gbps InfiniBand NIC

## Typical ML algorithms

- MF, LDA, BlockPG

## Large-scale dataset

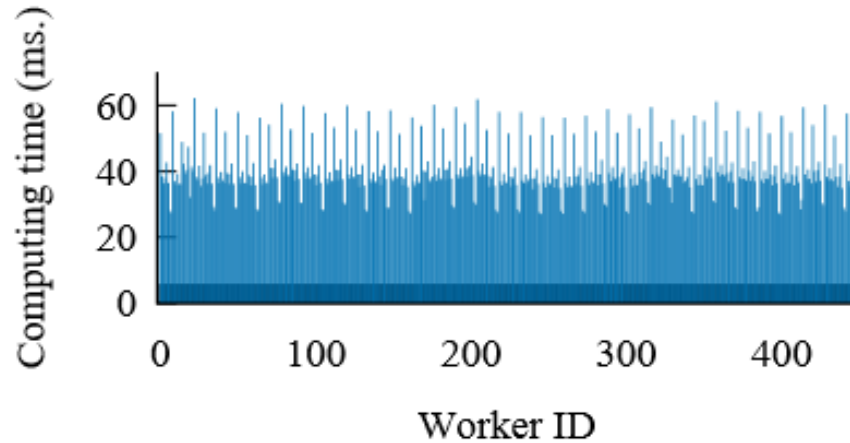
- Up to 64 billion edges graph

Dataset name	# of users/ docs/samples	# of items/ words/features	# of edges
NewsData(LDA)	7.3M	418.4K	1.4B
AdsData(BlockPG)	924.8M	209.3M	64.9B
Netflix(MF)	480.2K	17.8K	100.5M
Synthesized(MF)	30M	1M	6.3M

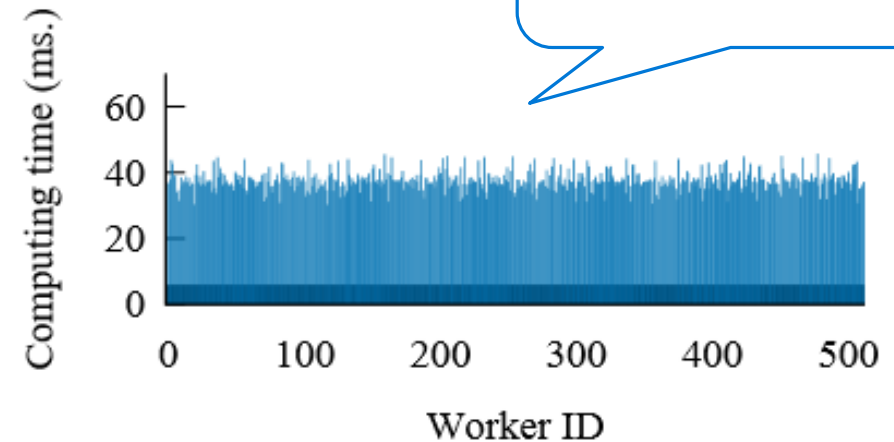
# Evaluation

## Compare to Parameter Server

- **48%** improvement on 32 servers!
- Algorithm: BlockPG
- Dataset: Microsoft private AdsData (64B edges)



Imbalance in Parameter Server

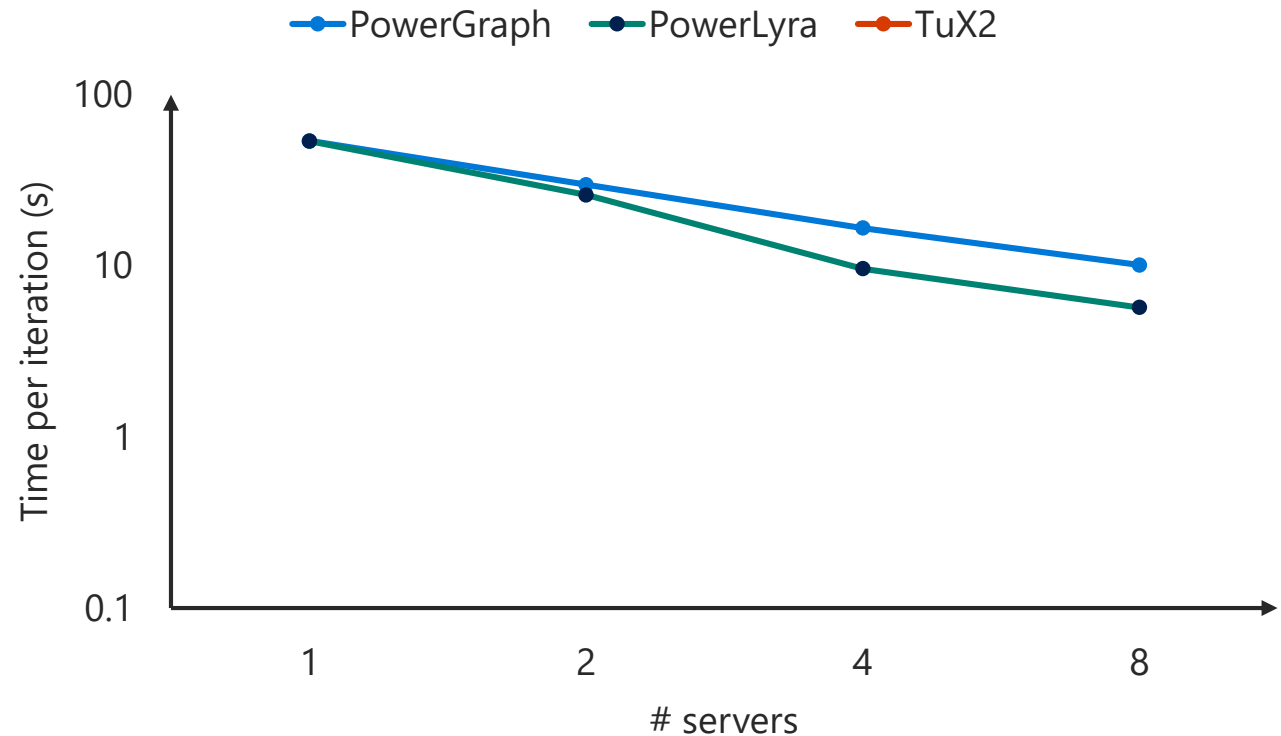


Balance in TuX2

# Evaluation

## Compare to PowerGraph, PowerLyra

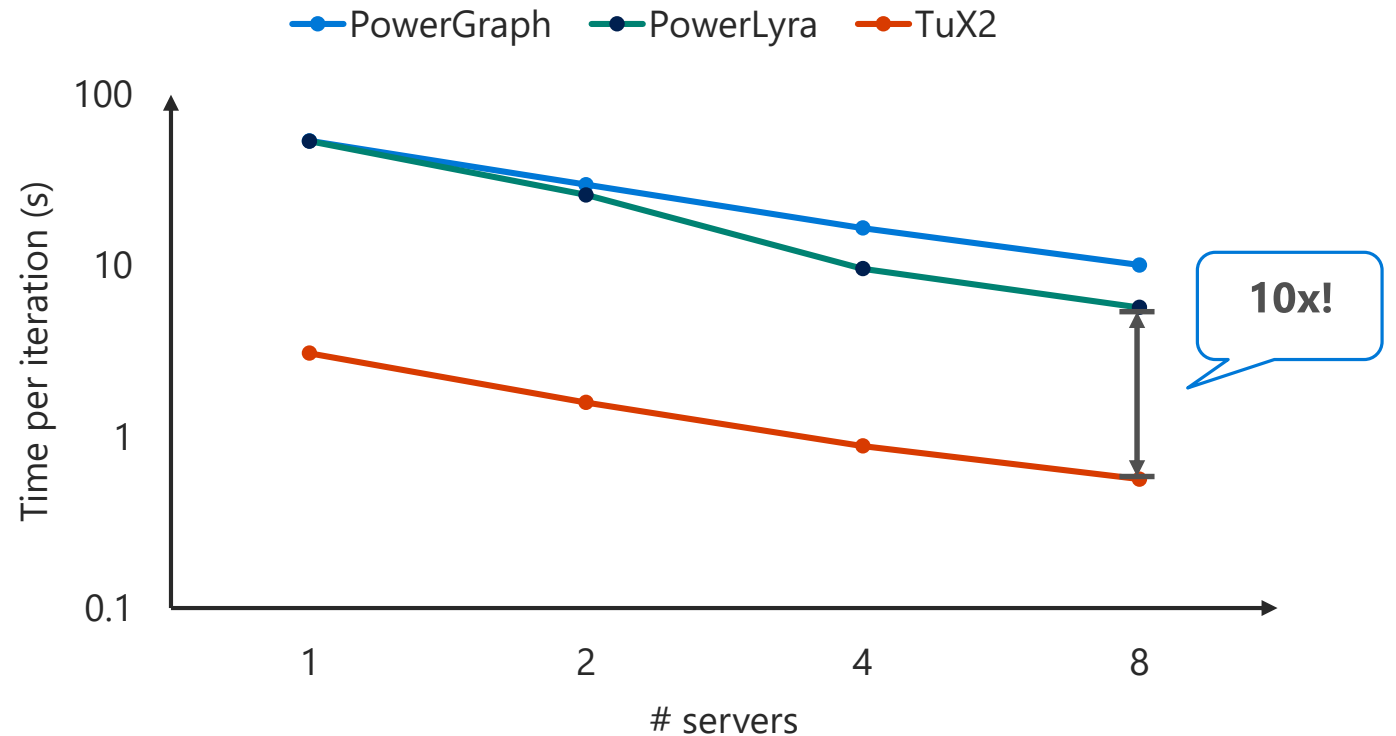
- Algorithm: Matrix Factorization
- Dataset: Netflix



# Evaluation

## Compare to PowerGraph, PowerLyra

- Algorithm: Matrix Factorization
- Dataset: Netflix





# Conclusion

TuX<sup>2</sup>: advocates the convergence of graph computation and distributed machine learning

- Introduce important machine learning concepts to graph computation
- Define a new, flexible graph model to express ML algorithms efficiently
- Demonstrate TuX<sup>2</sup> outperform existing Graph and ML systems in representative ML algorithms respectively

Thanks!  
Q&A