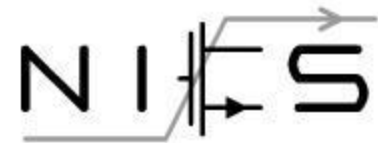# NXgraph:
# An Efficient Graph Processing System on a Single Machine

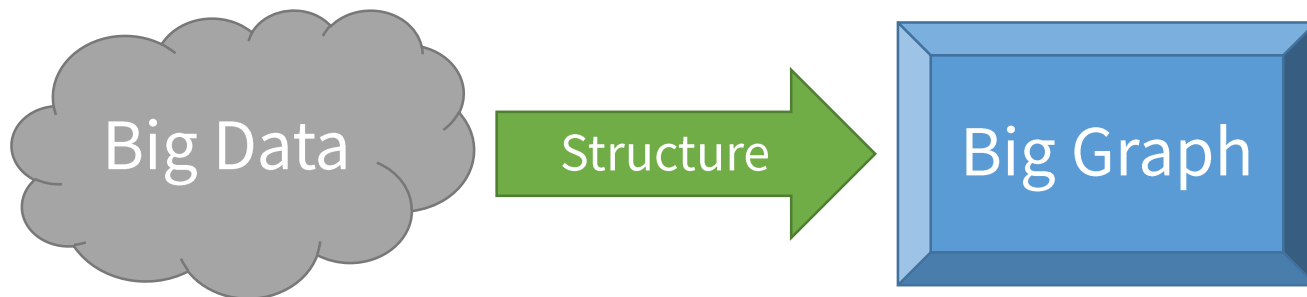Yuze Chi[1], Guohao Dai[1], Yu Wang[1], Guangyu Sun[2], Guoliang Li[1] and Huazhong Yang[1]

Tsinghua National Laboratory for Information Science and Technology[1], Tsinghua University

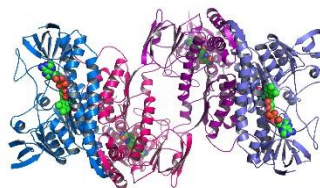Center for Energy Efficient Computing and Applications[2], Peking University

# Motivation

Big Data → Structure → Big Graph

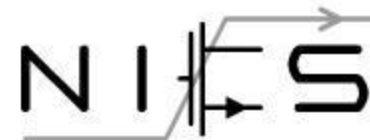social media

science

advertising

web

# Motivation

Can we exploit it well?

- Computation capacity of a single CPU
  - Intel i7-5820K: 12 hyper-threads, 3.3GHz
  - Assume 50 clock cycles/edge: **792MEPS** (Million Edges Per Second)

| System | Throughput/MEPS (PageRank on Twitter) | Notes |
|---|---|---|
| Spark[HotCloud2010] | 15 in total/0.15 each | 100 CPUs in 50 nodes |
| PowerGraph[OSDI2012] | 408 in total/6.4 each | 64 CPUs in 64 nodes |
| GraphChi[OSDI2012] | 50 | |
| TurboGraph[SIGKDD2013] | 108 | |
| X-stream[SOSP2013] | 20 | no pre-processing |
| VENUS[ICDE2015] | 15.4 | on HDD |
| GridGraph[ATC2015] | 61 | |

- Large gap!
- Our objective: **higher MEPS/CPU**

NICS

# Vertex-centric Model

- Graph G = (V,E)
  - vertex v = (id, at
  - edge e = (src vid, dst vid, onal attribute)
- "Think like a vertex"

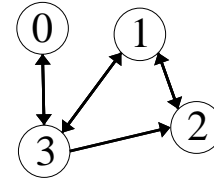> Update my neighbors...

- Example: Breadth-First Search (BFS)

```
for each dst in my.out_edges
    if dst.depth > my.depth+1
    then
        dst.depth = my.depth+1
```

- Vertices → Intervals
- Edges → Shards
  - Edges in each shard → Sub-Shards

- Objective
  - Limit memory access to a small region
  - Improve locality

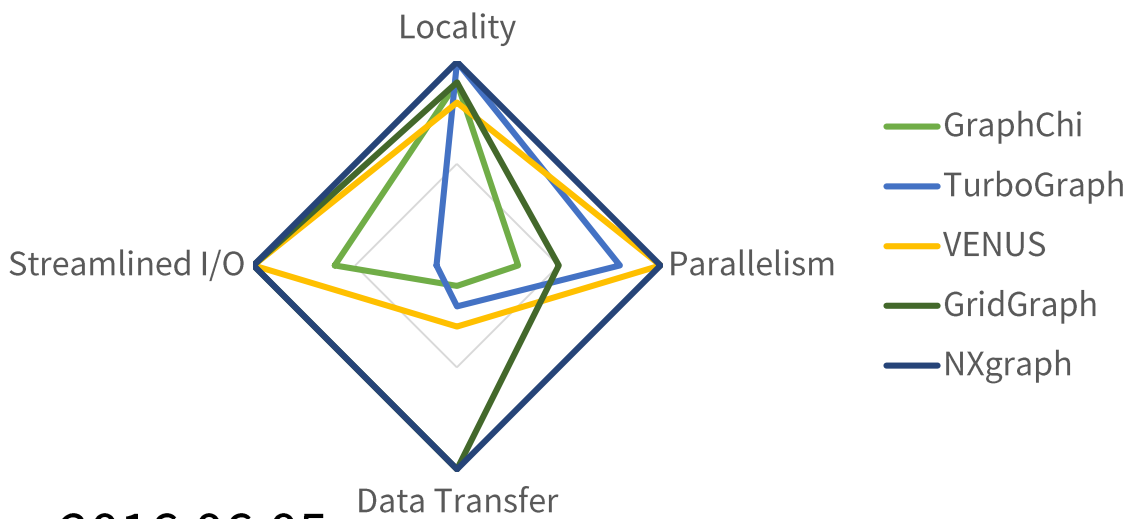| $I_1$ | $I_2$ |
|---|---|
| 0, 1 | 2, 3 |

| $S_1$ | $S_2$ |
|---|---|
| $SS_{1.1}$ | $SS_{1.2}$ |
|  | $1 \to 2$ <br> $0,1 \to 3$ |
| $SS_{2.1}$ | $SS_{2.2}$ |
| $3 \to 0$ <br> $2,3 \to 1$ | $3 \to 2$ |

# Four Optimizing Rules

| | GraphChi<br>OSDI2012 | TurboGraph<br>SIGKDD2013 | VENUS<br>ICDE2015 | GridGraph<br>ATC2015 | NXgraph<br>ICDE2016 |
|---|---|---|---|---|---|
| 1. Exploit the locality of graph data | ☺ | ☺ | ☺ | ☺ | ☺ |
| 2. Utilize the parallelism of multi-thread CPU | | ☺ | ☺ | | ☺ |
| 3. Reduce the amount of disk data transfer | | | | ☺ | ☺ |
| 4. Streamline the disk I/O | ☺ | | ☺ | ☺ | ☺ |



- GraphChi
- TurboGraph
- VENUS
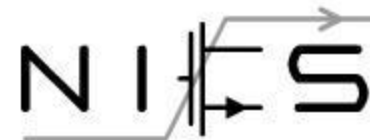- GridGraph
- NXgraph

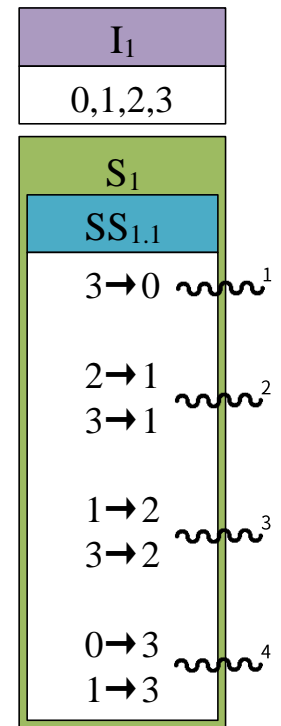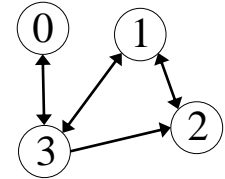No previous system considers them all

6

# Our Effort

- **Follows all four optimizing rules**

- According to optimizing rule 1 & 2, we design
  - Destination-Sorted Sub-Shard (DSSS) structure

- According to optimizing rule 3 & 4, we design
  - Adaptive updating strategies
    - Single-Phase Update (SPU)
    - Double-Phase Update (DPU)
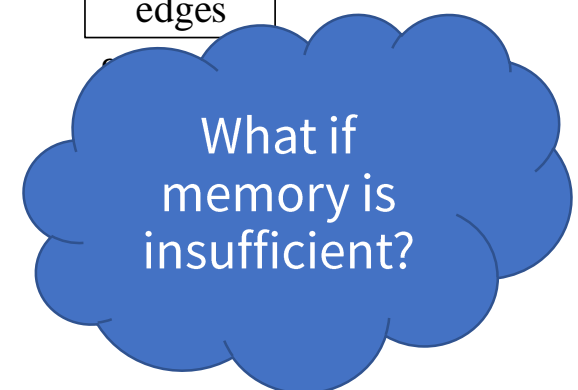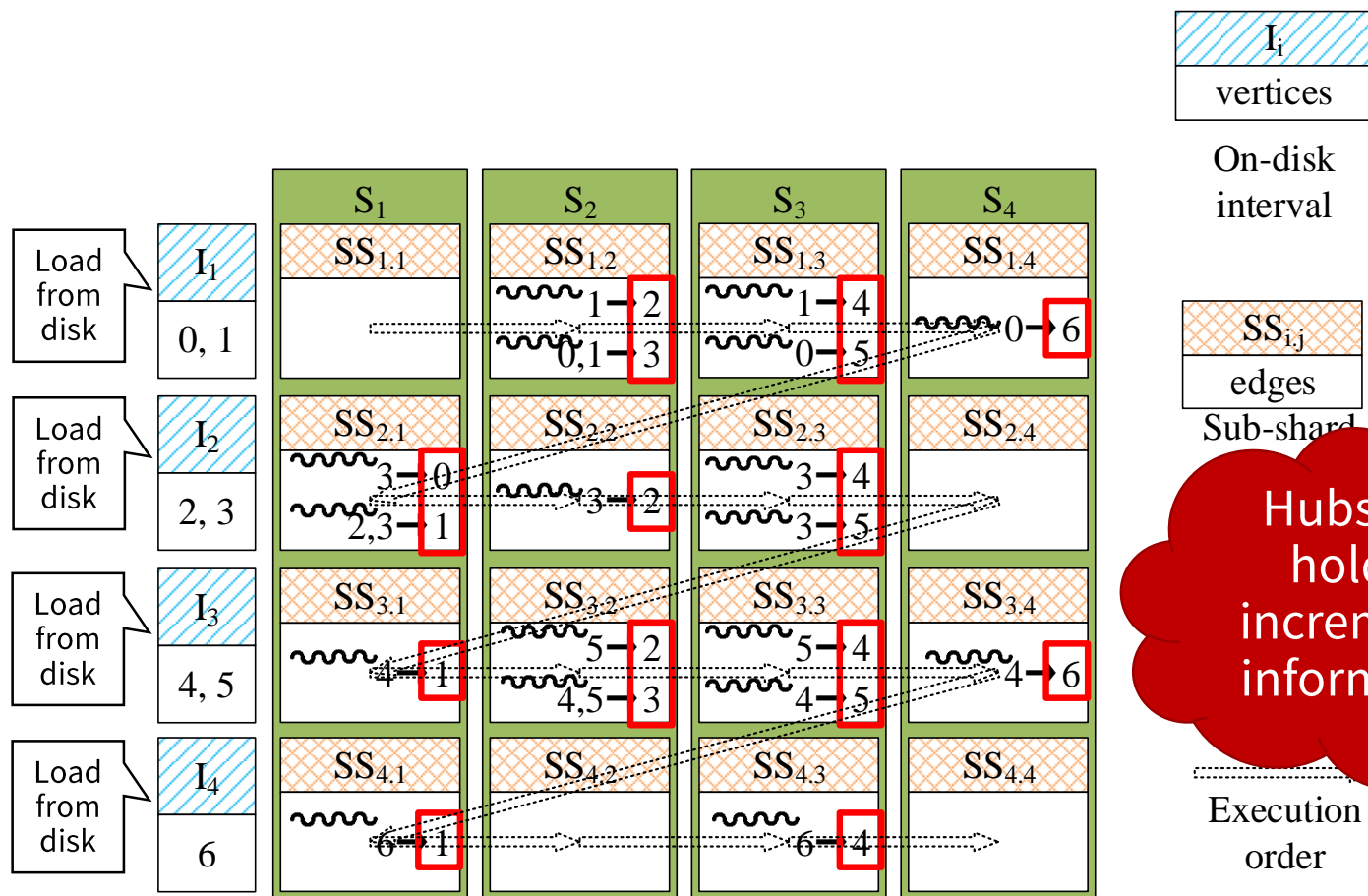    - Mixed-Phase Update (MPU)

# Destination-Sorted Sub-Shards

- Sort edges in each sub-shard

- Sequential reads from source interval (Rule 1)
  - Potentially improves cache hit rate

- Parallel writes to destination vertices (Rule 2)
  - No write conflict among threads

- Defines behavior **inside** each sub-shard
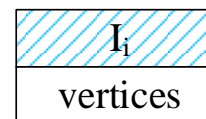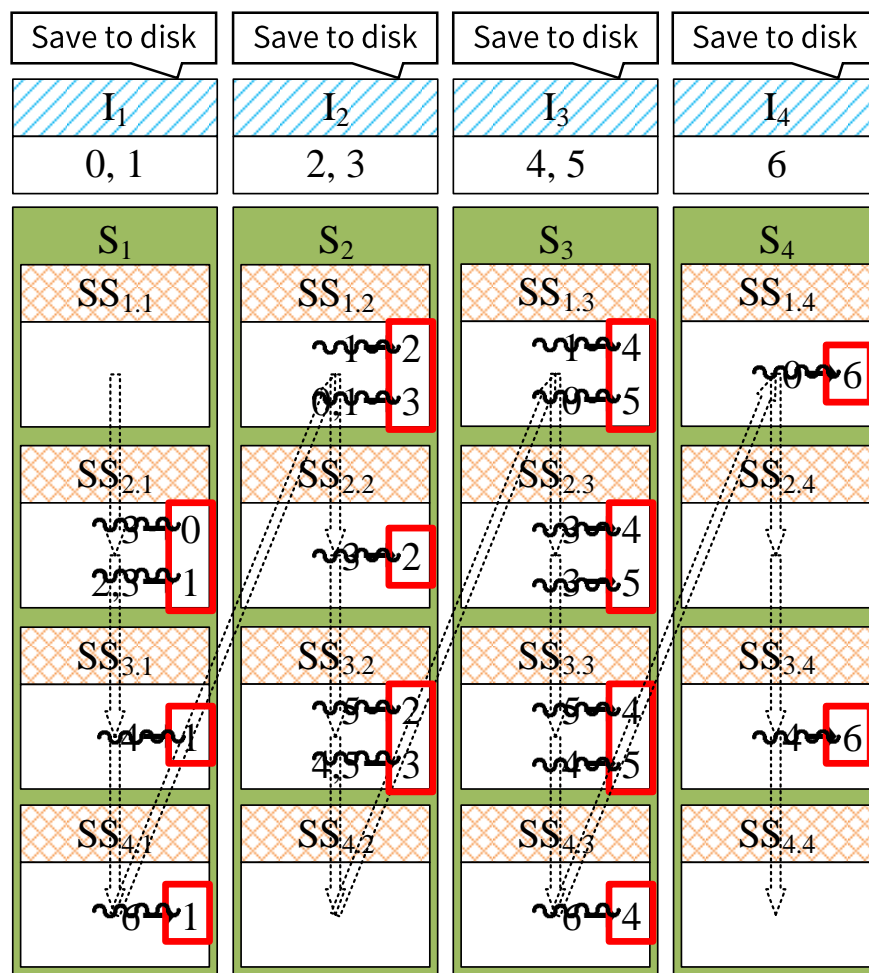


$I_1$

0,1,2,3

$S_1$

$SS_{1.1}$

$3 \rightarrow 0$ ~~~~ 1

$2 \rightarrow 1$ ~~~~ 2
$3 \rightarrow 1$

$1 \rightarrow 2$ ~~~~ 3
$3 \rightarrow 2$

$0 \rightarrow 3$ ~~~~ 4
$1 \rightarrow 3$

# Single-Phase Update

# Mixed-Phase Update

# System Architecture

NXgraph

preprocessing

updater

indexer    sharder

## indexer

| generate vertex mapping | generate binary edge list | calculate vertex degree |

## sharder

| put edges into each shard | sort edges in each shard |

## updater

calculate output in a user-defined format with user-defined input

# Results

- Evaluation platform

  - Hex-core hyper-threading Intel i7-5820K CPU @ 3.3GHz

  - 8x8G DDR4 RAM, 2x128G RAID0 SSD, 1T HDD

  - Ubuntu 14.04 LTS 64bit/Windows 10 Edu 64bit

# Results: Design Decisions

Performance with different sub-shard model

| Model | Elapsed Ti... | |
|---|---|---|
| | Live-journal | Tw... |
| src-sorted, coarse-grained | 1.44s | 72.... |
| dst-sorted, fine-grained | 1.00s | 20.50s |

> Destination-sorted is the right choice!

Performance with different numbers of interval on Twitter

> ~10s intervals are near-optimal choices

- PageRank
- BFS
- SCC



2016.06.05

15

## SPU vs DPU on performance

# Results: Different Environments

### PageRank on Live-journal
(Elapsed Time (sec./10iters) vs Memory Size (GB))

### PageRank on Twitter
(Elapsed Time vs Memory Size (GB))

### PageRank on Yahoo-web
(Elapsed Time vs Memory Size (GB))

Up to **13.64x** speedup over TurboGraph

Legend:
- NXgraph
- GraphChi
- Turbograph

### PageRank on Live-journal
(Elapsed Time (sec./10iters) vs Number of Threads)

### PageRank on Twitter
(Elapsed Time vs Number of Threads)

### PageRank on Yahoo-web
(Elapsed Time vs Number of Threads)

Legend:
- NXgraph
- GraphChi
- TurboGraph

Scalability

More tasks on Live-journal

More tasks on Twitter

More tasks on Yahoo-web

Best scalability

Average speedup of **74.6x**

2016.06.05

18

## System performance with limited resources

| System | Time (s) | Speedup | Evaluation environment |
|---|---|---|---|
| NXgraph | 7.13 | 1.00 | Intel i7 3.3GHz, 8t, 8G, SSD |
| GridGraph | 26.91 | 3.77 | AWS EC2 8t, 8G/30.5G, SSD |
| X-stream | 88.95 | 12.48 | |
| NXgraph | 12.55 | 1.00 | |
| VENUS | 95.48 | 7.60 | HDD |
| GridGraph | 24.11 | 1.92 | /30.5G, HDD |
| X-stream | 81.70 | 6.51 | AWS EC2, 8t, 8G/30.5G, HDD |

Average speedup of **6.6x** over various state-of-the-art single-machine systems with limited resources

Task: 1 iteration of PageRank on Twitter graph

2016.06.05

19

## System performance in the best case

| System | Time (s) | Speedup | Evaluation environment |
|--------|----------|---------|------------------------|
| NXgraph | 2.05 | 1.00 | Intel i7 3.3GHz, 8t, 16G, SSD |
| X-stream | 23.25 | 11.57 | ... SSD |
| GridGraph | 24.11 | 11.99 | |
| MMAP | 13.10 | 6.52 | |
| PowerGraph | 3.60 | 1.79 | (16t, 23G) |

**717MEPS** actual throughput vs **792MEPS** hypothetical limit

Task: 1 iteration of PageRank on Twitter graph

# Future Work

- NXgraph is still under development and subject to changes in data structures and APIs

v0.2 (current)

v0.3 (scheduled)

Contact us if you would like to see a pre-release!

timeline

✓ Faster sharding (done)
✓ More flexible partitioning (done)
➢ Faster indexing (active)
➢ More complicated algorithms (active)
  Dynamic graph support (scheduled)
  Open source (scheduled: **mid August**)

  Asynchronous updating (scheduled)
  More intelligent partitioning (pending)
  Search for practical applications (pending)

# Reference

- H. Kwak, C. Lee, H. Park, and S. Moon, "What is twitter, a social network or a news media?" in *WWW*. ACM, 2010, pp. 591–600.

- J. Gonzalez, Y. Low, and H. Gu, "Powergraph: Distributed graphparallel computation on natural graphs," in *OSDI*, 2012, pp. 17–30.

- J. Cheng, Q. Liu, Z. Li, W. Fan, J. C. S. Lui, and C. He, "VENUS: Vertex-Centric Streamlined Graph Computation on a Single PC," in *ICDE*, 2015, pp. 1131–1142.

- A. Kyrola, G. Blelloch, and C. Guestrin, "GraphChi: Large-Scale Graph Computation on Just a PC," in *OSDI*, 2012, pp. 31–46.

- W.-S. Han, S. Lee, K. Park, J.-H. Lee, M.-S. Kim, J. Kim, and H. Yu, "TurboGraph: A Fast Parallel Graph Engine Handling Billion-Scale Graphs in a Single PC," in *SIGKDD*, 2013, pp. 77–85.

- M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica, "Spark: cluster computing with working sets," in *HotCloud*, vol. 10, 2010, p. 10.

- L. Page, S. Brin, R. Motwani, and T. Winograd, "The pagerank citation ranking: bringing order to the web." 1999.

- X. Zhu, W. Han, and W. Chen, "GridGraph : Large-Scale Graph Processing on a Single Machine Using 2-Level Hierarchical Partitioning," in *ATC*, 2015, pp. 375–386.

- A. Roy, I. Mihailovic, and W. Zwaenepoel, "X-Stream: Edge-centric Graph Processing using Streaming Partitions," in *SOSP*, 2013, pp. 472–488.

# Reference

- Yahoo! altavisata web page hyperlink connectivity graph, circa 2002,"
- http://webscope.sandbox.yahoo.com/.
- "Livejournal social network," http://snap.stanford.edu/data/ soc-LiveJournal1.html.
- Z. Lin, M. Kahng, K. Sabrin, D. Horng, and P. Chau, "MMap : Fast Billion-Scale Graph Computation on a PC via Memory Mapping," in *ICBD*. IEEE, 2014.
- G. Malewicz, M. H. Austern, A. J. C. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski, "Pregel : A System for Large-Scale Graph Processing," in SIGMOD, 2010, pp. 135–145.

# Thank you!
## Q&A