# Callisto-RTS: Fine-Grain Parallel Loops

Tim Harris, Oracle Labs
Stefan Kaestle, ETH Zurich
Daniel Goodman, Oracle Labs

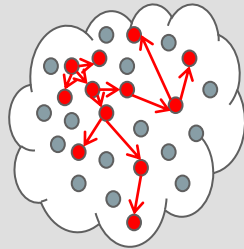15 July 2016

ORACLE®

**ORACLE®**

# In-memory graph analytics

Using a graph representation for your data

# In-memory graph analytics



Using a graph representation for your data

enables many interesting new analyses

Product Recommendation

Purchase Record

customer items

Influencer Identification

Communication

Community Detection

Pattern Matching

# In-memory graph analytics

Using a graph representation for your data



and eliminates repeated join operations, which is much more efficient when you have a lot of relationships to traverse

enables many interesting new analyses

**Product Recommendation**

Purchase Record

customer   items

**Influencer Identification**

Communication

**Community Detection**

**Pattern Matching**

# PageRank inner loop

# PageRank inner loop

# PageRank inner loop

# PageRank inner loop

# Hardware options

# Hardware options



My laptop
- Insufficient RAM
- Insufficient CPU capacity

Non-starter

# Hardware options



My laptop
- Insufficient RAM
- Insufficient CPU capacity

Non-starter



Cluster
- Enough RAM
- Enough CPU capacity
- Distributed memory model

Irregular memory accesses
make it hard to program

# Hardware options



**My laptop**
- Insufficient RAM
- Insufficient CPU capacity

Non-starter

**Cluster**
- Enough RAM
- Enough CPU capacity
- Distributed memory model

Irregular memory accesses
make it hard to program

**Large shared memory machine**
- Enough RAM
- Enough CPU capacity
- Shared memory model

# In-memory graph analytics

**Domain specific languages**
- Queries expressed in terms of graph concepts
- Tailor for different kinds of workload (e.g., sub-graph isomorphism)

**Generated code**
- Efficient in-memory data representations, e.g. compressed-sparse-row format
- Abundant parallelism

**Runtime system**
- Allocation of resources to a query
- Distribution of work and data within a machine

```
parallel_for<node_t>([&](node_t n) {
    …
});
```

| RAM | CPU | CPU | RAM |
|-----|-----|-----|-----|
| RAM | CPU | CPU | RAM |

# Batch size / load imbalance trade-off

# Batch size / load imbalance trade-off



Divide iteration space evenly between threads and get good load balancing

Fixed amount of work in each iteration

Iteration execution time

Iteration number

# Batch size / load imbalance trade-off



Iteration execution time

Iteration number

Variable amount of work per iteration

(Actual data – #out-edges of the top 1000 nodes in the SNAP Twitter dataset)

ORACLE®

# Batch size / load imbalance trade-off



Divide into large batches

Reduce contention distributing work
Risk load imbalance

Divide into small batches

Increase contention distributing work
Achieve better load balance

# Batch size / load imbalance trade-off

Typically, choose manually –
but getting this right
depends on (1) algorithm,
(2) machine, (3) data

Divide into large batches

Reduce contention distributing work
Risk load imbalance

Divide into small batches

Increase contention distributing work
Achieve better load balance

# PageRank – SNAP LiveJournal (4.8M vertices, 69M edges)

**OpenMP static & dynamic loops**



8-socket SPARC T5
16 cores per socket
8 h/w threads per core

PageRank
SNAP LiveJournal data set

Best performance: 0.26s

# My laptop

1 socket

# My laptop

1 socket

4 cores per socket

# My laptop



1 socket

4 cores per socket

2 h/w contexts per core

# My laptop

Counter

1 socket

4 cores per socket

2 h/w contexts per core

# My laptop

Counter

1 socket

Counter

4 cores per socket

2 h/w contexts per core

# My laptop

Counter

1 socket

Counter

4 cores per socket

2 h/w contexts per core

# My laptop

Counter

1 socket

Counter

4 cores per socket

2 h/w contexts per core

# T5-8

8 sockets

# T5-8



8 sockets

16 cores per socket

# T5-8

8 sockets

16 cores per socket

8 h/w contexts per core

# T5-8

Counter

8 sockets

16 cores per socket

8 h/w contexts per core

# T5-8

# T5-8

Counter

8 sockets

16 cores per socket

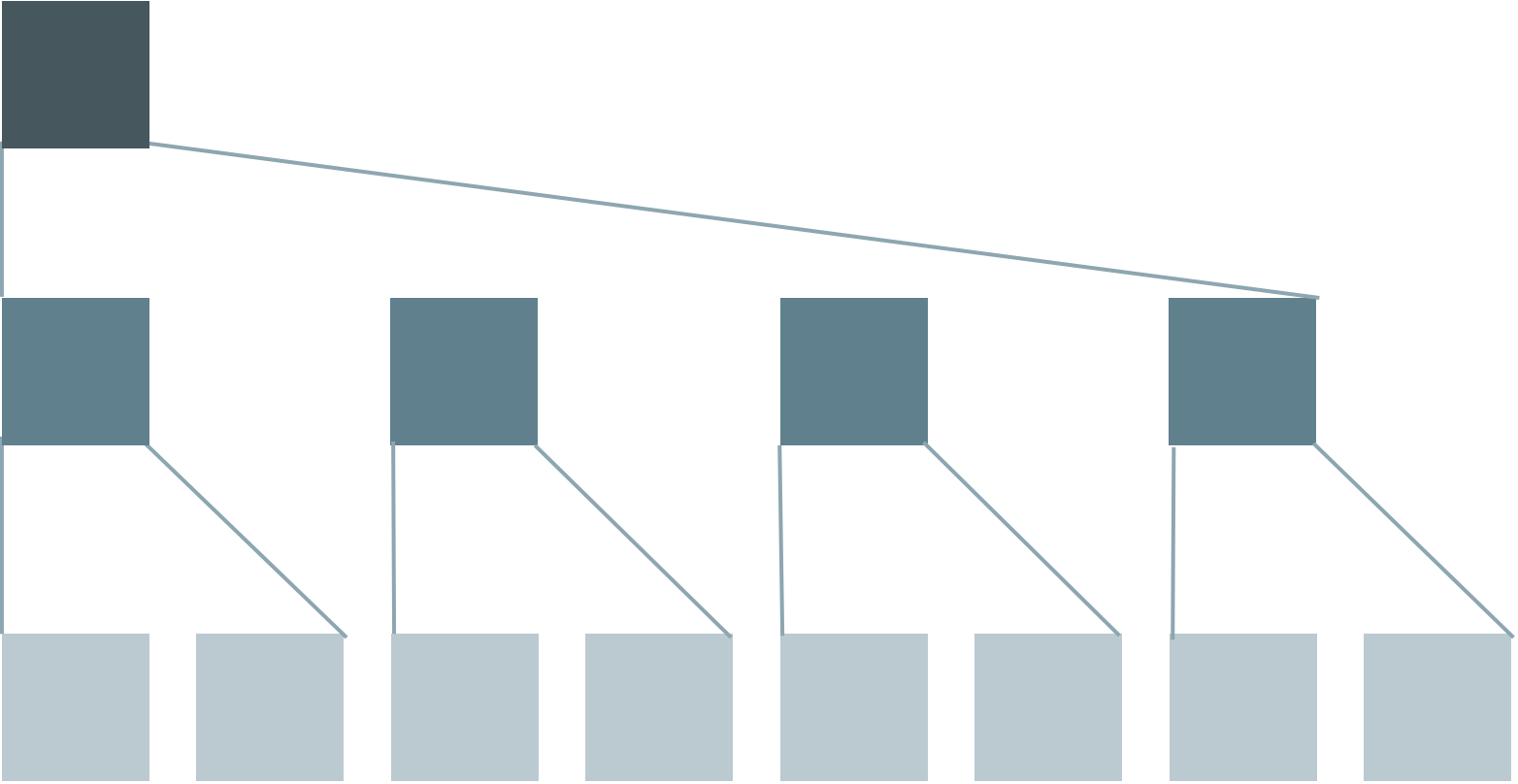8 h/w contexts per core

# T5-8

Counter

8 sockets

16 cores per socket

8 h/w contexts per core

# The problem

## My laptop

- 8 Threads accessing the counter

- The counter is always on the required socket

- 1 time in 4 the counter is on the required core

## T5-8

- 1024 Threads accessing the counter

- 1 time in 8 the counter is on the required socket

- 1 time in 128 the counter is on the required core

ORACLE®

# PageRank – SNAP LiveJournal (4.8M vertices, 69M edges)

OpenMP

Callisto-RTS

Threads: 32, 64, 128, 256, 512, 1024

Batch size: 1024, 256, 64, 16, 4

Normalized execution time: 4.0, 3.5, 3.0, 2.5, 2.0, 1.5, 1.0

# PageRank – SNAP LiveJournal (4.8M vertices, 69M edges)

OpenMP

Callisto-RTS

Threads: 32, 64, 128, 256, 512, 1024

Batch size: 1024, 256, 64, 16, 4

Normalized execution time: 4.0, 3.5, 3.0, 2.5, 2.0, 1.5, 1.0

17% improvement in best-case performance

# Batch size / load imbalance trade-off

Typically, choose manually –
but getting this right
depends on (1) algorithm,
(2) machine, (3) data

Our approach: support
efficient small batches

Divide into large batches

Reduce contention
Risk load imbalance

Divide into small batches

Increase contention distributing work
Achieve better load balance

**ORACLE®**

# Techniques

**1** ▶ Request combining

**2** ▶ Asynchronous work requests

ORACLE®

# Techniques

**1** ▶ Request combining

**2** ▶ Asynchronous work requests

# Approach



8 sockets

16 cores per socket

8 h/w contexts per core

# Approach

Per-socket iteration counters, reducing communication between sockets. Steal from other sockets when own work complete.

8 sockets

16 cores per socket

8 h/w contexts per core

# Approach

Per-socket iteration counters, reducing communication between sockets. Steal from other sockets when own work complete.

8 sockets

Aggregate requests within a socket, reducing contention on per-socket counter.

16 cores per socket

8 h/w contexts per core

# Approach

Per-socket iteration counters, reducing communication between sockets. Steal from other sockets when own work complete.

8 sockets

Aggregate requests within a socket, reducing contention on per-socket counter.

16 cores per socket

Further aggregation within a core, exploiting shared cache.

8 h/w contexts per core

# Approach, consider a loop 0..65536, batch size 8



8 sockets

Distribute iterations at start of loop down to per-core counters

16 cores per socket

8 h/w contexts per core

ORACLE®

# Approach, consider a loop 0..65536, batch size 8



8 sockets

Distribute iterations at start of loop down to per-core counters

0..512   512..1024

16 cores per socket

8 h/w contexts per core

# Approach, consider a loop 0..65536, batch size 8



8 sockets

**Distribute iterations at start of loop down to per-core counters**

0..512

512..1024

16 cores per socket

**Aggregate requests upwards within a core**

8 h/w contexts per core

Per-core lock

Per-thread request flags

# Approach, consider a loop 0..65536, batch size 8



8 sockets

Distribute iterations at start of loop down to per-core counters

0..512

512..1024

16 cores per socket

Aggregate requests upwards within a core

8 h/w contexts per core

Per-core lock

# Approach, consider a loop 0..65536, batch size 8



8 sockets

Distribute iterations at start of loop down to per-core counters

16 cores per socket

Aggregate requests upwards within a core

8 h/w contexts per core

0..512

512..1024

Per-core lock

# Approach, consider a loop 0..65536, batch size 8



8 sockets

Distribute iterations at start of loop down to per-core counters

16 cores per socket

0..512

512..1024

Aggregate requests upwards within a core

8 h/w contexts per core

Per-core lock

# Approach, consider a loop 0..65536, batch size 8



8 sockets

Distribute iterations at start of loop down to per-core counters

0..512    512..1024

16 cores per socket

Aggregate requests upwards within a core

8 h/w contexts per core

Per-core lock

# Approach, consider a loop 0..65536, batch size 8



8 sockets

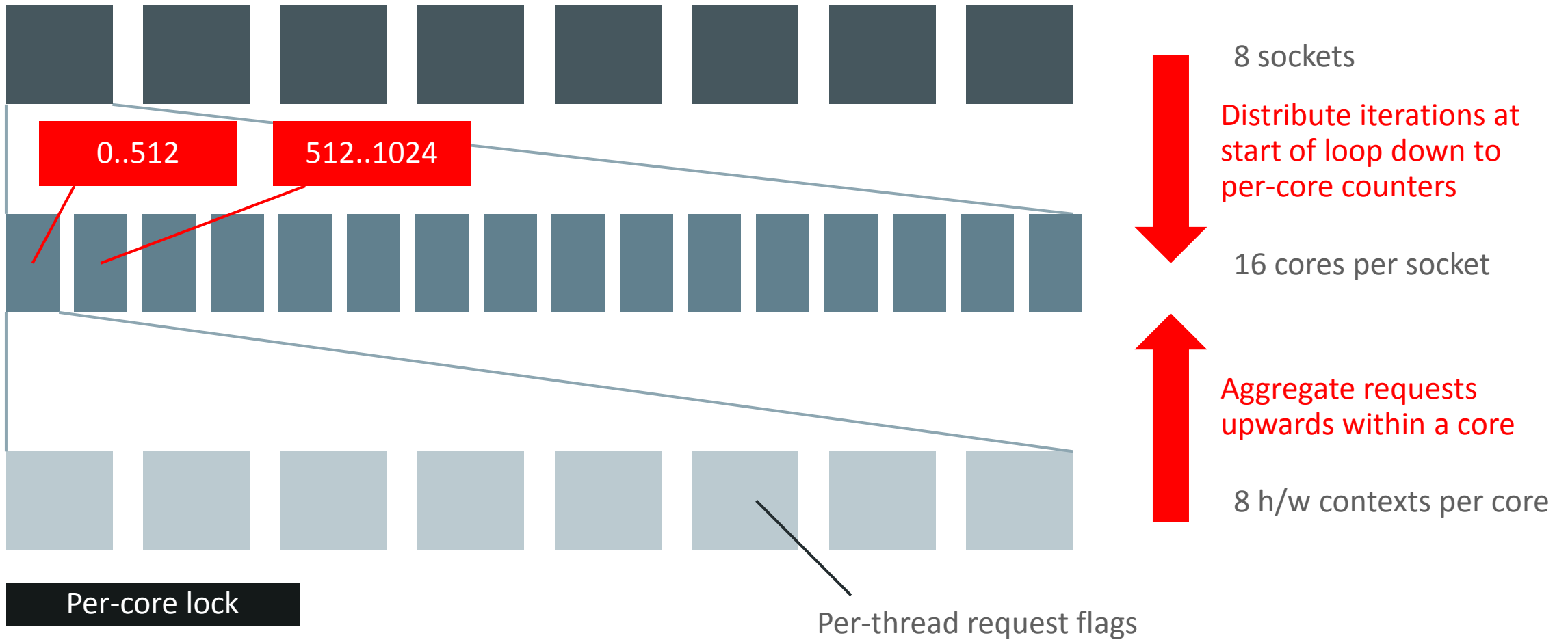Distribute iterations at start of loop down to per-core counters

0..512    512..1024

3*8

16 cores per socket

Aggregate requests upwards within a core

8 h/w contexts per core

Per-core lock

# Approach, consider a loop 0..65536, batch size 8



8 sockets

Distribute iterations at start of loop down to per-core counters

24..512    512..1024

16 cores per socket

Aggregate requests upwards within a core

0..8    8..16    16..24

8 h/w contexts per core

Per-core lock

# Approach, consider a loop 0..65536, batch size 8



**8 sockets**

Distribute iterations at start of loop down to per-core counters

24..512    512..1024

**16 cores per socket**

Aggregate requests upwards within a core

0..8    8..16    16..24

**8 h/w contexts per core**

Per-core lock

# Hierarchical distribution with request combining

- Combining implemented over flags in a single line in the shared L1 D$
- On TSO: no memory fences
- Synchronization remains core-local if work is evenly distributed
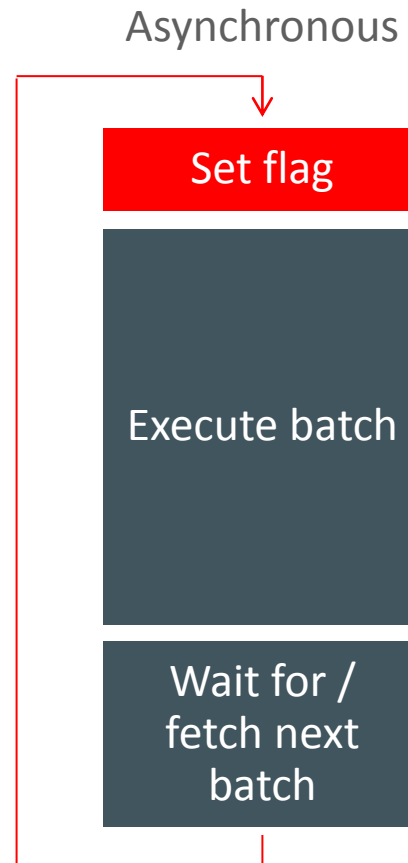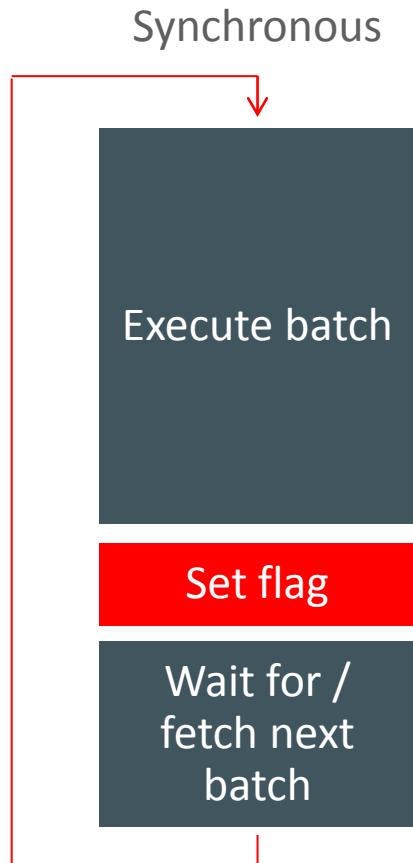- Threads waiting for combining can use mwait

# Techniques

**1** ▶ Request combining

**2** ▶ Asynchronous work requests
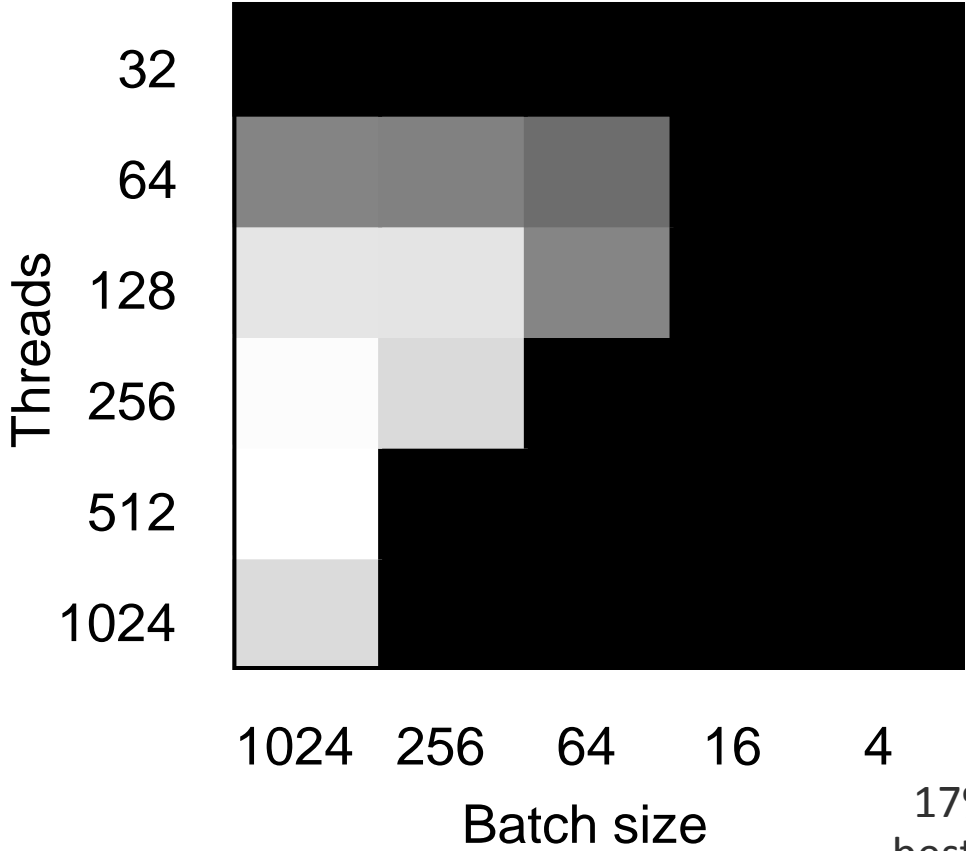
# Asynchronous combining of requests

Synchronous

# Asynchronous combining of requests



Synchronous

Execute batch

Set flag

Wait for / fetch next batch

Asynchronous
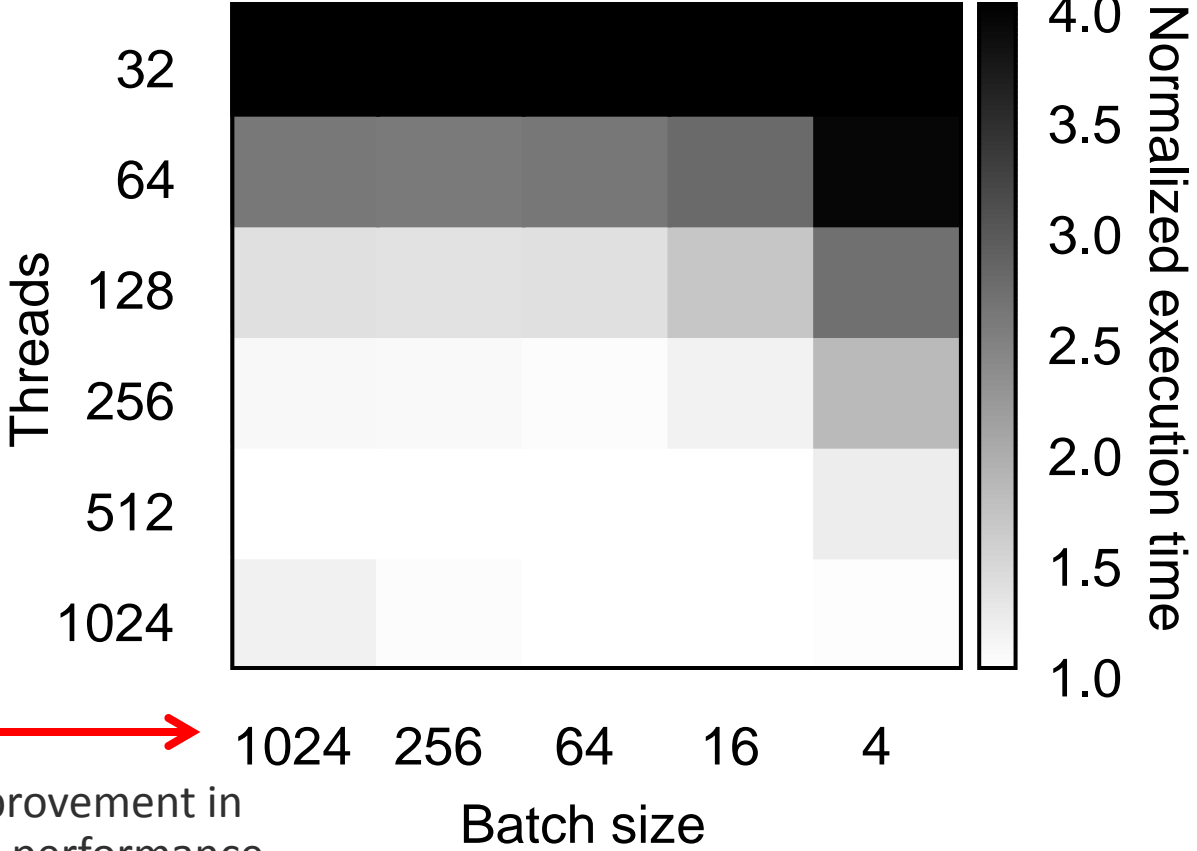
Set flag

Execute batch

Wait for / fetch next batch

Intuition: the time taken to execute the current batch provides an opportunity for other cores to service our request without us needing to wait, and the number of requests batched together will be larger.

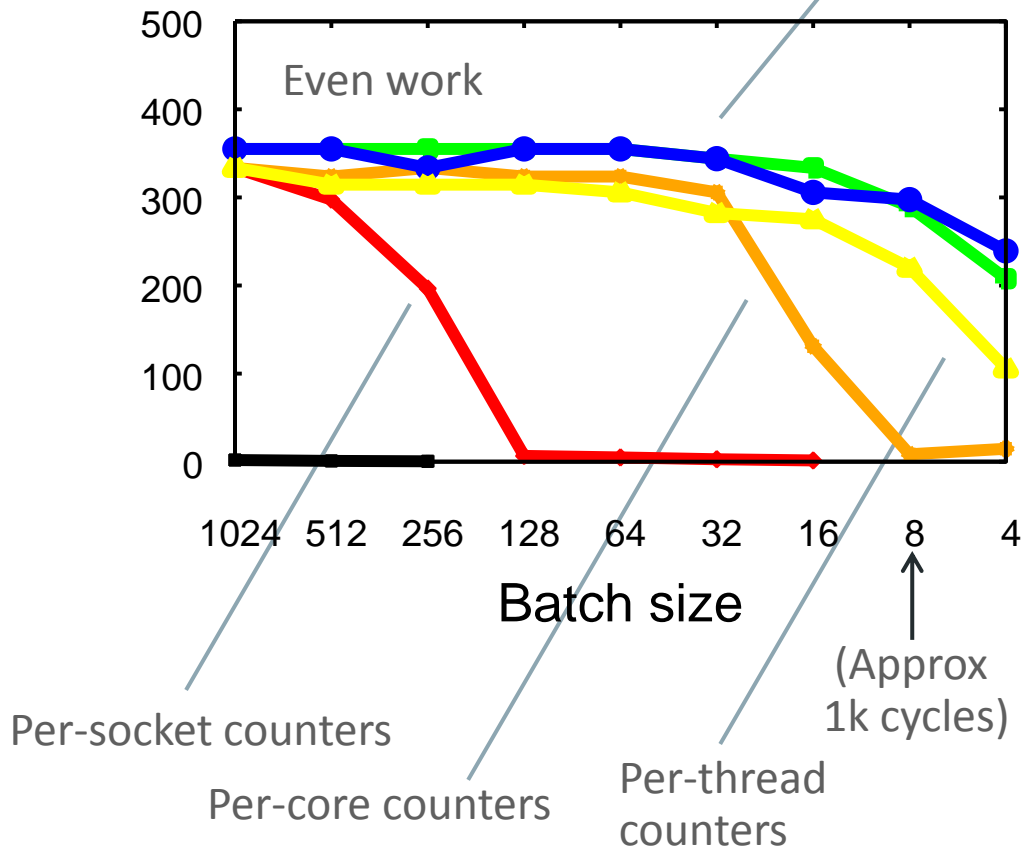# PageRank – SNAP LiveJournal (4.8M vertices, 69M edges)

OpenMP

Callisto-RTS



17% improvement in best-case performance
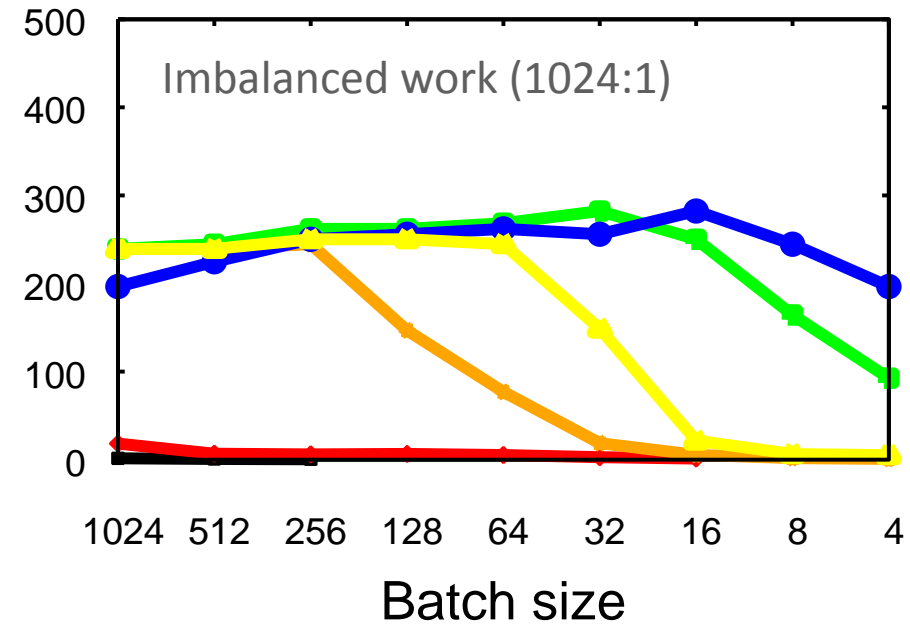
# Microbenchmark results

**SPARC T5-8, 1024 threads**
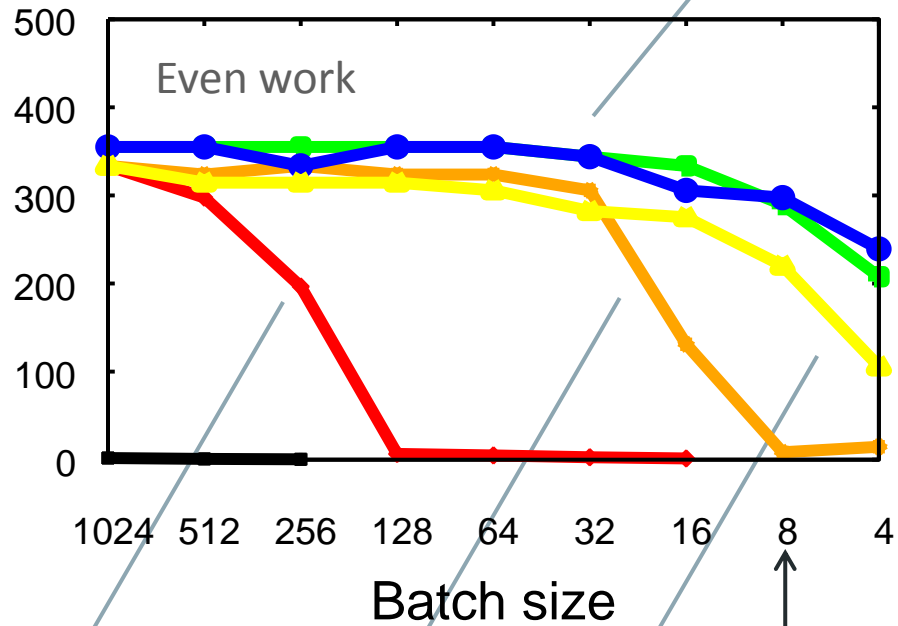
Per-core + asynchronous combining (blue)
Per-core + synchronous combining (green)



Even work

500
400
300
200
100
0

1024 512 256 128 64 32 16 8 4

Batch size

(Approx
1k cycles)

Per-socket counters

Per-core counters

Per-thread
counters

# Microbenchmark results

**SPARC T5-8, 1024 threads**

Per-core + asynchronous combining (blue)
Per-core + synchronous combining (green)



Even work

Imbalanced work (1024:1)

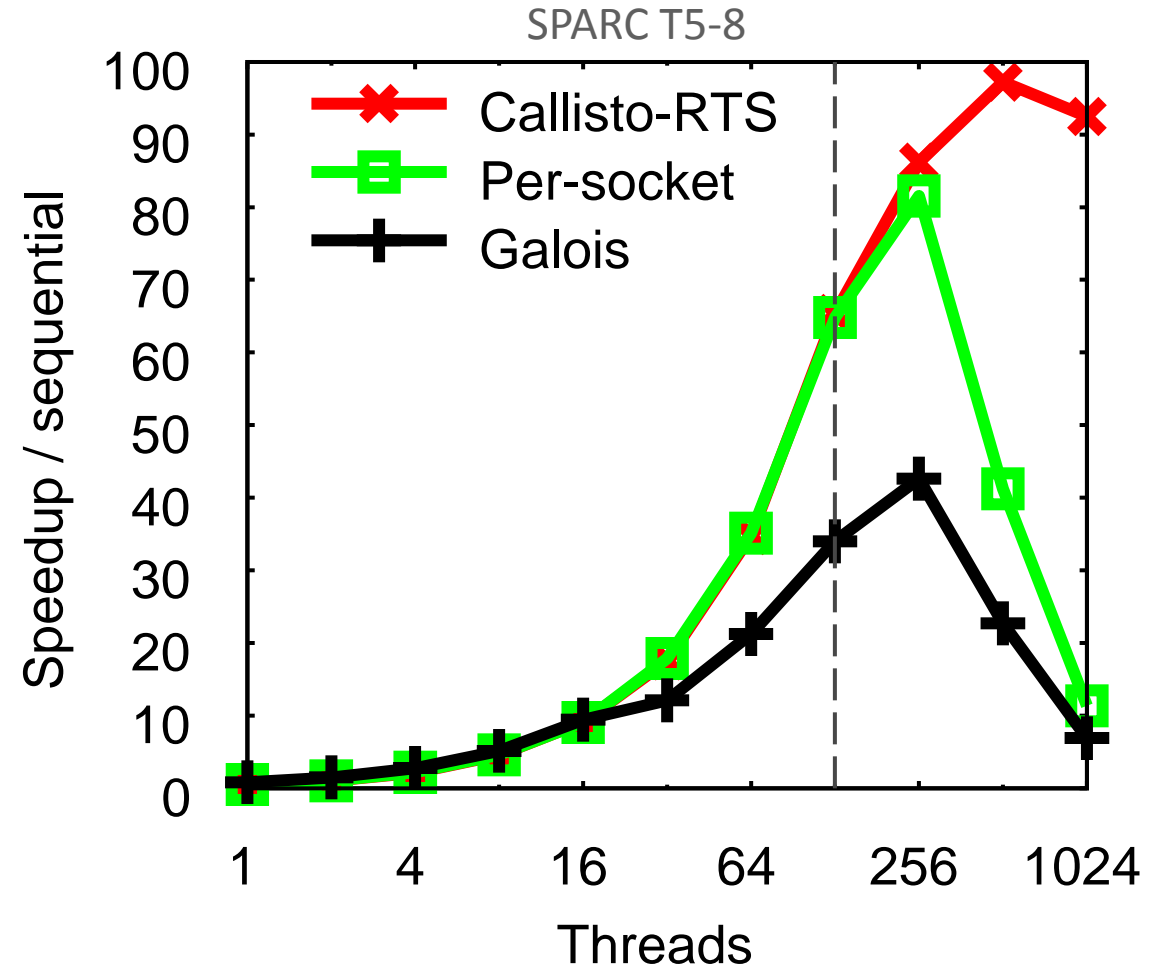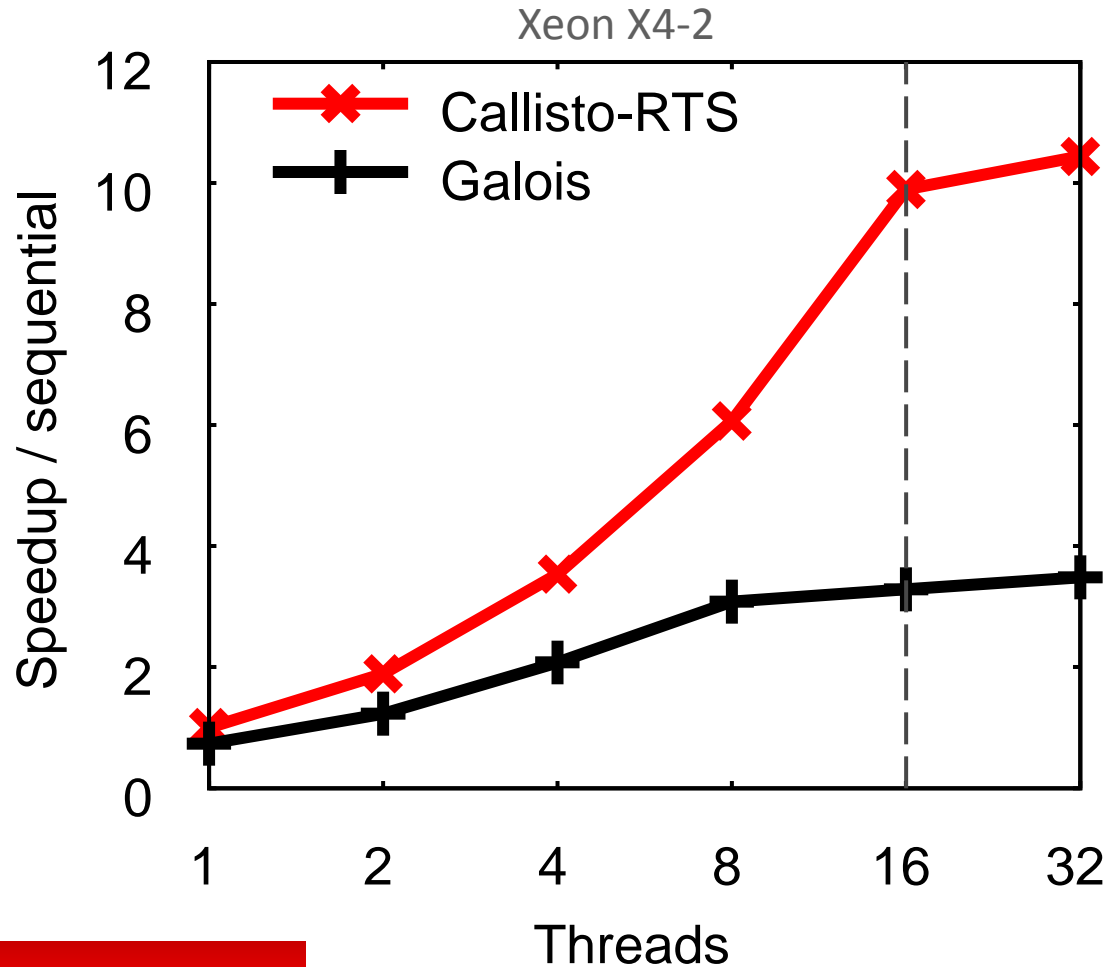Batch size

(Approx 1k cycles)

Per-socket counters

Per-core counters

Per-thread counters

# Comparison with Galois

## SNAP LiveJournal data set

# Nested loops

# Nested loops

- Abundant parallelism, why use nesting?

# Nested loops

- Abundant parallelism, why use nesting?

- Contention between iterations of an outer loop

- E.g., betweenness-centrality:
  - Iterate over vertices
  - BFS traversal from each vertex (plus additional work)

# Nested loops

- Abundant parallelism, why use nesting?

- Contention between iterations of an outer loop

- E.g., betweenness-centrality:
  – Iterate over vertices
  – BFS traversal from each vertex (plus additional work)

Better cache locality within each traversal
than between (unrelated) traversals
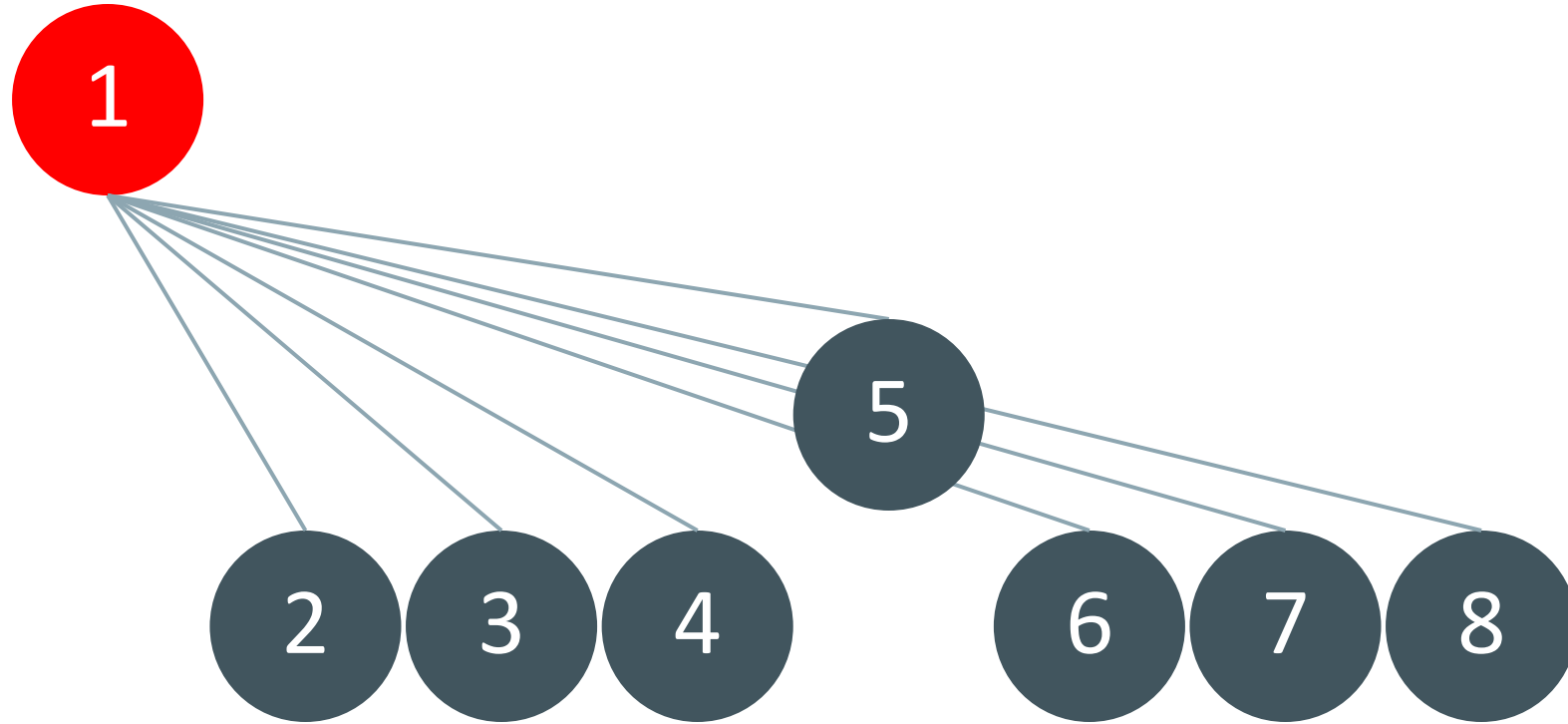
**ORACLE®**

# Nested loops

- Abundant parallelism, why use nesting?

- Contention between iterations of an outer loop

- E.g., betweenness-centrality:
  - Iterate over vertices
  - BFS traversal from each vertex (plus additional work)

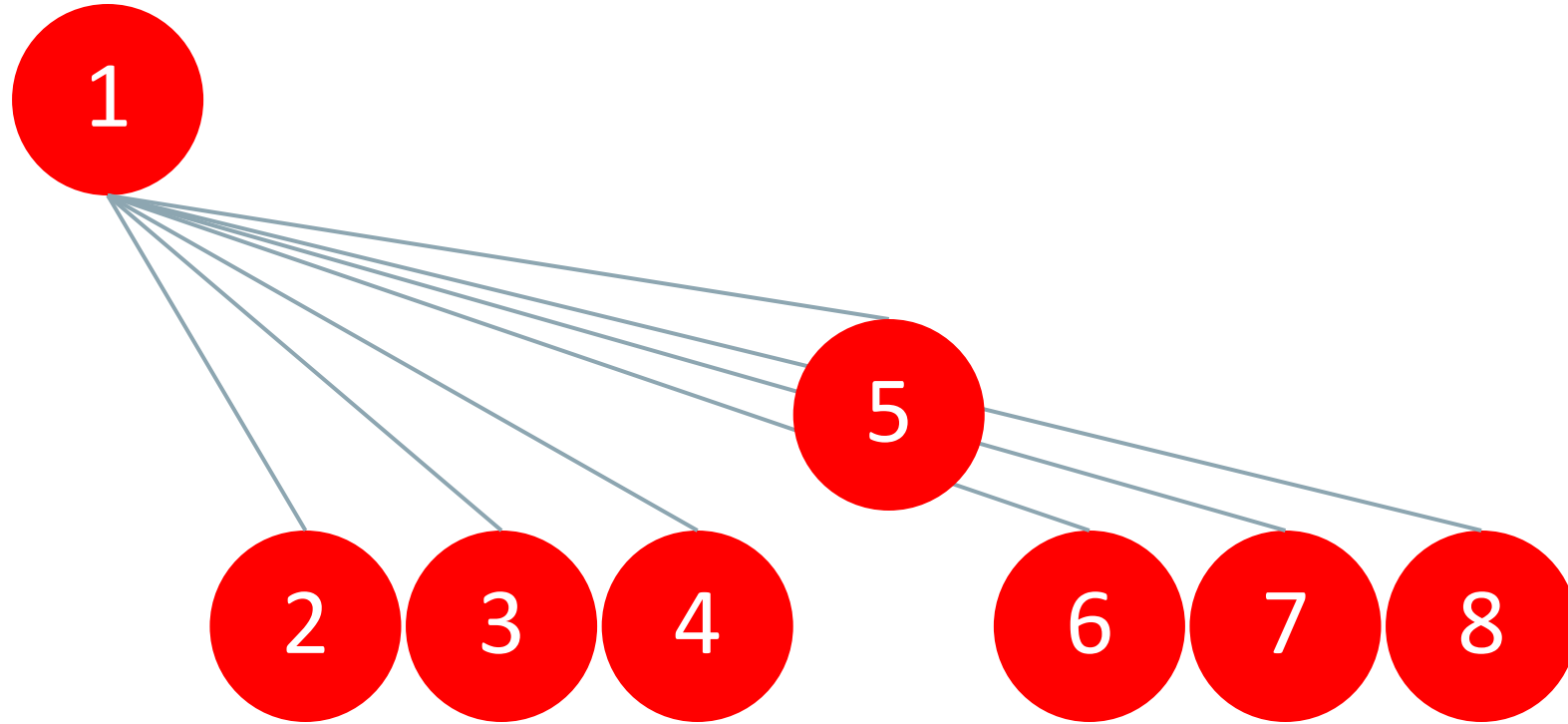Better cache locality within each traversal than between (unrelated) traversals → Run at most one of these per L2 D$
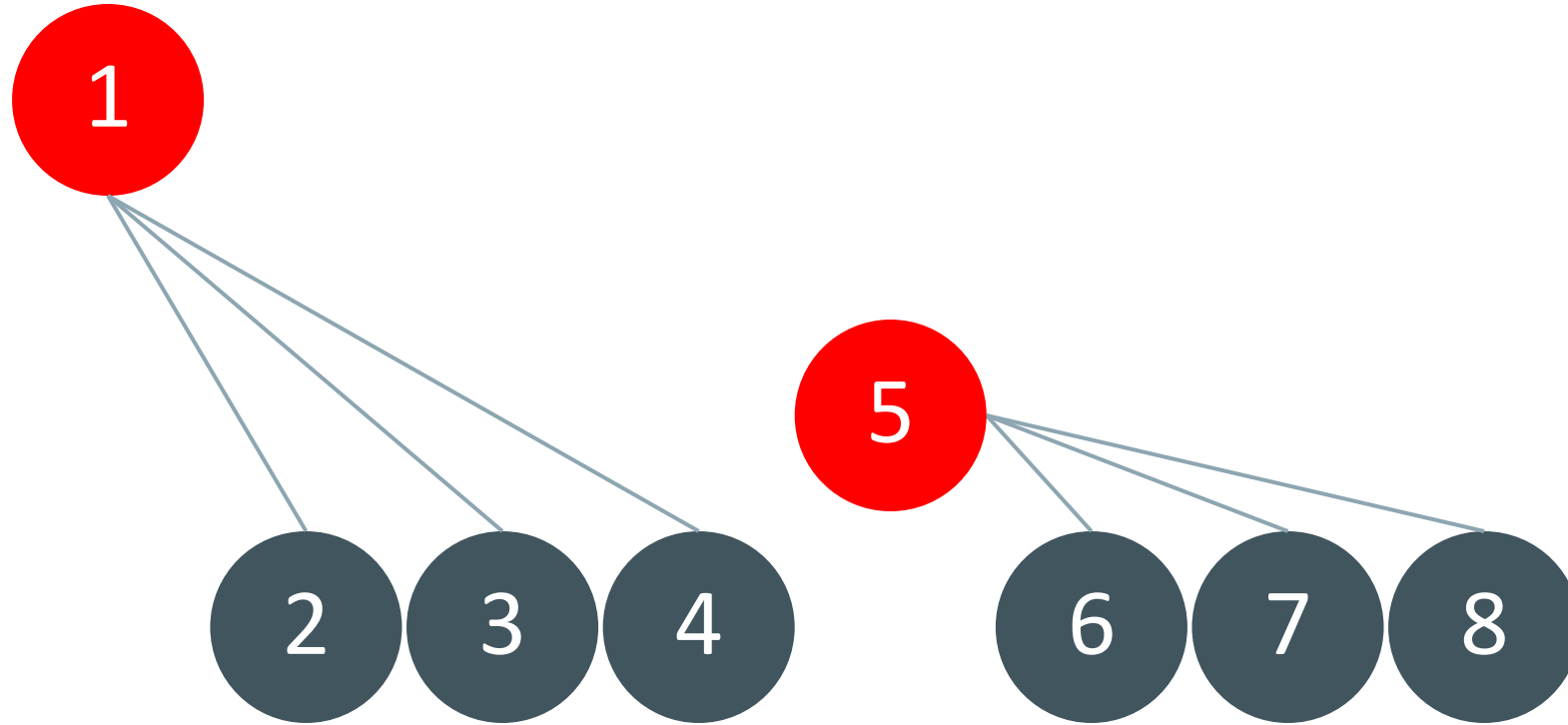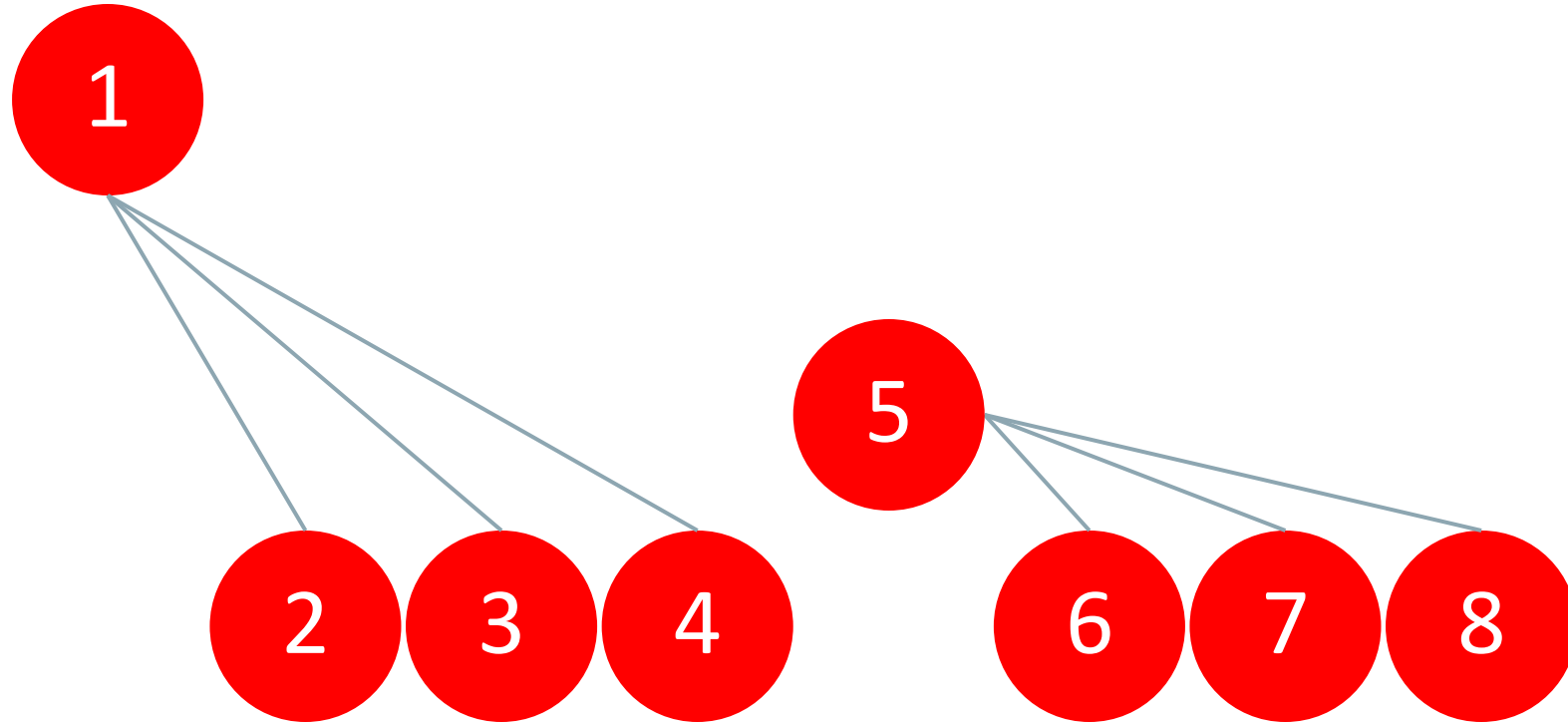
# Nested loops

# Nested loops: default, all threads participate

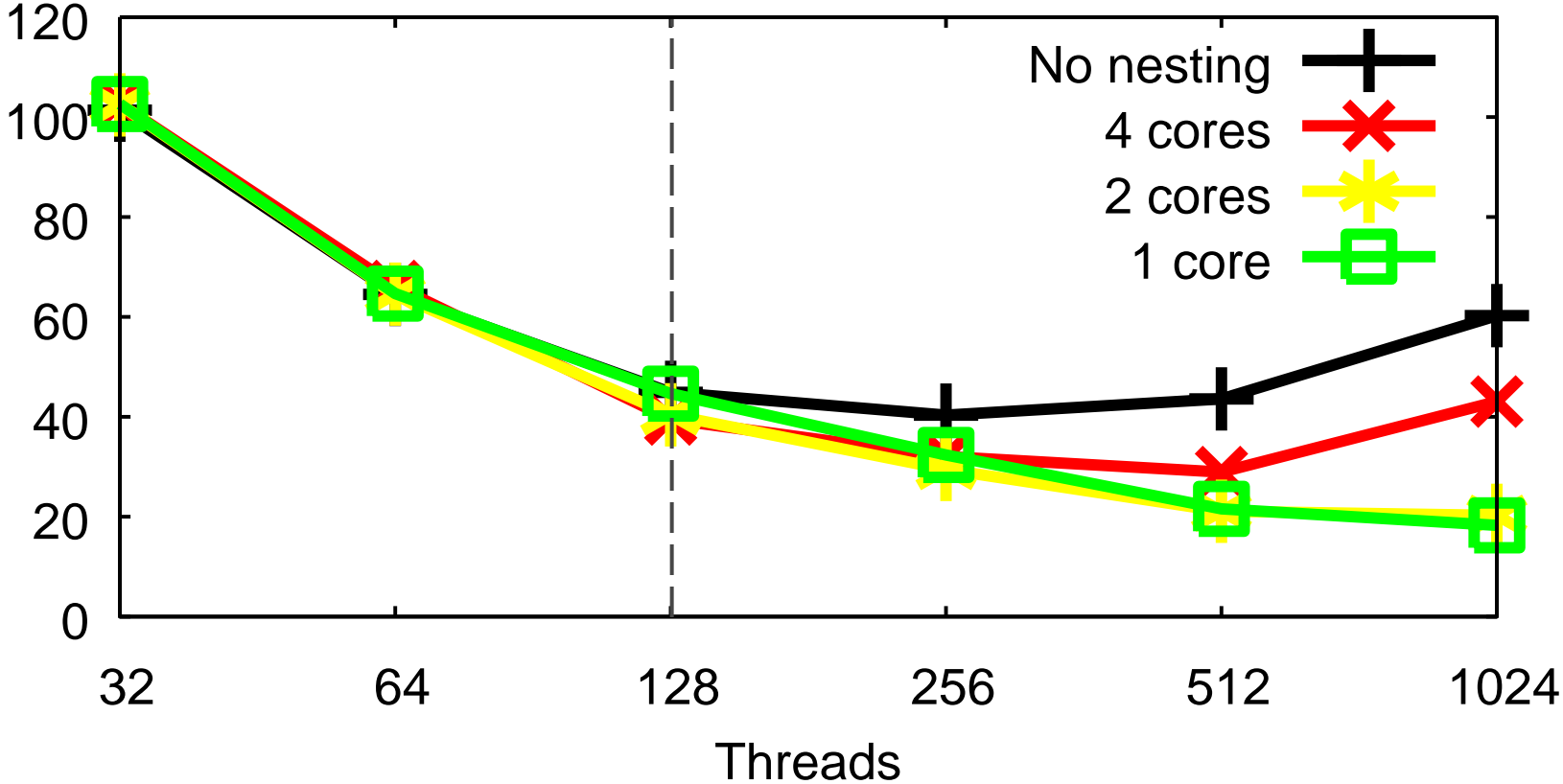# Nested loops: outer level – just 1+5 participate

# Nested loops: inner level –help respective leaders

# Betweenness-centrality

**SNAP Slashdot data set (82.1K nodes, 948K edges), T5-8**

# In-memory graph analytics
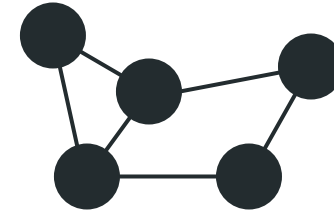
**Domain specific languages**
- Queries expressed in terms of graph concepts
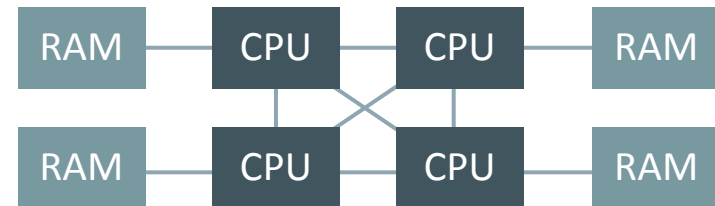- Tailor for different kinds of workload (e.g., sub-graph isomorphism)

**Generated code**
- Efficient in-memory data representations, e.g. compressed-sparse-row format
- Abundant parallelism

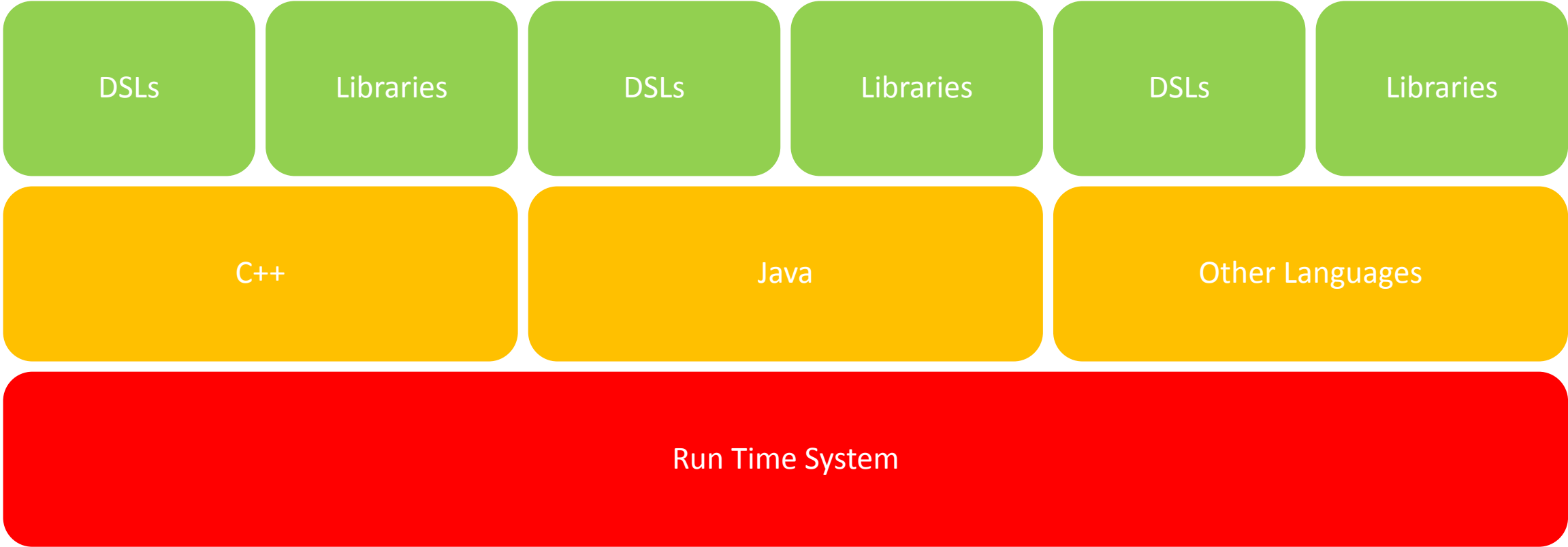**Runtime system**
- Allocation of resources to a query
- Distribution of work and data within a machine

```
parallel_for<node_t>([&](node_t n) {
  …
});
```

| RAM | CPU | CPU | RAM |
| --- | --- | --- | --- |
| RAM | CPU | CPU | RAM |

# RTS use cases

# Future work

- Continuing development of the programming model

- Control over data placement as well as threads
  - Initial examples from graph workloads generally have random accesses: spread data and threads widely in the machine
  - (See "Shoal", USENIX ATC 2015)

- Interactions between multiple parallel workloads
  - OS/runtime system interaction (ref our prior work at EuroSys 2014)
  - Placement in the machine
  - Control over degree of parallelism

# Integrated Cloud
## Applications & Platform Services