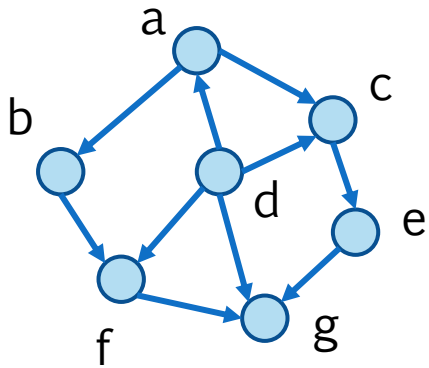# A Specialized B-tree for Concurrent Datalog Evaluation

**Herbert Jordan**[1], Pavle Subotić[3], David Zhao[2], and Bernhard Scholz[2]

PPoPP 2019, 16-20 February 2019, Washington, DC

universität innsbruck
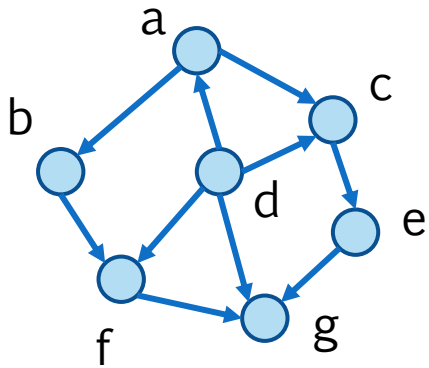
UCL PPLV

DPS
Distributed and Parallel Systems

1) University of Innsbruck
2) University of Sydney
3) University College London

# Datalog (by Example)



graph

| from | to |
| --- | --- |
| a | b |
| a | c |
| b | f |
| c | e |
| d | a |
| d | c |
| ... | ... |

edge relation

Which nodes
are connected?

# Datalog (by Example)



graph

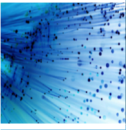| from | to |
|------|-----|
| a | b |
| a | c |
| b | f |
| c | e |
| d | a |
| d | c |
| ... | ... |

edge relation

```
path(X,Y) :- edge(X,Y).

path(X,Z) :- path(X,Y),
             edge(Y,Z).
```

Datalog
query

# Datalog

› Benefits:
  – a concise formalism for powerful data analysis
  – lately major performance improvements and tool support

› Applications:
  – data base queries
  – program analysis
  – security vulnerability analysis
  – network analysis

100s of relations and rules, billions of tuples, all in-memory

# Query Processing

relations ➡ set of integer tuples

rules ➡ sequence of relational algebra operations on sets

# Example

```
path(X,Z) :- path(X,Y), edge(Y,Z).
```

$$delta \leftarrow path$$
$$\text{while } ( \ delta \neq \emptyset \ ) \{$$
$$\quad new \leftarrow \pi(delta \bowtie edge) \setminus path$$
$$\quad path \leftarrow path \cup new$$
$$\quad delta \leftarrow new$$
$$\}$$

computational expensive and dominating part

# Example

$$new \leftarrow \pi(delta \bowtie edge) \setminus path$$

```
Relation new;
for t1 ∈ delta {
    auto l = edge.lower_bound( { t1[1], 0 } );
    auto u = edge.upper_bound( { t1[1]+1, 0 } );
    for t2 ∈ [l,u] {
        Tuple t3 = { t1[0], t2[1] };
         if ( t3 ∉ path ) {
             new.insert(t3);
         }
     }
}
```
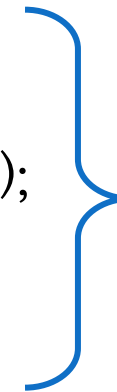
# Example

$$new \leftarrow \pi(delta \bowtie edge) \setminus path$$

```
Relation new;
#pragma omp parallel for
for t1 ∈ delta {
    auto l = edge.lower_bound( { t1[1], 0 } );
    auto u = edge.upper_bound( { t1[1]+1, 0 } );
    for t2 ∈ [l,u] {
        Tuple t3 = { t1[0], t2[1] };
        if ( t3 ∉ path ) {
            new.insert(t3);
        }
    }
}
```

all read accesses
(right hand side)

one write access
(assignment)

But: write target is
never read on
right hand side!

# Needed

› efficient   data structure for relations
- – maintain set of n-dimensional tuples

- – efficient support for
  - › insertion,
  - › scans,
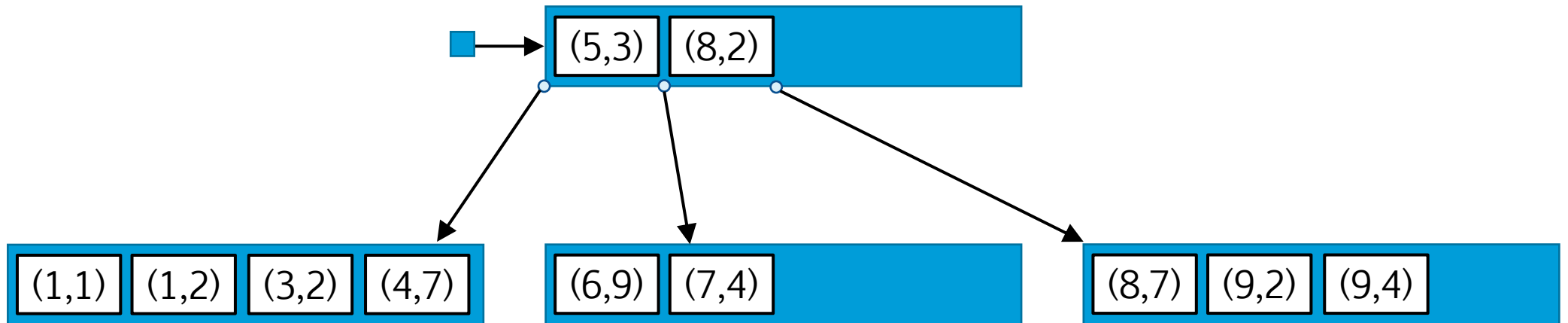  - › range queries,
  - › membership tests,
  - › emptiness checks

  well supported
  by B-trees

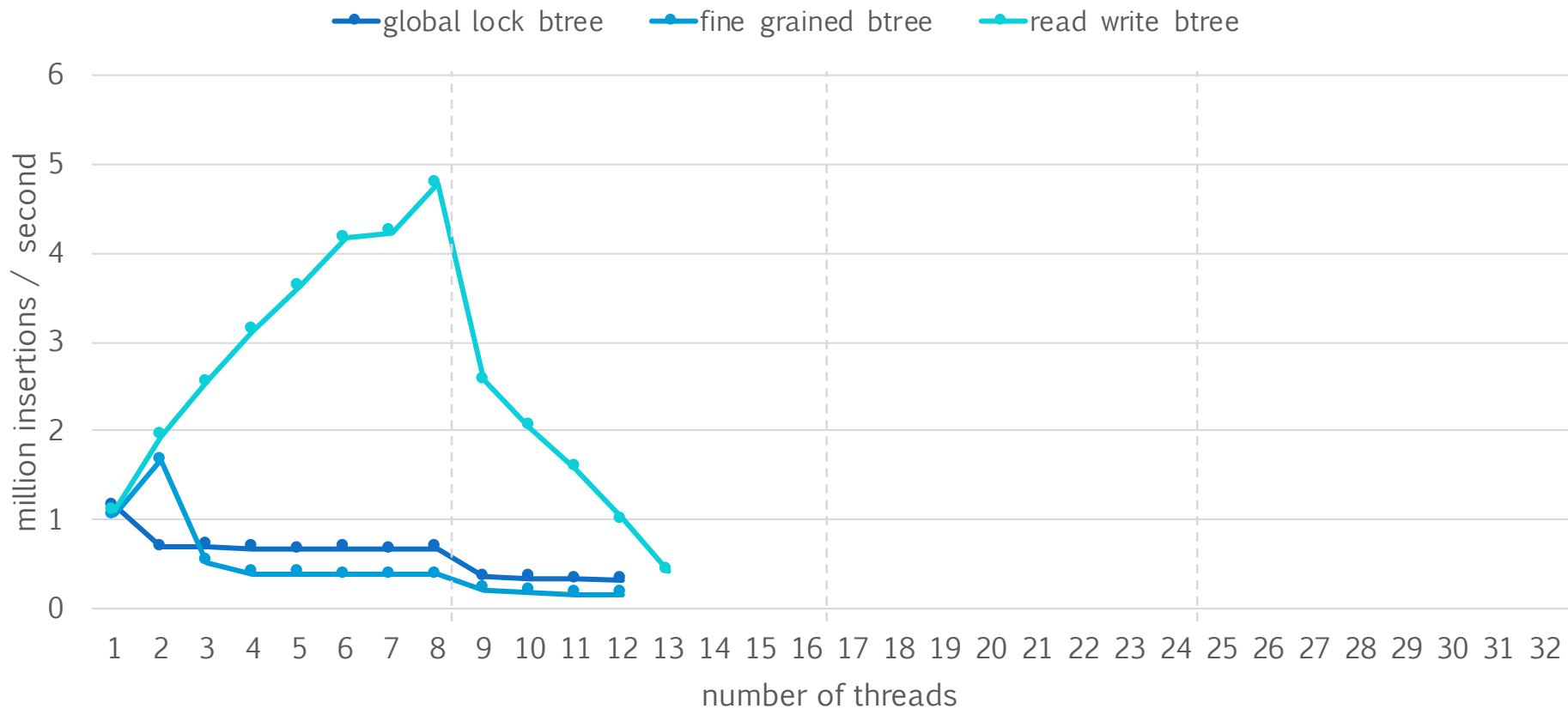- – efficient synchronization of concurrent inserts

  not so much …

# B-tree



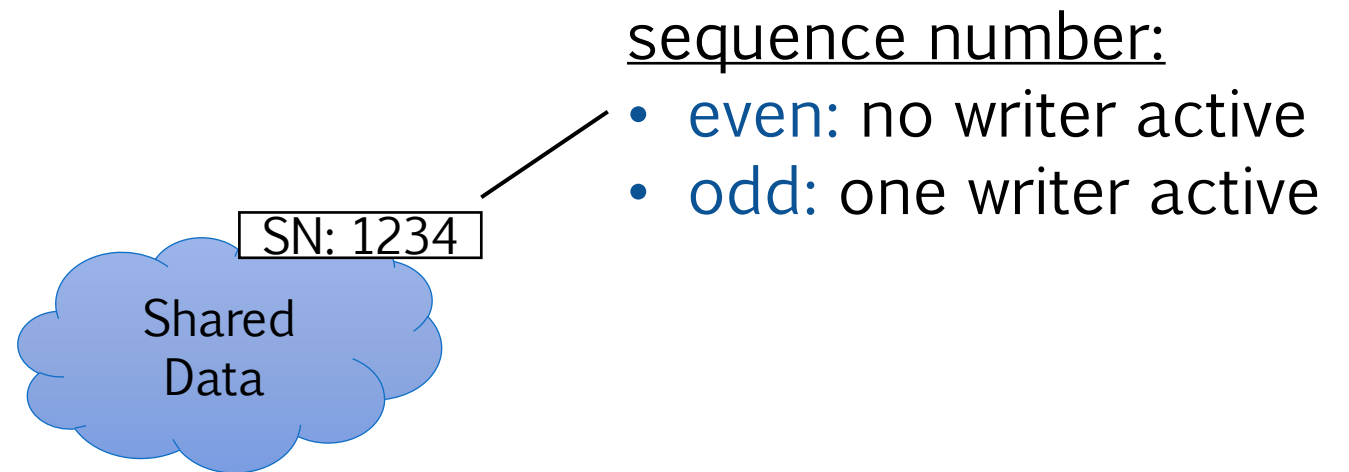The diagram shows a B-tree. The root node contains keys (5,3) and (8,2), with three child pointers. The left leaf contains (1,1), (1,2), (3,2), (4,7). The middle leaf contains (6,9), (7,4). The right leaf contains (8,7), (9,2), (9,4).

› Insertion:
 – locate target leaf node
 – split leaf node if necessary, may propagate up
 – insert element in sorted leaf-node element array

# B-tree Locking Strategies
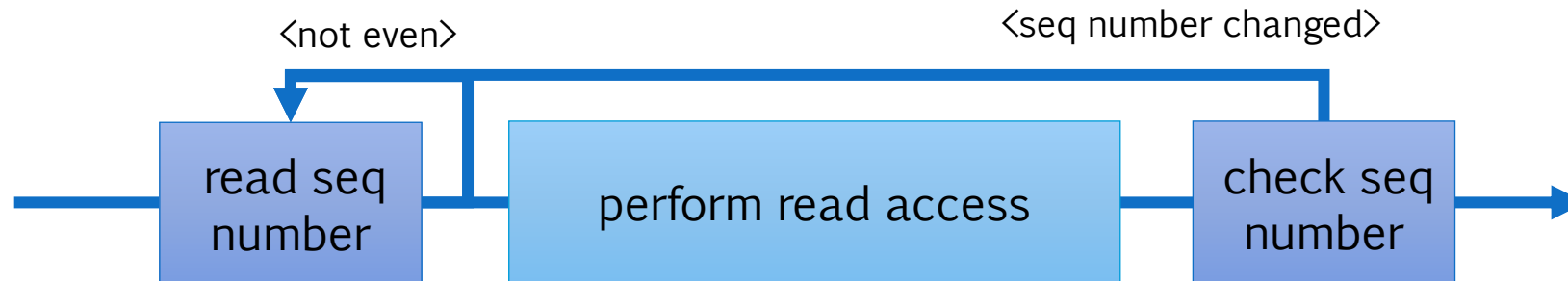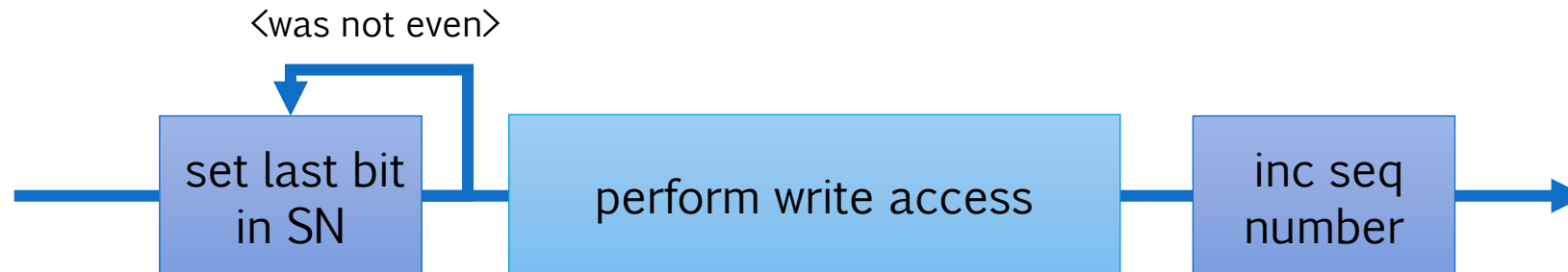


random order, on 4x8 core Intel Xeon E5-4650

# Seqlocks / Optimistic Locking

sequence number:
- even: no writer active
- odd: one writer active

SN: 1234

Shared
Data

# Seqlocks

SN: 1234

Shared Data

› Reader:

<not even>   <seq number changed>

| read seq number | perform read access | check seq number |

› Writer:

<was not even>

| set last bit in SN | perform write access | inc seq number |

13

# Optimistic Read/Write Lock

› Reader

Shared Data



14

# Optimistic B-tree

› Protect nodes and root pointer with optimistic R/W lock

› Synchronize insert operation
– read access on inner nodes, update to write when necessary

› Key challenge:
– pointer indirection
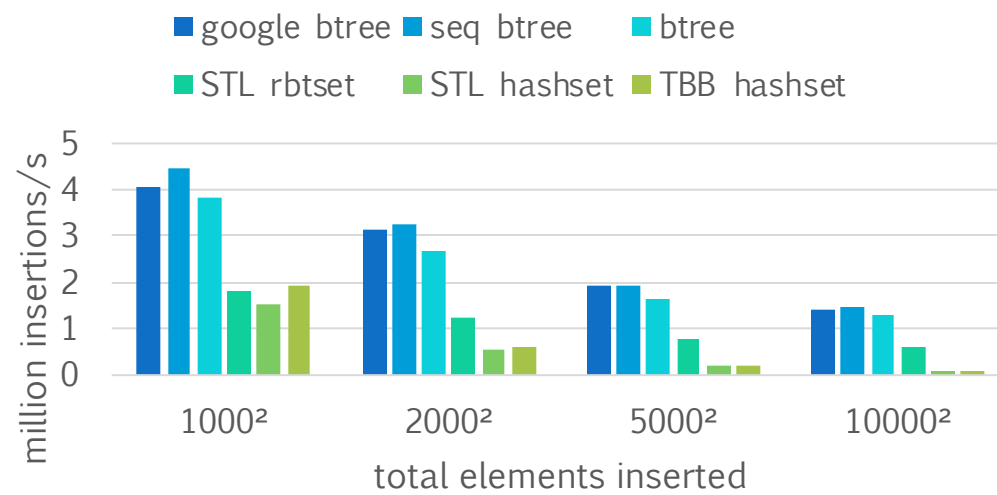– concurrency memory model

# B-tree Locking Strategies (cont)



random order, on 4x8 core Intel Xeon E5-4650
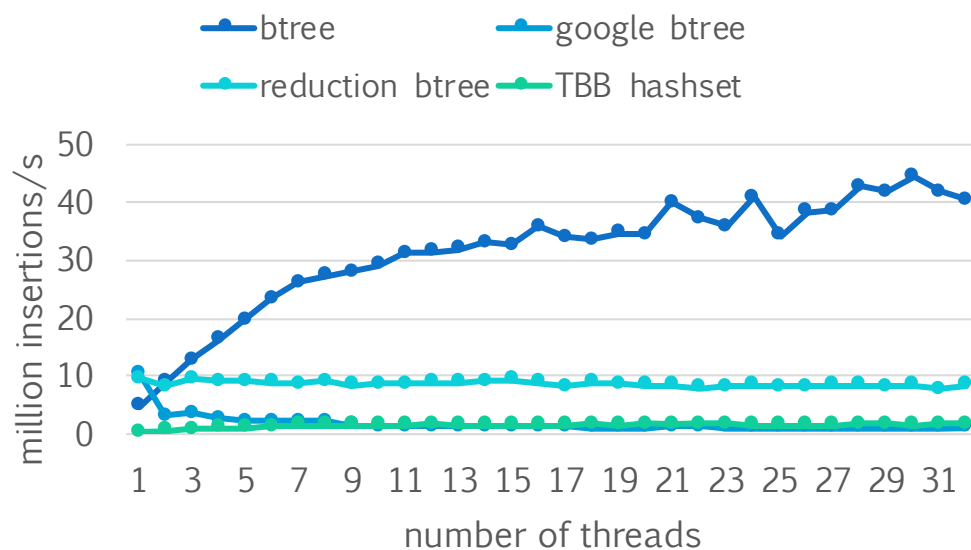
# Sequential Performance



ordered insertion

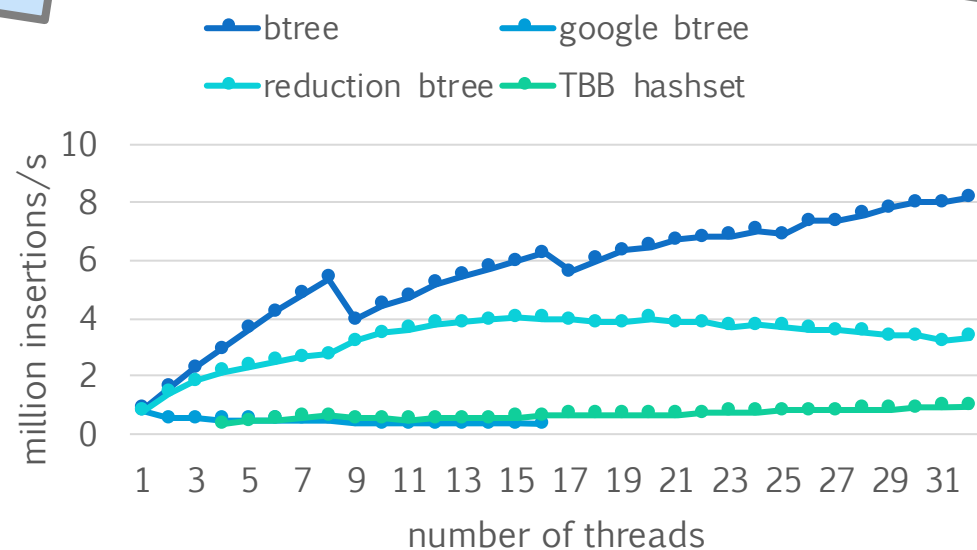random order insertion

(additional data structures covered in paper)

# Parallel Performance



ordered insertion

random order insertion
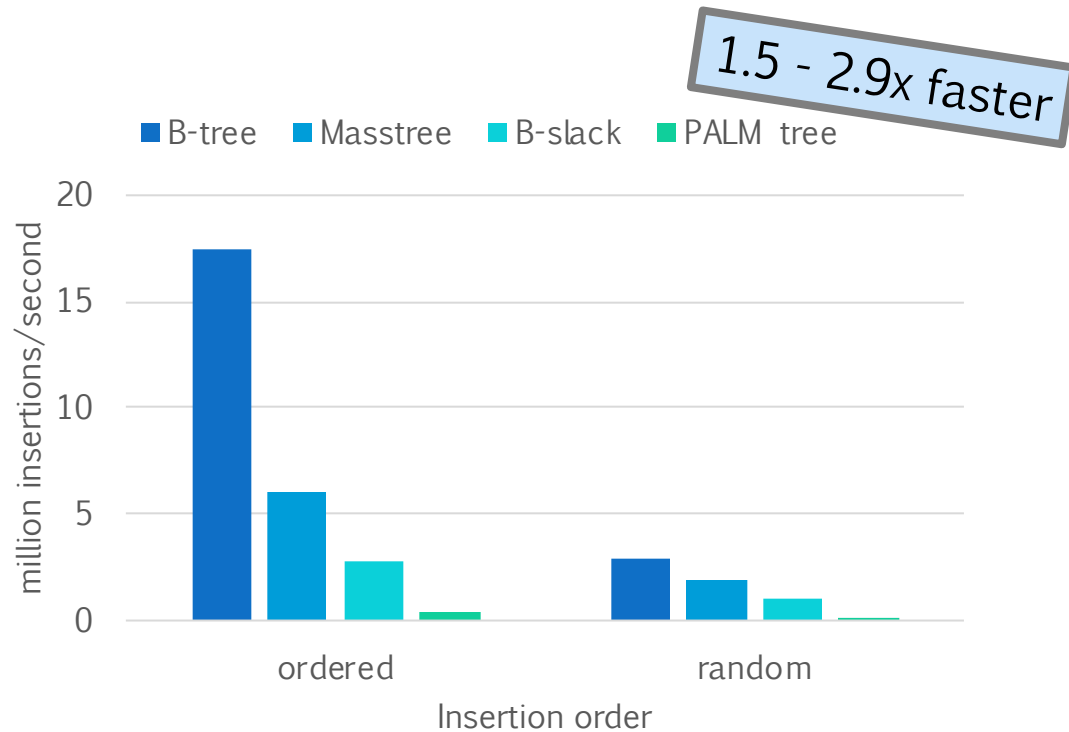
4x8 core Intel Xeon E5-4650

# Other Concurrent Tree Data Structures

1.5 - 2.9x faster

1.5 - 2.6x faster

■ B-tree  ■ Masstree  ■ B-slack  ■ PALM tree

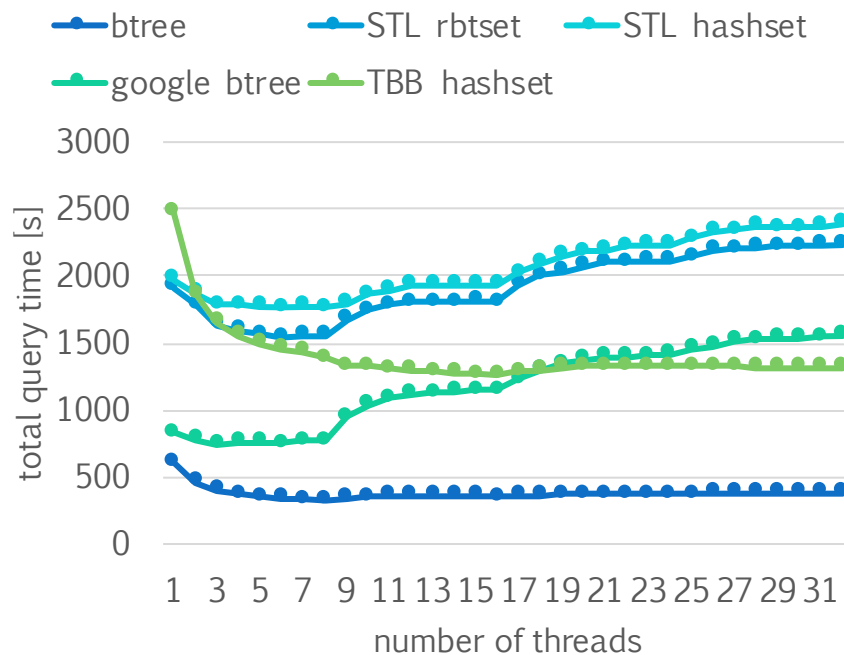■ B-tree  ■ Masstree  ■ B-slack  ■ PALM tree

single-threaded

8 threads
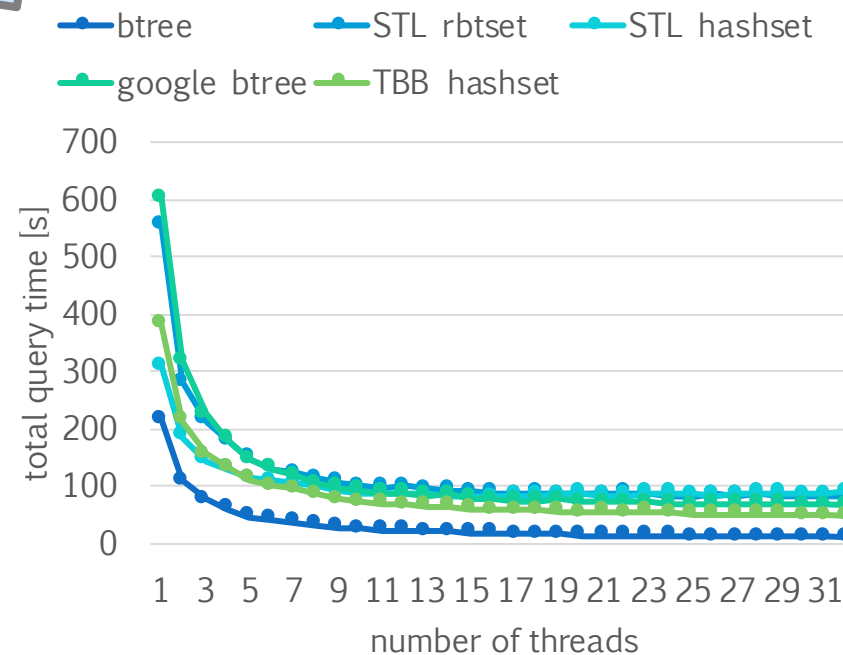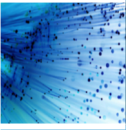
# Datalog Query Processing

1.4 - 6.3x faster

1.4 - 7.7x faster



context sensitive
var-points-to analysis

security vulnerability
analysis

# Conclusion

› Developed concurrent set for Datalog relations:

  – B-tree foundation
    › good sequential performance, cache friendly

  – Fine-grained synchronization
    › based on customized seqlock variant

› Results:
  – up to 59x faster than state-of-the-art hash based sets
  – up to 2.9x faster than state-of-the-art tree based sets
  – up to 7.7x faster for real-world query processing

› Future work:
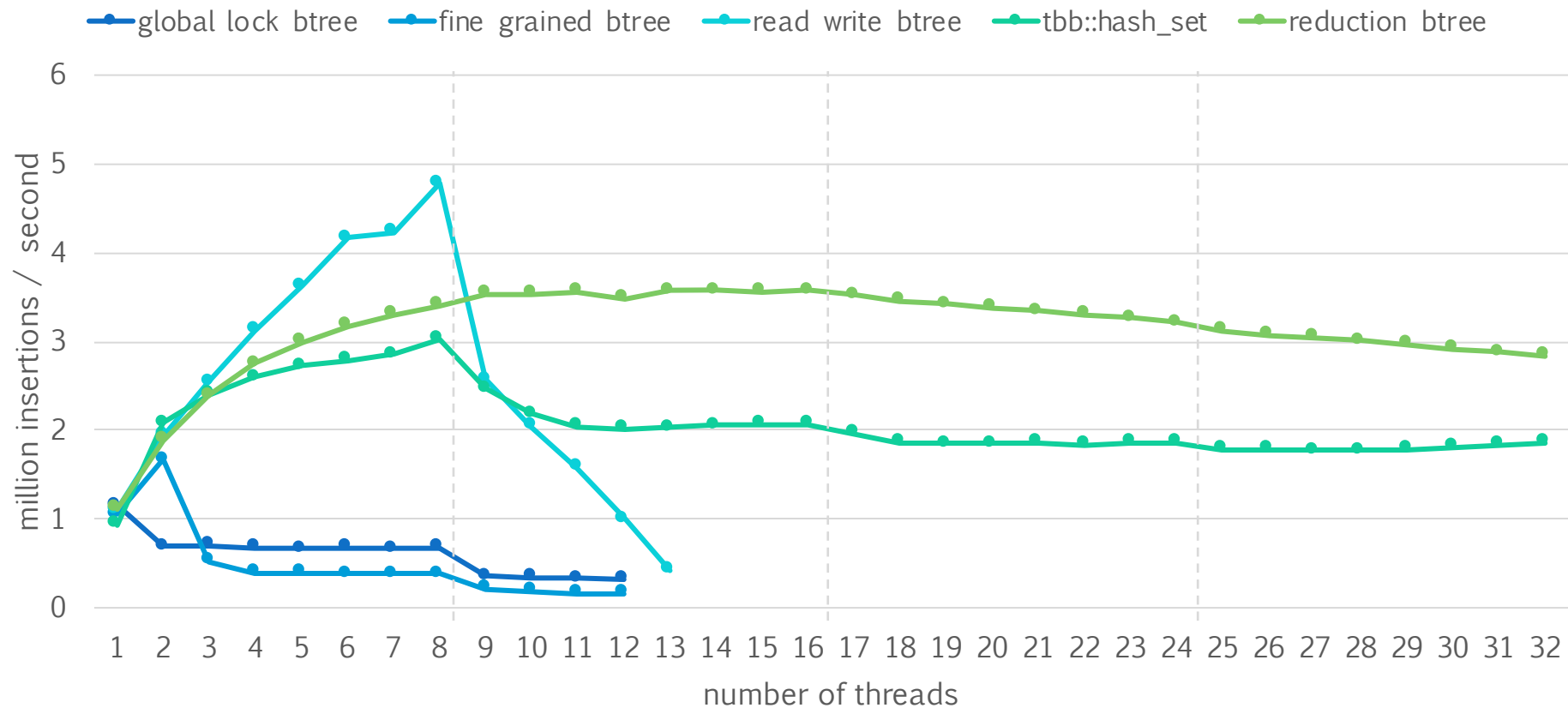  – investigate other data structures for specialized use cases

# Thank you!

visit us on https://souffle-lang.github.io
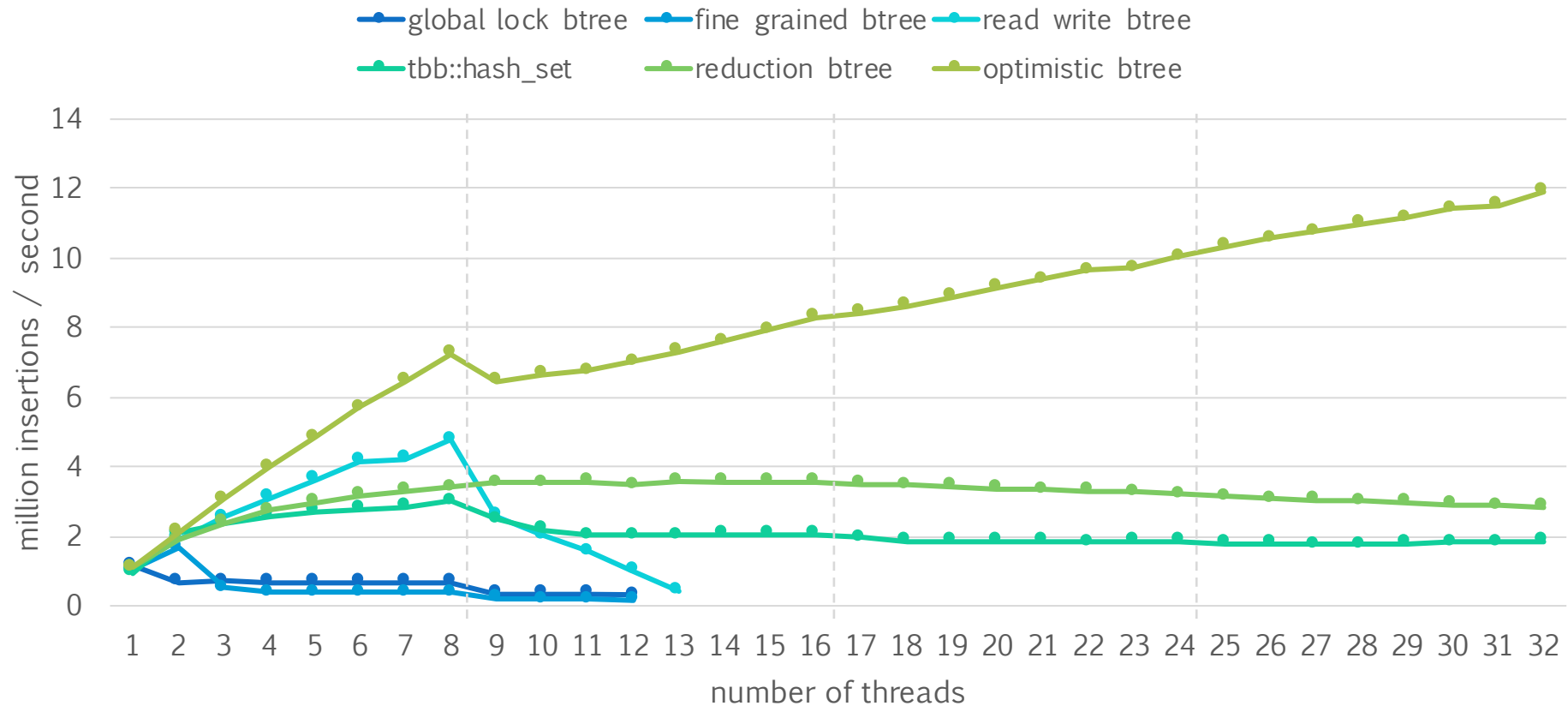
sources: https://github.com/souffle-lang/souffle

# B-tree Locking Strategies



random order, on 4x8 core Intel Xeon E5-4650

# B-tree Locking Strategies (cont)



random order, on 4x8 core Intel Xeon E5-4650