

Inspector

Data Provenance using Intel Processor Trace (PT)

Jörg Thalheim

Pramod Bhatotia & Christof Fetzer

Technical University of Dresden

Data Provenance

...records transformations made on data to explain how the computation was performed

Motivation: Use-case examples

Dependability:

Debugging
programs

Security:

Dynamic Information
Flow Tracking (DIFT)

Efficiency:

Memory management
for NUMA

Research gap

- Currently limited either to sequential programs

For parallel programs:

- Require manual annotations w/ new type systems
- Restrictive programming model & synchronization primitives

Design goals

- **Transparency**
 - Unmodified multithreaded programs
- **Generality**
 - Shared-memory model w/ POSIX sync. primitives
- **Efficiency**
 - Low overheads using a parallel provenance algorithm

Inspector: Easy to use!

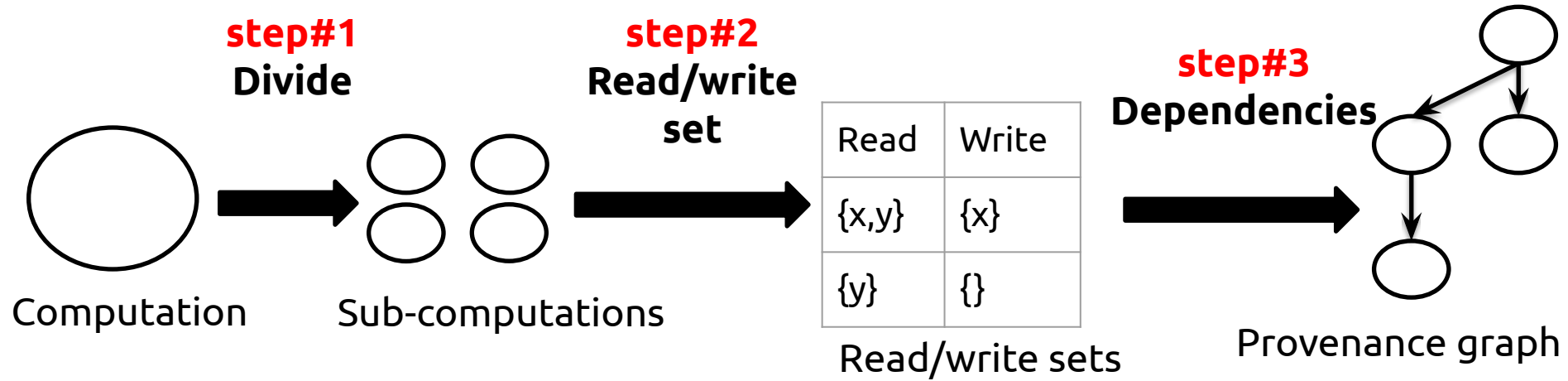
```
$ inspector -- ./<program> <arguments>
```

1. Preloads the Inspector library
LD_PRELOAD=libinspector.so
2. Executes “existing binaries” w/o re-compilation
3. Writes the provenance log to ./perf.data

Agenda

- ✓ Motivation
- ☐ Design
- ☐ Implementation
- ☐ Evaluation

Behind the scenes



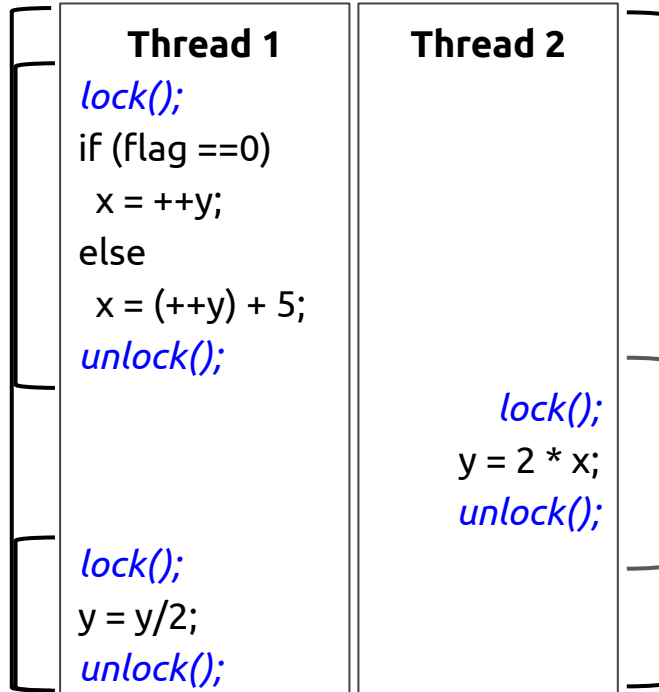
A simple example

Shared variables: x and y

Thread 1	Thread 2
<pre>lock(); if (flag == 0) x = ++y; else x = (++y) + 5; unlock();</pre>	<pre>lock(); y = 2 * x; unlock();</pre>
<pre>lock(); y = y/2; unlock();</pre>	

Step #1: Sub-computations

Shared variables: x and y



Approaches:

Coarsed Grained

- Whole Thread
- imprecise

Fine Grained

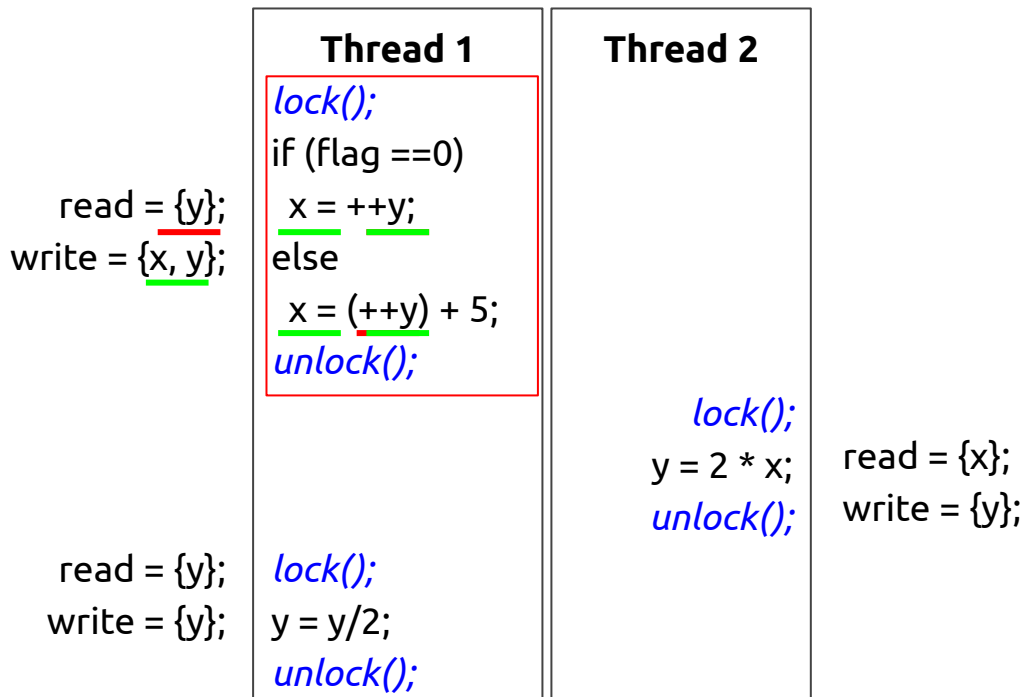
- Every Instruction
- expensive

Middle Ground

- Sub-Computations

Step #2: Read-write sets

Shared variables: x and y



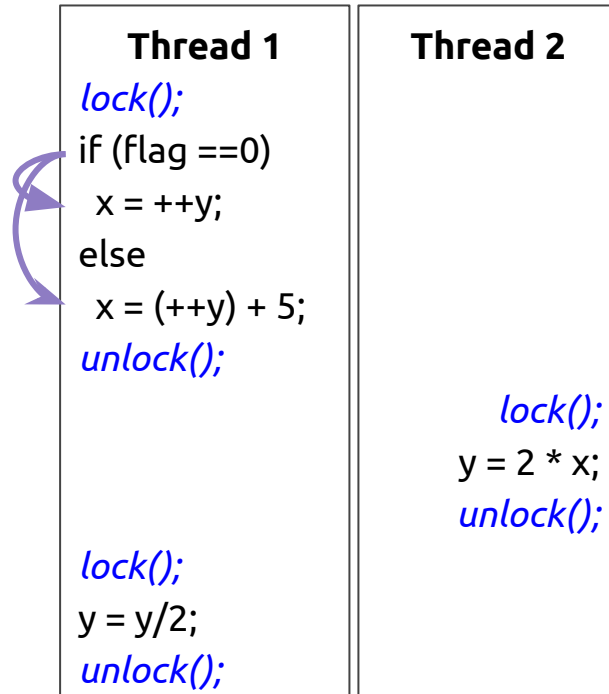
Step #3: Provenance graph

We record three dependencies:

- A. Control
- B. Schedule
- C. Data

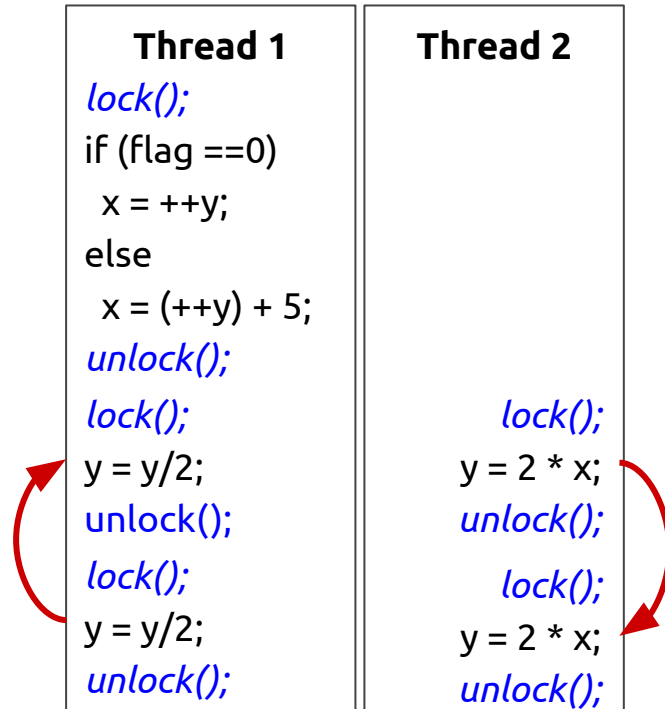
A: Control dependencies

Shared variables: x and y



B: Synchronization dependencies

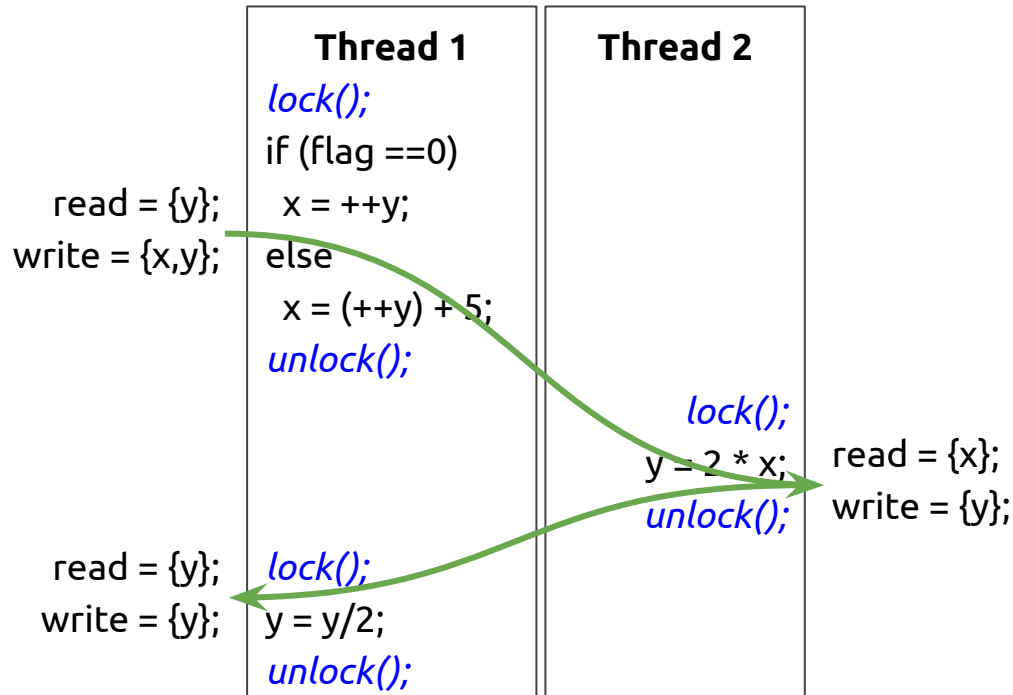
Shared variables: x and y



Change of Schedule

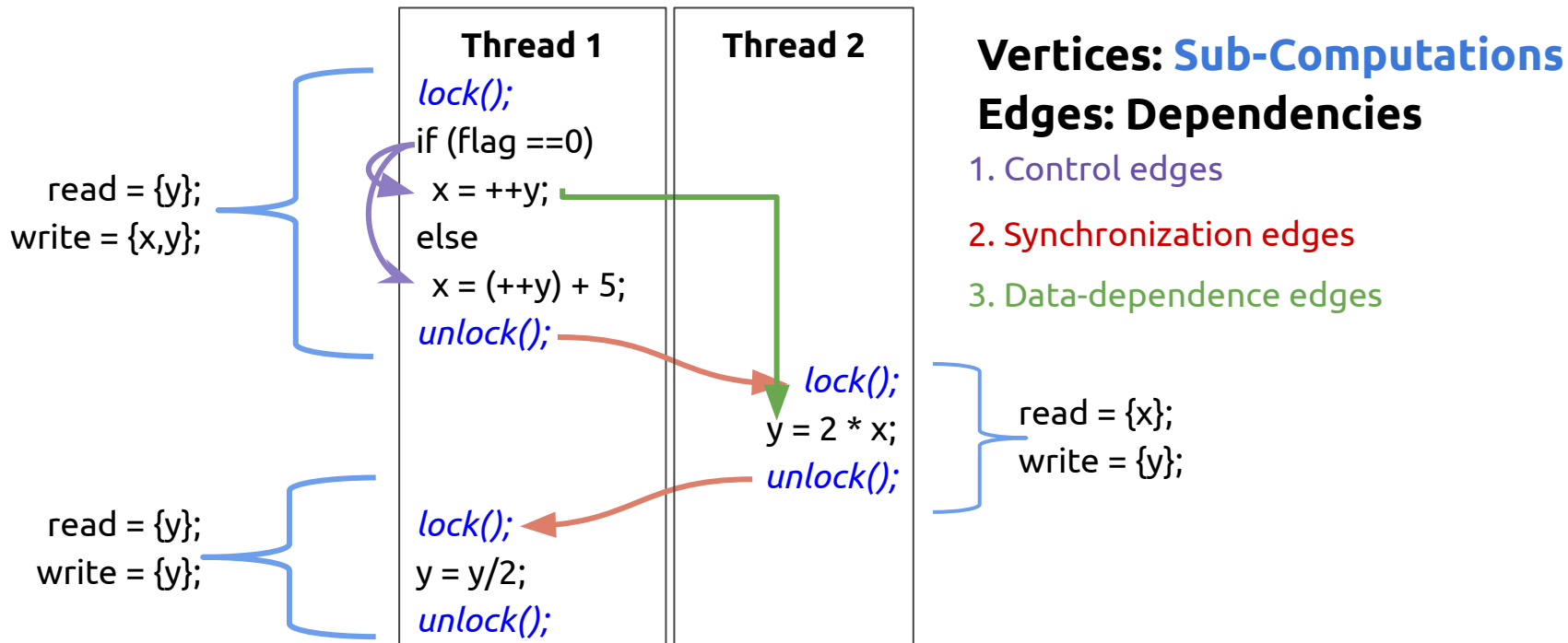
C: Data dependencies

Shared variables: x and y



Concurrent Provenance Graph (CPG)

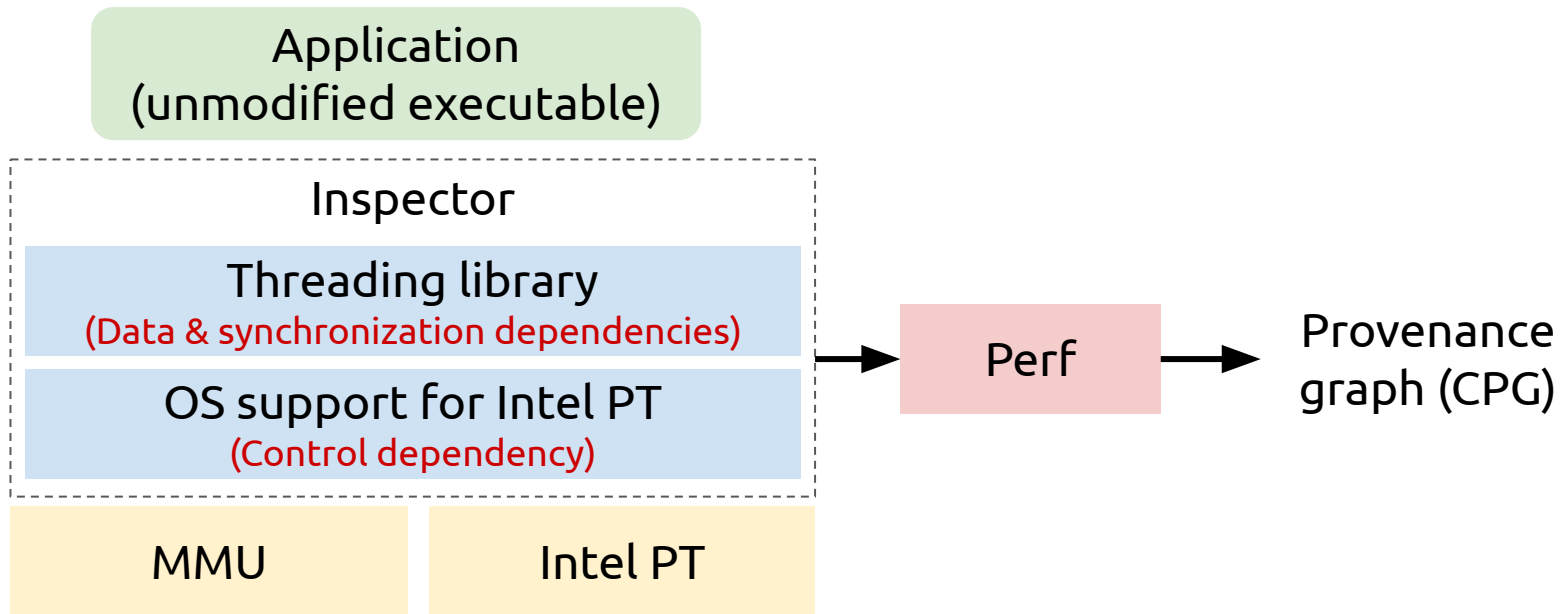
Shared variables: x and y



Agenda

- ✓ Motivation
- ✓ Design
- ☐ Implementation
- ☐ Evaluation

Inspector architecture



Agenda

- ✓ Motivation
- ✓ Design
- ✓ Implementation
- ☐ Evaluation

Evaluation

Questions:

1. Performance overheads
2. Sources for these overheads

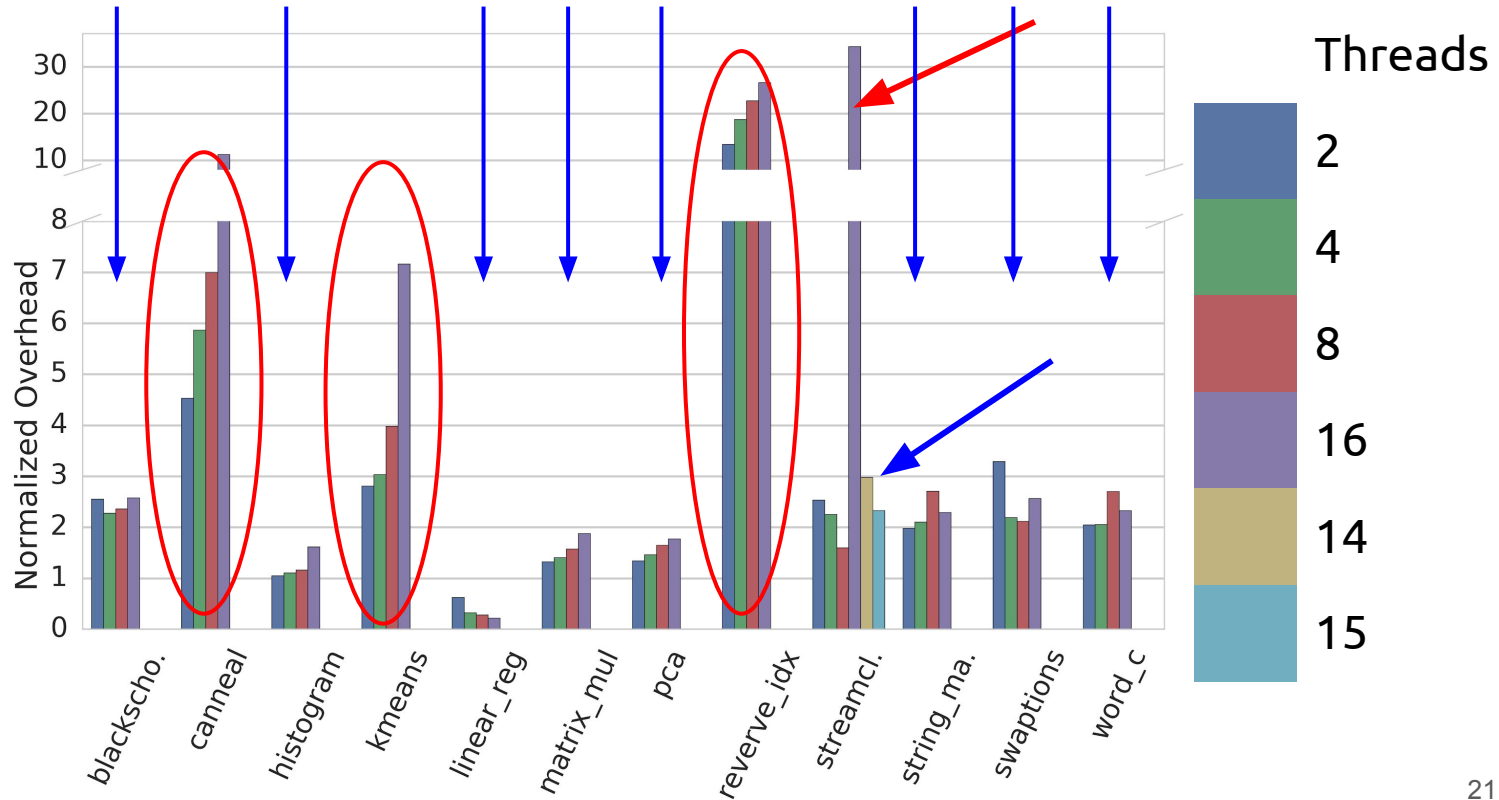


More results in the paper

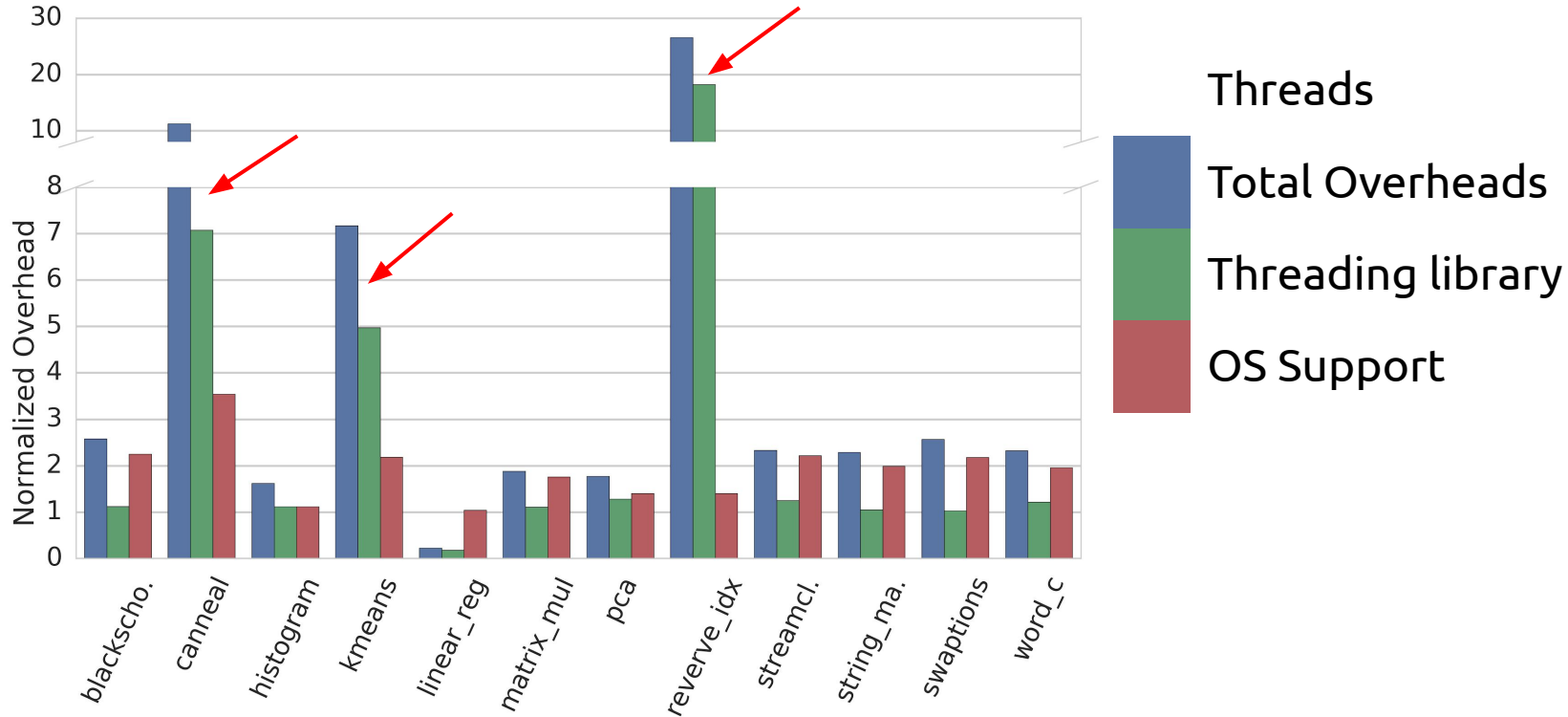
Experimental setup:

- Benchmarks: Phoenix 2.0 and PARSEC 3.0
- Platform: Intel Broadwell CPU with 8 cores (16 hyper-threads)

Q1: Performance overheads



Q2: Source of the overheads



Summary

Inspector: Data provenance using Intel Processor Trace (PT)

- **Transparent:** Targets unmodified multithreaded programs
- **General:** Supports the shared-memory model w/ POSIX sync primitives
- **Efficient:** Employs a parallel provenance algorithm

Usage: A dynamically linkable shared library

- **Source Code:** <https://github.com/Mic92/inspector>