# Dynamic Tracing: Memoization of Task Graphs for Dynamic Task-Based Runtimes

Wonchan Lee, Todd Warszawski, Alex Aiken
Elliott Slaughter
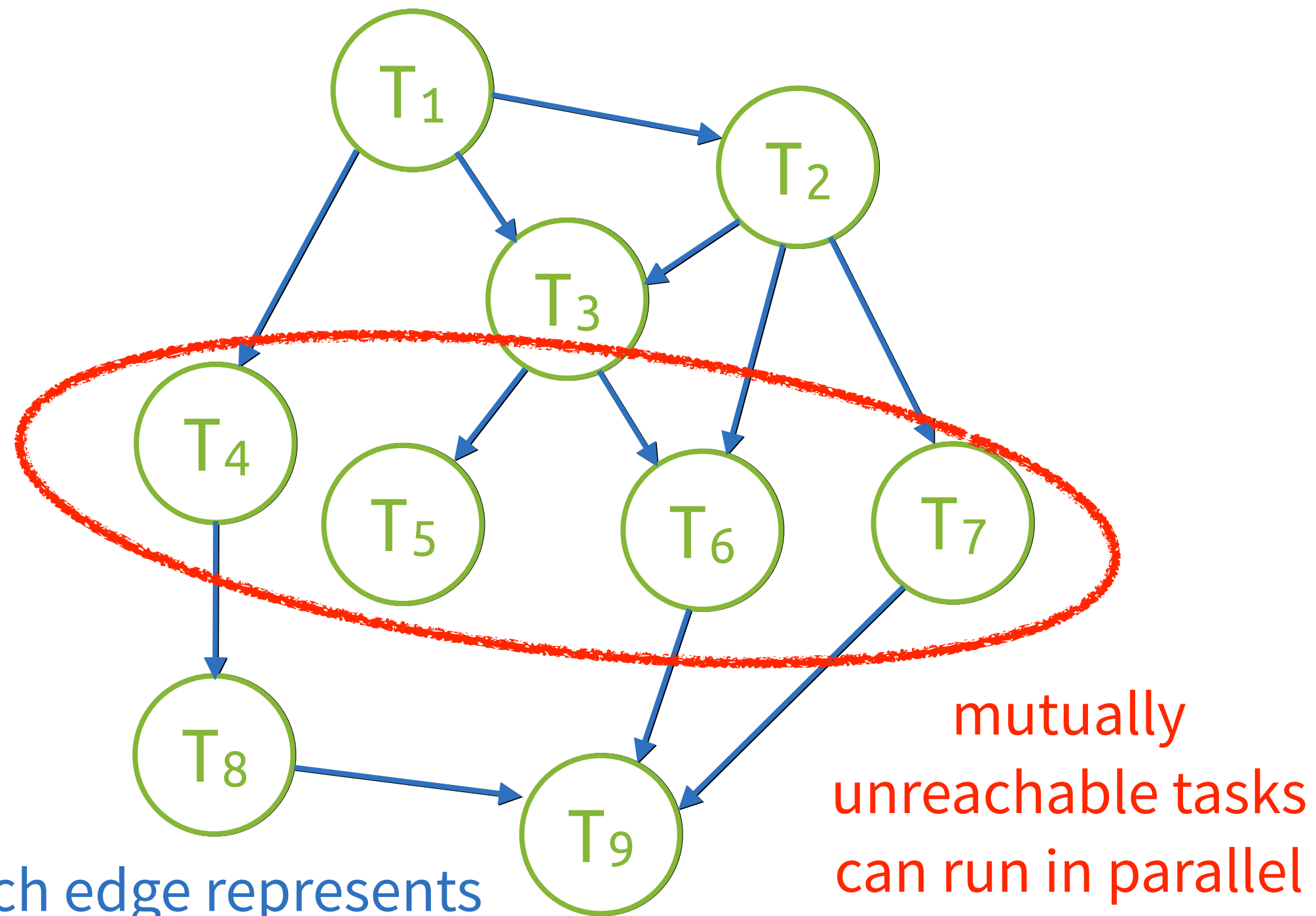Michael Bauer, Sean Treichler, Michael Garland

Stanford University
SLAC National Accelerator Laboratory
NVIDIA

November 14, 2018 @ SC'18

# Task Graphs Simplify Distributed Programming

Task graph is a DAG of tasks where

each task is an
opaque computation

$T_1$

$T_2$

$T_3$

$T_4$

$T_5$

$T_6$

$T_7$

$T_8$

$T_9$

mutually
unreachable tasks
can run in parallel

each edge represents
ordering between tasks

- Parallel execution is "straightforward" with task graphs

- Task graphs are most flexible when dynamically generated

  - Dynamic task graphs also facilitate fault recovery, load balancing, task (re-)mapping, resource allocation, etc.
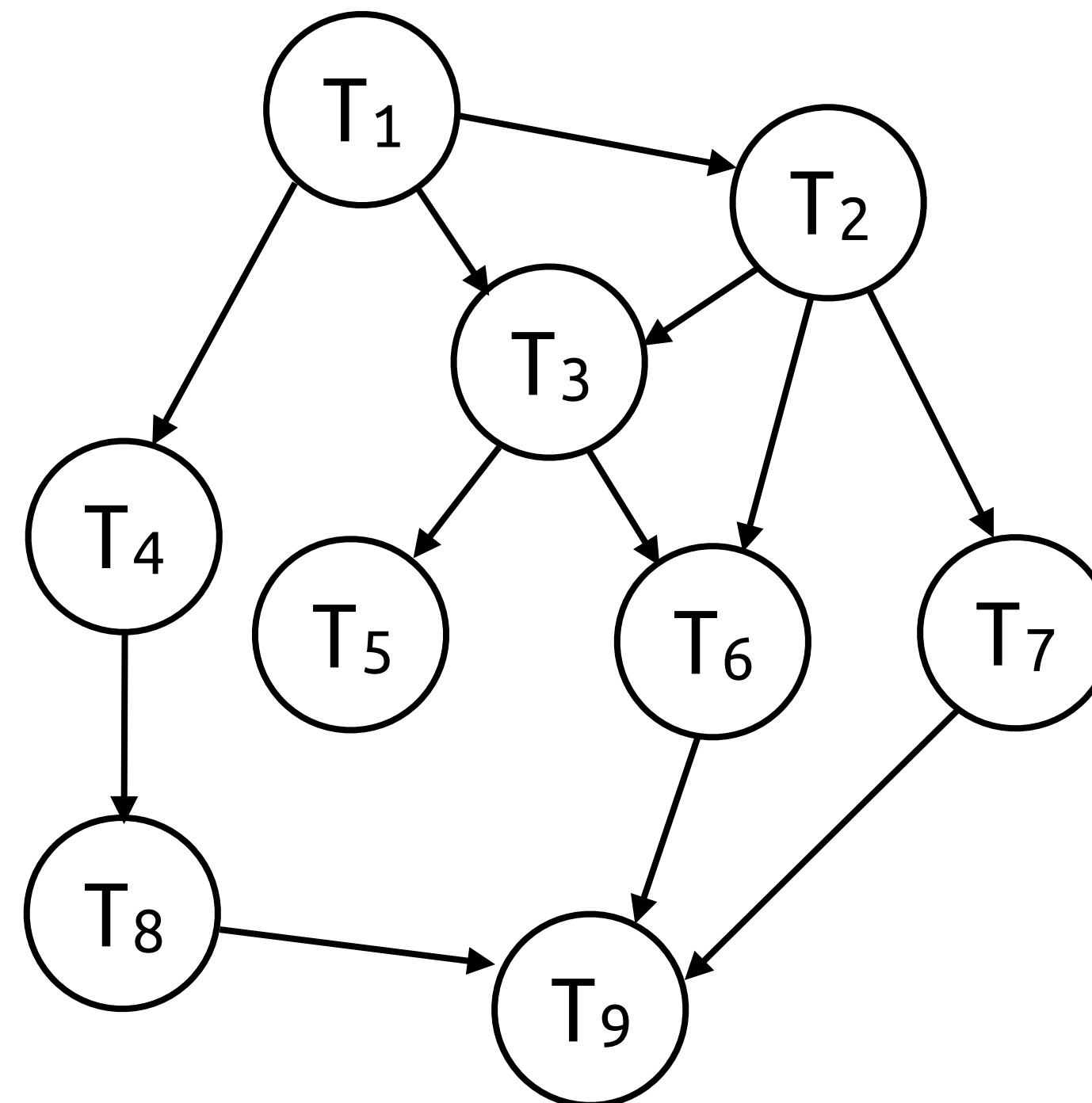
# Approaches to Dynamic Task Graph Construction

## Explicit construction

Program = Graph generator

```
g = new TaskGraph()
g.add(T₁)
g.add(T₂ ← T₁)
g.add(T₃ ← {T₁, T₂})
g.add(T₄ ← T₁)
…
```

**Exec.**

✔ Efficient
✘ Error-prone
✘ Not composable

e.g., Realm, CUDA

## Implicit construction

Program = Task generator

**Dep. Analysis**

$T_1(A,B)$ $T_2(A)$ $T_3(A)$ $T_4(B)$ ...

**Exec.**

```
T₁(A,B)  // writes(A),reads(B)
T₂(A)    // reads(A)
T₃(A)    // writes(A)
T₄(B)    // writes(B)
```

✔ Correct by construction
✔ Composable
✘ Runtime overhead

Dynamic task-based runtimes
(Legion, StarPU, PaRSEC, PyTorch, etc.)



**Is there a hybrid approach that enjoys benefits of both?**

# Dynamic Tracing: Memoizing Task Graphs

- Bring efficiency of explicit construction to dynamic task-based runtimes

- Key observation: programs often have *traces* of repetitive tasks

  - The same traces produce the same subgraph

Subgraphs are isomorphic!

```
task T(x,y) writes(x),reads(y)
task U(x,y) reads(x), reads(y)

while (*):
  T(A,B); T(C,D)
  U(A,D); U(C,B)
```
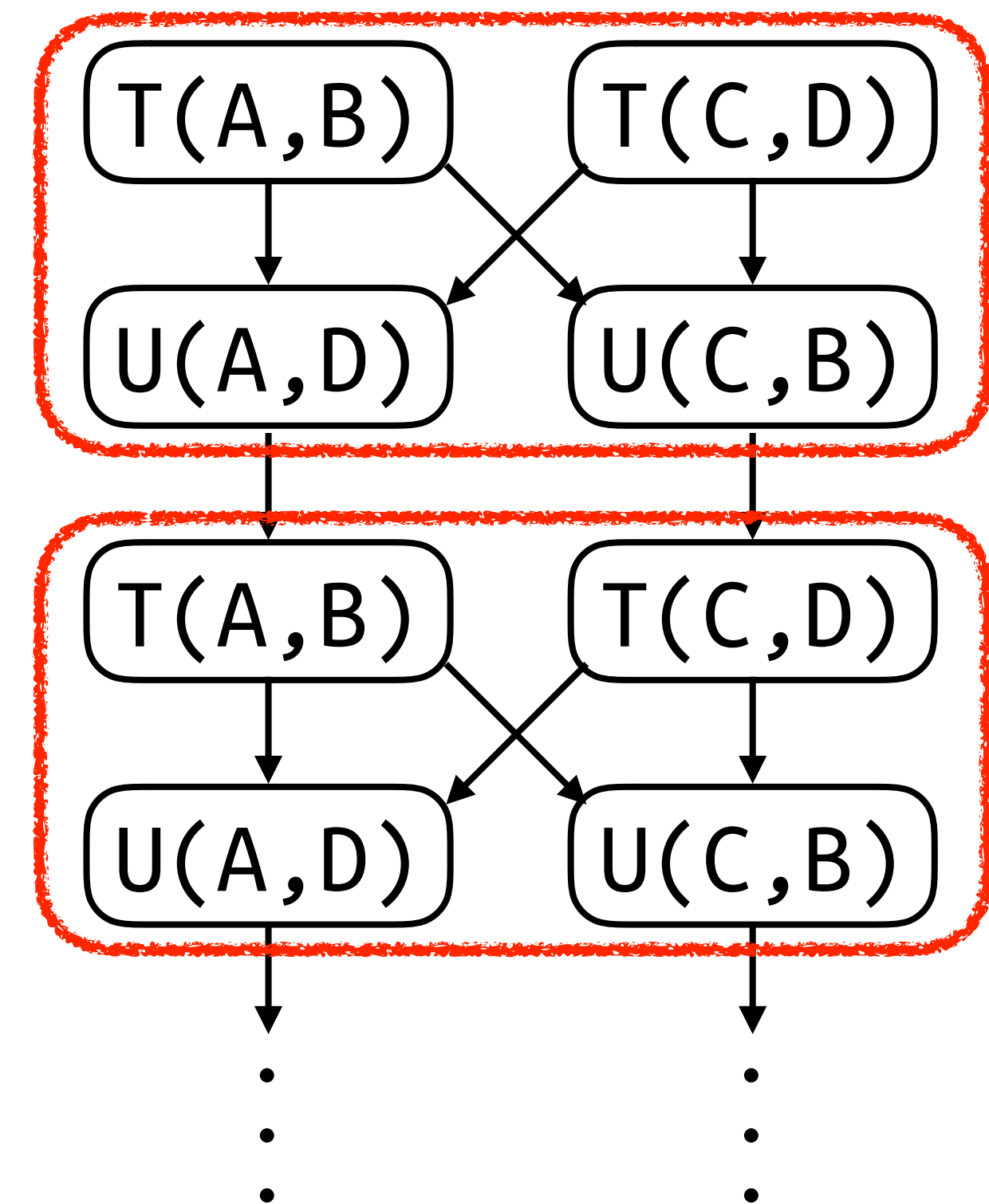
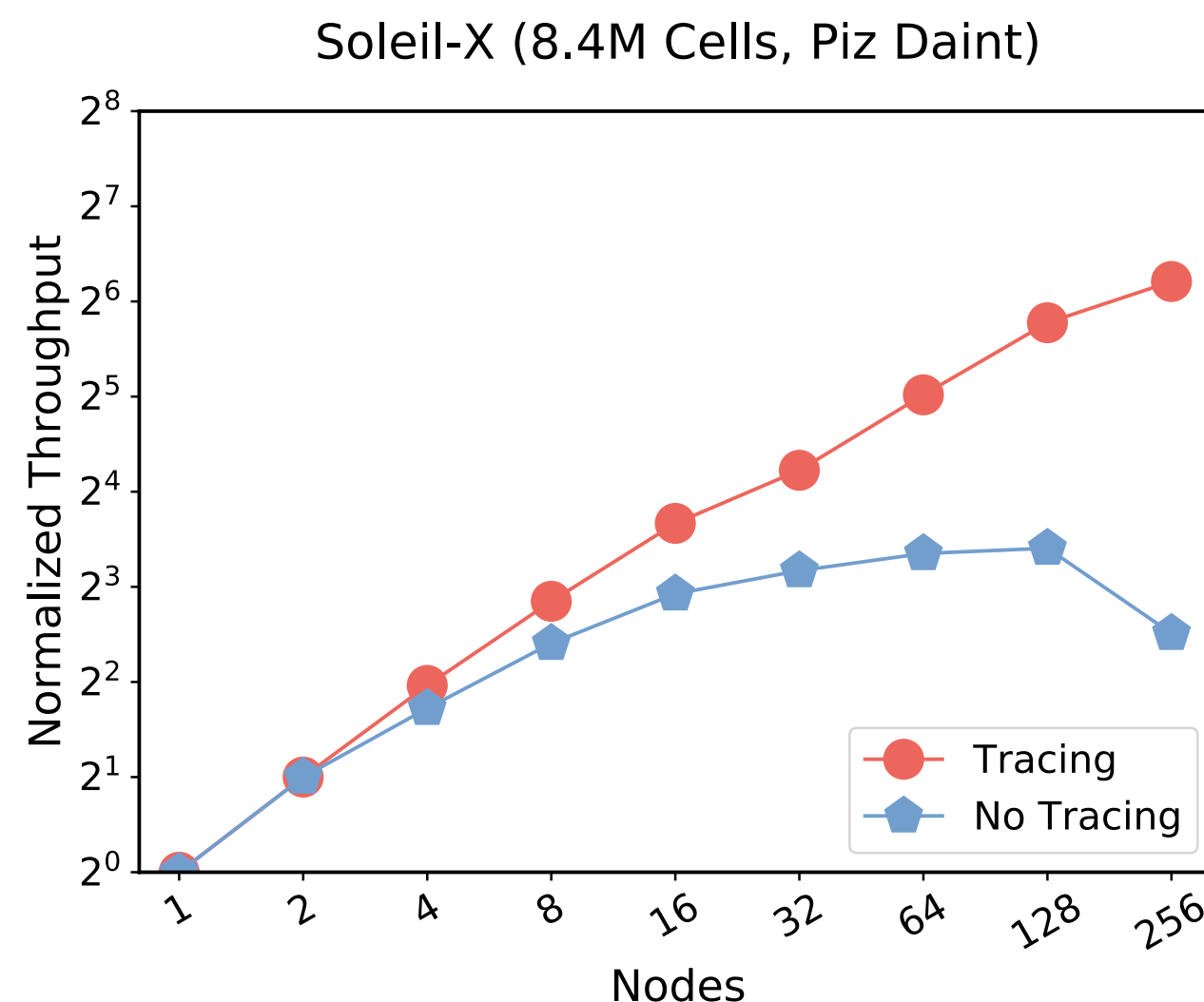**Exec.**

(T(A,B) T(C,D) U(A,D) U(C,B))*

Repetitive tasks

**Dep. Analysis**

# Dynamic Tracing: Memoizing Task Graphs

- Idea: record-and-replay

  - Record the subgraph once for a trace

  - Replay the recording whenever applicable

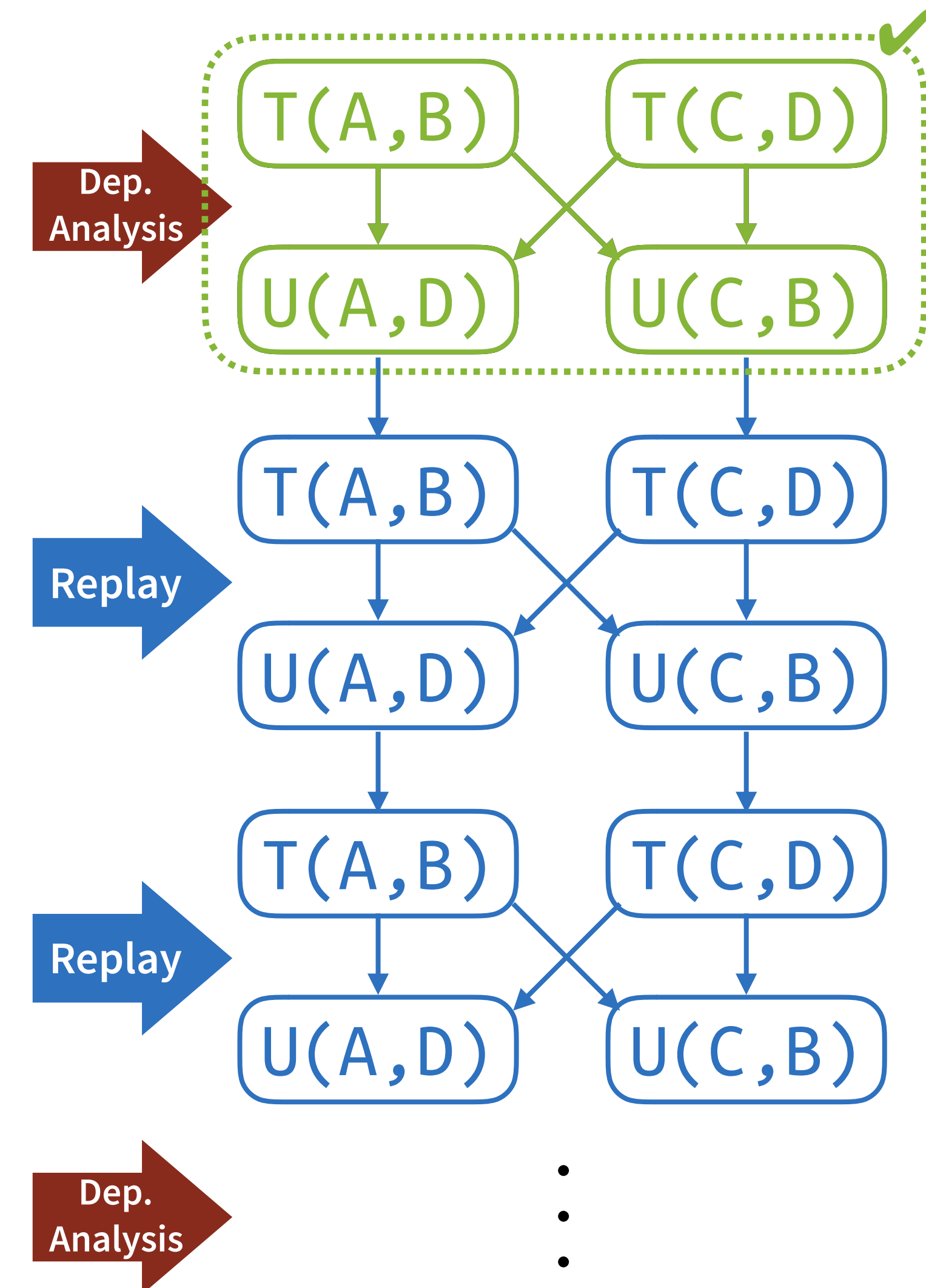- **Improves strong scaling performance by 4.9X**



Soleil-X (8.4M Cells, Piz Daint)

T(A,B)
T(C,D)
U(A,D)
U(C,B)

Dep. Analysis

T(A,B)    T(C,D)

U(A,D)    U(C,B)

T(A,B)
T(C,D)
U(A,D)
U(C,B)

Replay

T(A,B)    T(C,D)

U(A,D)    U(C,B)

T(A,B)
T(C,D)
U(A,D)
U(C,B)

Replay

T(A,B)    T(C,D)

U(A,D)    U(C,B)

S(A,C)

Dep. Analysis

# Contents

- Programming model

- Baseline dependence analysis

- Challenges in dynamic tracing

- Optimizations

- Experiment results

# Programming Model

- Task-based

  - Programs consist of tasks
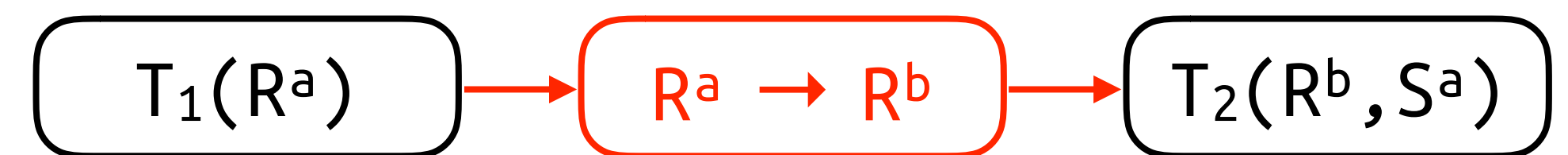
  - Tasks use regions and declare permissions

```
task T₁(x)   reads(x),writes(x)
task T₂(x,y) reads(x),writes(y)

T₁(R)
T₂(R,S)
```

- Distributed

  - Regions must be mapped to instances

  - One region can be mapped to multiple instances
    → *Coherence* must be maintained by the runtime

$T_1(R)$ $T_2(R,S)$ **Mapping** → $T_1(R^a)$ $T_2(R^b,S^a)$

**Dep. Analysis**

$T_1(R^a)$ → $R^a \rightarrow R^b$ → $T_2(R^b,S^a)$
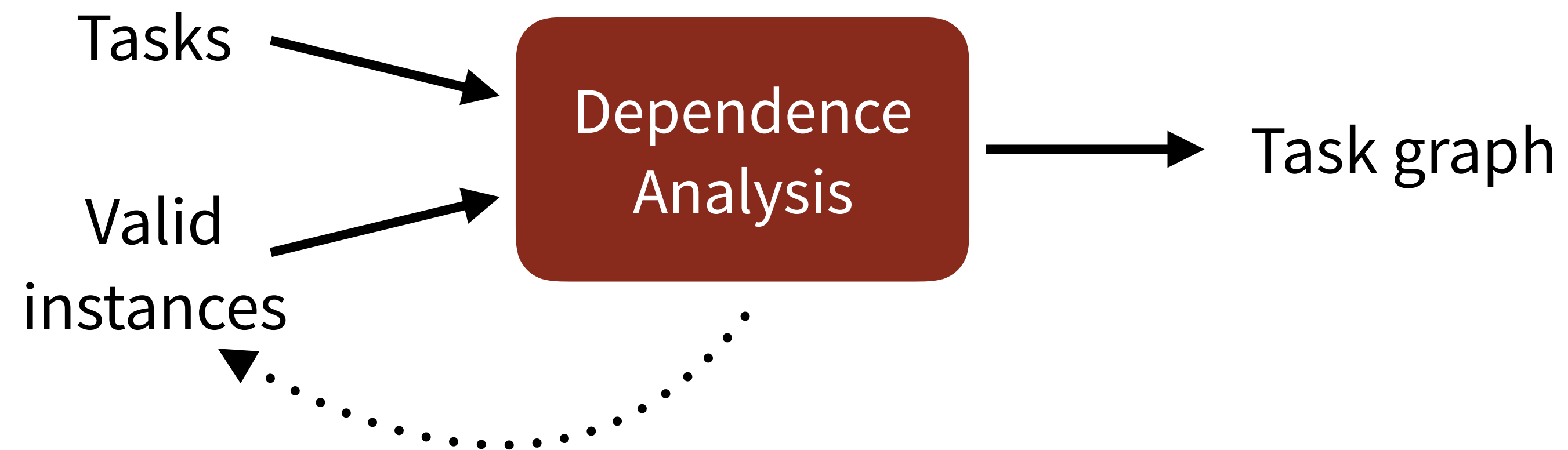
$R^a$ contains the last updates to R, or is *valid*

copy is issued for coherent access to R

$R^b$ is not valid yet for read access

# Baseline Dependence Analysis

Tasks

Valid
instances

Dependence
Analysis

Task graph

# Baseline Dependence Analysis

$T_1(R^a)$      **reads** $(R^a)$,**writes**$(R^a)$
$T_2(R^b,S^a)$    **reads** $(R^b)$,**writes**$(S^a)$
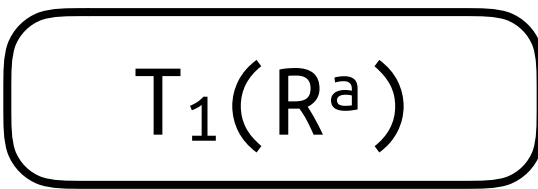$T_3(R^a,S^a)$    **writes**$(R^a)$,**reads** $(S^a)$

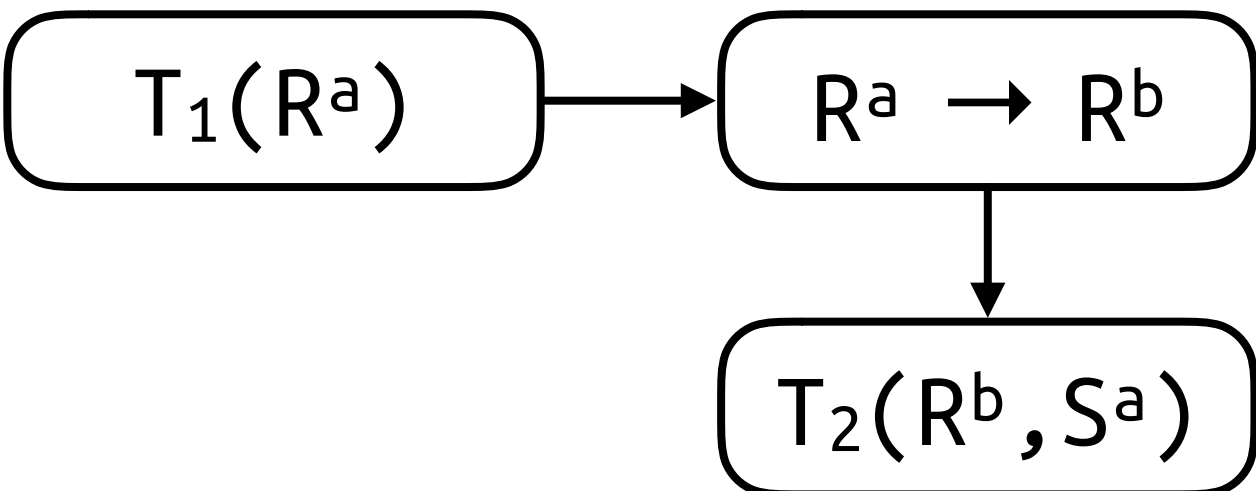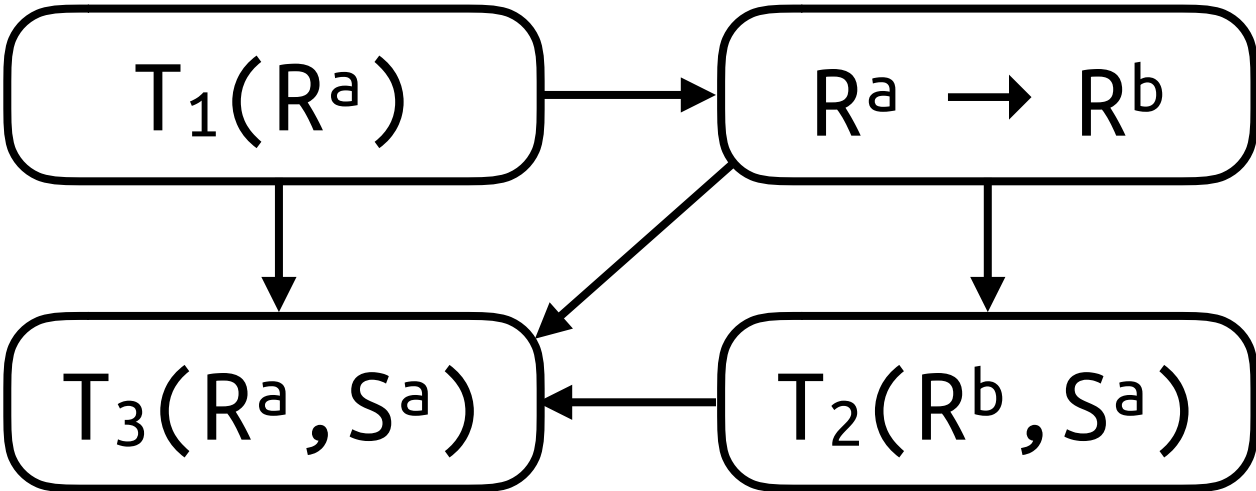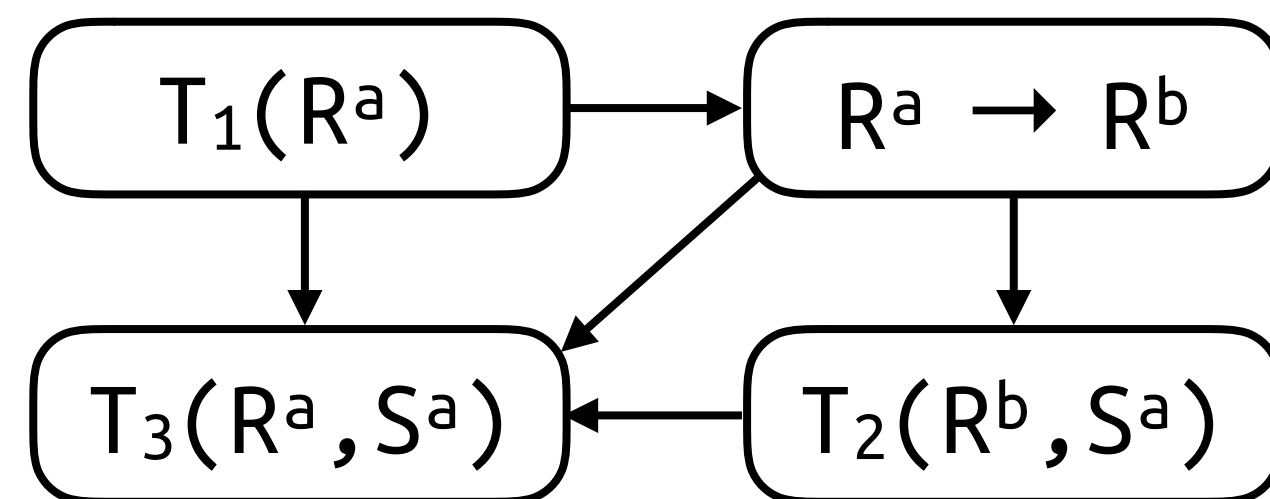| Task | Valid Instances | Task graph |
|------|-----------------|------------|
| $T_1(R^a)$ | R: $R^a$ | $T_1(R^a)$ |
| $T_2(R^b,S^a)$ | R: $R^b$,$R^a$ <br><br> S: $S^a$ | $T_1(R^a)$ → $R^a \rightarrow R^b$ → $T_2(R^b,S^a)$ |
| $T_3(R^a,S^a)$ | R: $R^a$ <br><br> S: $S^a$ | $T_1(R^a)$ → $R^a \rightarrow R^b$ ; $T_1 \rightarrow T_3(R^a,S^a)$ ; $R^a \rightarrow R^b \rightarrow T_2(R^b,S^a)$ → $T_3(R^a,S^a)$ |

# Challenges in Dynamic Tracing

- Challenge 1: transition from dep. analysis to subgraph replay
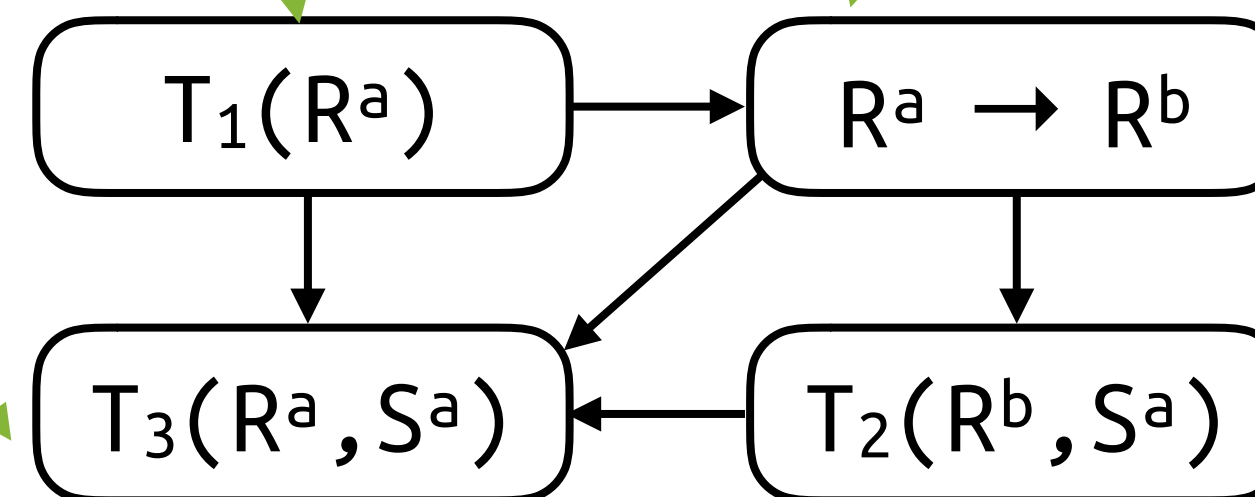
Captured subgraph **G**



Found the same trace
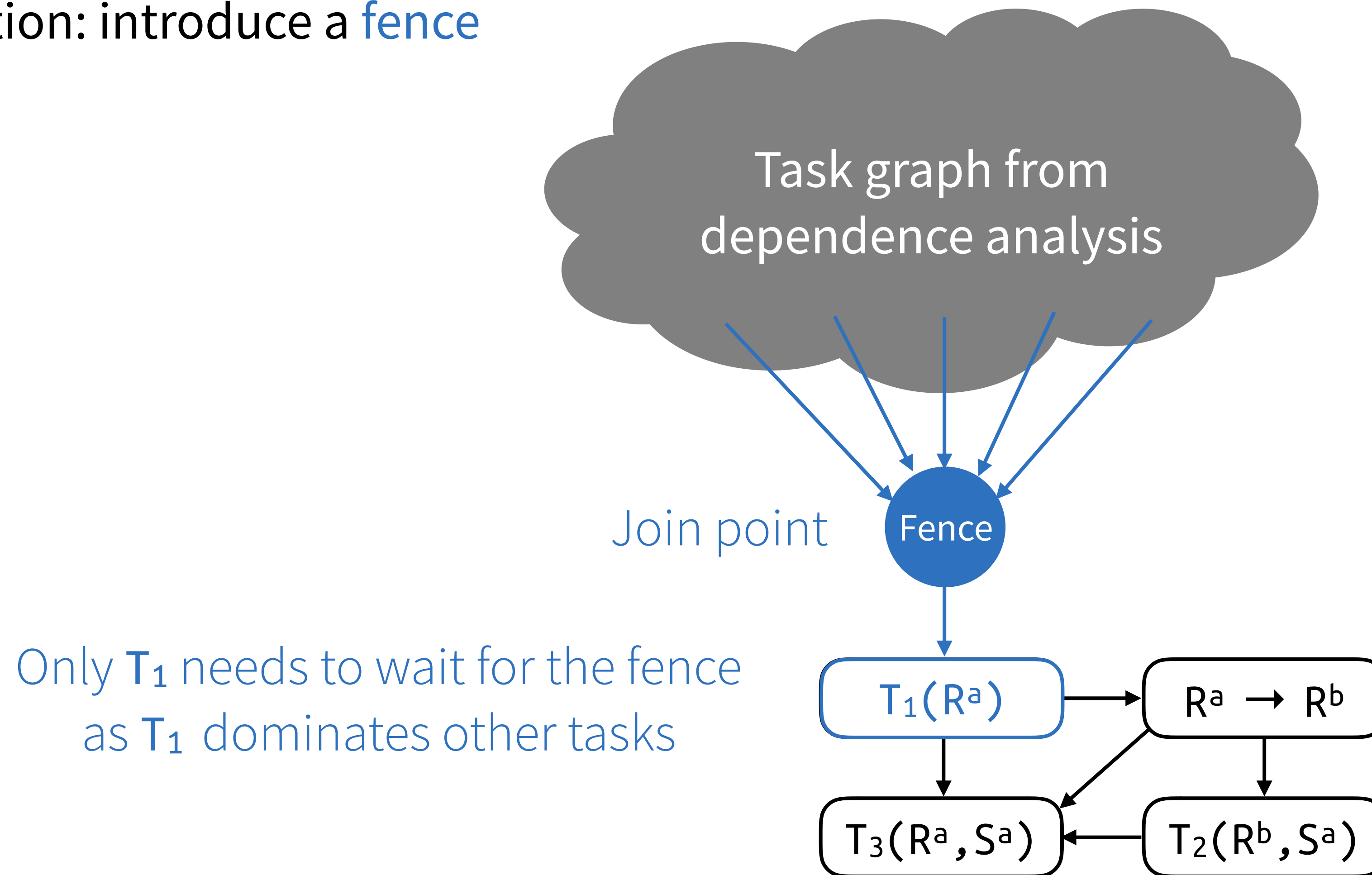
... $T_1(R^a)$  $T_2(R^b, S^a)$  $T_3(R^a, S^a)$  ...

Task graph from dependence analysis



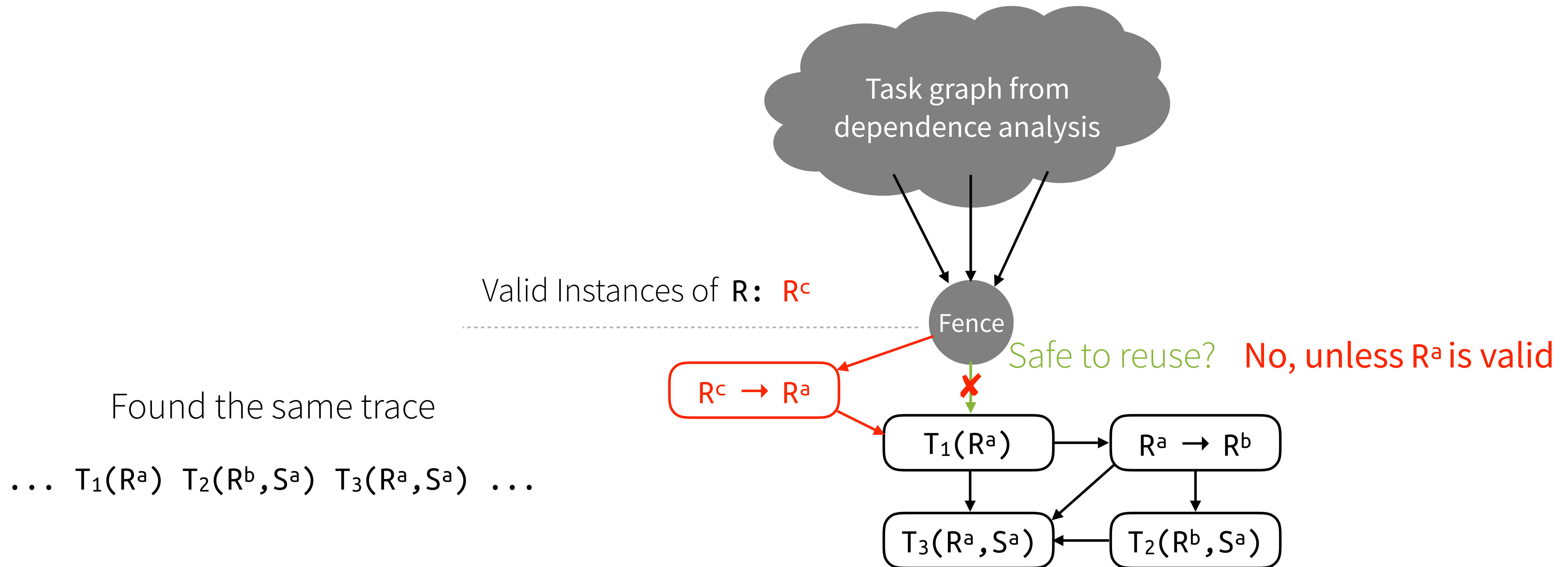How can we connect G to the graph from dep. analysis?

- Solution: introduce a fence

Task graph from dependence analysis

Join point — Fence

Only $T_1$ needs to wait for the fence as $T_1$ dominates other tasks

$T_1(R^a)$ → $R^a \rightarrow R^b$

$T_3(R^a, S^a)$ ← $T_2(R^b, S^a)$

# Challenges in Dynamic Tracing

- Challenge 2: coherence



Valid Instances of **R:** $R^c$

Found the same trace

$\ldots\ T_1(R^a)\ \ T_2(R^b,S^a)\ \ T_3(R^a,S^a)\ \ldots$

Task graph from dependence analysis

Fence

Safe to reuse?   No, unless $R^a$ is valid

$R^c \rightarrow R^a$

$T_1(R^a)$     $R^a \rightarrow R^b$

$T_3(R^a,S^a)$     $T_2(R^b,S^a)$

- Solution: remember precondition for a safe replay

Tasks ────→

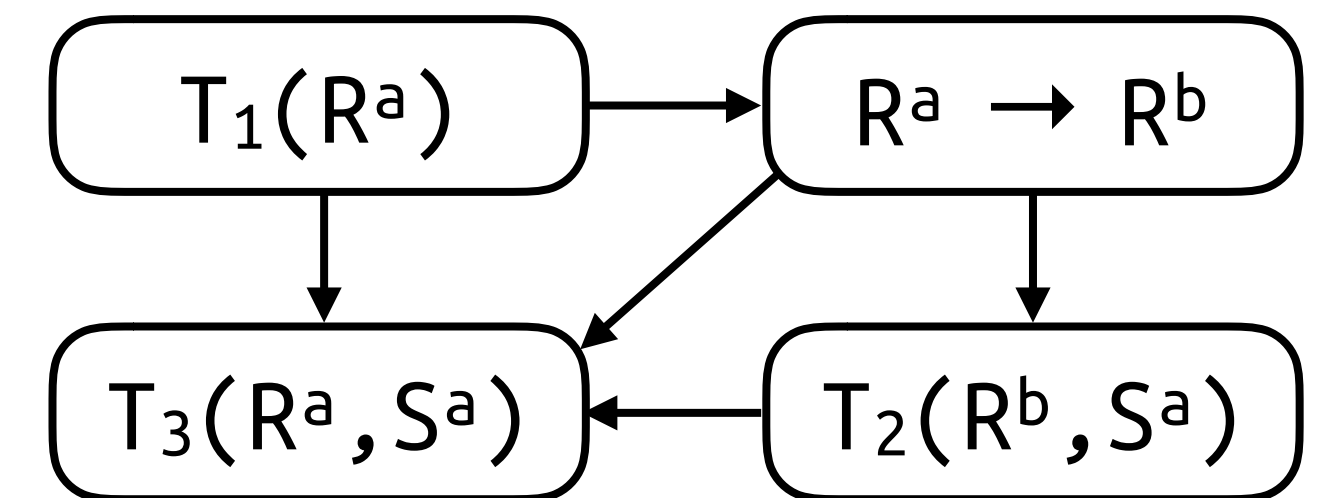Valid instances ────→

**Dependence Analysis** ────→ Task graph

Valid instances are input of dependence analysis

Tasks: $T_1(R^a)$ $T_2(R^b, S^a)$ $T_3(R^a, S^a)$

Precondition: $R^a$ is valid

Task graph:

$T_1(R^a)$ ──→ $R^a \rightarrow R^b$

$T_1(R^a)$ ──→ $T_3(R^a, S^a)$

$R^a \rightarrow R^b$ ──→ $T_3(R^a, S^a)$

$R^a \rightarrow R^b$ ──→ $T_2(R^b, S^a)$
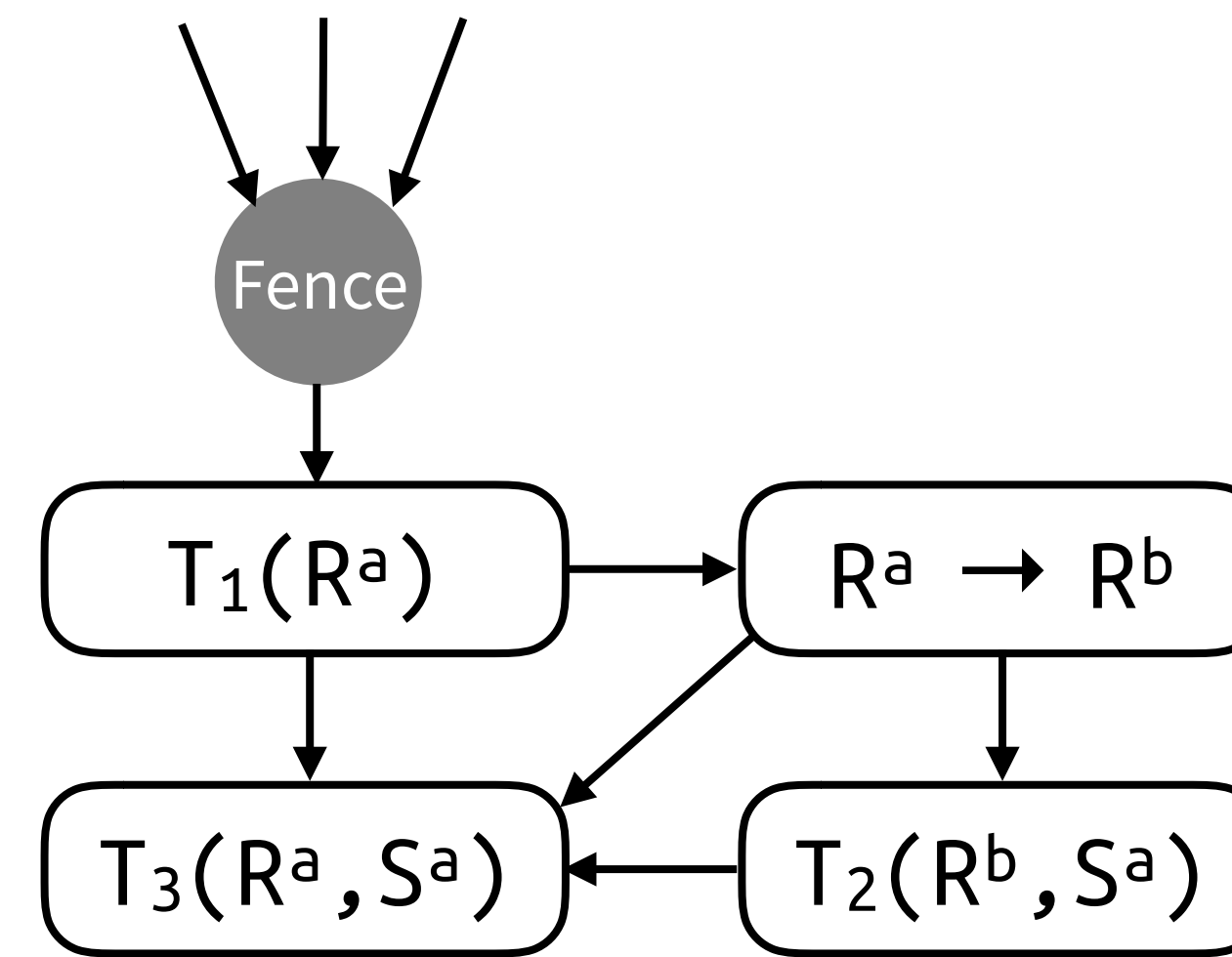
$T_2(R^b, S^a)$ ──→ $T_3(R^a, S^a)$

- Challenge 3: transition back to normal dep. analysis

Two inconsistencies
1. $T_1$, $T_2$, and $T_3$ are unknown to dep. analysis
2. The list of valid instances is stale
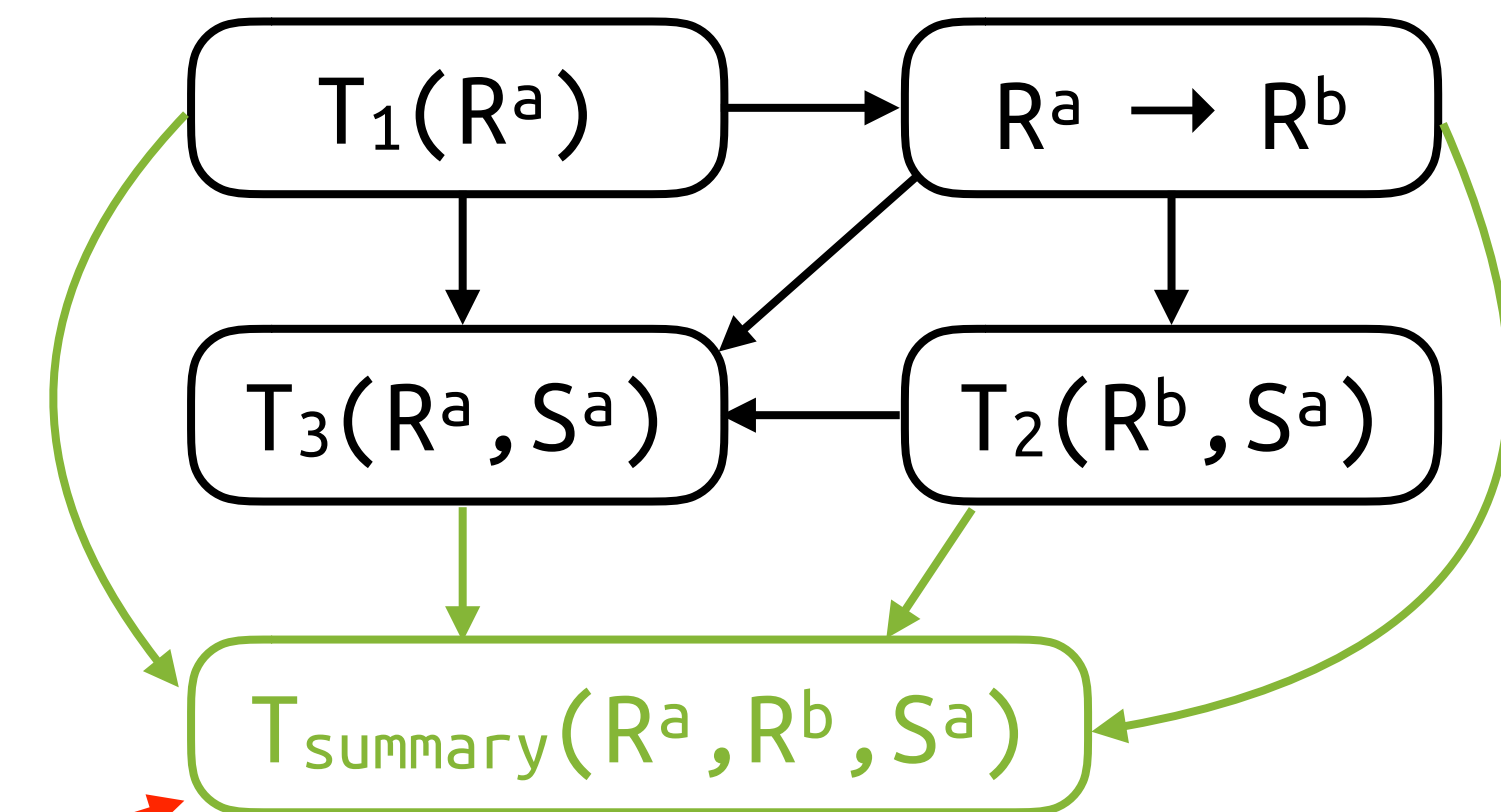
$$U(R^b, S^c)$$

**Dep. Analysis**

Fence

$T_1(R^a)$   →   $R^a → R^b$

$T_3(R^a, S^a)$   ←   $T_2(R^b, S^a)$

- Solution

  - Make a summary task

  - Compute postcondition to apply after each replay

Tasks: $T_1(R^a)$  $T_2(R^b,S^a)$  $T_3(R^a,S^a)$

Precondition: $R^a$ is valid

Task graph:


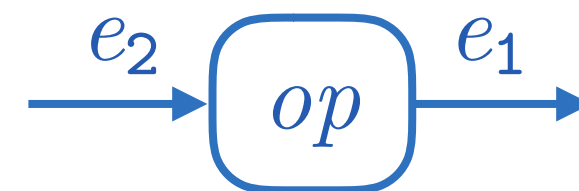
Summary task goes through normal dependence analysis

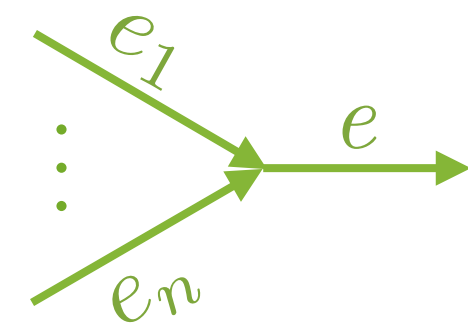Postcondition: $R^a$ and $S^a$ become valid

15

# Graph Calculus

- Simple graph construction language

  - Use *events* that signify termination of operations

  - Syntax:
    $$c ::= \quad e_1 := \mathtt{op}(op, e_2) \mid e := \mathtt{merge}(\overline{e_i}) \mid e := \mathtt{fence} \mid c; c$$
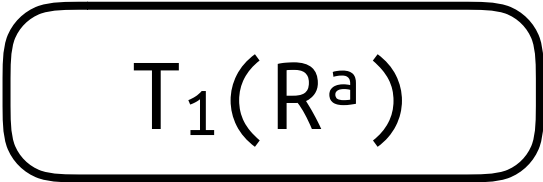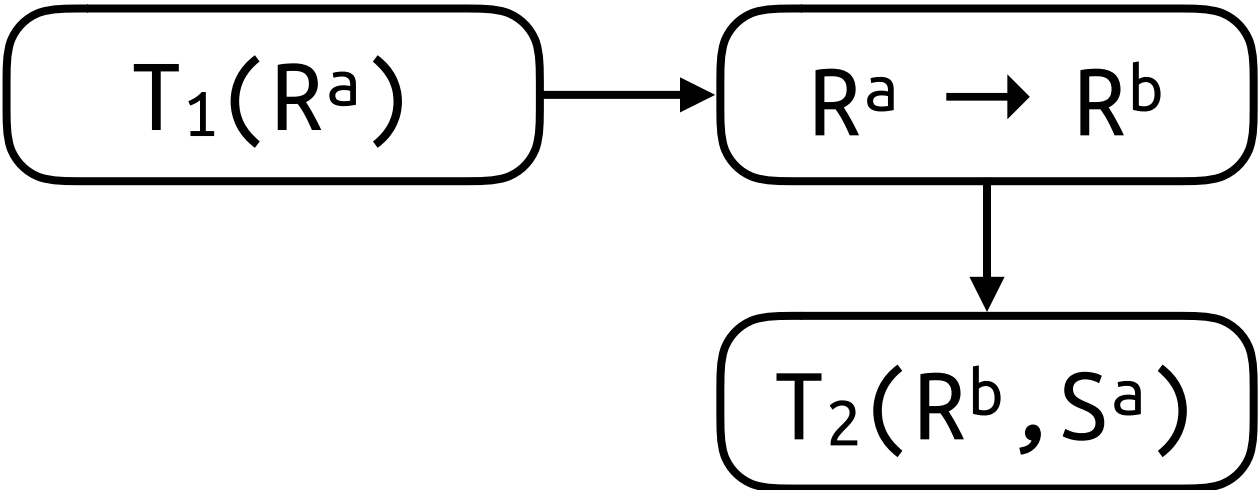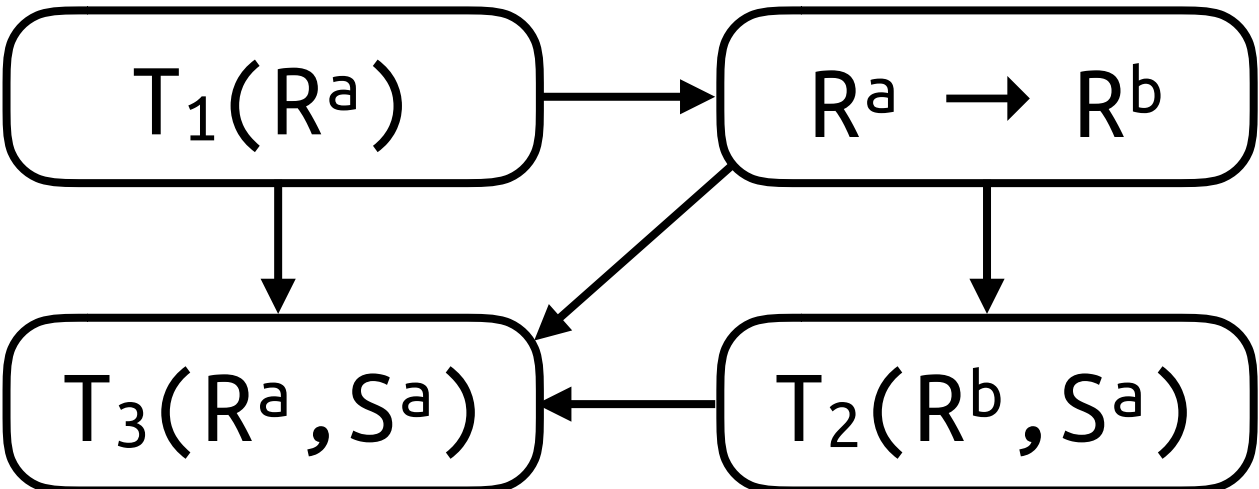
    $op$ waits until $e_2$ is signaled and triggers $e_1$ when it's done

    $e$ is triggered when $e_1, \ldots, e_n$ are

    issue a fence that signals $e$
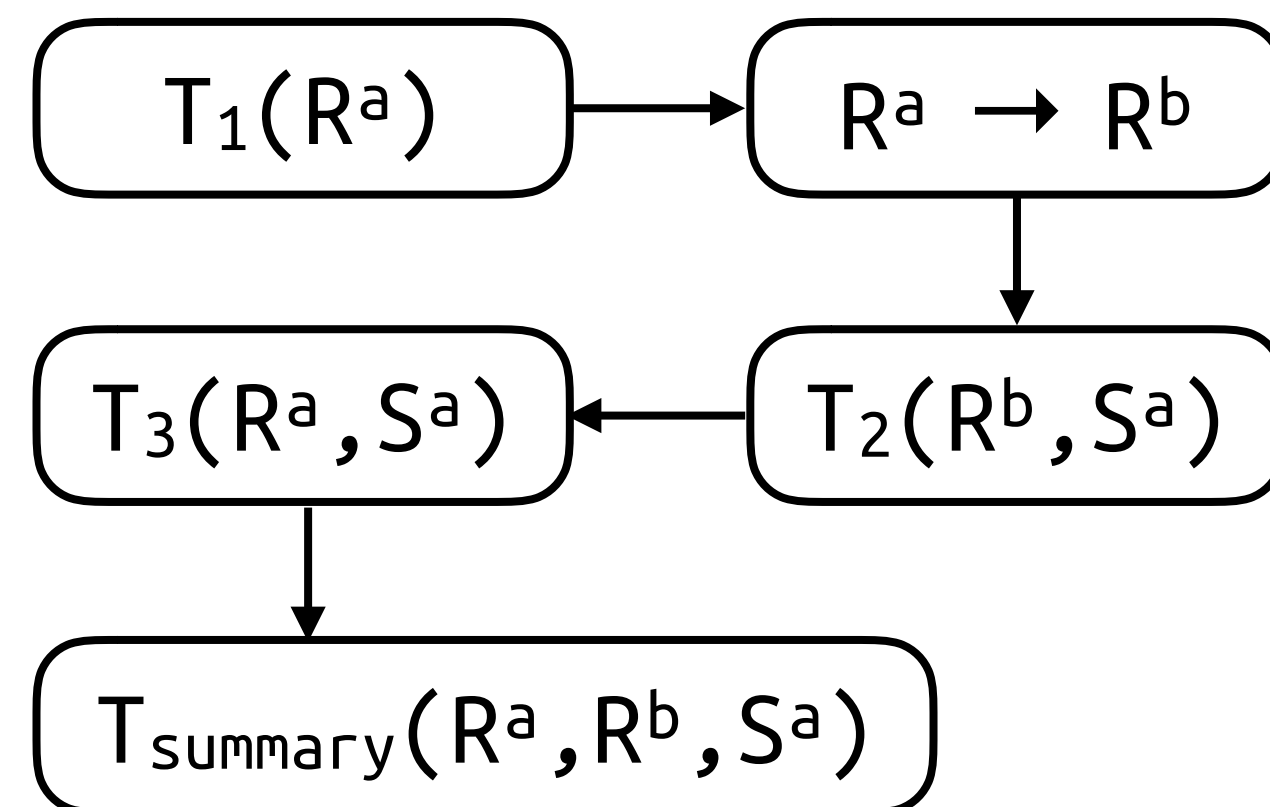
# Trace Recording Example

| Task | Task graph | Command | Recording state |
|------|-----------|---------|-----------------|
| $T_1(R^a)$ | $T_1(R^a)$ | `e₁ := fence`<br>`e₂ := op(T₁(Rᵃ), e₁)` | $T_1(R^a) = e_2$ |
| $T_2(R^b,S^a)$ | $T_1(R^a) \rightarrow R^a \rightarrow R^b$<br>$\downarrow$<br>$T_2(R^b,S^a)$ | `e₃ := op(Rᵃ ➜ Rᵇ, e₂)`<br>`e₄ := op(T₂(Rᵇ,Sᵃ), e₃)` | $R^a \rightarrow R^b = e_3$<br>$T_2(R^b,S^a) = e_4$ |
| $T_3(R^a,S^a)$ | $T_1(R^a) \rightarrow R^a \rightarrow R^b$<br>$T_3(R^a,S^a) \leftarrow T_2(R^b,S^a)$ | `e₅ := merge(e₂, e₃, e₄)`<br>`e₆ := op(T₃(Rᵃ,Sᵃ), e₅)` | $T_3(R^a,S^a) = e_6$ |
| Insert summary task | | `e₇ := merge(e₂, e₃, e₄, e₆)`<br>`e₈ := op(T_summary(Rᵃ,Rᵇ,Sᵃ), e₇)` | Pre: $R^a$<br>Post: $R^a,S^a$ |

# Idempotent Recordings

- When the postcondition subsumes the precondition

Task graph G:

$$T_1(R^a) \rightarrow R^a \rightarrow R^b$$

$$\downarrow$$

$$T_3(R^a, S^a) \leftarrow T_2(R^b, S^a)$$

$$\downarrow$$

$$T_{summary}(R^a, R^b, S^a)$$

Precondition: $R^a$ is valid

Postcondition: $R^a$ and $S^a$ become valid

→ Precondition is satisfied immediately after
   postcondition is applied

- Optimization: precondition check elision
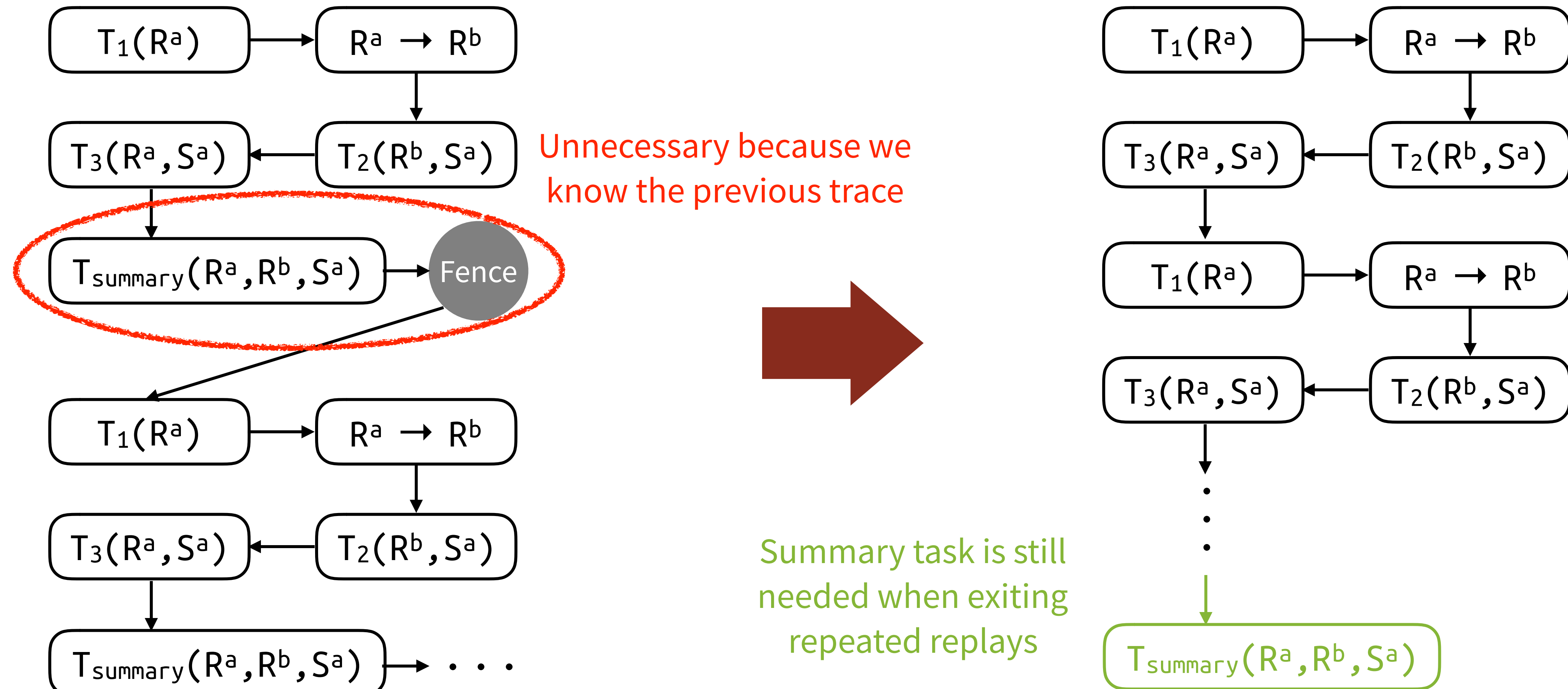  (when the same trace repeatedly appears)

$$(\text{Check pre.} \rightarrow \text{Replay} \rightarrow \text{Apply post.})^*$$

$$\Downarrow$$

$$\text{Check pre.} \rightarrow (\text{Replay})^* \rightarrow \text{Apply post.}$$

# Fence Elision

- We can remove summary tasks and fences when we replay the same trace repeatedly



Unnecessary because we know the previous trace

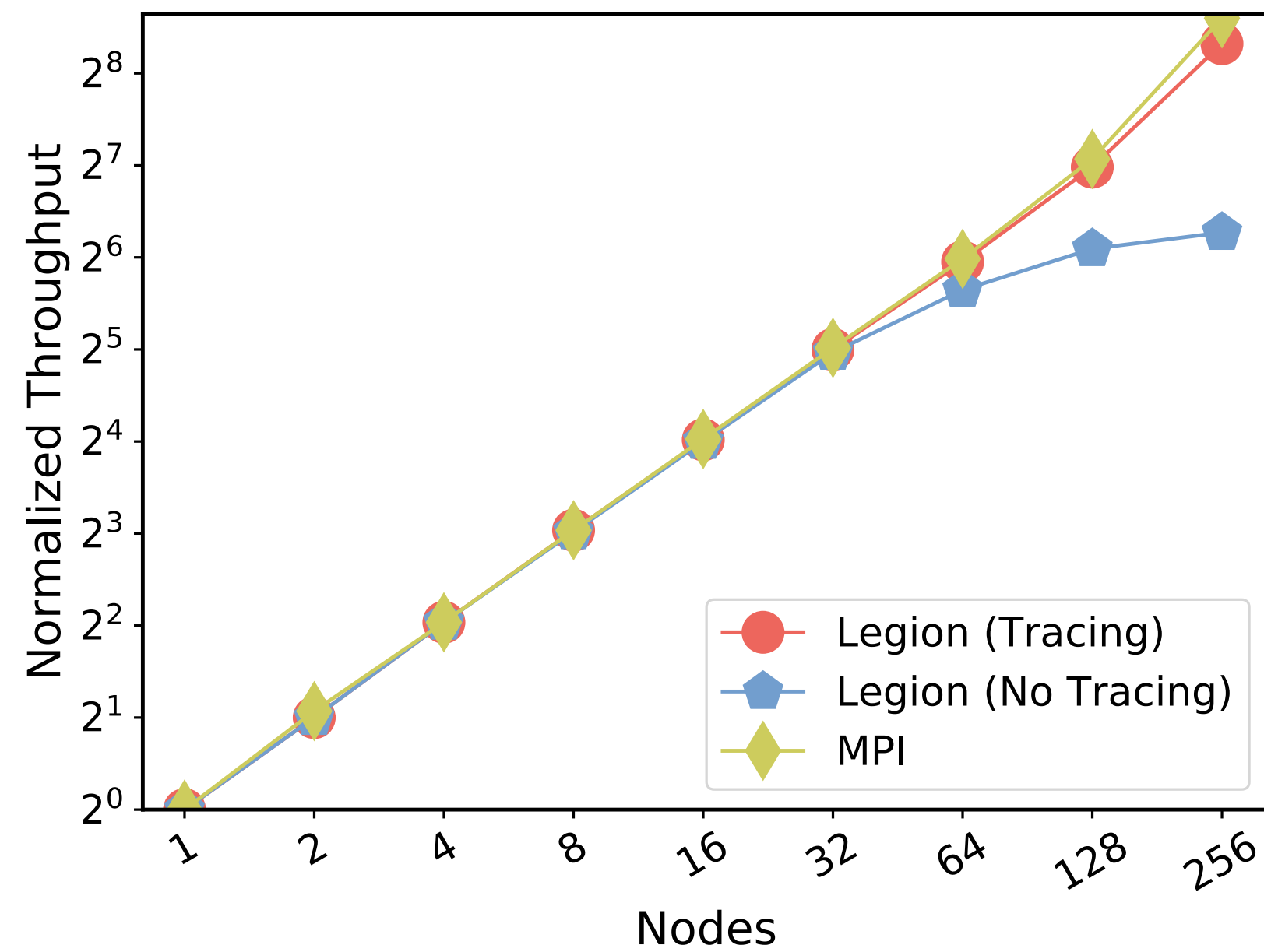Summary task is still needed when exiting repeated replays

# Experiment Results

- Implemented dynamic tracing in Legion

- Measure strong scaling performance of five Legion applications

  - Varying complexity (from 9-point stencil to multi-physics solver)

  - Already optimized for weak scaling performance[†]

  - Machine: Piz Daint (Cray XC50, Xeon E5-2690 with 12 cores & 64 GB memory per node)

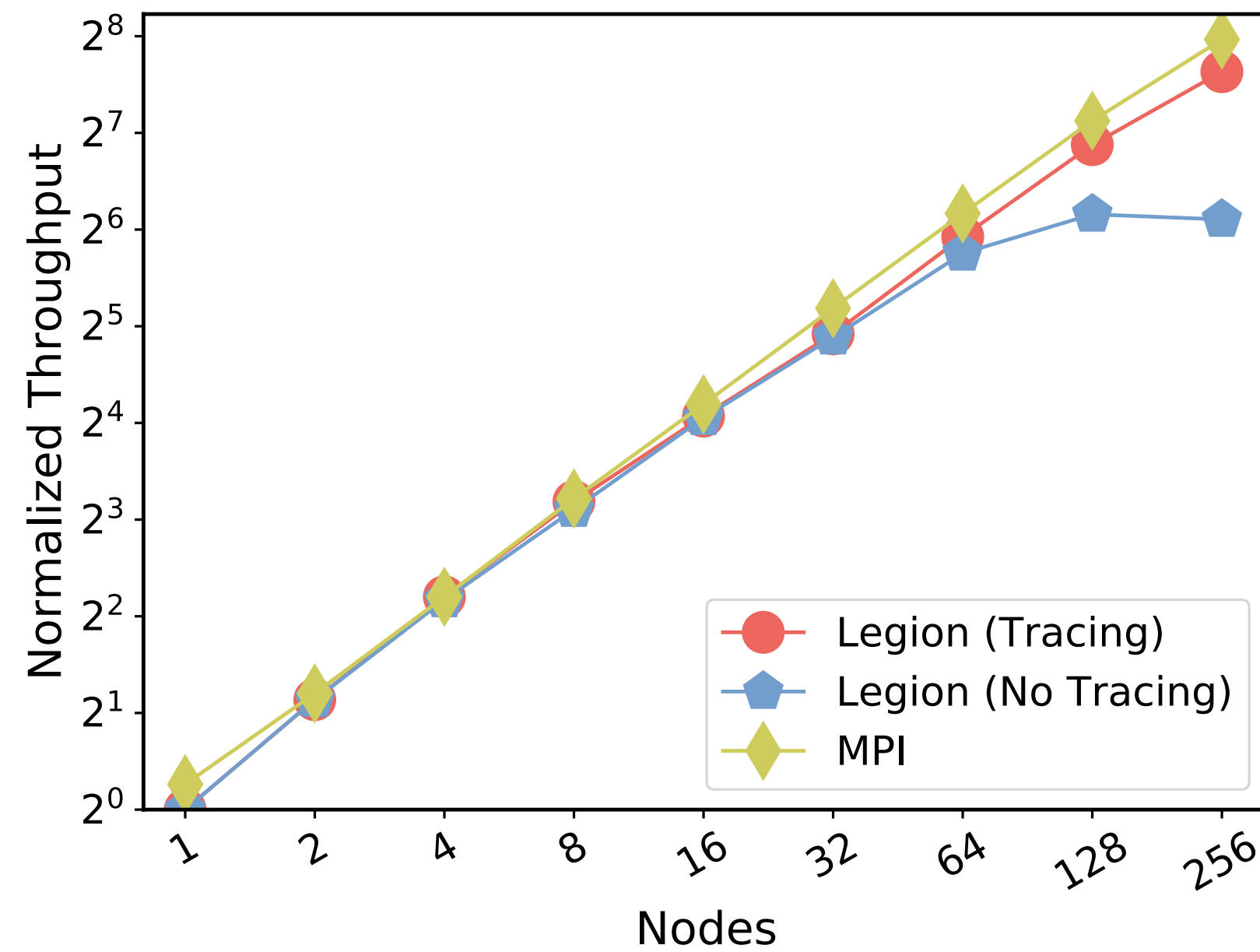- Compare with MPI references for Stencil, MiniAero, and PENNANT

† E. Slaughter, W. Lee, S. Treichler, W. Zhang, M. Bauer, G. Shipman, P. McCormick, and A. Aiken, "Control replication: Compiling implicit parallelism to efficient SPMD with logical regions," SC'17
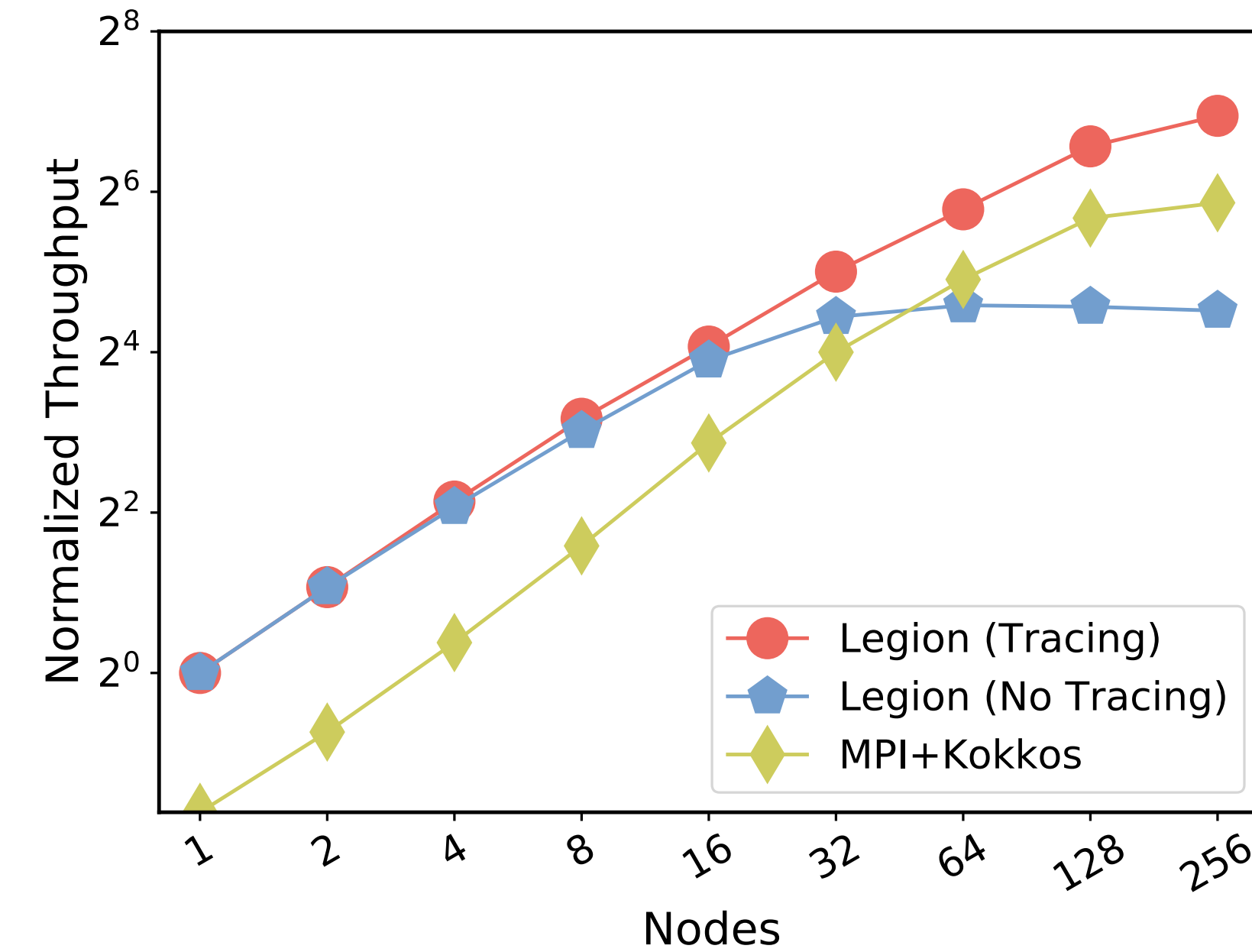
# Strong Scaling Performance



Stencil (0.4B Cells, Piz Daint)

PENNANT (29M Zones, Piz Daint)

MiniAero (1M Cells, Piz Daint)

| | | | |
|---|---|---|---|
| Trace size / node | 47 | 121 | 210 |
| Improvement | 4.2X | 2.8X | 5.1X |
| DT / MPI | 82% | 79% | 212% |

Dynamic time stepping blocks runtime analysis every iteration

Legion spares 3 cores for runtime

Legion version uses a better data layout

# Strong Scaling Performance

### Circuit (74K Wires, Piz Daint)



### Soleil-X (8.4M Cells, Piz Daint)



| | Circuit | Soleil-X |
|---|---|---|
| Trace size / node | 76 | 344 |
| Improvement | 5.3X | 7.0X |

# Conclusion

- Dynamic tracing brings performance of explicit task graph construction to dynamic task-based runtimes

    - Strong scaling performance is improved by 4.9X on average

- Feel free to try out Dynamic Tracing!

    - Checked in to the Legion repository: https://github.com/StanfordLegion/legion

    - Experiment scripts are here: https://gitlab.com/StanfordLegion/legion/tree/tracing-sc18

# Acknowledgment

# Questions?

# Programming Model

- Traces are annotated in programs

  - Places where tracing is beneficial are often obvious

  - Finding such places is important, but an orthogonal issue
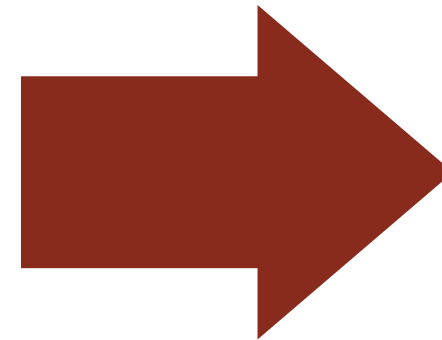
```
task T(x,y) writes(x),reads(y)
task U(x,y) reads(x), reads(y)

while (*):
  begin_trace
  T(A,B); T(C,D)
  U(A,D); U(C,B)
  end_trace
```

# Optimizing Graph Calculus Commands

- Two standard optimizations: transitive reduction and copy propagation

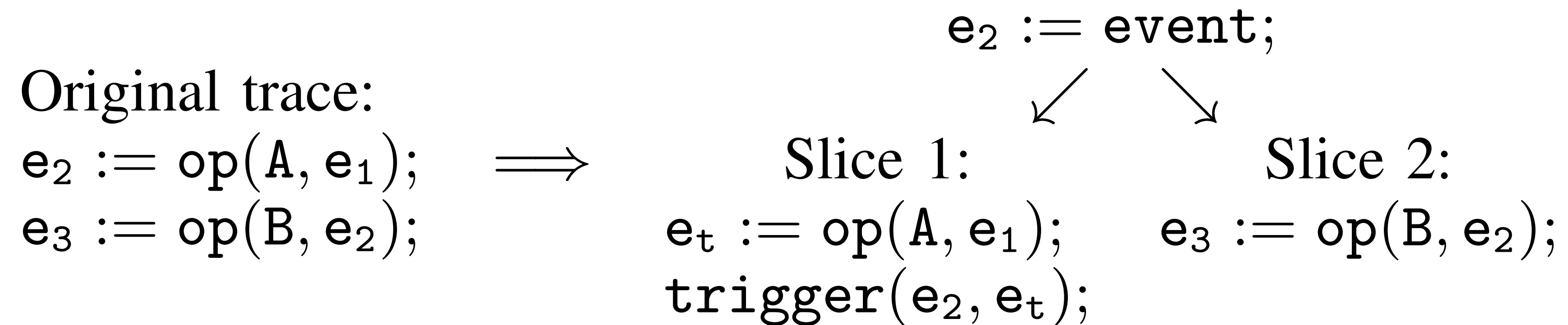  - The overhead is amortized by repeated replays

```
e₁ := fence
e₂ := op(T₁(Rᵃ), e₁)
e₃ := op(Rᵃ → Rᵇ, e₂)
e₄ := op(T₂(Rᵇ,Sᵃ), e₃)
e₅ := merge(e₂, e₃, e₄)
e₆ := op(T₃(Rᵃ,Sᵃ), e₅)
e₇ := merge(e₂, e₃, e₄, e₆)
e₈ := op(T_summary(Rᵃ,Rᵇ,Sᵃ), e₇)
```

➡

```
e₁ := fence
e₂ := op(T₁(Rᵃ), e₁)
e₃ := op(Rᵃ → Rᵇ, e₂)
e₄ := op(T₂(Rᵇ,Sᵃ), e₃)
e₅ := merge(e₂, e₃, e₄)
e₆ := op(T₃(Rᵃ,Sᵃ), e₄)
e₇ := merge(e₂, e₃, e₄, e₆)
e₈ := op(T_summary(Rᵃ,Rᵇ,Sᵃ), e₆)
```

# Parallel Replays

- Trace replay can be a bottleneck if the trace is long

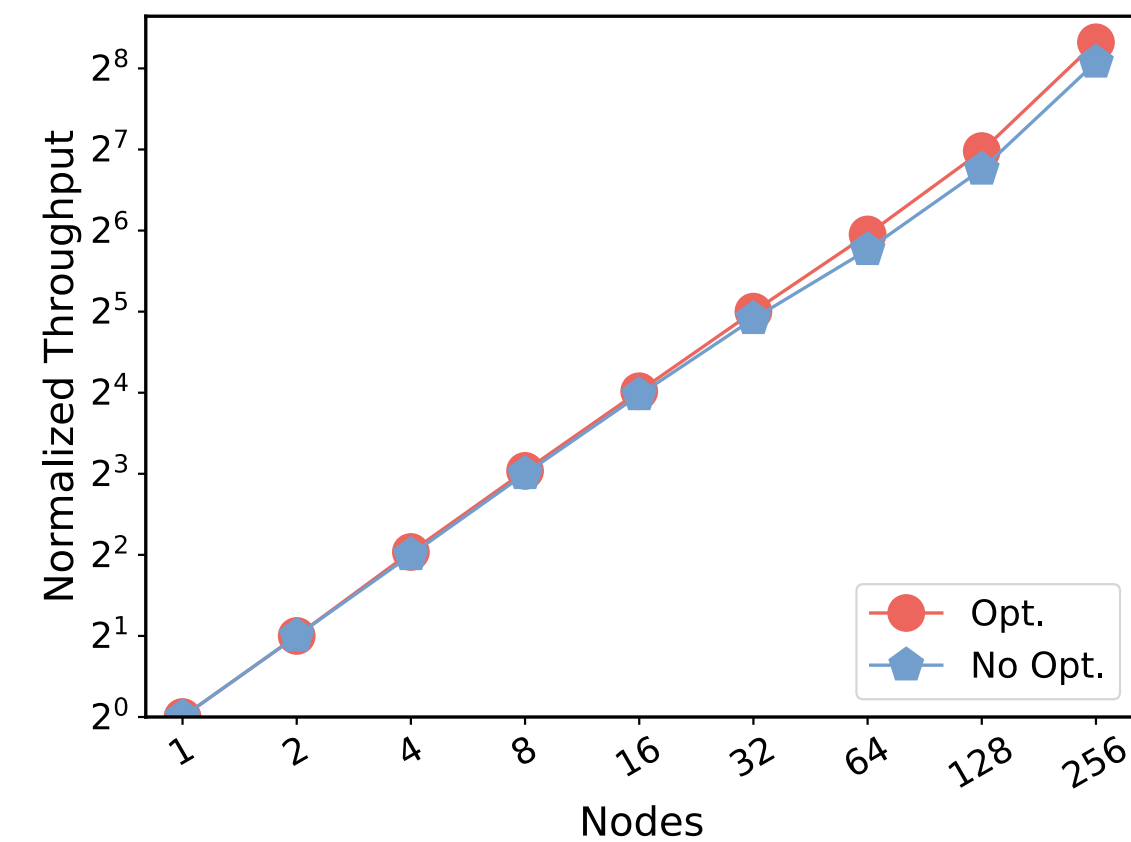- We can parallelize trace replay by slicing the trace

$$\texttt{e}_2 := \texttt{event};$$

$$\swarrow \qquad \searrow$$

Original trace:

$$\texttt{e}_2 := \texttt{op}(\texttt{A}, \texttt{e}_1); \quad \Longrightarrow \qquad \text{Slice 1:} \qquad\qquad \text{Slice 2:}$$

$$\texttt{e}_3 := \texttt{op}(\texttt{B}, \texttt{e}_2); \qquad\qquad \texttt{e}_t := \texttt{op}(\texttt{A}, \texttt{e}_1); \quad \texttt{e}_3 := \texttt{op}(\texttt{B}, \texttt{e}_2);$$

$$\texttt{trigger}(\texttt{e}_2, \texttt{e}_t);$$

Extended graph calculus $\quad c ::= \cdots \mid e := \texttt{event} \mid \texttt{trigger}(e, e)$
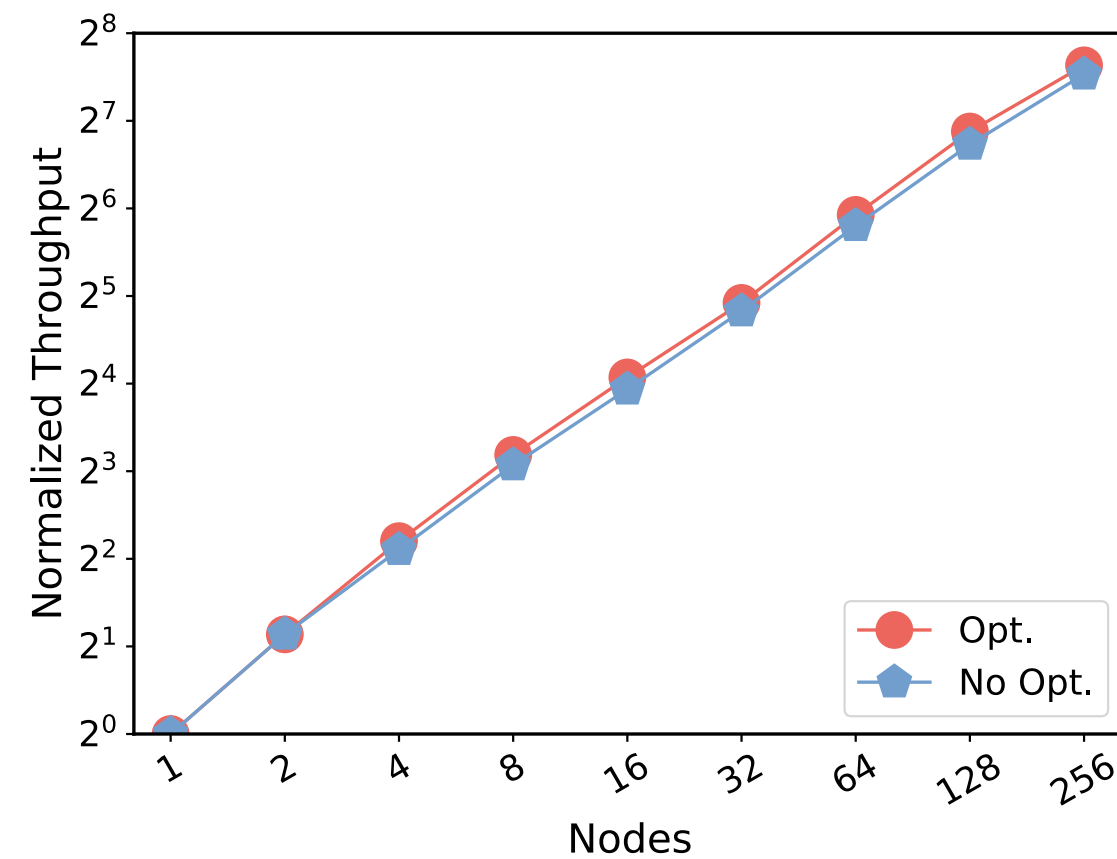
- Balanced slicing uses the implicit knowledge encoded in the application's mapping

# Effect of Idempotent Trace Optimizations
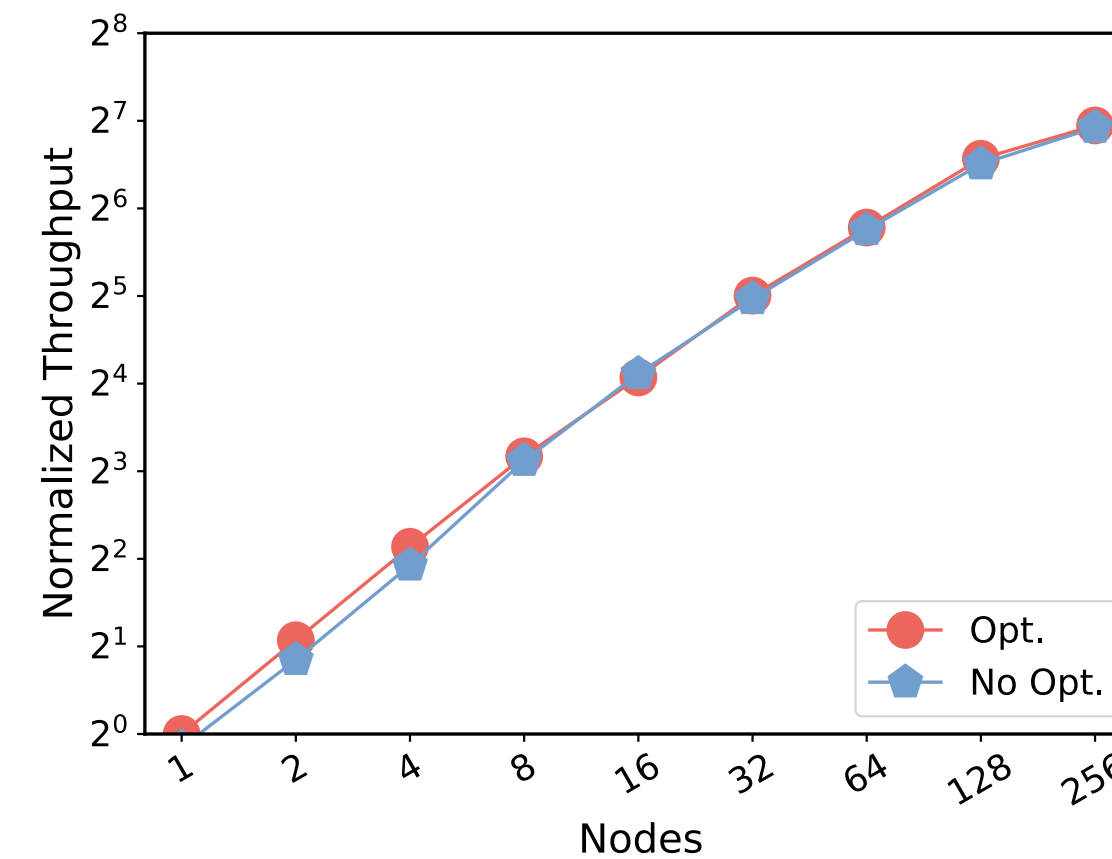


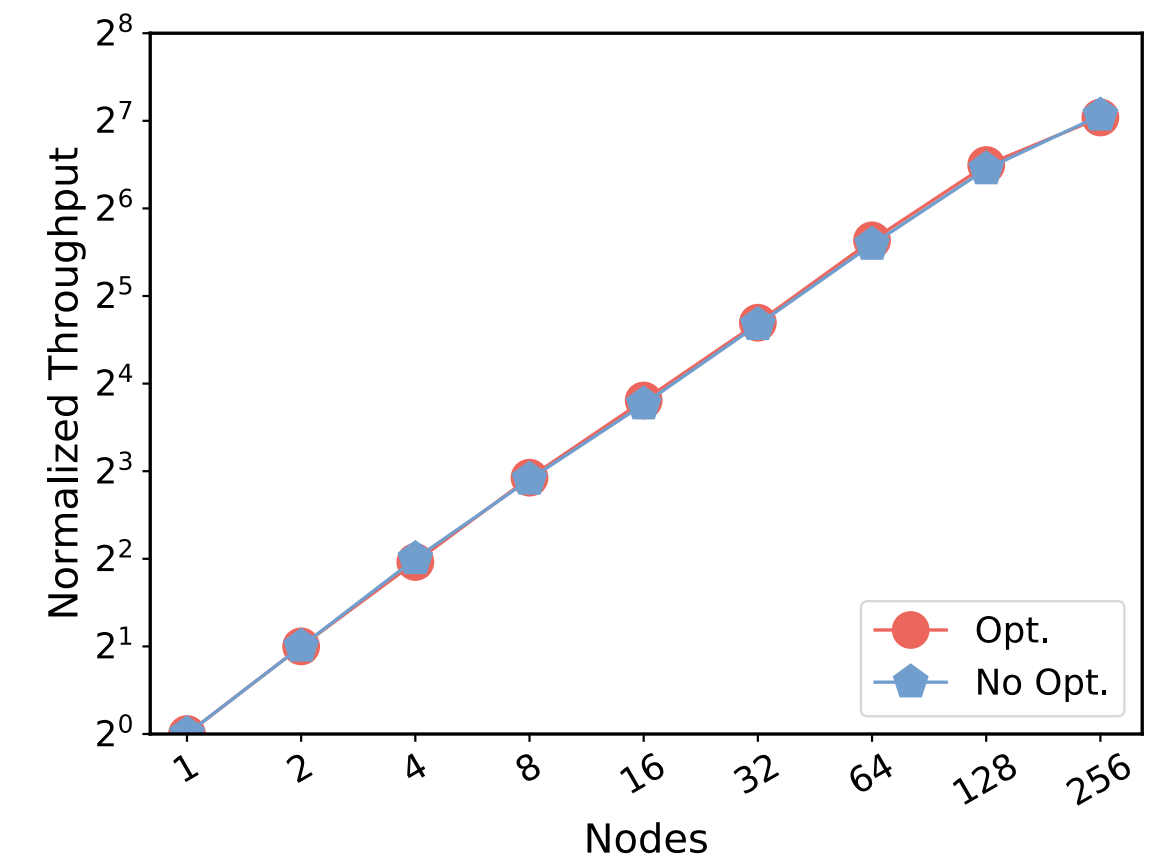Stencil (0.4B Cells, Piz Daint)



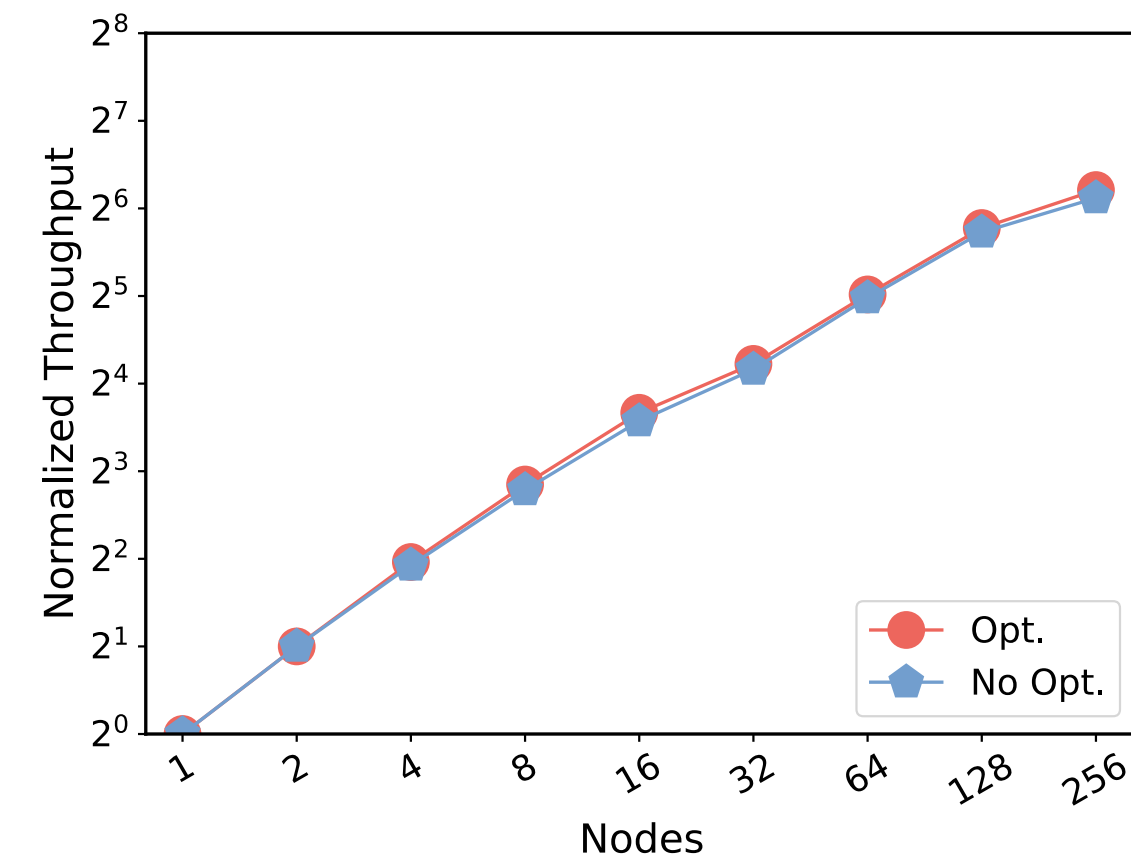PENNANT (29M Zones, Piz Daint)



MiniAero (1M Cells, Piz Daint)



Circuit (74K Wires, Piz Daint)



Soleil-X (8.4M Cells, Piz Daint)

- Idempotent trace optimizations improve performance by an average of 5% and a maximum of 19%

  - Fence elision removes spurious task dependencies, thereby improving performance considerably

  - No benefit on Circuit as it has all-to-all task dependencies on each node
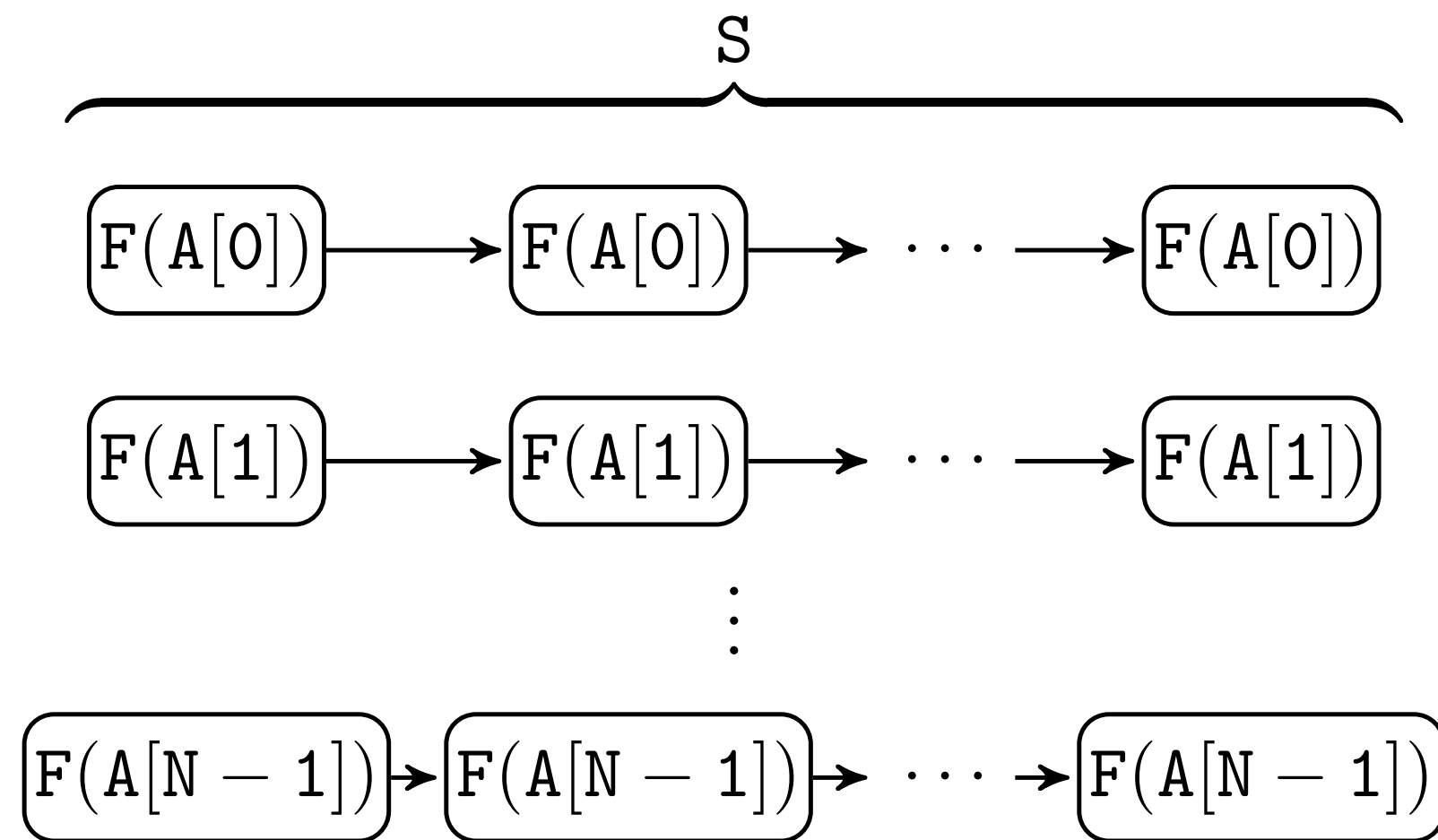
# Average Task Granularity

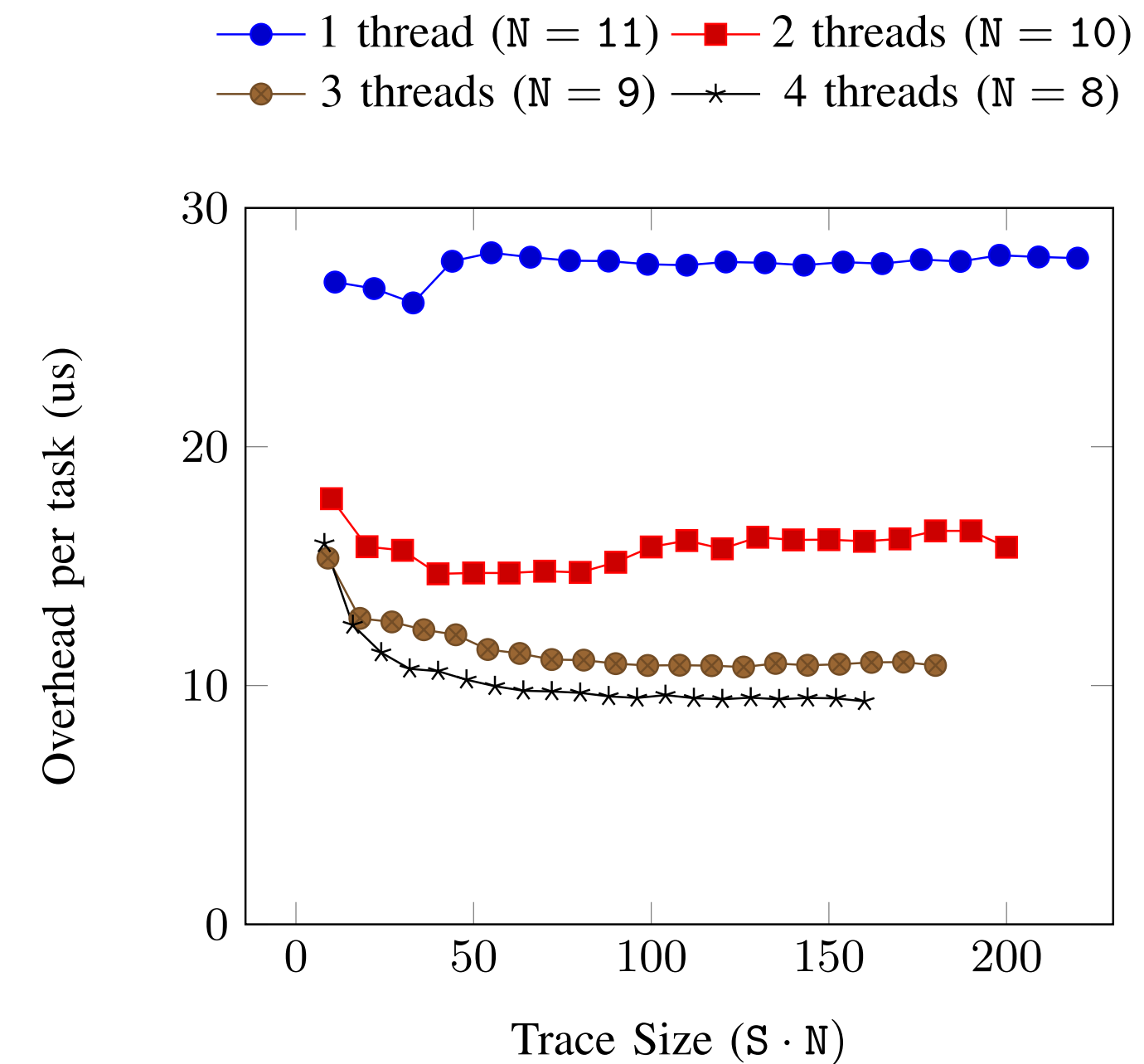| | MiniAero | | Soleil-X | |
|---|---|---|---|---|
| | Tr. | No Tr. | Tr. | No Tr. |
| Num. tasks per processor | 36 | | 56 | |
| Min. time per iteration | 6.6ms | 33.8ms | 23.1ms | 161.2ms |
| Avg. task granularity | 183us | 940us | 413us | 2,879us |

- Achieves sub-millisecond task granularity with dynamic tracing

- Soleil-X tasks take twice more steps on average per replay than MiniAero tasks

Task graph for benchmarking

Trace replay overhead per task



- Using more runtime threads has diminishing return

- Longer traces better amortize the replay overhead

# Tracing Overhead

|              | Stencil | Circuit | PENNANT | MiniAero | Soleil-X |
|--------------|---------|---------|---------|----------|----------|
| No Tracing   | 2.23    | 10.29   | 10.47   | 4.99     | 19.41    |
| Tracing      | 0.29    | 0.53    | 0.86    | 0.68     | 2.26     |
| Improv.      | 7.6×    | 19.5×   | 12.2×   | 7.4×     | 8.6×     |
| Trace size   | 47      | 76      | 121     | 210      | 344      |
| Trace opt.   | 0.72    | 1.70    | 3.90    | 1.75     | 5.86     |

TABLE IV: Runtime overhead per trace (all in milliseconds)