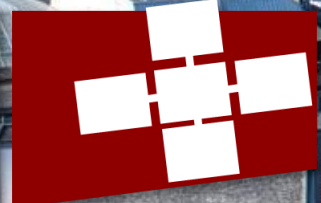
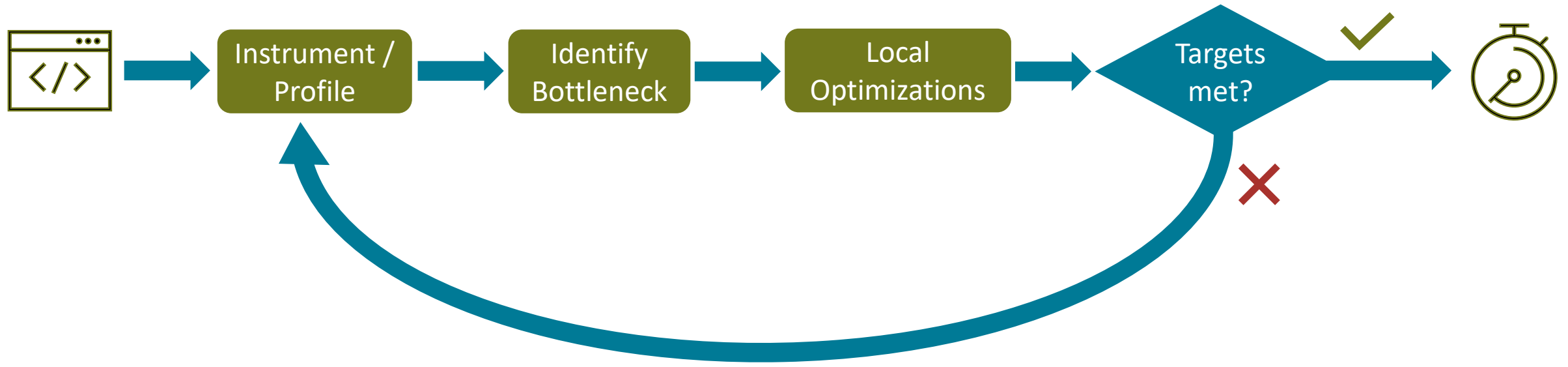


PHILIPP SCHAAD, TAL BEN-NUN, TORSTEN HOEFLER

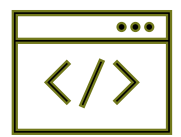
C.A.T.S.: Memory and Control Flow Tracing for Whole-Program Performance Analysis



Typical Optimization Workflows



Typical Optimization Workflows



Instrument / Profile

Identify Bottleneck



The collage shows several performance analysis tools:

- Vampir**: A timeline-based profiler showing OS runtime libraries and process execution.
- AppViewer**: A call stack statistics viewer showing a call stack for a thread.
- Cube GUI**: A 3D visualization tool for performance data, showing a cube with various views (Metric tree, Call tree, Peer distribution).



Typical Optimization Workflows

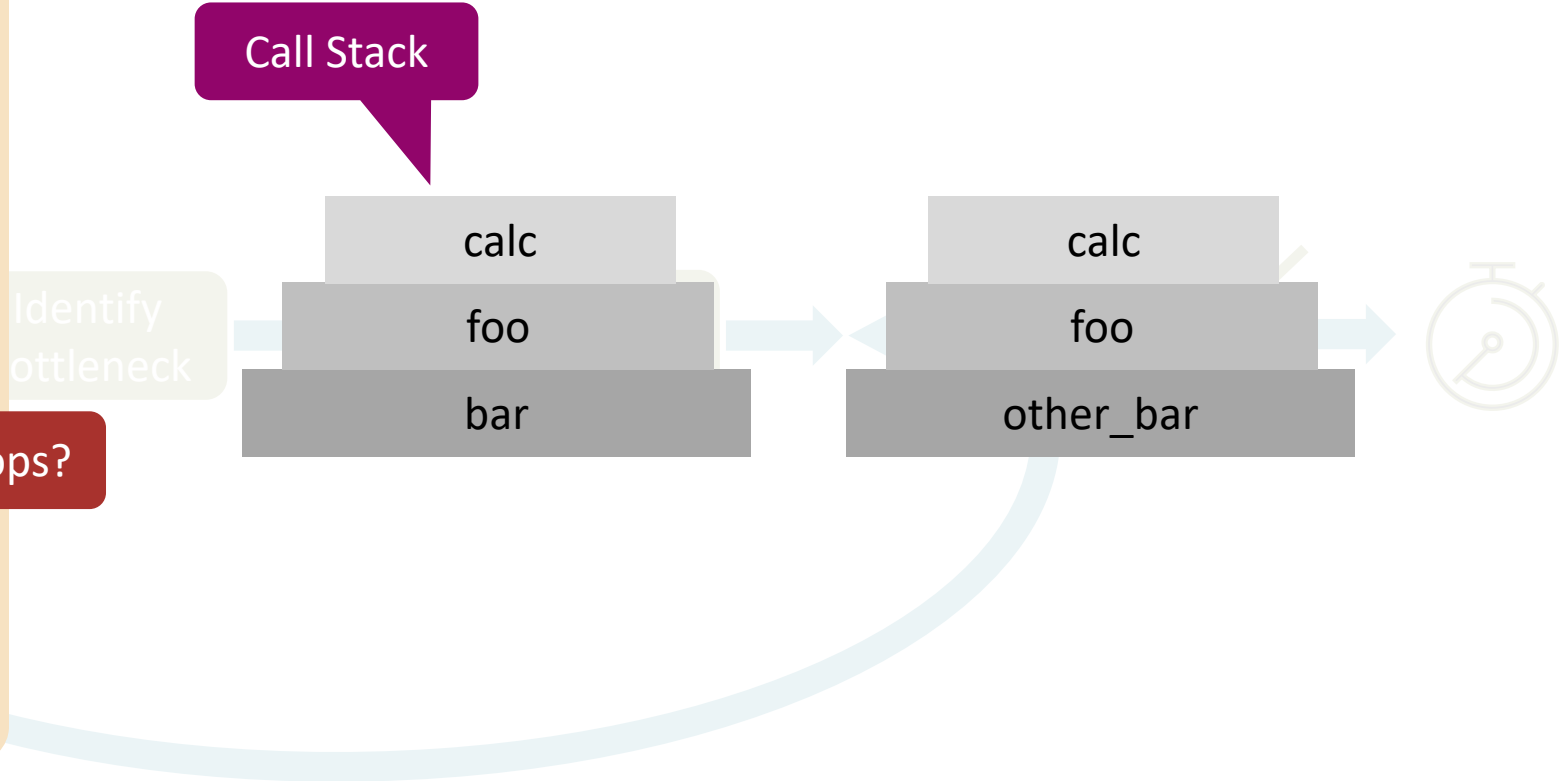
```

void foo() {
  for (int i = 0; i < N; i++) {
    for (int j = 0; j < M; j++) {
      A[i][j] = calc(i, j);
    }
  }
}

void bar() {
  foo();
}

void other_bar() {
  for (int k = 0; k < K; k++) {
    foo();
  }
}
  
```

But what about loops?



Typical Optimization Workflows

```

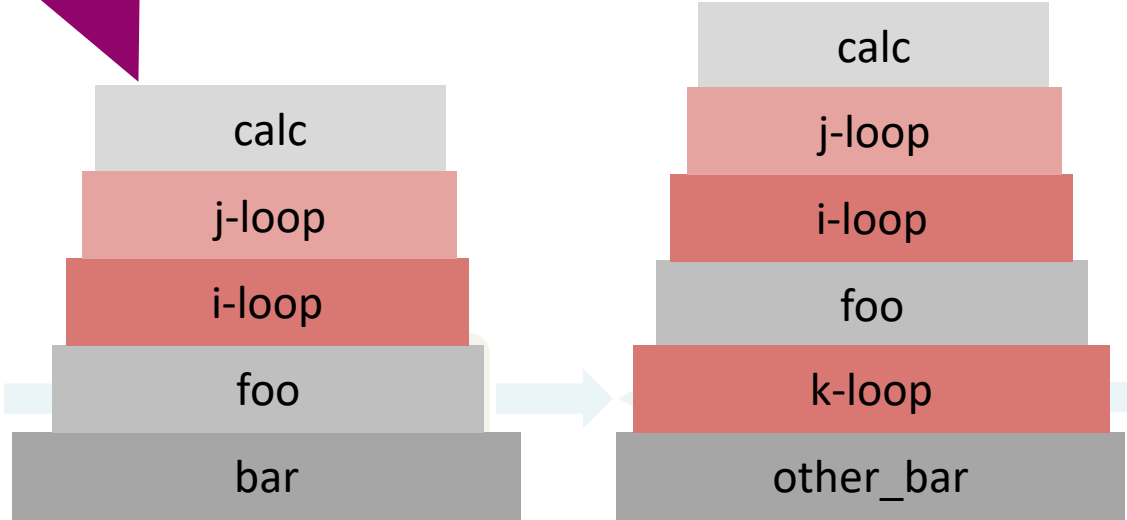
void foo() {
  for (int i = 0; i < N; i++) {
    for (int j = 0; j < M; j++) {
      A[i][j] = calc(i, j);
    }
  }
}

void bar() {
  foo();
}

void other_bar() {
  for (int k = 0; k < K; k++) {
    foo();
  }
}
  
```

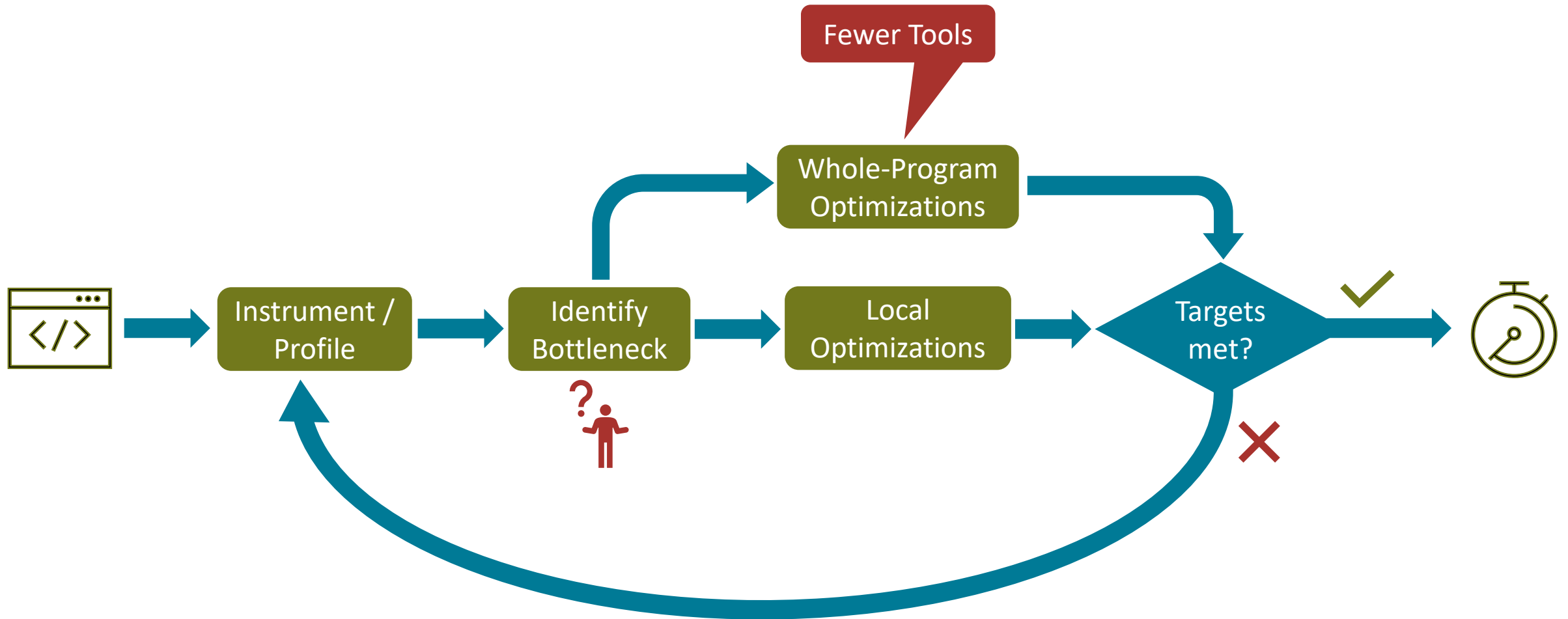
Control Flow Stack

Identify bottleneck



C.A.T.S.: Control Augmented Trace Structure

Typical Optimization Workflows



Typical Optimization Workflows

3 Whole-Program Performance Visualizations

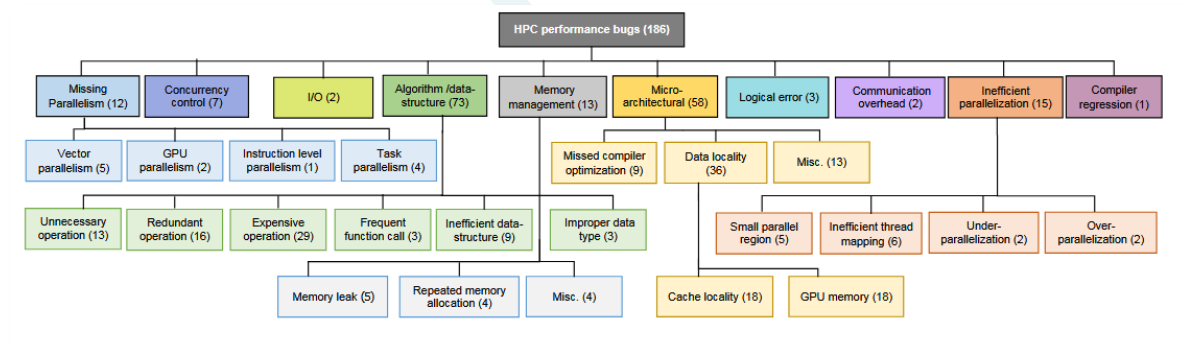
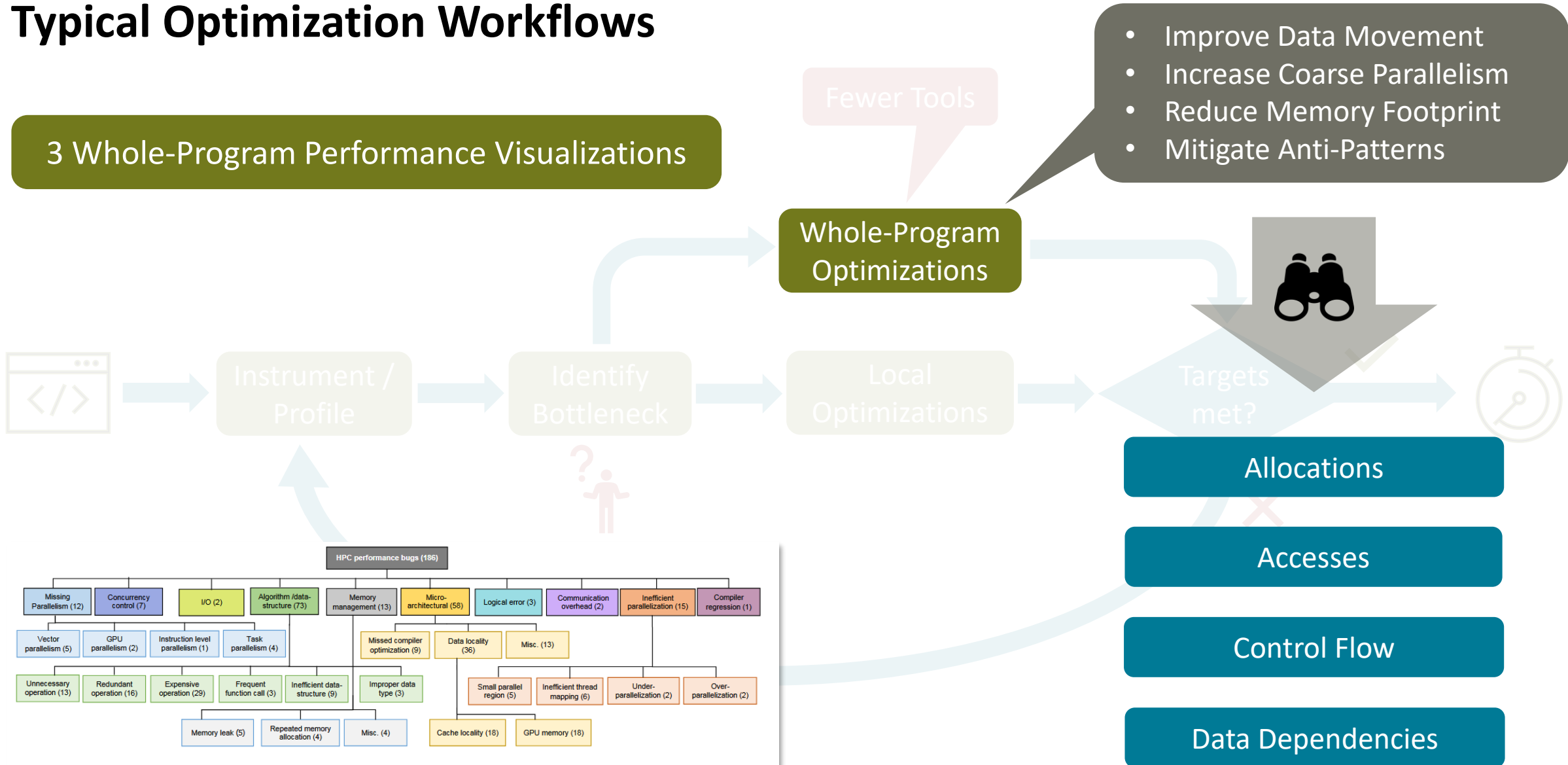
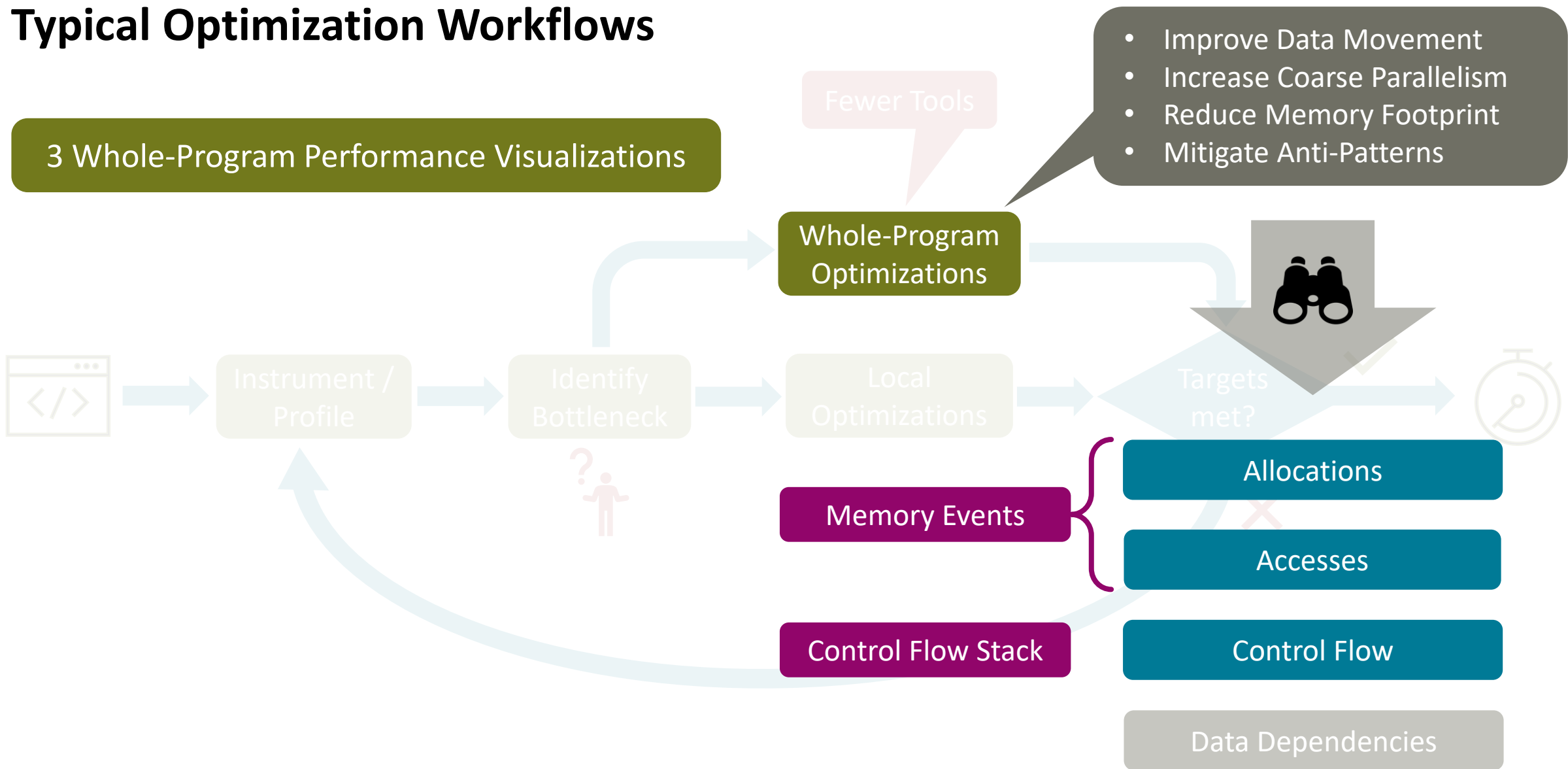


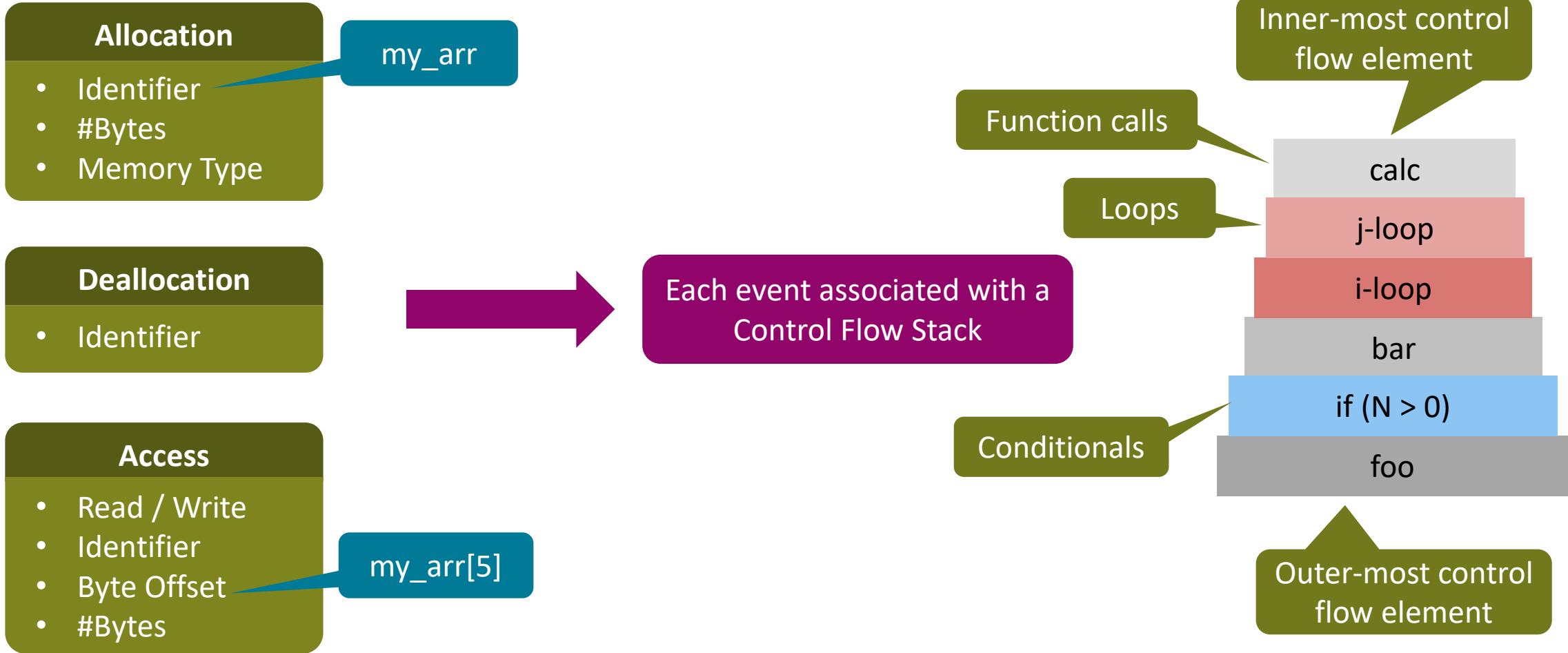
Fig. 2. Taxonomy of HPC Performance Bugs [1]

[1] Azad et al., An Empirical Study of High Performance Computing (HPC) Performance Bugs

Typical Optimization Workflows







C.A.T.S.: Control Augmented Trace Structure



C.A.T.S.: Trace Construction




Runtime Tracing

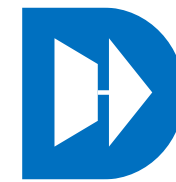
-  Observes Reality
-  Data-Dependent Execution
-  High Capture Time
-  Memory Overhead



All LLVM compatible
source languages

Static Tracing

-  Low Capture Time
-  Captures Symbolic Information
-  Infeasible for Data-Dependent Execution



C.A.T.S.: Trace Compression

```

void foo() {
  for (int i = 0; i < N; i++) {
    for (int j = 0; j < M; j++) {
      A[i][j] = calc(i, j);
    }
  }
}

```

Control flow stack captures fact that this is in a loop nest

Only really care about the event once per loop nest

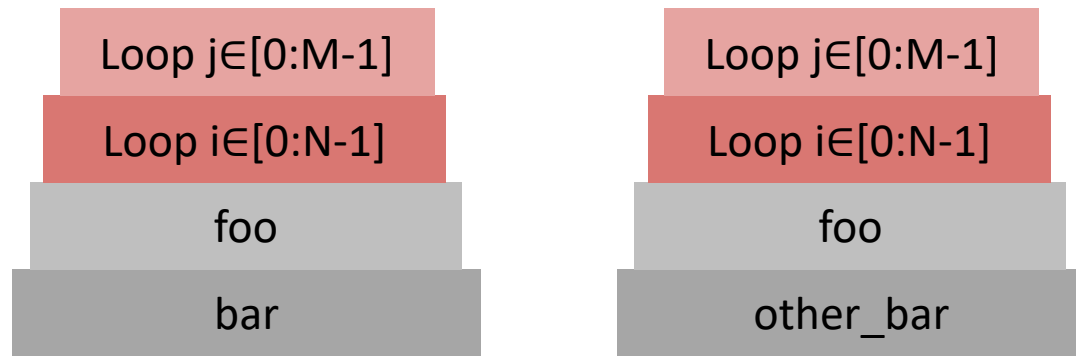
Instrumenting this access causes $2 * N * M$ events

```

void bar() {
  foo();
}

void other_bar() {
  foo();
}

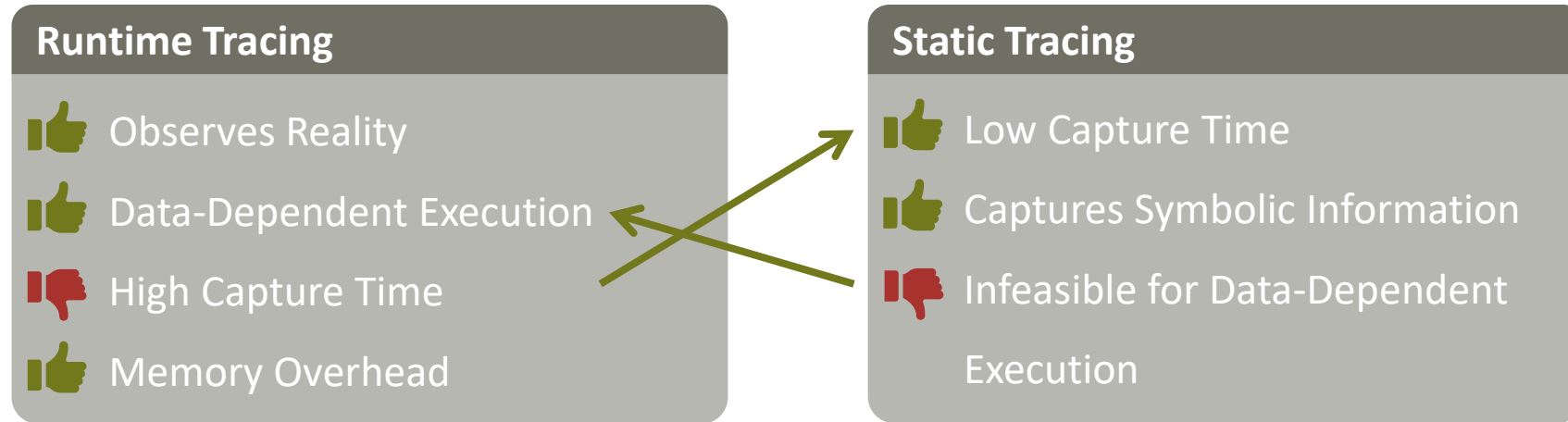
```



→ Events only recorded once per *unique* control flow stack

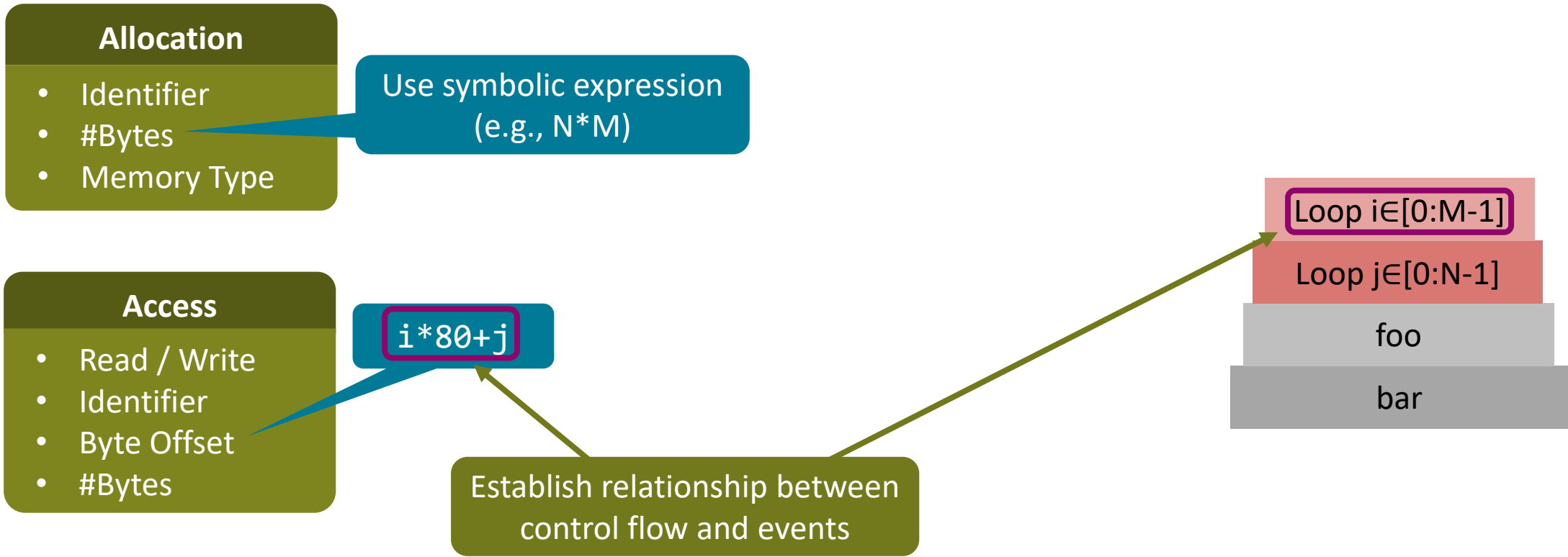
Only 2 unique stacks → 2 events

C.A.T.S.: Implementations



C.A.T.S.: Static Tracing

Tracing via symbolic execution gives symbolic information



C.A.T.S.: Novel Performance Visualizations

Reduce Data Movement

Increase Coarse Parallelism

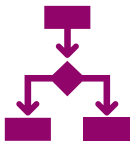
Reduce Memory Footprint

Mitigate Anti-Patterns

Control Flow Data Dependency View

Memory Timeline View

Data Access Statistics View



C.A.T.S.: Novel Performance Visualizations

Reduce Data Movement

Increase Coarse Parallelism

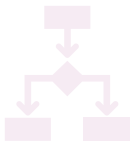
Reduce Memory Footprint

Mitigate Anti-Patterns

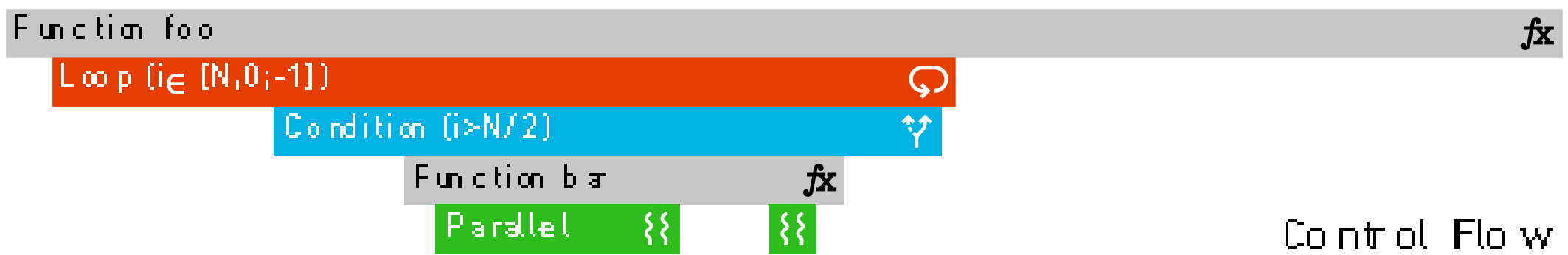
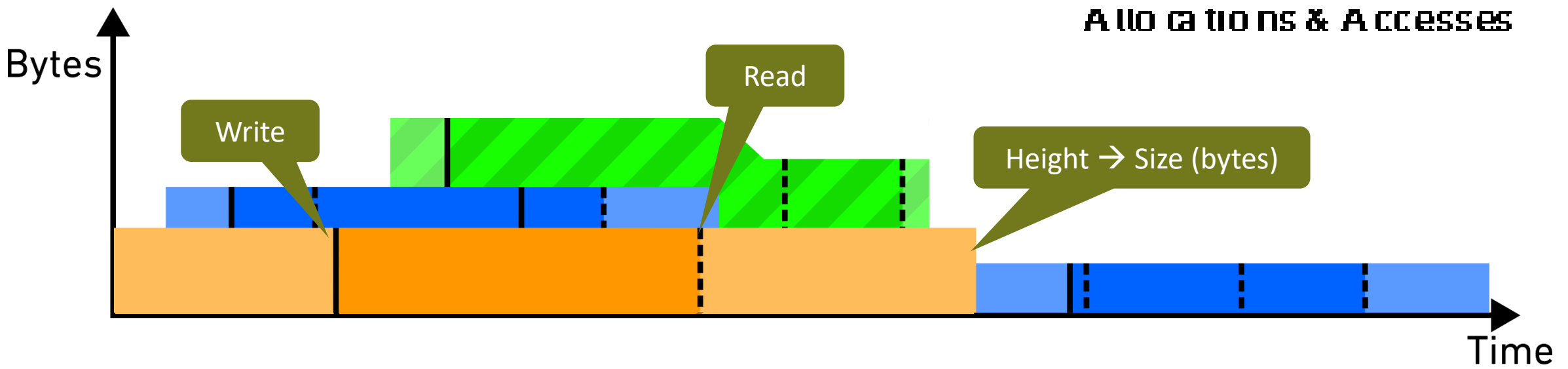
Control Flow Data Dependency View

Memory Timeline View

Data Access Statistics View



Memory Timeline View

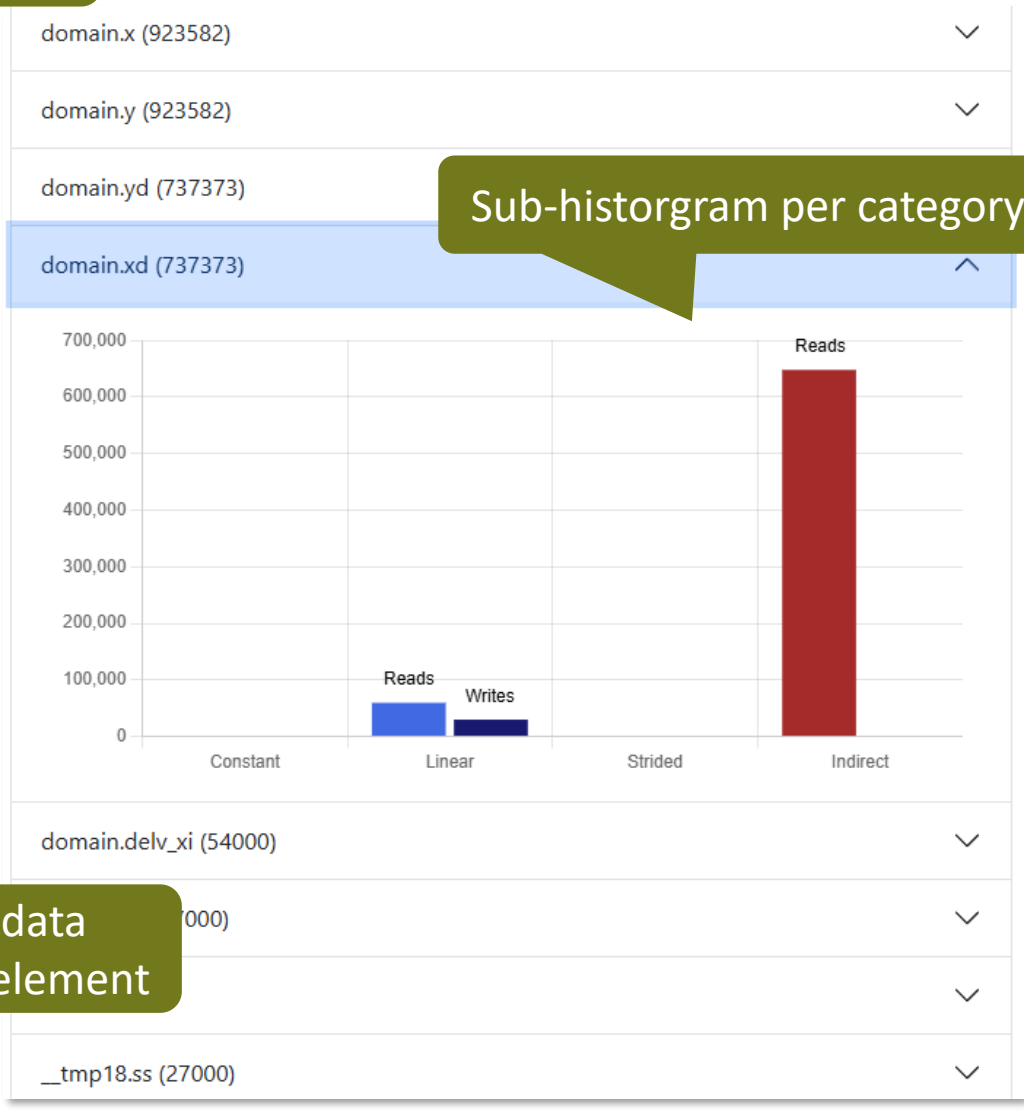


Control Flow Stack

Access Statistics View



Histogram of data accesses by type

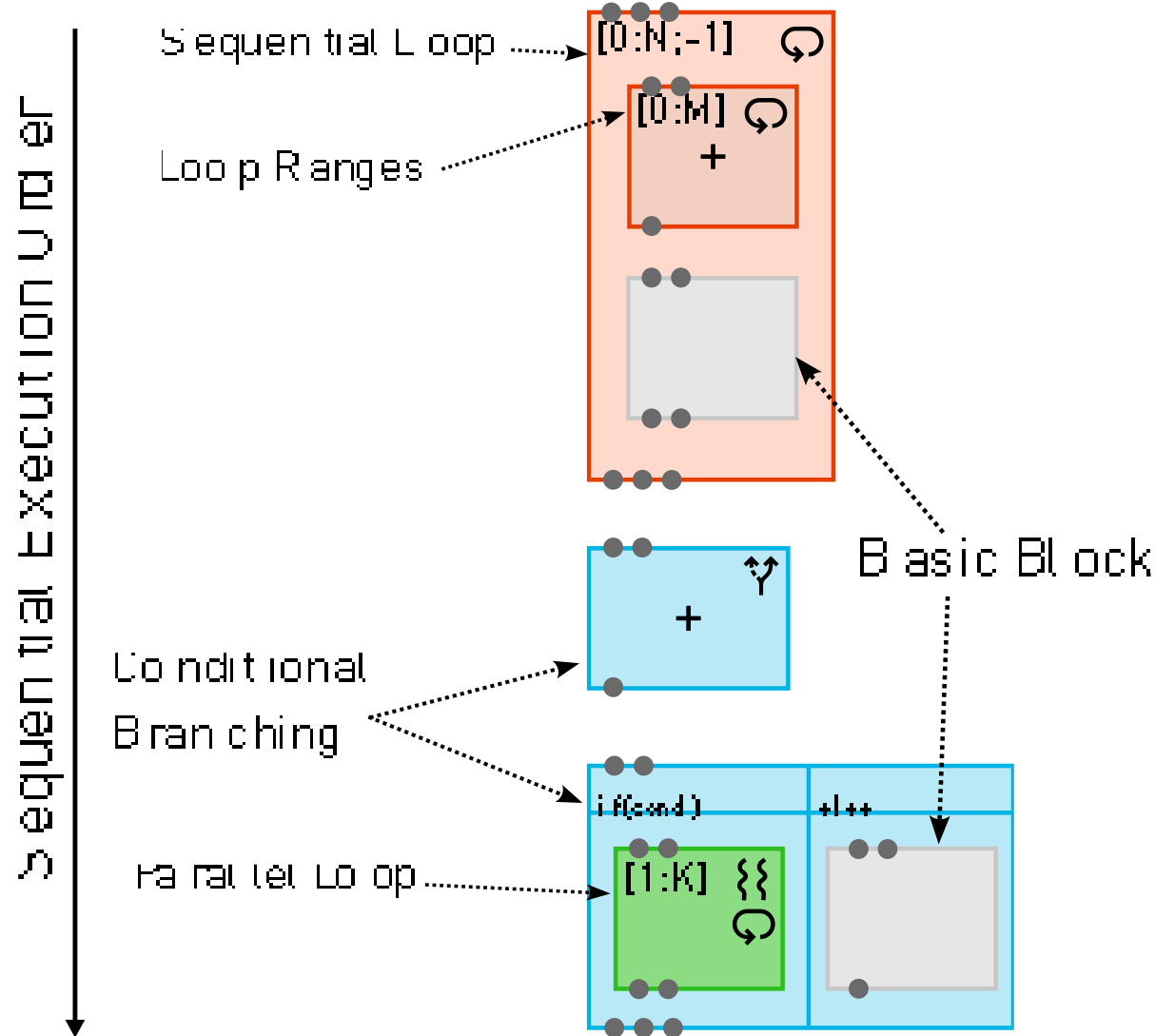
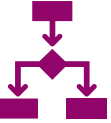


Sub-histogram per category

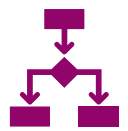
Breakdown based on data container / control flow element



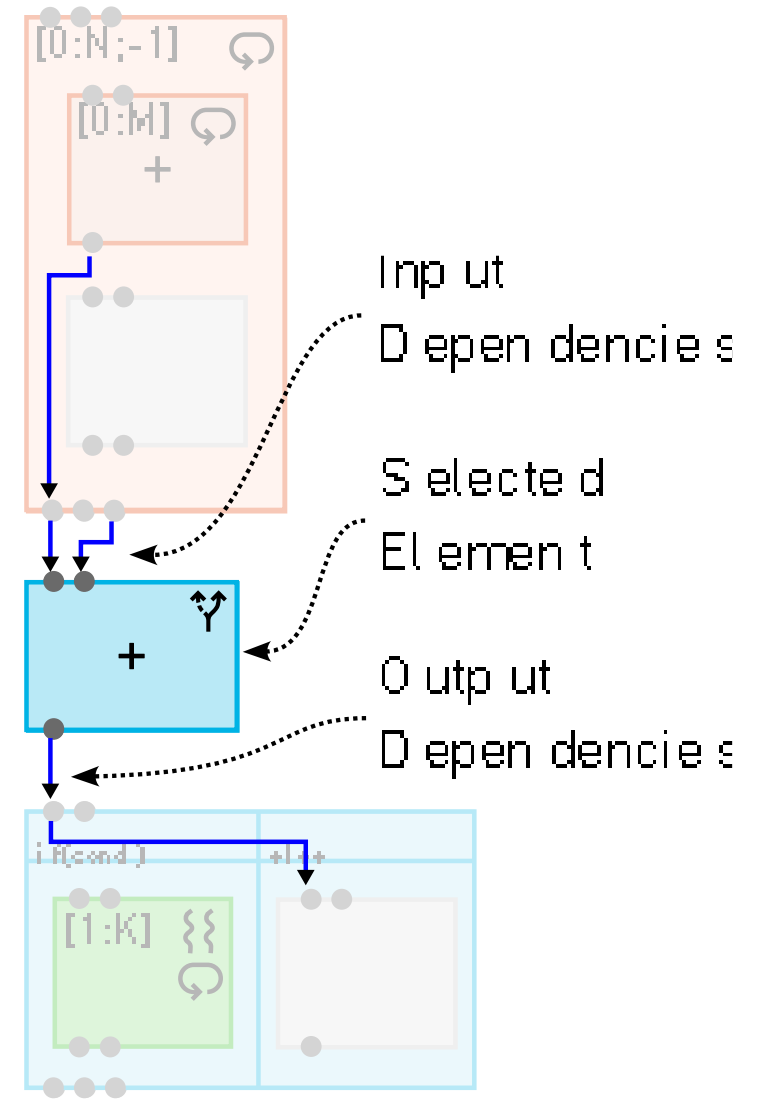
Control Flow Data Dependency View



Control Flow Data Dependency View



Sequential Execution Order



Shock Hydrodynamics – LULESH

Valgrind: 41 minutes / 11.8 MB trace



Static Tracing



24 seconds / 395 KB trace

30x smaller trace

Shock Hydrodynamics – LULESH

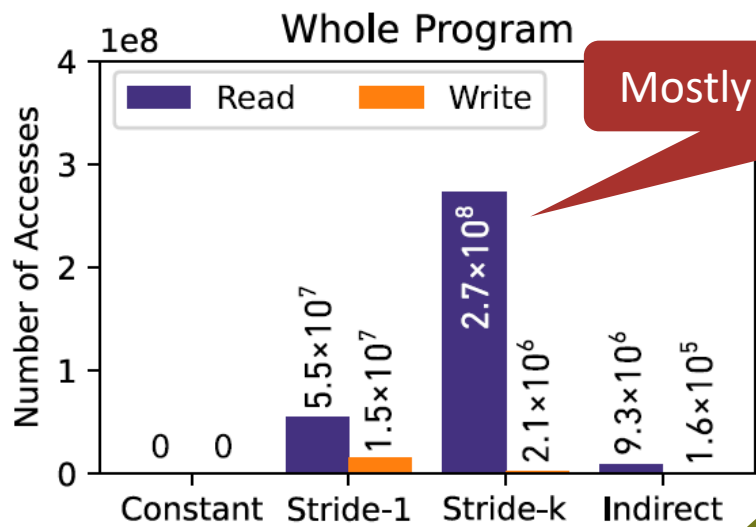


Static Tracing



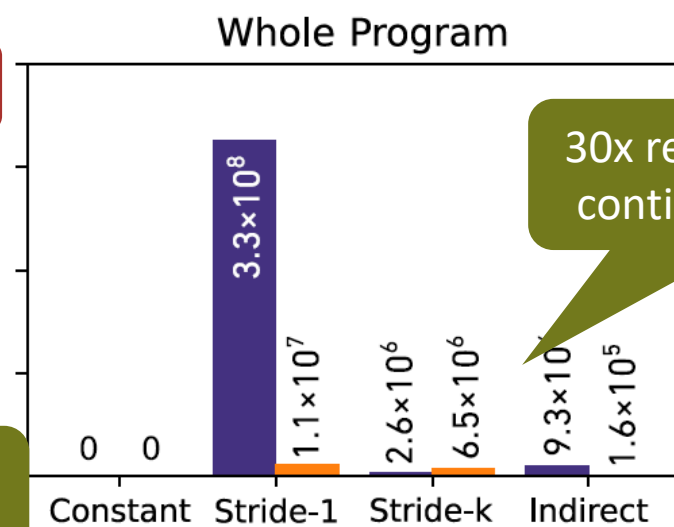
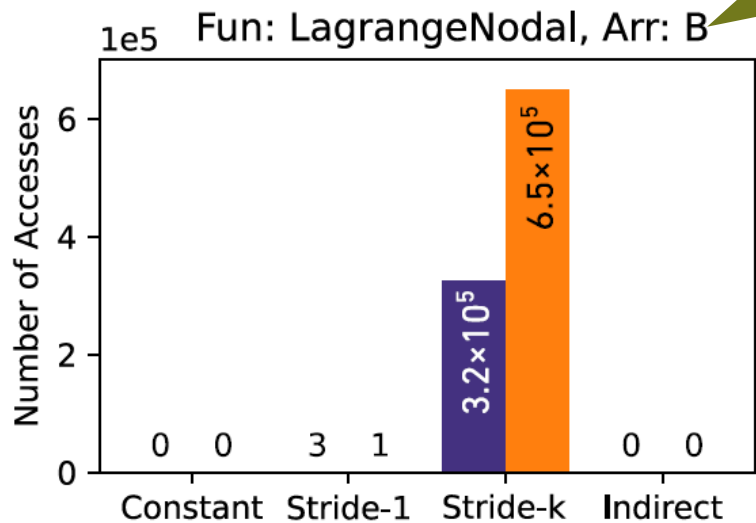
24 seconds / 395 KB trace

1.18x single-node speedup



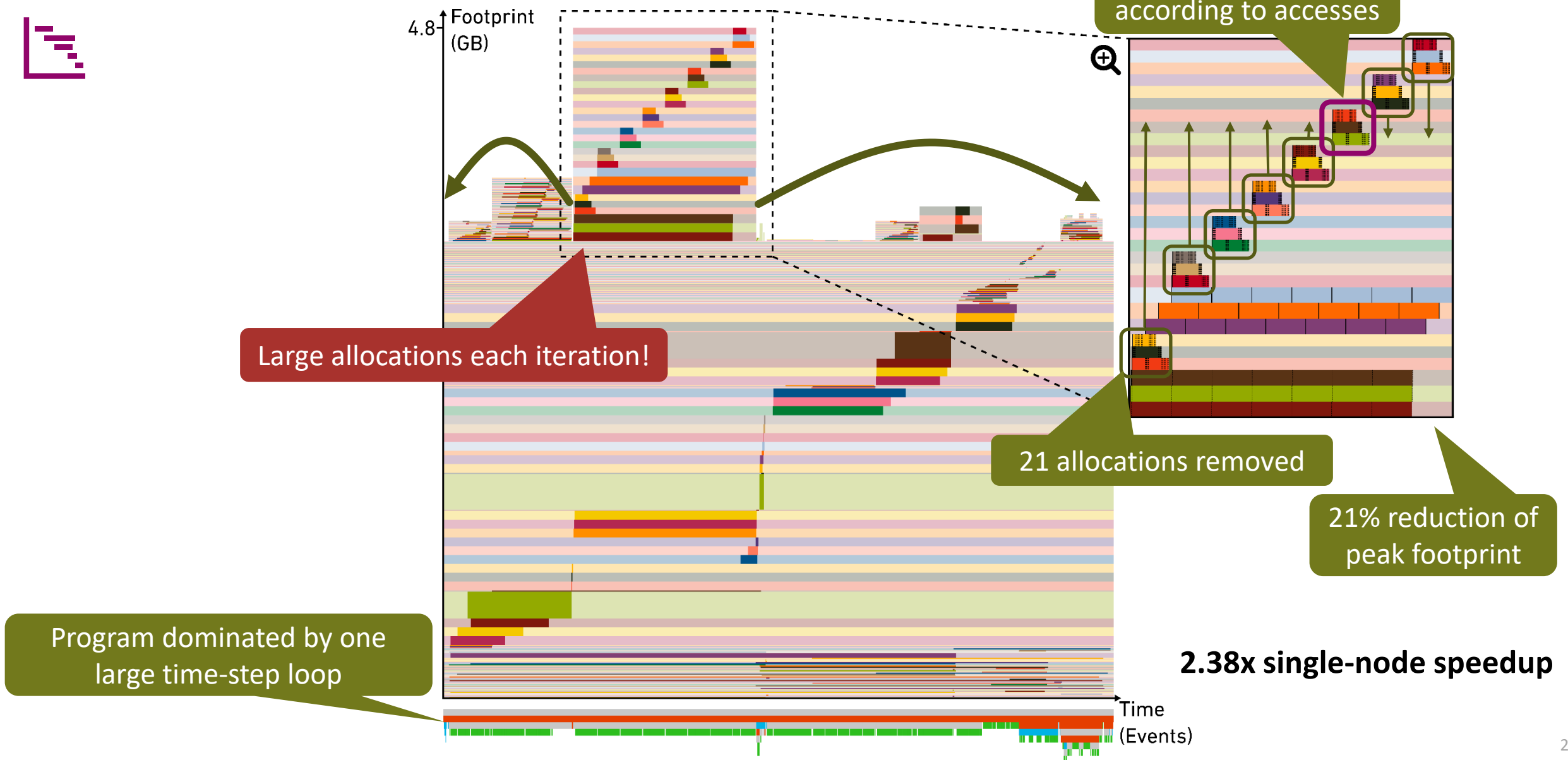
Mostly non-contiguous accesses

Container with most non-contiguous accesses



30x reduction in non-contiguous accesses

Shock Hydrodynamics – LULESH



Large allocations each iteration!

Reuse is possible according to accesses

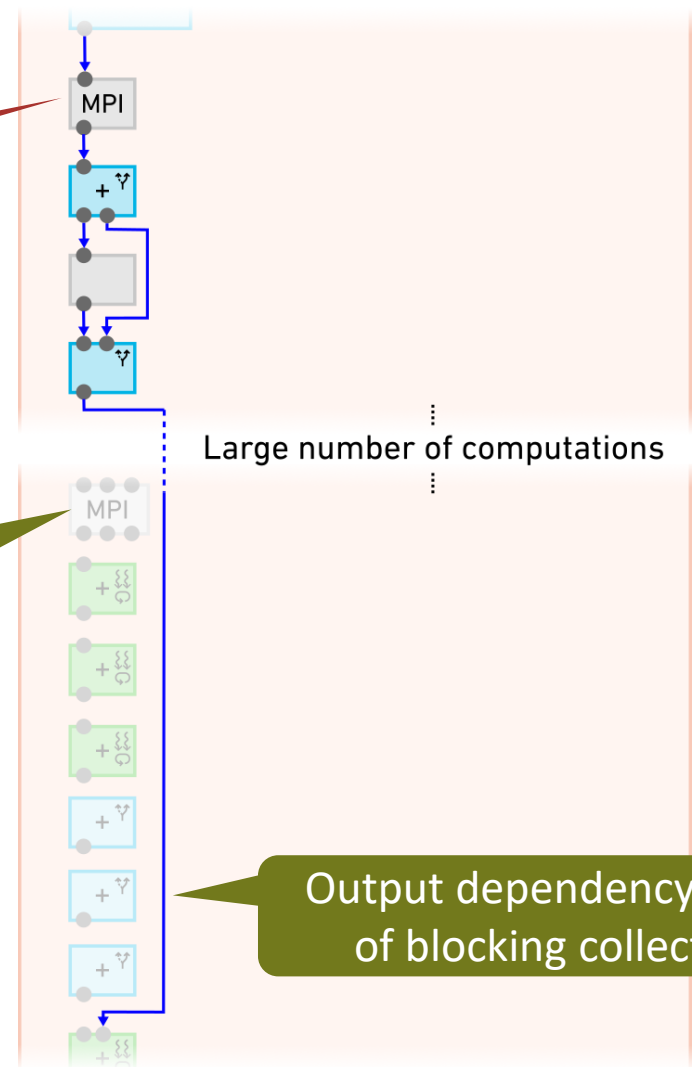
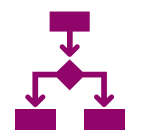
21 allocations removed

21% reduction of peak footprint

Program dominated by one large time-step loop

2.38x single-node speedup

Shock Hydrodynamics – LULESH



Blocking MPI_Allreduce

Use non-blocking version!

Next blocking MPI synchronization point

Large number of computations

Output dependency chain of blocking collective

Time-step loop



256 nodes with 4 GH200 each, 2197 ranks

5% speedup / 10% communication overhead reduction

Kilometer-Scale Weather Modeling

Finite-Volume Cubed-Sphere Dynamical Core (FV3)
used by NASA and NOAA



Static Tracing

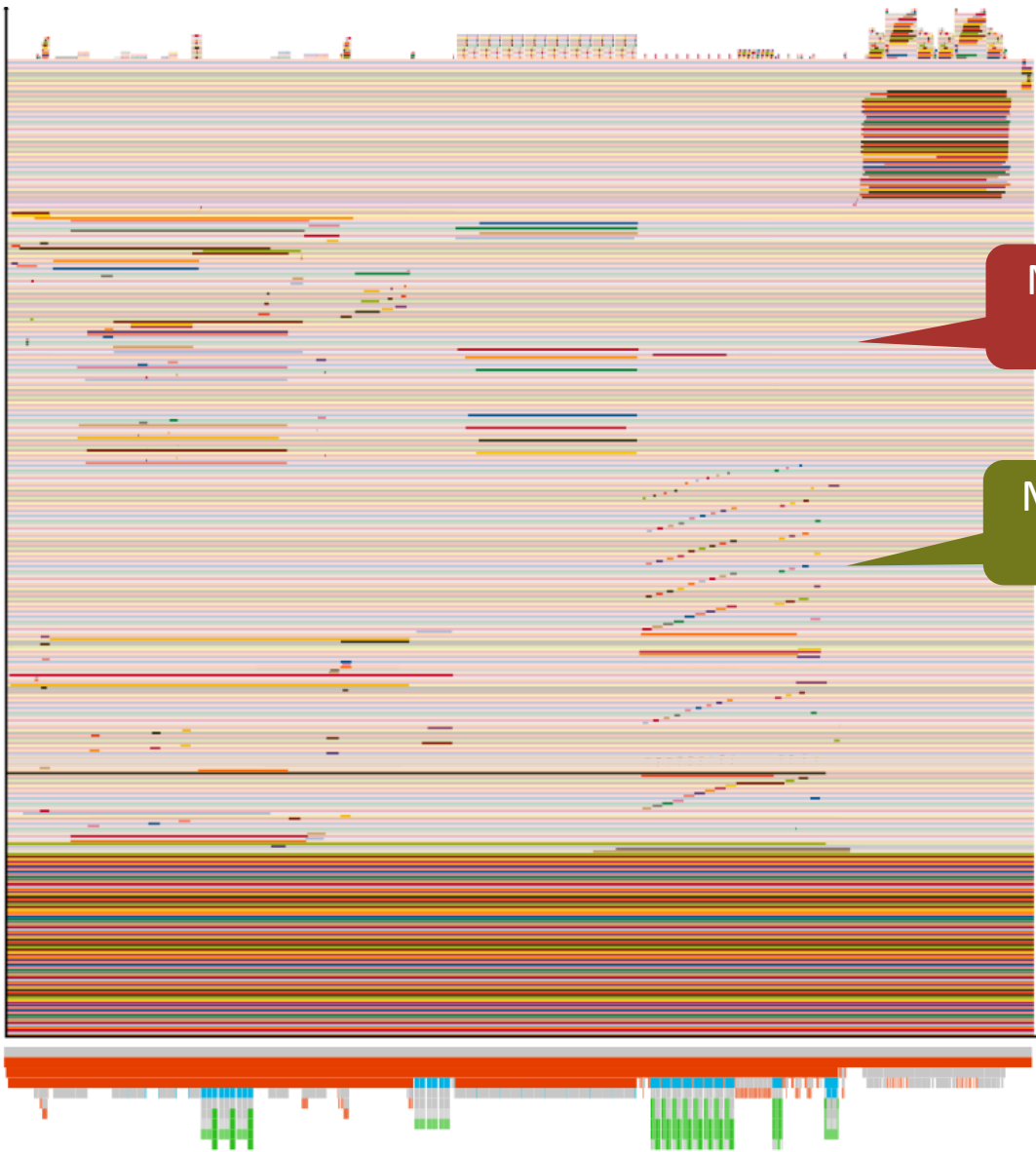


65 seconds / 1.2 MB trace

Memory footprint on GPU
restricts KM-scale modeling

More than 3100 H100 GPUs

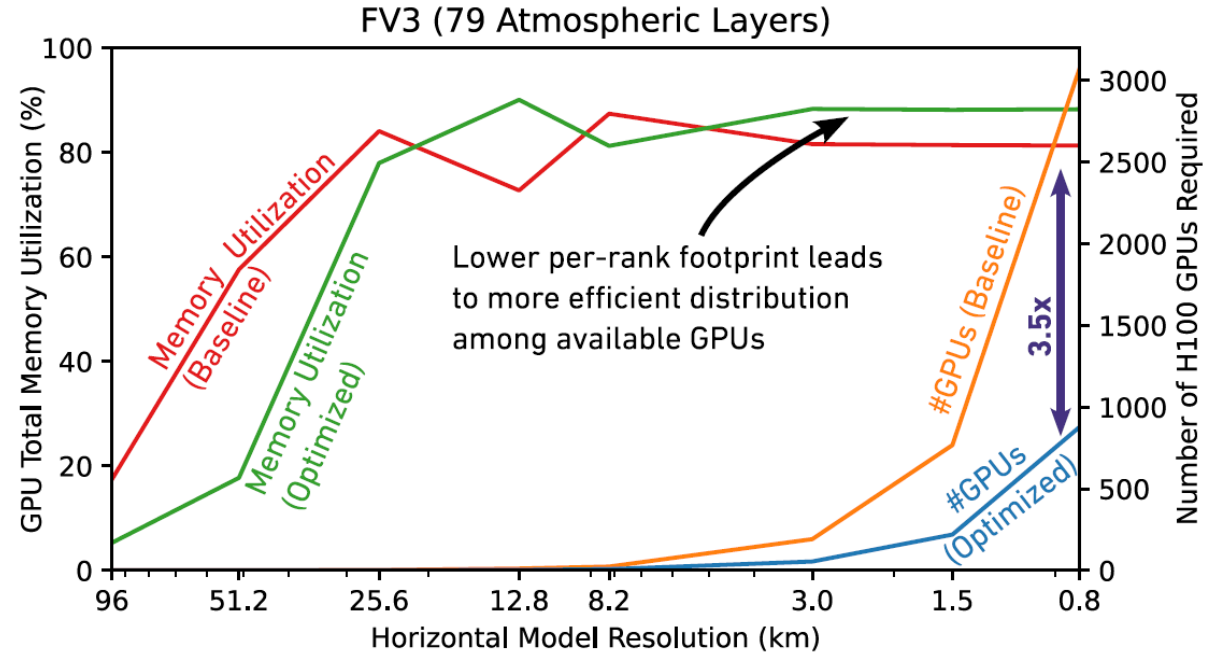
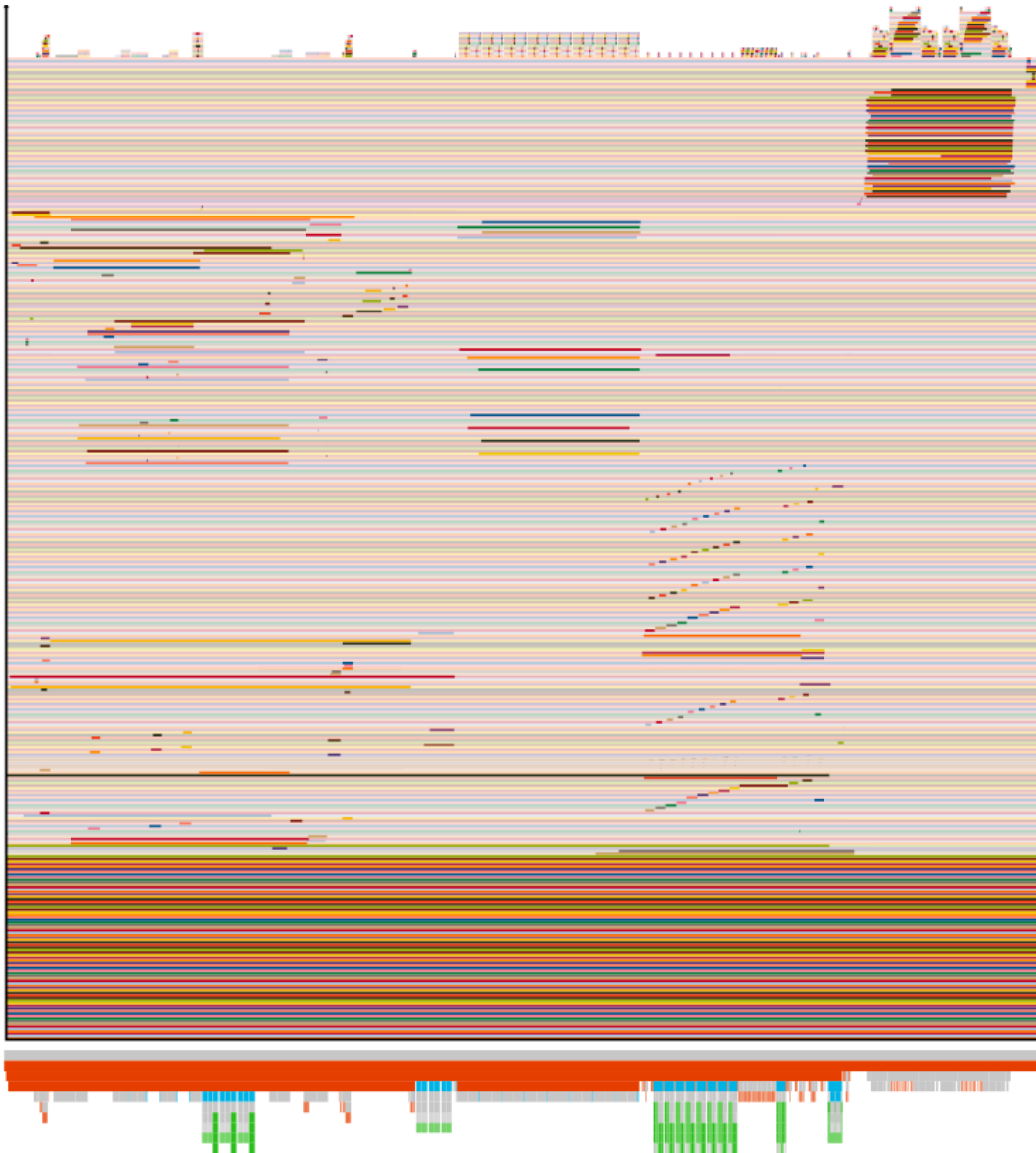
Kilometer-Scale Weather Modeling



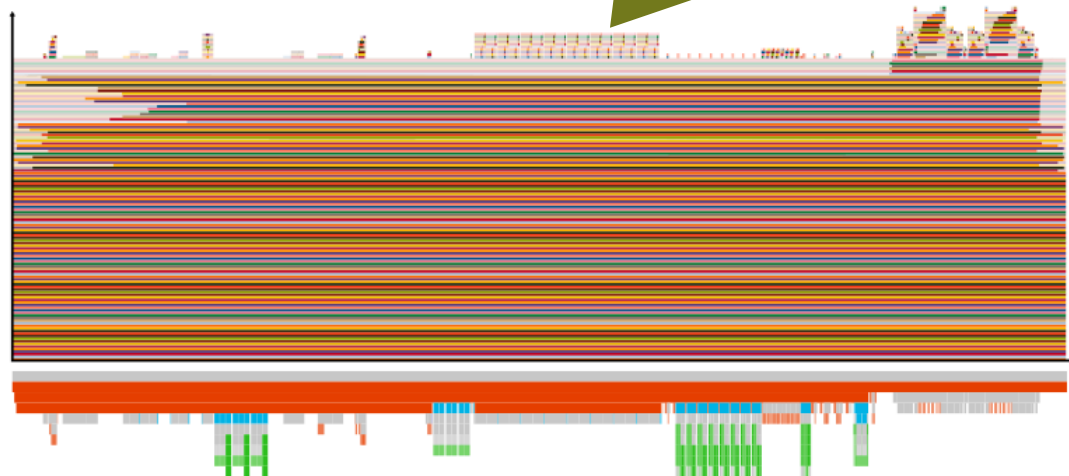
Mostly unused GPU memory

Most temporaries can be reused!

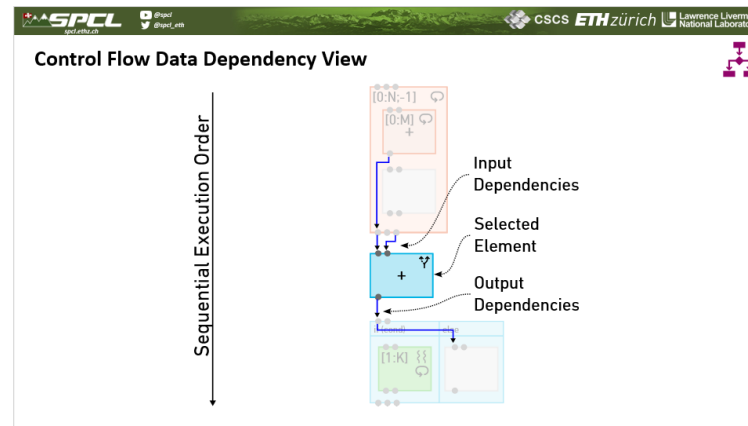
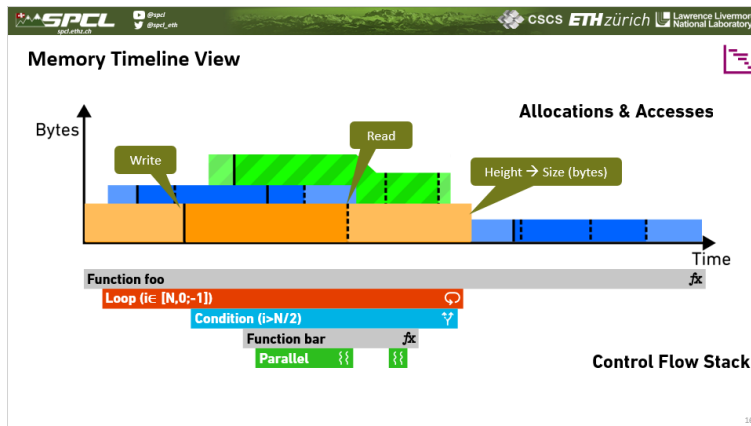
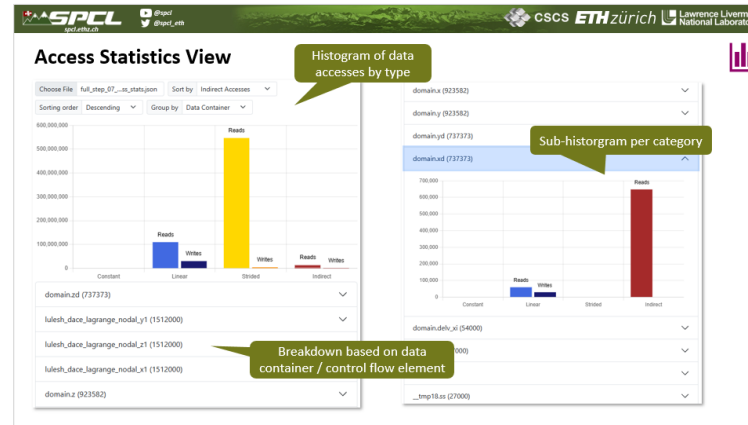
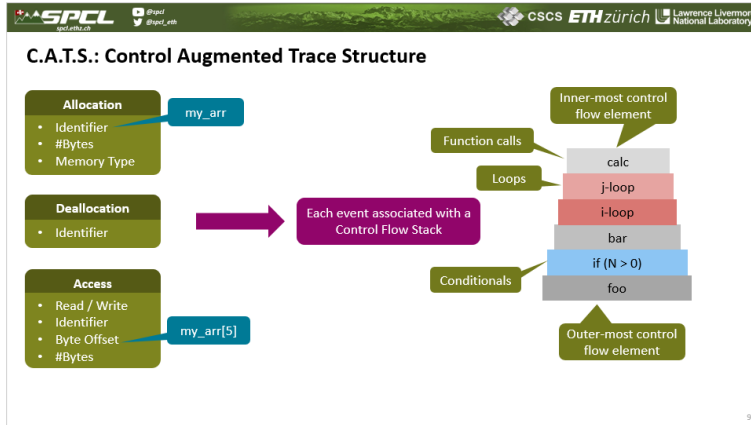
Kilometer-Scale Weather Modeling



More than 3x footprint reduction




Conclusions



More of SPCL's research:

 youtube.com/@spcl **240+ Talks**

 x.com/spcl_eth **1.7K+ Followers**

 github.com/spcl **7.2K+ Stars**




... or spcl.ethz.ch



Conclusions



More of SPCL's research:

-  youtube.com/@spcl **240+ Talks**
-  x.com/spcl_eth **1.7K+ Followers**
-  github.com/spcl **7.2K+ Stars**

... or spcl.ethz.ch



Data-Dependent Program: HPCG

