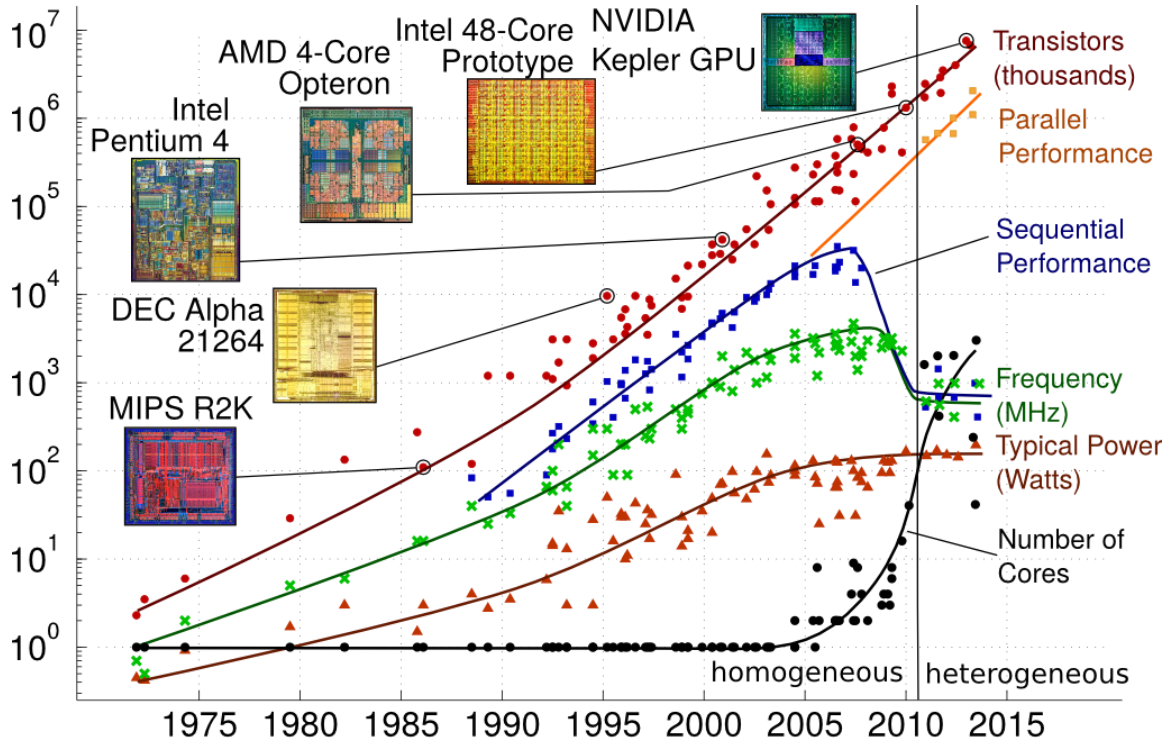# Polly-ACC: Transparent Compilation to Heterogeneous Hardware
## Torsten Hoefler (with Tobias Grosser)

# Evading various "ends" – the hardware view



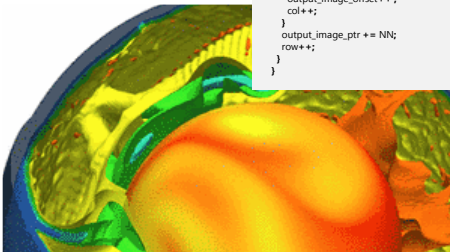Data partially collected by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond
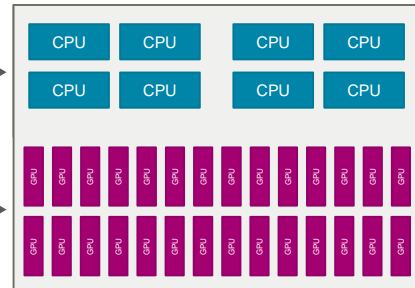
# Sequential Software

# Parallel Hardware

## Multi-Core CPU

## Accelerator

# Design Goals



Automatic

**Automatic accelerator mapping
-
How close can we get?**

"Regression Free"    High Performance

**ETH**zürich

# Tool: Polyhedral Modeling

*Program Code*

*Iteration Space*

```
for (i = 0; i <= N; i++)
  for (j = 0; j <= i; j++)
    S(i,j);
```

N = 4

(i, j) = (4,4)

Polly -- Performing Polyhedral
Optimizations on a Low-Level
Intermediate Representation
Tobias Grosser  et al,
Parallel Processing Letter, 2012

i ≤ N = 4

0 ≤ j

j ≤ i

0 ≤ i

D = { (i,j) | 0 ≤ i ≤ N ∧ 0 ≤ j ≤ i }

# Mapping Computation to Device

*Iteration Space*

*Device Blocks & Threads*



$$BID = \{(i,j) \rightarrow \left( \left\lfloor \frac{i}{4} \right\rfloor \% \ 2, \left\lfloor \frac{j}{3} \right\rfloor \% \ 2 \right)\}$$

$$TID = \{(i,j) \rightarrow (i \% \ 4, j \% \ 3)\}$$

# Memory Hierarchy of a Heterogeneous System



Main Memory

Device Memory

Shared Memory

Registers

CPU  CPU

CPU  CPU

GPU GPU  GPU GPU  GPU GPU

GPU GPU  GPU GPU  GPU GPU

GPU GPU  GPU GPU  GPU GPU

GPU GPU  GPU GPU  GPU GPU

# Host-device date transfers



Main Memory

Device Memory

Shared Memory

Registers

# Host-device date transfers



Main Memory

Device Memory

Shared Memory

Registers

# Mapping onto fast memory



Main Memory

Device Memory

Shared Memory

Registers

CPU  CPU

CPU  CPU

GPU GPU  GPU GPU  GPU GPU
GPU GPU  GPU GPU  GPU GPU
GPU GPU  GPU GPU  GPU GPU
GPU GPU  GPU GPU  GPU GPU

# Mapping onto fast memory



Polyhedral parallel code generation for CUDA, Verdoolaege, Sven et. al, ACM Transactions on Architecture and Code Optimization, 2013

# Profitability Heuristic

Modeling Execution GPU

All
Loop
Nests

static dynamic

Insufficient Compute

Unsuitable

Trivial

T. Grosser, TH: Polly-ACC: Transparent compilation to heterogeneous hardware, ACM ICS'16

# From kernels to program – data transfers

```
void heat(int n, float A[n], float hot, float cold) {

    float B[n] = {0};

    initialize(n, A, cold);
    setCenter(n, A, hot, n/4);

    for (int t = 0; t < T; t++) {
        average(n, A, B);
        average(n, B, A);
        printf("Iteration %d done", t);
    }
}
```

# Data Transfer – Per Kernel

```
void heat(int n, float A[n], ...) {
    initialize(n, A, cold);
    setCenter(n, A, hot, n/4);
    for (int t = 0; t < T; t++) {
        average(n, A, B);
        average(n, B, A);
        printf("Iteration %d done", t);
    } }
```

Host Memory                     Device Memory

initialize()        $D \rightarrow H$



setCenter()                     $D \rightarrow H$



average()        $H \rightarrow D \ D \rightarrow H$



average()                     $H \rightarrow D \ \ D \rightarrow H$



average()        $H \rightarrow D \ D \rightarrow H$



time

# Data Transfer – Inter Kernel Caching

```
void heat(int n, float A[n], ...) {
  initialize(n, A, cold);
  setCenter(n, A, hot, n/4);
  for (int t = 0; t < T; t++) {
    average(n, A, B);
    average(n, B, A);
    printf("Iteration %d done", t);
  } }
```
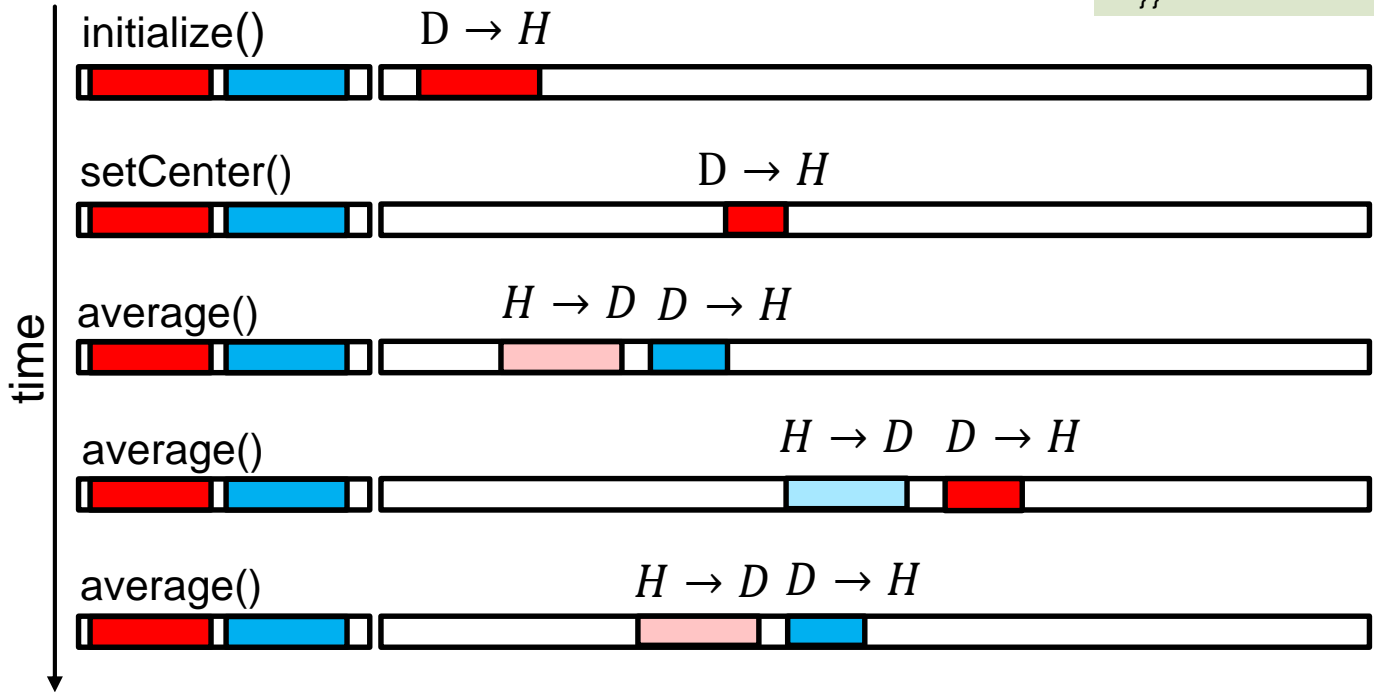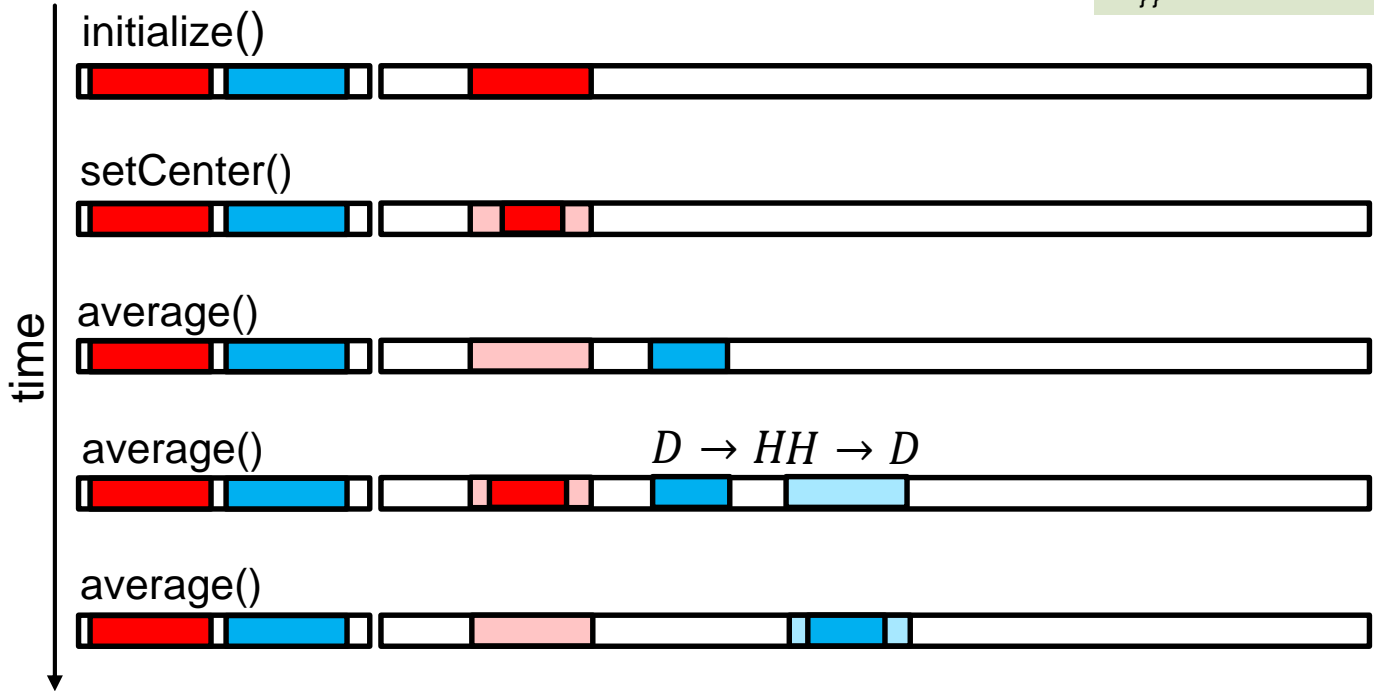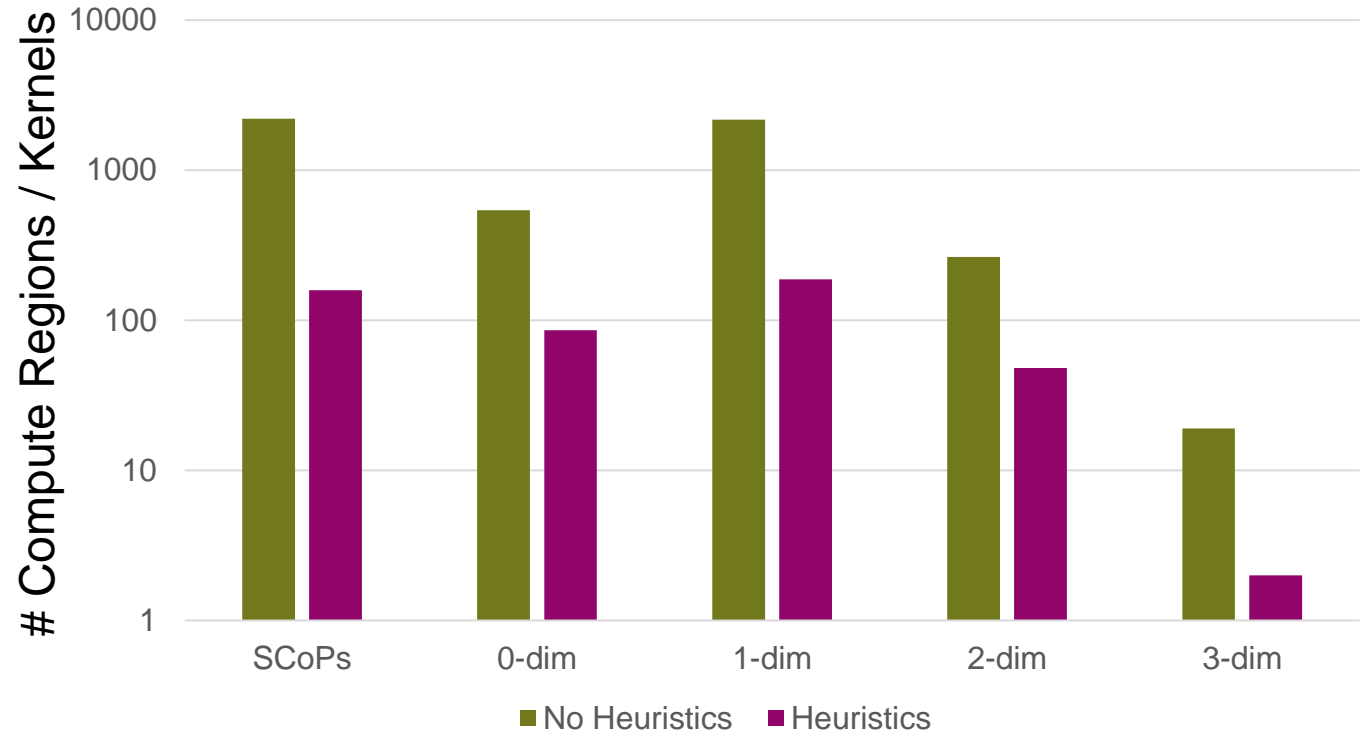
Host Memory                    Device Memory

initialize()

setCenter()

average()

average()          $D \rightarrow H H \rightarrow D$

average()

time

# Evaluation

Workstation: 10 core SandyBridge    NVIDIA Titan Black (Kepler)
Mobile:          4 core Haswell        NVIDIA GT730M (Kepler)

# LLVM Nightly Test Suite



T. Grosser, TH: Polly-ACC: Transparent compilation to heterogeneous hardware, ACM ICS'16

# Some results: Polybench 3.2



Speedup over icc –O3

geomean: ~6x

arithmean: ~30x

legend: icc.parallel | polly | polly-gpu
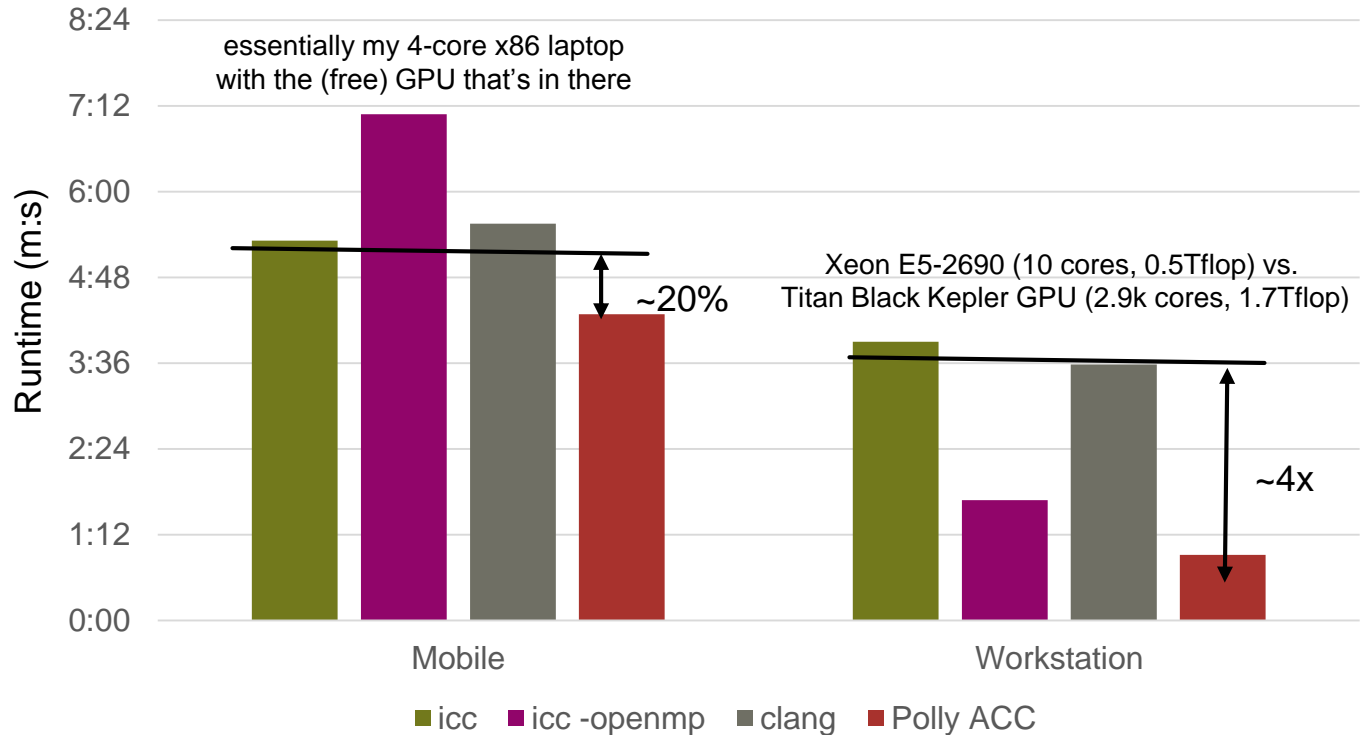
Xeon E5-2690 (10 cores, 0.5Tflop) vs. Titan Black Kepler GPU (2.9k cores, 1.7Tflop)

T. Grosser, TH: Polly-ACC: Transparent compilation to heterogeneous hardware, ACM ICS'16

18

**ETH***zürich*

# Compiles all of SPEC CPU 2006 – Example: LBM



essentially my 4-core x86 laptop
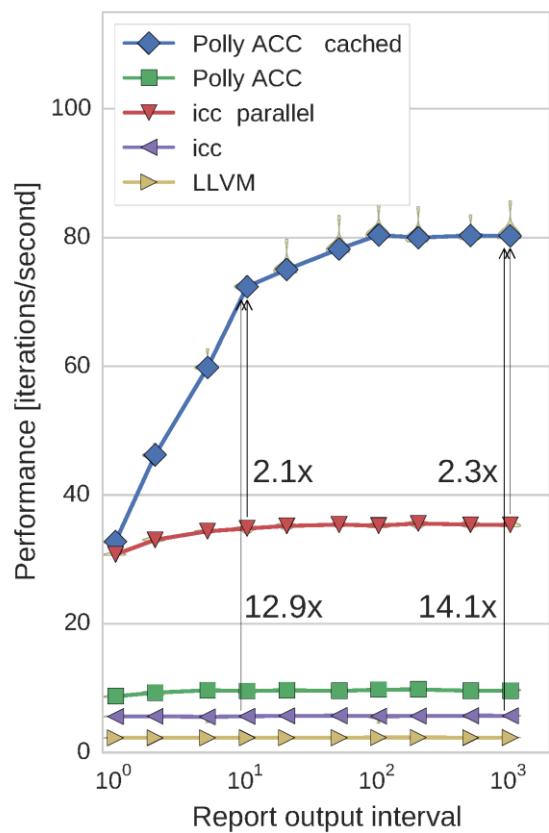with the (free) GPU that's in there

Xeon E5-2690 (10 cores, 0.5Tflop) vs.
Titan Black Kepler GPU (2.9k cores, 1.7Tflop)

~20%

~4x

**Runtime (m:s)**

8:24
7:12
6:00
4:48
3:36
2:24
1:12
0:00

Mobile          Workstation

■ icc  ■ icc -openmp  ■ clang  ■ Polly ACC

T. Grosser, TH: Polly-ACC: Transparent compilation to heterogeneous hardware, ACM ICS'16

# Cactus ADM (SPEC 2006)

# Cactus ADM (SPEC 2006) - Data Transfer



T. Grosser, TH: Polly-ACC: Transparent compilation to heterogeneous hardware, ACM ICS'16
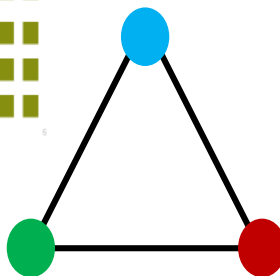
# Polly-ACC

http://spcl.inf.ethz.ch/Polly-ACC

**Mapping Computation to Device**

**Automatic**

**"Regression Free"**

**High Performance**



T. Grosser, TH: Polly-ACC: Transparent compilation to heterogeneous hardware, ACM ICS'16

# Brave new compiler world!?

- **Unfortunately not …**
  - Limited to affine code regions
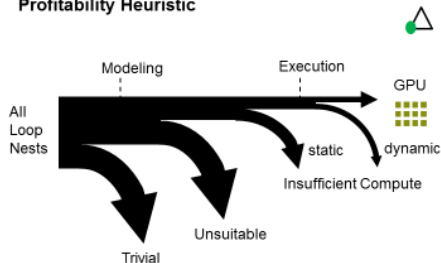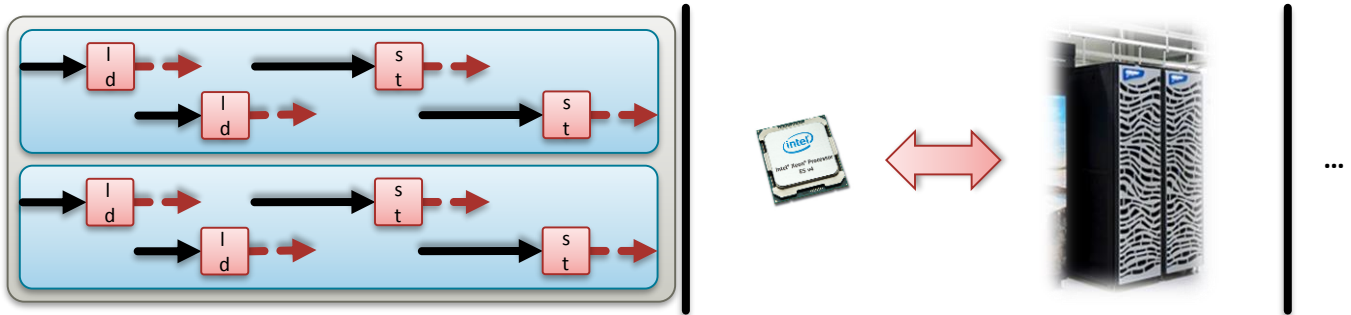  - Maybe generalizes to control-restricted programs
  - No distributed anything!!

- **Good news:**
  - Much of traditional HPC fits that model
  - Infrastructure is coming along




Platform for Advanced Scientific Computing

- **Bad news:**
  - Modern data-driven HPC and Big Data fits less well
  - Need a programming model for **distributed** heterogeneous machines!

# How do we program GPUs today?



**CUDA**
- over-subscribe hardware
- use spare parallel slack for latency hiding

**MPI**
- host controlled
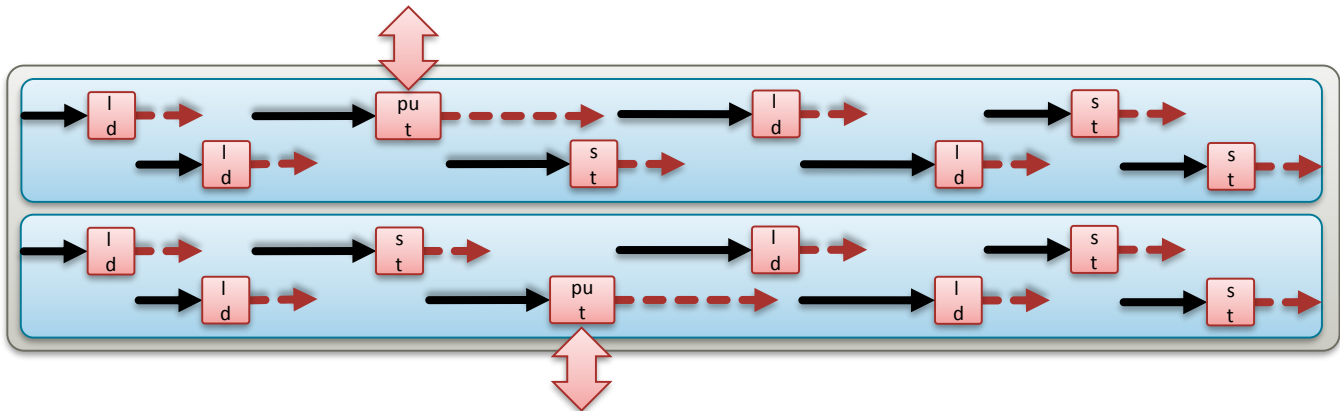- full device synchronization

device · compute core · active thread · instruction latency

# Latency hiding at the cluster level?



**dCUDA (distributed CUDA)**

- unified programming model for GPU clusters
- avoid unnecessary device synchronization to enable system wide latency hiding

device · compute core · active thread · instruction latency

# Talk on Wednesday

**Tobias Gysi, Jeremiah Baer, TH: "dCUDA: Hardware Supported Overlap of Computation and Communication"**

**Wednesday, Nov. 16th**

**4:00-4:30pm**

**Room 355-D**