



Polly Labs



U.S. DEPARTMENT OF
ENERGY

Code Analysis and Optimization using Loop Hierarchies

Workshop on Program Synthesis for Scientific Computing

Michael Kruse, Hal Finkel

Argonne Leadership Computing Facility
Argonne National Laboratory

2020-08-05

Table of Contents

1 The Idea

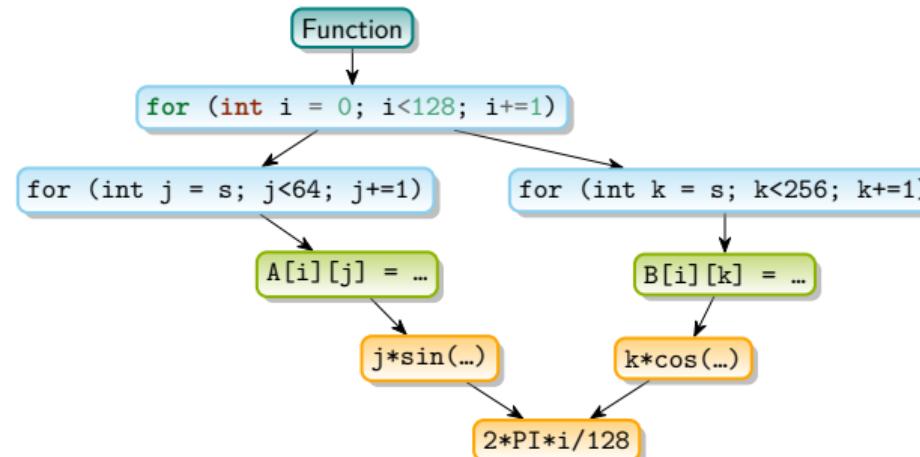
2 Impact

3 Compute Hierarchy

4 Conclusion

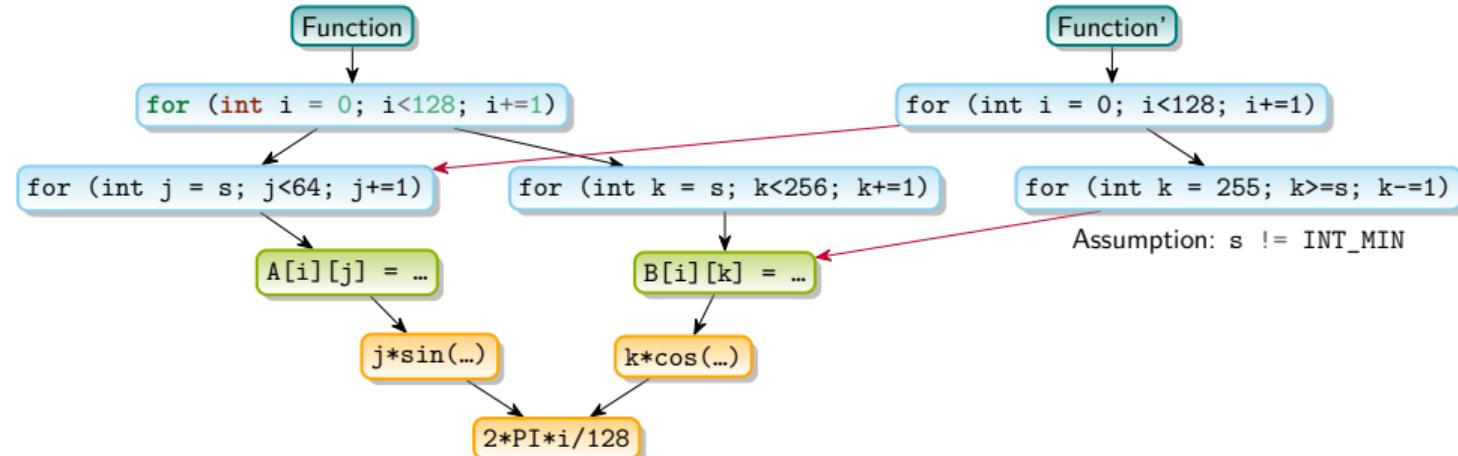
Loop Hierarchy DAG

```
void Function(int s) {  
    for (int i = 0; i < 128; i+=1) {  
        for (int j = s; j < 64; j+=1) A[i][j] = j*sin(2*PI*i/128);  
        for (int k = s; k < 256; k+=1) B[i][k] = k*cos(2*PI*i/128);  
    }  
}
```

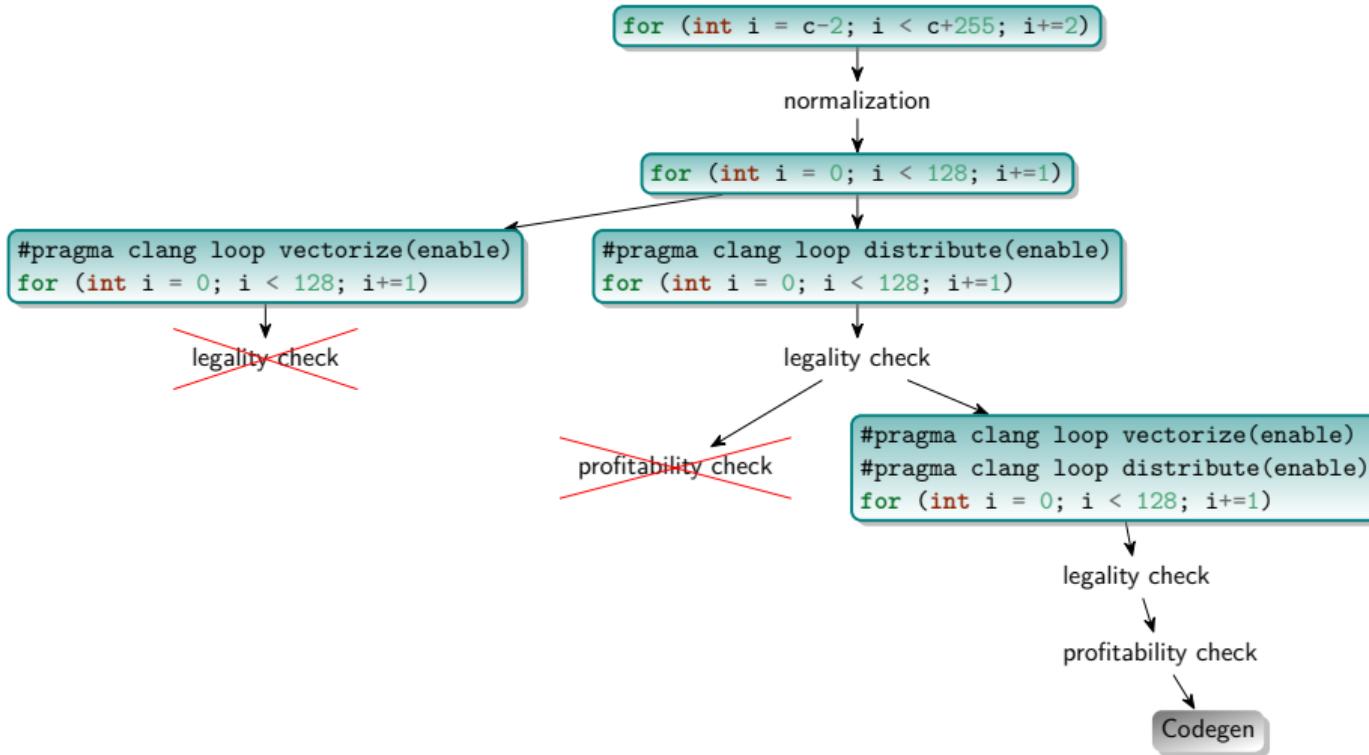


Loop Hierarchy DAG

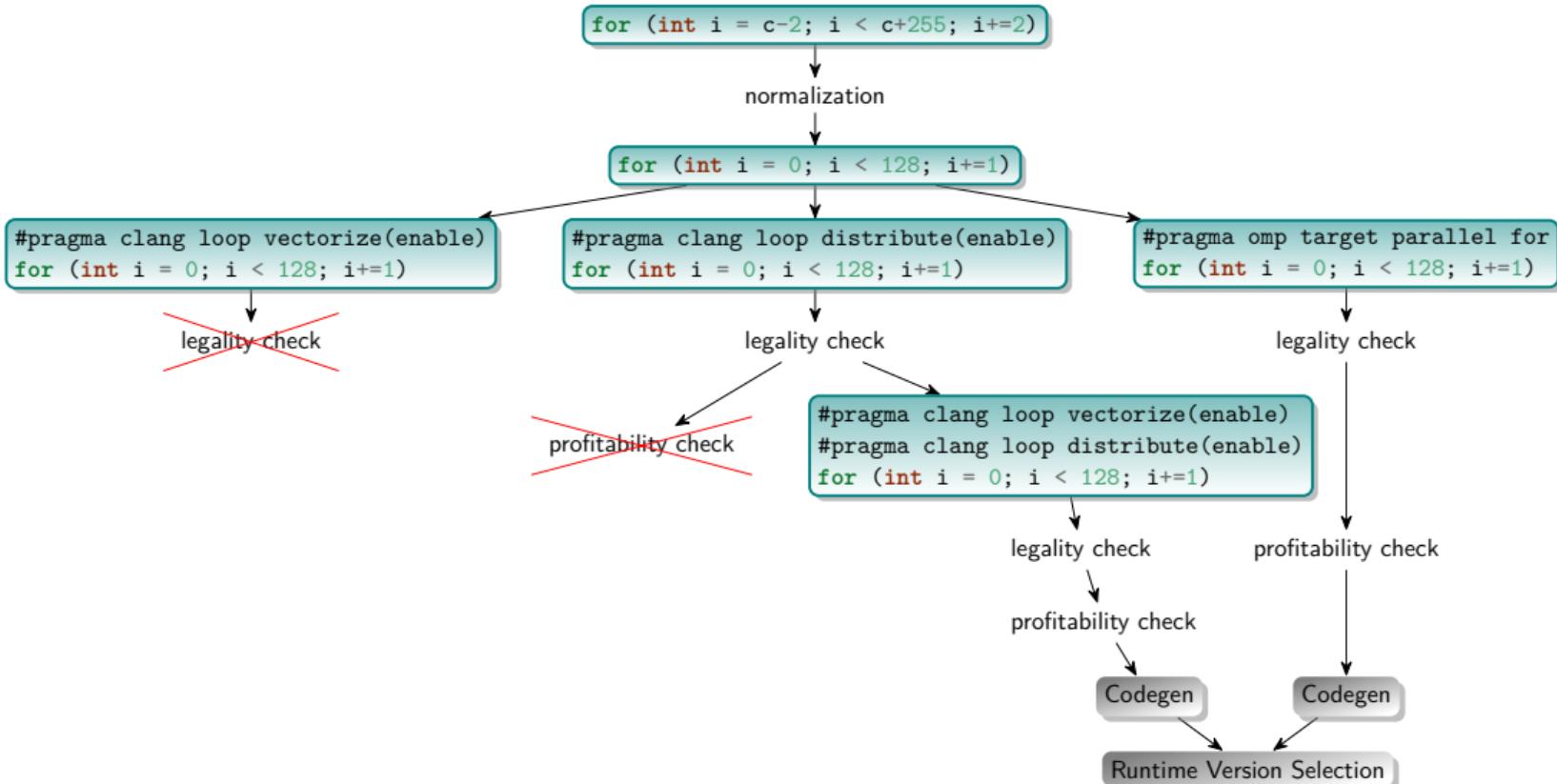
```
void Function(int s) {
    for (int i = 0; i < 128; i+=1) {
        for (int j = s ; j < 64; j+=1) A[i][j] = j*sin(2*PI*i/128);
        for (int k = 255; k >= s ; k-=1) B[i][k] = k*cos(2*PI*i/128);
    }
}
```



Shallow Copy Mechanics



Shallow Copy Mechanics



Legality Check

Comparison between known-good (original) and transformed loop tree

- Executes all statement instances
- No additional instances
- No dependency reversals

Profitability Check

Execution Time Model

- Estimate cycles of straight-line code
 - `llvm-mca`
 - Memory access latency
- Estimate trip count
 - Arbitrary constant ("100")
 - "infinity" (only innermost kernel counts)
 - From user annotations (`#pragma loop count(n)`)
 - From PGO / previous JIT stage

Profitability Check

Execution Time Model

- Estimate cycles of straight-line code
 - `llvm-mca`
 - Memory access latency
- Estimate trip count
 - Arbitrary constant ("100")
 - "infinity" (only innermost kernel counts)
 - From user annotations (`#pragma loop count(n)`)
 - From PGO / previous JIT stage

Autotuning

- Select most-promising not-yet-evaluated
 - or –
- Select known-fastest

Profitability Check

Execution Time Model

- Estimate cycles of straight-line code
 - `llvm-mca`
 - Memory access latency
- Estimate trip count
 - Arbitrary constant ("100")
 - "infinity" (only innermost kernel counts)
 - From user annotations (`#pragma loop count(n)`)
 - From PGO / previous JIT stage

Autotuning

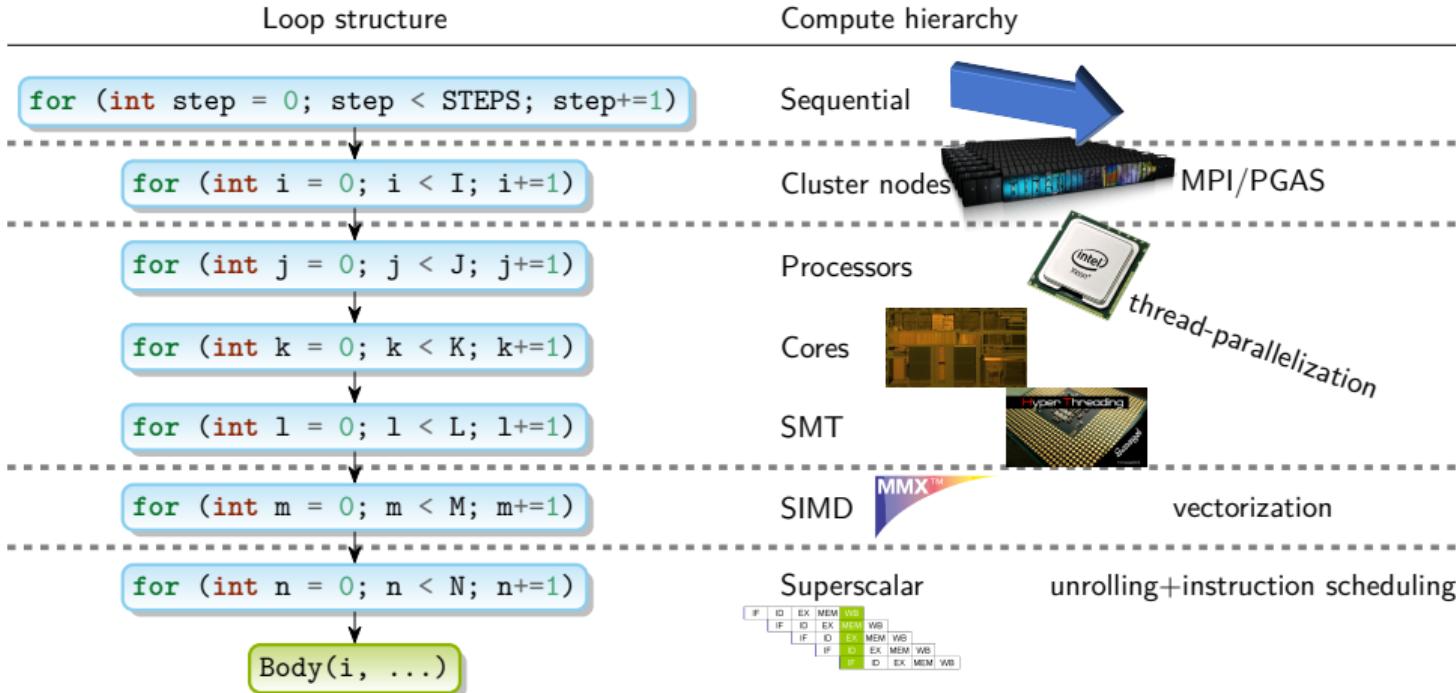
- Select most-promising not-yet-evaluated
 - or –
- Select known-fastest

Machine Learning

- Apply a per-architecture pre-trained model

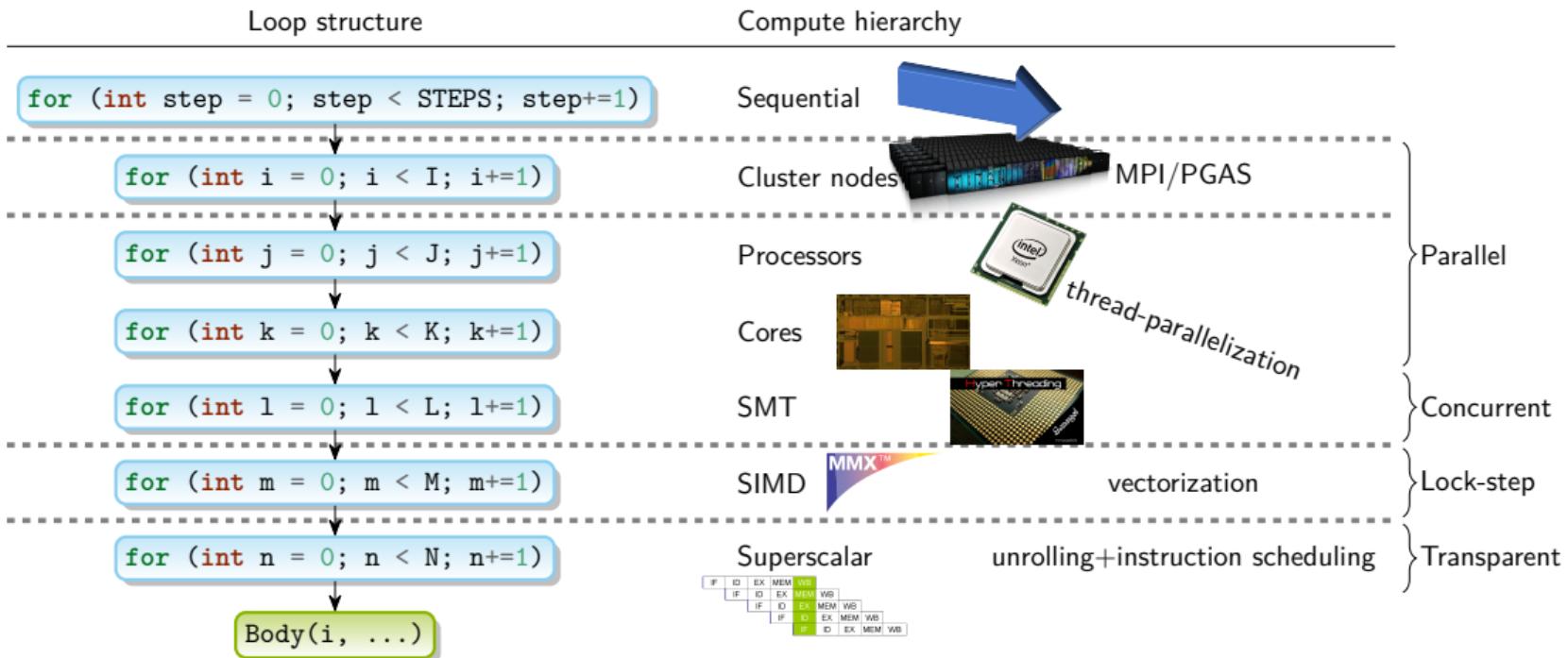
Static Compute Hierarchy

CPU-based



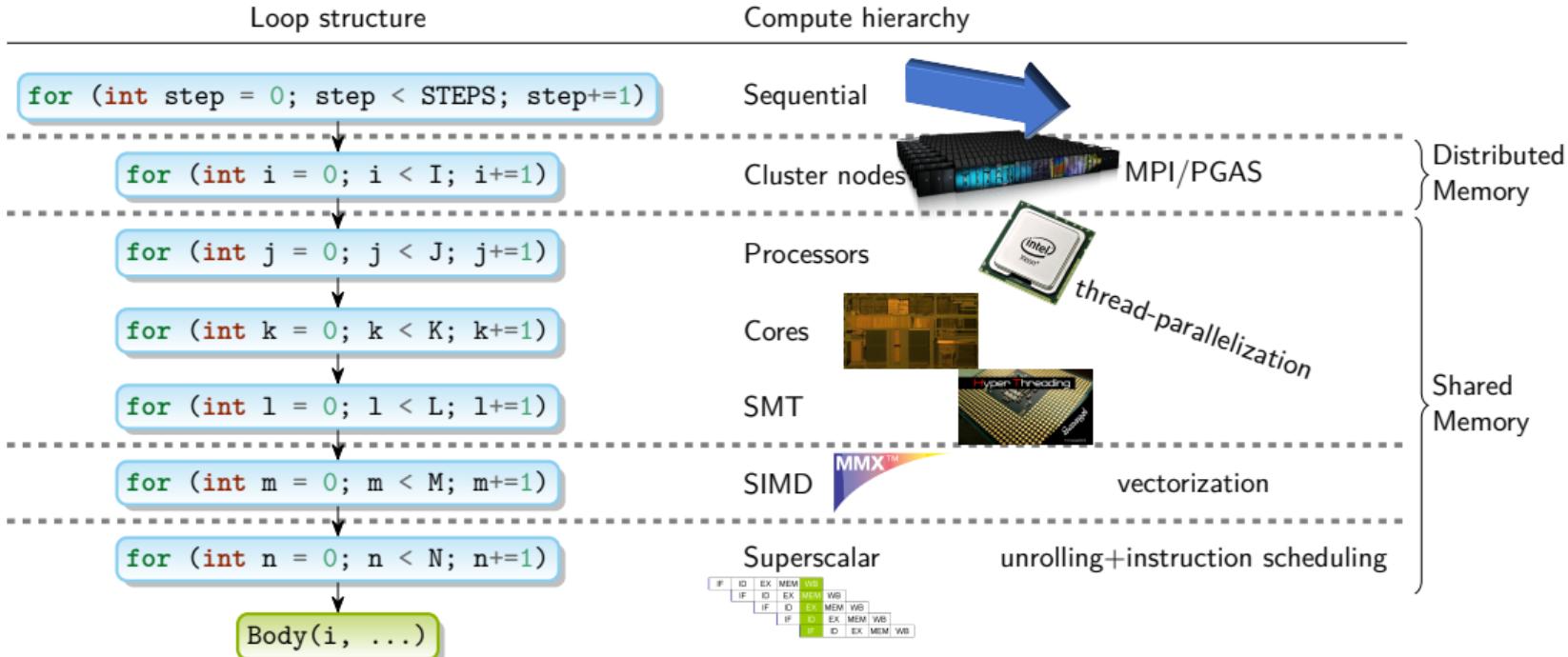
Static Compute Hierarchy

CPU-based



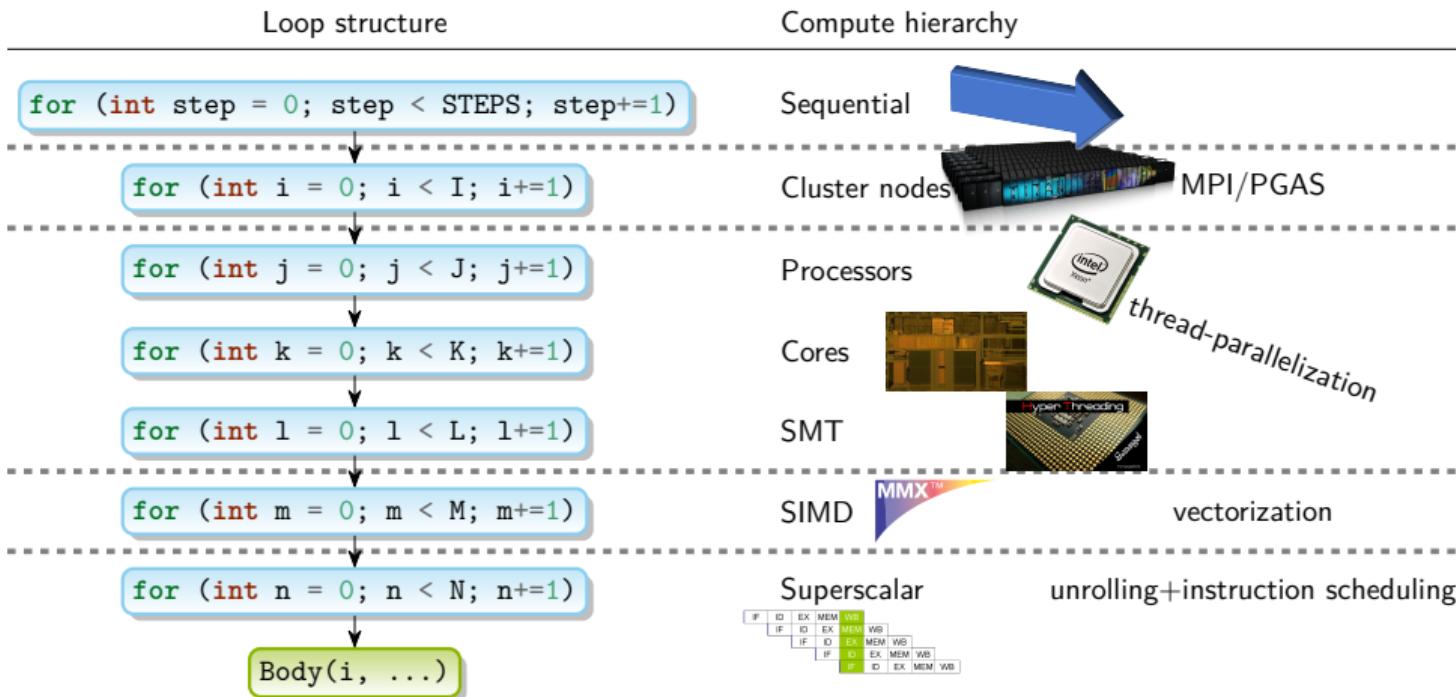
Static Compute Hierarchy

CPU-based



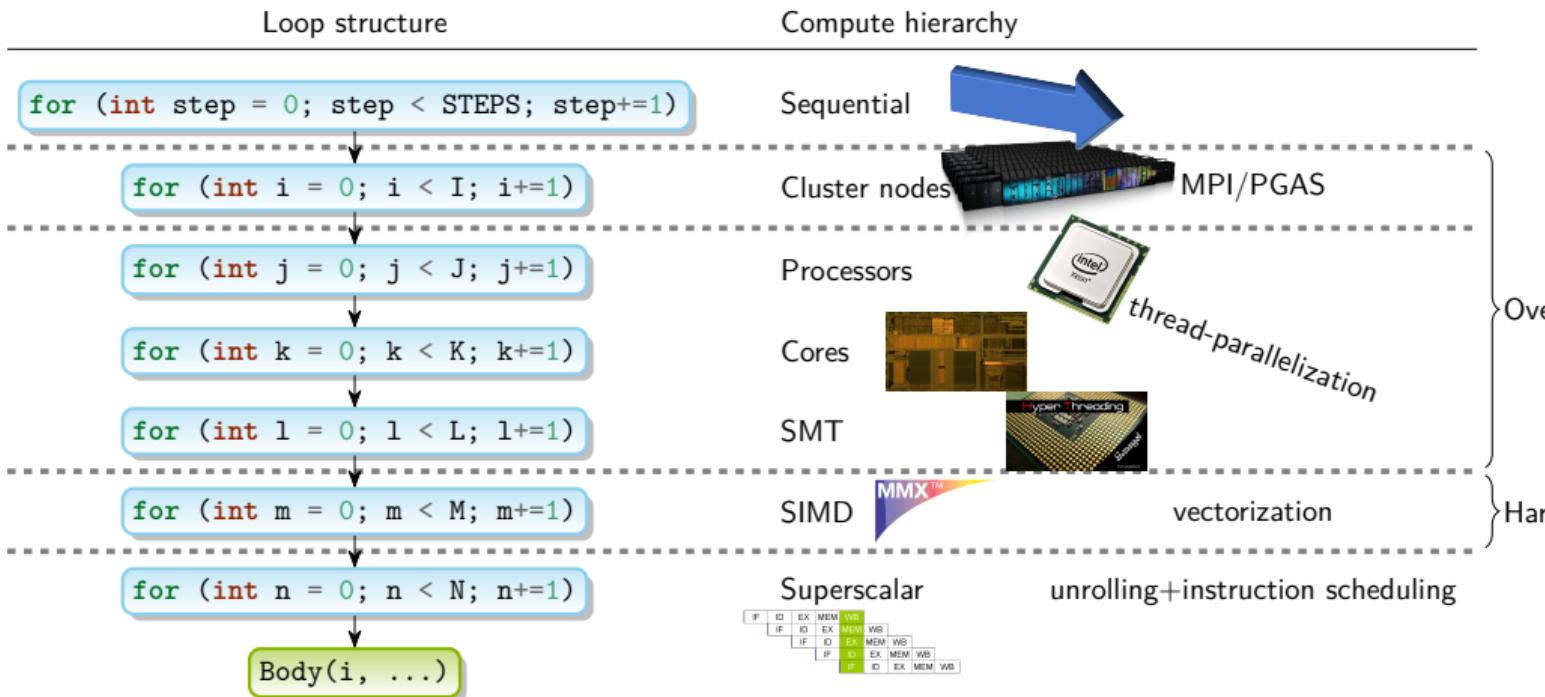
Static Compute Hierarchy

CPU-based



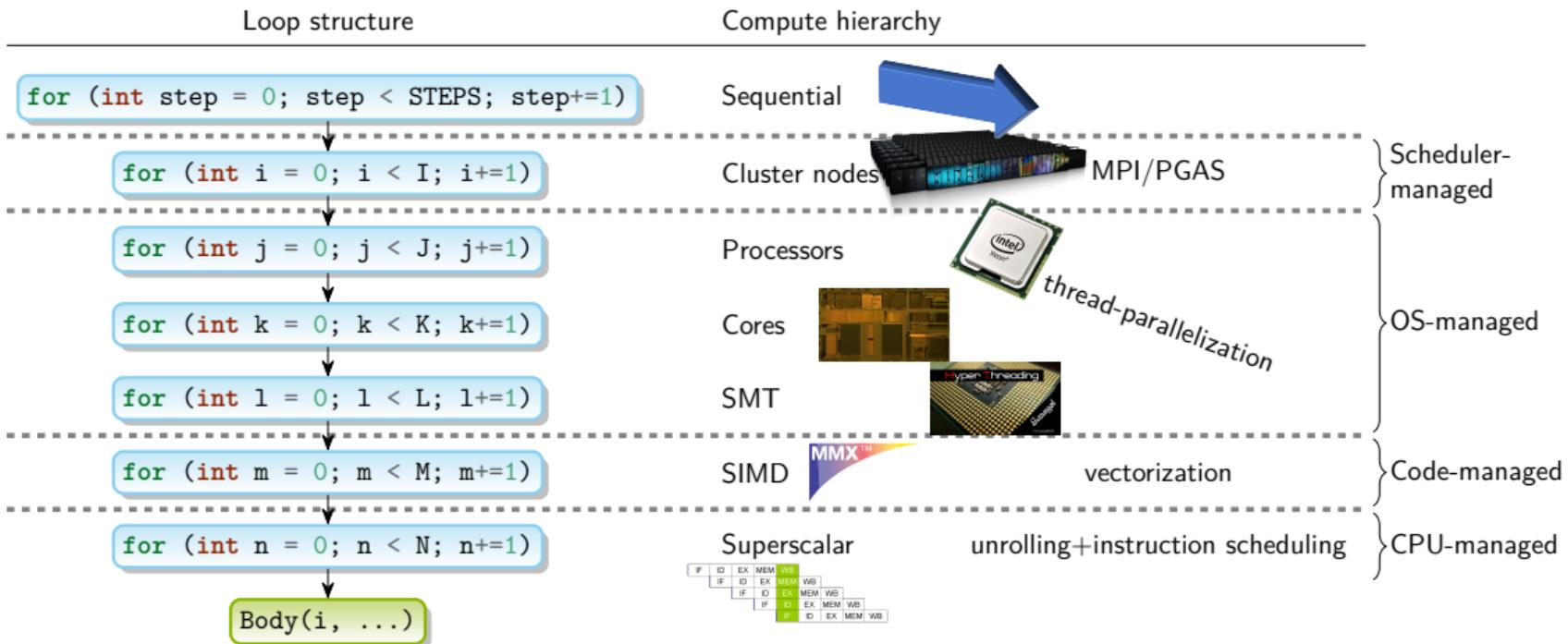
Static Compute Hierarchy

CPU-based



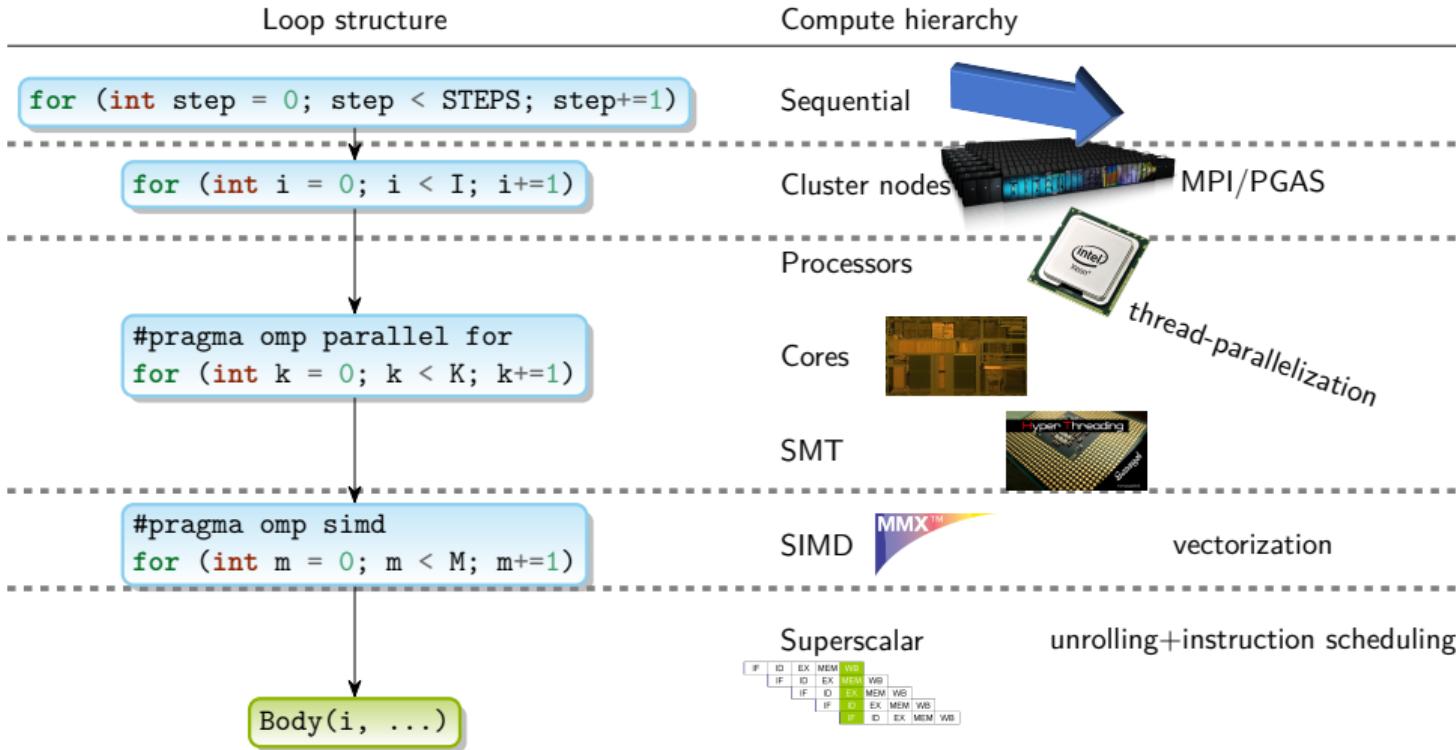
Static Compute Hierarchy

CPU-based



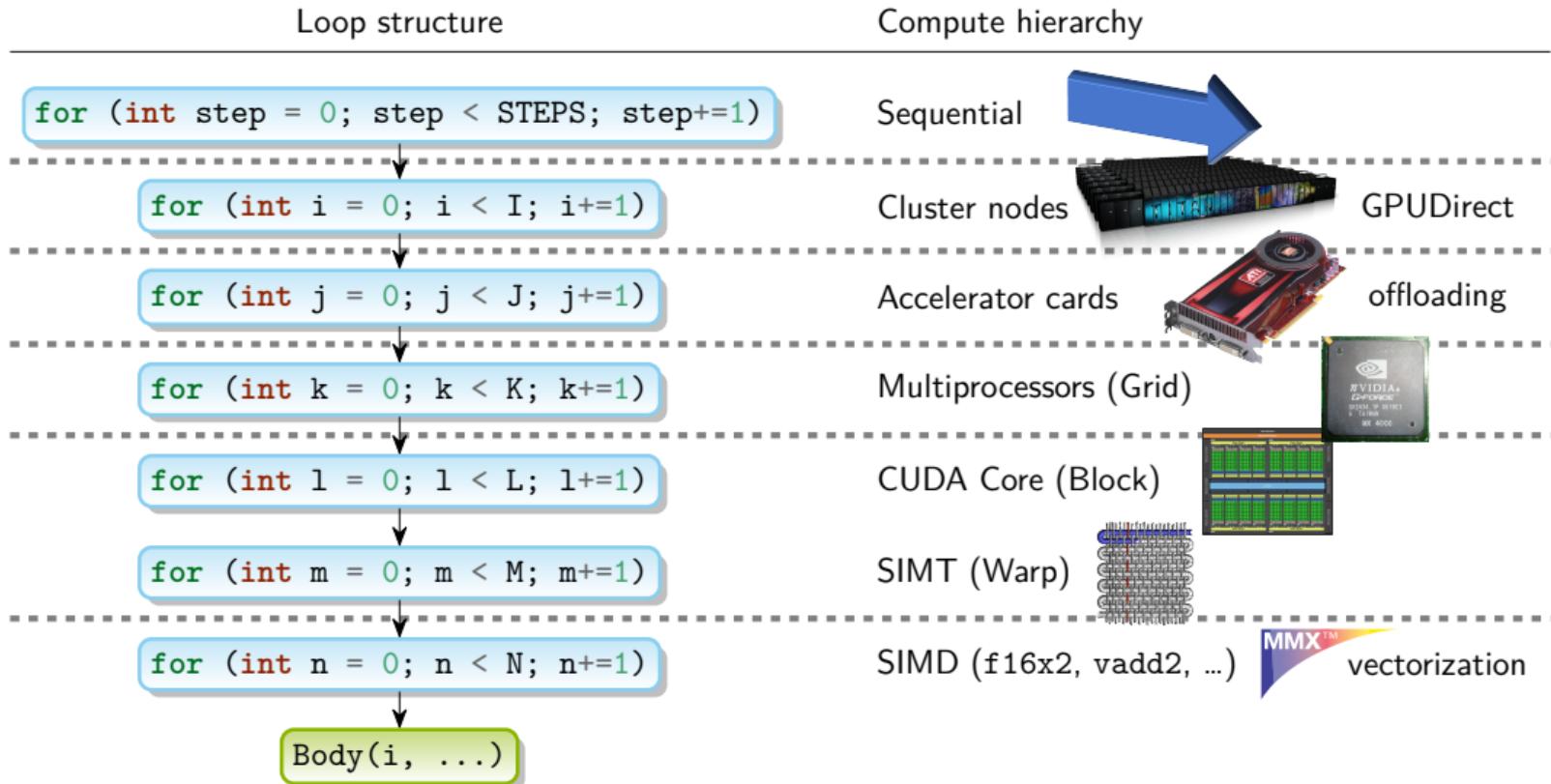
Static Compute Hierarchy

CPU-based (OpenMP)



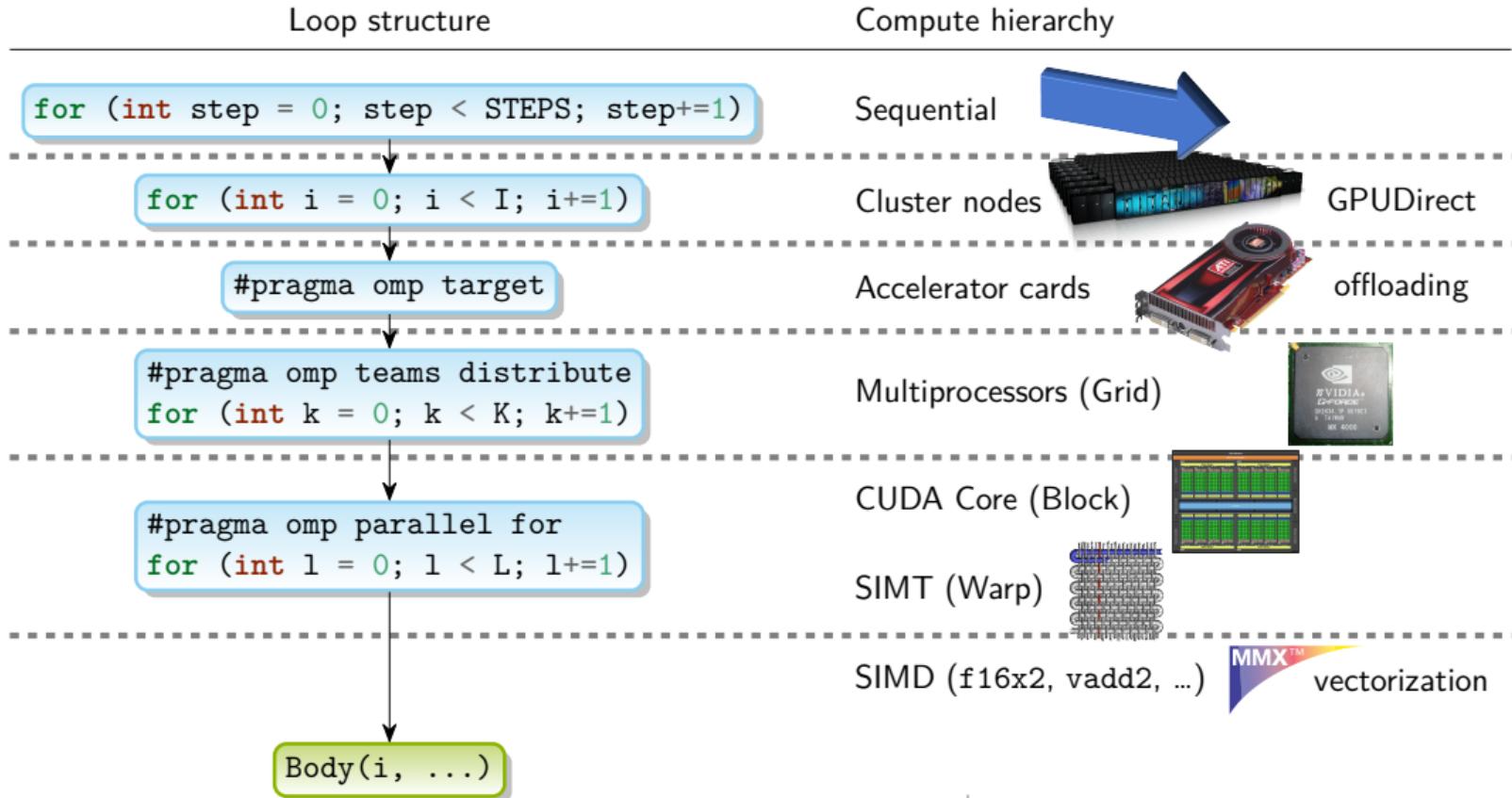
Static Compute Hierarchy

GPU-based (CUDA)

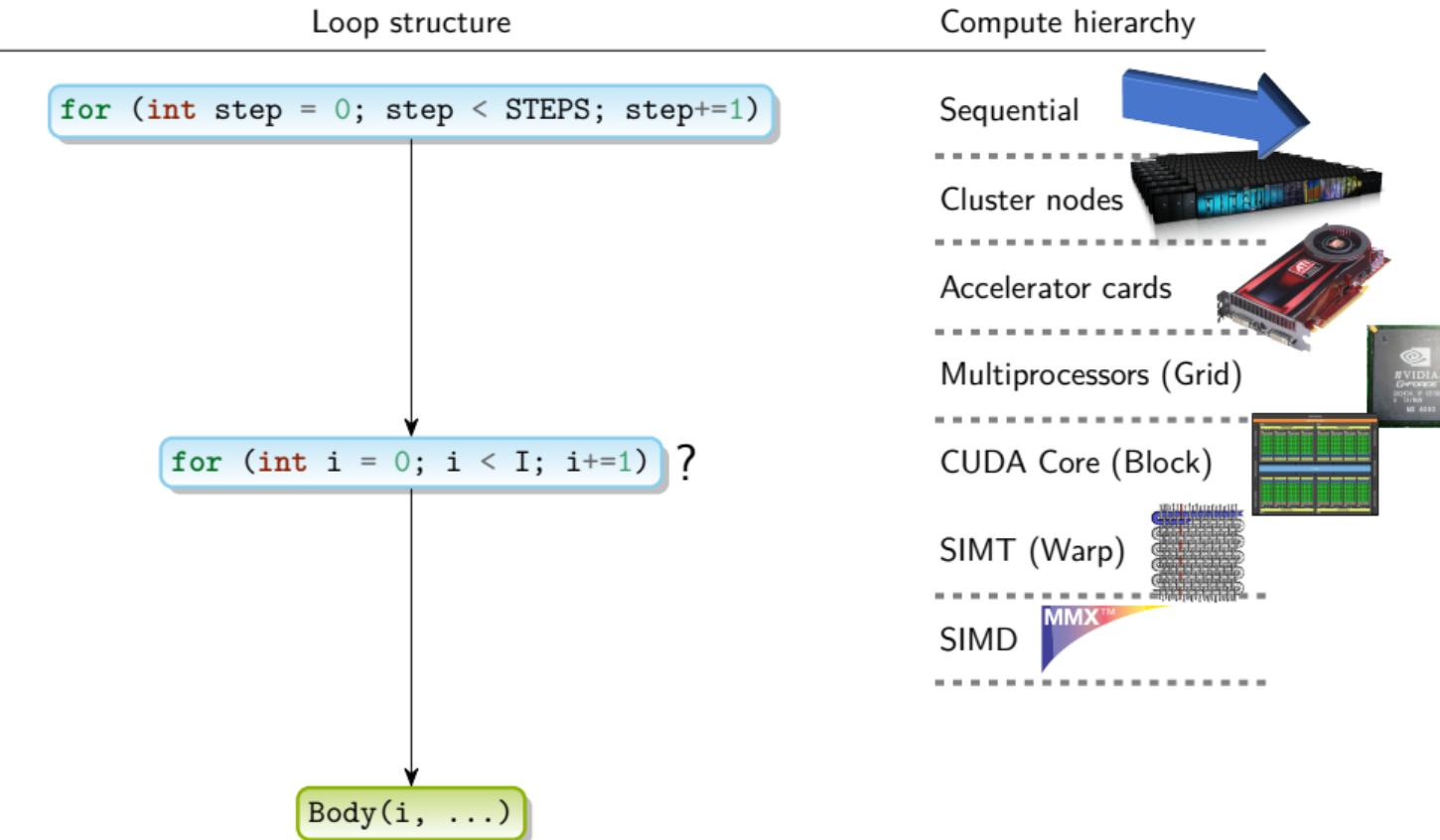


Static Compute Hierarchy

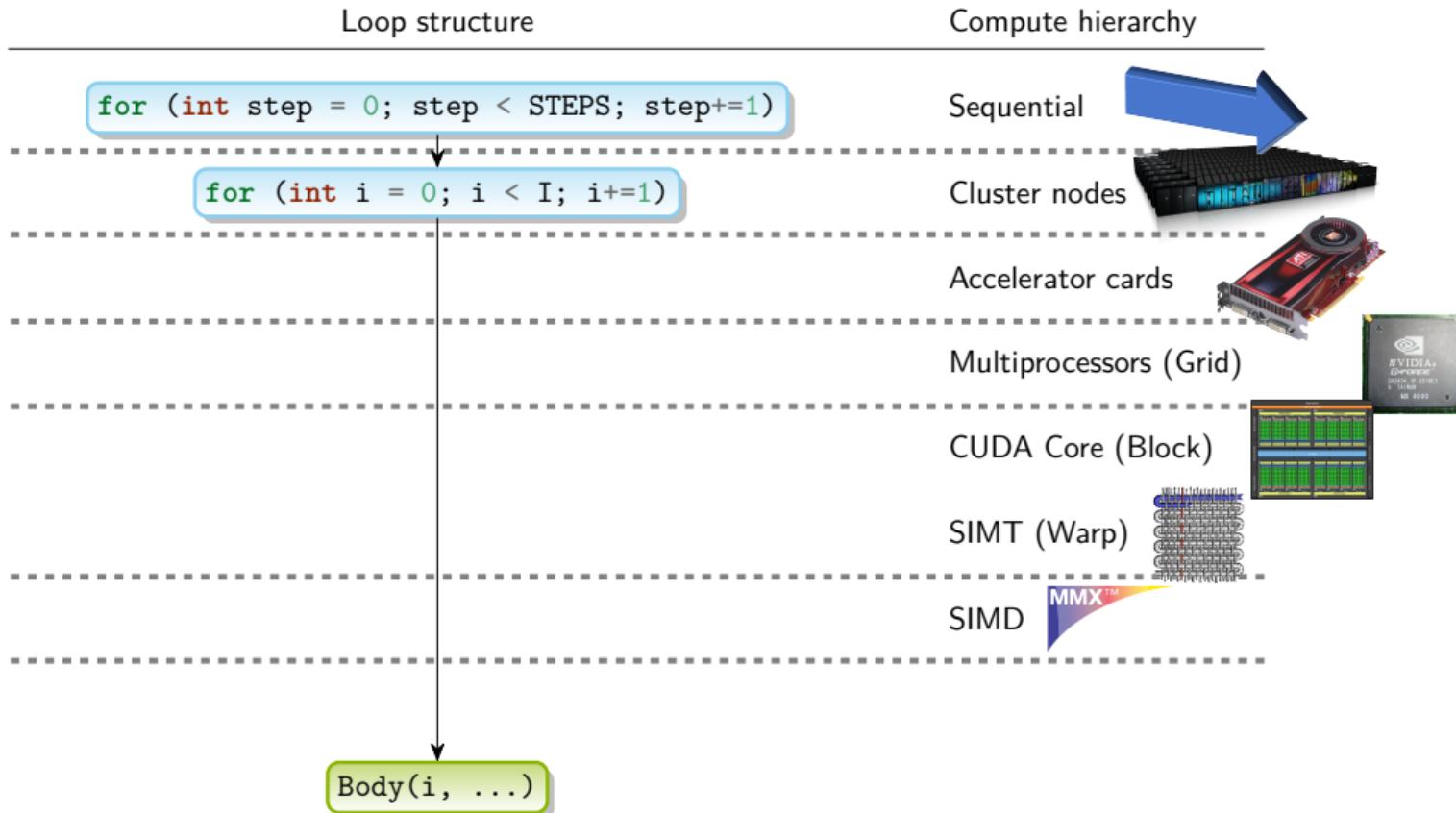
GPU-based (OpenMP-to-ptx)



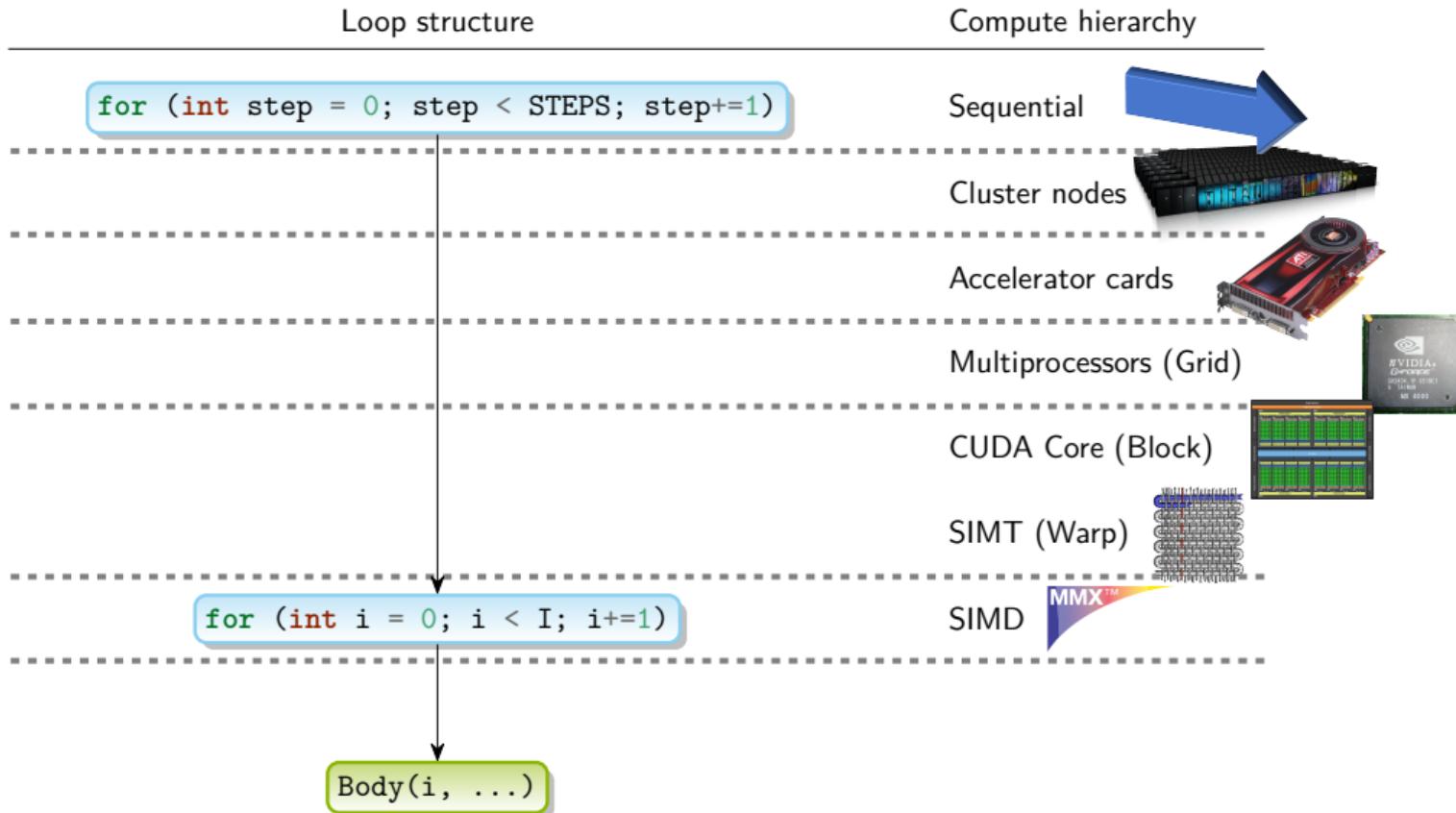
Compute Hierarchy Mapping



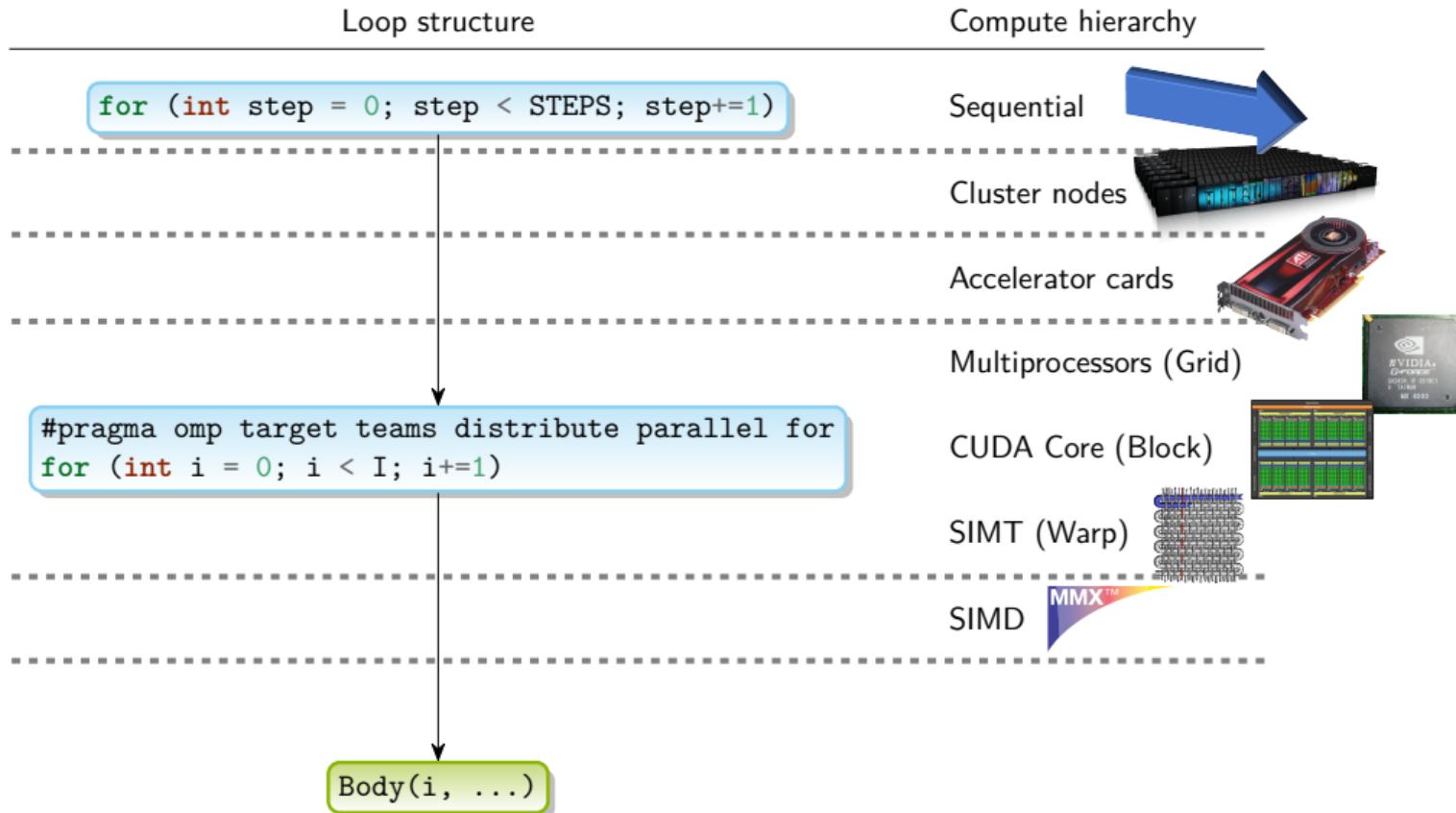
Compute Hierarchy Mapping



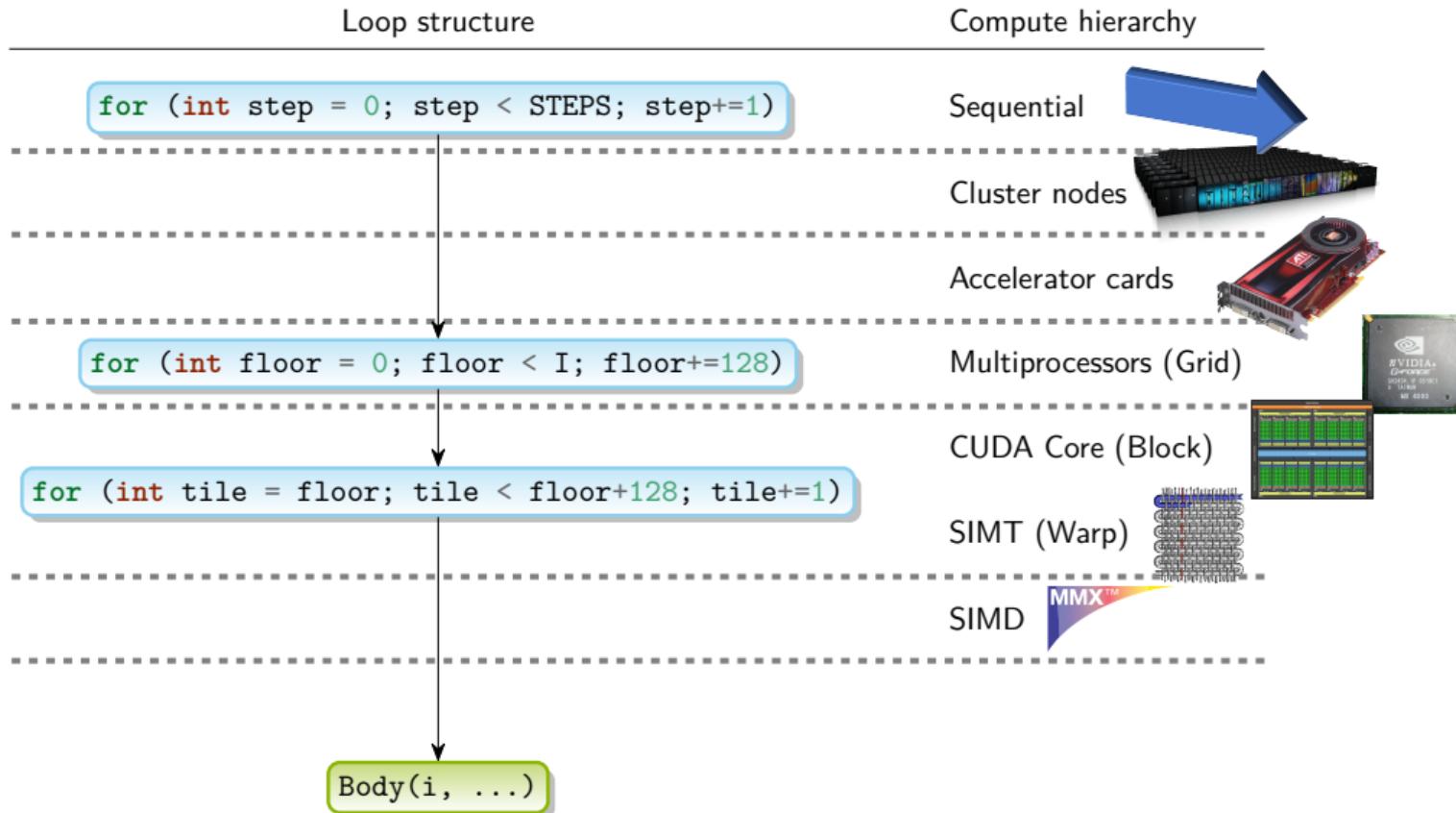
Compute Hierarchy Mapping



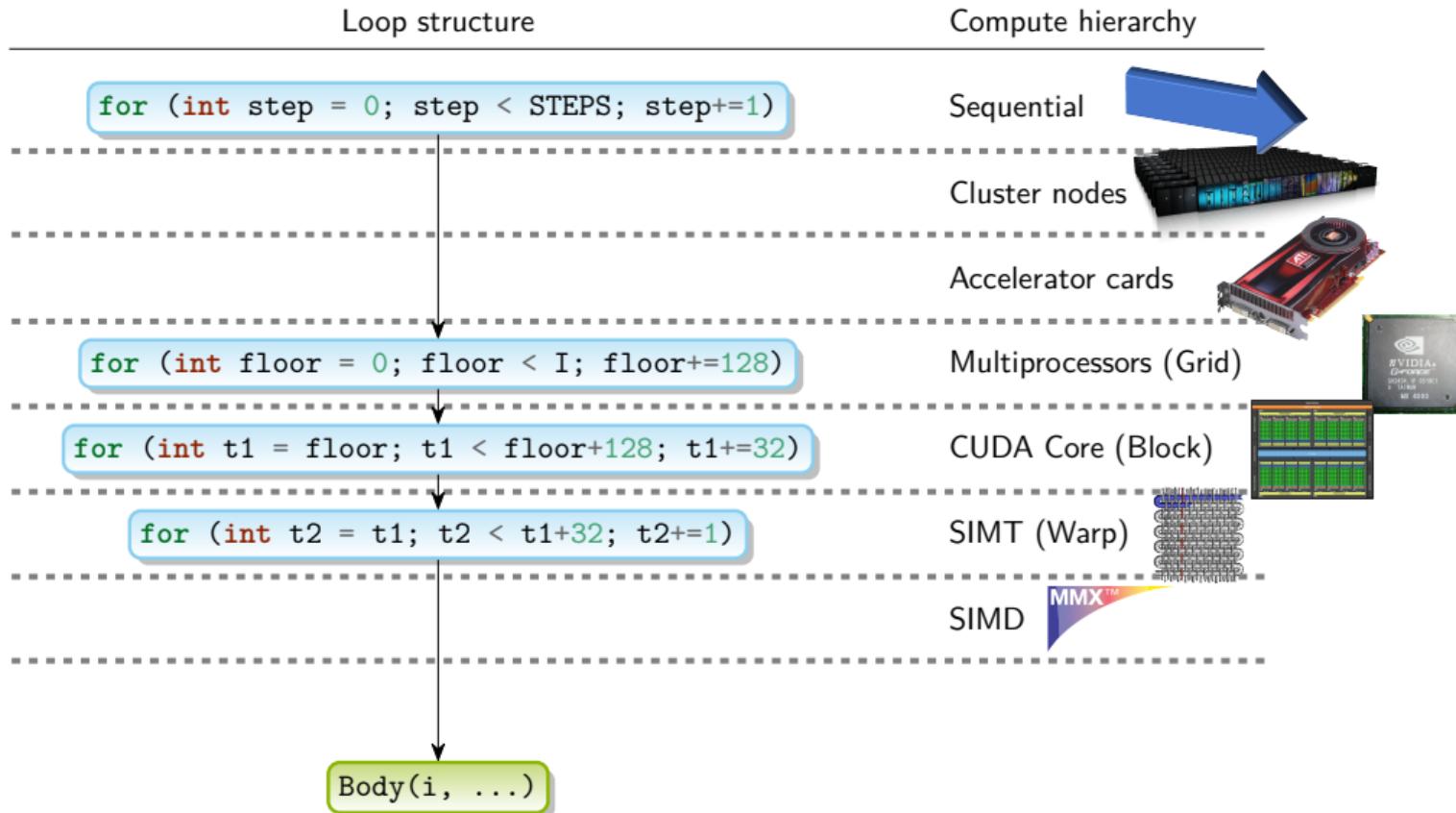
Compute Hierarchy Mapping



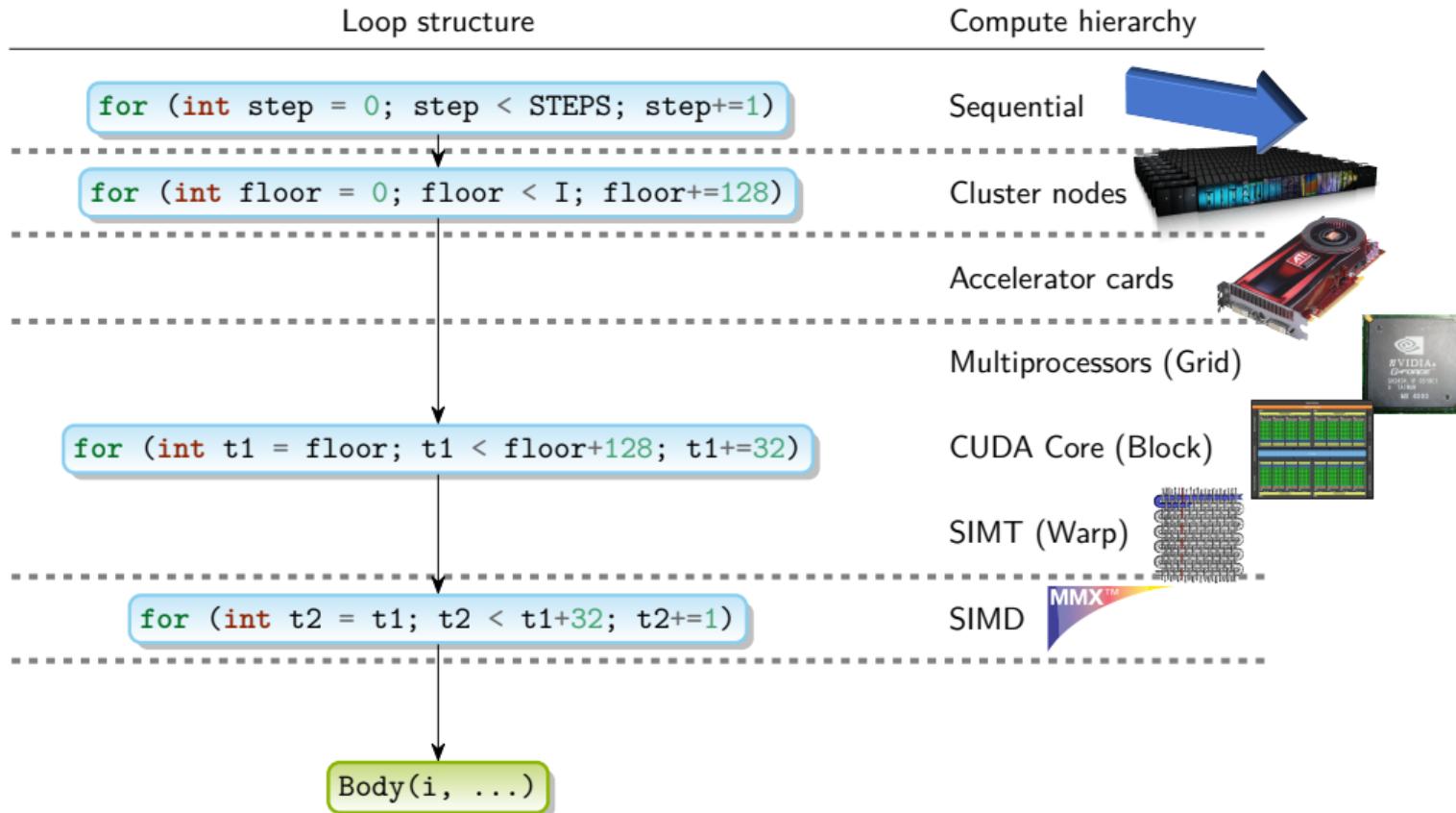
Compute Hierarchy Mapping



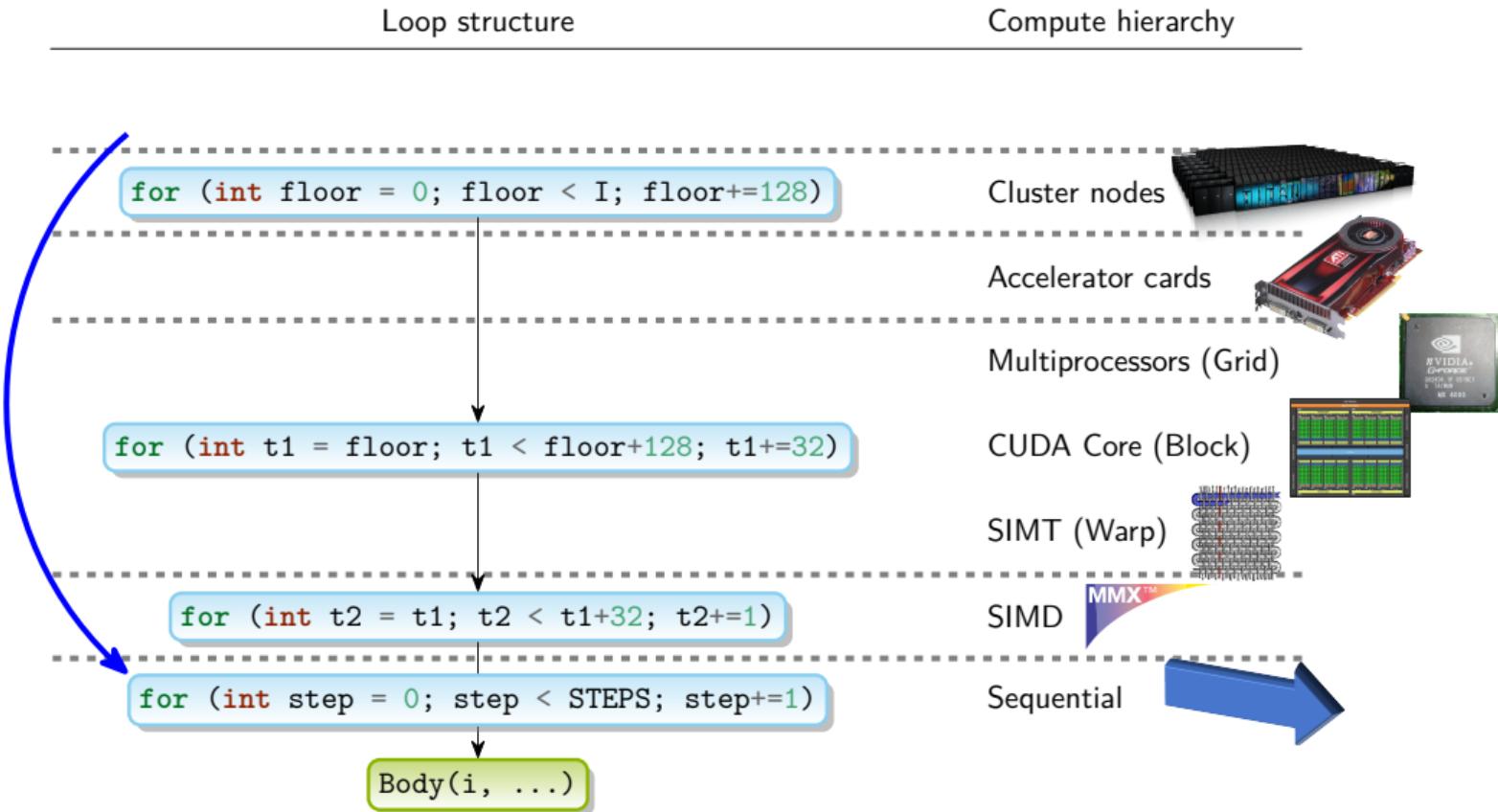
Compute Hierarchy Mapping



Compute Hierarchy Mapping



Compute Hierarchy Mapping



(Dis-)advantages

Advantages

- Speculatively apply transformations
 - Generalized legality check
 - Generalized profitability check
- Loop canonicalization only applied to DAG
- Assumptions not part of core structure
 - No per-transformation code versioning
- Independent of underlying IR

Disadvantages

- Cannot directly reference siblings/parents
 - Def/Use chains
- LLVM-IR, MachineIR, VPlan, MLIR, ...



That's all Folks!