

# COMPILER OPTIMIZATION FOR (OPENMP) ACCELERATOR OFFLOADING

---

EuroLLVM — April 8, 2019 — Brussels, Belgium

Johannes Doerfert and Hal Finkel

Leadership Computing Facility  
Argonne National Laboratory  
<https://www.alcf.anl.gov/>



## ACKNOWLEDGMENT

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration) responsible for the planning and preparation of a capable exascale ecosystem, including software, applications, hardware, advanced system engineering, and early testbed platforms, in support of the nation's exascale computing imperative.



## Original Program

```
int y = 7;  
  
for (i = 0; i < N; i++) {  
    f(y, i);  
}  
g(y);
```

## After Optimizations

```
for (i = 0; i < N; i++) {  
    f(7, i);  
}  
g(7);
```



## Original Program

```
int y = 7;  
#pragma omp parallel for  
for (i = 0; i < N; i++) {  
    f(y, i);  
}  
g(y);
```

## After Optimizations



## Original Program

```
int y = 7;
#pragma omp parallel for
for (i = 0; i < N; i++) {
    f(y, i);
}
g(y);
```

## After Optimizations

```
int y = 7;
#pragma omp parallel for
for (i = 0; i < N; i++) {
    f(y, i);
}
g(y);
```



# CURRENT COMPILER OPTIMIZATION FOR PARALLELISM



None<sup>\*†</sup>

---

\*At least for LLVM/Clang up to 8.0

†And not considering smart **runtime libraries!**

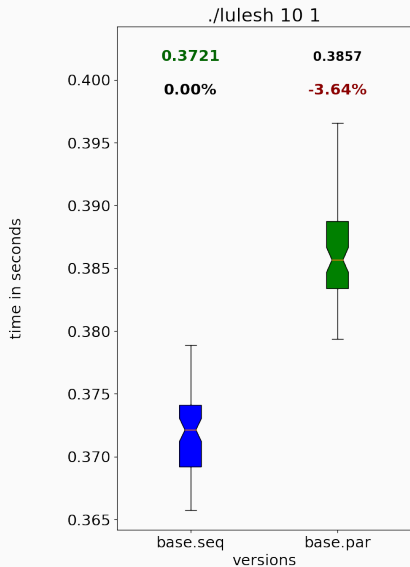


Why is this important?

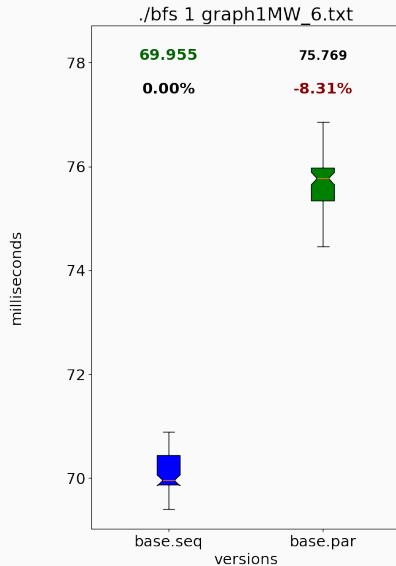
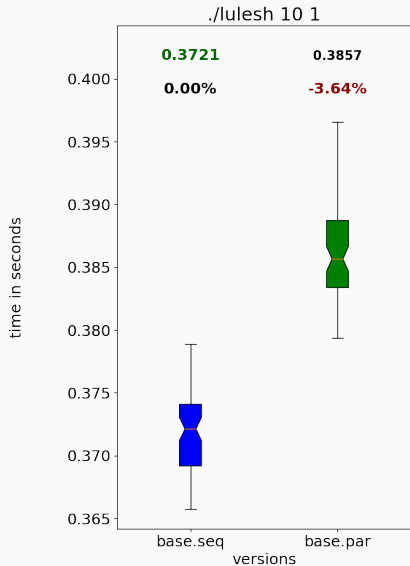




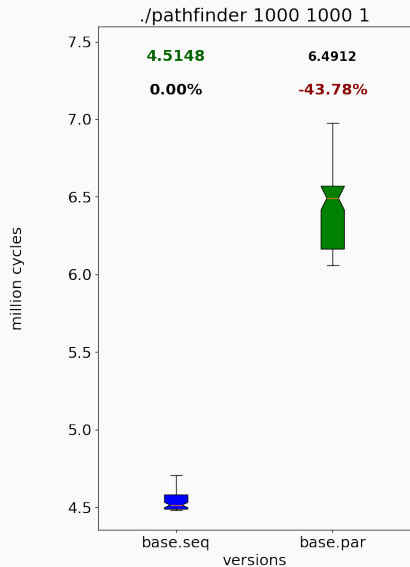
# PERFORMANCE IMPLICATIONS



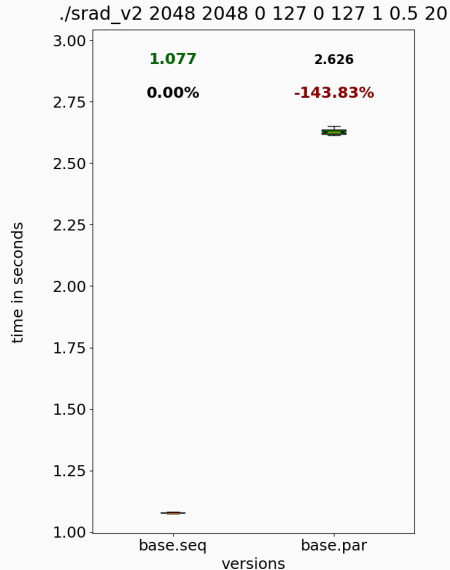
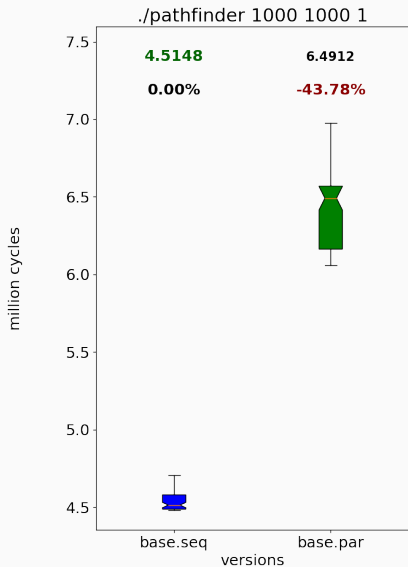
# PERFORMANCE IMPLICATIONS



# PERFORMANCE IMPLICATIONS



# PERFORMANCE IMPLICATIONS



Optimizations for *sequential* aspects

Optimizations for *parallel* aspects

---

$a$

$b$



## Optimizations for *sequential* aspects

- May reuse *existing* transformations

(patches up for review!)

## Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations*

---

*a*

*b*



## Optimizations for *sequential* aspects

- May reuse *existing* transformations (patches up for review!)
- ⇒ Introduce *suitable abstractions* to bridge the indirection (DONE!)

## Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations*

---

*a*

*b*



## Optimizations for *sequential* aspects

- May reuse *existing* transformations (patches up for review!)
- ⇒ Introduce *suitable abstractions* to bridge the indirection (DONE!)

## Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations*
- ⇒ Introduce a unifying abstraction layer

---

*a*

*b*





## Optimizations for *sequential* aspects

- May reuse *existing* transformations (patches up for review!)
- ⇒ Introduce *suitable abstractions* to bridge the indirection (DONE!)

## Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations* (see IWOMP'18<sup>a</sup>)
- ⇒ Introduce a unifying abstraction layer (see EuroLLVM'18<sup>b</sup>)

---

<sup>a</sup>Compiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

<sup>b</sup>A Parallel IR in Real Life: Optimizing OpenMP, H. Finkel, J. Doerfert, X. Tian, G. Stelle, Euro



## Optimizations for *sequential* aspects

- May reuse *existing* transformations (patches up for review!)
- ⇒ Introduce *suitable abstractions* to bridge the indirection (DONE!)

## Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations* (see IWOMP'18<sup>a</sup>)
- ⇒ Introduce a unifying abstraction layer (see EuroLLVM'18<sup>b</sup>)

---

<sup>a</sup>Compiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

<sup>b</sup>A Parallel IR in Real Life: Optimizing OpenMP, H. Finkel, J. Doerfert, X. Tian, G. Stelle, Euro



```
{  
  #pragma omp target teams parallel  
  work1();  
  
}
```



```
>> clang -O3 -fopenmp-targets=...
```

```
{  
  #pragma omp target teams parallel  
  work1();  
}
```



```
>> clang -O3 -fopenmp-targets=...
```

```
{  
  #pragma omp target teams parallel  
  work1();  
}
```

“relatively” good performance :)



```
>> clang -O3 -fopenmp-targets=...
```

```
{  
    #pragma omp target teams parallel  
    work1();  
  
    #pragma omp target teams  
    #pragma omp parallel  
    work2();  
}
```



```
>> clang -O3 -fopenmp-targets=...
```

```
{  
  #pragma omp target teams parallel  
  work1();  
  
  #pragma omp target teams  
  #pragma omp parallel  
  work2();  
}
```

“relatively” good performance :)



```
>> clang -O3 -fopenmp-targets=...
```

```
#pragma omp target teams
```

```
{
```

```
#pragma omp parallel
```

```
work1();
```

```
#pragma omp parallel
```

```
work2();
```

```
}
```





```
>> clang -O3 -fopenmp-targets=...
```

```
#pragma omp target teams
```

```
{
```

```
#pragma omp parallel
```

```
work1();
```

```
#pragma omp parallel
```

```
work2();
```

```
}
```

probably poor performance :(



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
{  
  #pragma omp target teams  
  foo();  
  #pragma omp target teams parallel  
  work();                               // <- Hotspot  
  #pragma omp target teams  
  bar();  
}
```



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
{  
  #pragma omp target teams  
  foo();  
  #pragma omp target teams parallel  
  work();                               // <- Hotspot  
  #pragma omp target teams  
  bar();  
}
```

N teams,



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
{  
  #pragma omp target teams  
  foo();  
  #pragma omp target teams parallel  
  work();                               // <- Hotspot  
  #pragma omp target teams  
  bar();  
}
```

N teams, with M threads each,



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
{  
  #pragma omp target teams  
  foo();  
  #pragma omp target teams parallel  
  work();                               // <- Hotspot  
  #pragma omp target teams  
  bar();  
}
```

N teams, with M threads each, all executing **work** concurrently.



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
{

    foo();
    #pragma omp parallel
    work();                               // <- Hotspot

    bar();
}
```



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
{

    foo();
    #pragma omp parallel
    work();                               // <- Hotspot

    bar();
}
```

1 master and N-1 worker teams,



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
{

    foo();
    #pragma omp parallel
    work();                               // <- Hotspot

    bar();
}
```

1 master and N-1 worker teams, worker teams M threads:





## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
{

    foo();
    #pragma omp parallel
    work();                               // <- Hotspot

    bar();
}
```

1 master and N-1 worker teams, worker teams M threads:  
Masters execute **foo** concurrently, workers idle.



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
{

    foo();
    #pragma omp parallel
    work();                               // <- Hotspot

    bar();
}
```

1 master and N-1 worker teams, worker teams M threads:  
Masters execute **foo** concurrently, workers idle.  
Masters delegate **work** for concurrent execution.



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
{

    foo();
    #pragma omp parallel
    work();                               // <- Hotspot

    bar();
}
```

1 master and N-1 worker teams, worker teams M threads:

Masters execute **foo** concurrently, workers idle.

Masters delegate **work** for concurrent execution.

Masters execute **bar** concurrently, workers idle.



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
```

```
{
```

Problems:

```
}
```

```
1 m
```

Masters execute **foo** concurrently, workers idle.

Masters delegate **work** for concurrent execution.

Masters execute **bar** concurrently, workers idle.



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
```

```
{
```

Problems:

- a separate master team costs resources

```
}
```

```
1 m
```

Masters execute **foo** concurrently, workers idle.

Masters delegate **work** for concurrent execution.

Masters execute **bar** concurrently, workers idle.



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
```

```
{
```

Problems:

- a separate master team costs resources
- synchronization has overhead

```
}
```

```
1 m
```

Masters execute **foo** concurrently, workers idle.

Masters delegate **work** for concurrent execution.

Masters execute **bar** concurrently, workers idle.



## THE COMPILER BLACK BOX — BEHIND THE CURTAIN (OF CLANG)

```
#pragma omp target teams
```

```
{
```

Problems:

- a separate master team costs resources
- synchronization has overhead
- currently impossible to optimization

```
}
```

```
1 m
```

Masters execute **foo** concurrently, workers idle.

Masters delegate **work** for concurrent execution.

Masters execute **bar** concurrently, workers idle.



OpenMP

Clang

LLVM-IR

LLVM

Assembly





OpenMP

Clang

LLVM-IR

LLVM

Assembly

Code



# OPENMP OFFLOAD — OVERVIEW

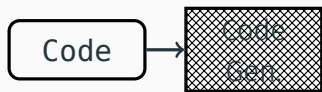
OpenMP

Clang

LLVM-IR

LLVM

Assembly



# OPENMP OFFLOAD — OVERVIEW

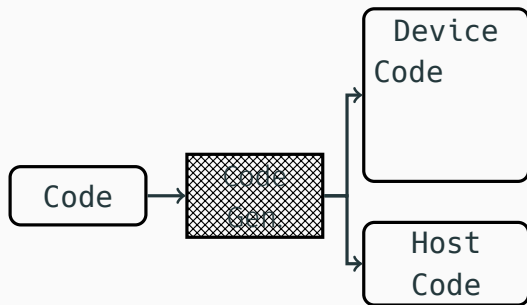
OpenMP

Clang

LLVM-IR

LLVM

Assembly



# OPENMP OFFLOAD — OVERVIEW

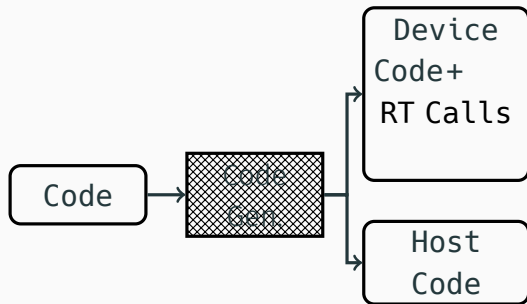
OpenMP

Clang

LLVM-IR

LLVM

Assembly



# OPENMP OFFLOAD — OVERVIEW

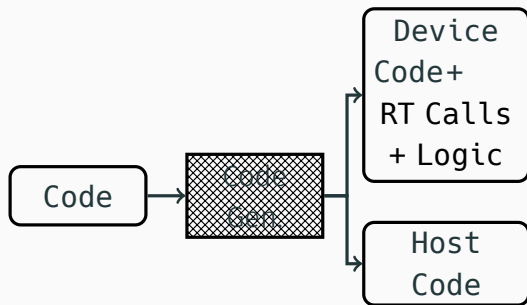
OpenMP

Clang

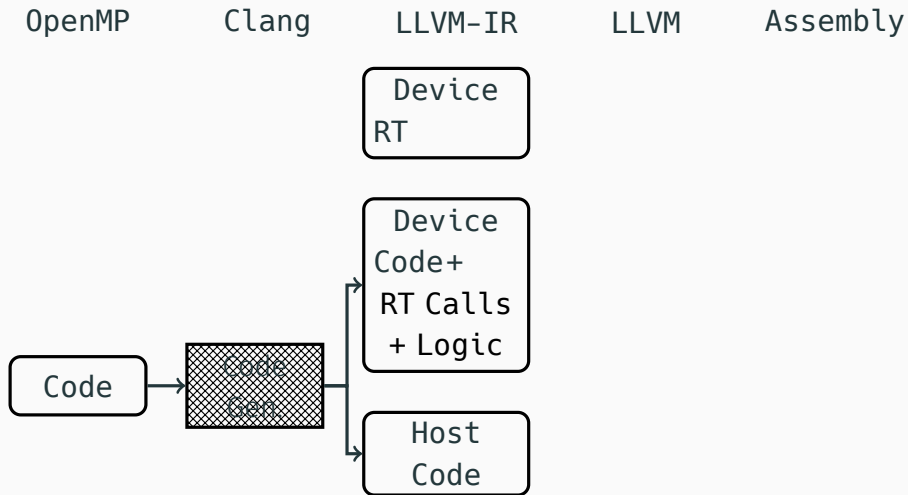
LLVM-IR

LLVM

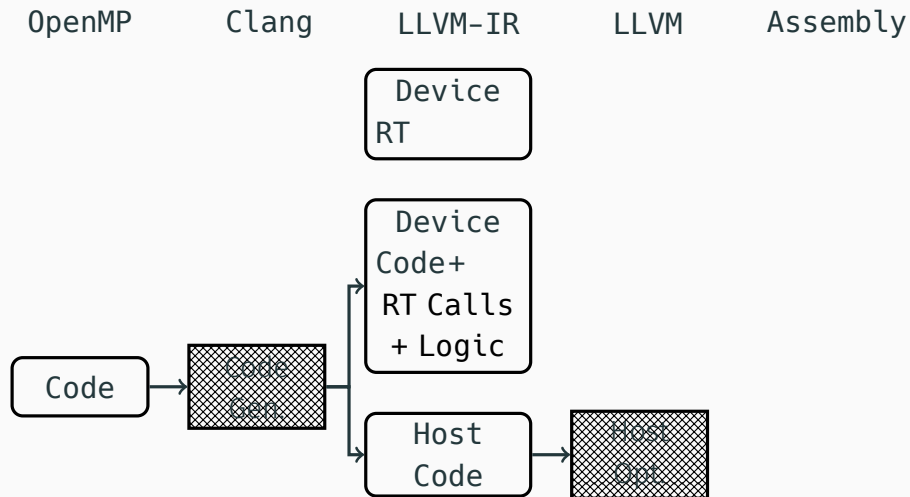
Assembly



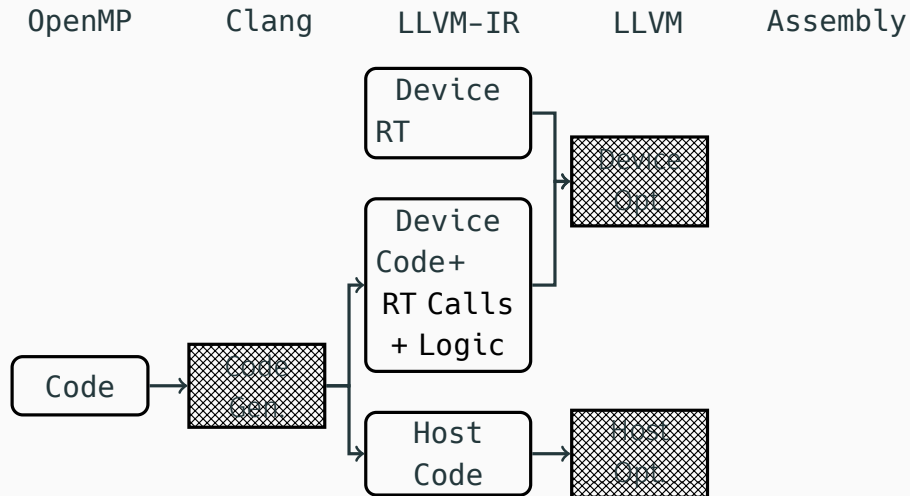
# OPENMP OFFLOAD — OVERVIEW



# OPENMP OFFLOAD — OVERVIEW

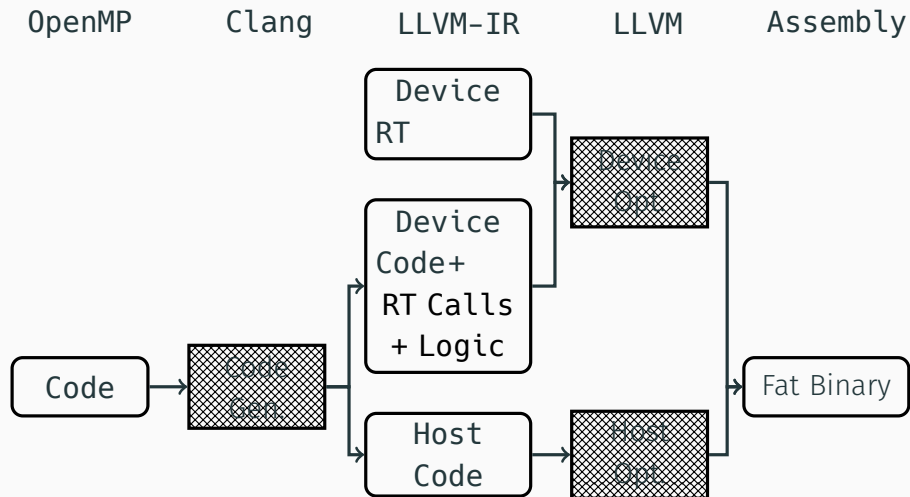


# OPENMP OFFLOAD — OVERVIEW

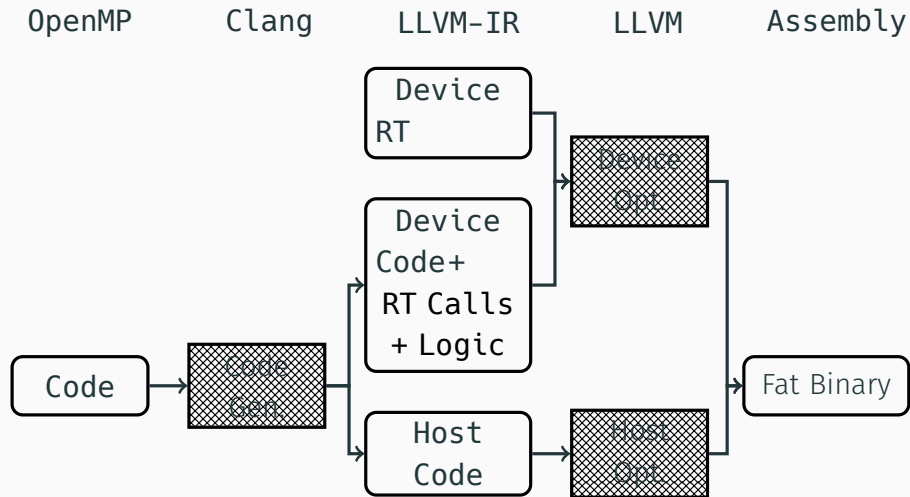




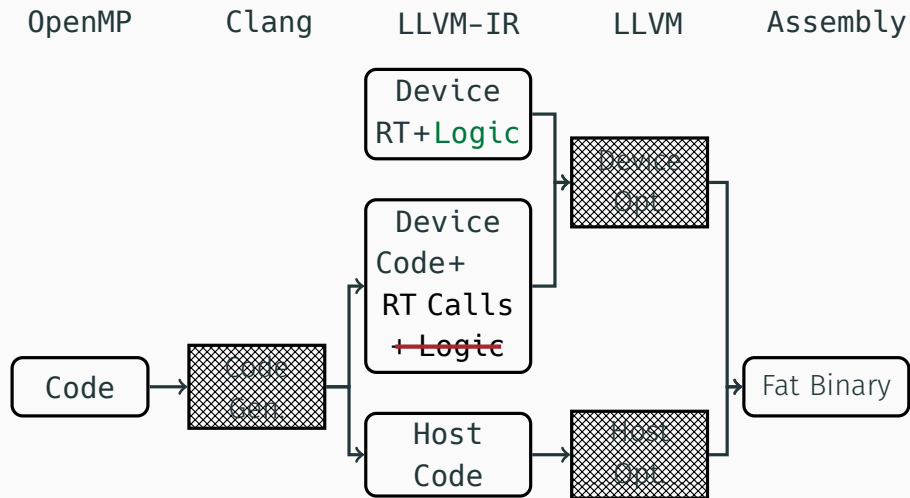
# OPENMP OFFLOAD — OVERVIEW



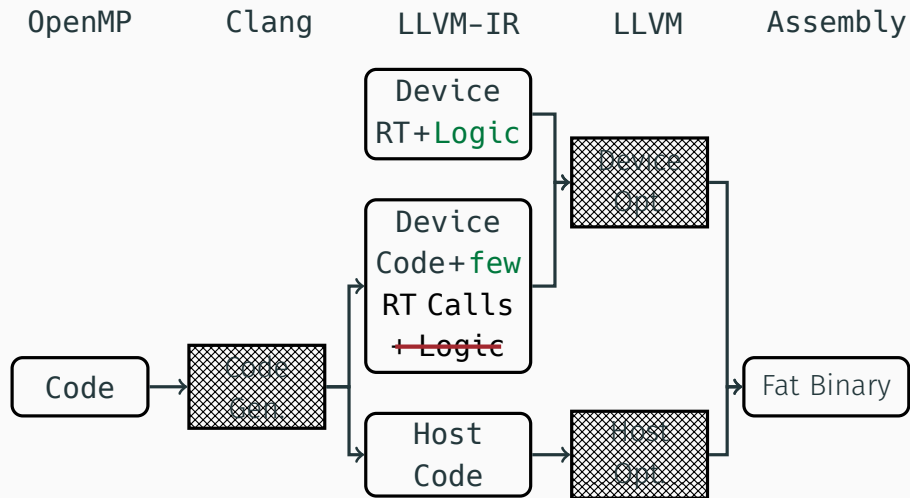
## 1. Offload-Specific Optimizations on Device Code



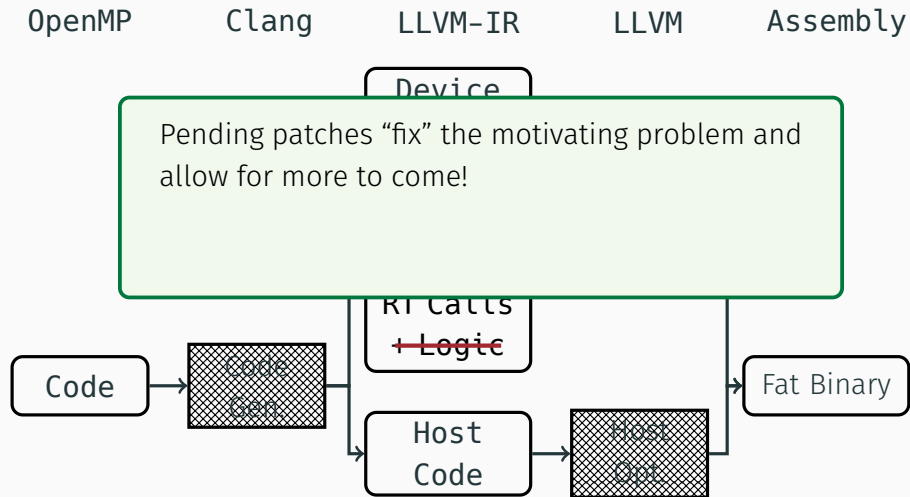
## 1. Offload-Specific Optimizations on Device Code



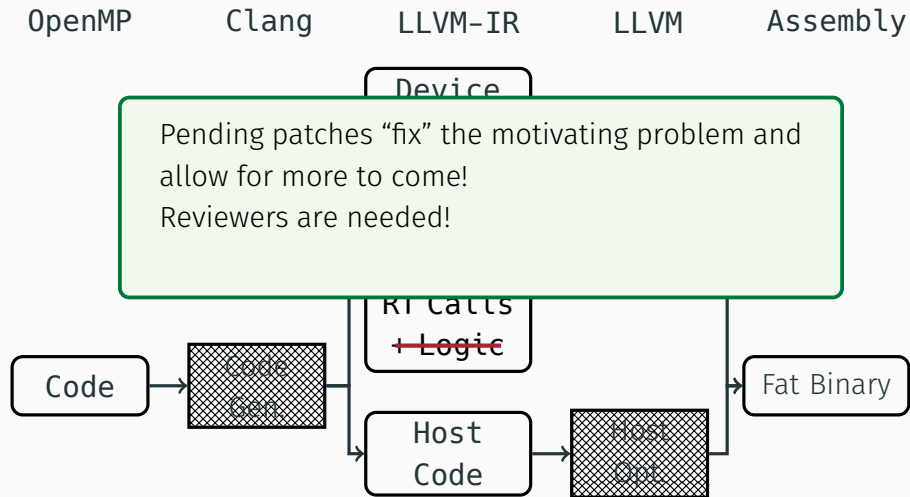
## 1. Offload-Specific Optimizations on Device Code



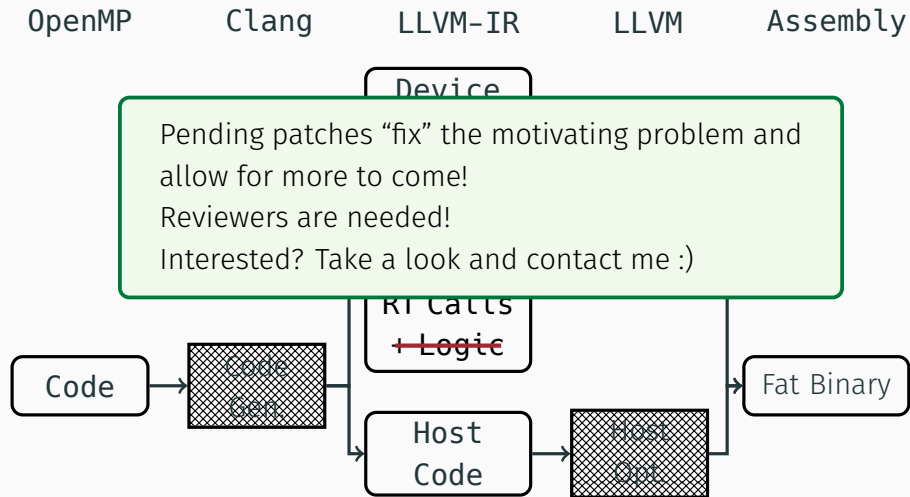
## 1. Offload-Specific Optimizations on Device Code



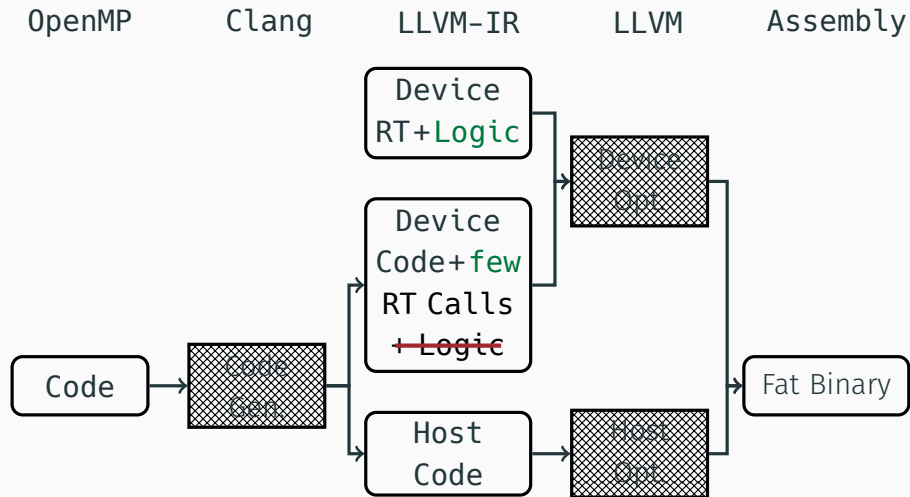
## 1. Offload-Specific Optimizations on Device Code



## 1. Offload-Specific Optimizations on Device Code

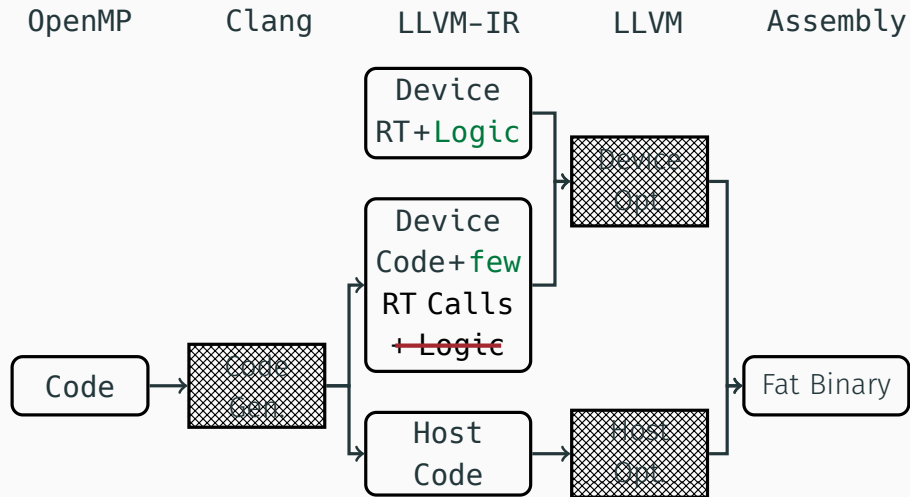


# OPENMP OFFLOAD — OVERVIEW & DIRECTIONS

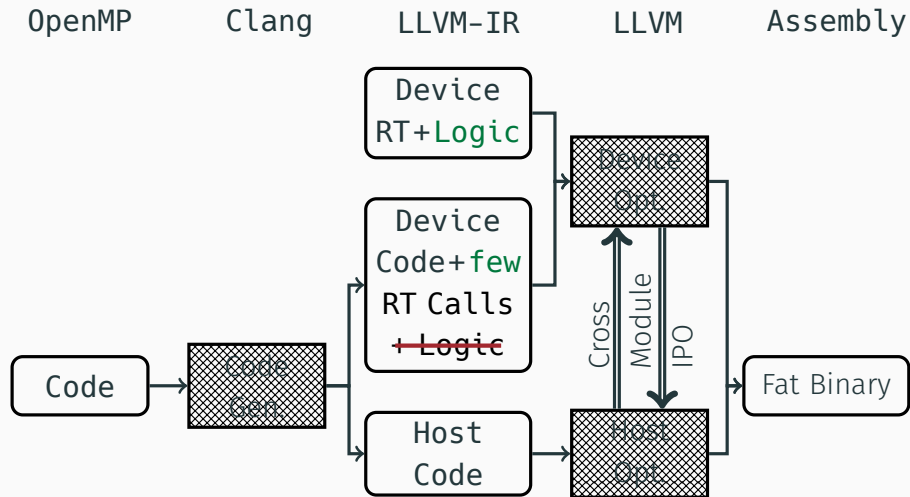




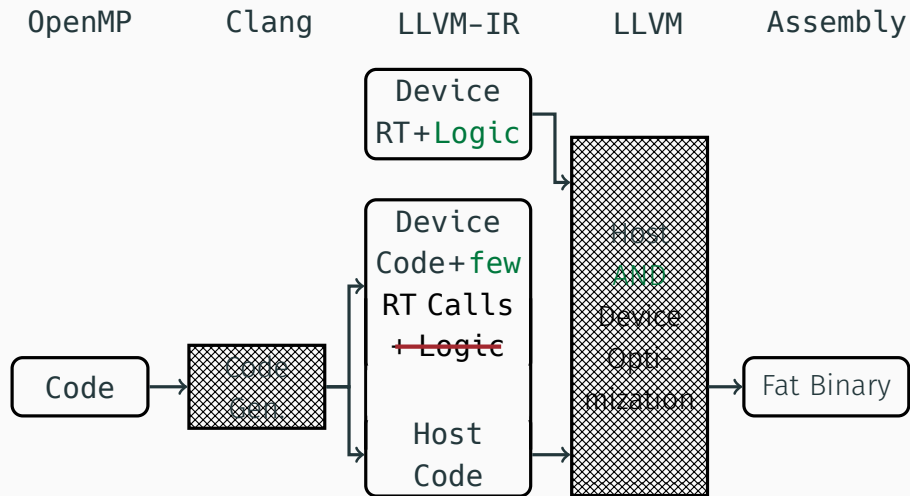
## 2. Optimize Device and Host Code Together



## 2. Optimize Device and Host Code Together



## 2. Optimize Device and Host Code Together





- A straight-forward `#pragma omp target` front-end:
  - ◇ simplified implementation
  - ◇ improved reusability (F18, ...)



- A straight-forward `#pragma omp target` front-end:
  - ◇ simplified implementation CGOpenMPRuntimeNVPTX.cpp ~5.0k loc
  - ◇ improved reusability (F18, ...) CGOpenMPRuntimeTRegion.cpp ~0.5k loc



- A straight-forward `#pragma omp target` front-end:
  - ◇ simplified implementation CGOpenMPRuntimeNVPTX.cpp ~5.0k loc
  - ◇ improved reusability (F18, ...) CGOpenMPRuntimeTRegion.cpp ~0.5k loc
- Interface exposes information and implementation choices:
  - ◇ “smartness” is moved in the compiler middle-end
  - ◇ simplifies analyses and transformations in LLVM



- A straight-forward `#pragma omp target` front-end:
  - ◇ simplified implementation CGOpenMPRuntimeNVPTX.cpp ~5.0k loc
  - ◇ improved reusability (F18, ...) CGOpenMPRuntimeTRegion.cpp ~0.5k loc
- Interface exposes information and implementation choices:
  - ◇ “smartness” is moved in the compiler middle-end
  - ◇ simplifies analyses and transformations in LLVM
- Device RT interface & implementation are separated:
  - ◇ simplifies generated LLVM-IR
  - ◇ most LLVM & Clang parts become target agnostic





# 1. OFFLOAD-SPECIFIC OPTIMIZATIONS — “SPMD-ZATION”

- use inter-procedural reasoning to place minimal guards/synchronization



# 1. OFFLOAD-SPECIFIC OPTIMIZATIONS — “SPMD-ZATION”

- use inter-procedural reasoning to place minimal guards/synchronization
- if legal, switch all boolean **UseSPMDMode** flags to **true**



# 1. OFFLOAD-SPECIFIC OPTIMIZATIONS — “SPMD-ZATION”

- use inter-procedural reasoning to place minimal guards/synchronization
- if legal, switch all boolean **UseSPMDMode** flags to **true**
- currently, no (unknown) global side-effects allowed outside parallel regions.



# 1. OFFLOAD-SPECIFIC OPTIMIZATIONS — CUSTOM STATE MACHINES

- use optimized state-machines when unavoidable



# 1. OFFLOAD-SPECIFIC OPTIMIZATIONS — CUSTOM STATE MACHINES

- use optimized state-machines when unavoidable
- reachability & post-dominance restrict the set of potential next parallel regions to work on



# 1. OFFLOAD-SPECIFIC OPTIMIZATIONS — CUSTOM STATE MACHINES

- use optimized state-machines when unavoidable
- reachability & post-dominance restrict the set of potential next parallel regions to work on
- reuse already communicated/shared values if possible

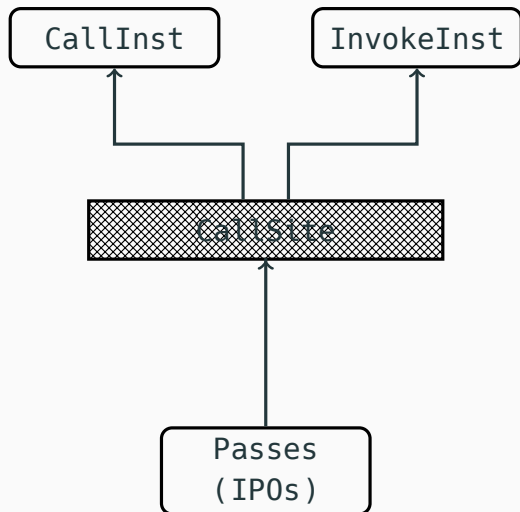


# 1. OFFLOAD-SPECIFIC OPTIMIZATIONS — CUSTOM STATE MACHINES

- use optimized state-machines when unavoidable
- reachability & post-dominance restrict the set of potential next parallel regions to work on
- reuse already communicated/shared values if possible
- currently, a simple state machine is generated with explicit conditionals for all known parallel regions in the module

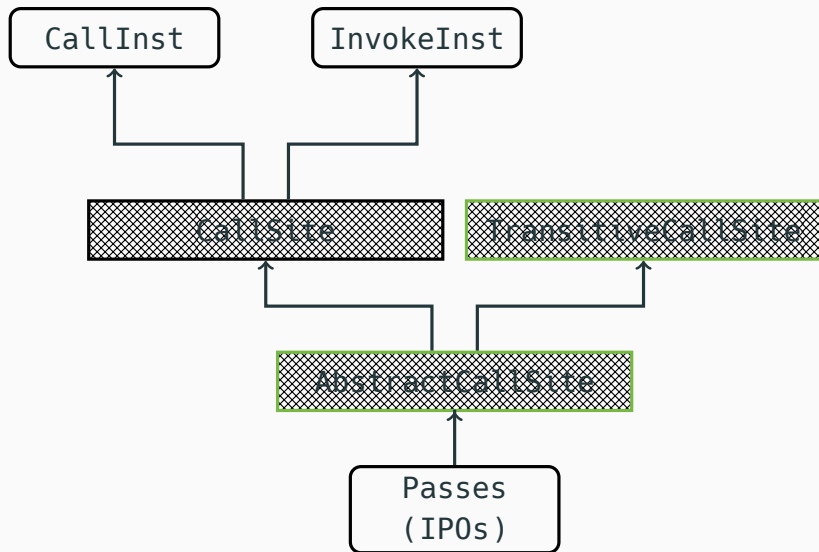


## 2. OPTIMIZE DEVICE AND HOST TOGETHER — ABSTRACT CALL SITES





## 2. OPTIMIZE DEVICE AND HOST TOGETHER — ABSTRACT CALL SITES



## 2. OPTIMIZE DEVICE AND HOST TOGETHER — ABSTRACT CALL SITES

Call

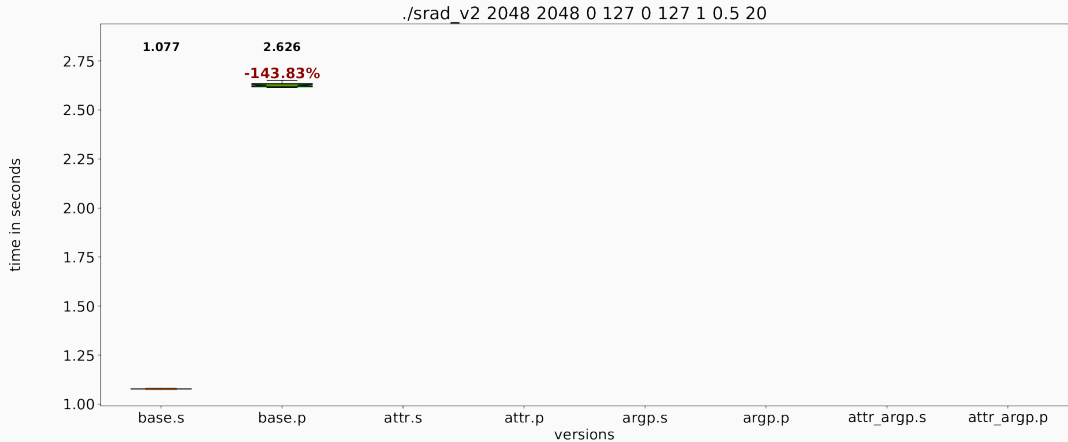
Functional changes required for  
Inter-procedural Constant Propagation:

```
for (int i = 0; i < NumArgs; i++) {  
    Value *ArgOp = ACS.getArgOperand(i);  
    if (!ArgOp) {  
        // handle non-constant  
        continue;  
    }  
    ...  
}
```

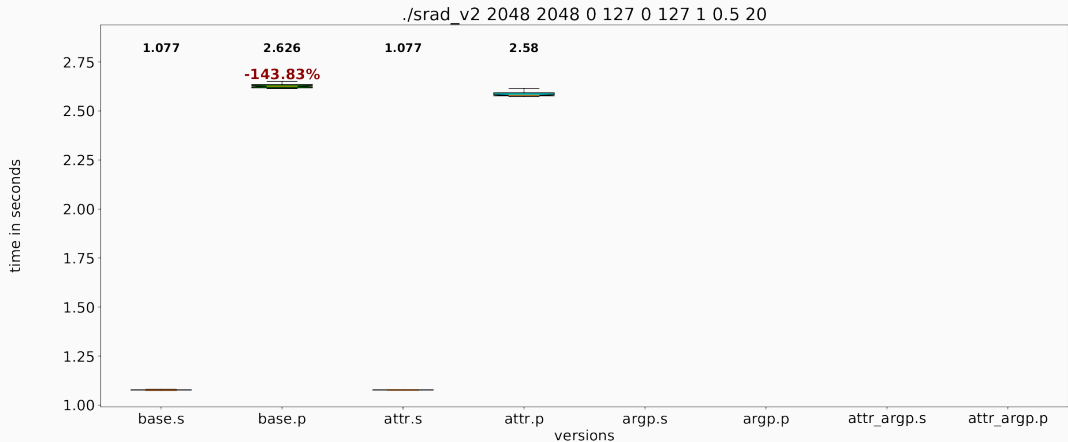
Passes  
(IPOs)



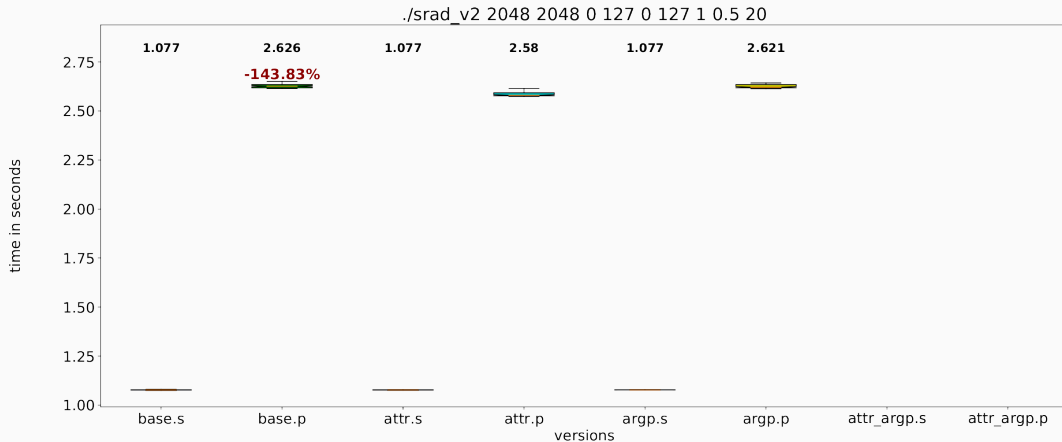
# ABSTRACT CALL SITES — PERFORMANCE RESULTS



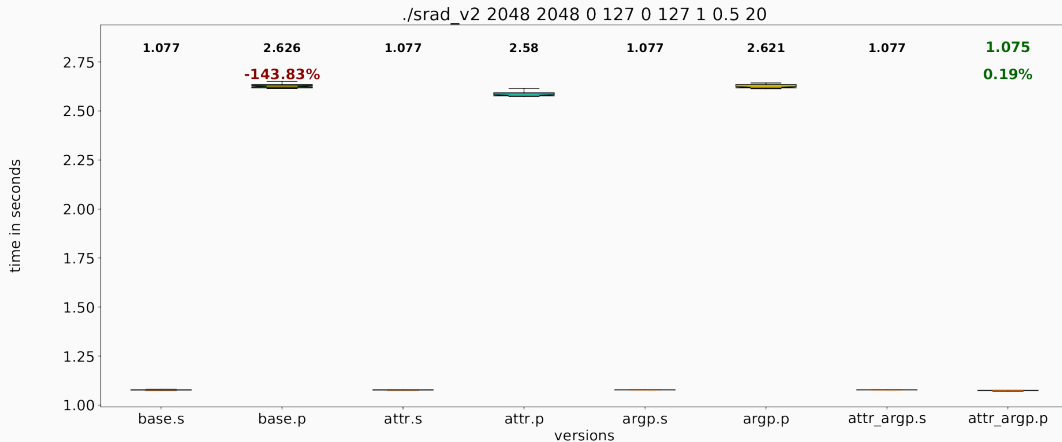
# ABSTRACT CALL SITES — PERFORMANCE RESULTS



# ABSTRACT CALL SITES — PERFORMANCE RESULTS



# ABSTRACT CALL SITES — PERFORMANCE RESULTS



# CONCLUSION



# CONCLUSION

## THE COMPILER BLACK BOX

```
>> clang -O3 -fopenmp-targets=...
```

```
#pragma omp target teams  
{  
  #pragma omp parallel  
  work1();  
}
```

```
#pragma omp parallel  
work2();  
}
```

probably poor performance :(





# CONCLUSION

## THE COMPILER BLACK BOX

```
>> clang -O3 -fopenmp-targets=...
```

```
#pragma omp target teams
{
  #pragma omp parallel
  work1();

  #pragma omp parallel
  work2();
}
```

probably poor performance :(

## OPTIMIZATION CATEGORIES

### Optimizations for *sequential* aspects

- May reuse *existing* transformations (patches up for review!)
- ⇒ Introduce *suitable abstractions* to bridge the indirection (DONE!)

### Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations* (see IWOMP'18<sup>a</sup>)
- ⇒ Introduce a unifying abstraction layer (see EuroLLVM'18<sup>b</sup>)

<sup>a</sup>Compiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

<sup>b</sup>A Parallel IR in Real Life: Optimizing OpenMP, H. Finkel, J. Doerfert, X. Tian, G. Stelle, Euro LLVM Meeting 2018



# CONCLUSION

## THE COMPILER BLACK BOX

```
>> clang -O3 -fopenmp-targets=...
```

```
#pragma omp target teams  
{  
  #pragma omp parallel  
  work1();  
  
  #pragma omp parallel  
  work2();  
}
```

probably poor performance :(

## OPTIMIZATION CATEGORIES

### Optimizations for *sequential* aspects

- May reuse *existing* transformations (patches up for review!)
- ⇒ Introduce *suitable abstractions* to bridge the indirection (DONE!)

### Optimizations for *parallel* aspects

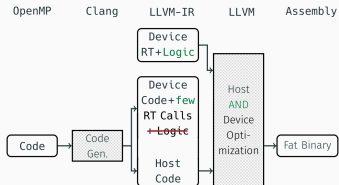
- New *explicit parallelism-aware transformations* (see IWOMP'18<sup>a</sup>)
- ⇒ Introduce a unifying abstraction layer (see EuroLLVM'18<sup>b</sup>)

<sup>a</sup>Compiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

<sup>b</sup>A Parallel IR in Real Life: Optimizing OpenMP, H. Finkel, J. Doerfert, X. Tian, G. Stelle, Euro LLVM Meeting 2018

## OPENMP OFFLOAD — OVERVIEW & DIRECTIONS

### 2. Optimize Device and Host Code Together



# CONCLUSION

## THE COMPILER BLACK BOX

```
>> clang -O3 -fopenmp-targets=...
```

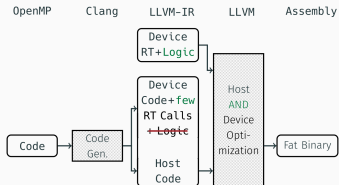
```
#pragma omp target teams  
{  
  #pragma omp parallel  
  work1();  
  
  #pragma omp parallel  
  work2();  
}
```

probably poor performance :(



## OPENMP OFFLOAD — OVERVIEW & DIRECTIONS

### 2. Optimize Device and Host Code Together



## OPTIMIZATION CATEGORIES

### Optimizations for *sequential* aspects

- May reuse *existing* transformations (patches up for review)
- ⇒ Introduce *suitable abstractions* to bridge the indirection (DONE)

### Optimizations for *parallel* aspects

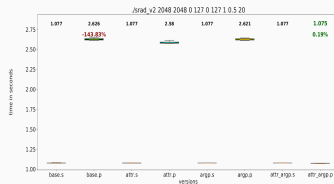
- New *explicit parallelism-aware transformations* (see IWOMP'18<sup>a</sup>)
- ⇒ Introduce a unifying abstraction layer (see EuroLLVM'18<sup>b</sup>)

<sup>a</sup>Compiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

<sup>b</sup>A Parallel IR in Real Life: Optimizing OpenMP, H. Finkel, J. Doerfert, X. Tian, G. Stelle, Euro LLVM Meeting 2018



## ABSTRACT CALL SITES — PERFORMANCE RESULTS







I: Attribute Propagation — Bidirectional Information Transfer:  
read/write-only, **restrict**/**noalias**, ...



I: Attribute Propagation — Bidirectional Information Transfer:

`read/write-only`, **`restrict`**/**`noalias`**, ...

II: Variable Privatization — Limit Variable Lifetimes:

`shared(var)`  $\rightarrow$  `firstprivate(var)`  $\rightarrow$  `private(var)`



I: Attribute Propagation — Bidirectional Information Transfer:  
`read/write-only`, **`restrict/noalias`**, ...

II: Variable Privatization — Limit Variable Lifetimes:  
`shared(var)` → `firstprivate(var)` → `private(var)`

III: Parallel Region Expansion — Maximize Parallel Contexts:  
⇒ reduce start/stop overheads and expose barriers





I: Attribute Propagation — Bidirectional Information Transfer:  
`read/write-only`, **`restrict/noalias`**, ...

II: Variable Privatization — Limit Variable Lifetimes:  
`shared(var)` → `firstprivate(var)` → `private(var)`

III: Parallel Region Expansion — Maximize Parallel Contexts:  
⇒ reduce start/stop overheads and expose barriers

IV: Barrier Elimination — Eliminate Redundant Barriers



- I: Attribute Propagation — Bidirectional Information Transfer:  
`read/write-only`, **`restrict/noalias`**, ...
- II: Variable Privatization — Limit Variable Lifetimes:  
`shared(var)` → `firstprivate(var)` → `private(var)`
- III: Parallel Region Expansion — Maximize Parallel Contexts:  
⇒ reduce start/stop overheads and expose barriers
- IV: Barrier Elimination — Eliminate Redundant Barriers
- V: Communication Optimization — Move Computations Around: seq. compute & result comm. **vs.** operand comm. & par. compute



I: Attribute Propagation — Bidirectional Information Transfer:  
`read/write-only`, **`restrict/noalias`**, ...

II: Variable Privatization — Limit Variable Lifetimes:  
`shared(var)` → `firstprivate(var)` → `private(var)`

III: Parallel Region Expansion — Maximize Parallel Contexts:  
⇒ reduce start/stop overheads and expose barriers

IV: Barrier Elimination — Eliminate Redundant Barriers

V: Communication Optimization — Move Computations Around: seq. compute & result comm. vs. operand comm. & par. compute



- I: Attribute Propagation — In LLVM: *Attribute Deduction (IPO!)*  
read/write-only, **restrict/noalias**, ...
- II: Variable Privatization — In LLVM: *Argument Promotion (IPO!)*  
**shared(var) → firstprivate(var) → private(var)**
- III: Parallel Region Expansion — Maximize Parallel Contexts:  
⇒ reduce start/stop overheads and expose barriers
- IV: Barrier Elimination — Eliminate Redundant Barriers
- V: Communication Optimization — Move Computations Around: seq. compute & result comm. vs. operand comm. & par. compute



```
#pragma omp parallel for  
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---



```
OpenMP Input: #pragma omp parallel for  
for (int i = 0; i < N; i++)  
    Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.  
omp_rt_parallel_for(0, N, &body_fn, &N, &In, &Out);
```



## EARLY OUTLINING

```
#pragma omp parallel for
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.
```

```
omp_rt_parallel_for(0, N, &body_fn, &N, &In, &Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
static void body_fn(int tid, int *N,  
                    float** In, float** Out) {  
    int lb = omp_get_lb(tid), ub = omp_get_ub(tid);  
    for (int i = lb; i < ub; i++)  
        (*Out)[i] = (*In)[i] + (*In)[i + (*N)]  
}
```



## EARLY OUTLINING

```
#pragma omp parallel for
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.
```

```
omp_rt_parallel_for(0, N, &body_fn, &N, &In, &Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
static void body_fn(int tid, int* N,  
                   float** In, float** Out) {  
    int lb = omp_get_lb(tid), ub = omp_get_ub(tid);  
    for (int i = lb; i < ub; i++)  
        (*Out)[i] = (*In)[i] + (*In)[i + (*N)]  
}
```





## AN ABSTRACT PARALLEL IR

```
#pragma omp parallel for
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by an annotated loop
```

```
for /* parallel */ (int i = 0; i < N; i++)  
    body_fn(i, &N, &In, &Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
static void body_fn(int i, int* N,  
                    float** In, float** Out) {
```

```
    (*Out)[i] = (*In)[i] + (*In)[i + (*N)]
```

```
}
```



## EARLY OUTLINED + TRANSITIVE CALLS

```
#pragma omp parallel for
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.
```

```
omp_rt_parallel_for(0, N, &body_fn, &N, &In, &Out);
```

```
// Model transitive call: body_fn(?, &N, &In, &Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
static void body_fn(int tid, int* N,  
                   float** In, float** Out) {  
    int lb = omp_get_lb(tid), ub = omp_get_ub(tid);  
    for (int i = lb; i < ub; i++)  
        (*Out)[i] = (*In)[i] + (*In)[i + (*N)]  
}
```



## EARLY OUTLINED + TRANSITIVE CALLS

```
#pragma omp parallel for
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.
```

```
omp_rt_parallel_for(0, N, &body_fn, &N, &In, &Out);
```

```
// Model transitive call: body_fn(?, &N, &In, &Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
static void body_fn(int tid, int* N,  
                   float** In, float** Out) {
```

```
    int lb = 0;
```

```
    for (int i
```

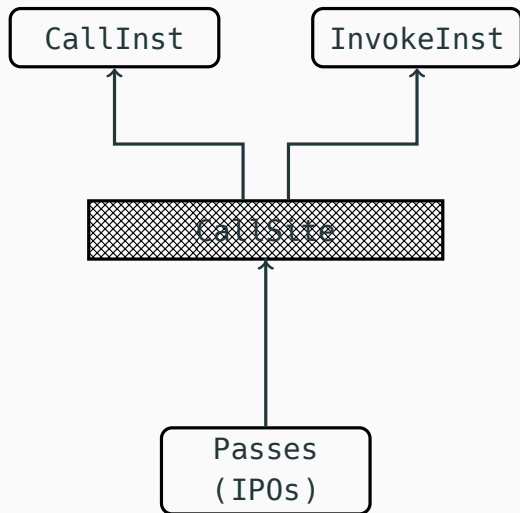
```
        (*Out)[i
```

```
    }
```

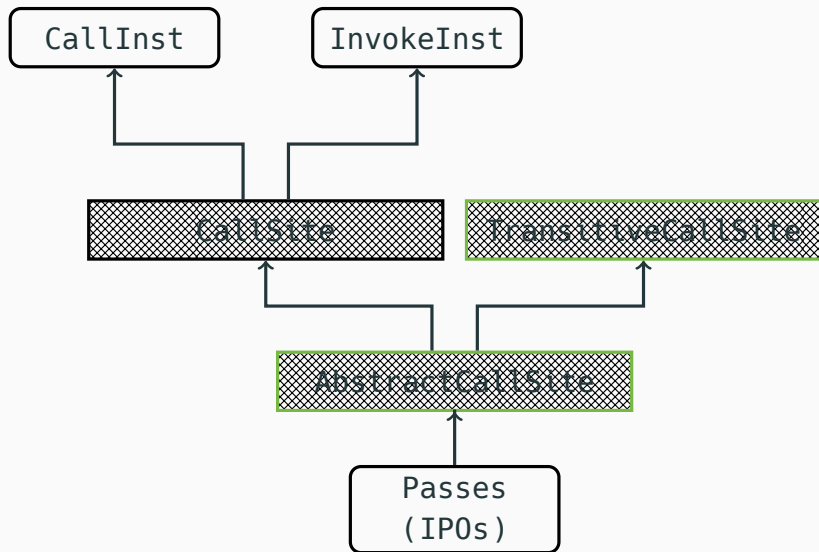
- + valid and executable IR
- integration cost per IPO
- + no unintended interactions



## IPO IN LLVM



## TRANSITIVE CALL SITES IN LLVM



Call

Functional changes required for  
Inter-procedural Constant Propagation:

```
for (int i = 0; i < NumArgs; i++) {  
    Value *ArgOp = ACS.getArgOperand(i);  
    if (!ArgOp) {  
        // handle non-constant  
        continue;  
    }  
    ...  
}
```

(IPOs)



Version	Description	Opt.
<i>base</i>	plain “-O3”, thus no parallel optimizations	
<i>attr</i>	attribute propagation through attr. deduction (IPO)	I
<i>argp</i>	variable privatization through arg. promotion (IPO)	II
<i>n/a</i>	constant propagation (IPO)	







## SOME CONTEXT



### Examples

Examples are given in a C-like language with OpenMP annotations.

### Transformations

Our transformations work on the LLVM intermediate representation (LLVM-IR), thus take and produce LLVM-IR.

### OpenMP Runtime Library

We experience OpenMP annotations as OpenMP runtime library calls and the situation is most often more complicated than presented here.



- Run with 1 Thread<sup>2</sup>
- Median and variance of 51 runs is shown
- Rodiana 3.1 benchmarks and LULESH v1.0 (OpenMP)
- Only time in parallel constructs was measured

---

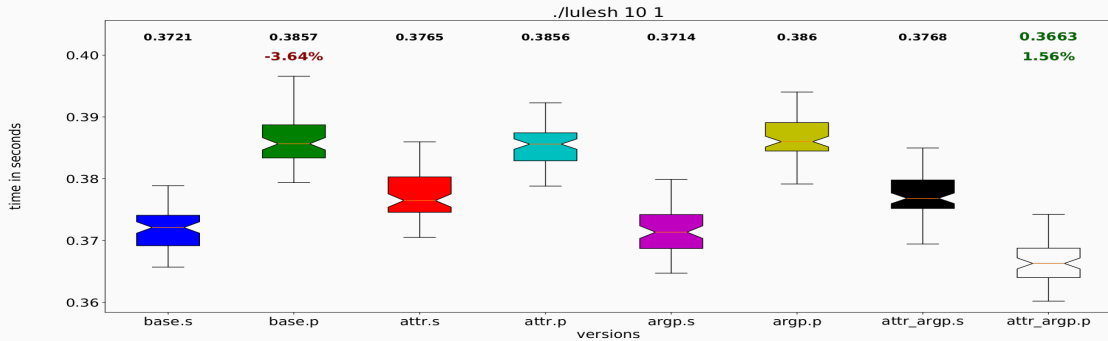
<sup>2</sup>Intel(R) Core(TM) i7-4800MQ CPU @ 2.70GHz



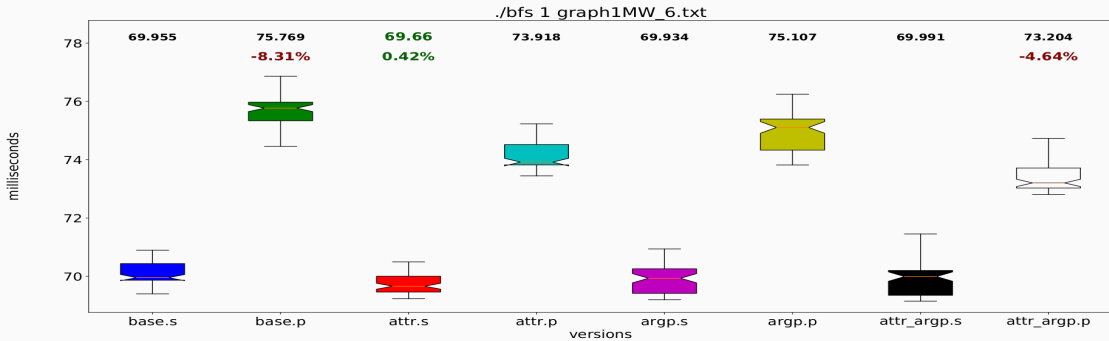
# PERFORMANCE RESULTS



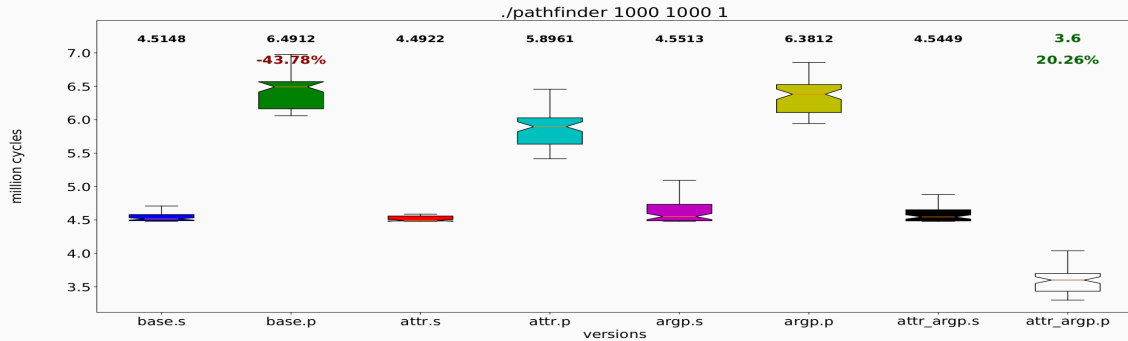
# PERFORMANCE RESULTS



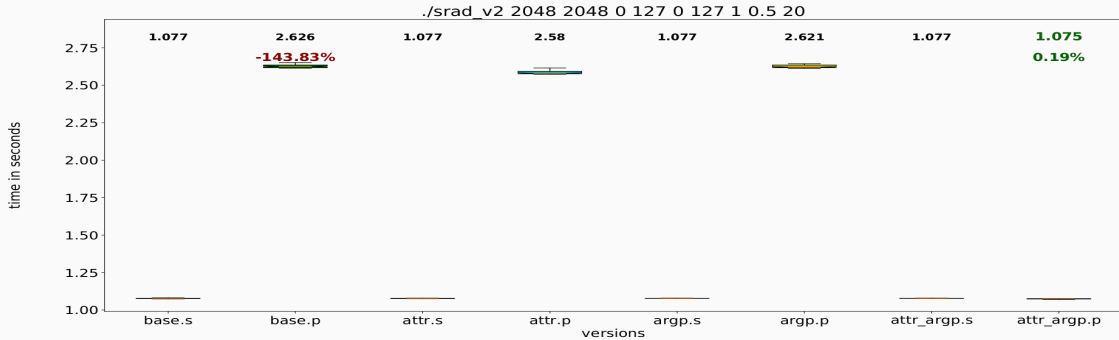
# PERFORMANCE RESULTS



# PERFORMANCE RESULTS



# PERFORMANCE RESULTS





# ACTION ITEM I

---

\*

†



1) Run *your* OpenMP code sequentially<sup>\*</sup>,  
with and without OpenMP.

---

<sup>\*</sup>export OMP\_NUM\_THREADS=1  
<sup>†</sup>



1) Run *your* OpenMP code sequentially<sup>\*</sup>,  
with and without OpenMP.

2) Email me<sup>†</sup> the results!

---

<sup>\*</sup> `export OMP_NUM_THREADS=1`

<sup>†</sup> `jdoerfert@anl.gov`



# ACTION ITEM II

---

\*



1) Always\* use `default(none)` and `firstprivate(...)`

---

\* For scalars/pointers if you do not have explicit synchronization.



1) Always\* use `default(none)` and `firstprivate(...)`

2) Revisit ACTION ITEM I

---

\*For scalars/pointers if you do not have explicit synchronization.



NO need to “share” the *variable A!*

```
#pragma omp parallel for shared(A)  
for (int i = 0; i < N; i++)  
    A[i] = i;
```

---

\*For scalars/pointers if you do not have explicit synchronization.



## CONSTANT PROPAGATION EXAMPLE

```
double gamma[4][8];
```

```
gamma[0][0] = 1;
```

```
// ... and so on till ...
```

```
gamma[3][7] = -1;
```

```
Kokkos::parallel_for(
```

```
    "CalcFBHourglassForceForElems A",
```

```
    numElem, KOKKOS_LAMBDA(const int &i2) {
```

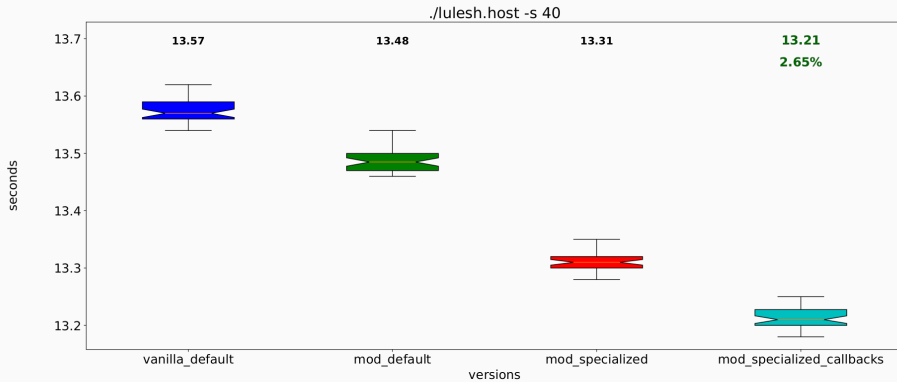
```
        // Use gamma[0][0] ... game[3][7]
```

```
    }
```





# CONSTANT PROPAGATION PERFORMANCE



## OPTIMIZATION I: ATTRIBUTE PROPAGATION

```
#pragma omp parallel for firstprivate(...)
```

OpenMP Input: `for (int i = 0; i < N; i++)`  
`Out[i] = In[i] + In[i+N];`

---



## OPTIMIZATION I: ATTRIBUTE PROPAGATION

```
#pragma omp parallel for firstprivate(...)
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.  
omp_rt_parallel_for(0, N, &body_fn, N, In, Out);
```

```
// Parallel region outlined in the front-end (clang)!  
void body_fn(int i, int N,  
             float* In,  
             float* Out) {  
    Out[i] = In[i] + In[i + N];  
}
```



## OPTIMIZATION I: ATTRIBUTE PROPAGATION

```
#pragma omp parallel for firstprivate(...)
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.  
omp_rt_parallel_for(0, N, &body_fn, N, In, Out);
```

```
// Parallel region outlined in the front-end (clang)!  
void body_fn(int i, int N,  
             float* /* read-only & no-escape */ In,  
             float* /* write-only & no-escape */ Out) {  
    Out[i] = In[i] + In[i + N];  
}
```



## OPTIMIZATION I: ATTRIBUTE PROPAGATION

```
#pragma omp parallel for firstprivate(...)
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
               Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.
```

```
omp_rt_parallel_for(0, N, &body_fn, N,  
    /* ro & no-esc */ In, /* wo & no-esc */ Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
void body_fn(int i, int N,  
    float* /* read-only & no-escape */ In,  
    float* /* write-only & no-escape */ Out) {  
    Out[i] = In[i] + In[i + N];  
}
```



```
int foo() {  
    int a = 0;  
  
    #pragma omp parallel  
    {  
        #pragma omp critical  
        { a += 1; }  
        bar();  
        #pragma omp critical  
        { a *= 2; }  
    }  
    return a;  
}
```

```
int foo() {  
    int a = 0;  
  
#pragma omp parallel  
    {  
#pragma omp critical  
        { a += 1; }  
        bar();  
#pragma omp critical  
        { a *= 2; }  
    }  
    return a;  
}
```

```
int foo() {  
    int a = 0;  
    int *restrict p = &a;  
    omp_rt_parallel_for(pwork, p);  
    return a;  
}  
void pwork(int tid, int *p) {  
    if (omp_critical(tid)) {  
        *p = *p + 1;  
        omp_critical_end(tid);  
    }  
    bar();  
    if (omp_critical(tid)) {  
        *p = *p * 2;  
        omp_critical_end(tid);  
    }  
}
```

```
void pwork(int tid,
           int *restrict p) {
    if (omp_critical(tid)) {
        omp_critical_end(tid);
    }
    bar();
    if (omp_critical(tid)) {
        *p = 2 * (*p + 1);
        omp_critical_end(tid);
    }
}
```

```
int foo() {
    int a = 0;
    int *restrict p = &a;
    omp_rt_parallel_for(pwork, p);
    return a;
}

void pwork(int tid, int *p) {
    if (omp_critical(tid)) {
        *p = *p + 1;
        omp_critical_end(tid);
    }
    bar();
    if (omp_critical(tid)) {
        *p = *p * 2;
        omp_critical_end(tid);
    }
}
```



```
void pwork(int tid,
           int *restrict p) {
    if (omp_critical(tid)) {
        *p = *p + 1;
        omp_critical_end(tid);
    }
    bar()[p]; // May "use" p.
    if (omp_critical(tid)) {
        *p = *p * 2;
        omp_critical_end(tid);
    }
}
```

```
int foo() {
    int a = 0;
    int *restrict p = &a;
    omp_rt_parallel_for(pwork, p);
    return a;
}

void pwork(int tid, int *p) {
    if (omp_critical(tid)) {
        *p = *p + 1;
        omp_critical_end(tid);
    }
    bar();
    if (omp_critical(tid)) {
        *p = *p * 2;
        omp_critical_end(tid);
    }
}
```

## OPTIMIZATION II: VARIABLE PRIVATIZATION

```
#pragma omp parallel for shared(...)
```

OpenMP Input: `for (int i = 0; i < N; i++)`  
`Out[i] = In[i] + In[i+N];`

---



## OPTIMIZATION II: VARIABLE PRIVATIZATION

```
#pragma omp parallel for shared(...)
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.
```

```
omp_rt_parallel_for(0, N, &body_fn, &N, &In, &Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
void body_fn(int i, int* N,  
             float** In,  
             float** Out) {  
    (*Out)[i] = (*In)[i] + (*In)[i + (*N)];  
}
```



## OPTIMIZATION II: VARIABLE PRIVATIZATION

```
#pragma omp parallel for shared(...)
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.
```

```
omp_rt_parallel_for(0, N, &body_fn, &N, &In, &Out);
```

```
// Parallel region outlined in the front-end (clang)!
```

```
void body_fn(int i, int* /* ro & ne */ N,  
             float** /* ro & ne */ In,  
             float** /* ro & ne */ Out) {  
    (*Out)[i] = (*In)[i] + (*In)[i + (*N)];  
}
```



## OPTIMIZATION II: VARIABLE PRIVATIZATION

```
#pragma omp parallel for firstprivate(...)
```

```
OpenMP Input: for (int i = 0; i < N; i++)  
                Out[i] = In[i] + In[i+N];
```

---

```
// Parallel region replaced by a runtime call.  
omp_rt_parallel_for(0, N, &body_fn, N, In, Out);
```

```
// Parallel region outlined in the front-end (clang)!  
void body_fn(int i, int N,  
             float* In,  
             float* Out) {  
    Out[i] = In[i] + In[i + N];  
}
```



## OPTIMIZATION III: PARALLEL REGION EXPANSION



## OPTIMIZATION III: PARALLEL REGION EXPANSION

```
void copy(float* dst, float* src, int N) {  
    #pragma omp parallel for  
    for(int i = 0; i < N; i++) {  
        dst[i] = src[i];  
    } // implicit barrier!  
}
```

```
void compute_step_factor(int nelr, float* vars,  
                        float* areas, float* sf) {  
    #pragma omp parallel for  
    for (int blk = 0; blk < nelr / block_length; ++blk) {  
        ...  
    } // implicit barrier!  
}
```



## OPTIMIZATION III: PARALLEL REGION EXPANSION

```
for (int i = 0; i < iterations; i++) {  
    copy(old_vars, vars, nelr * NVAR);  
  
    compute_step_factor(nelr, vars, areas, sf);  
  
    for (int j = 0; j < RK; j++) {  
        compute_flux(nelr, ese, normals, vars, fluxes, ff_vars,  
                    ff_m_x, ff_m_y, ff_m_z, ff_denergy);  
  
        time_step(j, nelr, old_vars, vars, sf, fluxes);  
    }  
}
```






## OPTIMIZATION III: PARALLEL REGION EXPANSION

```
for (int i = 0; i < iterations; i++) {  
    #pragma omp parallel for          // copy  
    for (...) {  
        /* write old_vars, read vars */  
    } // implicit barrier!  
    compute_step_factor(nelr, vars, areas, sf);  
  
    for (int j = 0; j < RK; j++) {  
        compute_flux(nelr, ese, normals, vars, fluxes, ff_vars,  
                    ff_m_x, ff_m_y, ff_m_z, ff_denergy);  
  
        time_step(j, nelr, old_vars, vars, sf, fluxes);  
    }  
}
```



## OPTIMIZATION III: PARALLEL REGION EXPANSION

```
for (int i = 0; i < iterations; i++) {  
    #pragma omp parallel for           // copy  
    for (...) {  
        /* write old_vars, read vars */  
    } // implicit barrier!  
    #pragma omp parallel for         // compute_step_factor  
    for (...) {  
        /* write sf, read vars & area */  
    } // implicit barrier!  
    for (int j = 0; j < RK; j++) {  
        #pragma omp parallel for     // compute_flux  
        for (...) {  
            /* write fluxes, read vars & ... */  
        } // implicit barrier!  
        ...  
    }  
}
```



## OPTIMIZATION III: PARALLEL REGION EXPANSION

```
#pragma omp parallel
for (int i = 0; i < iterations; i++) {
    #pragma omp for                // copy
    for (...) {
        /* write old_vars, read vars */
    } // explicit barrier in LLVM-IR!
    #pragma omp for                // compute_step_factor
    for (...) {
        /* write sf, read vars & area */
    } // explicit barrier in LLVM-IR!
    for (int j = 0; j < RK; j++) {
        #pragma omp for            // compute_flux
        for (...) {
            /* write fluxes, read vars & ... */
        } // explicit barrier in LLVM-IR!
        ...
    }
}
```



## OPTIMIZATION IV: BARRIER ELIMINATION

```
#pragma omp parallel
for (int i = 0; i < iterations; i++) {
  #pragma omp for                                // copy
  for (...) {
    /* write old_vars, read vars */
  } // explicit barrier in LLVM-IR!
  #pragma omp for                                // compute_step_factor
  for (...) {
    /* write sf, read vars & area */
  } // explicit barrier in LLVM-IR!
  for (int j = 0; j < RK; j++) {
    #pragma omp for                                // compute_flux
    for (...) {
      /* write fluxes, read vars & ... */
    } // explicit barrier in LLVM-IR!
    ...
  }
}
```



## OPTIMIZATION IV: BARRIER ELIMINATION

```
#pragma omp parallel
for (int i = 0; i < iterations; i++) {
    #pragma omp for // copy
    for (...) {
        /* write old_vars, read vars */
    } // explicit barrier in LLVM-IR!
    #pragma omp for // compute_step_factor
    for (...) {
        /* write sf, read vars & area */
    } // explicit barrier in LLVM-IR!
    for (int j = 0; j < RK; j++) {
        #pragma omp for // compute_flux
        for (...) {
            /* write fluxes, read vars & ... */
        } // explicit barrier in LLVM-IR!
    }
    ...
}
```



## OPTIMIZATION IV: BARRIER ELIMINATION

```
#pragma omp parallel
for (int i = 0; i < iterations; i++) {
  #pragma omp for nowait // copy
  for (...) {
    /* write old_vars, read vars */
  }
  #pragma omp for nowait // compute_step_factor
  for (...) {
    /* write sf, read vars & area */
  }
  for (int j = 0; j < RK; j++) {
    #pragma omp for // compute_flux
    for (...) {
      /* write fluxes, read vars & ... */
    } // explicit barrier in LLVM-IR!
  }
  ...
}
```



# OPTIMIZATION V: COMMUNICATION OPTIMIZATION



## OPTIMIZATION V: COMMUNICATION OPTIMIZATION

```
void f(int *X, int *restrict Y) {  
    int L = *X;          // immovable  
    int N = 512;        // movable  
  
    int A = N + L;      // movable  
    #pragma omp parallel for |  
        firstprivate(X, Y, N, L, A)  
    for (int i = 0; i < N; i++) {  
        int K = *Y;     // movable  
        int M = N * K;  // movable  
        X[i] = M+A*L*i; // immovable  
    }  
}
```





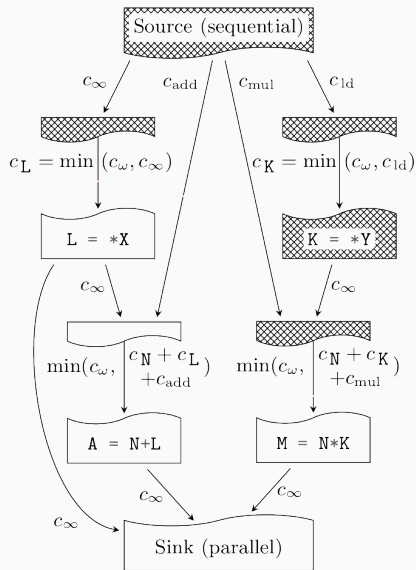


# OPTIMIZATION V: COMMUNICATION OPTIMIZATION

```
void f(int *X, int *restrict Y) {
    int L = *X;           // immovable
    int N = 512;         // movable
```

```
    int A = N + L;      // movable
    #pragma omp parallel for firstprivate(X, Y, N, L, A)
    for (int i = 0; i < N; i++) {
        int K = *Y;     // movable
        int M = N * K;  // movable
        X[i] = M + A * L * i; // immovable
    }
}
```




 $c_\infty = \infty$     $c_{\text{add}} = 5$     $c_\omega = 15$     $c_{\text{mul}} = 10$   
 $c_{\text{ld}} = 20$     $c_N = c_{\text{cst}}$     $c_{\text{cst}} = 0$     cut

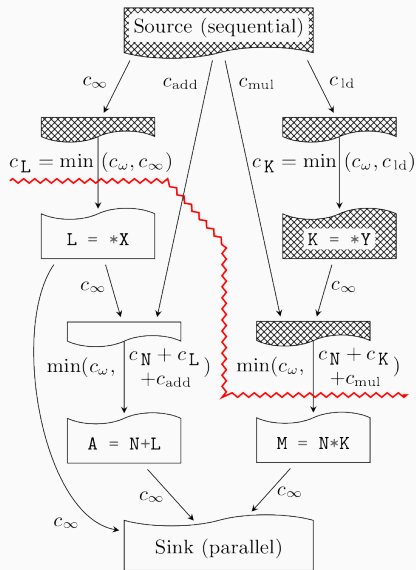


# OPTIMIZATION V: COMMUNICATION OPTIMIZATION

```
void f(int *X, int *restrict Y) {
    int L = *X;           // immovable
    int N = 512;         // movable
```

```
    int A = N + L;      // movable
    #pragma omp parallel for firstprivate(X, Y, N, L, A)
    for (int i = 0; i < N; i++) {
        int K = *Y;    // movable
        int M = N * K; // movable
        X[i] = M + A * L * i; // immovable
    }
}
```


 $c_\infty = \infty$     $c_{add} = 5$     $c_\omega = 15$     $c_{mul} = 10$   
 $c_{ld} = 20$     $c_N = c_{cst}$     $c_{cst} = 0$     cut



## OPTIMIZATION V: COMMUNICATION OPTIMIZATION

```
void f(int *X, int *restrict Y) {
    int L = *X;          // immovable
    int N = 512;        // movable

    int A = N + L;      // movable
    #pragma omp parallel for |
        firstprivate(X, Y, N, L, A)
    for (int i = 0; i < N; i++) {
        int K = *Y;    // movable
        int M = N * K; // movable
        X[i] = M+A*L*i; // immovable
    }
}
```



```
void g(int *X, int *restrict Y) {
    int L = *X;          // immovable
    int K = *Y;         //  $c_{ld} > c_{\omega}$ 
    int M = 512 * K;    //  $c_{mul} + c_K > c_{\omega}$ 
    #pragma omp parallel |
        firstprivate(X, M, L)
    {
        int A = 512 + L; //  $c_{add} < c_{\omega}$ 
        #pragma omp for |
            firstprivate(X, M, A, L)
        for (int i = 0; i < 512; i++) {

            X[i] = M+A*L*i; // immovable
        }
    }
}
```

Information Transfer

Value Transfer



**NO** Information Transfer:  
*outlined function*  $\iff$  *runtime library call site*

Value Transfer



**NO** Information Transfer:  
*outlined function*  $\iff$  *runtime library call site*

Declaration	Value Transfer OpenMP Clause	Communication Type
T var;	<i>default = shared</i>	
T var;	<b>shared</b> (var)	
T var;	<b>lastprivate</b> (var)	



**NO** Information Transfer:  
*outlined function*  $\iff$  *runtime library call site*

Declaration	Value Transfer OpenMP Clause	Communication Type
T var;	<i>default = shared</i>	&var of type T*
T var;	<b>shared</b> (var)	&var of type T*
T var;	<b>lastprivate</b> (var)	&var of type T*



**NO** Information Transfer:  
*outlined function*  $\iff$  *runtime library call site*

Declaration	Value Transfer OpenMP Clause	Communication Type
T var;	<i>default = shared</i>	&var of type T*
T var;	<b>shared</b> (var)	&var of type T*
T var;	<b>lastprivate</b> (var)	&var of type T*
T var;	<b>firstprivate</b> (var)	var of type T





NO Information Transfer:  
*outlined function*  $\iff$  *runtime library call site*

Declaration	Value Transfer OpenMP Clause	Communication Type
T var;	<i>default = shared</i>	&var of type T*
T var;	<b>shared</b> (var)	&var of type T*
T var;	<b>lastprivate</b> (var)	&var of type T*
T var;	<b>firstprivate</b> (var)	var of type T
T var;	<b>private</b> (var)	<i>none</i>





# TARGET REGION — THE INTERFACE

```
void kernel(...) {  
  
  init:  
    char ThreadKind = __kmpc_target_region_kernel_init(...);  
    if (ThreadKind == -1) {           // actual worker thread  
      if (!UsedLibraryStateMachine)  
        user_code_state_machine();  
      goto exit;  
    } else if (ThreadKind == 0) {     // surplus worker thread  
      goto exit;  
    } else {                          // team master thread  
      goto user_code;  
    }  
  
  user_code:  
    // User defined kernel code, parallel regions are replaced by  
    // by __kmpc_target_region_kernel_parallel(...) calls.  
  
    // Fallthrough to de-initialization  
  deinit:  
    __kmpc_target_region_kernel_deinit(...);  
  
  exit:  
    /* exit the kernel */  
}
```



# TARGET REGION — THE INTERFACE

```
// Initialization
int8_t __kmpc_target_region_kernel_init(ident_t *Ident,
                                         bool UseSPMDMode,
                                         bool RequiresOMPRuntime,
                                         bool UseStateMachine,
                                         bool RequiresDataSharing);
```

```
// De-Initialization
void __kmpc_target_region_kernel_deinit(ident_t *Ident,
                                         bool UseSPMDMode,
                                         bool RequiredOMPRuntime);
```

```
// Parallel execution
typedef void (*ParallelWorkFnTy)(void * /* SharedValues */,
                                  void * /* PrivateValues */)
```

```
CALLBACK(ParallelWorkFnTy, SharedValues, PrivateValues)
void __kmpc_target_region_kernel_parallel(ident_t *Ident,
                                         bool UseSPMDMode, bool RequiredOMPRuntime,
                                         ParallelWorkFnTy ParallelWorkFn, void *SharedValues,
                                         uint16_t SharedValuesBytes, void *PrivateValues,
                                         uint16_t PrivateValuesBytes, bool SharedMemPointers);
```



## TARGET REGION — THE IMPLEMENTATION

- (almost) the same as with the current NVPTX backend, except for shared/firstprivate variables
- implemented in Cuda as part of the library, not generated into the user code/module/TU by Clang
- the boolean flags are commonly constant, after inlining all target region abstractions is gone



# ACTION ITEM III

---

†



1) Review your OpenMP target code.

---

†



- 1) Review your OpenMP target code.
- 2) Email me<sup>†</sup> if you use the “bad” pattern!

---

<sup>†</sup> [jdoerfert@anl.gov](mailto:jdoerfert@anl.gov)





- started the review process



- started the review process
- more test-cases needed to determine benefit



## CURRENT WORK — REVIEWS, EVALUATION, FEATURES, HARDENING

- started the review process
- more test-cases needed to determine benefit
- more developers needed to add missing features



## CURRENT WORK — REVIEWS, EVALUATION, FEATURES, HARDENING

- started the review process
- more test-cases needed to determine benefit
- more developers needed to add missing features
- more users/developers needed to improve test coverage



Interested? Please let me know!

- started the review process
- more test-cases needed to determine benefit
- more developers needed to add missing features
- more users/developers needed to improve test coverage



- improve and extend the LLVM's OpenMP optimizations: connection to abstract callsites, memory placement, ...



## FUTURE WORK — OPTIMIZATIONS, FRONT-ENDS, TARGETS

- improve and extend the LLVM's OpenMP optimizations: connection to abstract callsites, memory placement, ...
- use target regions in other “front-ends”: F18, Polly, Rust?, ...



## FUTURE WORK — OPTIMIZATIONS, FRONT-ENDS, TARGETS

- improve and extend the LLVM's OpenMP optimizations: connection to abstract callsites, memory placement, ...
- use target regions in other “front-ends”: F18, Polly, Rust?, ...
- implement the interface for other targets: GPUs, FPGAs?, ...





Interested? Please let me know!

- improve and extend the LLVM's OpenMP optimizations: connection to abstract callsites, memory placement, ...
- use target regions in other “front-ends”: F18, Polly, Rust?, ...
- implement the interface for other targets: GPUs, FPGAs?, ...

