

OPTIMIZING INDIRECTIONS, OR USING ABSTRACTIONS WITHOUT REMORSE

LLVMDev'18 — October 18, 2018 — San Jose, California, USA

Johannes Doerfert, Hal Finkel

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



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.

CONTEXT & MOTIVATION

Optimizations for *sequential* aspects

Optimizations for *parallel* aspects

a

b



Optimizations for *sequential* aspects

- Can reuse (improved) *existing* transformations

Optimizations for *parallel* aspects

a

b



Optimizations for *sequential* aspects

- Can reuse (improved) *existing* transformations

Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations* (see IWOMP'18^a)

^aCompiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

^b



Optimizations for *sequential* aspects

- Can reuse (improved) *existing* transformations
- ⇒ Introduce *suitable abstractions and transformations* to bridge the indirection

Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations* (see IWOMP'18^a)

^aCompiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

^b



Optimizations for *sequential* aspects

- Can reuse (improved) *existing* transformations
- ⇒ Introduce *suitable abstractions and transformations* to bridge the indirection

Optimizations for *parallel* aspects

- New *explicit parallelism-aware transformations* (see IWOMP'18^a)
- ⇒ Introduce a unifying abstraction layer (see EuroLLVM'18 Talk^b)

^aCompiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

^bA Parallel IR in Real Life: Optimizing OpenMP, H. Finkel, J. Doerfert, X. Tian, G. Stelle, Euro-LLVM



Optimizations for *sequential* aspects

- Can reuse (improved) *existing* transformations
- ⇒ Introduce *suitable abstractions and transformations* to bridge the indirection

Optimizations for *parallel* aspects

Interested?

Contact me and come to our BoF!

- New *explicit parallelism-aware transformations* (see IWOMP'18^a)
- ⇒ Introduce a unifying abstraction layer (see EuroLLVM'18 Talk^b)

^aCompiler Optimizations For OpenMP, J. Doerfert, H. Finkel, IWOMP 2018

^bA Parallel IR in Real Life: Optimizing OpenMP, H. Finkel, J. Doerfert, X. Tian, G. Stelle, Euro-LLVM

Original Program

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

After Optimizations



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

MOTIVATION — COMPILER OPTIMIZATION FOR PARALLELISM

Original Program

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

After Optimizations

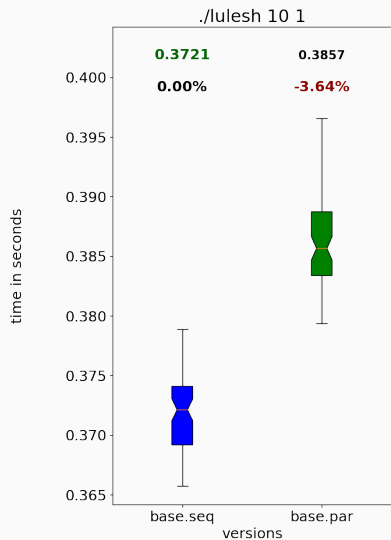
```
int y = 7;
#pragma omp parallel for
for (i = 0; i < N; i++) {
    f(y, i);
}
g(y);
```



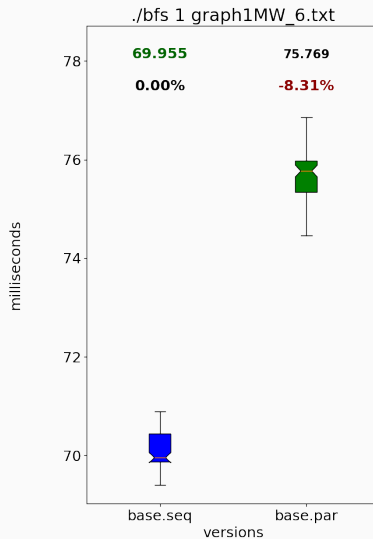
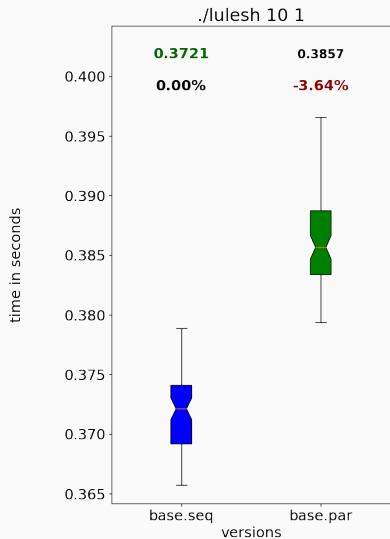
Why is this important?



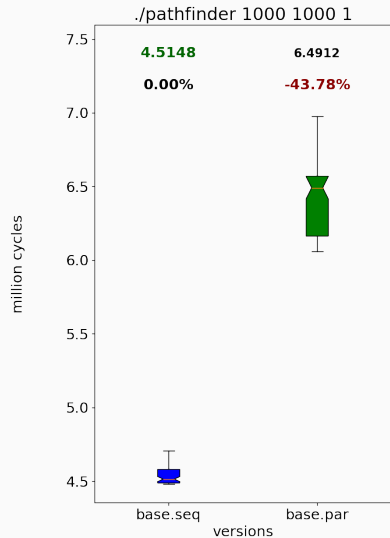
SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS



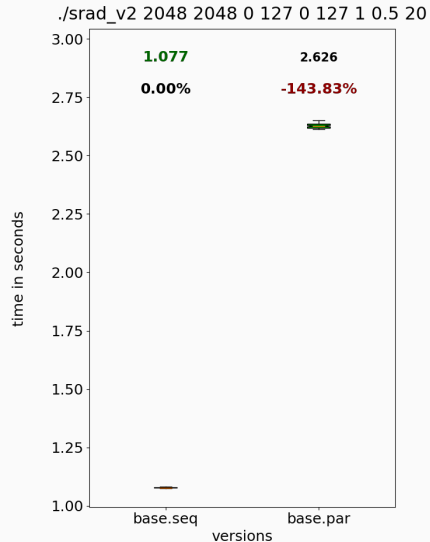
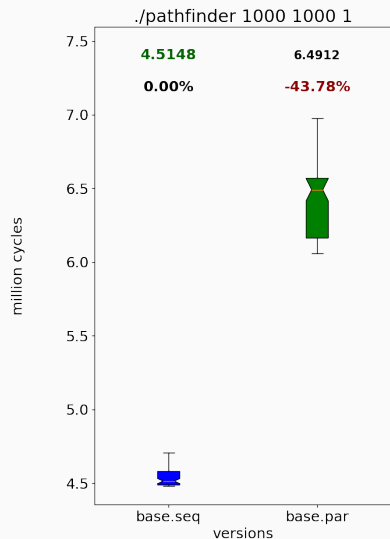
SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS



SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS



SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS



OpenMP Input: `#pragma omp parallel for`
`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);`

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);`

`// 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

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);`

`// 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

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

// Parallel region replaced by an annotated loop

```
parfor (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 OUTLINING

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);`

`// 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 + TRANSITIVE CALLS

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);`

// 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 OUTLINING + TRANSITIVE CALLS

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

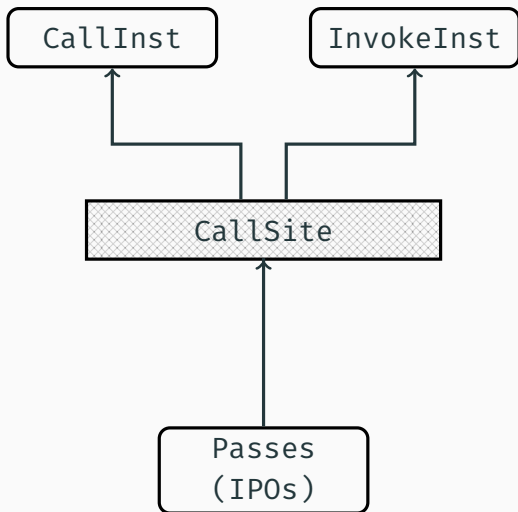
`// Parallel`
`static void`
`int lb =`
`for (int`
`(*Out)`
`}`

- + *valid and executable IR*
- + *no unintended interactions*
- + *>1k function pointers arguments in LLVM-TS + SPEC*
- *integration cost per IPO*

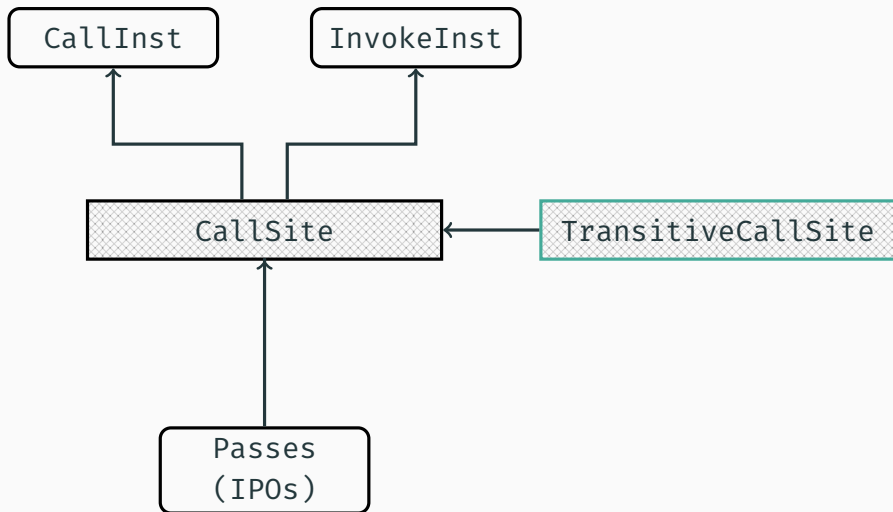
`Out) {`



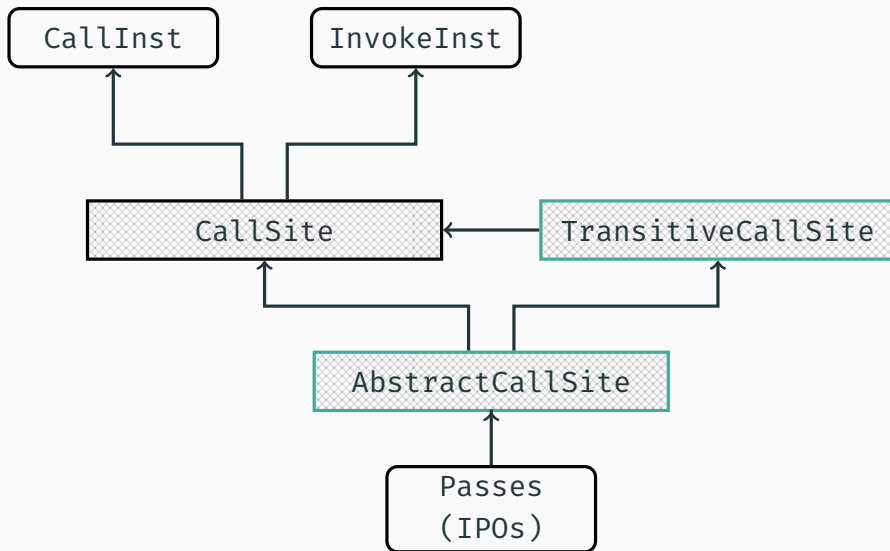
CALL ABSTRACTION IN LLVM



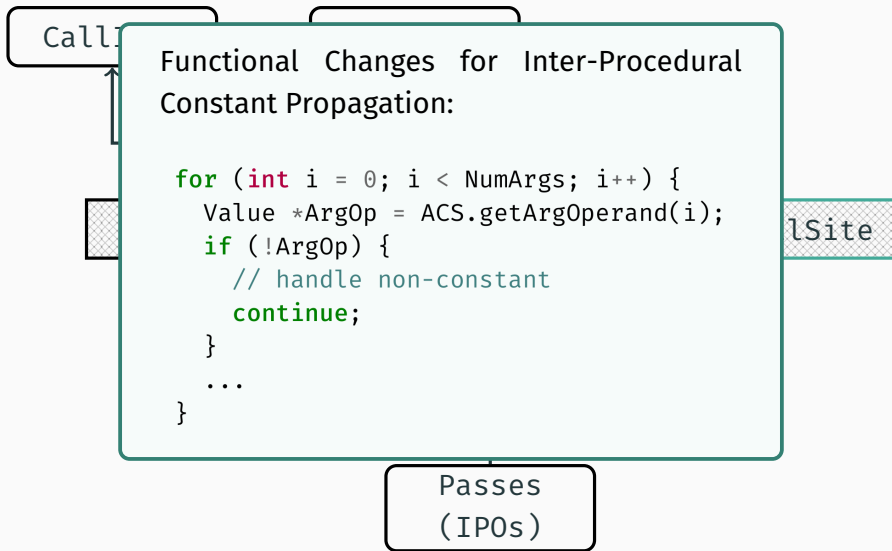
CALL ABSTRACTION IN LLVM + TRANSITIVE CALL SITES



CALL ABSTRACTION IN LLVM + TRANSITIVE CALL SITES



CALL ABSTRACTION IN LLVM + TRANSITIVE CALL SITES



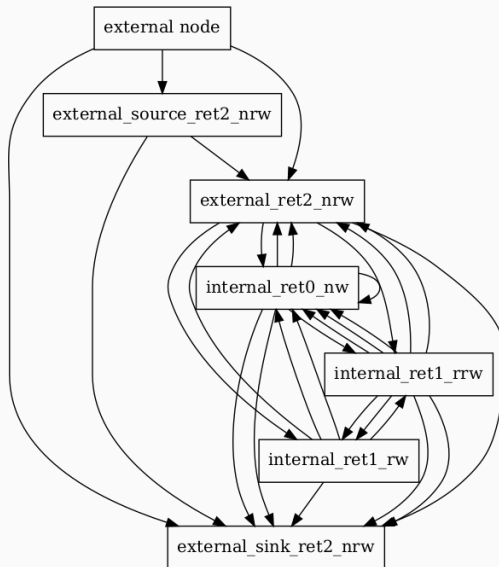
INTER-PROCEDURAL OPTIMIZATION (IPO) IN LLVM

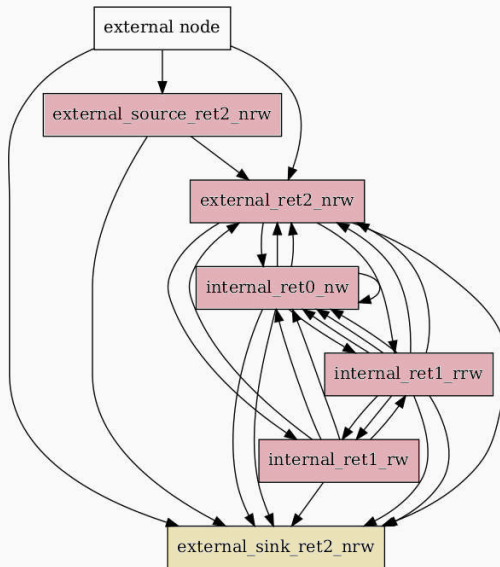
```
static int* internal_ret1_rrw(int *r0, int *r1, int *w0);  
static int* internal_ret0_nw(int *n0, int *w0);  
static int* internal_ret1_rw(int *r0, int *w0);  
int* external_source_ret2_nrw(int *n0, int *r0, int *w0);  
int* external_sink_ret2_nrw(int *n0, int *r0, int *w0);  
int* external_ret2_nrw(int *n0, int *r0, int *w0);
```

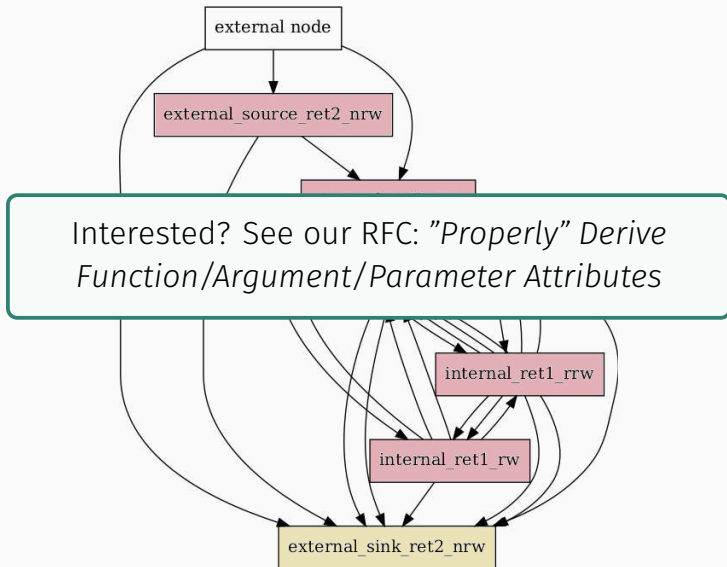



```
static int* internal_ret1_rrw(int *r0, int *r1, int *w0) {  
    if (!*r0)  
        return r1;  
    internal_ret1_rw(r0, w0);  
    *w0 = *r0 + *r1;  
    internal_ret1_rw(r1, w0);  
    internal_ret0_nw(r0, w0);  
    internal_ret0_nw(w0, w0);  
    external_ret2_nrw(r0, r1, w0);  
    external_ret2_nrw(r1, r0, w0);  
    external_sink_ret2_nrw(r0, r1, w0);  
    external_sink_ret2_nrw(r1, r0, w0);  
    return internal_ret0_nw(r1, w0);  
}
```









```
static int foo(int a, int b) {  
    return a + b; // 5?  
}
```

```
int bar() {  
    return foo(2, 3);  
}
```



```
static int foo(int a, int b) {  
    return a + b; // 5?  
}
```

```
int bar() {  
    return foo(2, 3);  
}
```

```
struct Pair {  
    int a, b;  
};
```

```
static int foo(struct Pair p) {  
    return p.a + p.b; // 5?  
}
```

```
int bar() {  
    struct Pair p = {2, 3};  
    return foo(p);  
}
```



```
static int foo(int a, int b) {  
    return 5;  
}
```

```
int bar() {  
    return foo(2, 3);  
}
```

```
struct Pair {  
    int a, b;  
};
```

```
static int foo(struct Pair p) {  
    return 5;  
}
```

```
int bar() {  
    struct Pair p = {2, 3};  
    return foo(p);  
}
```



```
struct Pair {  
    int a, b;  
};
```

```
static int foo(struct Pair *p) {  
    return p->a + p->b; // 5?  
}
```

```
int bar() {  
    struct Pair p = {2, 3};  
    return foo(&p);  
}
```

```
struct Tuple {  
    int a, b, c, d;  
};
```

```
static int foo(struct Tuple t) {  
    return t.a + t.b + t.c + t.d; // 5?  
}
```

```
int bar() {  
    struct Tuple t = {2, 3, 0, 0};  
    return foo(t);  
}
```




```
struct Pair {  
    int a, b;  
};
```

```
static int foo(struct Pair *p) {  
    return p->a + p->b;  
}
```

```
int bar() {  
    struct Pair p = {2, 3};  
    return foo(&p);  
}
```

```
struct Tuple {  
    int a, b, c, d;  
};
```

```
static int foo(struct Tuple t) {  
    return t.a + t.b + t.c + t.d;  
}
```

```
int bar() {  
    struct Tuple t = {2, 3, 0, 0};  
    return foo(t);  
}
```



```
struct Pair {  
    int a, b;  
};
```

```
static int foo(struct Pair p) {  
    return p->a + p->b;  
}
```

```
int bar() {  
    struct Pair p = {2, 3};  
    return foo(&p);  
}
```

```
struct Tuple {  
    int a, b, c, d;  
};
```

```
static int foo(struct Tuple t) {  
    return t.a + t.b + t.c + t.d;  
}
```

```
int bar() {  
    struct Tuple t = {2, 3, 0, 0};  
    return foo(t);  
}
```

Why? Pipeline is less tuned and passes are conservative for IPO.



```
struct Tuple {  
    int a, b, c, d, e, f, g, *h;  
};  
  
static  
int f(struct Tuple *t) {  
    return t->a+t->c+t->e+t->g;  
}  
  
int bar(struct Tuple *t) {  
    t->a = 3; t->e = 7; /* ... */  
    f(t);  
    // ... t->h does escape in f!  
}
```



```
struct Tuple {  
    int a, b, c, d, e, f, g, *h;  
};  
static  
int f(struct Tuple *t) {  
    return t->a+t->c+t->e+t->g;  
}
```

```
int bar(struct Tuple *t) {  
    t->a = 3; t->e = 7; /* ... */  
    f(t);  
    // ... t->h does escape in f!  
}
```

```
struct Tuple {  
    int a, b, c, d, e, f, g, *h;  
};  
static  
int f(int a, int c, int e, int g) {  
    return a + c + e + g;  
}
```

```
int bar(struct Tuple *t) {  
    t->a = 3; t->e = 7; /* ... */  
    f(3, t->c, 7, t->g);  
    // ... t->h does not escape in f!  
}
```



```
struct Tuple {  
    int a, b, c, d, e, f, g, *h;  
};  
static  
int f(struct Tuple *t) {  
    return t->a+t->c+t->e+t->g;  
}
```

```
int bar(struct Tuple *t) {  
    t->a = 3; t->e = 7; /* ... */  
    f(t);  
    // ... t->h does escape in f!  
}
```

```
struct Tuple {  
    int a, b, c, d, e, f, g, *h;  
};  
static  
int f(int c, int g) {  
    return 3 + c + 3 + g;  
}
```

```
int bar(struct Tuple *t) {  
    t->a = 3; t->e = 7; /* ... */  
    f(t->c, t->g);  
    // ... t->h does escape in f!  
}
```



```
struct Tuple {  
    int a, b, c, d, e, f, g, *h;  
};  
static  
int f(struct Tuple *t) {  
    return t->a+t->c+t->e+t->g;  
}
```

```
int bar(struct Tuple *t) {  
    t->a = 3; t->e = 7; /* ... */  
    f(t);  
    // ... t->h does escape in f!  
}
```

```
struct Tuple {  
    int a, b, c, d, e, f, g, *h;  
};  
static  
int f(struct Tuple *t) {  
    return 3 + t->c + 7 + t->g;  
}
```

```
int bar(struct Tuple *t) {  
    t->a = 3; t->e = 7; /* ... */  
    f(t);  
    // ... t->h does escape in f!  
}
```



```
struct Tuple {
    int a, b, c, d, e, f, g, *h;
};
```

```
static
```

```
int f(struct Tuple t) {
    return t->a+t->g;
}
```

```
int bar(struct Tuple t) {
    t->a = 3; t->e = 7; /* ... */
    f(t);
    // ... t->h does escape in f!
}
```

```
struct Tuple {
    int a, b, c, d, e, f, g, *h;
};
```

```
static
```

```
int f(struct Tuple t) {
    return 7 + t->g;
}
```

```
int bar(struct Tuple t) {
    t->a = 3; t->e = 7; /* ... */
    f(t);
    // ... t->h does escape in f!
}
```

Aggressively unpack object arguments early and condense arguments late as an alternative/substitution for inlining.

- track values of fields across function calls, e.g., closure initialization [prototype]
- determine performance impact of missing static information [ongoing]
- export attributes for libraries, e.g., add `__attribute__((const))` [planned]

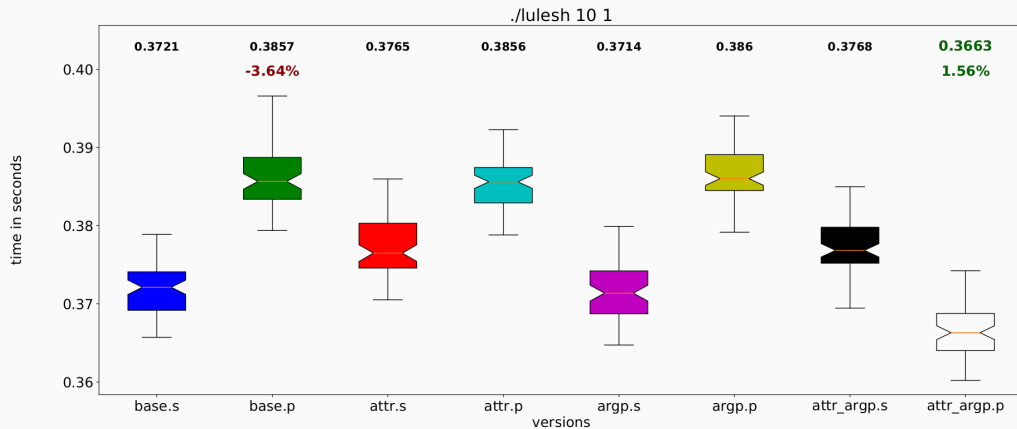
- track values of fields across function calls, e.g., closure initialization [prototype]
- determine performance impact of missing static information [ongoing]
- export attributes for libraries, e.g., add `__attribute__((const))` [planned]

Interested? Contact me!

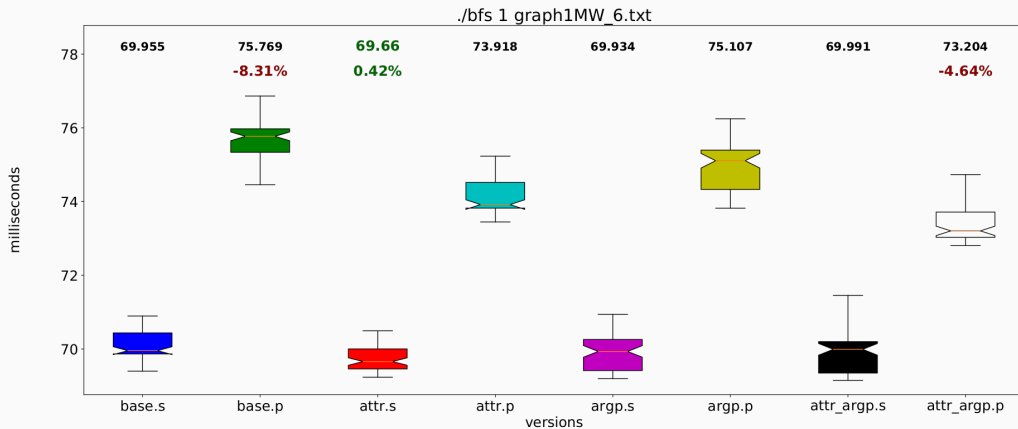
EVALUATION

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)	

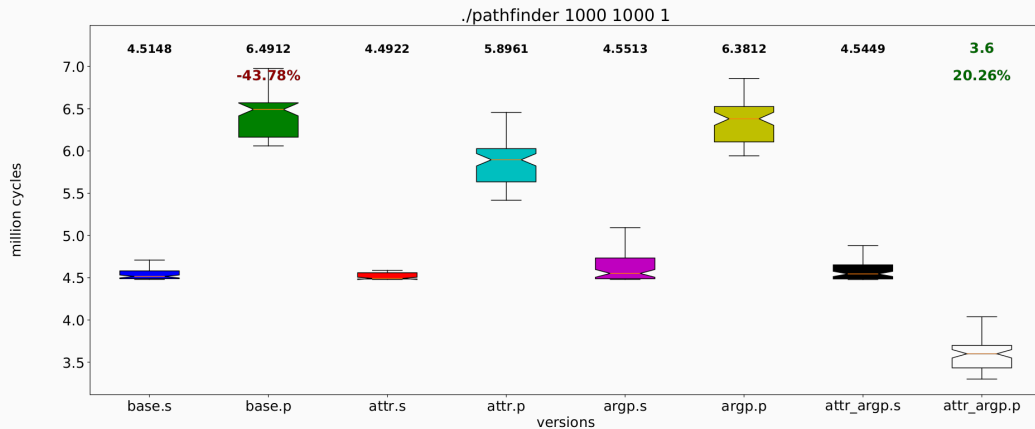
OPENMP OPTIMIZATIONS — PERFORMANCE RESULTS



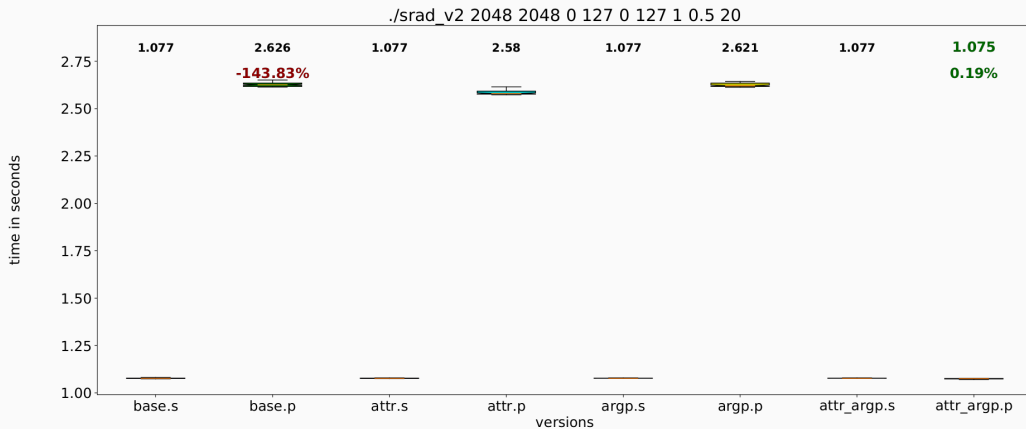
OPENMP OPTIMIZATIONS — PERFORMANCE RESULTS



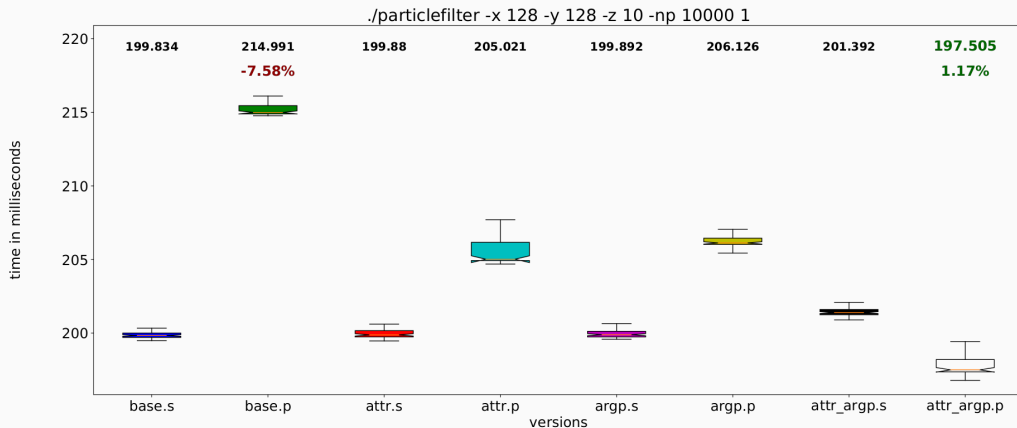
OPENMP OPTIMIZATIONS — PERFORMANCE RESULTS



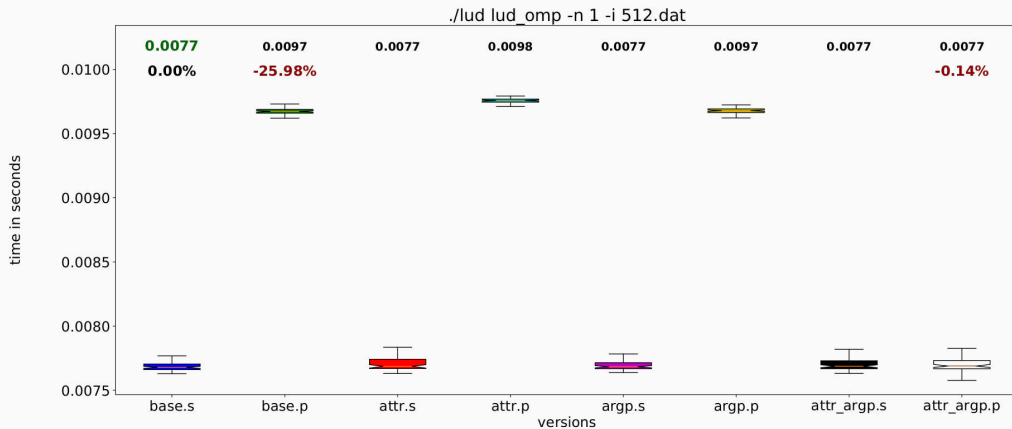
OPENMP OPTIMIZATIONS — PERFORMANCE RESULTS



OPENMP OPTIMIZATIONS — PERFORMANCE RESULTS



OPENMP OPTIMIZATIONS — PERFORMANCE RESULTS

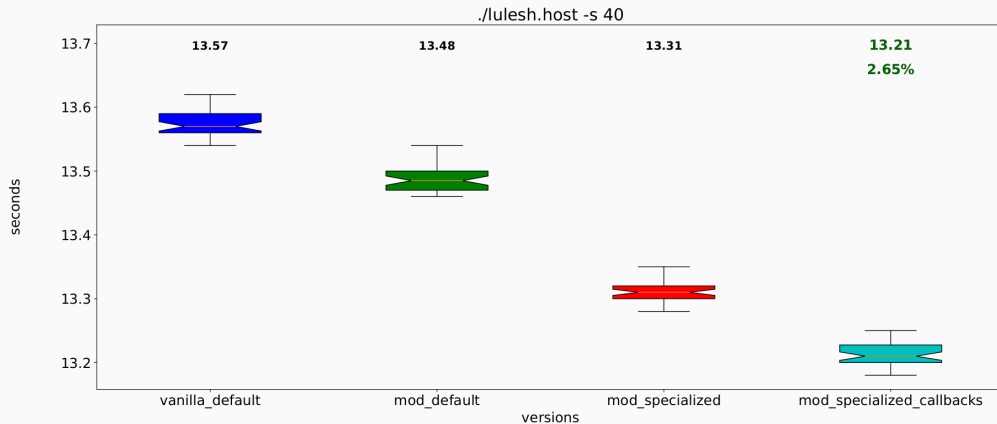


ARRAY CONSTANT PROPAGATION EXAMPLE

```
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]  
    }  
)
```



ARRAY CONSTANT PROPAGATION PERFORMANCE



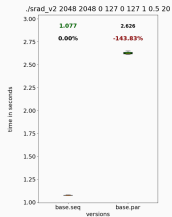
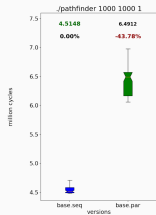
CONCLUSION

CONCLUSION



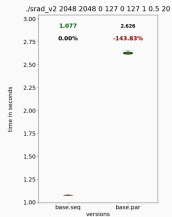
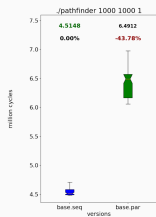
CONCLUSION

SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS

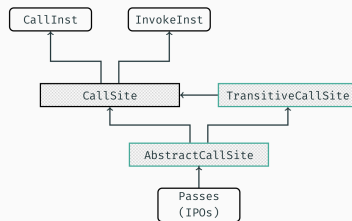


CONCLUSION

SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS

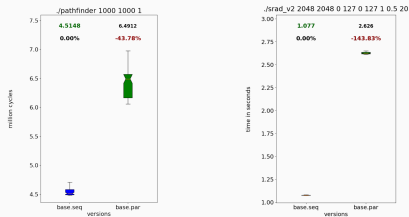


CALL ABSTRACTION IN LLVM + TRANSITIVE CALL SITES

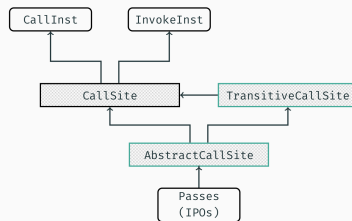


CONCLUSION

SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS

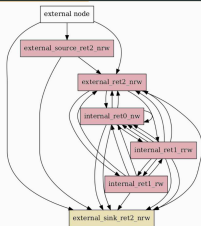


CALL ABSTRACTION IN LLVM + TRANSITIVE CALL SITES



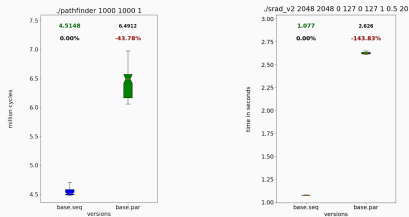
IPO — ATTRIBUTE INFERENCE

-O3 -DISABLE-INLINING

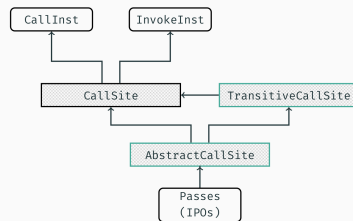


CONCLUSION

SEQUENTIAL PERFORMANCE OF PARALLEL PROGRAMS

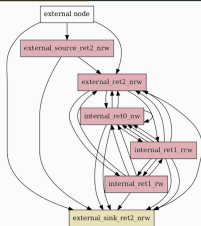


CALL ABSTRACTION IN LLVM + TRANSITIVE CALL SITES



IPO — ATTRIBUTE INFERENCE

-O3 -DISABLE-INLINING



OPENMP OPTIMIZATIONS — PERFORMANCE RESULTS

