

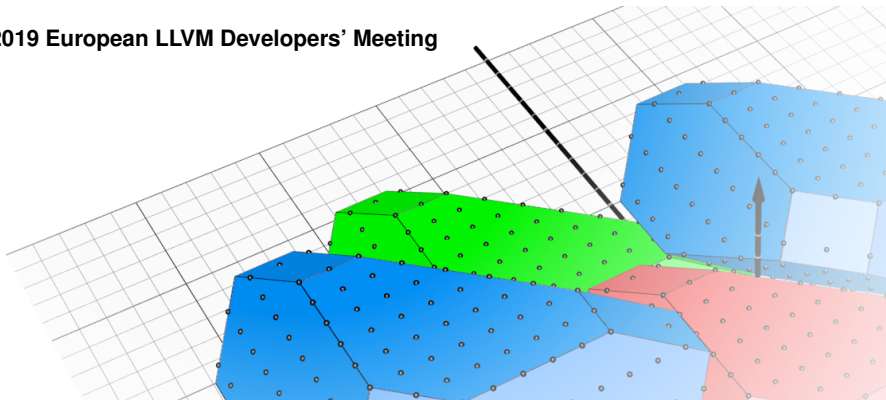
An alternative OpenMP Backend for Polly



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Michael Halkenhäuser

2019 European LLVM Developers' Meeting





- ▶ Polyhedral framework on LLVM-IR

- ▶ Polyhedral framework on LLVM-IR
 - ▶ Efficient analyses and transformations
 - ▶ Code generation

- ▶ Polyhedral framework on LLVM-IR
 - ▶ Efficient analyses and transformations
 - ▶ Code generation
- ▶ Example transformations
 - ▶ Loop interchange / fission / fusion
 - ▶ Strip mining (Vectorization)

- ▶ Polyhedral framework on LLVM-IR
 - ▶ Efficient analyses and transformations
 - ▶ Code generation

- ▶ Example transformations
 - ▶ Loop interchange / fission / fusion
 - ▶ Strip mining (Vectorization)
 - ▶ Automatic parallelization

- ▶ Automatic parallelization
 - ▶ No need for **manual annotation**

- ▶ Automatic parallelization
 - ▶ No need for **manual annotation**

```
// "matvect" -- Sequential
// (Simplified dependencies)

for (i = 0; i <= n; i++) {
    for (j = 0; j <= n; j++)
        s[i] = s[i] + a[i][j] * x[j];
}
```

Input

- ▶ Automatic parallelization
 - ▶ No need for manual annotation

```
// "matvect" -- Sequential
// (Simplified dependencies)

for (i = 0; i <= n; i++) {
    for (j = 0; j <= n; j++)
        s[i] = s[i] + a[i][j] * x[j];
}
```

Input

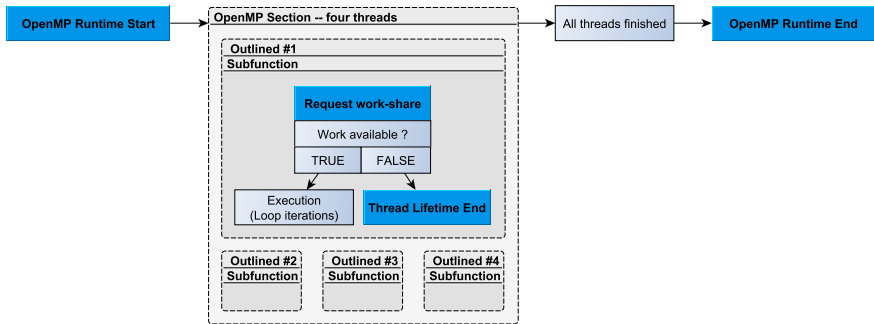
```
// "matvect" -- OpenMP parallelized
// Equivalent to the LLVM-IR output

#pragma omp parallel for [...] \
    schedule (dynamic, 1) num_threads(N)
for (i = 0; i <= n; i++) {
    for (j = 0; j <= n; j++)
        s[i] = s[i] + a[i][j] * x[j];
}
```

Output

Polly – Parallelization Scheme

- ▶ Polly detects parallelizable code regions
 - ▶ Moved into an *outlined* function
 - ▶ Executed using **OpenMP API**





- ▶ Limited influence on OpenMP execution
 - ▶ Increase number of user options
 - ▶ Improve fine-tuning possibilities



- ▶ Limited influence on OpenMP execution
 - ▶ Increase number of user options
 - ▶ Improve fine-tuning possibilities
- ▶ Dependent on GNU OpenMP API
 - ▶ Expand the scope of application

Motivation for an alternative OpenMP Backend



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- ▶ Limited influence on OpenMP execution
 - ▶ Increase number of user options
 - ▶ Improve fine-tuning possibilities
- ▶ Dependent on GNU OpenMP API
 - ▶ Expand the scope of application
- ▶ LLVM OpenMP implementation available
 - ▶ Enable direct use of LLVM's OpenMP runtime
 - ▶ Support automated testing

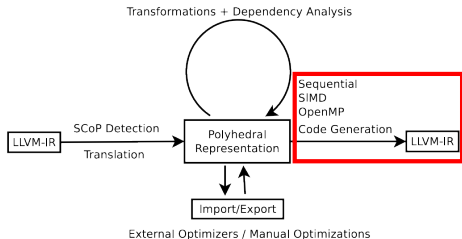
- ▶ Extension of the preexisting backend
 - ▶ Reused common functionalities
 - ▶ Moved into abstract base class



- ▶ Extension of the preexisting backend
 - ▶ Reused common functionalities
 - ▶ Moved into abstract base class
 - ▶ API-specific call creation and placement
 - ▶ Implemented in derived class per backend

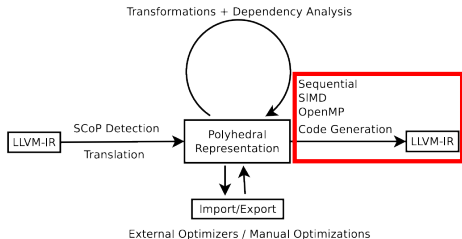
- ▶ Extension of the preexisting backend
 - ▶ Reused common functionalities
 - ▶ Moved into abstract base class
 - ▶ API-specific call creation and placement
 - ▶ Implemented in derived class per backend

- ▶ User may **choose backend**
 - ▶ Via CL switch, similar to
 - ▶ *Number of threads*



- ▶ Extension of the preexisting backend
 - ▶ Reused common functionalities
 - ▶ Moved into abstract base class
 - ▶ API-specific call creation and placement
 - ▶ Implemented in derived class per backend

- ▶ User may **choose backend**
 - ▶ Via CL switch, similar to
 - ▶ *Number of threads*
 - ▶ Additional options
 - ▶ *Scheduling type*
 - ▶ *Chunk size*



- ▶ *Scheduling type* determines work distribution

<i>static</i>	<i>dynamic</i>	<i>guided</i>
Predetermined, uniform distribution of iterations	Threads request work shares of <i>chunk size</i>	Hybrid scheduling of <i>static</i> and <i>dynamic</i> , using a <i>minimum chunk size</i>

- ▶ *Scheduling type* determines work distribution

	<i>static</i>	<i>dynamic</i>	<i>guided</i>
Load Balancing	–	+	○
Organization Overhead	+	–	○

- ▶ *Scheduling type* determines work distribution

	<i>static</i>	<i>dynamic</i>	<i>guided</i>
Load Balancing	–	+	○
Organization Overhead	+	–	○

- ▶ *static* suited for constant computational demands
- ▶ *dynamic* suited for shifting computational demands
- ▶ *guided* suited for "both"

- ▶ PolyBench¹
 - ▶ Provides multiple datasets
 - ▶ Triggers auto-parallelization in 18 benchmarks

¹<https://sourceforge.net/projects/polybench/>

- ▶ PolyBench¹
 - ▶ Provides multiple datasets
 - ▶ Triggers auto-parallelization in 18 benchmarks
- ▶ Runtime results
 - ▶ Average from 50 out of 60 runs (10% trimmed-mean)
 - ▶ Utilized CPU: AMD R5 1600X

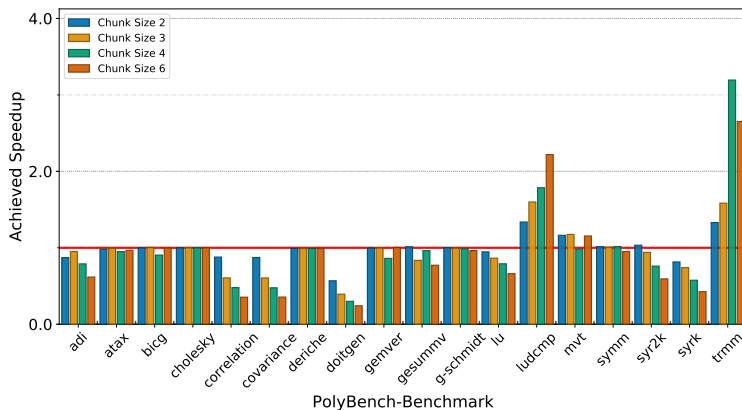
¹<https://sourceforge.net/projects/polybench/>

- ▶ PolyBench¹
 - ▶ Provides multiple datasets
 - ▶ Triggers auto-parallelization in 18 benchmarks
- ▶ Runtime results
 - ▶ Average from 50 out of 60 runs (10% trimmed-mean)
 - ▶ Utilized CPU: AMD R5 1600X
- ▶ Plots show *relative speedup*
 - ▶ $speedup = \frac{runtime\ of\ baseline}{runtime\ of\ competitor}$

¹<https://sourceforge.net/projects/polybench/>

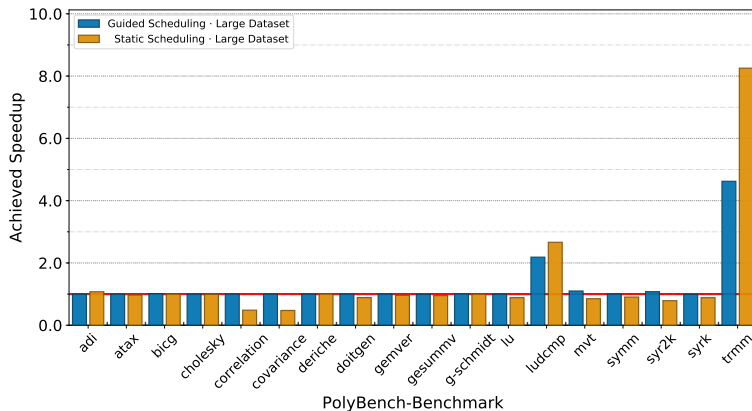
Performance Impact of *chunk size*

LLVM OpenMP Chunk Size Comparison
Large Dataset · No Vectorization · Dynamic Scheduling · 12 Threads · Baseline: Chunk Size 1



Performance Impact of *scheduling type*

LLVM OpenMP Scheduling Comparison
No Vectorization · 12 Threads · Baseline: Dynamic Scheduling





▶ *Chunk size*

- ▶ 1 is usually a reasonable choice
- ▶ Very beneficial in particular cases
 - ▶ More than $3\times$ speedup possible

▶ *Chunk size*

- ▶ 1 is usually a reasonable choice
- ▶ Very beneficial in particular cases
 - ▶ More than $3\times$ speedup possible

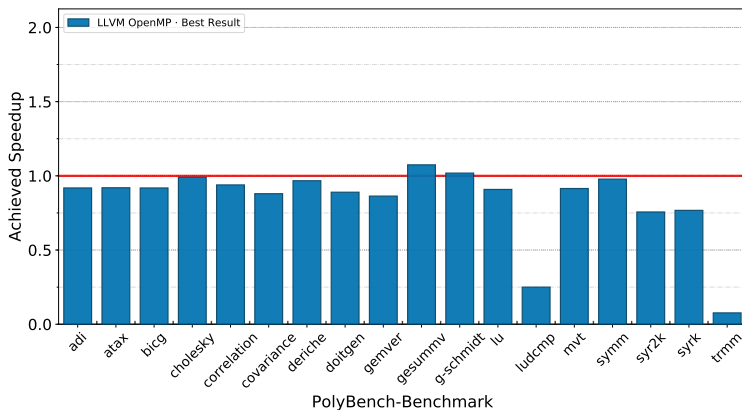
▶ *Scheduling type*

- ▶ *Dynamic*: Good overall performance
- ▶ *Guided*: Performs at least as good as *dynamic*
- ▶ *Static*: Problem-dependent
 - ▶ May achieve $8\times$ speedup compared to *dynamic*

Backend Comparison

LLVM versus GNU OpenMP Backend

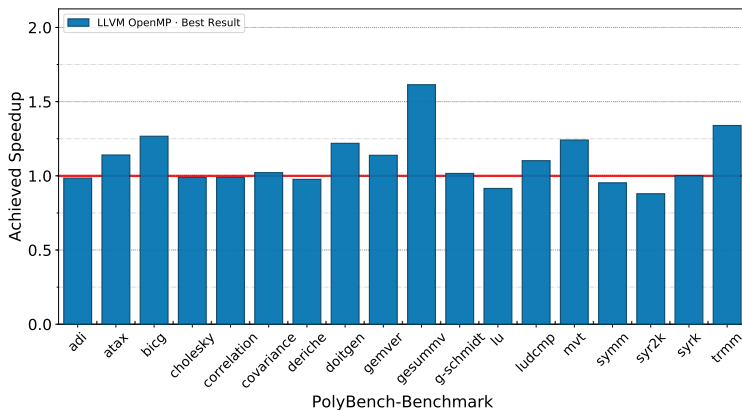
GNU & LLVM Backend Comparison
Large Dataset · No Vectorization · 4 Threads · Baseline: GNU Backend



Backend Comparison

LLVM versus GNU OpenMP Backend

GNU & LLVM Backend Comparison
Large Dataset · No Vectorization · 12 Threads · Baseline: GNU Backend





- ▶ Using the maximum number of available threads
- ▶ Our "LLVM" backend
 - ▶ Achieves comparable performance
 - ▶ Performs significantly faster than "GNU" in seven cases
 - ▶ Reaches up to $1.6\times$ speedup

- ▶ Using the maximum number of available threads
- ▶ Our "LLVM" backend
 - ▶ Achieves comparable performance
 - ▶ Performs significantly faster than "GNU" in seven cases
 - ▶ Reaches up to $1.6\times$ speedup
- ▶ GNU backend
 - ▶ Only a single, considerable lead

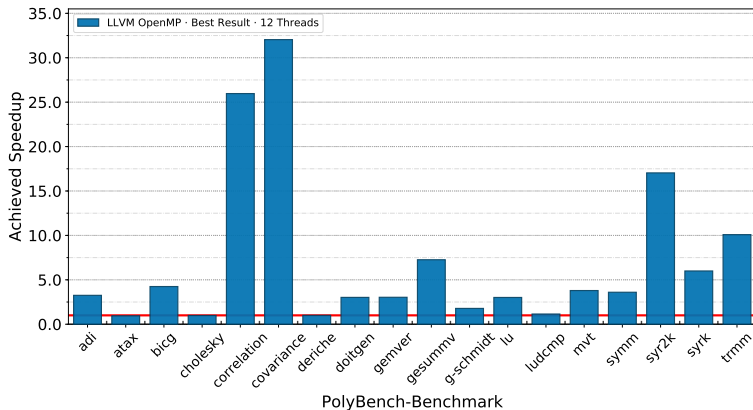


- ▶ Using the maximum number of available threads
- ▶ Our "LLVM" backend
 - ▶ Achieves comparable performance
 - ▶ Performs significantly faster than "GNU" in seven cases
 - ▶ Reaches up to $1.6\times$ speedup
- ▶ GNU backend
 - ▶ Only a single, considerable lead
- ▶ Additional switches
 - ▶ Allow problem-specific adjustments
 - ▶ ... without depending on env. variable

General Comparison

LLVM OpenMP Backend versus clang

clang Comparison
Large Dataset · With Vectorization · Baseline: clang-8 -O3



Conclusion

- ▶ Our "LLVM" OpenMP backend for Polly
 - ▶ Represents a superior alternative

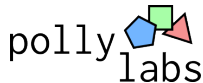
- ▶ Our "LLVM" OpenMP backend for Polly
 - ▶ Represents a superior alternative
 - ▶ Acts as drop-in replacement

- ▶ Our "LLVM" OpenMP backend for Polly
 - ▶ Represents a superior alternative
 - ▶ Acts as drop-in replacement
 - ▶ Provides more customization options

- ▶ Our "LLVM" OpenMP backend for Polly
 - ▶ Represents a superior alternative
 - ▶ Acts as drop-in replacement
 - ▶ Provides more customization options
 - ▶ Carries no clear drawbacks, but instead ...
 - ▶ Reaches up to $1.6\times$ speedup

- ▶ Our "LLVM" OpenMP backend for Polly
 - ▶ Is publicly available
 - ▶ Review accepted on March 19th
<https://reviews.lvm.org/D59100>
 - ▶ Currently on Polly's master branch
<https://github.com/llvm/llvm-project/commit/89251ed>
- ▶ References:
 - ▶ Title graphic: <https://polly.llvm.org/images/header-background.png>
 - ▶ T. Grosser, H. Zheng, R. Aloor, A. Simbürger, A. Größlinger, and L. - N. Pouchet, "Polly - Polyhedral optimization in LLVM," in Proceedings of the First International Workshop on Polyhedral Compilation Techniques (IMPACT), vol. 2011, 2011, p. 1.

Questions ?



- ▶ Ask them now, or ...
- ▶ Find me tomorrow, at the poster session
 - ▶ 09:00 am - 10:00 am (Foyer)