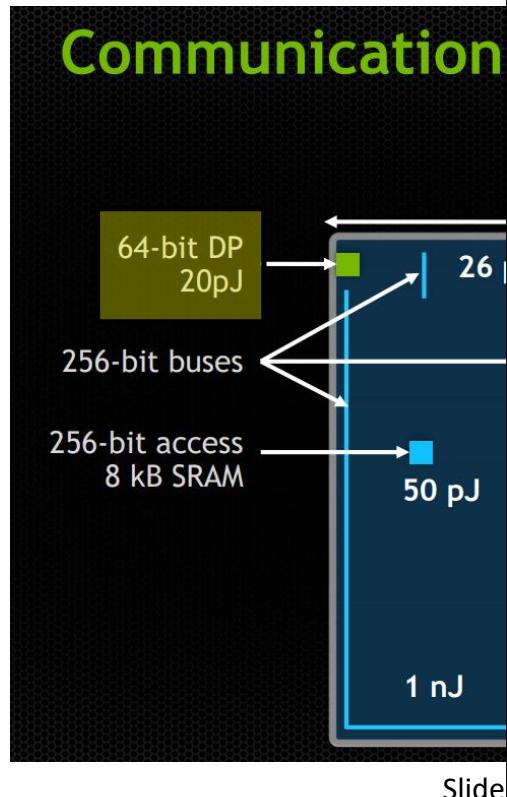


Alexandru Calotoiu, Tal Ben-Nun, Grzegorz Kwasniewski, Johannes de Fine Licht,
Timo Schneider, Philipp Schaad, Torsten Hoefer

Lifting C Semantics for Dataflow Optimization

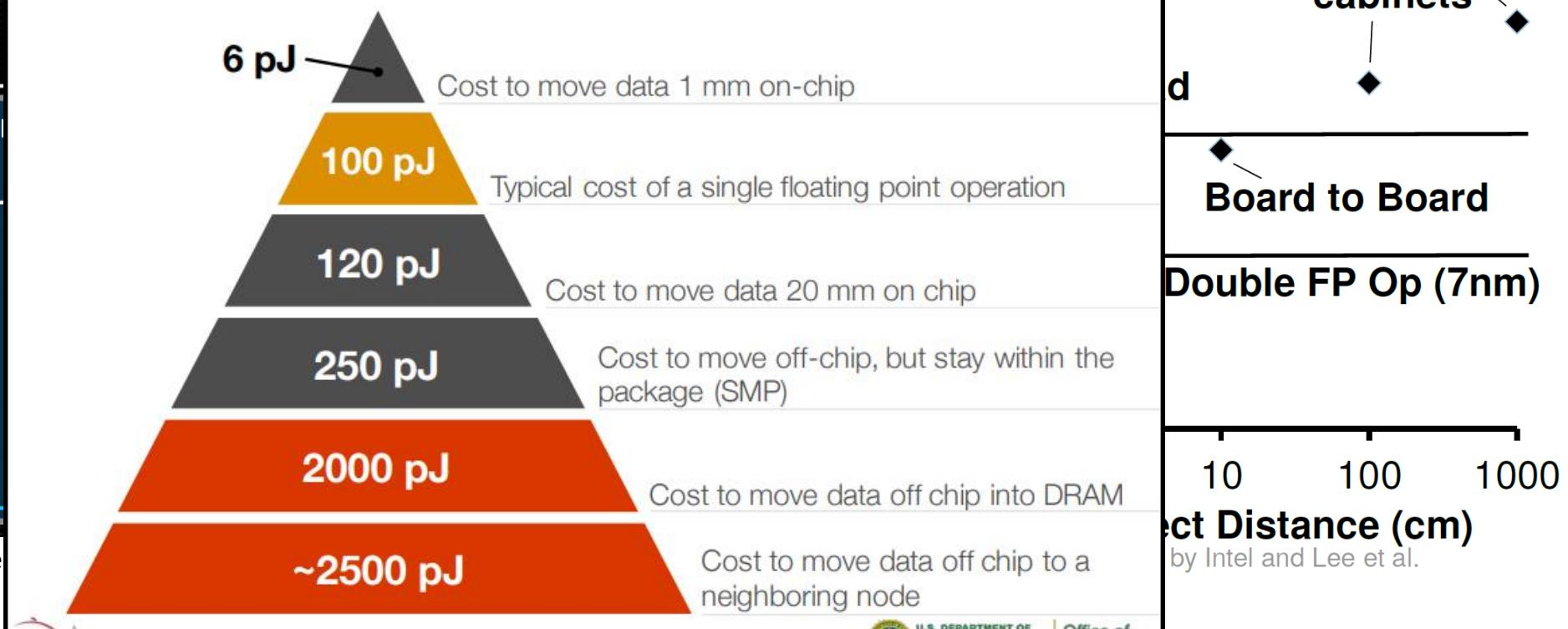


Data movement dominates energy costs!



Hierarchical Power Costs

Data Movement is the Dominant Power Cost



Efficient data movement is hard!

Halide: Decoupling Schedules from Data Movement

MLIR: A Compiler Infrastructure for Machine Learning

LIFT: A Functional Data-Parallel IR for High-Performance GPU Code Generation

Michel Steuwer Toomas Remmelg Christophe Dubach
 University of Edinburgh, United Kingdom
 {michel.steuwer, toomas.remmelg, christophe.dubach}@ed.ac.uk

Abstract

Parallel patterns (*e.g.*, map, reduce) have gained traction as an abstraction for targeting parallel accelerators and are a promising answer to the performance portability problem. However, compiling high-level programs into efficient low-level parallel code is challenging. Current approaches start from a high-level parallel IR and proceed to emit GPU code directly in one big step. Fixed strategies are used to optimize and map parallelism exploiting properties of a particular GPU generation leading to performance portability issues.

We introduce the LIFT IR, a new data-parallel IR which encodes OpenCL-specific constructs as functional patterns. Our prior work has shown that this functional nature simplifies the exploration of optimizations and mapping of parallelism from portable high-level programs using rewrite-rules.

This paper describes how LIFT IR programs are compiled into efficient OpenCL code. This is non-trivial as many performance sensitive details such as memory allocation, array

Polyhedral-Based Data Reuse Optimization for Configurable Computing

Louis-Noël Pouchet,¹ Peng Zeng,²
¹ University of California, Los Angeles
² Ohio State University

MAPS: Optimizing Massively Parallel Applications Using Device-Level Memory Abstraction

EPI DURIN ELY LEVY, AMNON BARAK, and TAL BEN-NUN,
 University of Jerusalem



reasingly important role in high-performance computing. While developing naive code is optimizing massively parallel applications requires deep understanding of the underlying developer must struggle with complex index calculations and manual memory transfers. identifies memory access patterns used in most parallel algorithms, based on Berkeley's Parthenon proposes the MAPS framework, a device-level memory abstraction that facilitates GPUs, alleviating complex indexing using on-device containers and iterators. This article implementation of MAPS and shows that its performance is comparable to carefully optimized real-world applications.

Subject Descriptors: C.1.4 [Parallel Architectures]: GPU Memory Abstraction
 Parallelism, Abstraction, Performance

Keywords and Phrases: GPGPU, memory abstraction, heterogeneous computing architectures, patterns

Format:

Levy, Amnon Barak, and Tal Ben-Nun. 2014. MAPS: Optimizing massively parallel applications using device-level memory abstraction. ACM Trans. Architec. Code Optim. 11, 4, Article 44 (December

Spatial: A Language and Compiler for Application Accelerators

David Koeplinger[†] Matthew Feldman[†] Raghu Prabhakar[†] Yaqi Zhang[†]
 Michael Hadjis[†] Ruben Fiszel[‡] Tian Zhao[†] Luigi Nardi[†] Ardavan Pedram[†]
 Christos Kozyrakis[†] Kunle Olukotun[†]

[†] Stanford University, USA

ne

air
Science
at
n
du

Science
at
n
i

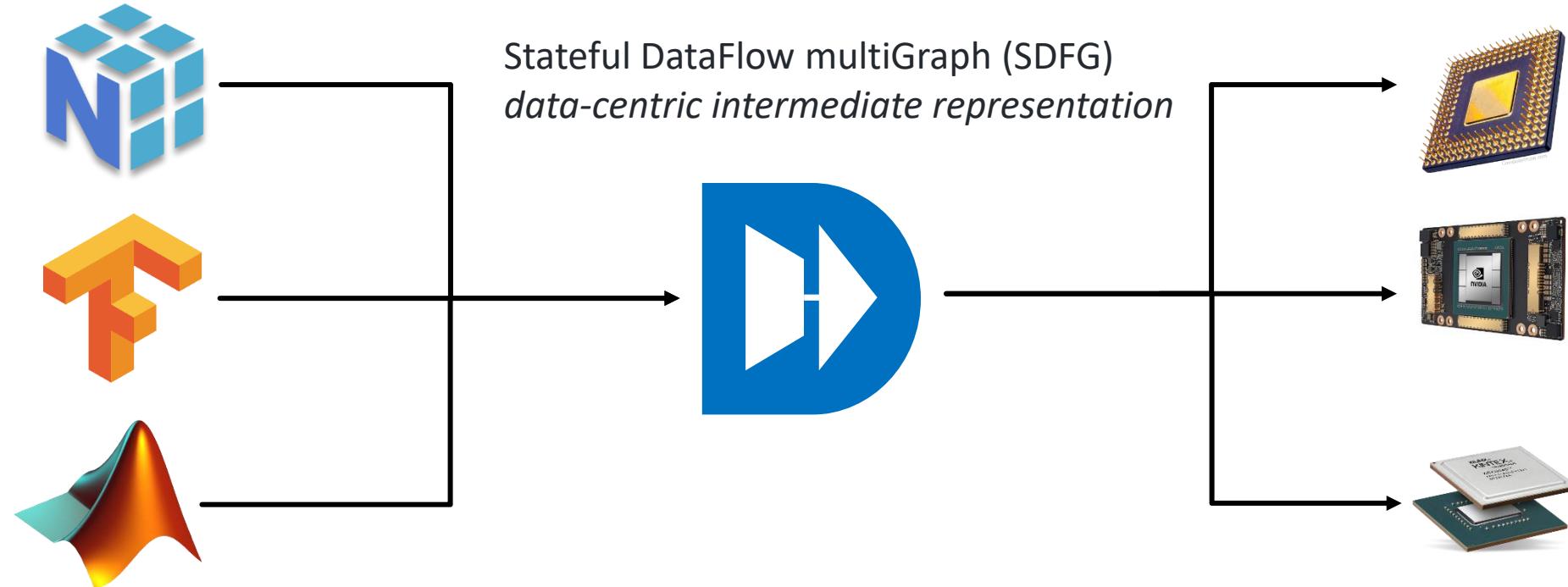
es can make
to exploit
onclude that
for achieving
parallelizing

nization →

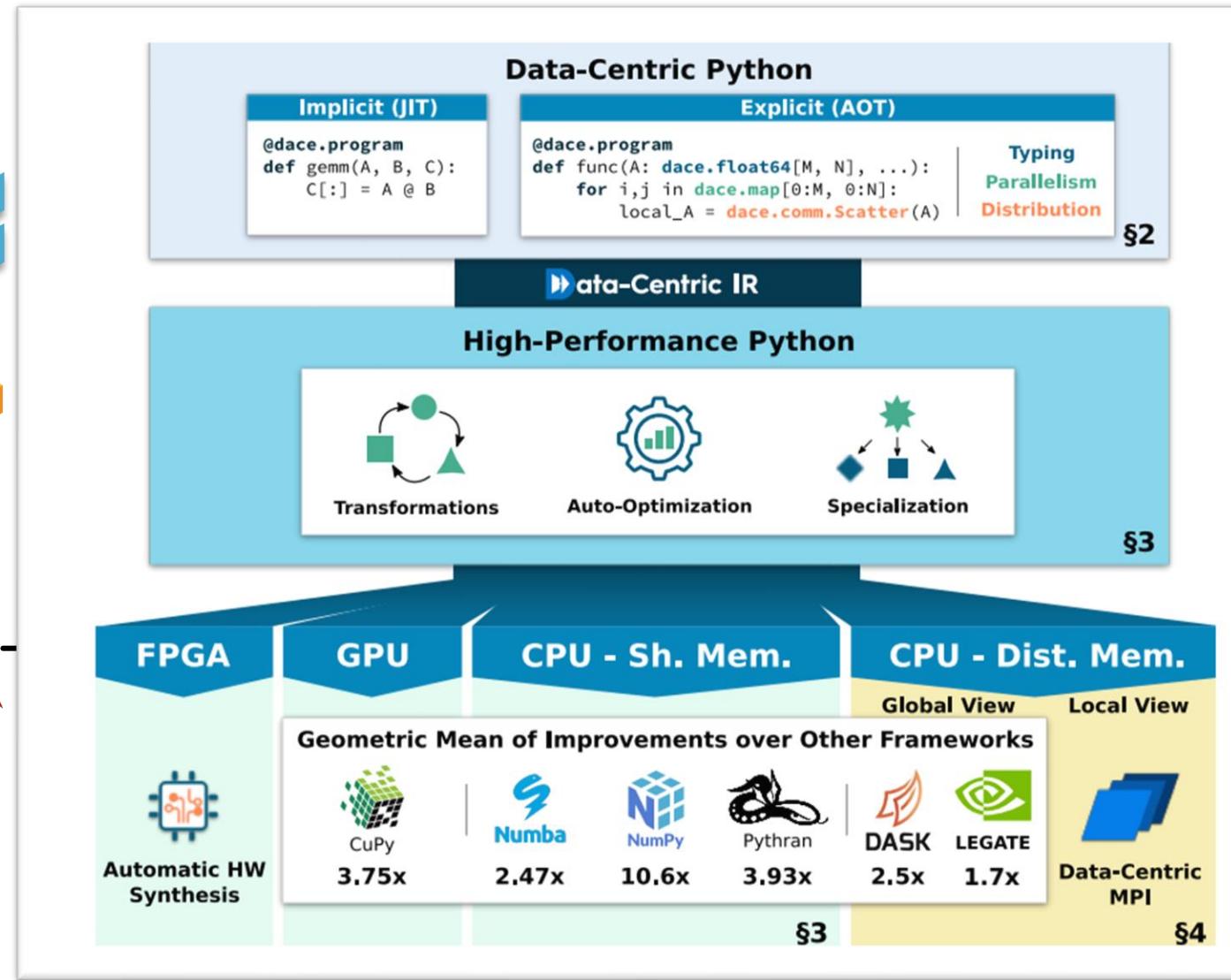
heterogeneou

increasingly
ile devices to
systems are
ing elements,
and domain-
e energy ef

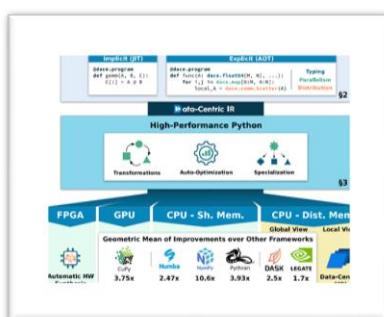
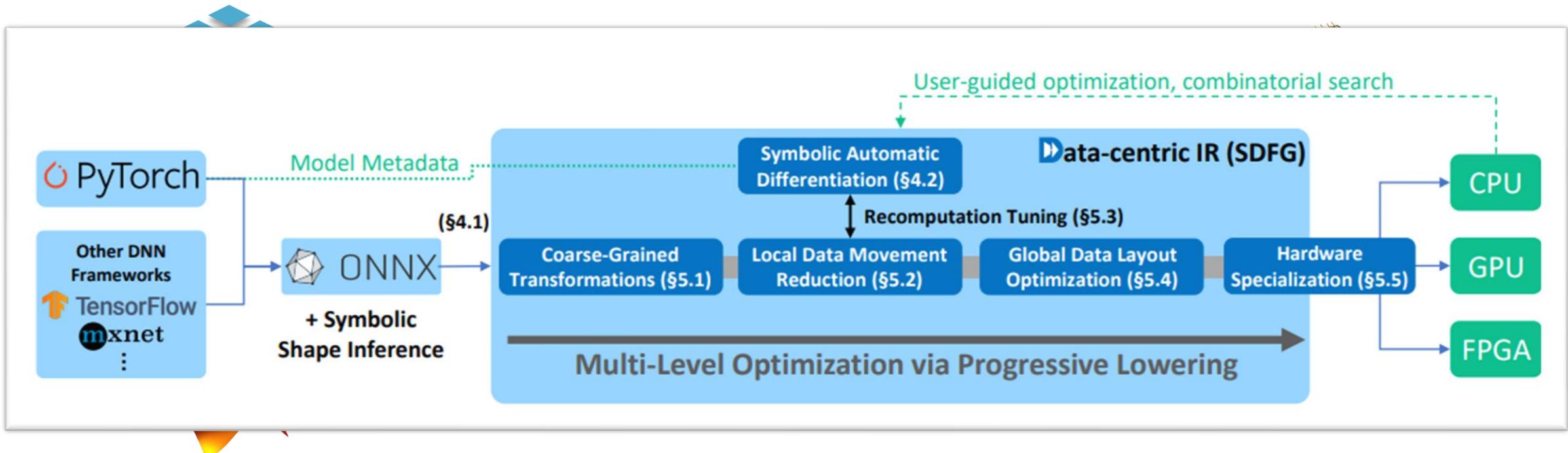
Data Centric Parallel Programming



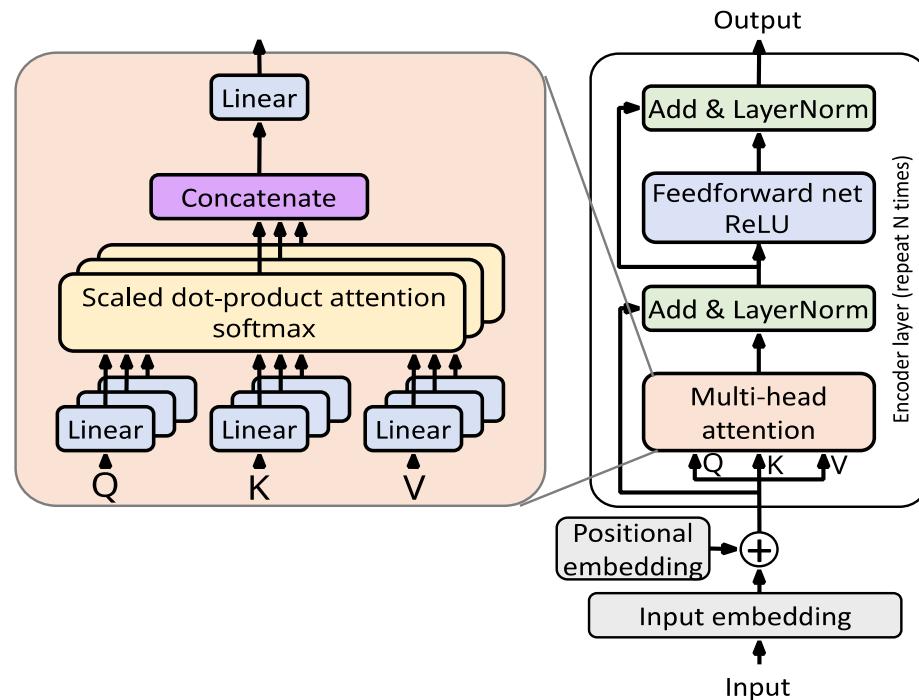
Data Centric Parallel Programming



Data Centric Parallel Programming



Data Centric Parallel Programming



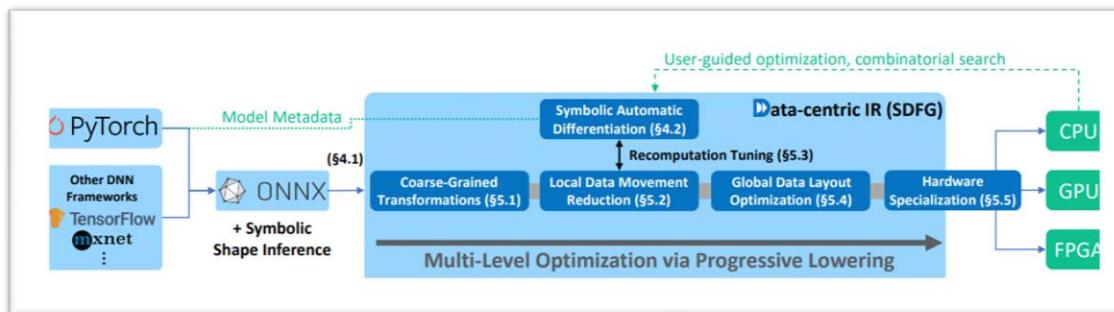
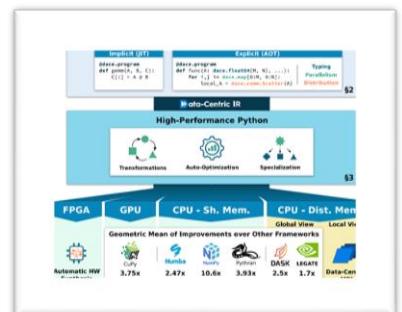
GPT-3 has a training cost of \$12M

Our performance improvement for BERT-large

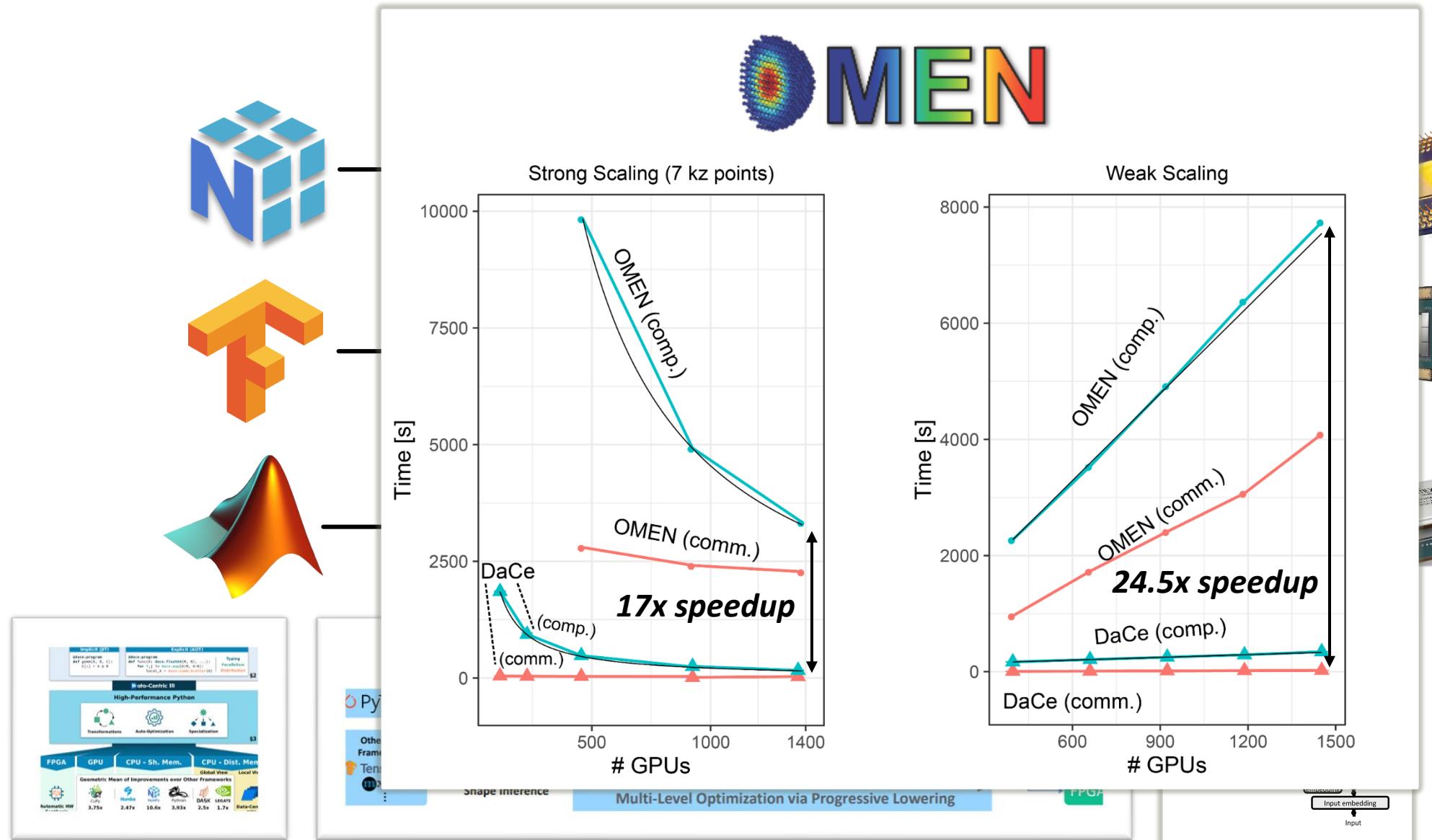
- 30% over PyTorch
- 20% over Tensorflow + XLA
- 8% over DeepSpeed

Estimated savings

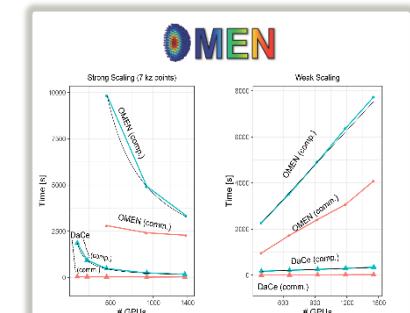
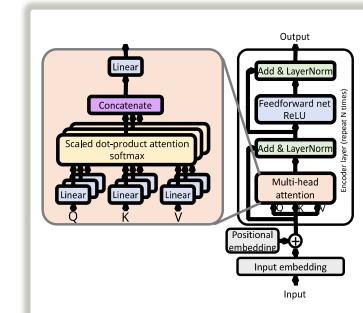
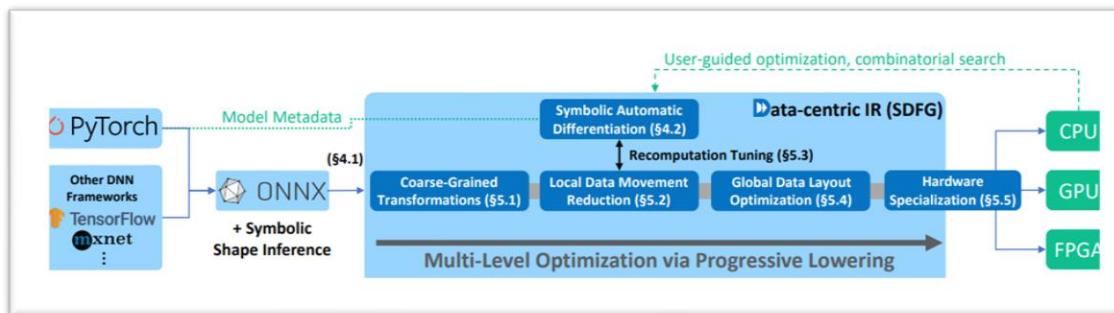
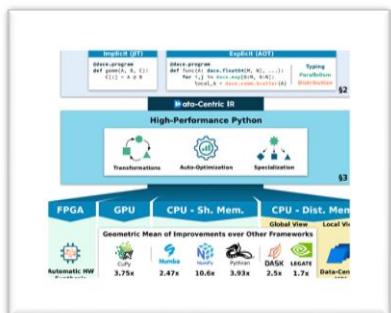
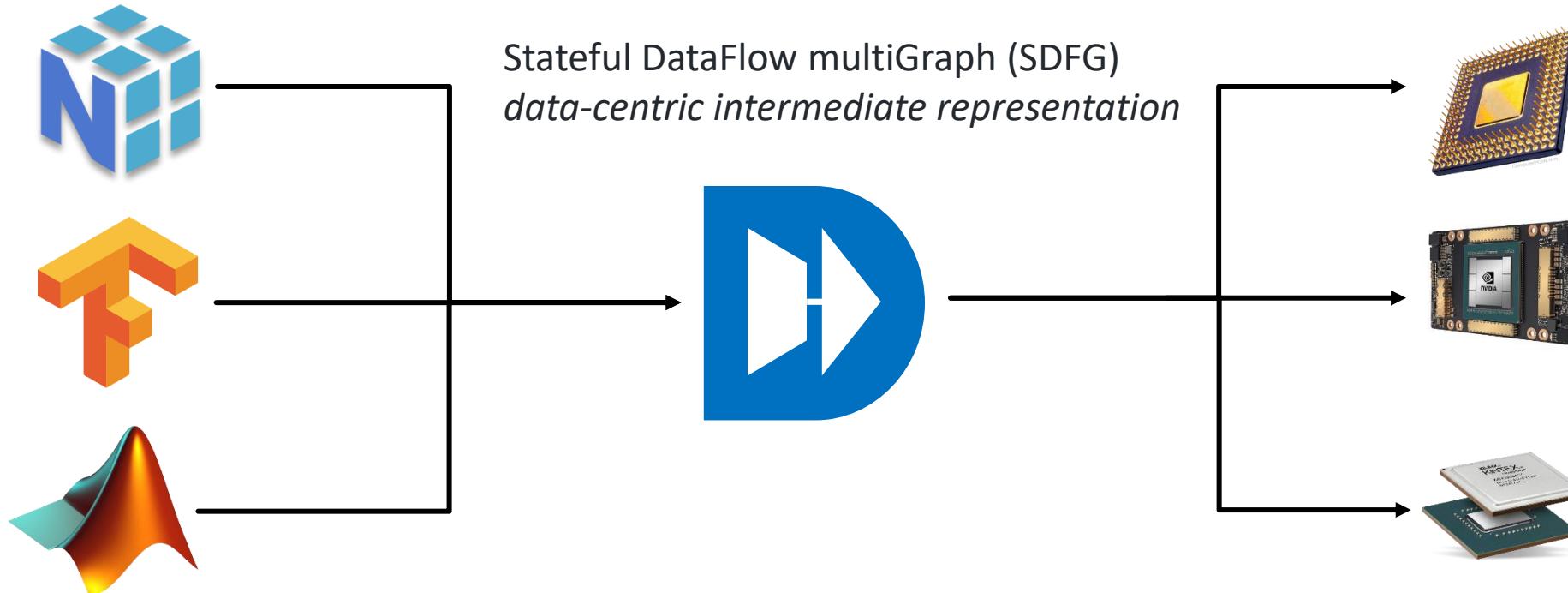
- \$85,000 for BERT on AWS using PyTorch
- \$3.6M and 120 MWh for GPT-3



Data Centric Parallel Programming



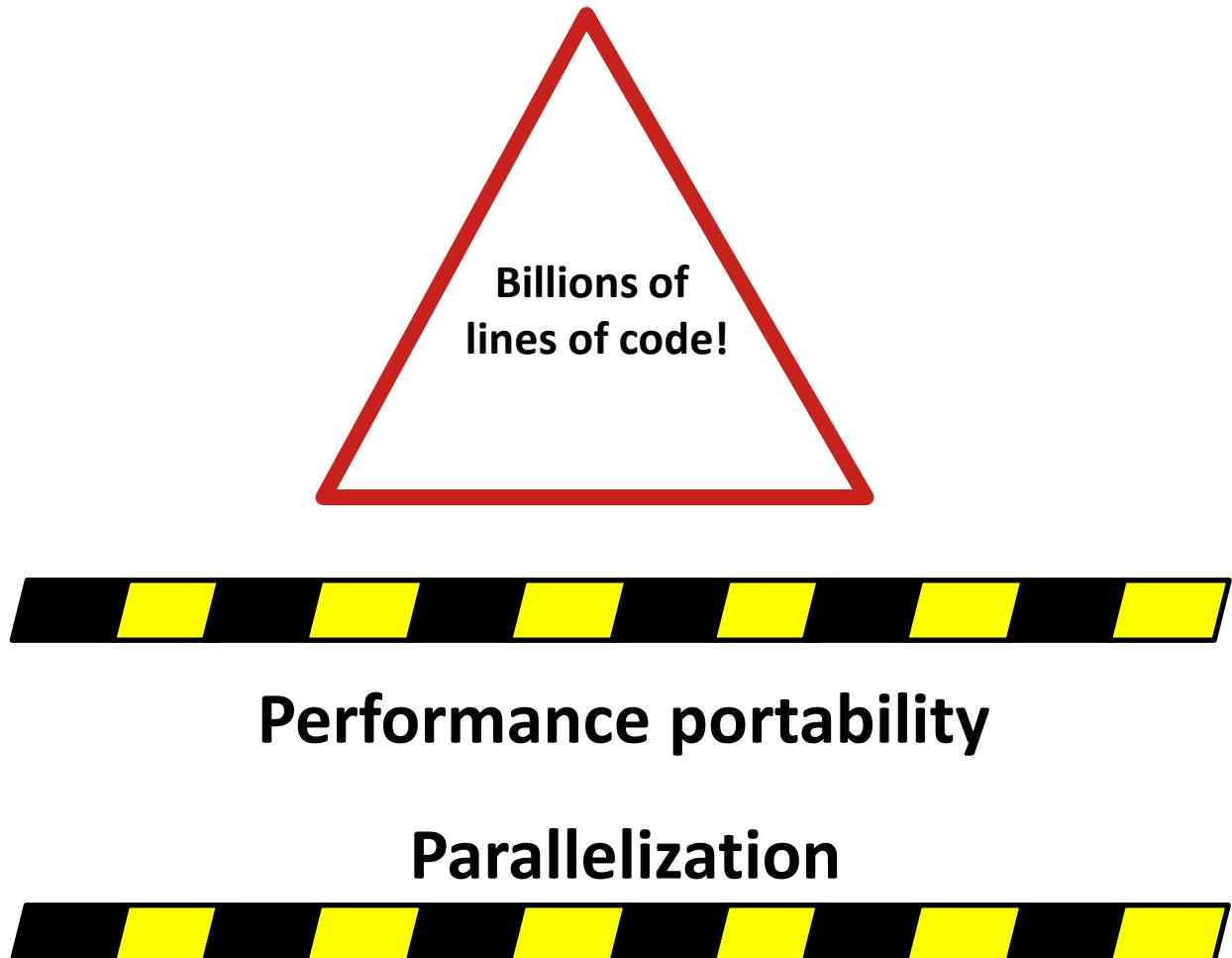
Data Centric Parallel Programming



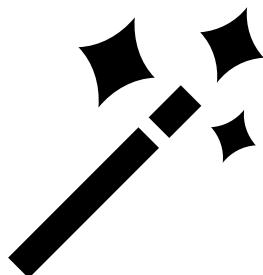
What about C?



TIOBE Index 2022	#2	#4
TIOBE Index 2021	#1	#4
TIOBE Index 2020	#1	#4
TIOBE Index 2019	#2	#4
TIOBE Index 2018	#2	#3
TIOBE Index 2017	#2	#3
TIOBE Index 2012	#2	#3
TIOBE Index 2007	#2	#3
TIOBE Index 2002	#2	#3
TIOBE Index 1997	#1	#2
TIOBE Index 1992	#1	#2
TIOBE Index 1987	#1	#5



What about C?

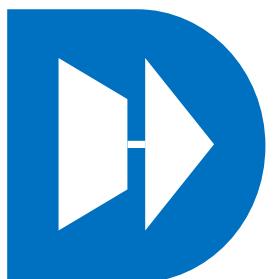
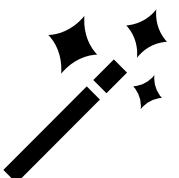


- Autoparallelization
- Generate GPU code
- Generate FPGA code
- Improve data movement
 - Apply tiling
 - Reorder loops

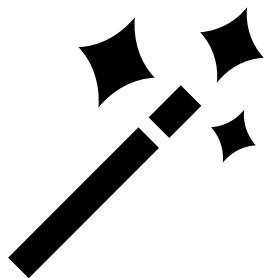
```
static void kernel_jacobi_1d(int tsteps,
                             int n,
                             double A[2000 + 0],
                             double B[2000 + 0])

{
    int t, i;

    for (t = 0; t < tsteps; t++)
    {
        #pragma omp parallel for
        for (i = 1; i < n - 1; i++)
        {
            B[i] = 0.33333 * (A[i - 1] + A[i] + A[i + 1]);
        }
        #pragma omp parallel for
        for (i = 1; i < n - 1; i++)
        {
            A[i] = 0.33333 * (B[i - 1] + B[i] + B[i + 1]);
        }
    }
}
```



What about C?

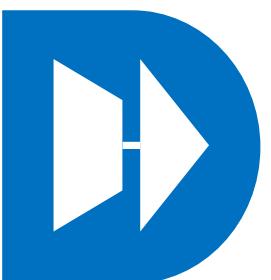
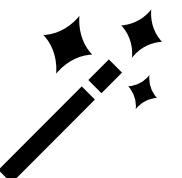


- Autoparallelization
- Generate GPU code ✓
- Generate FPGA code ✓
- Improve data movement
 - Apply tiling ✓
 - Reorder loops ✓

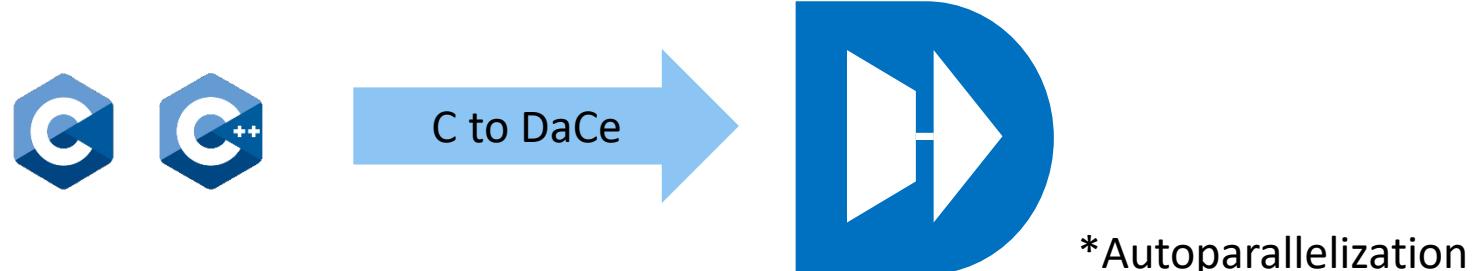
```
static void kernel_jacobi_1d(int tsteps,
                             int n,
                             double A[2000 + 0],
                             double B[2000 + 0])

{
    int t, i;

    for (t = 0; t < tsteps; t++)
    {
        #pragma omp parallel for
        for (i = 1; i < n - 1; i++)
        {
            B[i] = 0.33333 * (A[i - 1] + A[i] + A[i + 1]);
        }
        #pragma omp parallel for
        for (i = 1; i < n - 1; i++)
        {
            A[i] = 0.33333 * (B[i - 1] + B[i] + B[i + 1]);
        }
    }
}
```



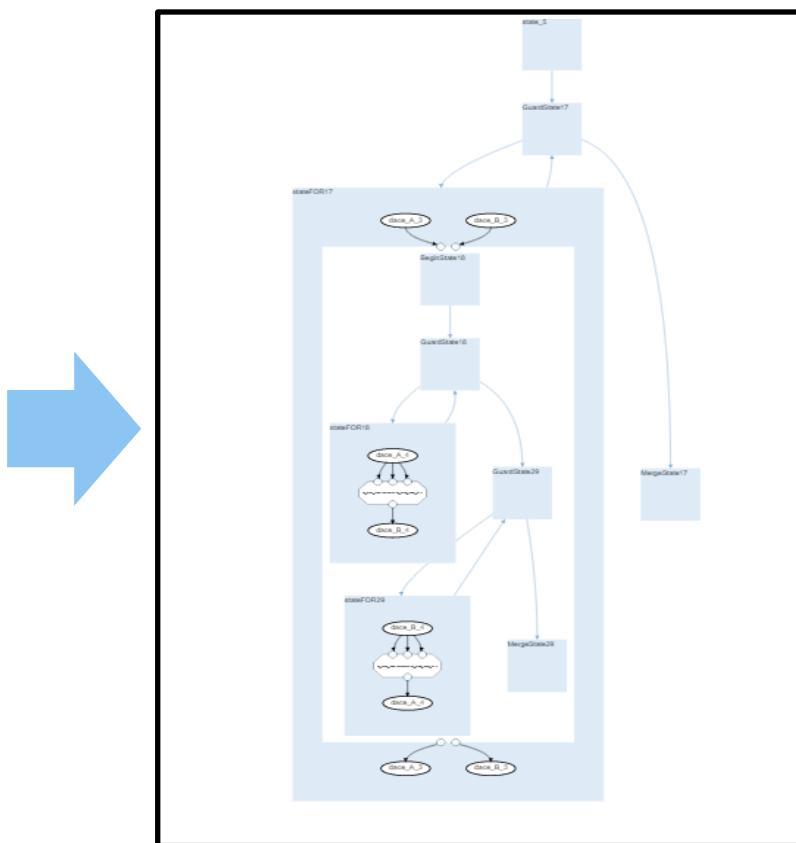
What about C?



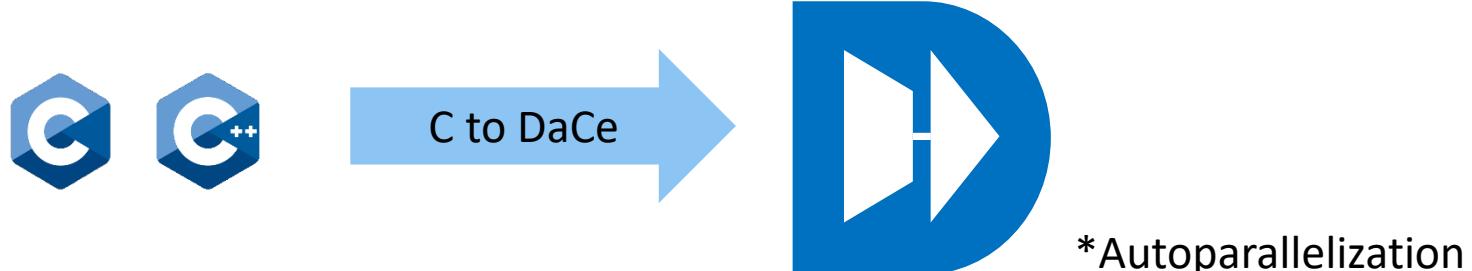
```
static void kernel_jacobi_1d(int tsteps,
                            int n,
                            double A[2000 + 0],
                            double B[2000 + 0])
{
    int t, i;

    for (t = 0; t < tsteps; t++)
    {
        for (i = 1; i < n - 1; i++)
        {
            B[i] = 0.33333 * (A[i - 1] + A[i] + A[i + 1]);
        }

        for (i = 1; i < n - 1; i++)
        {
            A[i] = 0.33333 * (B[i - 1] + B[i] + B[i + 1]);
        }
    }
}
```



What about C?

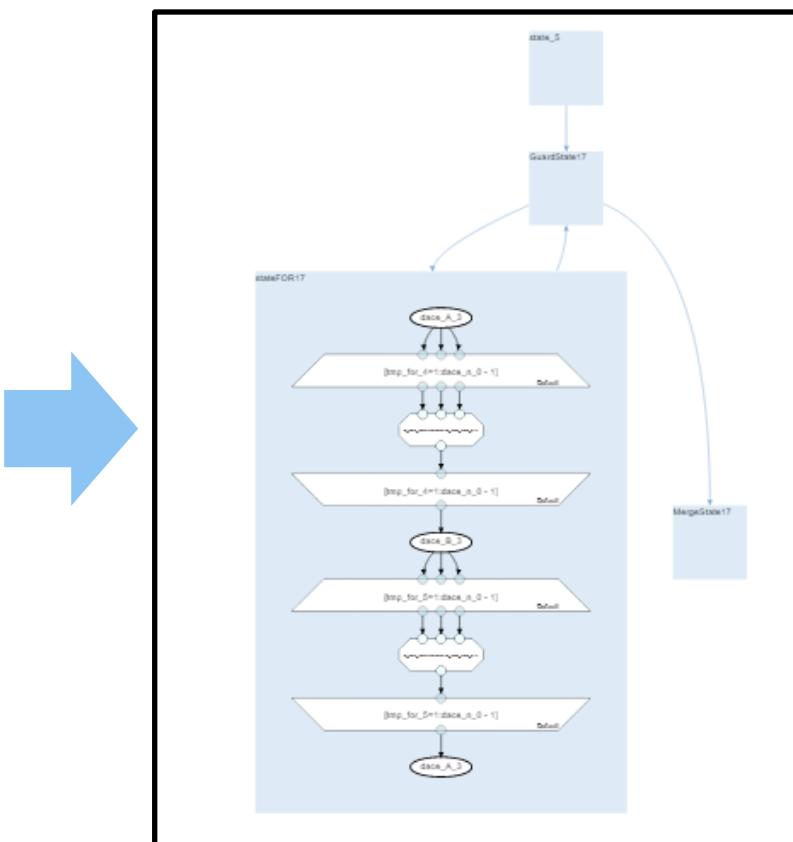


*Autoparallelization

```
static void kernel_jacobi_1d(int tsteps,
                            int n,
                            double A[2000 + 0],
                            double B[2000 + 0])
{
    int t, i;

    for (t = 0; t < tsteps; t++)
    {
        for (i = 1; i < n - 1; i++)
        {
            B[i] = 0.33333 * (A[i - 1] + A[i] + A[i + 1]);
        }

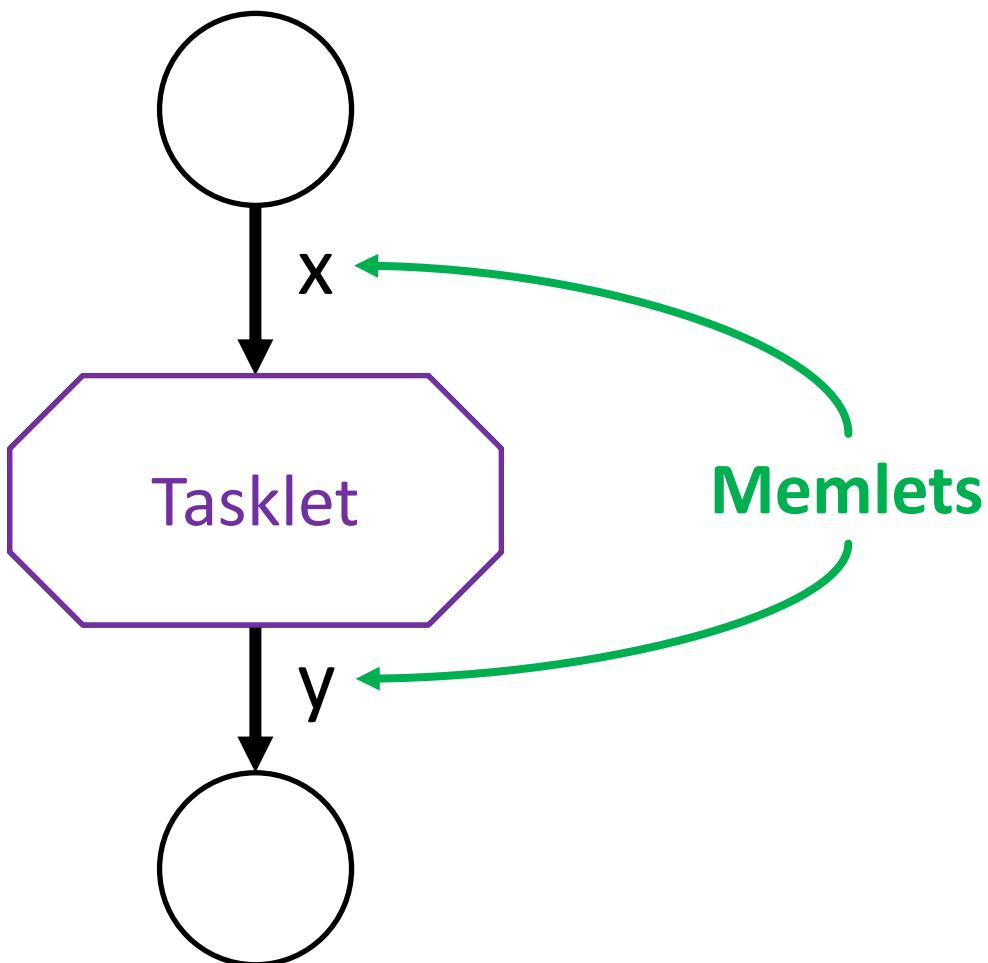
        for (i = 1; i < n - 1; i++)
        {
            A[i] = 0.33333 * (B[i - 1] + B[i] + B[i + 1]);
        }
    }
}
```



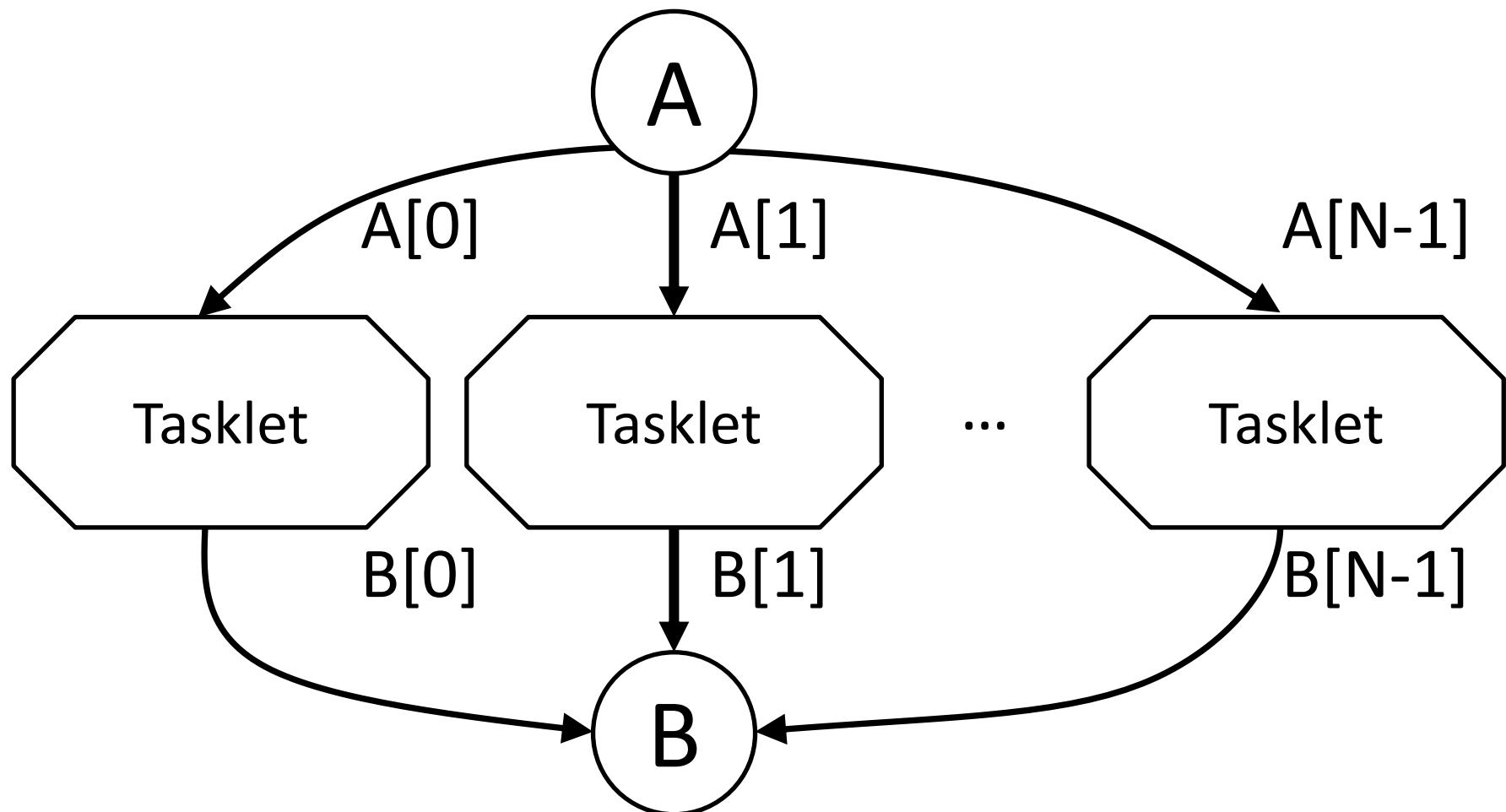
Dataflow Programming in DaCe

$$y = x^2 + \sin \frac{x}{\pi}$$

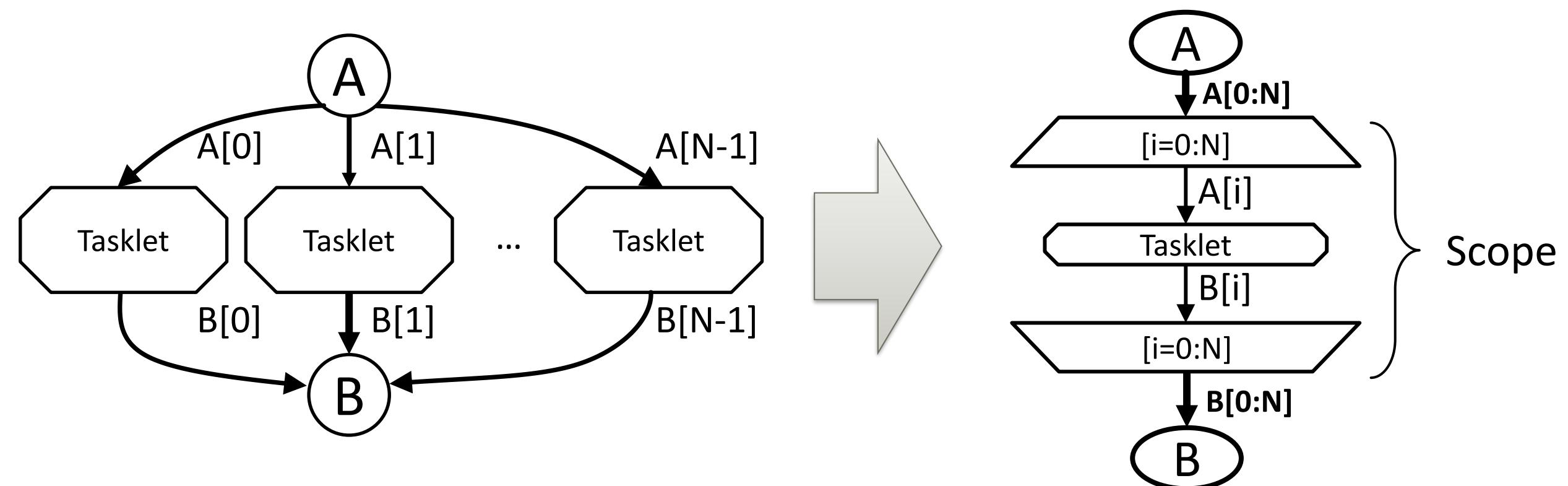
Dataflow Programming in DaCe



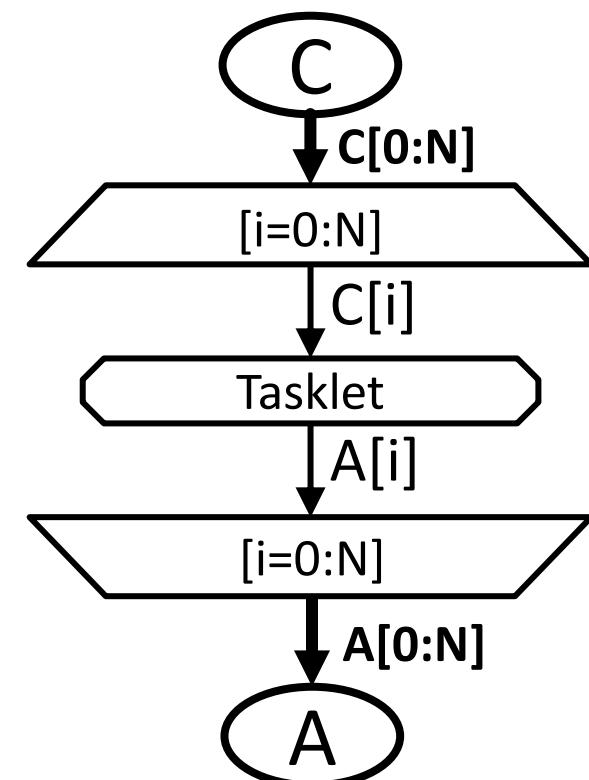
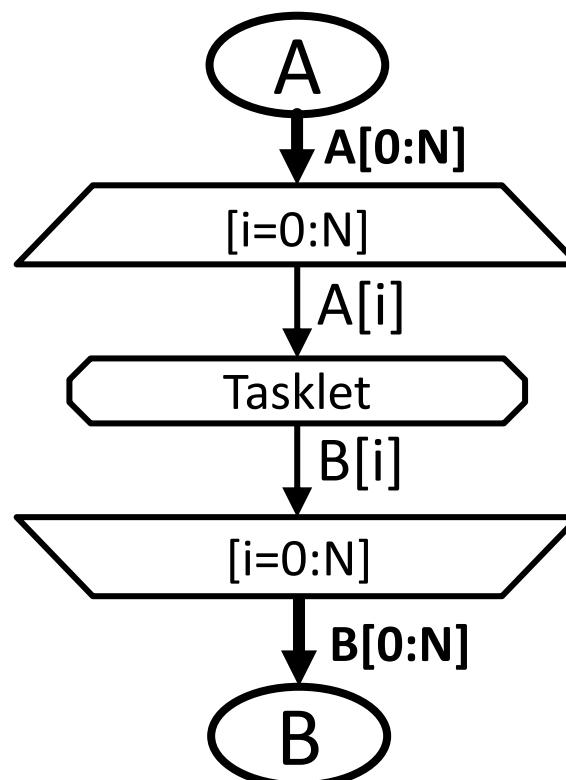
Parallel Dataflow Programming



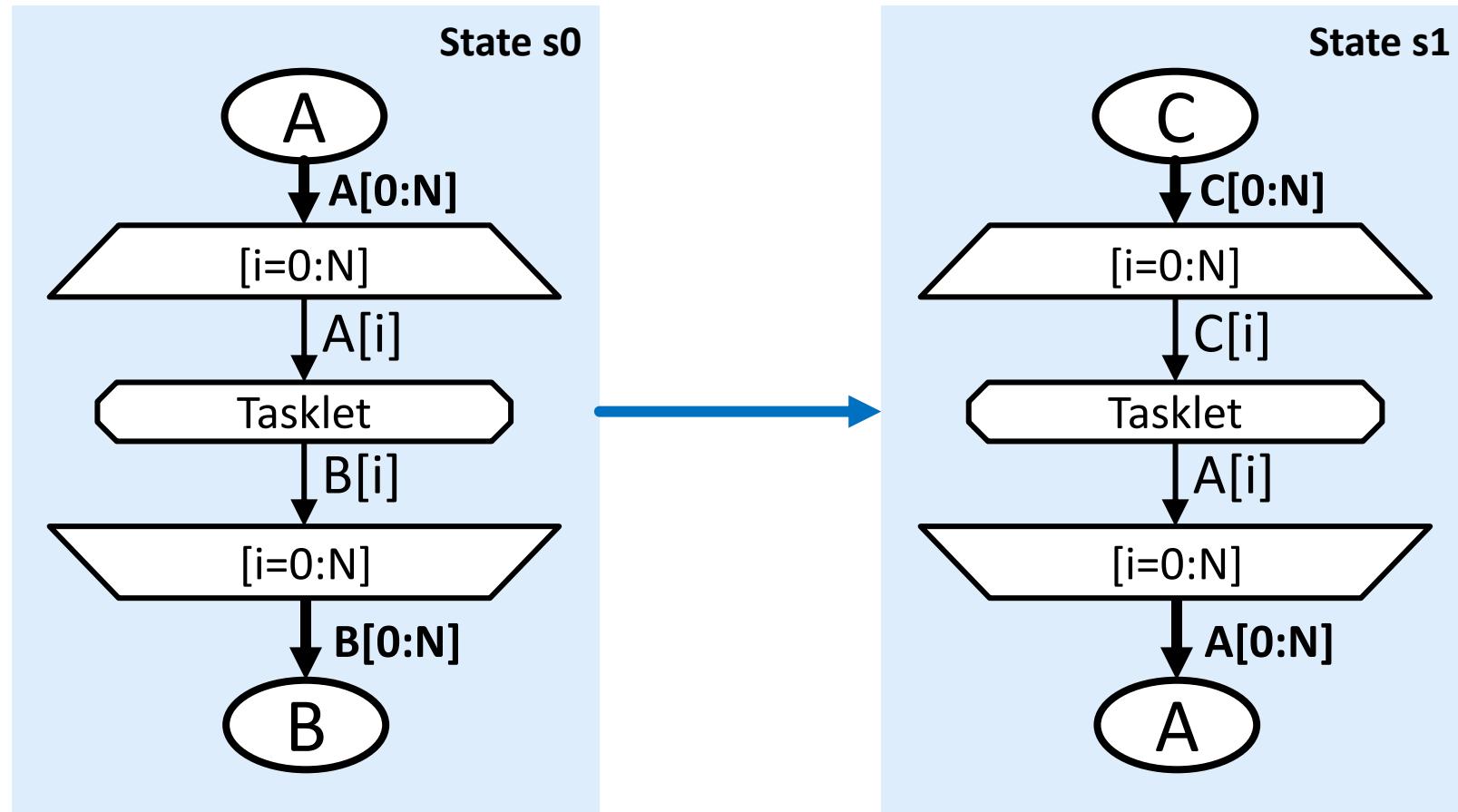
Parallel Dataflow Programming



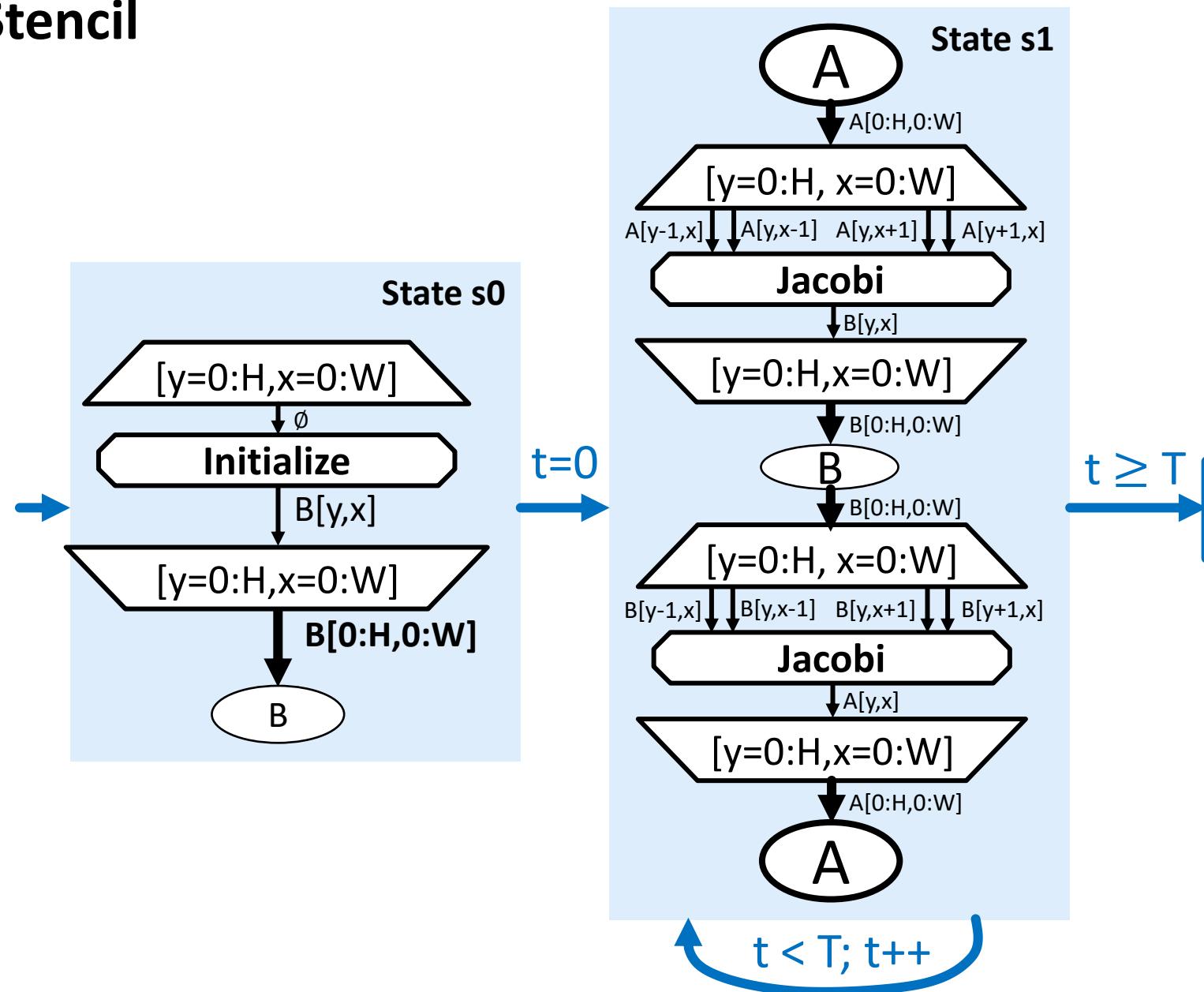
Stateful Parallel Dataflow Programming



Stateful Parallel Dataflow Programming



Example: 2D Stencil



C to DaCe: SpMV

1. AST Transformations

```
for (int i = 0; i < N ; i++)  
  
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++)  
  
        y[i] += A[col_idx[j]] * x[j];
```

C to DaCe: SpMV

1. AST Transformations

- Make basic blocks explicit

```
for (int i = 0; i < N ; i++){  
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){  
  
        y[i] += A[col_idx[j]] * x[j];  
    }  
}
```

C to DaCe: SpMV

1. AST Transformations

- Make basic blocks explicit
- Extract array indices*
- And many others...

```
for (int i = 0; i < N ; i++){  
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){  
        int idx = col_idx[j];  
        y[i] += A[ idx ] * x[j];  
    }  
}
```

AST Transformations *canonicalize* the program representation and allows the translation to SDFG to make simplifying assumptions.

*All array indices are extracted, including i , $i+1$, j . They are omitted here for simplicity and space

C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG

```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

```
for (int i = 0; i < N ; i++){

    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG

```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

```
for (int i = 0; i < N ; i++){
```

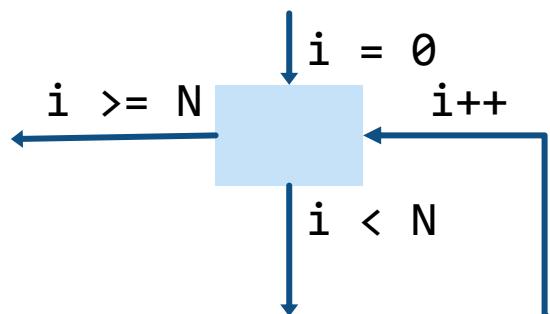
```
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

```
}
```

C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG

```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

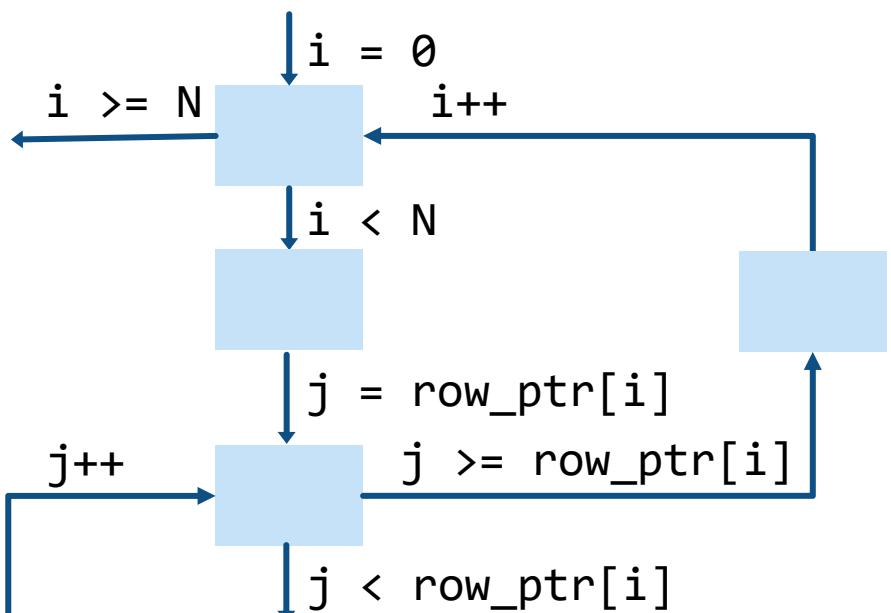


```
for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
    int idx = col_idx[j];
    y[i] += A[ idx ] * x[j];
}
```

C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG

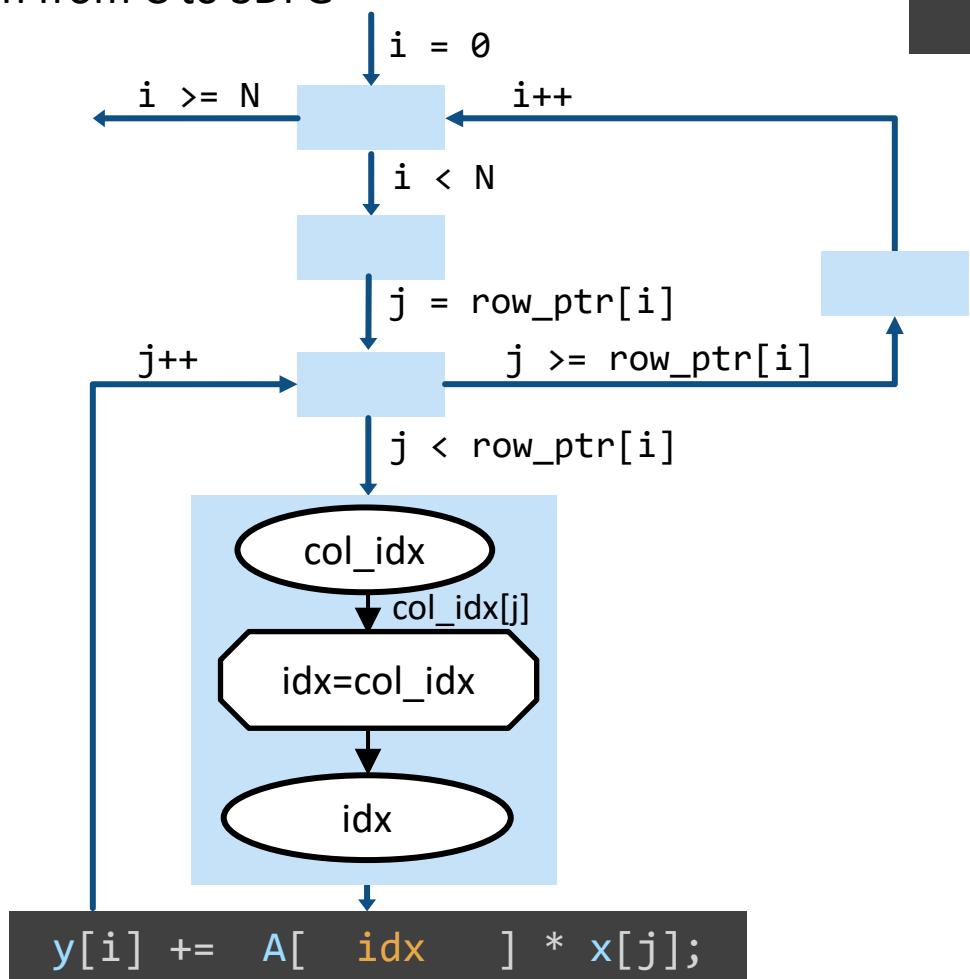
```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```



```
int idx = col_idx[j];
y[i] += A[ idx ] * x[j];
```

C to DaCe: SpMV

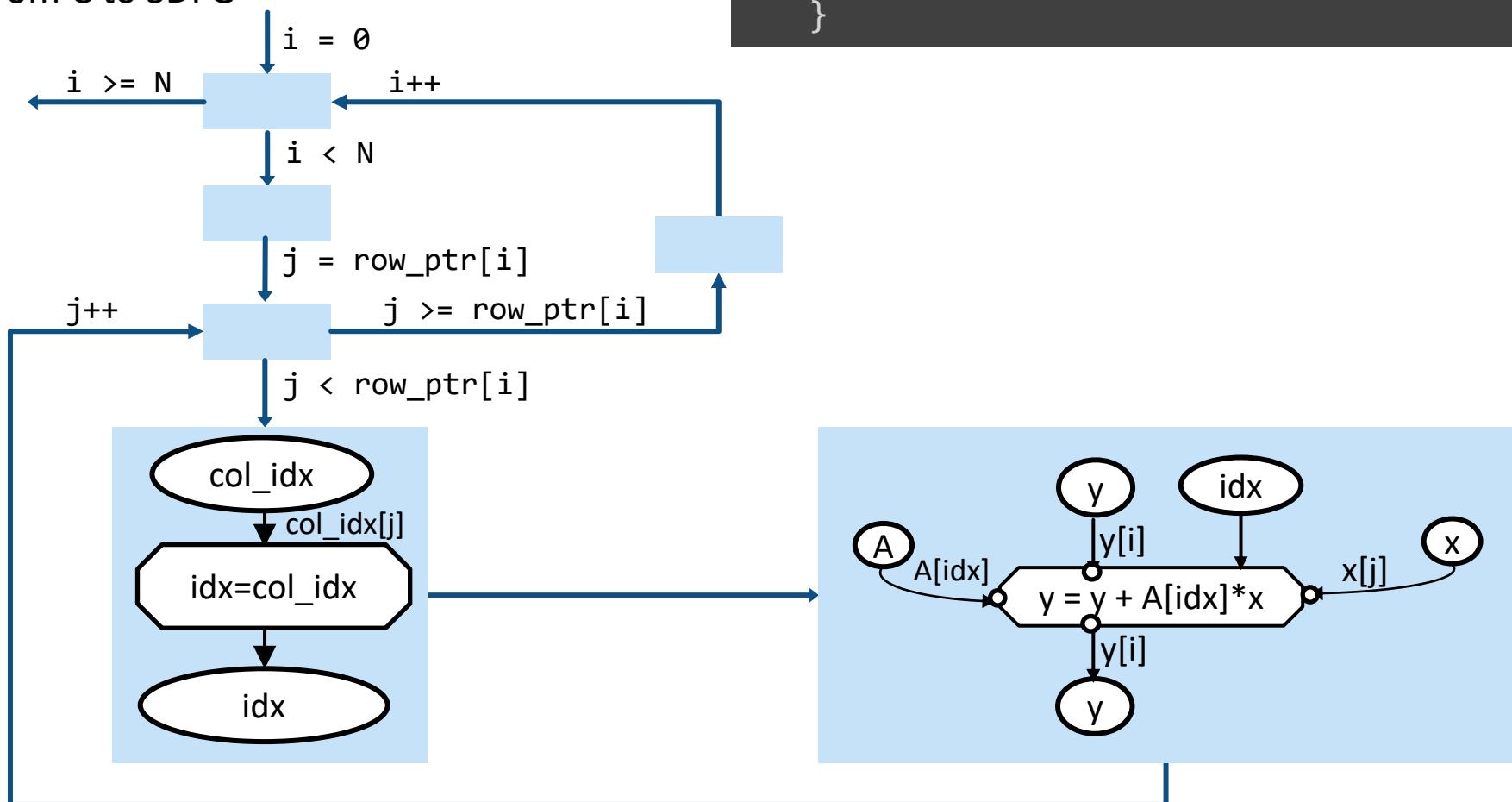
1. AST Transformations
2. Translation from C to SDFG



```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

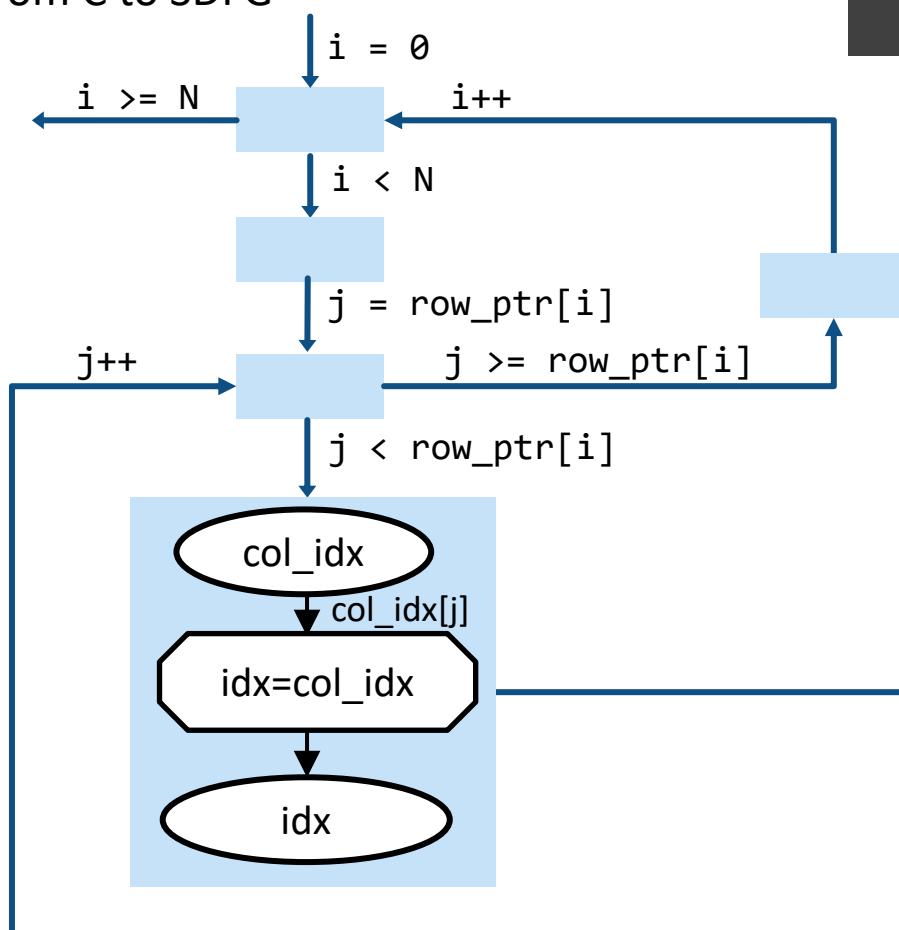
C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG



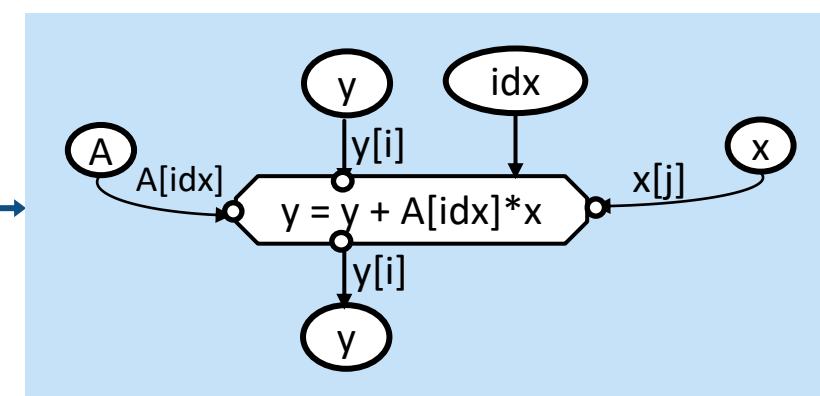
C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG



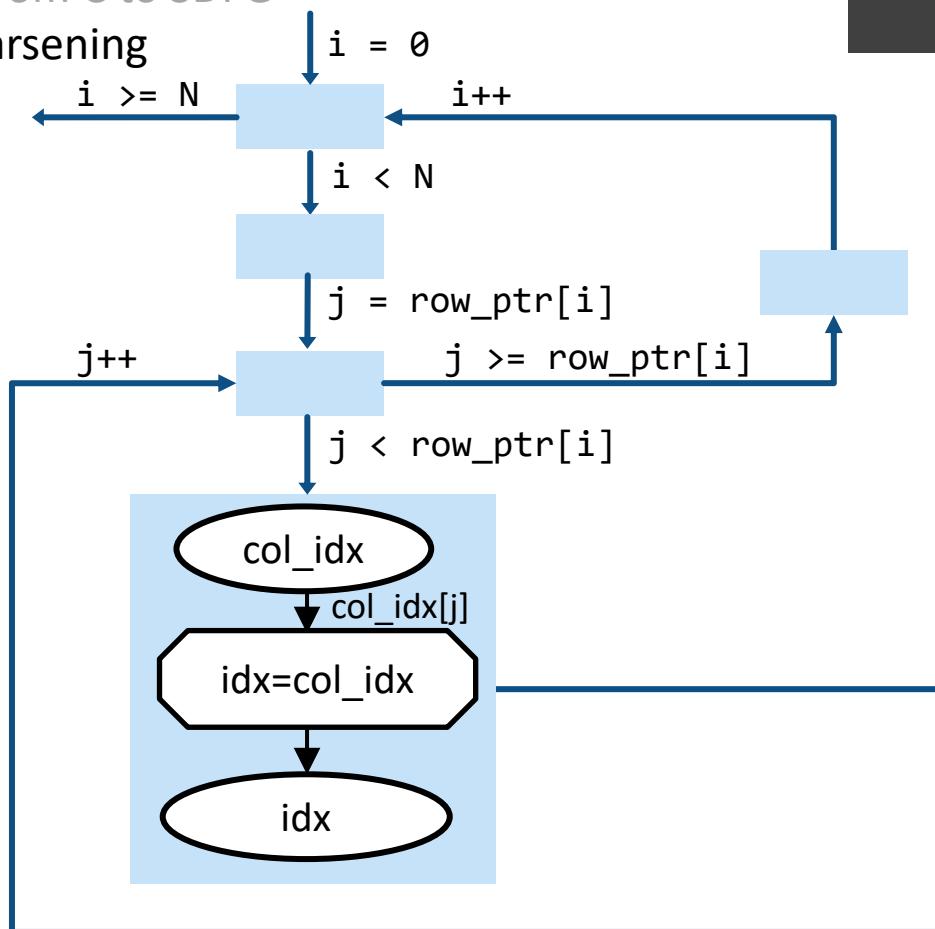
```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

The SDFG representation is at this point a direct representation of the *control-flow centric* C code.



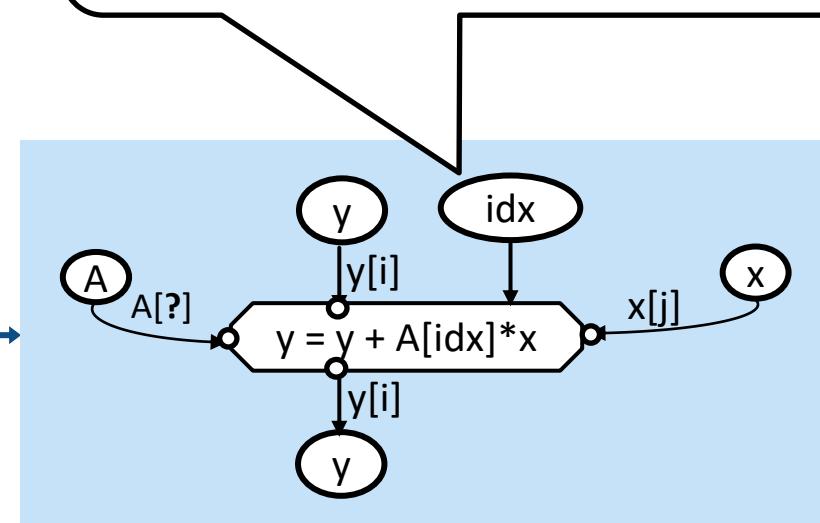
C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG
3. Dataflow coarsening

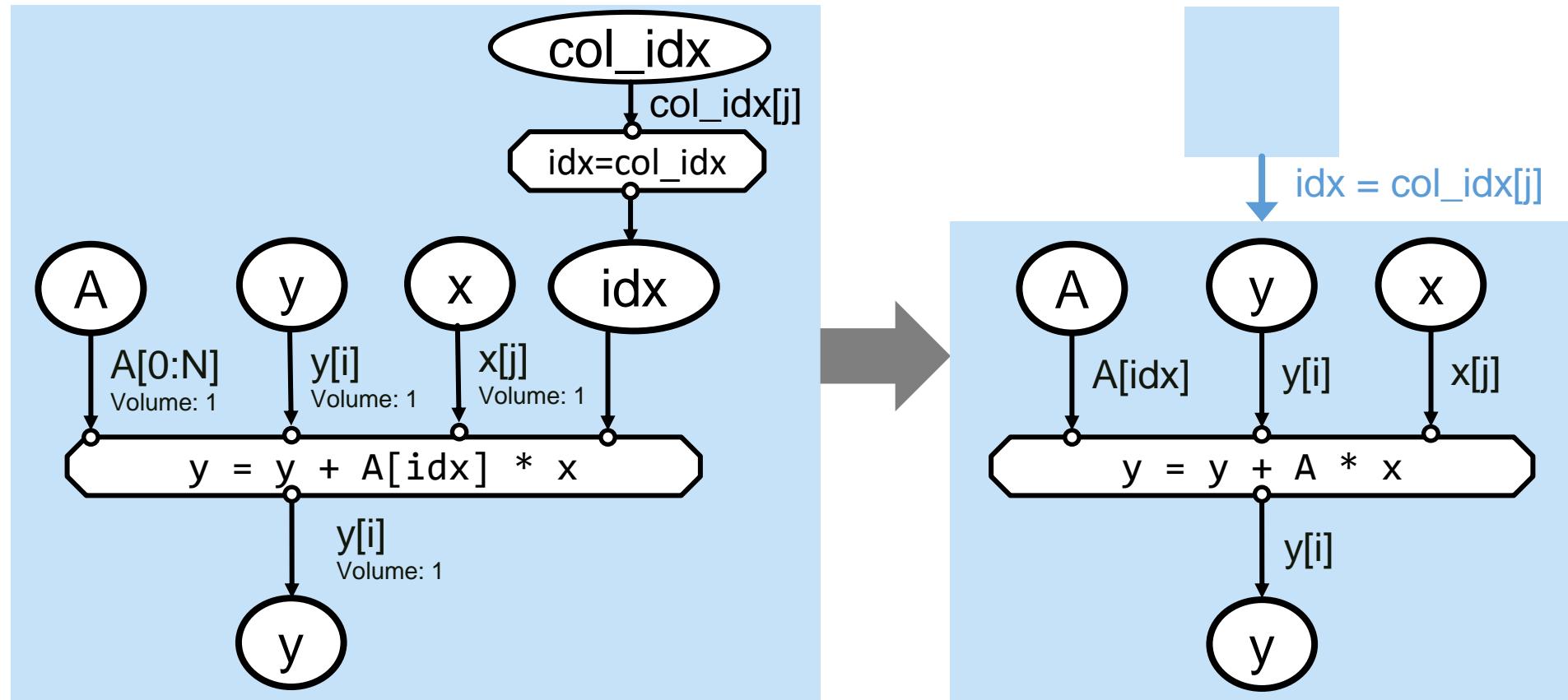


```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

Symbolic scalar analysis:
Improves understanding of data access patterns by allowing symbolic manipulation

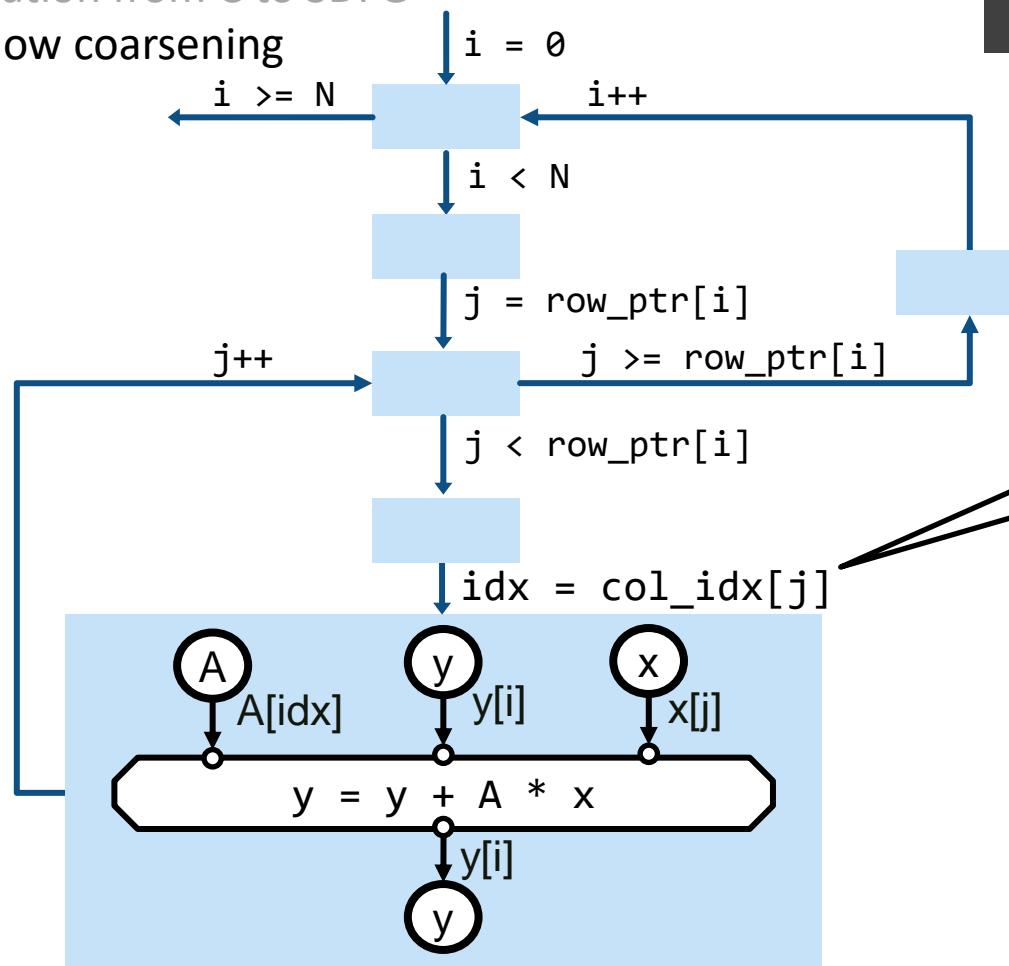


Symbolic scalar analysis



C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG
3. Dataflow coarsening

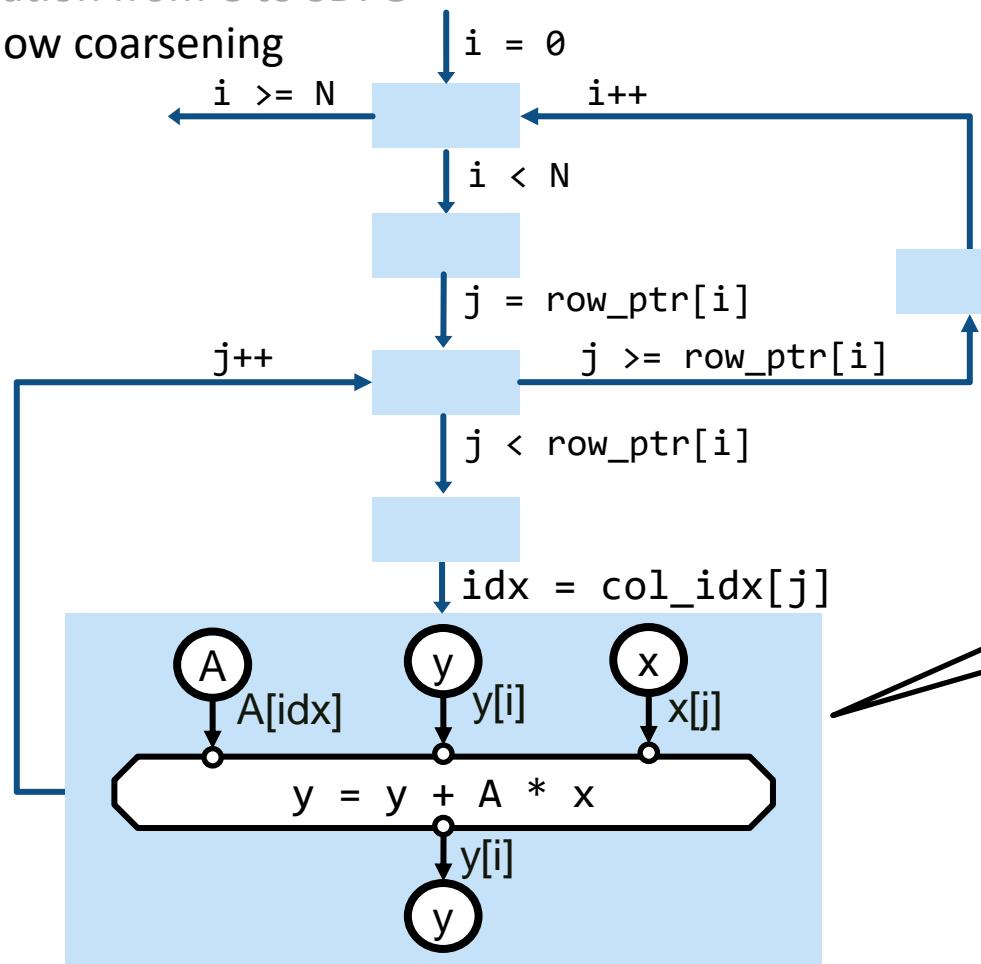


```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

Symbolic scalar analysis:
Improves understanding of data access patterns by allowing symbolic manipulation

C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG
3. Dataflow coarsening

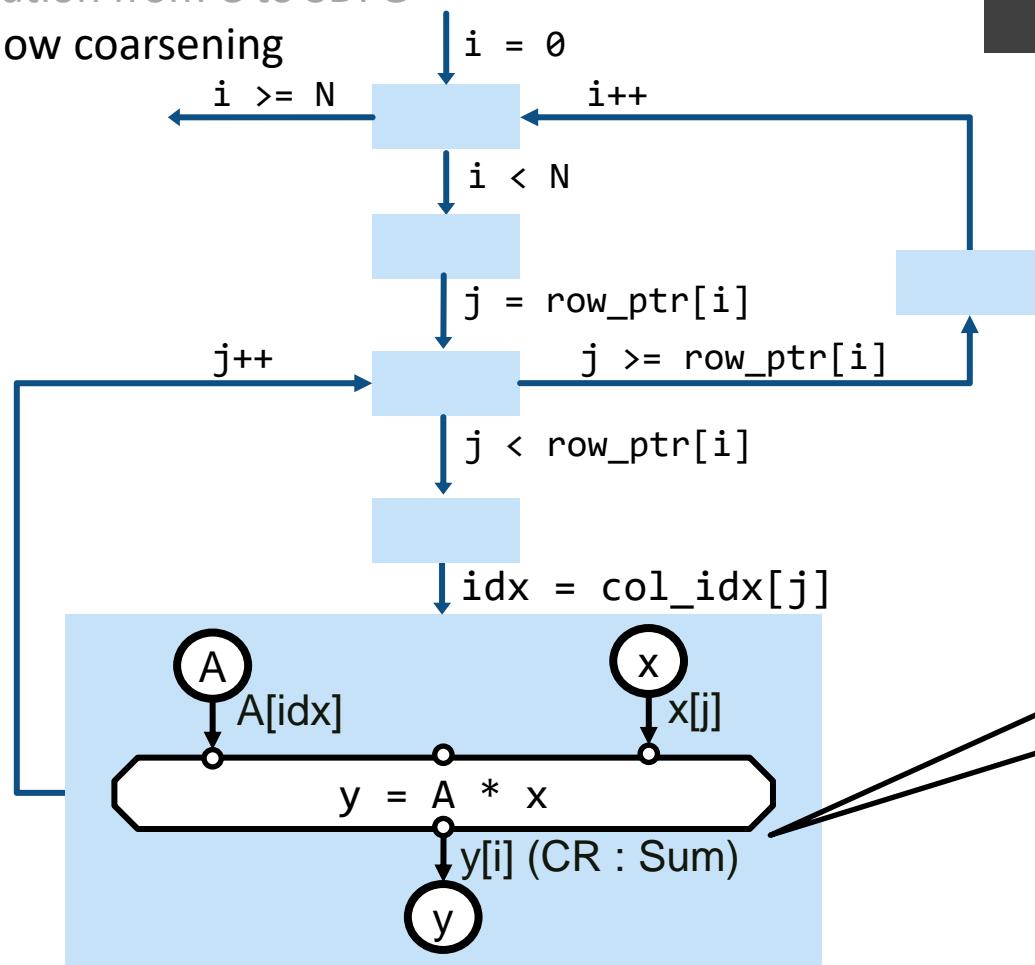


```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

Update detection

C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG
3. Dataflow coarsening

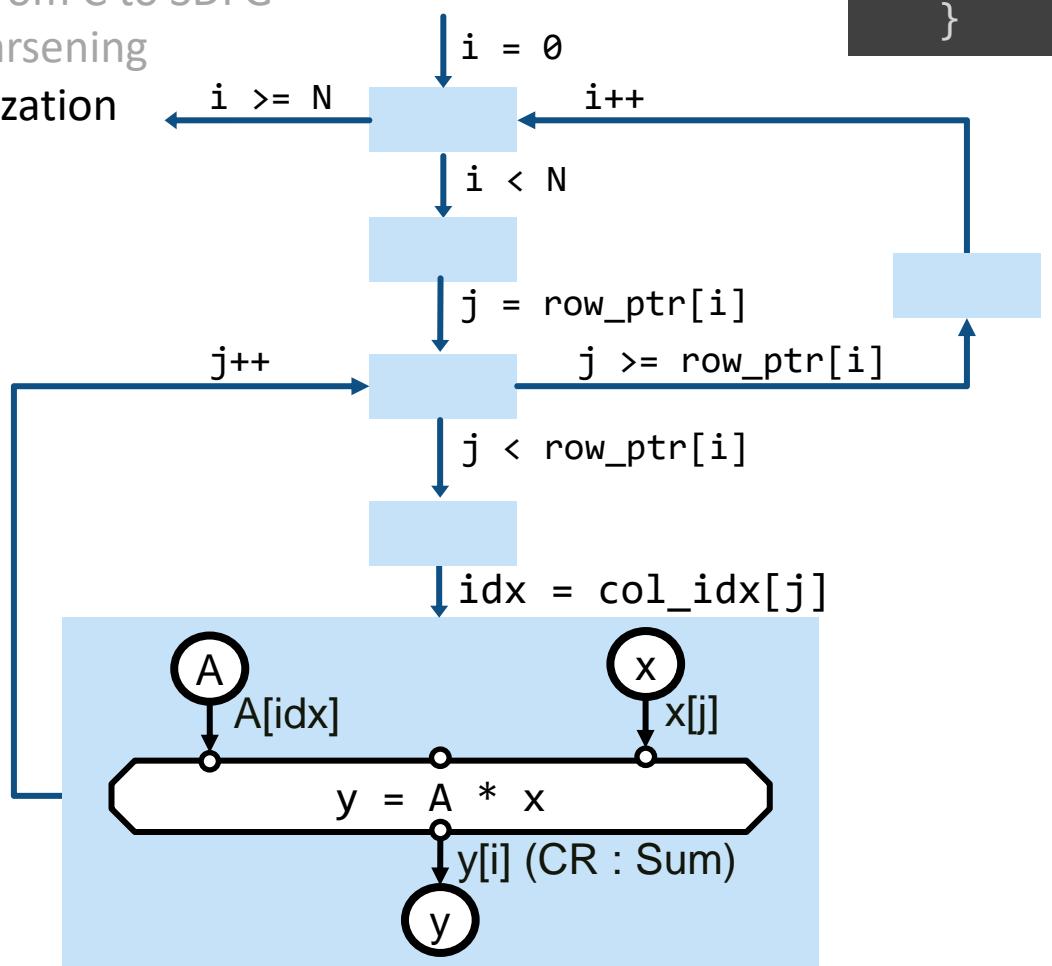


```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

Update detection:
Allows more parallelization by creating the specialized *update* type of assignment, supporting conflict resolution for multiple operations

C to DaCe: SpMV

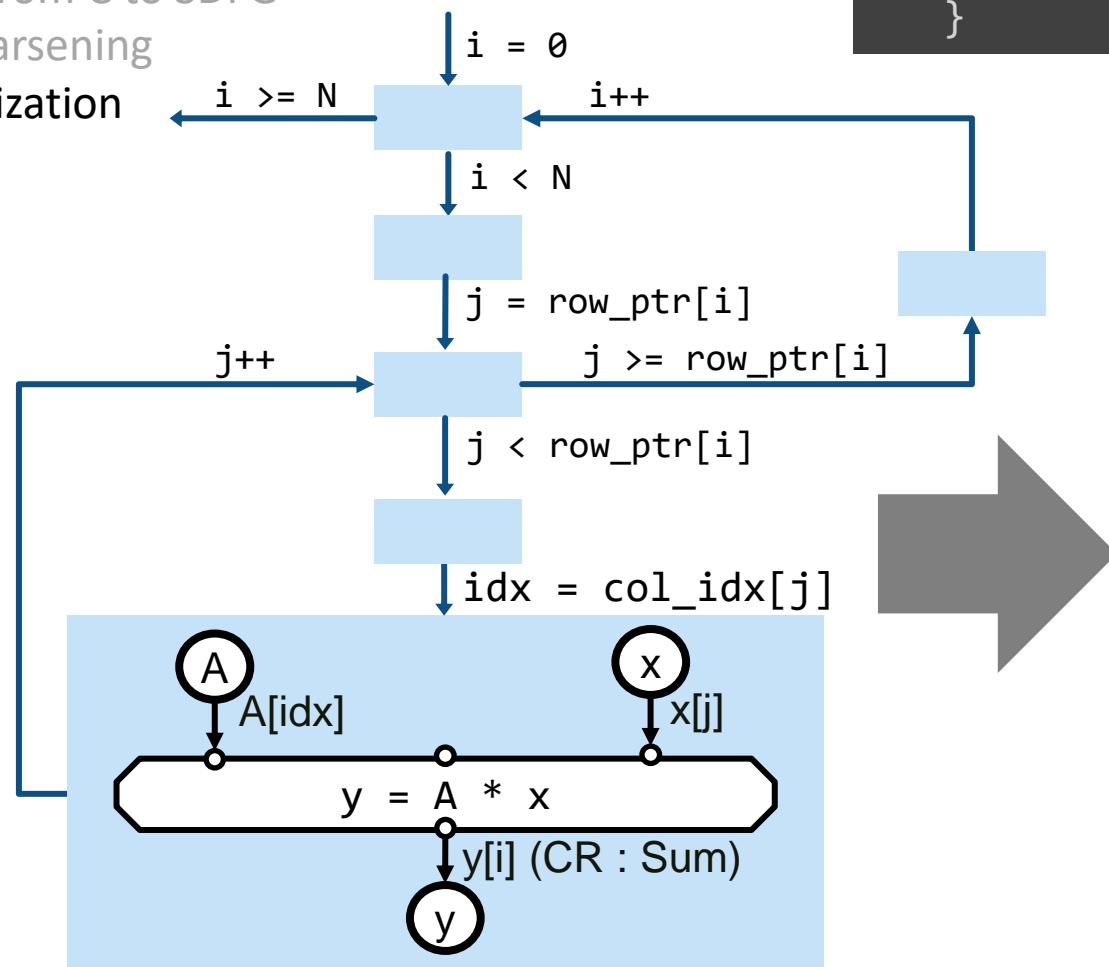
1. AST Transformations
2. Translation from C to SDFG
3. Dataflow coarsening
4. Autoparallelization



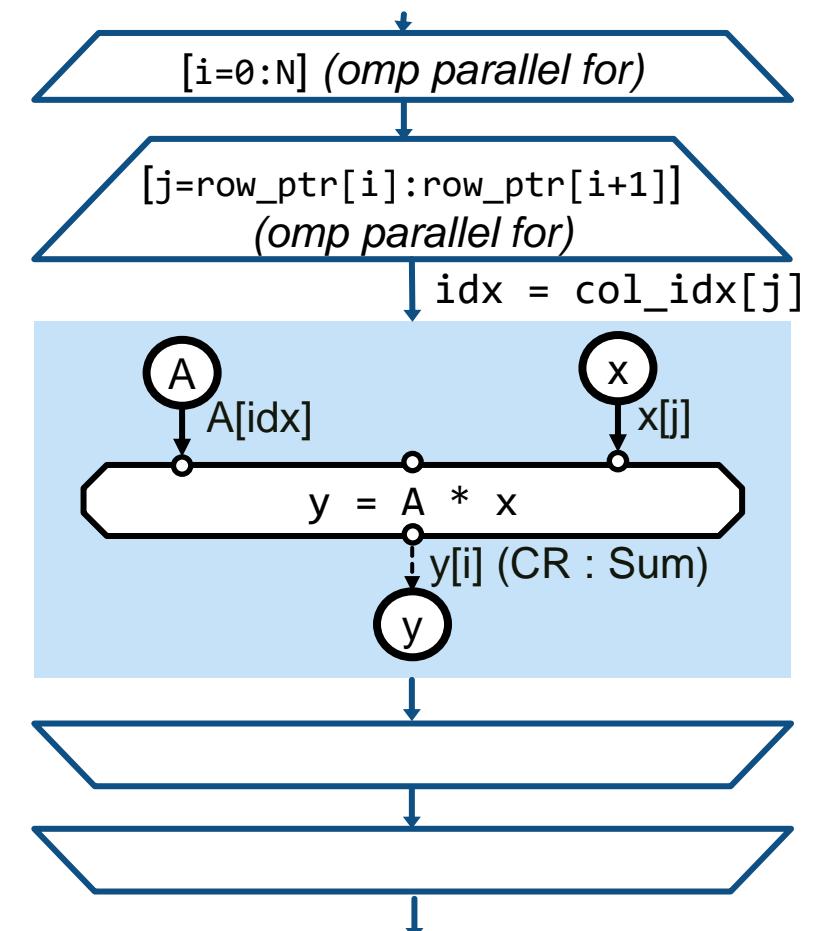
```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```

C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG
3. Dataflow coarsening
4. Autoparallelization



```
for (int i = 0; i < N ; i++){
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){
        int idx = col_idx[j];
        y[i] += A[ idx ] * x[j];
    }
}
```



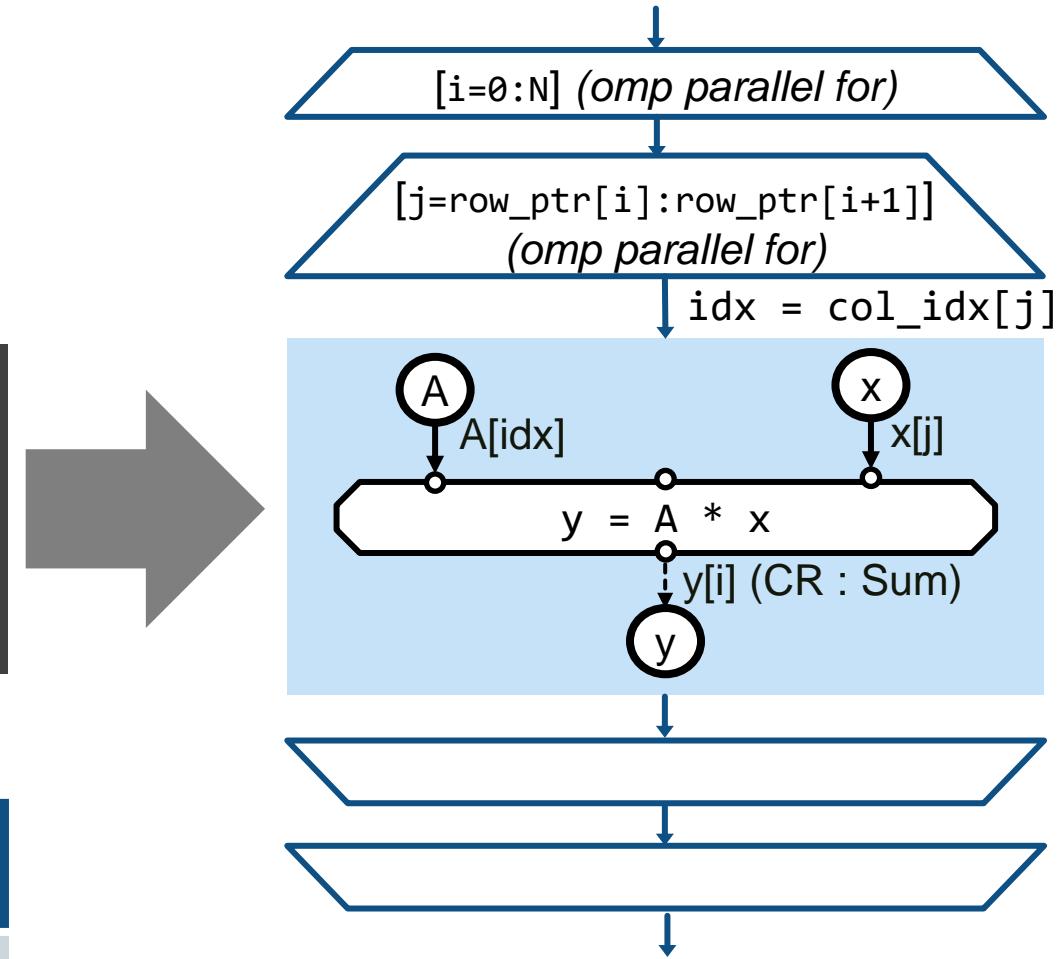
C to DaCe: SpMV

1. AST Transformations
2. Translation from C to SDFG
3. Dataflow coarsening
4. Autoparallelization

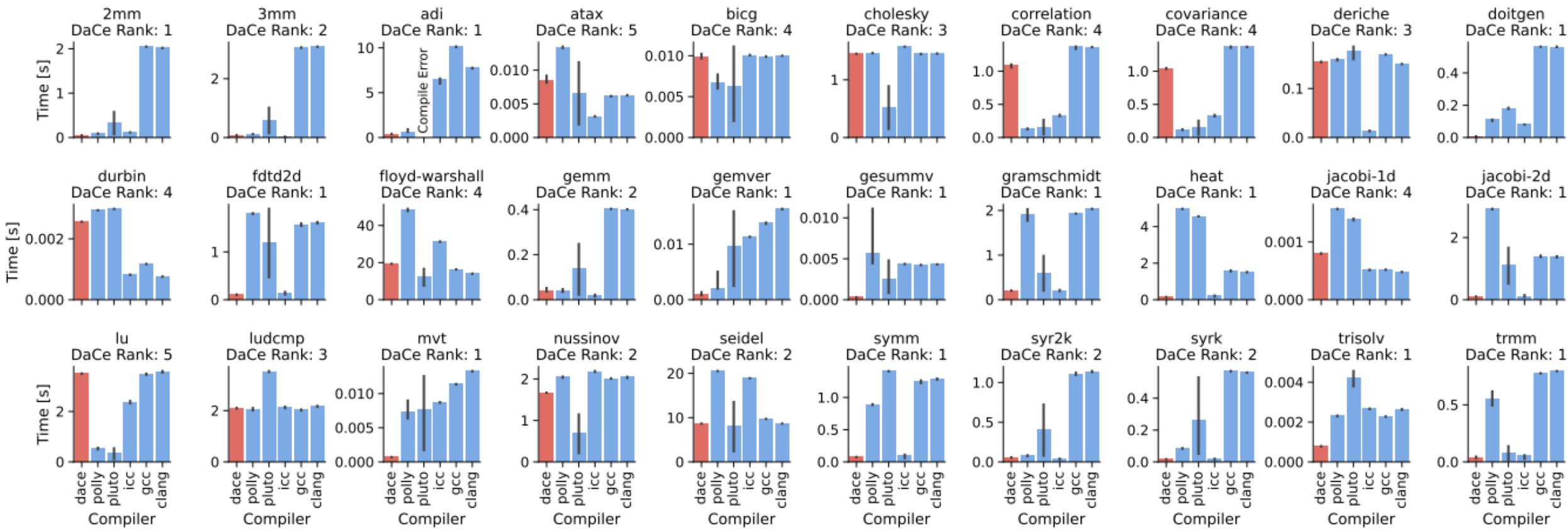
```
for (int i = 0; i < N ; i++){  
    for (int j = row_ptr[i]; j < row_ptr[i+1]; j++){  
        int idx = col_idx[j];  
        y[i] += A[ idx ] * x[j];  
    }  
}
```

Performance on up to 36 cores

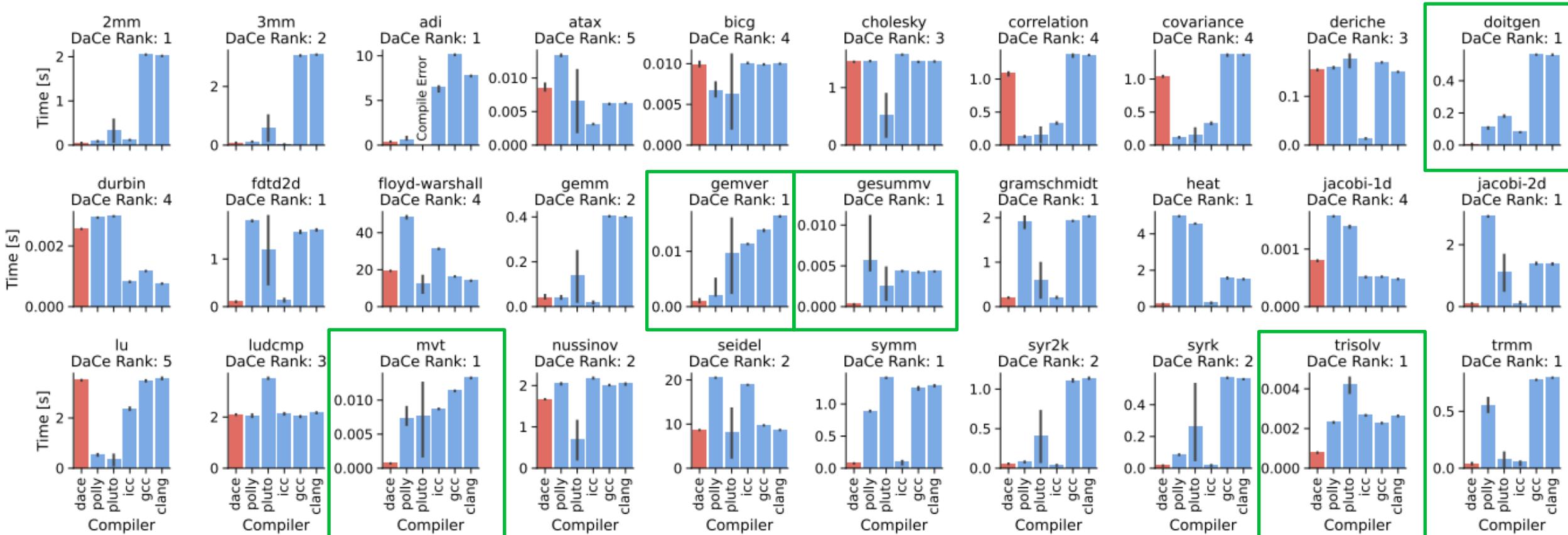
DaCe	polly	pluto	icc - parallel	icc	gcc	clang
0.54s	3.95s	failed	0.45s	4.98s	3.55s	4.19s



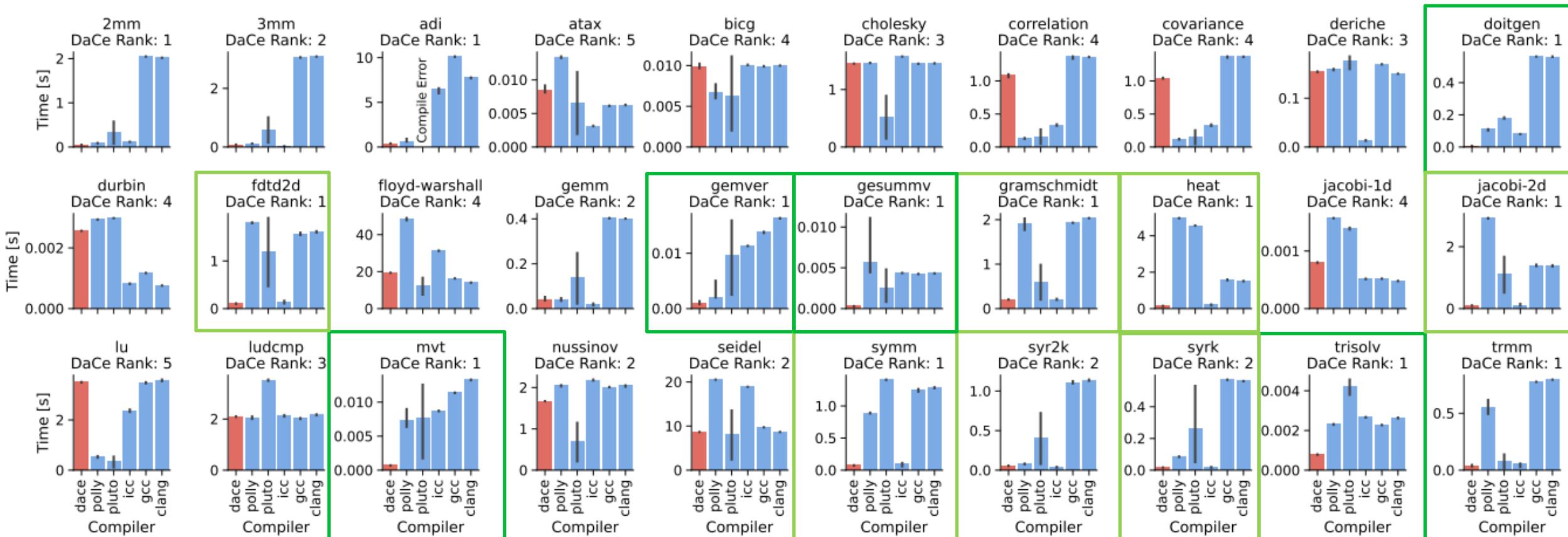
Results - Polybench



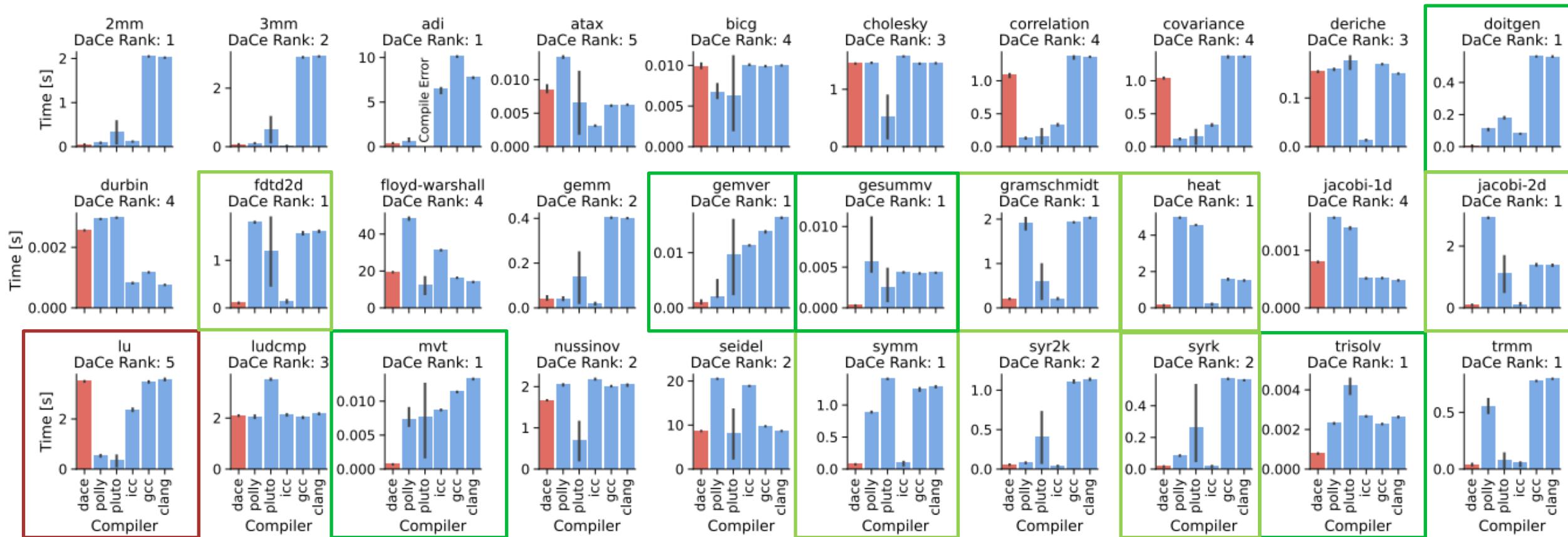
Results - Polybench



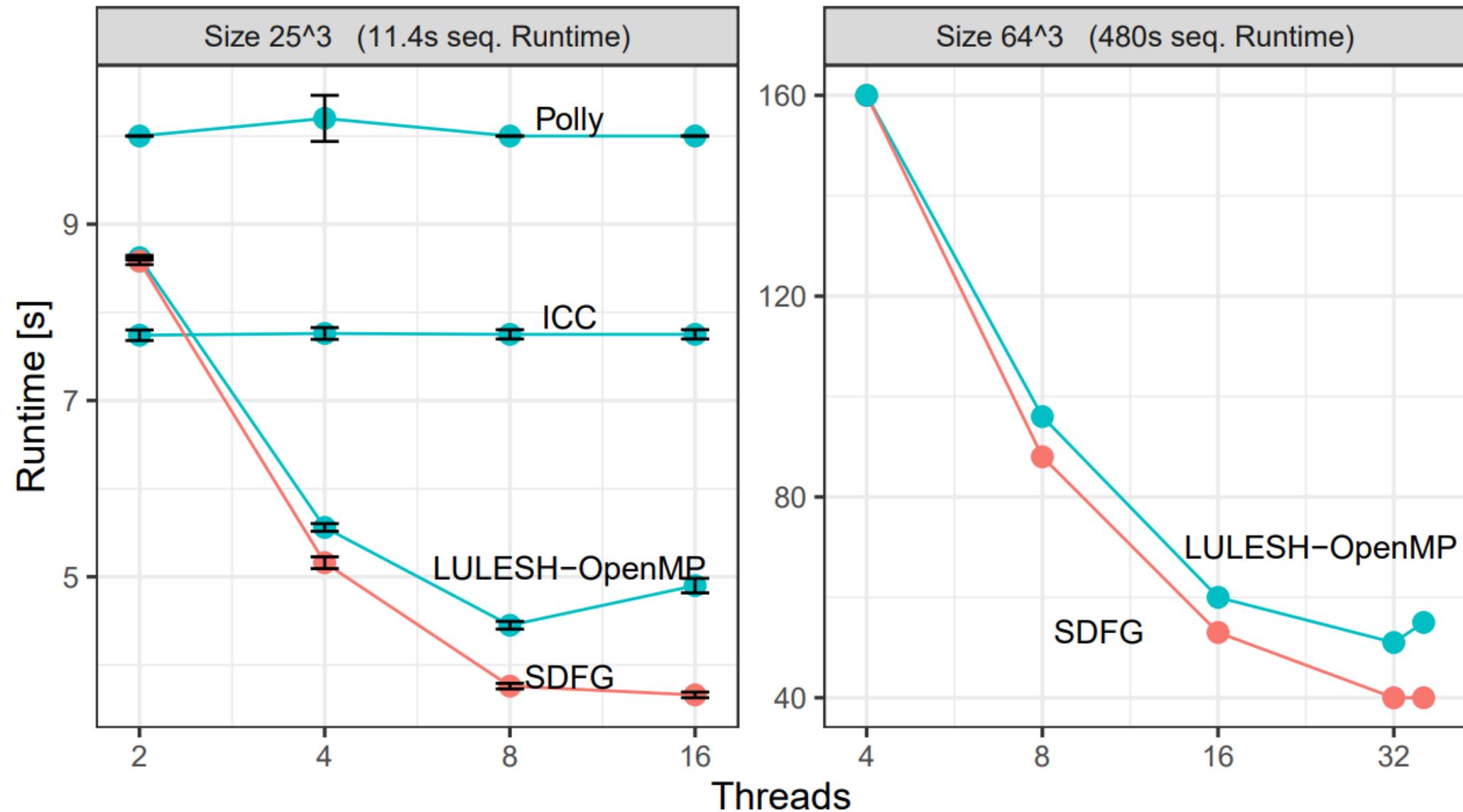
Results - Polybench



Results - Polybench

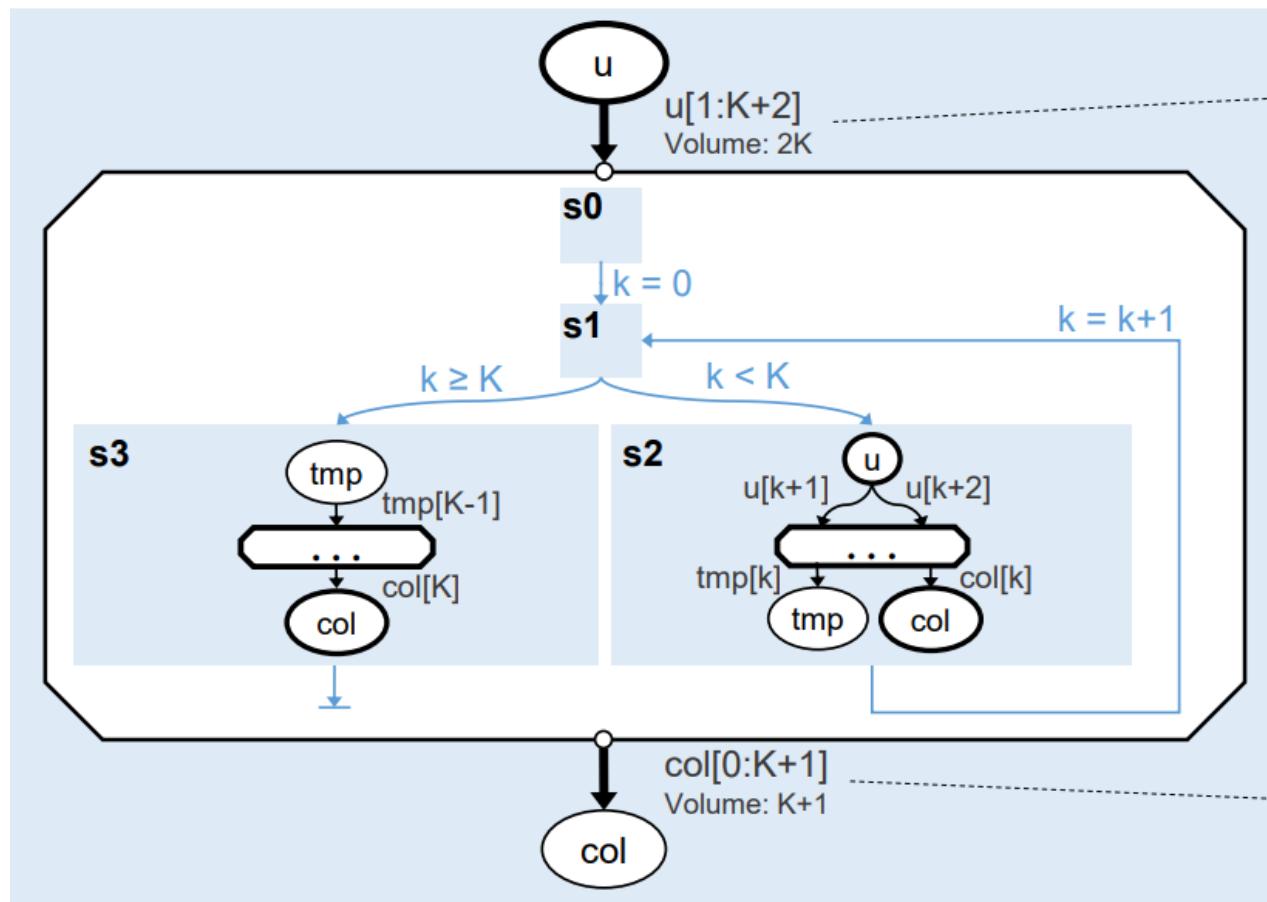


Results - LULESH



Thank you!

Access pattern propagation



Subset: $\bigcup_{k=0}^K \{k + 1, k + 2\}$
Volume: $2 \cdot \text{exec}(s_2)$

s_0 :
Executions = 1
Symbol Ranges = \emptyset

s_1, s_2 :
Executions = K
Symbol Ranges = $\{k: 0:K\}$

s_3 :
Executions = 1
Symbol Ranges = $\{k: K\}$

Subset: $\bigcup_{k=0}^K \{k\} \cup \{K\}$
Volume: $1 \cdot \text{exec}(s_2) + 1 \cdot \text{exec}(s_3)$

Enforcing order for data-less dependencies

- Global dictionary of states
- By default, all calls within a library should share a state
- Lists of stateless functions/libraries avoid needless complication of the SDFG

