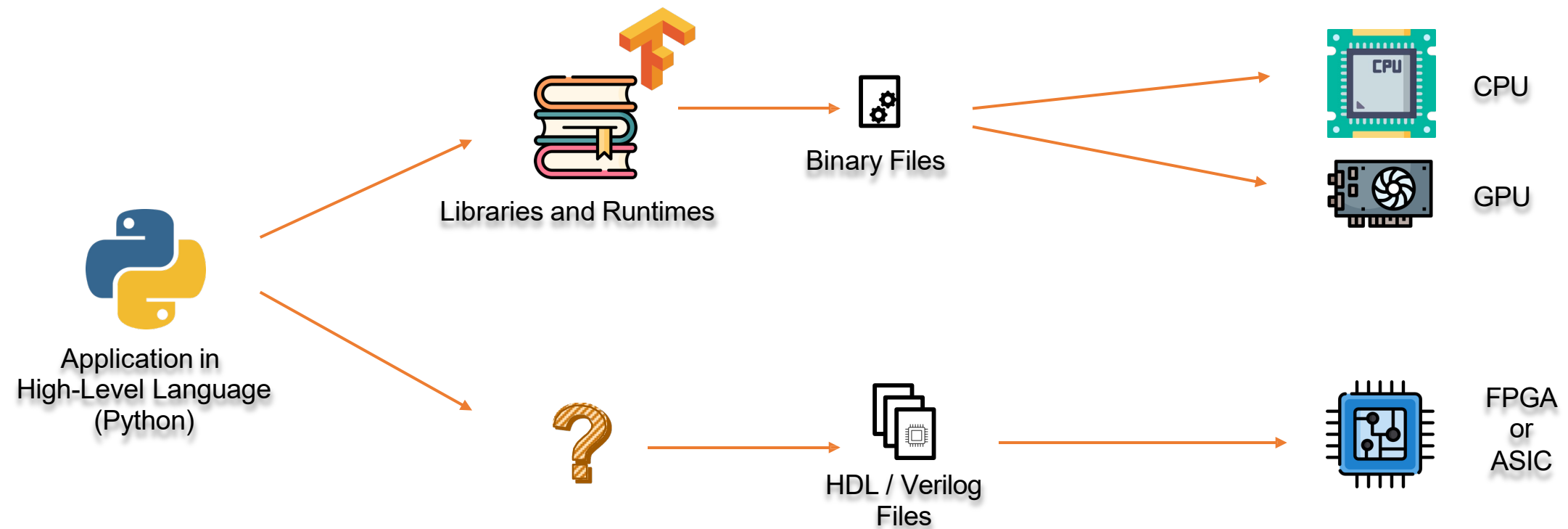


MLIR-based Compiler Flow for System-Level Design and Hardware Acceleration

Based on ICCAD '22 paper and ISCA '24 Tutorial
<https://dl.acm.org/doi/10.1145/3508352.3549424>
<https://hpc.pnl.gov/SODA/tutorials/2024/ISCA24.html>

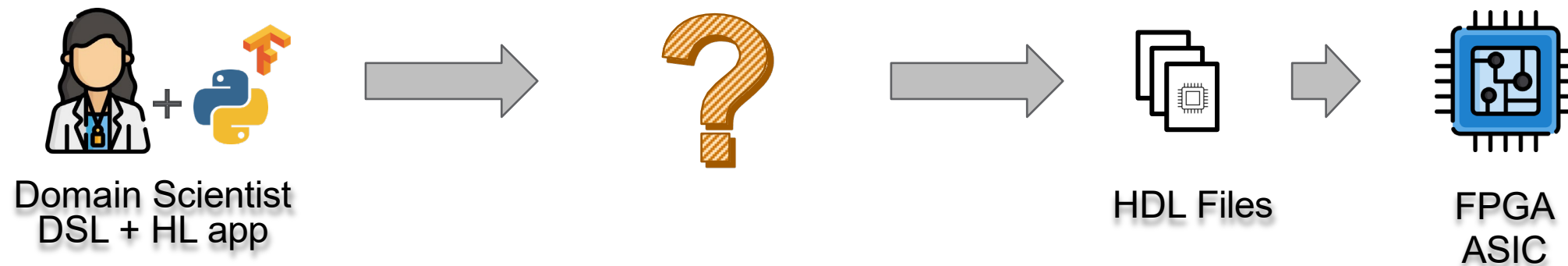
Qucheng Jiang

Agile Hardware Design and Prototyping



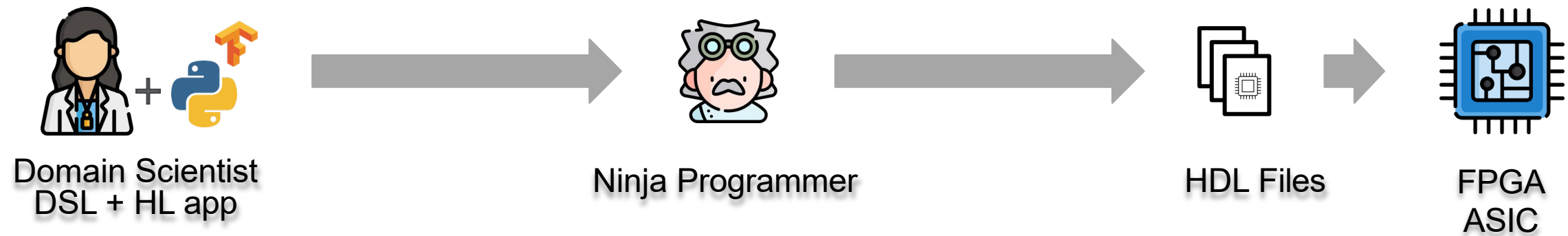
Current challenge

- It is challenging **to map** applications into custom hardware
- It is challenging **to extract performance** of the custom hardware
- Can we transform the Domain Scientist in a "**Lead User**" for Custom Domain Specific Accelerators?

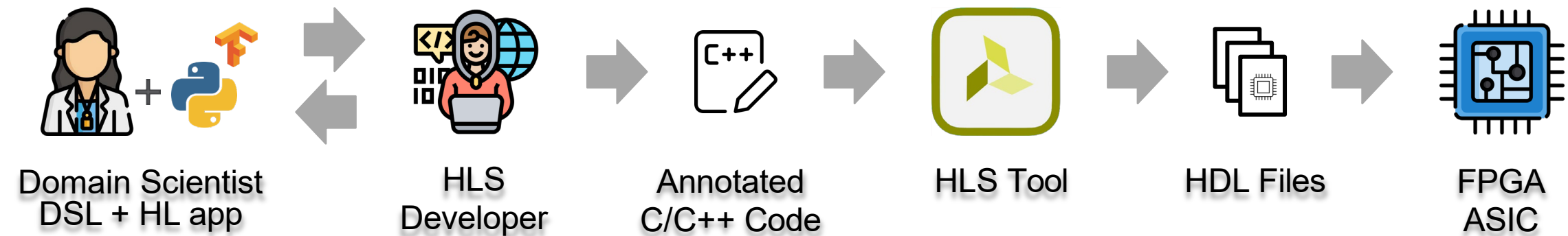


Current workflow

Conventional approach : Ninja Programmer



Conventional approach : HLS Developer



Related Work

Tool	Input	Backend tools	F1	F2	F3	F4	F5
ScaleHLS	MLIR (Affine or Higher)	Vivado HLS	✓	✗	✓	✓	✗
CIRCT HLS	MLIR (Affine or Higher)	CIRCT	✓	✗	✗	✗	✓
FROST	Halide, Tiramisu	Vivado HLS	✗	✗	✗	✓	✗
Hot & Spicy	Annotated Python	SDSoC	✗	✗	✗	✓	✗
HeteroCL	Annotated C, OpenCL	Intel or Vitis HLS	✗	✗	✗	✓	✗
HPVM2FPGA	HeteroC++	Intel HLS	✓	✓	✓	✓	✗
Phism	C/C++	Vitis HLS	✓	✗	✗	✓	✗

F1

Optimizations Without Manual Annotations

F2

Automatic Partitioning of Host and Kernel Code

F3

Design Space Exploration of Optimization Strategies

F4

FPGA Support

F5

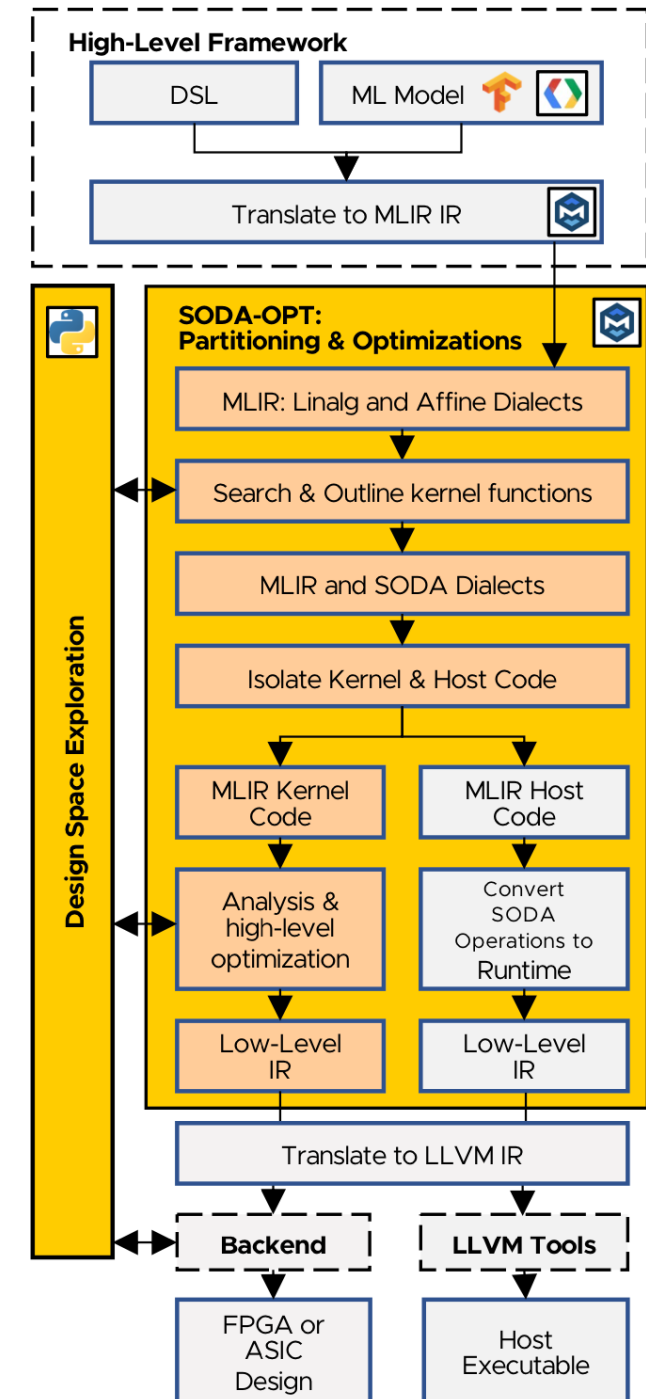
ASIC Support

SODA-OPT contributions / scope

- **MLIR-based inputs**
 - Supports all high-level application once it can be converted into `linalg`, `affine` dialects
- Configurable pipeline of **high-level optimization** passes
- **Automated flow** from high-level outlined kernels into custom hardware impl. through HLS backends
- **DSE** of compiler options: pipeline param. tuning



+



The Multi-Level Intermediate Representation Compiler Infrastructure

- **Open-source**
- **Progressive lowering** between existing and new operations
- **Reuse** of abstractions and compiler transformations
- Enables **co-existence** of different abstractions

```
%results:2 = "d.operation"(%arg0, %arg1) ({  
  // Regions belong to Ops and can have multiple blocks. Region  
  ^block(%argument: !d.type): Block  
    // Ops have function types (expressing mapping).  
    %value = "nested.operation"() ({  
      // Ops can contain nested regions. Region  
      "d.op"() : () -> ()  
    }) : () -> (!d.other_type)  
    "consume.value"(%value) : (!d.other_type) -> ()  
  ^other_block: Block  
    "d.terminator"() [^block(%argument : !d.type)] : () -> ()  
  })  
  // Ops can have a list of attributes.  
  {attribute="value" : !d.type} : () -> (!d.type, !d.other_type)
```

Figure from: Lattner, Chris, et al. "MLIR: Scaling compiler infrastructure for domain specific computation." *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. IEEE, 2021.

MLIR Dot Product Example and Lowering

Linalg abstraction

```
1 func.func @dot(%A: memref<100xf32>, %B: memref<100xf32>,  
2               %out: memref<f32>) {  
3   linalg.dot ins(%A, %B: memref<100xf32>, memref<100xf32>)  
4             outs(%out: memref<f32>)  
5   return  
6 }
```

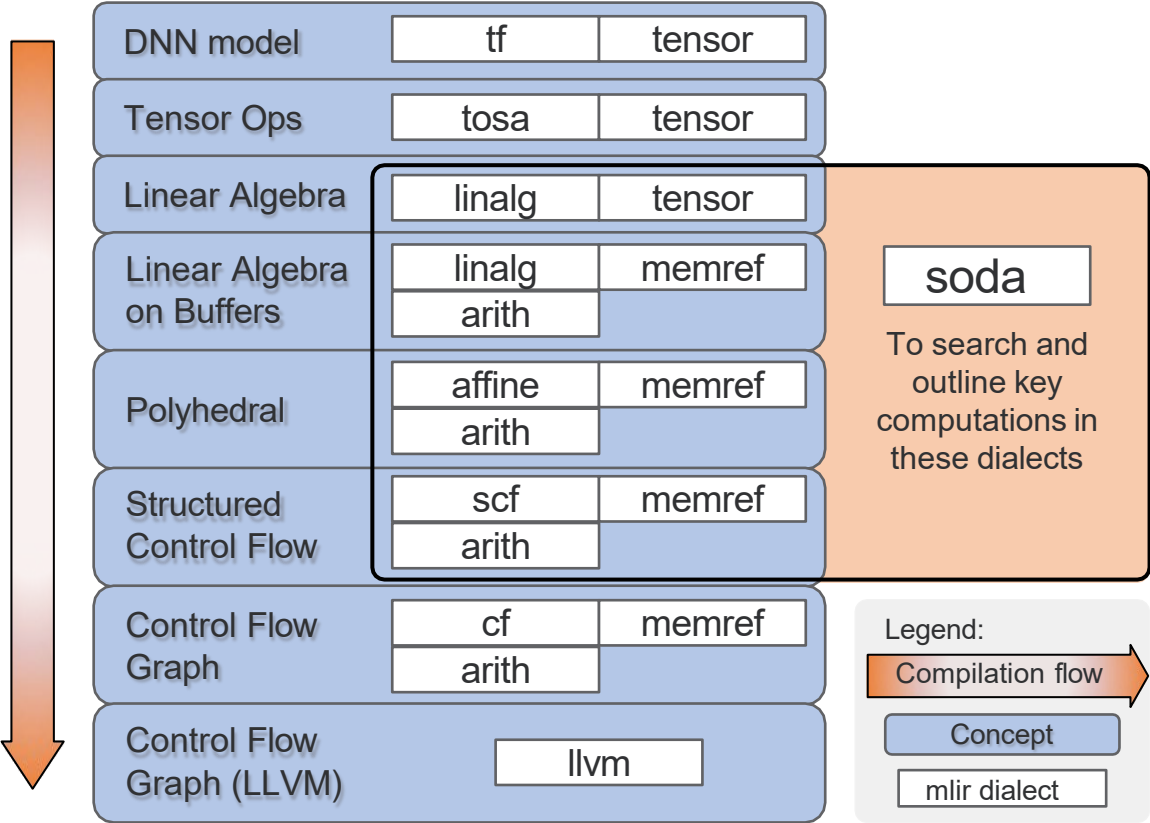
SCF abstraction

```
1 func.func @dot(%A: memref<100xf32>, %B: memref<100xf32>,  
2               %out: memref<f32>) {  
3   %c0 = arith.constant 0 : index  
4   %c100 = arith.constant 100 : index  
5   %c1 = arith.constant 1 : index  
6   scf.for %arg3 = %c0 to %c100 step %c1 {  
7     %0 = memref.load %A[%arg3] : memref<100xf32>  
8     %1 = memref.load %B[%arg3] : memref<100xf32>  
9     %2 = memref.load %out[] : memref<f32>  
10    %3 = arith.mulf %0, %1 : f32  
11    %4 = arith.addf %2, %3 : f32  
12    memref.store %4, %out[] : memref<f32>  
13  }  
14  return  
15 }
```

CF abstraction

```
1 func.func @dot(%A: memref<100xf32>, %B: memref<100xf32>,  
2               %out: memref<f32>) {  
3   %c0 = arith.constant 0 : index  
4   %c100 = arith.constant 100 : index  
5   %c1 = arith.constant 1 : index  
6   cf.br ^bb1(%c0 : index) // 2 preds: ^bb0, ^bb2  
7 ^bb1(%0: index): // 2 preds: ^bb0, ^bb2  
8   %1 = arith.cmpi slt, %0, %c100 : index  
9   cf.cond_br %1, ^bb2, ^bb3  
10 ^bb2: // pred: ^bb1  
11   %2 = memref.load %A[%0] : memref<100xf32>  
12   %3 = memref.load %B[%0] : memref<100xf32>  
13   %4 = memref.load %out[] : memref<f32>  
14   %5 = arith.mulf %2, %3 : f32  
15   %6 = arith.addf %4, %5 : f32  
16   memref.store %6, %out[] : memref<f32>  
17   %7 = arith.addi %0, %c1 : index  
18   cf.br ^bb1(%7 : index)  
19 ^bb3: // pred: ^bb1  
20   return  
21 }
```


The SODA Dialect



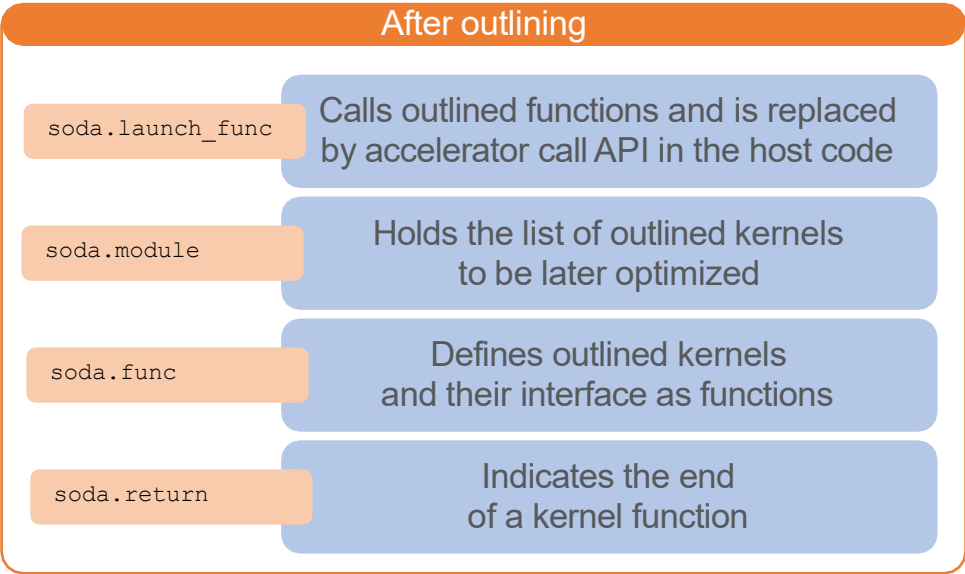
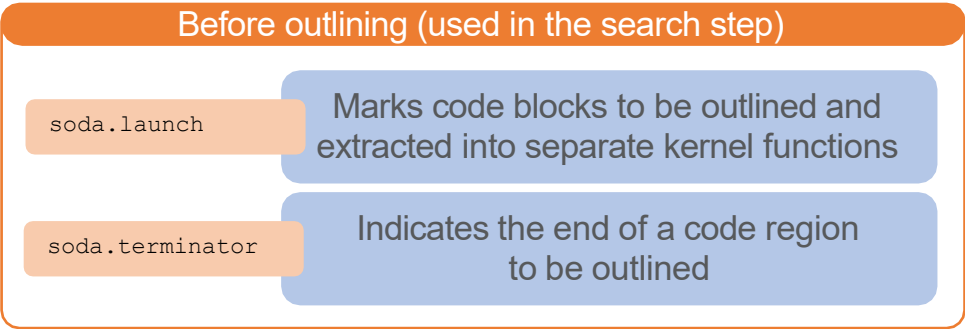
Before outlining (used in the search step)

<code>soda.launch</code>	Marks code blocks to be outlined and extracted into separate kernel functions
<code>soda.terminator</code>	Indicates the end of a code region to be outlined

After outlining

<code>soda.launch_func</code>	Calls outlined functions and is replaced by accelerator call API in the host code
<code>soda.module</code>	Holds the list of outlined kernels to be later optimized
<code>soda.func</code>	Defines outlined kernels and their interface as functions
<code>soda.return</code>	Indicates the end of a kernel function

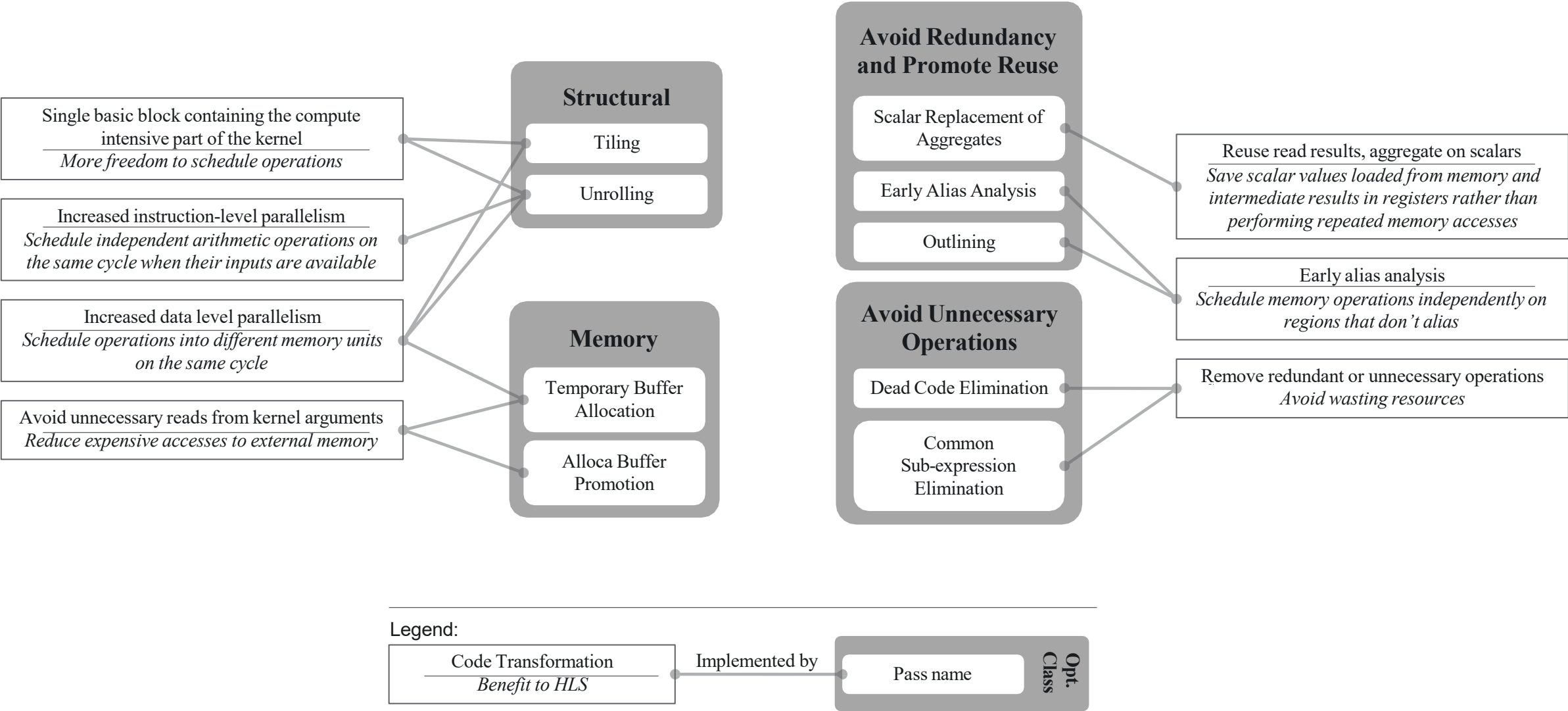
Search and Outlining



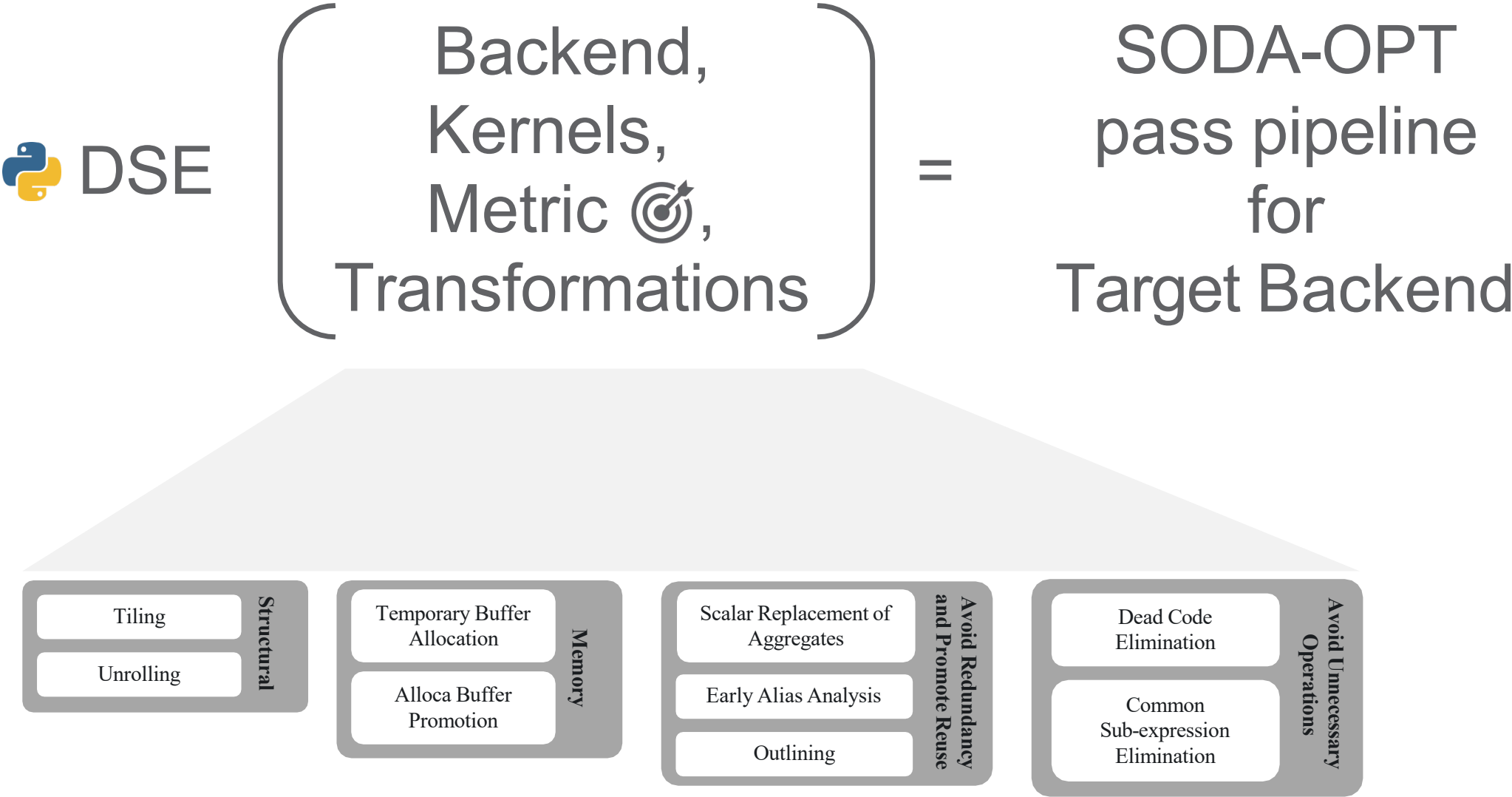
```
1 module {
2   func.func @main(%A: memref<42x42xf32>, %B: memref<42x42xf32>,
3                 %C: memref<42x42xf32>) {
4     soda.launch {
5       linalg.matmul ins(%A, %B : memref<42x42xf32>, memref<42x42xf32>)
6         outs(%C : memref<42x42xf32>)
7       soda.terminator
8     }
9     return
10  }
```

```
1 module attributes {soda.container_module} {
2   func.func @my_matmul(%A: memref<42x42xf32>, %B: memref<42x42xf32>,
3                     %C: memref<42x42xf32>) {
4     soda.launch_func @accelerators::@my_matmul_kernel args(
5       %A : memref<42x42xf32>, %B : memref<42x42xf32>, %C : memref<42x42xf32>)
6     return
7   }
8
9   soda.module @accelerators {
10     soda.func @my_matmul_kernel(%A_kernel: memref<42x42xf32>,
11                               %B_kernel: memref<42x42xf32>,
12                               %C_kernel: memref<42x42xf32>) kernel{
13       linalg.matmul ins(%A_kernel, %B_kernel : memref<42x42xf32>, memref<42x42xf32>)
14         outs(%C_kernel : memref<42x42xf32>)
15       soda.return
16     }
17   }
```

Optimizations for High-Level Synthesis

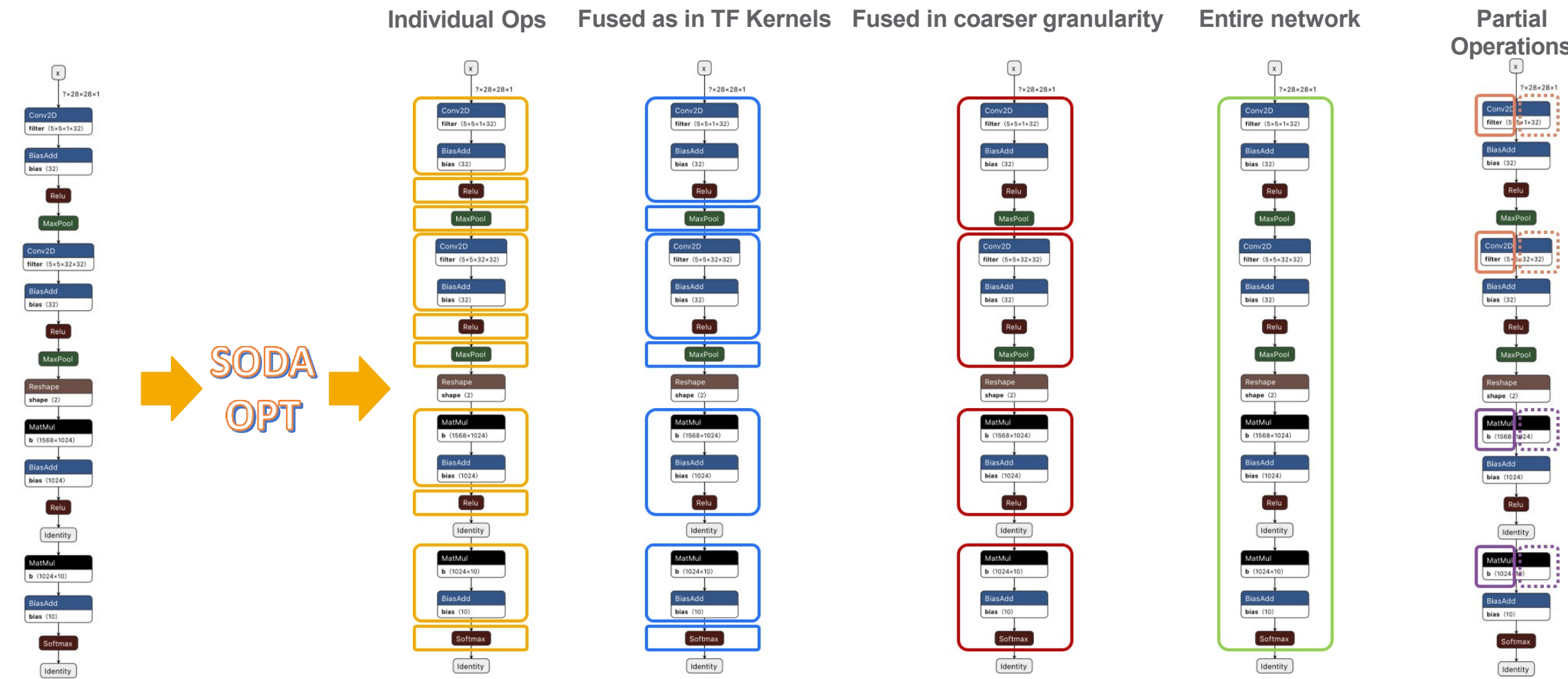


DSE Engine



Outlining in Different Granularities

- The impact of outlining and generating accelerators for different granularities of a DNN model



Outlining in Different Granularities

Execution times of LeNet with different outlining strategies.

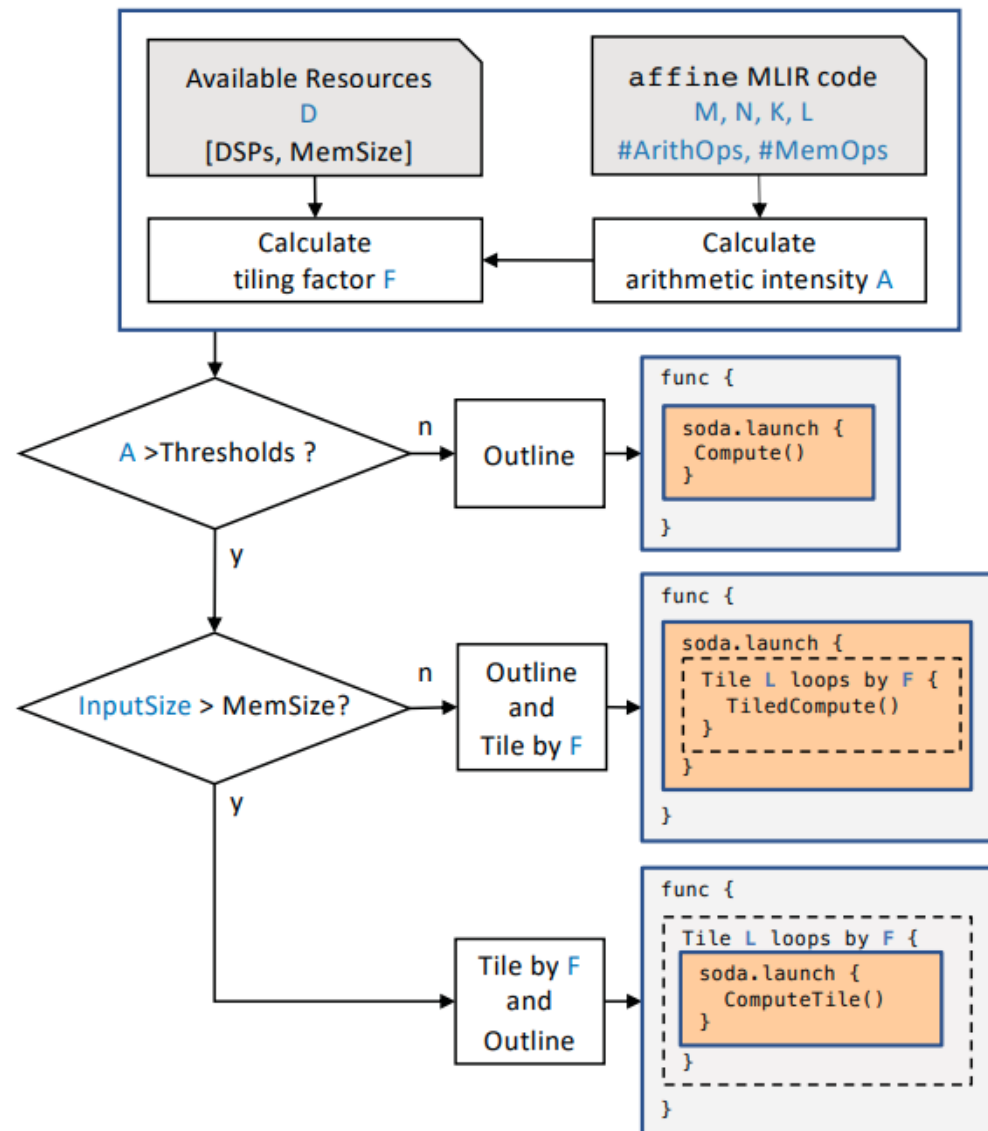
Layer type	Layer params	Individual Ops		Fused as in TF kernels		Fused in coarser granularity		Entire network	
		Baseline	Unrolled	Baseline	Unrolled	Baseline	Unrolled	Baseline	Unrolled
Conv2D	5x5,6,same	2,423,374	2,353,598	2,462,602	2,388,122	2,526,225	2,443,073	6,806,682	5,965,844
Activation	Relu	39,230	34,526						
AveragePooling2D	2x2, s2x2	88,244	84,338						
Conv2D	5x5,16,same	4,917,022	4,835,522	4,917,022	4,835,522	5,175,970	4,853,938		
Activation	Relu	12,912	11,312						
AveragePooling2D	2x2, s2x2	30,092	28,842						
Dense	120 units	1,010,522	926,402	1,011,602	927,362	1,011,602	927,362		
Activation	Relu	962	842						
Dense	84 units	213,446	195,722	214,202	196,394	214,202	196,394		
Activation	Relu	674	590						
Dense	10 units	17,841	16,372	17,841	16,372	19,084	16,731		
Activation	Softmax	454	410						
Total		8,754,773	8,488,476	8,742,059	8,477,362	8,947,083	8,437,498	6,806,682	5,965,844

* FPGA Target: Xilinx xc7vx690t-ffg1930-3 @ 100MHz
Merged cells correspond to a single accelerator.

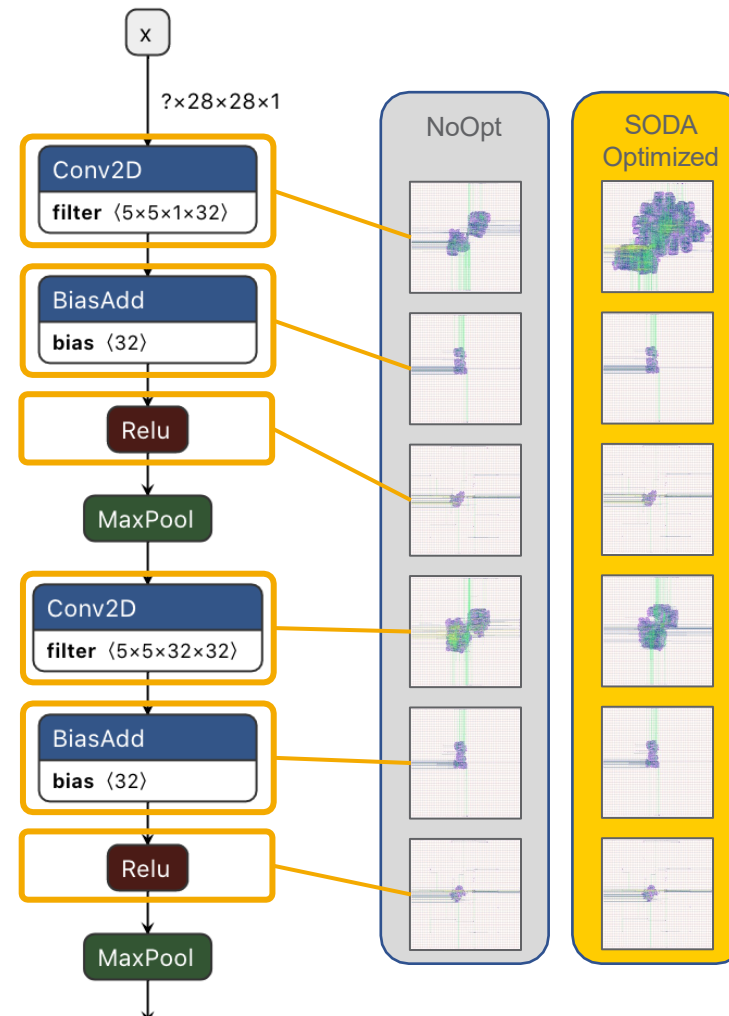
1.47x Speedup



Tiling, Outlining and Optimizing



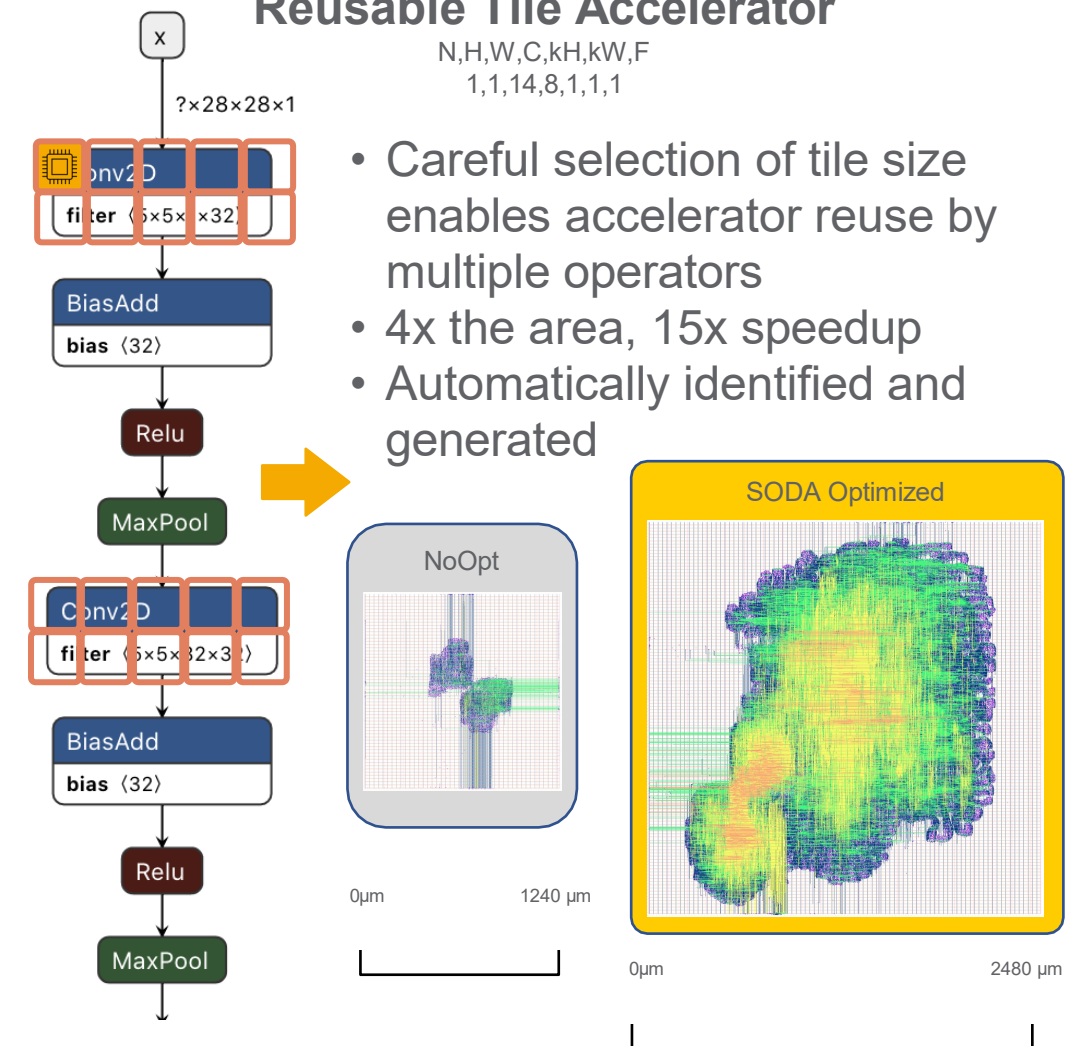
Individual Ops



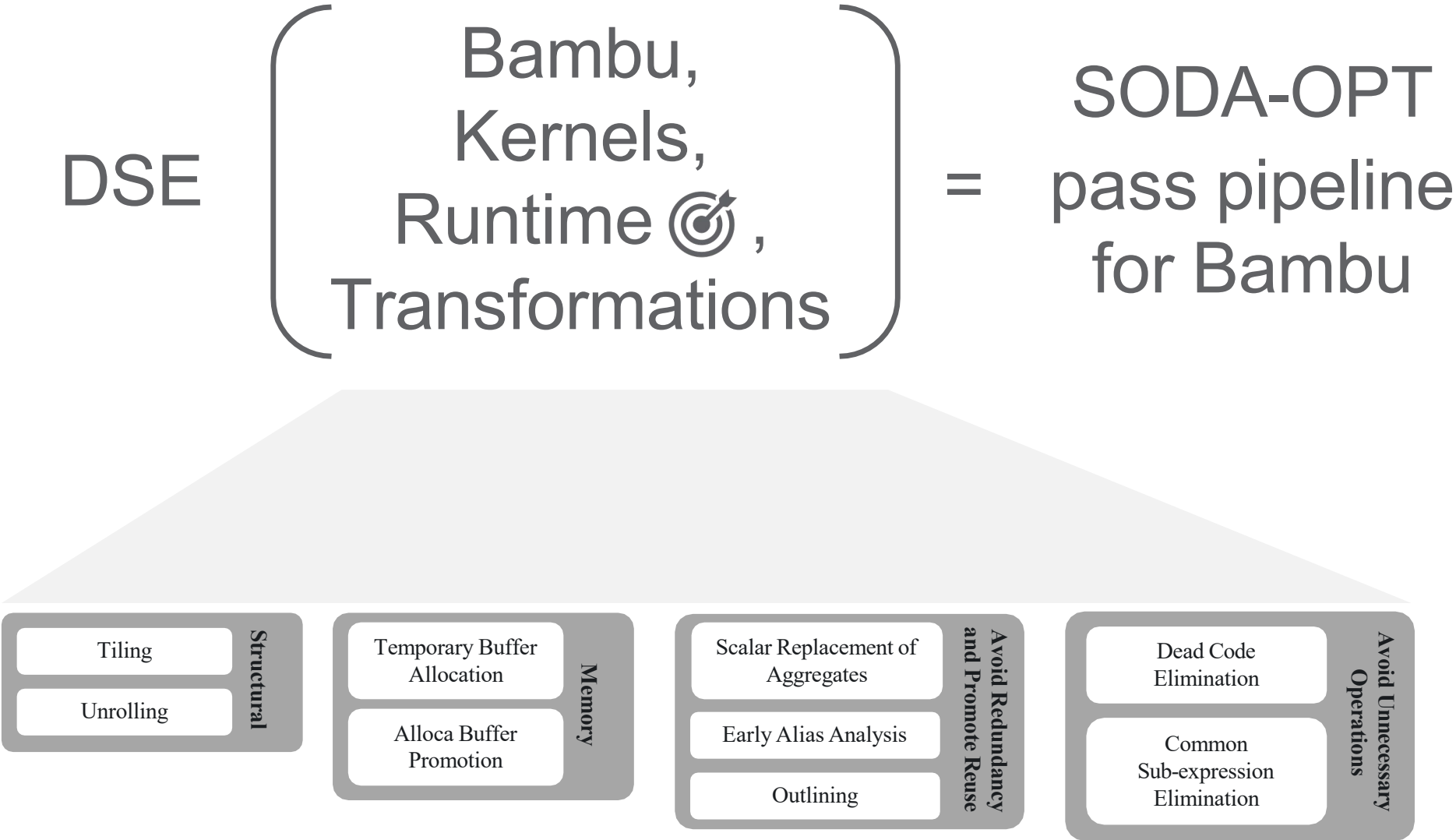
Reusable Tile Accelerator

N,H,W,C,kH,kW,F
1,1,14,8,1,1,1

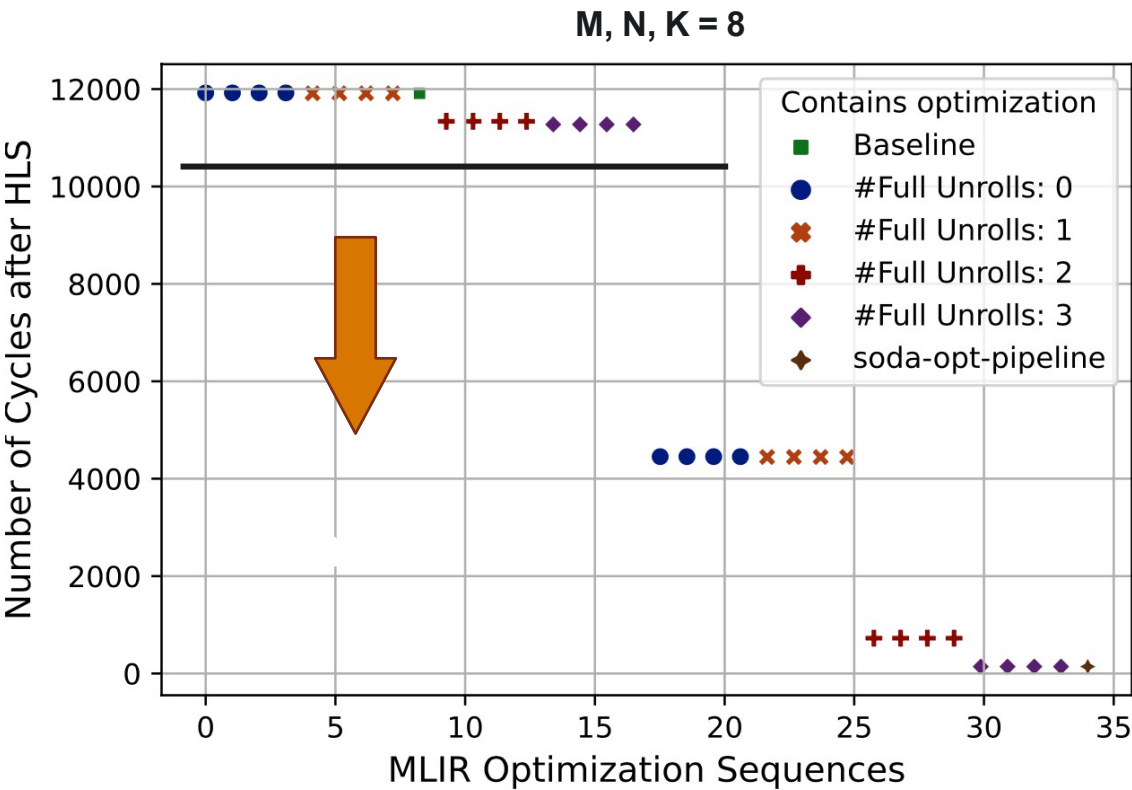
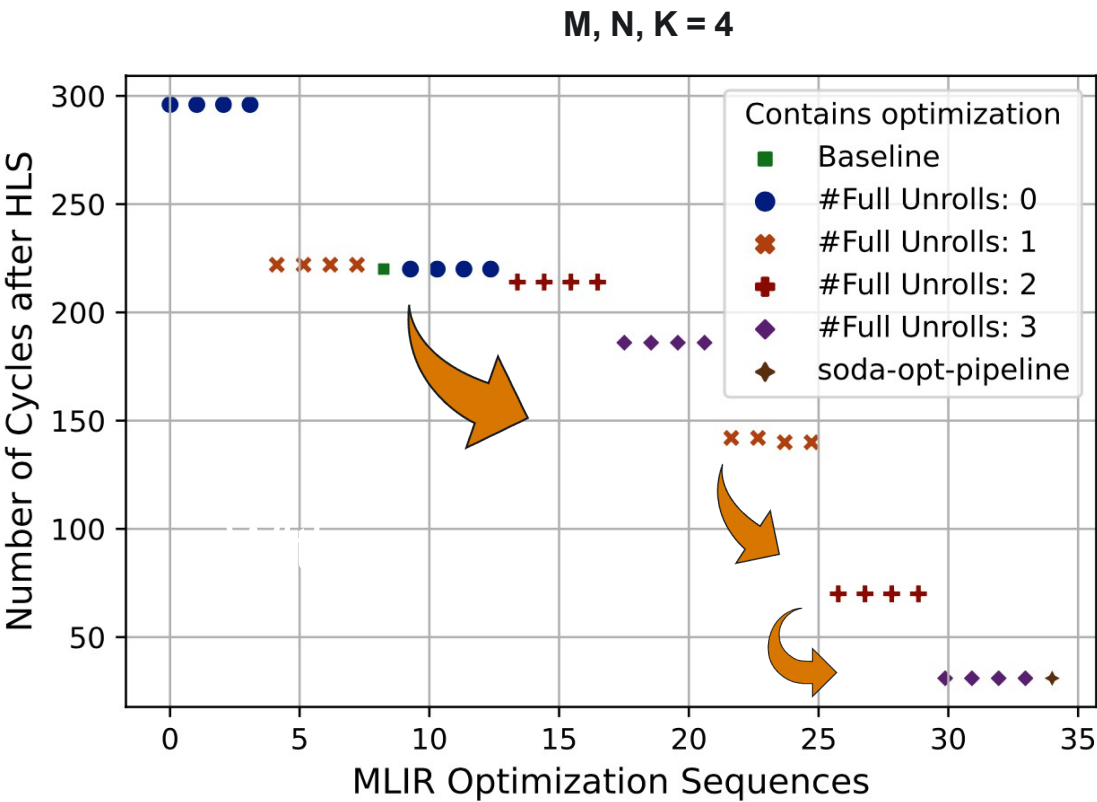
- Careful selection of tile size enables accelerator reuse by multiple operators
- 4x the area, 15x speedup
- Automatically identified and generated



Effects of DSE Engine



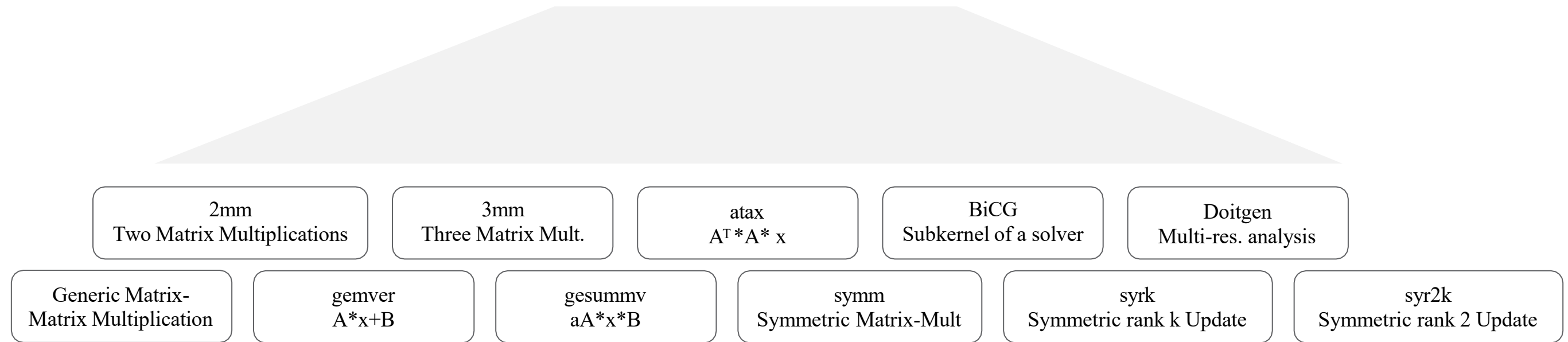
Effects of DSE Engine on GEMM Kernels



Structural Optimizations and Memory Optimizations

Benchmarking the SODA-OPT Pass Pipeline

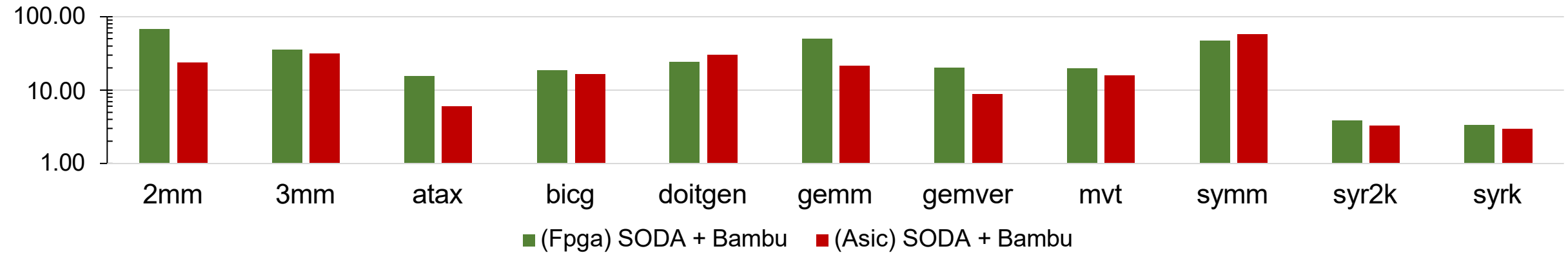
SODA-OPT pass pipeline $\left(\begin{array}{c} \text{For } \langle \text{HLS Tool} \rangle, \\ \text{Polybench Kernels} \end{array} \right) = \text{Runtime} \quad \text{🎯}$



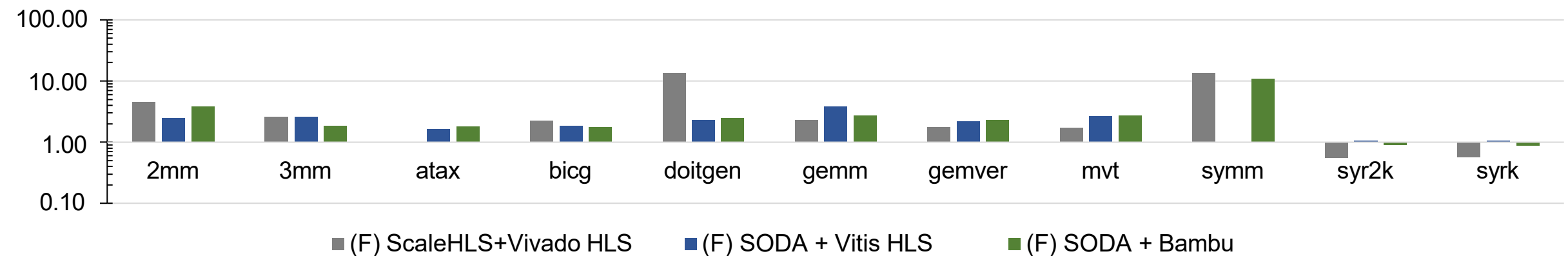
Kernels from: Louis-Noel Pouchet and Tomofumi Yuki. [n.d.]. PolyBench/C 4.2.1. <http://polybench.sourceforge.net>.

Performance Improvement vs Others

Average of Speedups
Over **Bambu**
(F:FPGA or A:ASIC)



Average of Speedups
Over **Vitis HLS**
(F:FPGA)

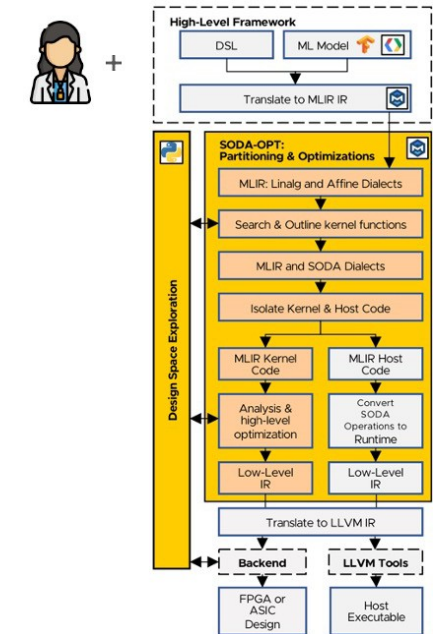


SODA-OPT derived accelerators **outperform**:

- **100%** of Bambu or Vitis HLS accelerators generated with no manual annotations
- **70%** of accelerators generated with current state of the art ScaleHLS

Contributions

- An **MLIR** based compiler flow for high-level synthesis
- An **MLIR dialect to search and outline** MLIR code at suitable abstractions
 - Kernels from HL applications can be outlined **at different granularities**
- Compiler **passes and pipelines** that enhance HLS of arbitrary regions of an application for **any target** (HLS Backend and Platform)
 - **Up to 60x speedup** over baseline designs that only leverage HLS optimizations
- Puts the domain scientist **in control** of custom accelerator generation with compiler passes. Transforming the scientist into a **Lead User** and potential **source of novel concepts**
- New HLS capabilities, essential to enable an “**agile hardware design**” approach



F1

Optimizations Without
Manual Annotations

F2

Automatic Partitioning of
Host and Kernel Code

F3

Design Space Exploration
of Optimization Strategies

F4

FPGA
Support

F5

ASIC
Support

Reference



<https://github.com/pnnl/soda-opt>



<https://hpc.pnl.gov/SODA/tutorials/2024/ISCA24.html>