



高层综合与AIE阵列映射的MLIR应用

罗国杰

gluo@pku.edu.cn

北京大学

Why Compiler Infrastructures? “ $p \times n \times q$ ” → “ $p + n + q$ ”



image source: lei.chat

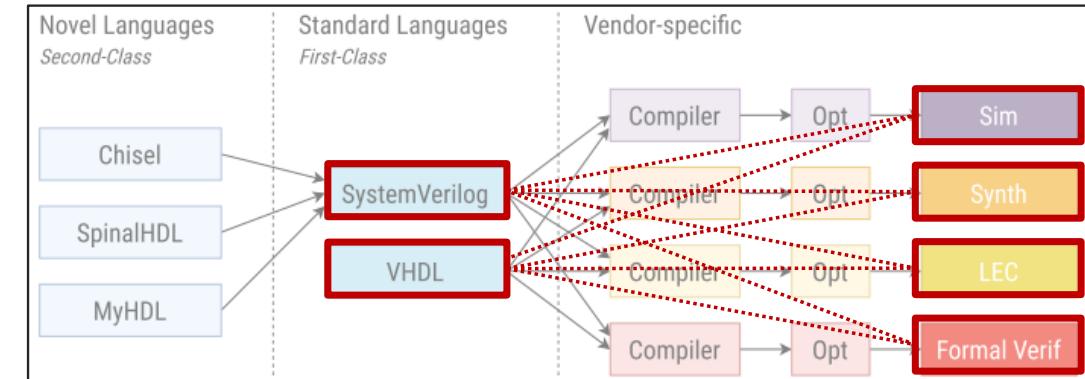
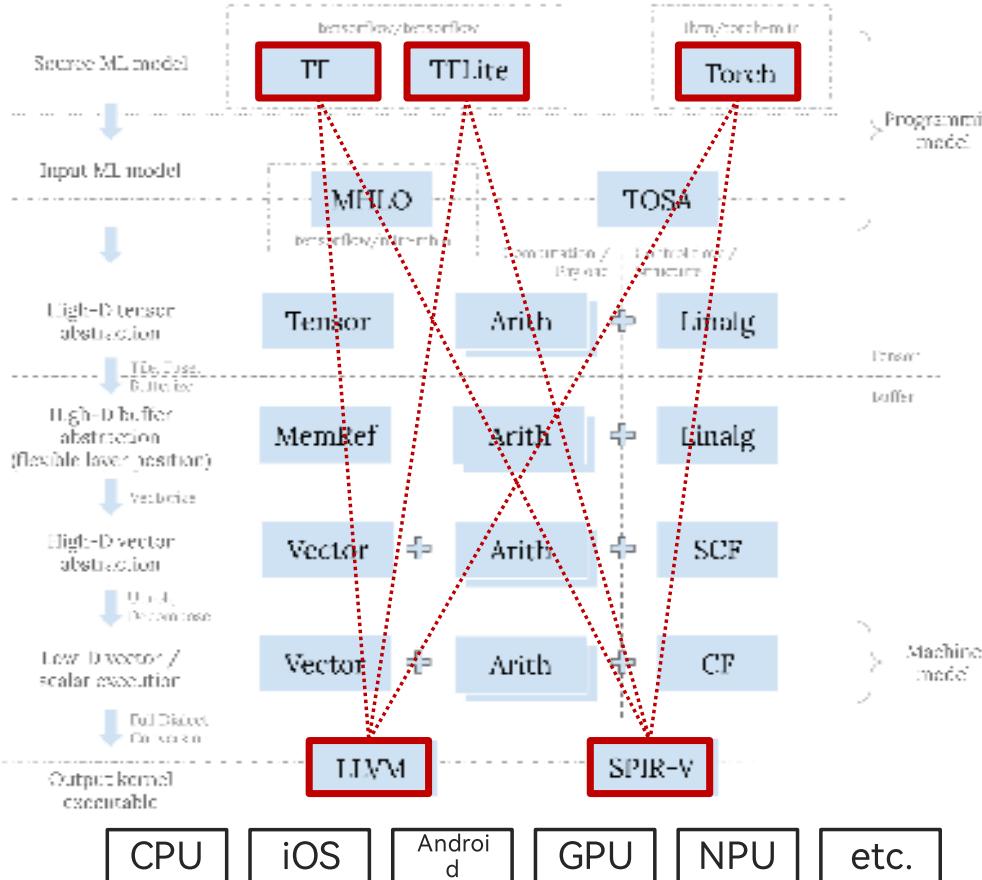


image source: [LLHD@PLDI'20]

Refactor: $x_1y_1z_1 + x_1y_1z_2 + x_1y_2z_1 + x_1y_2z_2 + x_1y_3z_1 + x_1y_3z_2 + x_2y_1z_1 + x_2y_1z_2 + x_2y_2z_1 + x_2y_2z_2 + x_2y_3z_1 + x_2y_3z_2 = (x_1 + x_2) \cdot (y_1 + y_2 + y_3) \cdot (z_1 + z_2)$

How MLIR and CIRCT Fit In

Fixed processor

Custom processor

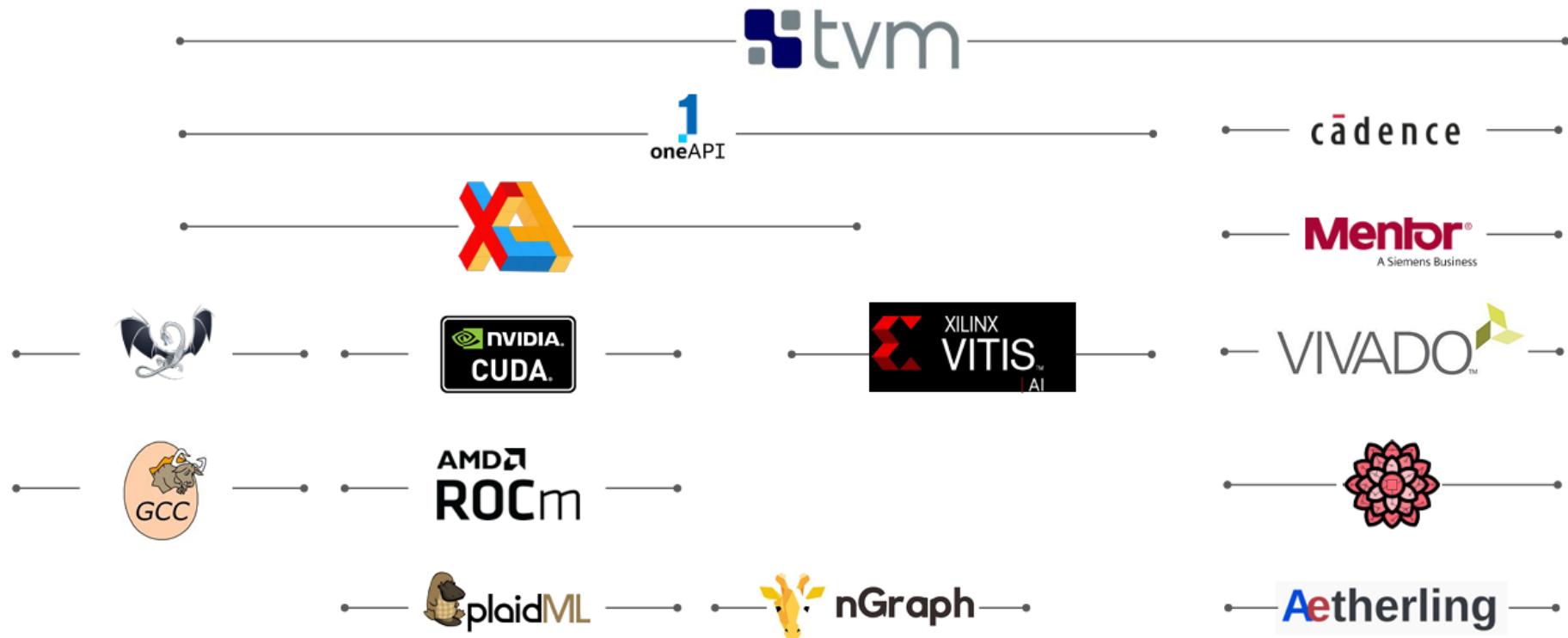
CPU, etc.

GPU, etc.

TPU, NPU, etc.

FPGA, CPLD, etc.

ASIC



MLIR and CIRCT Ecosystem

↔ Exists today
↔ Hypothetical

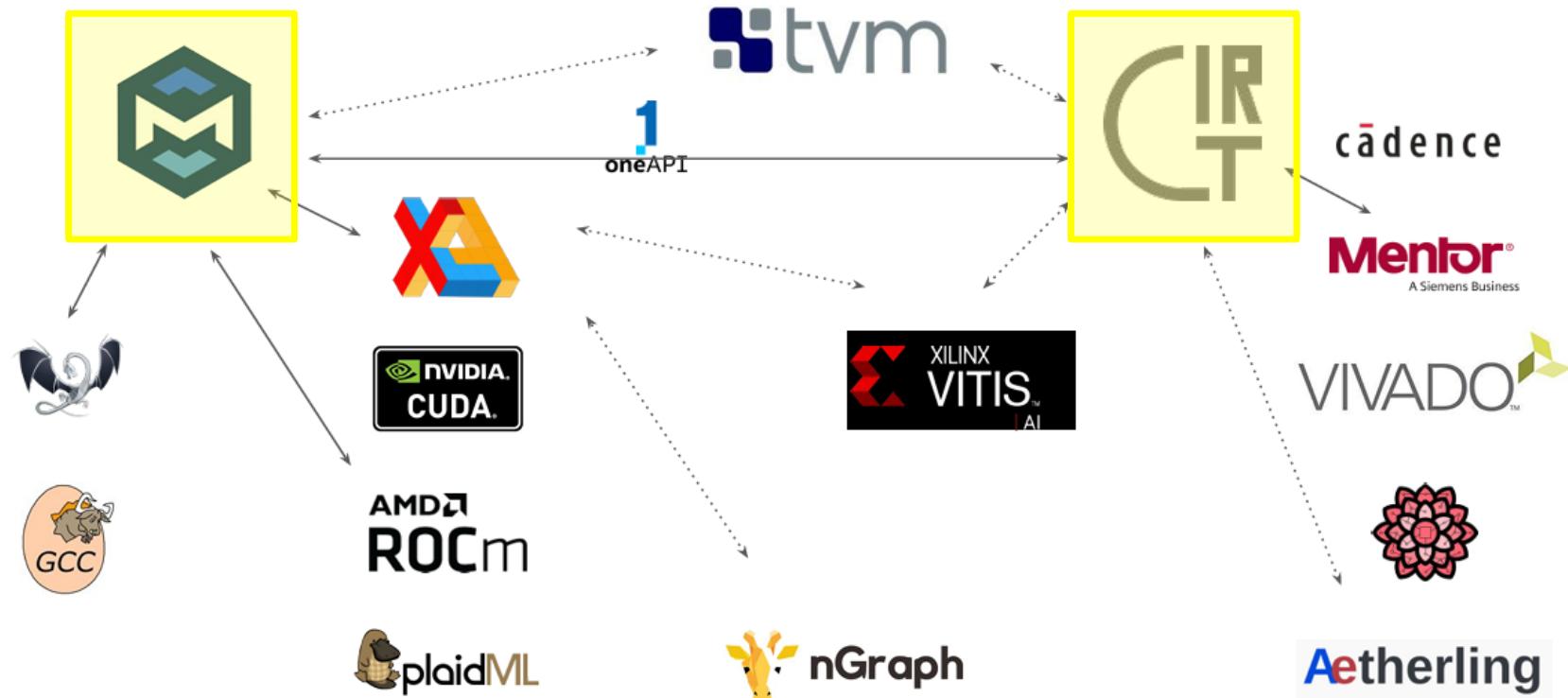
CPU, etc.

GPU, etc.

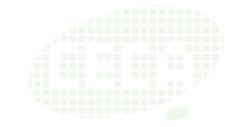
TPU, NPU, etc.

FPGA, CPLD, etc.

ASIC



What is MLIR for?



A hybrid IR to support multiple different requirements in a unified infrastructure:

- ▶ The ability to represent dataflow graphs (such as in TensorFlow)
- ▶ Optimizations and transformations typically done on such graphs (e.g. in Grappler).
- ▶ Ability to host high-performance loop optimizations and to transform memory layouts of data.
- ▶ Code generation “lowering” transformations (e.g., DMA insertion, memory tiling, vectorization).
- ▶ Ability to represent target-specific operations, e.g. accelerator-specific high-level operations.
- ▶ Quantization and other graph transformations done on a Deep-Learning graph.
- ▶ Polyhedral primitives.
- ▶ **Hardware Synthesis Tools / HLS.**

Non-goals

- ▶ not for low level machine code generation algorithms (like reg. allocation and instr. scheduling).
- ▶ not a source language to write kernels (analogous to CUDA C++).

* <https://mlir.llvm.org/>

* C. Lattner et al., “MLIR: Scaling Compiler Infrastructure for Domain Specific Computation,” CGO 2021.

Users of MLIR



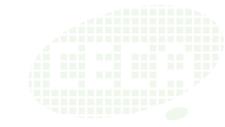
CIRCT: Circuit IR Compilers and Tools

- ▶ <https://circt.llvm.org/>
- ▶ The CIRCT project is an (experimental!) effort looking to apply MLIR and the LLVM development methodology to the domain of hardware design tools.

MLIR-AIE: Toolchain for AMD AI Engine devices

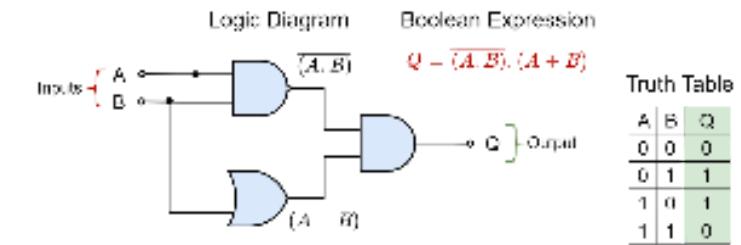
- ▶ <https://xilinx.github.io/mlir-aie/>
- ▶ MLIR-AIE is a toolchain providing low-level device configuration for Versal AIEngine-based devices. Support is provided to target the AIEngine portion of the device, including processors, stream switches, TileDMA and ShimDMA blocks. Backend code generation is included, targeting the LibXAIE library, along with some higher-level abstractions enabling higher-level design.

Background: Functionality of Digital Circuits

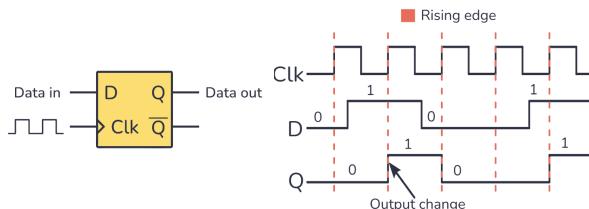


- ▶ Combinational (stateless) circuits
 - ▶ with input x and output y

$$y = f(x)$$

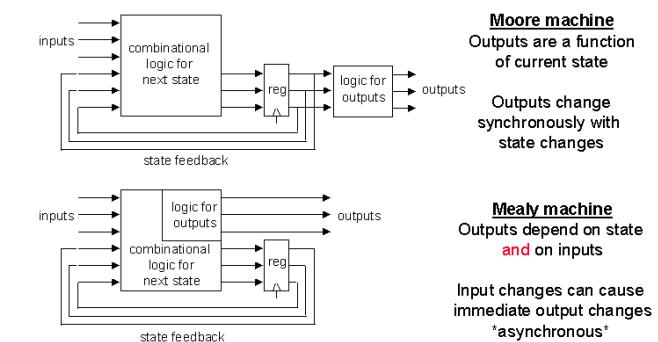


- ▶ Sequential (stateful) circuits
 - ▶ with input x , state S , and output y



$$\begin{cases} y = f(S) \\ S' = g(S, x) \end{cases} \quad (\text{Moore machine})$$

$$\begin{cases} y = f(S, x) \\ S' = g(S, x) \end{cases} \quad (\text{Mealy machine})$$



- ▶ Formally, $f(\cdot)$ and $g(\cdot)$ are 4-valued $\{0,1,X,Z\}$ logic, instead of Boolean, functions

What is missing? Timing, Power, Languages (SystemVerilog/VHDL), etc.

Background: High-Level Synthesis (HLS) and HLS C



HLS C

- ▶ A hardware synthesis language with automatic scheduling, array partitioning, interface generation, etc., in C syntax

HLS vs. compilation

- ▶ Compilation: **C** \Rightarrow instruction stream
 - To compute, a CISC/RISC processor executes the **instructions**
- ▶ HLS: **C** \Rightarrow hardware description language (HDL)
 - To compute, the HDL-described **circuit** executes Boolean switching

Background: HLS C (in AMD/Xilinx Vitis as example)



HLS C is a synthesizable and functional hardware description language
(functional: 功能性 not 函数式)

► Part I: C as the host language

- functions ↔ modules; arguments ↔ I/O ports; control flows ↔ control logic;
- operators ↔ functional units; scalar variables ↔ wires or regs; array variables ↔ on-chip memories;

► Part II: HLS pragmas <https://docs.amd.com/r/en-US/ug1399-vitis-hls/HLS-Pragmas>

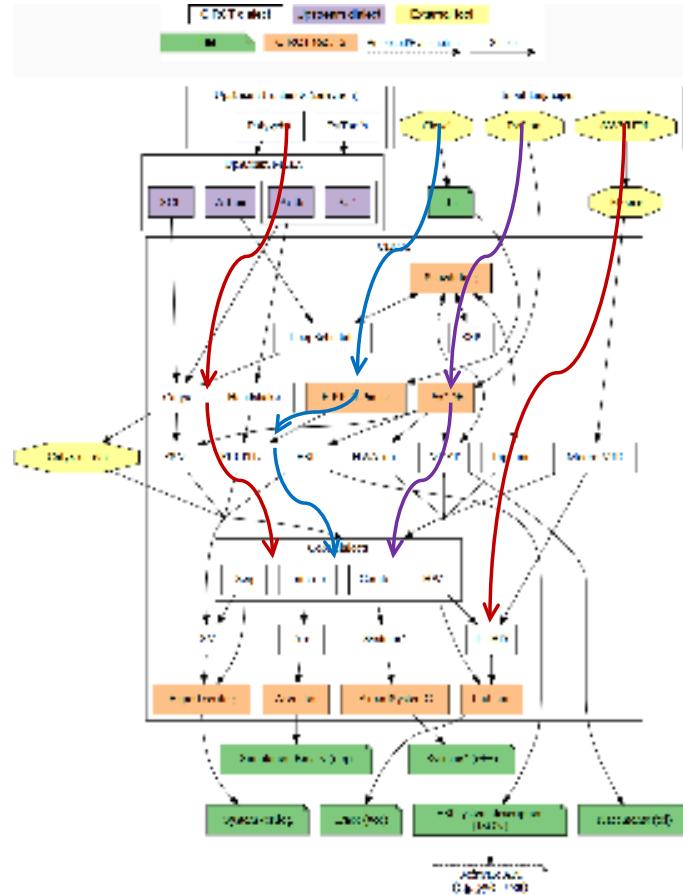
- interface synthesis: interface, stream
- schedule-related: pipeline, dataflow
- resource-related: unroll, array_partition, allocation

► Part III: HLS library

- HLS Arbitrary Precision Types Library https://github.com/Xilinx/HLS_arbitrary_Precision_Types/
- HLS Stream Library <https://github.com/Xilinx/hls-lib-stream/>

(The two repositories above is for simulation only. The synthesis is likely hardcoded in Vitis HLS tools.)

Circuit IR Compilers and Tools (CIRCT)



High-level IR

- ▶ Calyx
- ▶ Scheduling
- ▶ SSP
- ▶ LoopSchedule
- ▶ Pipeline
- ▶ Handshake
- ▶ ESI

RTL-synth IR

- ▶ FSM
- ▶ FIRRTL
- ▶ HWAarith
- ▶ Core dialects
- ▶ SV
- ▶ MSFTs

RTL-sim IR

- ▶ Arc
- ▶ SystemC
- ▶ LLHD

High-Level Synthesis (HLS) in CIRCT



doHLSFlowDynamic()

- ▶ Software lowering

- lower affine
- Convert SCF to CF
- DHLS pipeline
- handshake tranforms pipeline

- ▶ HW path

- simple canonicalizer
- handshake to HW
- simple canonicalizer
- ESI lowering pipeline
- HW lowering pipeline
- export Verilog

doHLSFlowCalyx()

- ▶ Software transformations

- SCF to Calyx; simple canonicalizer
- remove comb. groups; simple canonicalizer
- Calyx to FSM; simple canonicalizer
- materialize Calyx to FSM; simple canonicalizer
- remove groups from FSM; simple canonicalizer

- ▶ HW path

- Calyx to HW; simple canonicalizer
- convert FSM to SV
- HW lowering pipeline
- export Verilog

HLS in CIRCT



```
cgeist -S --function=* --memref-fullrank --scal-rep ${design}.c | \
mlir-opt --scf-for-to-while | \
circt-opt --lower-scf-to-calyx --canonicalize | \
circt-translate --export-calyx | \
calyx -b verilog
```

[lower-scf-to-fsm] error: Operation 'calyx.par' not supported for
FSM compilation #5154

<https://github.com/llvm/circt/issues/5154>

[CalyxToFSM] Add par FSM lowering #4606
<https://github.com/llvm/circt/issues/4606>

doHLSFlowCalyx()

► Software transformations

- SCF to Calyx; simple canonicalizer
- remove comb. groups; simple canonicalizer
- **Calyx to FSM**; simple canonicalizer
- materialize Calyx to FSM; simple canonicalizer
- remove groups from FSM; simple canonicalizer

► HW path

- Calyx to HW; simple canonicalizer
- convert FSM to SV
- HW lowering pipeline
- export Verilog

HLS in CIRCT



doHLSFlowDynamic()

- ▶ Software lowering
 - lower affine
 - Convert SCF to CF
 - **DHLS pipeline**
 - **handshake tranforms pipeline**
- ▶ HW path
 - simple canonicalizer
 - handshake to HW
 - simple canonicalizer
 - **ESI lowering pipeline**
 - **HW lowering pipeline**
 - export Verilog

```
cgeist -S --function=* --memref-fullrank -scal-rep ${design}.c | \
mlir-opt --lower-affine | \
mlir-opt --convert-scf-to-cf | \
circt-opt --flatten-memref | \
circt-opt --handshake-legalize-memrefs | \
mlir-opt --convert-scf-to-cf --canonicalize | \
circt-opt --lower-cf-to-handshake | \
circt-opt --handshake-lower-extmem-to-hw --canonicalize | \
circt-opt --handshake-materialize-forks-sinks --canonicalize | \
circt-opt --handshake-insert-buffers --canonicalize | \
circt-opt --lower-handshake-to-hw --canonicalize | \
circt-opt --lower-esi-ports | \
circt-opt --lower-esi-to-physical | \
circt-opt --lower-esi-to-hw --canonicalize | \
circt-opt --lower-seq-hlmem | \
circt-opt --hw-memory-sim | \
circt-opt --lower-seq-to-sv | \
circt-opt --hw-cleanup | \
circt-opt --hw-legalize-modules --canonicalize | \
circt-opt --prettyfy-verilog | \
circt-opt --export-verilog
```

MLIR: Affine Dialect

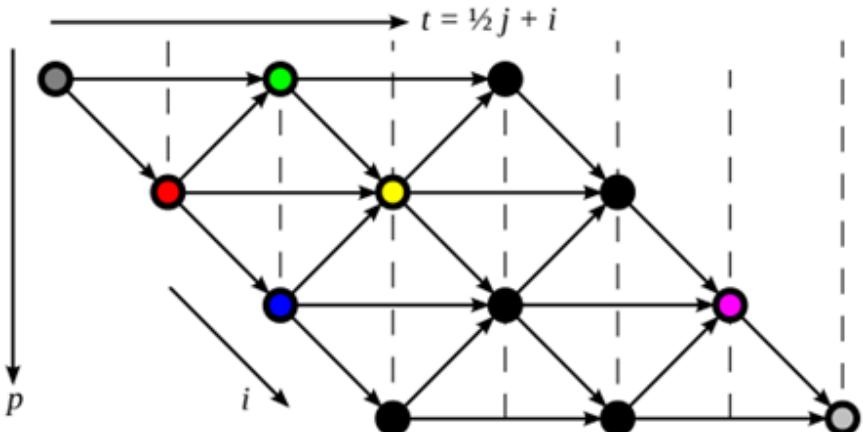
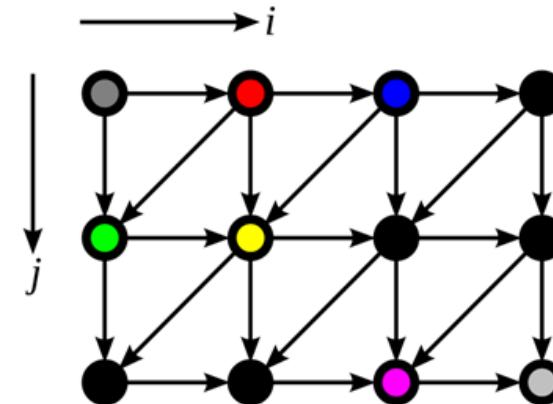


Polyhedral compilation toolbox

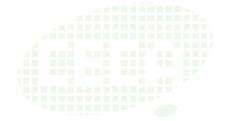
Affine expressions, maps, and integer sets

Affine loops and transformations

```
func @matmul(%A: memref<5x10xf32>, %B: memref<10x10xf32>,
             %C: memref<10x10xf32>) -> memref<10x10xf32> {
    affine.for %arg3 = 0 to 5 {
        affine.for %arg4 = 0 to 10 {
            affine.for %arg5 = 0 to 10 { ... }
        }
    }
    return %0 : memref<10x10xf32>
}
```



MLIR: SCF and CF Dialects



Structured Control Flow (SCF) :

- ▶ Contains operations that represent control flow constructs such as if and for
- ▶ The control flow has a structure unlike, for example, gotos or asserts

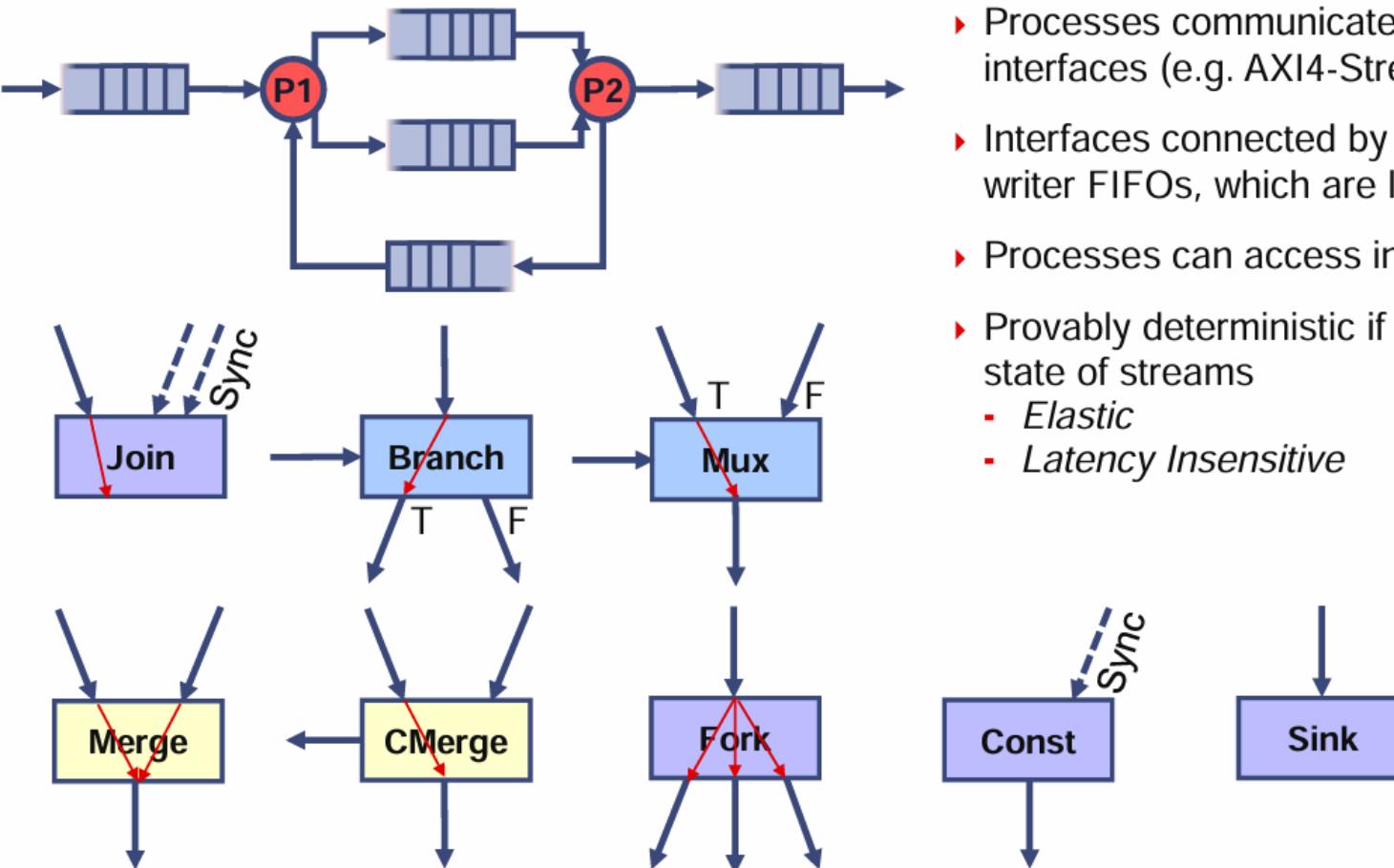
```
scf.for %iv = %lb to %ub step %step {  
    ... // body  
}
```

Control Flow (CF):

- ▶ Contains low-level, i.e. non-region based, control flow constructs
- ▶ Represents control flow directly on SSA blocks

```
^bb2:  
    %2 = call @someFn()  
    cf.br ^bb3(%2 : tensor<*xf32>)  
^bb3(%3: tensor<*xf32>):
```

CIRCT: Handshake Dialect



- Processes communicate through stream interfaces (e.g. AXI4-Stream)
- Interfaces connected by single-reader single-writer FIFOs, which are logically unbounded
- Processes can access interfaces in any order
- Provably deterministic if processes cannot test state of streams
 - Elastic*
 - Latency Insensitive*

CIRCT: ESI Dialect

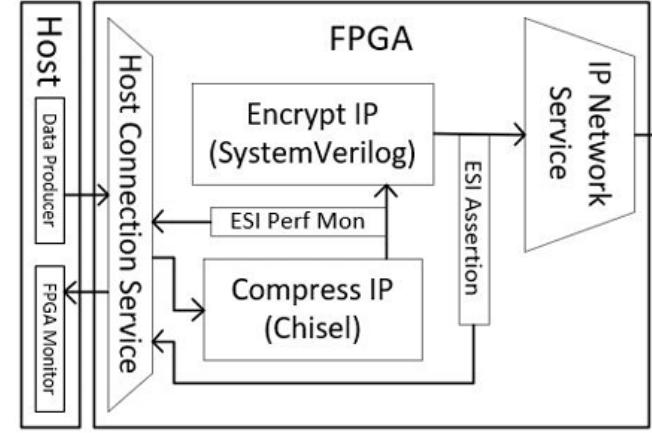


Type system for channels in hardware

Connect on- and off-chip components

Generates elastic connector circuits

Supports cosimulation



```

rtl.module @esi(%clock: i1, %reset: i1) {
    %esiChannel, %0 = rtl.instance "sender" @Sender (%clock) : (i1) -> (!esi.channel<i4>, i8)
    %bufferedChannel = esi.buffer %clock, %reset, %esiChannel { stages = 4 } : i4
    rtl.instance "receiver" @Reciever (%bufferedChannel, %clock) : (!esi.channel<i4>, i1) -> ()
}
  
```

CIRCT: Core Dialects



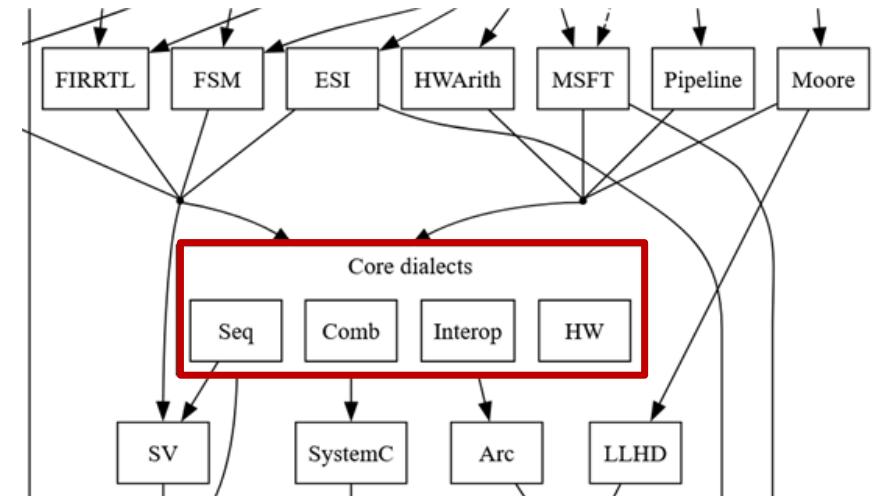
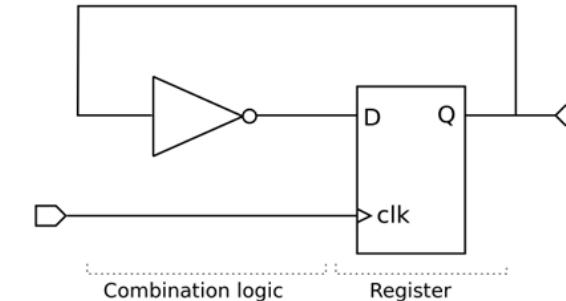
HW: structure (modules and ports)

Comb: Combinational logic (and, or, xor)

Seq: Sequential logic (registers)

Interop: between backends and tools

```
rtl.module @test1(%a: i12) -> (i32){
    %b = comb.add %a, %a : i12
    %c = comb.mul %a, %b : i12
    %d = comb.sextr %c : (i12) -> i32
    rtl.output %d : i32
}
```



CIRCT: SV Dialect



Represents System Verilog syntax

Designed for export, not transformation

Readability and other syntactical improvements happen here

Guides towards compatibility, but offers escape hatch

```
rtl.module @systemverilog(%clock: i1, %reset: i1) {
    %c0 = rtl.constant 0 : i32
    %c42 = rtl.constant 42 : i32
    %reg = sv.reg : !rtl.inout<i32>
    sv.alwaysff(posedge %clock) {
        sv.passign %reg, %c42 : i32
    } (asyncreset : posedge %reset) {
        sv.passign %reg, %c0 : i32
    }
}
```

CIRCT: Calyx Dialect

Components: the basic building block of Calyx programs, which has three sections: cells, wires, and control

Cells: instantiates sub-components like adders and registers

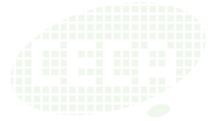
Assignments: The wires section, works like non-blocking assignments in RTL

Groups: a collection of assignments that are only active when run from the control program

Control: specifies the order to run groups

```
component main() -> () {  
    cells {  
        a = std_reg(32);  
        b = std_reg(32);  
    }  
    wires {  
        group write_a { ... }  
        group read_a { ... }  
        group write_b { ... }  
    }  
    control {  
        seq {  
            write_a;  
            par {  
                read_a;  
                write_b;  
            }  
        }  
    }  
}
```

CIRCT: FSM Dialect



A set of abstractions for finite-state machine

- ▶ Provide explicit and structural representations of states, transitions, and internal variables of an FSM, allowing convenient analysis and transformation.
- ▶ Provide a target-agnostic representation of FSM, allowing the state machine to be instantiated and attached to other dialects from different domains.
- ▶ allow to lower the FSM abstraction into HW+Comb+SV (Hardware) and Standard+SCF+MemRef (Software) dialects for the purposes of simulation, code generation, etc.

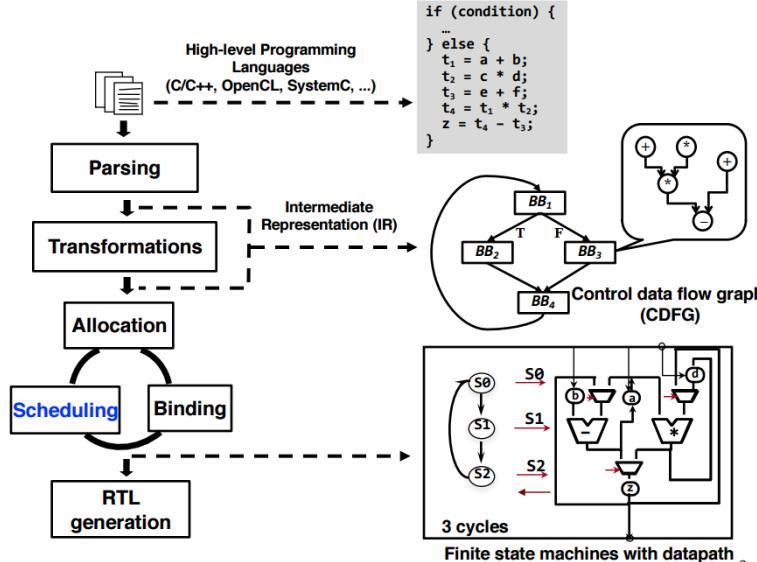
```
fsm.machine @foo(%arg0: i1) -> i1 attributes {stateType = i1} {  
    fsm.state "IDLE" output {  
        %true = constant true  
        fsm.output %true : i1  
    } transitions {  
        fsm.transition @BUSY ...  
    }  
  
    fsm.state "BUSY" output {  
        %false = constant false  
        fsm.output %false : i1  
    } transitions {  
        fsm.transition @BUSY ...  
        fsm.transition @IDLE ...  
    }  
}
```

Key Algorithms in HLS



1) Scheduling; 2) Binding; 3)

All three:



Code repository Q&A

- "Chat with Copilot" in [github.com](#)
- [deepwiki.com](#) by Devin

an scheduling example in CIRCT

```

36     /// This class provides a framework to model certain scheduling problems as
37     /// lexicographic linear programs (LP), which are then solved with an
38     /// extended version of the dual simplex algorithm.
39     ///
40     /// The approach is described in:
41     /// [1] B. D. de Dinechin, "Simplex Scheduling: More than Lifetime Sensitive
42     /// Instruction Scheduling", PRISM 1994.22, 1994.
43     /// [2] B. D. de Dinechin, "Fast Module Scheduling Under the Simplex Scheduling
44     /// Framework", PRISM 1995.23, 1995.
45
46     std::unique_ptr<mlir::Pass> createMapArithToCombPass();
47     std::unique_ptr<mlir::Pass> createFlattenMemRefPass();
48     std::unique_ptr<mlir::Pass> createFlattenMemRefCallsPass();
49     std::unique_ptr<mlir::Pass> createStripDebugInfoWithPredPass(
50         const std::function<bool(mlir::Location)> &pred);
51     std::unique_ptr<mlir::Pass> createMaximizeSSAPass();
52     std::unique_ptr<mlir::Pass> createInsertMergeBlocksPass();
53     std::unique_ptr<mlir::Pass> createPrintOpCountPass();
54     std::unique_ptr<mlir::Pass>
55     createMemoryBankingPass(ArrayRef<unsigned> bankingFactors = {},
56                           ArrayRef<unsigned> bankingDimensions = {});
57     std::unique_ptr<mlir::Pass> createIndexSwitchToIfPass();
58     std::unique_ptr<mlir::Pass> createHierarchicalRunner(
59         const std::string &topName,
60         llvm::function_ref<void(mlir::OpPassManager &)>> pipeline,
61         bool includeDoundInstances = false);
  
```

memory banking (kind of resource alloc.) in CIRCT

Short Summary of the CIRCT Project



CIRCT project is “**codified**” digital design experiences

“Large Circuit Model”: learned design experiences?

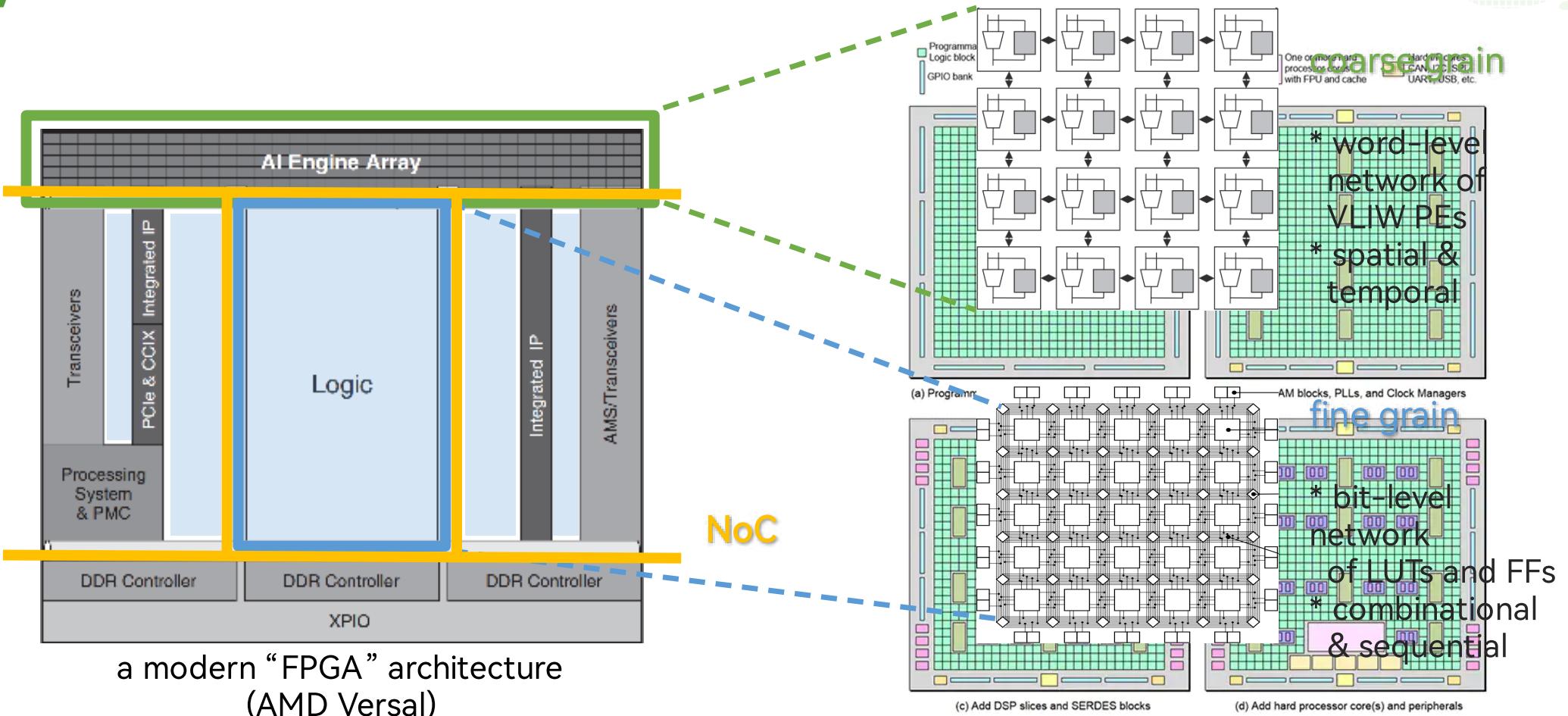
- ▶ L. Chen et al., “Large circuit models: opportunities and challenges,” SCIS vol. 67, no. 10, Oct. 2024.

Experiences in “good-for-hardware” patterns

James C. Hoe, "Lecture 6: Good-for-HW Computation Models," CMU ECE 643, Fall 2023.

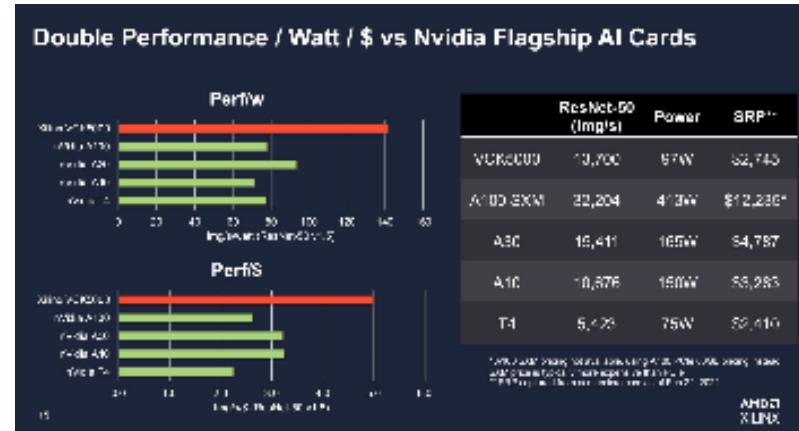
- ▶ FSM-Datapath
- ▶ Systolic Array
- ▶ Data Parallel
- ▶ Dataflow
- ▶ Stream Processing

可重构架构的演化 FPGA+CGRA

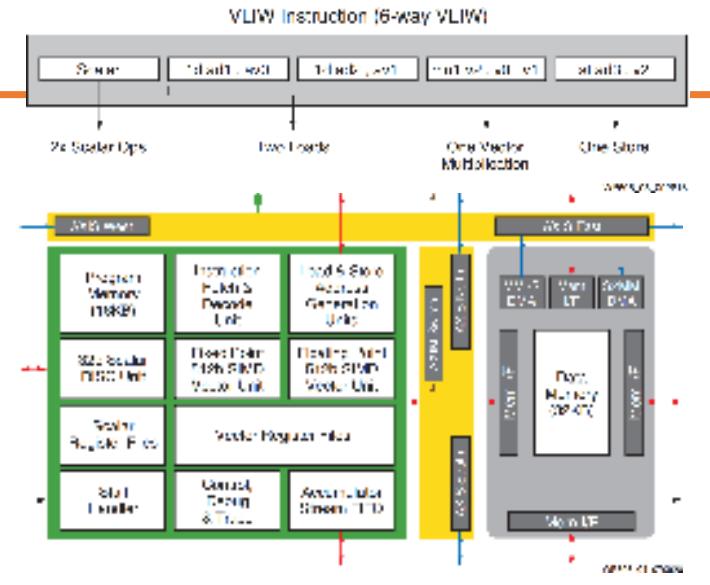


现代可重构计算的微架构

现代 PE 阵列架构的 VLIW 单元



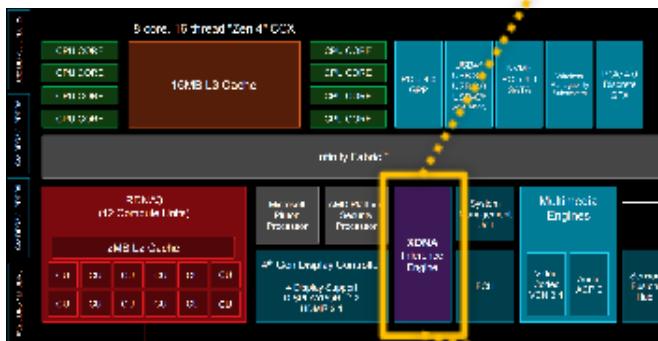
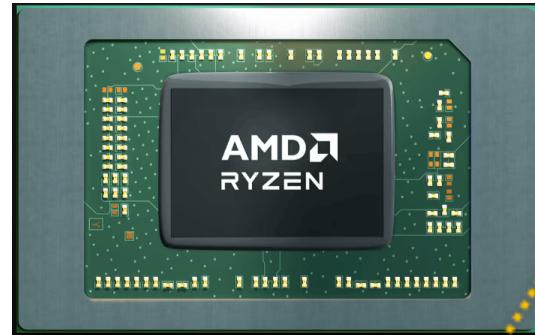
ResNet-50 推理任务：现代可重构架构（如 AMD Versal）的能效和性价比均优于 A100



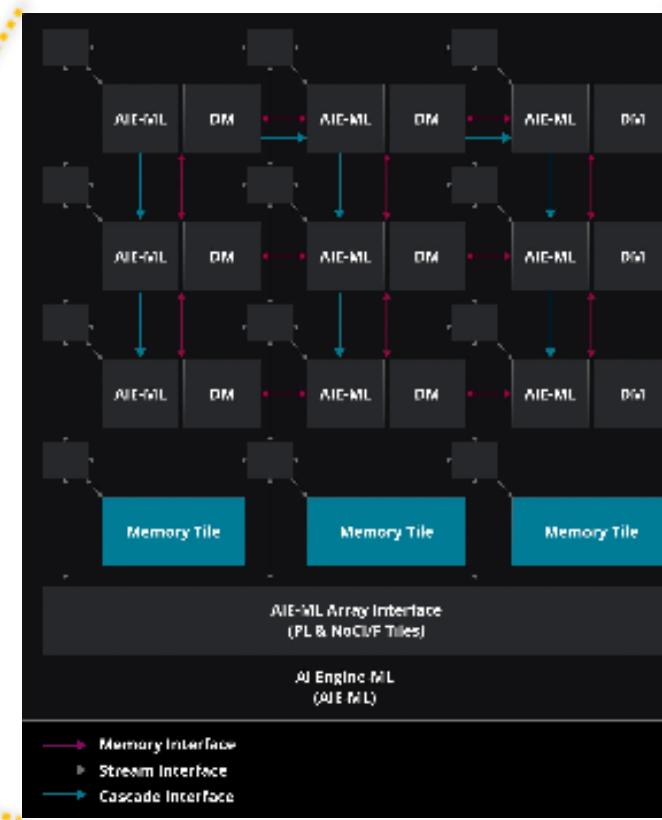
AI 引擎 (AIE) PU) 在 CPU 和 FPGA 的集成



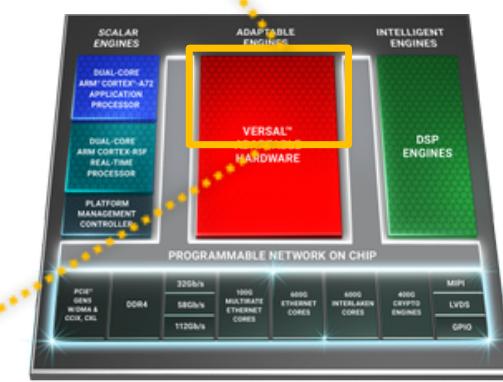
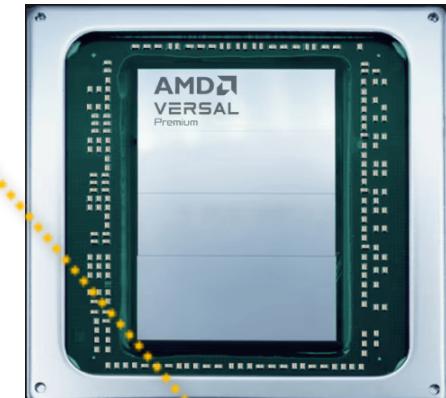
AMD Ryzen AI series
(CPU + AIE)



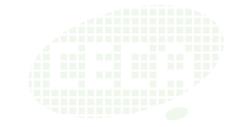
AI Engine 阵列



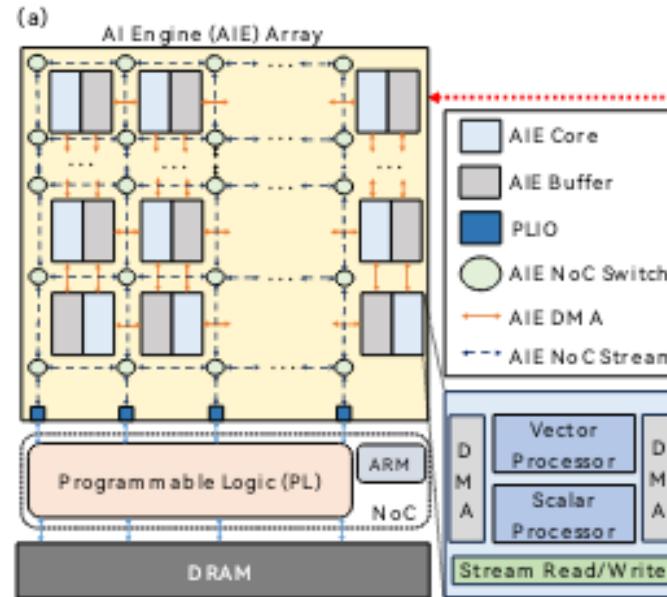
AMD Versal VU9P
(FPGA + AIE)



ACAP架构应用的编程模型



硬件结构	软件编程模型
AIE 阵列	ADF 图描述
AI Engine	AIE 计算核代码



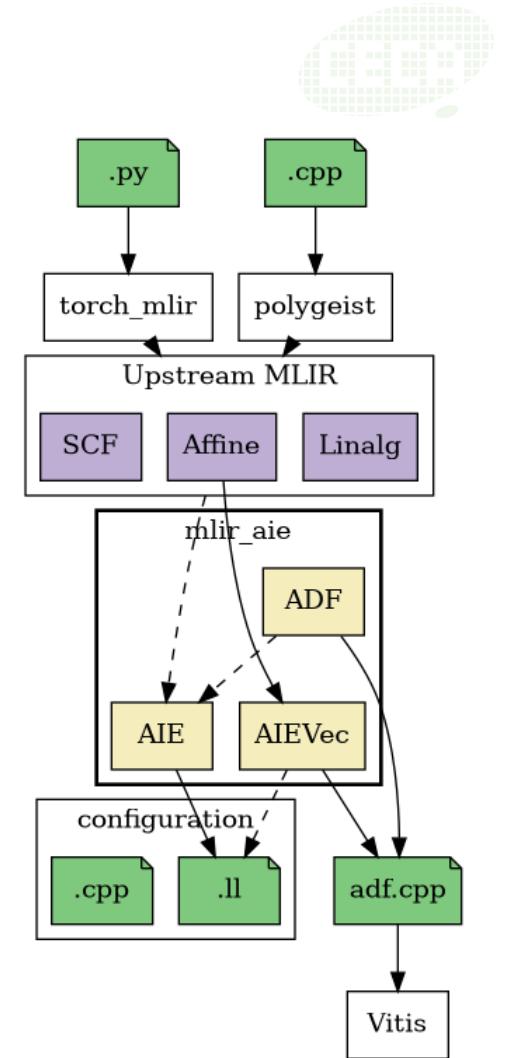
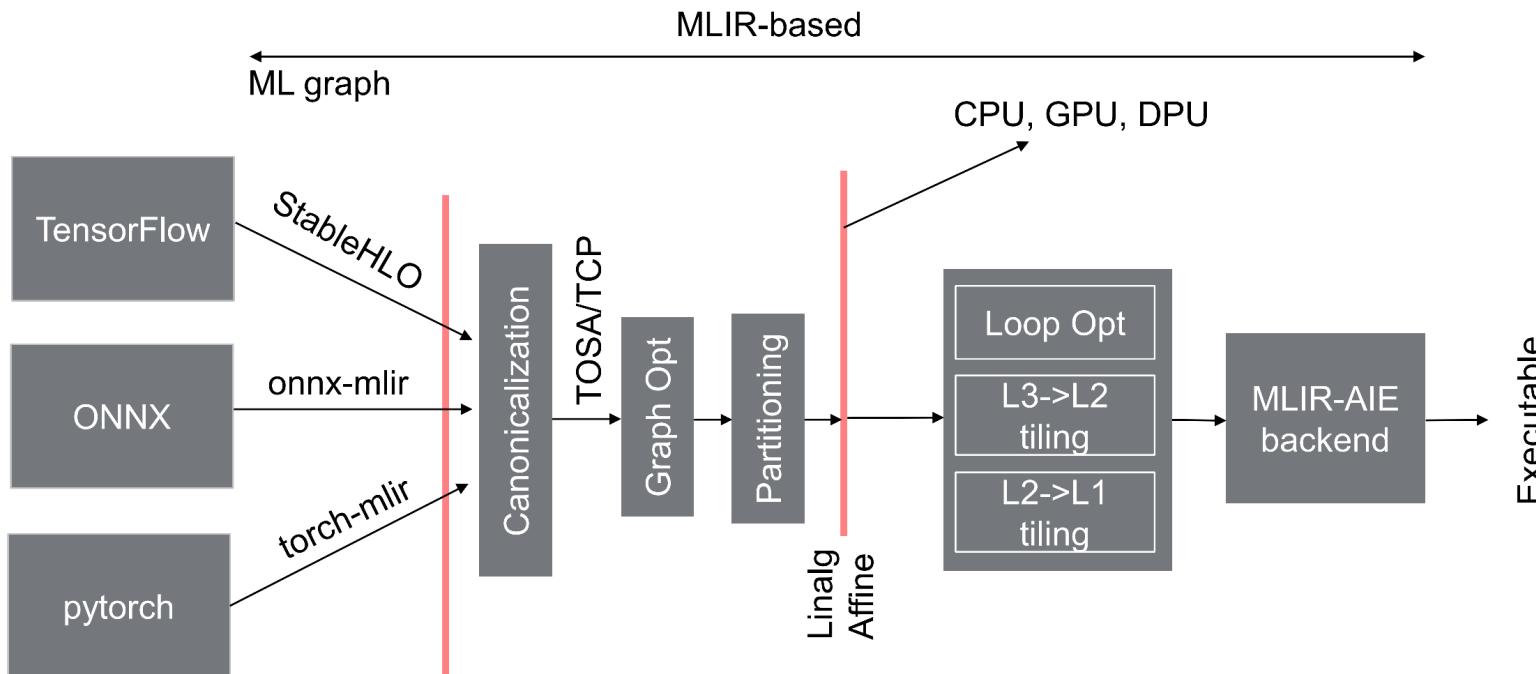
(b)

```
// AIE Kernel Creations          Adaptive Dataflow (ADF)
kernels = kernel::create(funcs);      Graph Programming
source<kernels> = "kernels.cpp"
// Inter-Core Connections via NoC
connect<pktstream, window<SIZE>> (src, dst);
// AIE Core Location Allocation
location<kernel>(kernels) = tile(pe_x, pe_y);
// Stack and Buffer Allocation
location<stack>(kernels) = location<kernel>(kernels);
location<buffer>(port) = location<kernel>(kernels);
.....
graph.h
```

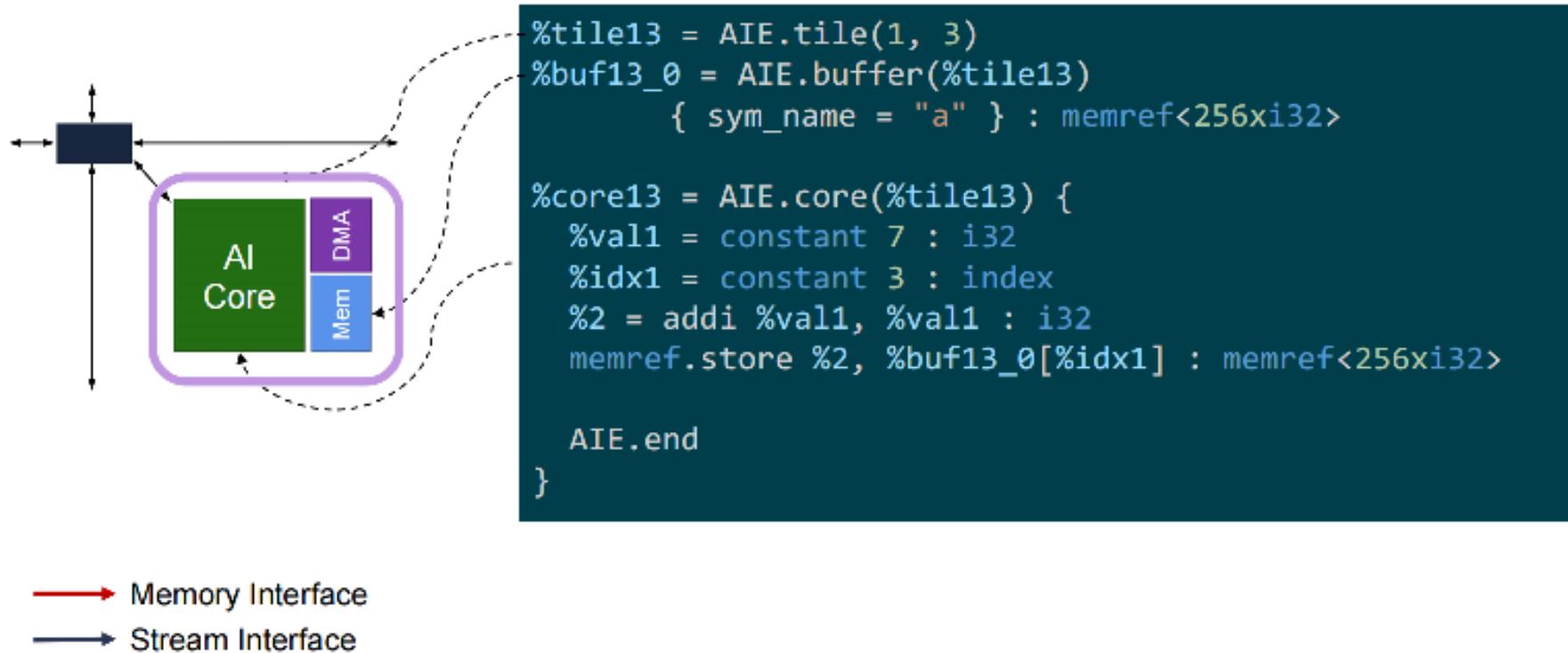


```
void funcs(IN, ..., OUTPUT){           AIE Kernel Programming
// AIE Kernel Description           with Intrinsic
for(int i = 0; i < M; ++i){
    BUF = fpmac(BUF, IN2, IDX1, ADDR, IN1, IDX0, 0x0);
    BUF = upd_v(BUF, IDX2, window_read_v4(IN1));
    window_incr(IN, M);
    window_write(OUTPUT, BUF);  }
kernels.cpp
```

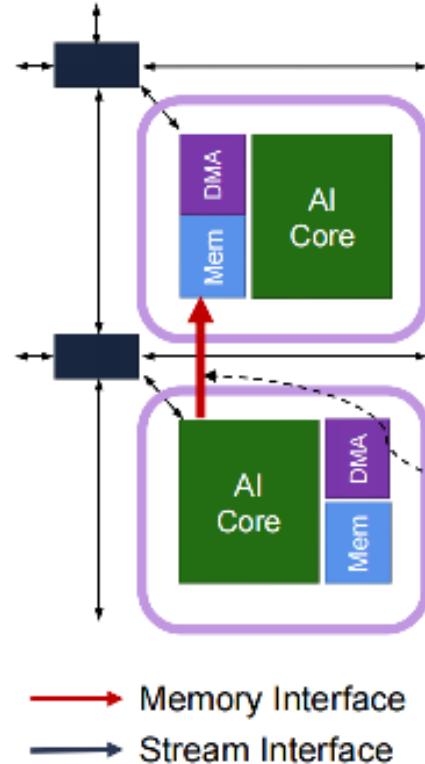
MLIR-AIE



Running Code on a Core



Moving Buffers



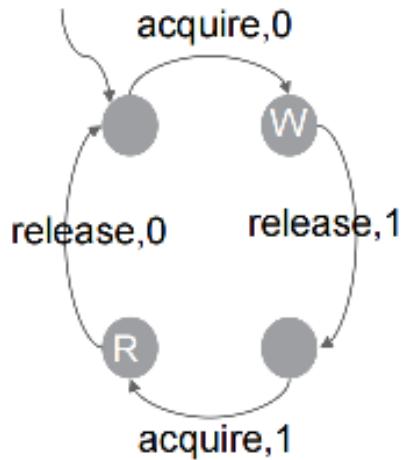
```
%tile14 = AIE.tile(1, 4)
%buf = AIE.buffer(%tile14)
  { sym_name = "a" } : memref<256xi32>

%tile13 = AIE.tile(1, 3)
%core13 = AIE.core(%tile13) {
  %val1 = constant 7 : i32
  %idx1 = constant 3 : index
  %2 = addi %val1, %val1 : i32

  memref.store %2, %buf[%idx1] : memref<256xi32>

AIE.end
}
```

Synchronization with Locks



```

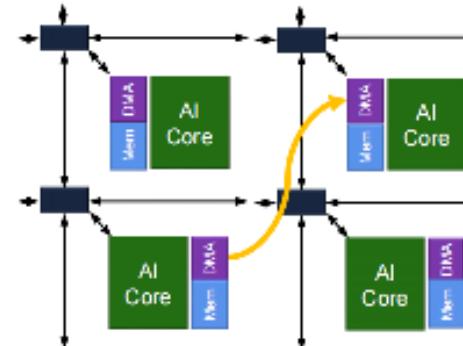
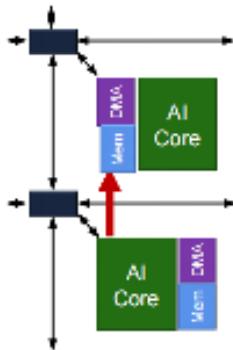
%tile14 = AIE.tile(1, 4)
%lock = AIE.lock(%tile14) { sym_name = "a_lock" }
%buf = AIE.buffer(%tile14)
{ sym_name = "a" } : memref<512xi32>

%tile13 = AIE.tile(1, 3)
%core13 = AIE.core(%tile13) {
  AIE.useLock(%lock, "Acquire", 0)
  memref.store %2, %buf[%idx1] : memref<512xi32>
  AIE.useLock(%lock, "Release", 1)
}

%core14 = AIE.core(%tile14) {
  AIE.useLock(%lock, "Acquire", 1)
  %data = memref.load %buf[%idx1] : memref<512xi32>
  AIE.useLock(%lock, "Release", 0)
}
  
```

Object FIFO Lowers to Communication Components Based on the Endpoints

aie-opt --aie-objectFifo-stateful-transform

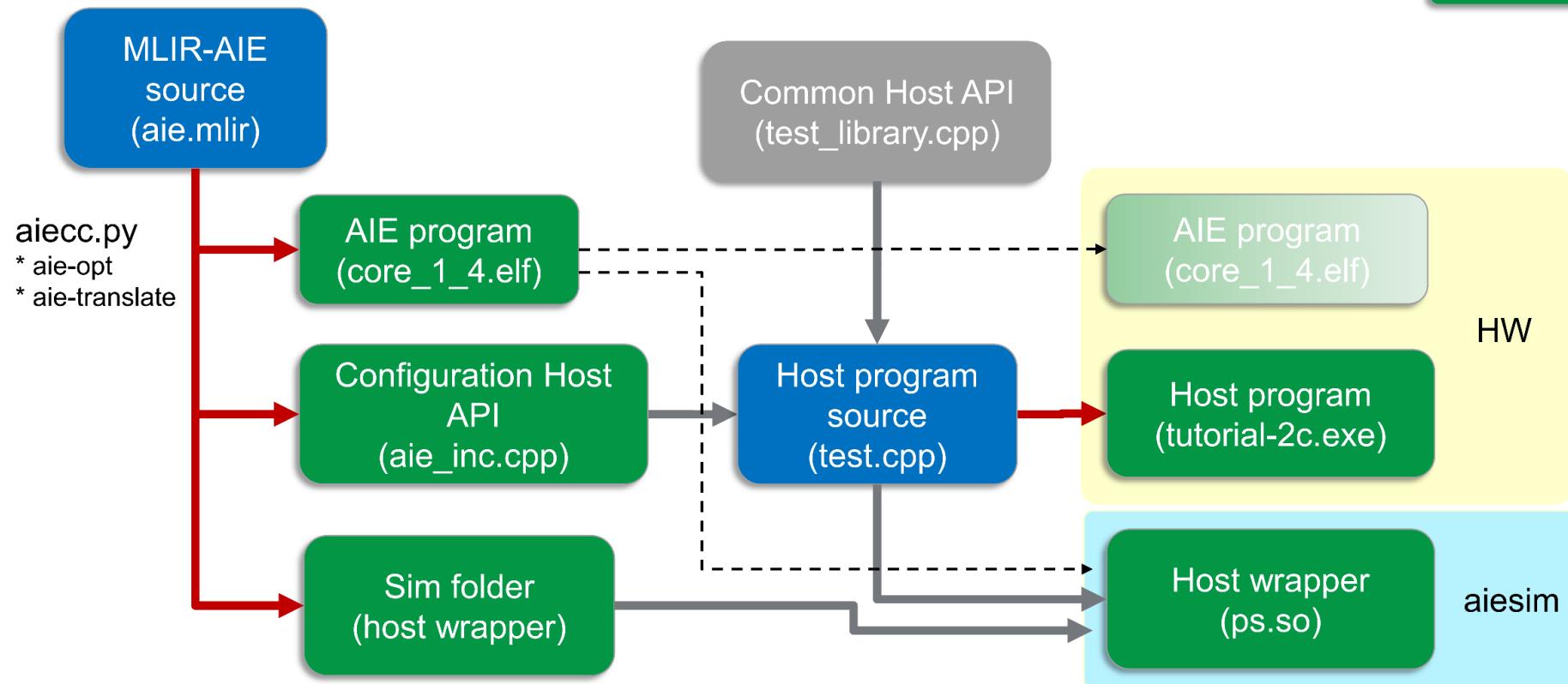


```
%4 = AIE.buffer(%0) : memref<16xi32>
%5 = AIE.lock(%0, 0)
%6 = AIE.buffer(%0) : memref<16xi32>
%7 = AIE.lock(%0, 1)
%8 = AIE.core(%0) { ... }
%9 = AIE.core(%1) { ... }
```

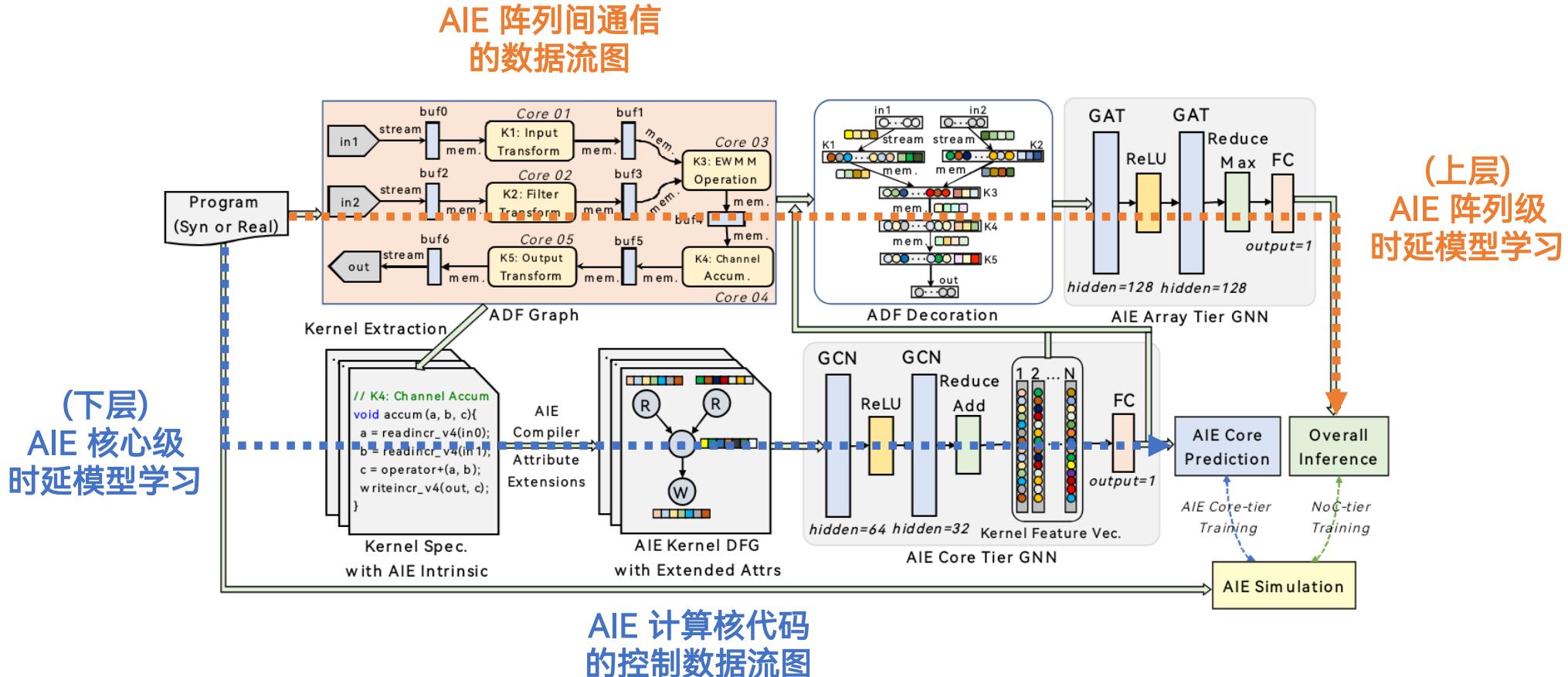
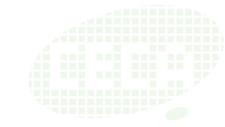
```
%6 = AIE.buffer(%0) : memref<16xi32>
%7 = AIE.lock(%0, 0)
%8 = AIE.buffer(%0) : memref<16xi32>
%9 = AIE.lock(%0, 1)
%10 = AIE.buffer(%2) : memref<16xi32>
%11 = AIE.lock(%2, 1)
%12 = AIE.buffer(%2) : memref<16xi32>
%13 = AIE.lock(%2, 2)
%14 = AIE.core(%0) { ... }
%15 = AIE.core(%2) { ... }
%16 = AIE.mem(%0) { ... }
%17 = AIE.mem(%2) { ... }
AIE.flow(%0, DMA : 0, %1, DMA : 0)
```



Lay of the Land (Summary)



相关工作：G2PM 基于双层图学习的性能模型



相关工作：WideSA - AIE阵列映射器（结果节选）



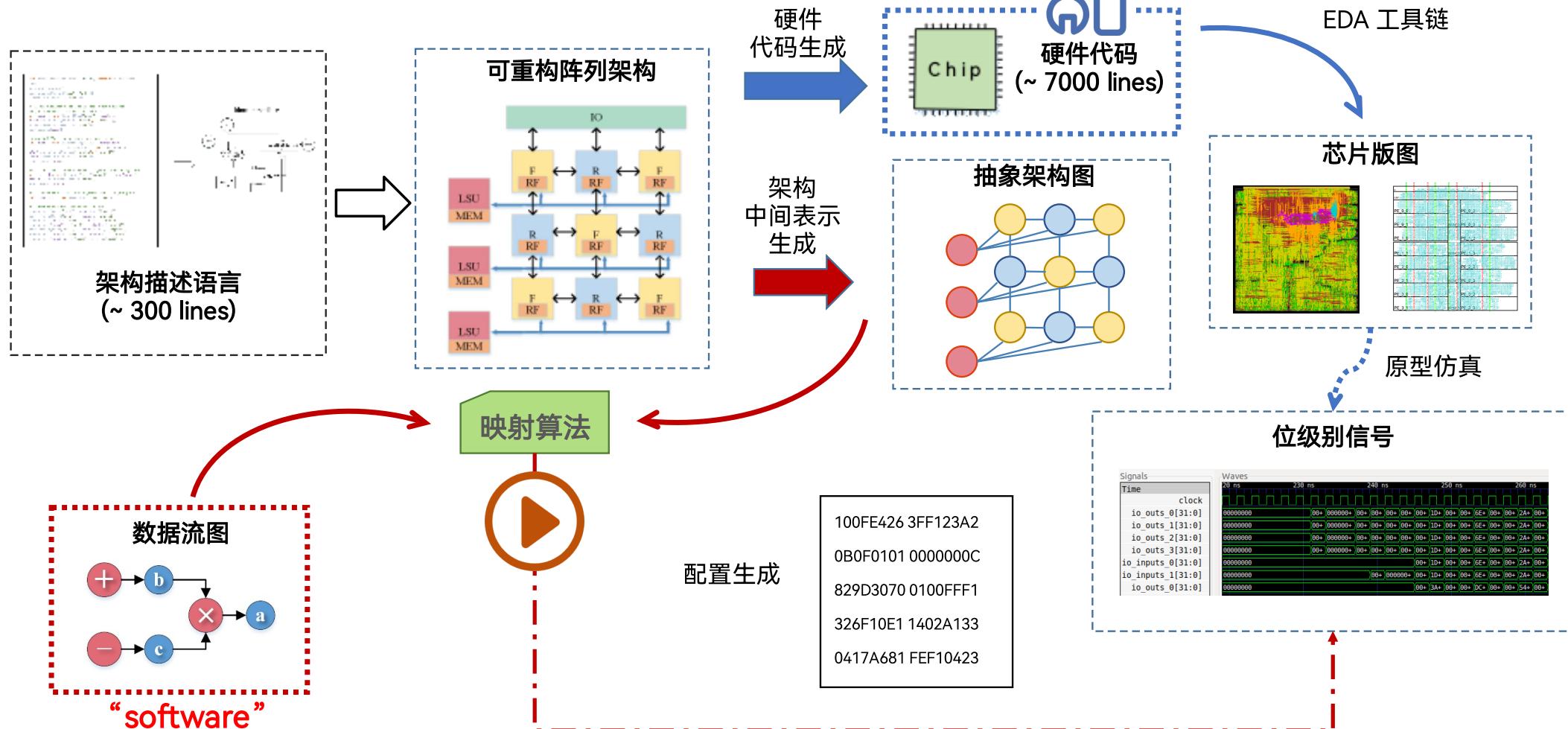
整体性能比较：达到100% AIE 利用率和36.02 TOPS 峰值性能

方法	指标	MM				2D-Conv				2D-FFT		FIR Filter			
		Float	Int8	Int16	Int32	Float	Int8	Int16	Int32	Cfloat	Cint16	Float	Int8	Int16	Cfloat
基线	#AIEs	384	384	384	384	-	256	-	-	10	10	10	10	10	10
	TOPS	3.73	29.78	7.82	3.72	-	31.40	-	-	0.04	0.13	0.15	2.56	0.62	0.15
	TOPS/#AIEs	0.010	0.077	0.020	0.010	-	0.123	-	-	0.004	0.013	0.015	0.256	0.062	0.015
WideSA	#AIEs	400	400	400	400	400	400	400	400	320	320	256	256	256	256
	TOPS	4.15	32.49	8.10	3.92	4.50	36.02	10.35	4.48	1.10	3.83	2.92	39.3	9.47	2.89
	TOPS/#AIEs	0.010	0.081	0.020	0.010	0.011	0.090	0.025	0.011	0.003	0.012	0.012	0.100	0.037	0.011

矩阵乘法与常规 FPGA 加速器比较：2.25倍的能效比提升

	PL-only				WideSA			
数据类型	Float	Int8	Int16	Int32	Float	Int8	Int16	Int32
DSPs	1536	1528	1516	1536	152	60	67	65
#AIEs	0	0	0	0	400	400	400	400
吞吐量	0.59	5.77	2.16	0.60	4.15	32.49	8.10	3.92
功率(W)	19.5	18.8	18.6	19.5	55.8	54.4	54.9	55.6
能效	0.03	0.31	0.12	0.03	0.07	0.60	0.15	0.07
能效比	1×	1×	1×	1×	2.25×	1.94×	1.29×	2.25×

Tile Arch. Generator: Workflow



总结：CIRCT 和 MLIR-AIE



CIRCT: Circuit IR Compilers and Tools

- ▶ <https://circt.llvm.org/>
- ▶ The CIRCT project is an (experimental!) effort looking to apply MLIR and the LLVM development methodology to the domain of hardware design tools.

MLIR-AIE: Toolchain for AMD AI Engine devices

- ▶ <https://xilinx.github.io/mlir-aie/>
- ▶ MLIR-AIE is a toolchain providing low-level device configuration for Versal AIEngine-based devices. Support is provided to target the AIEngine portion of the device, including processors, stream switches, TileDMA and ShimDMA blocks. Backend code generation is included, targeting the LibXAIE library, along with some higher-level abstractions enabling higher-level design.