

Accelerating Edge AI with Morpher: An Integrated Design, Compilation and Simulation Framework for CGRAs

Dhananjaya Wijerathne, Zhaoying Li , Tulika Mitra
School of Computing, National University of Singapore
{dmd, zhaoying, tulika } @comp.nus.edu.sg



Domain Specific vs Reconfigurable Accelerators

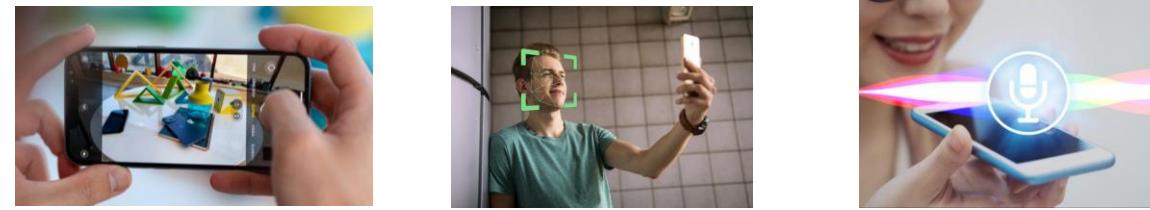
Tiny IoT devices have ultra-low area and power footprint need

Several domain-specific accelerators plus multiple general-purpose cores on SoC are not tenable for such devices

Domain Specific Accelerators



Apple iPhone 12's A14 Bionic SoC
with
Domain Specific Accelerators



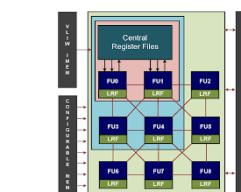
Video

ML

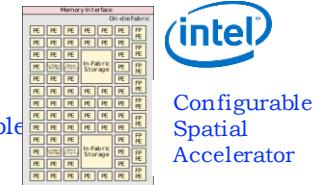
Audio



Field Programmable
Gate Arrays (FPGA)



SAMSUNG
Samsung
Reconfigurable
Processor

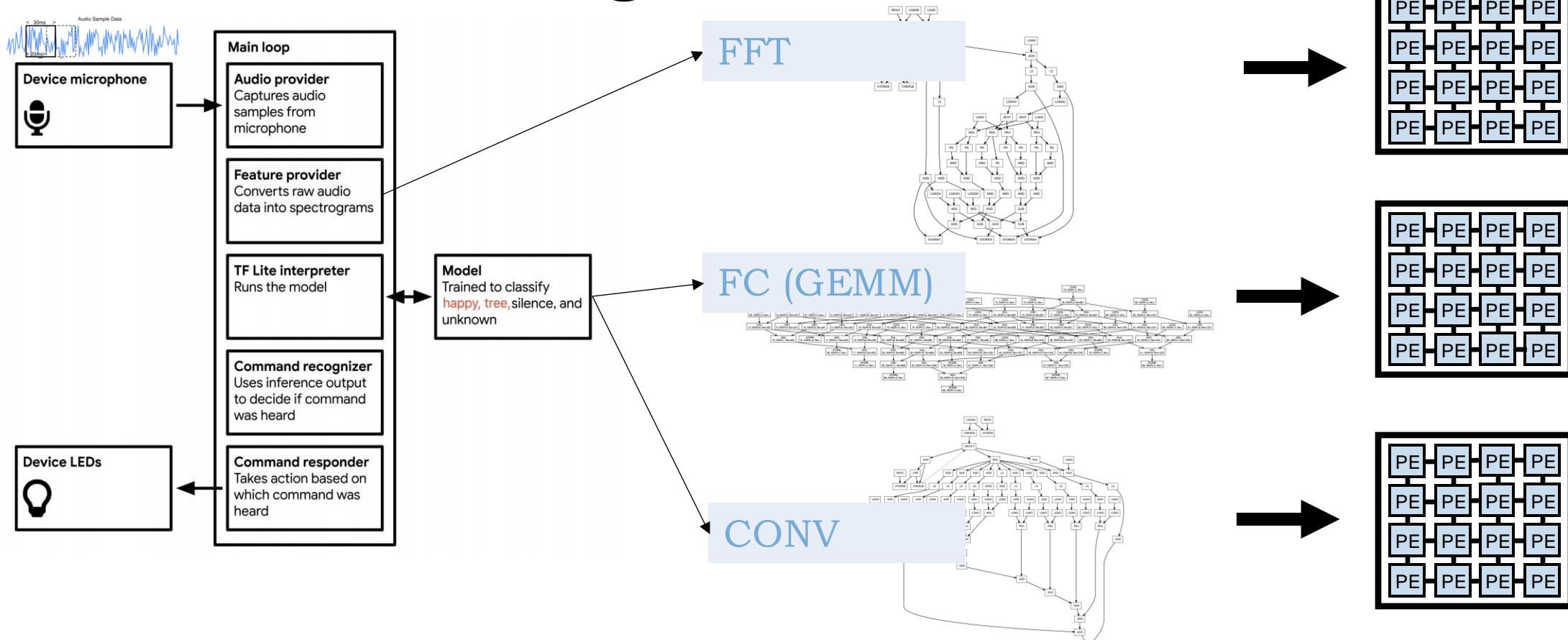


intel
Configurable
Spatial
Accelerator

Coarse-Grained Reconfigurable
Array (CGRA)



CGRA: Supporting Diverse Dataflows



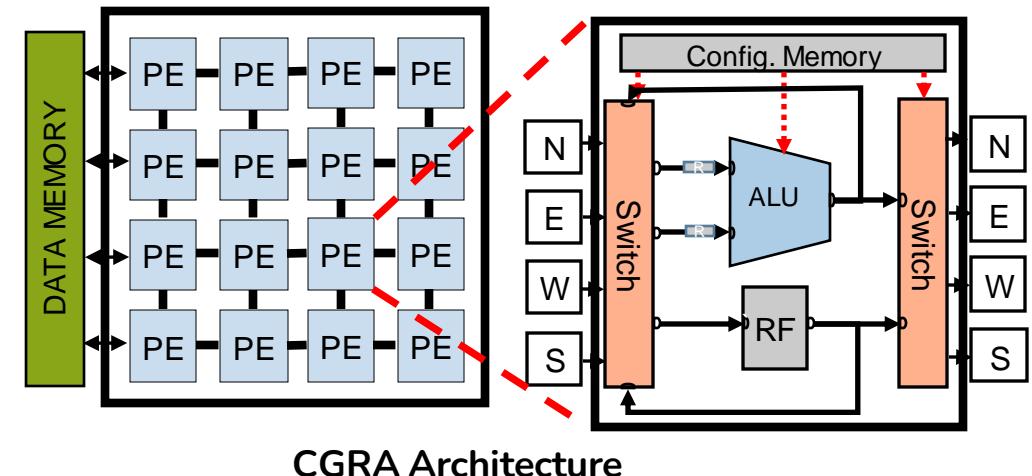
Speech Recognition Model

Application Kernels

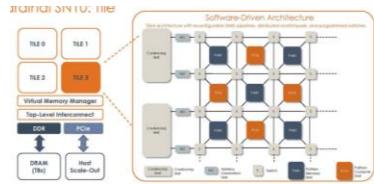
Different Dataflows

Coarse-Grained Reconfigurable Array (CGRA)

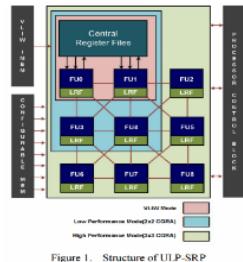
- **CGRA Architecture:** A power-efficient mesh of processing elements, each equipped with an ALU, register file, configuration memory, switches, and on-chip memory.
- **Execution Model:** CGRAs' flexibility lies in the compiler-generated execution schedules and configurations, adapting to different application kernels mapped onto the CGRAs.



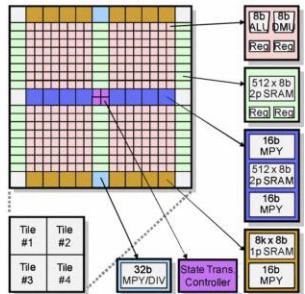
Coarse-Grained Reconfigurable Array (CGRA)



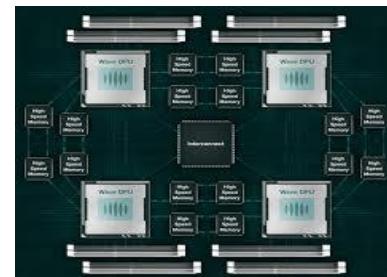
Sambanova
Reconfigurable
Dataflow Unit (RDU)



Samsung
Reconfigurable
Processor (SRP)

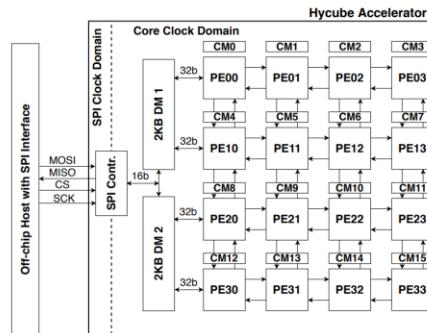


Renesas Dynamic
Reconfigurable Processor
(DRP)

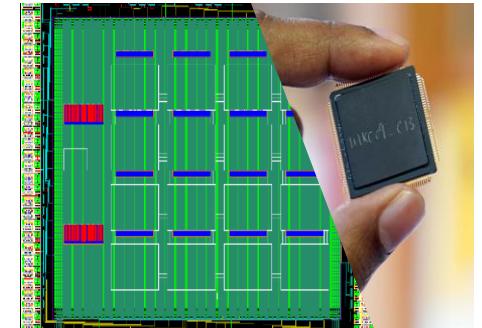


Wave DPU

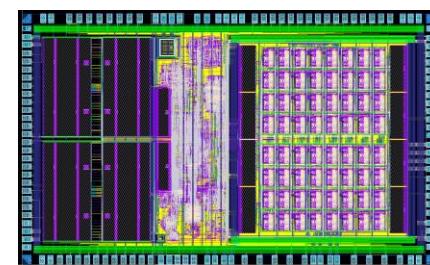
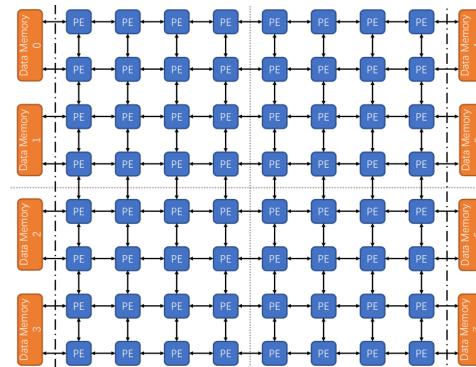
Commercial CGRAs



Hycube Accelerator

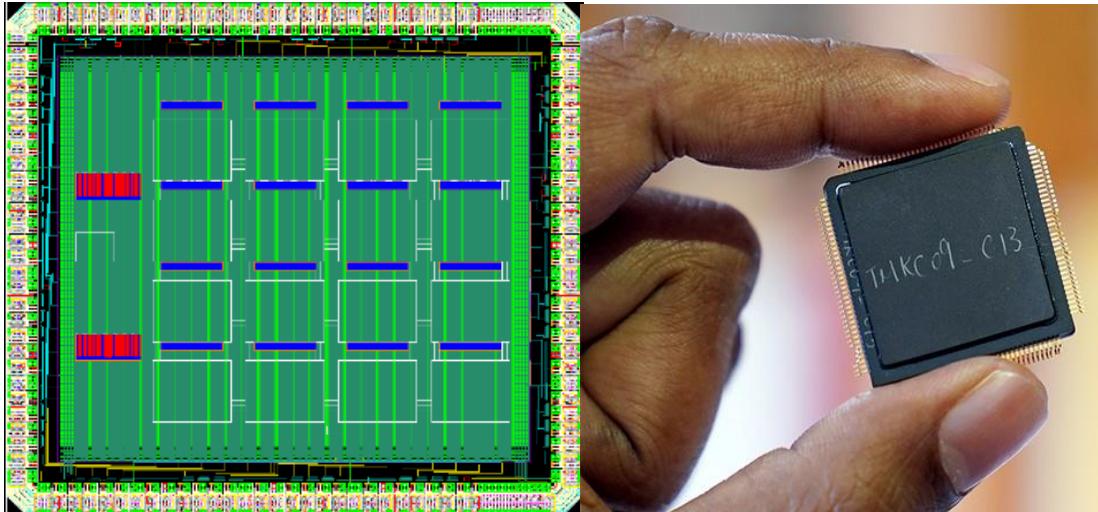


NUS HyCUBE



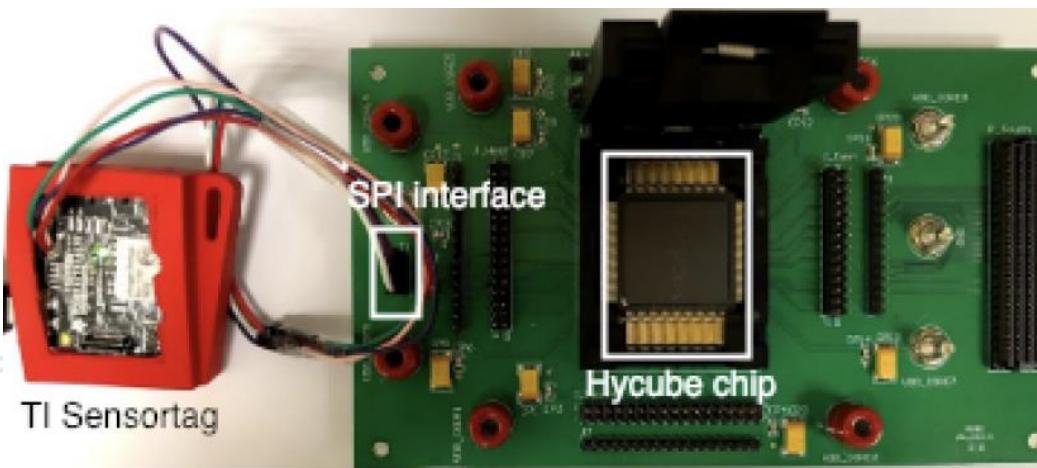
NUS PACE

HyCUBE CGRA @ NUS 2019



4x4 CGRA, TSMC 40nm

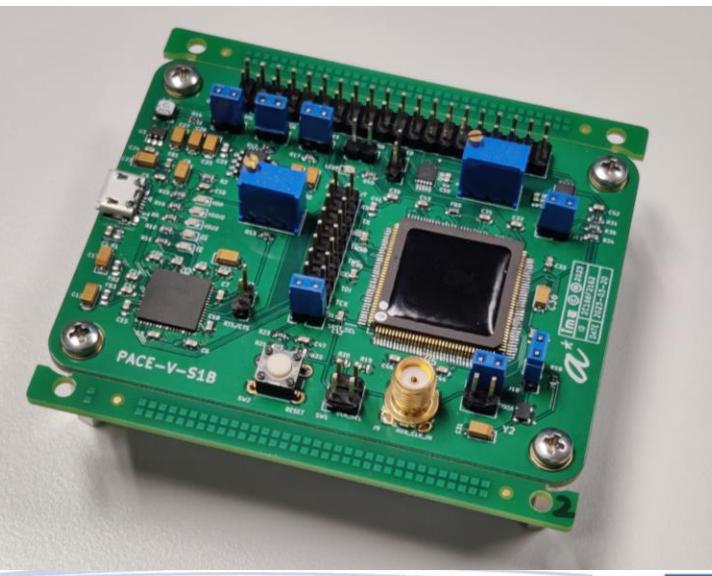
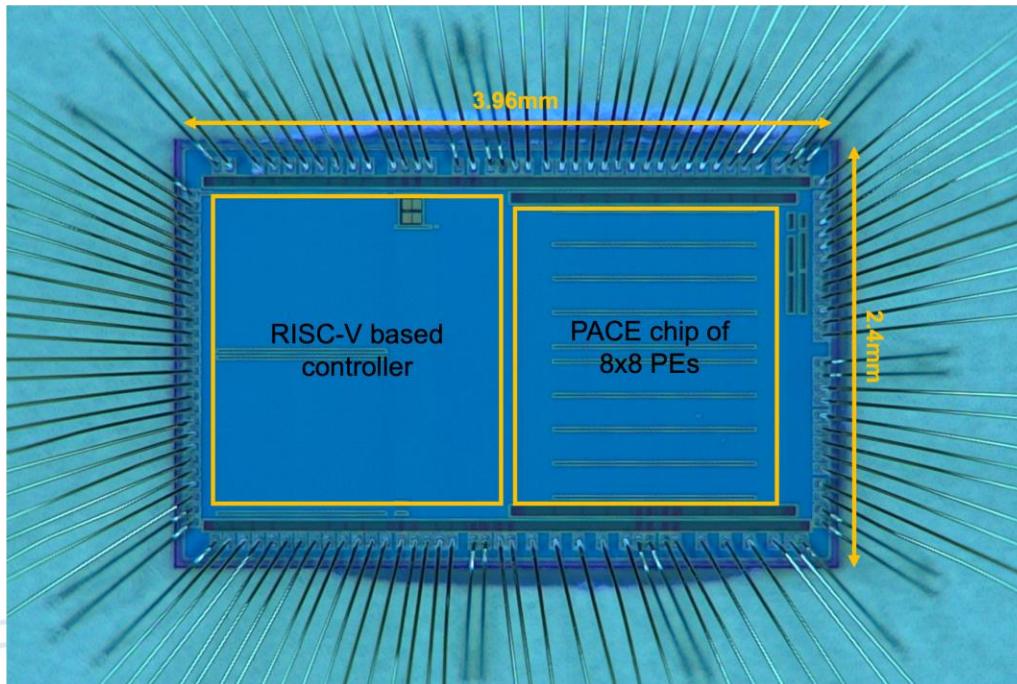
NUS HyCUBE: 90 MOPS/mW
Samsung SRP: 22 MOPS/mW
ARM CPU: 2.6 MOPS/mW
Xilinx FPGA: 25 MOPS/mW



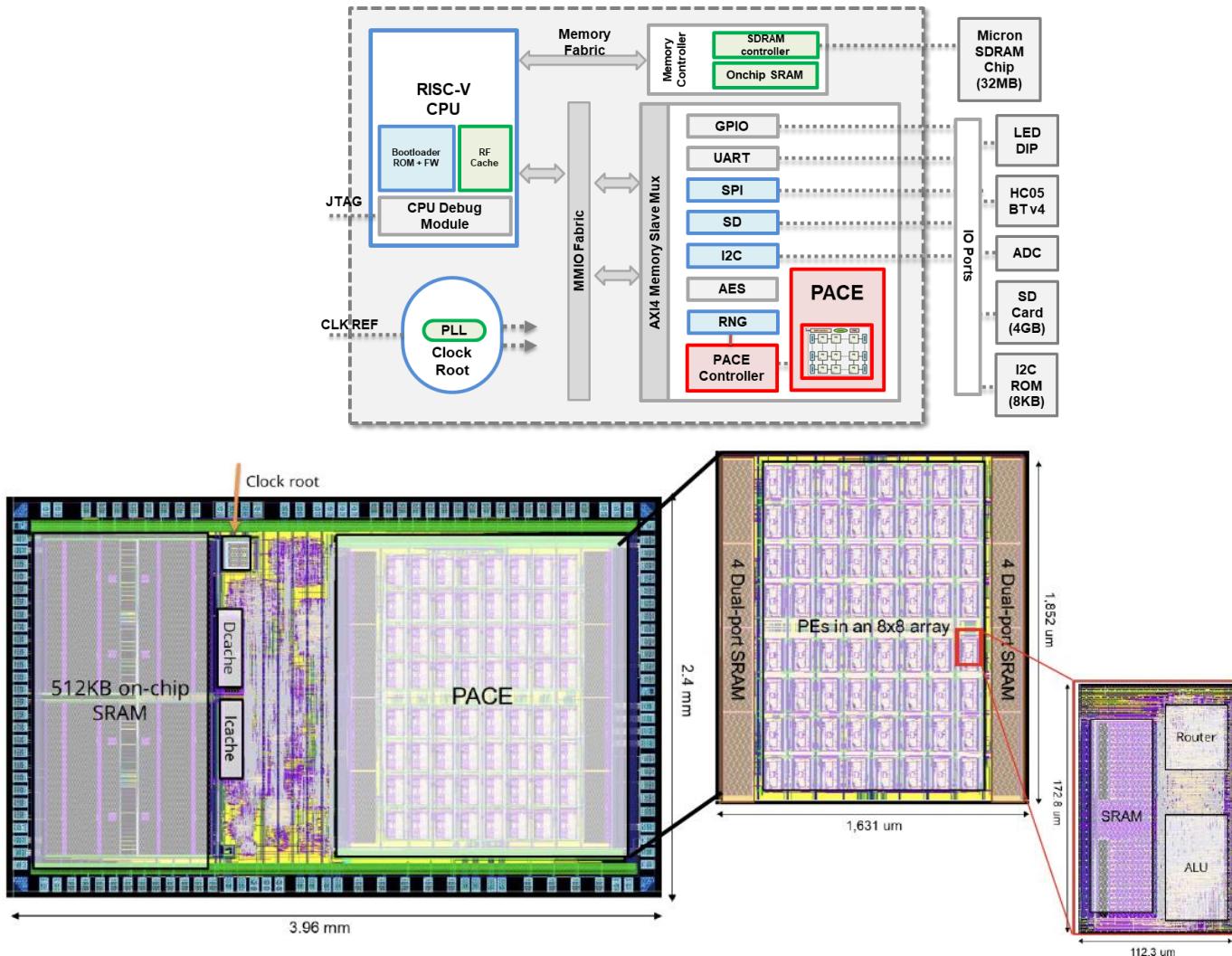
PACE CGRA+RISC-V SoC: NUS + IME 2023



8x8 CGRA:
582 GOPS/W at 0.45V, 40nm ULP
1.1 mW at 10MHz
Estimated 1 TOPS/W at 22nm

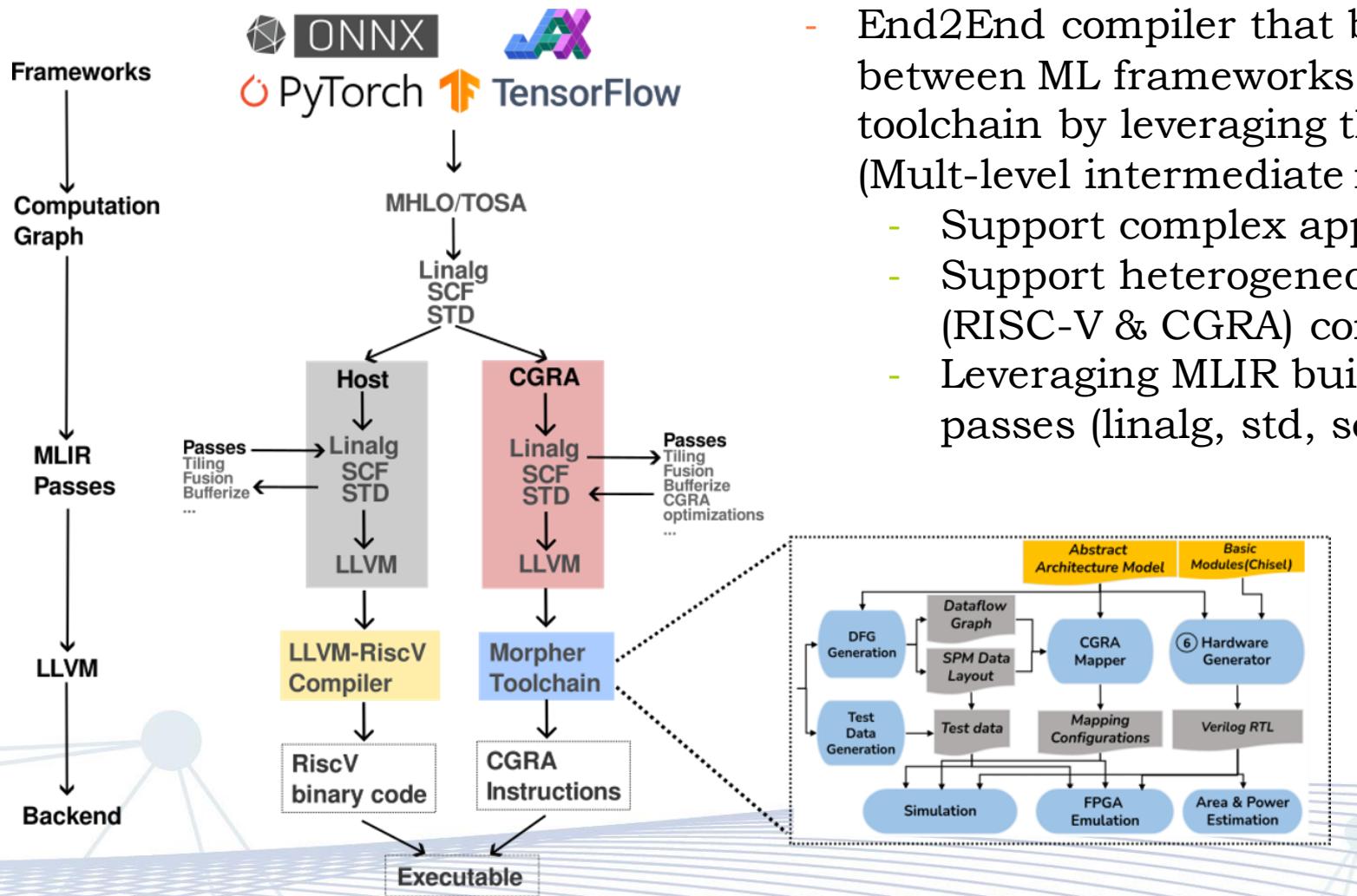


PACE CGRA SoC: NUS + IME



PACE SoC SPECIFICATIONS (SIMULATION)	
RISC-V based controller	
Tech. Node	UMC 40nm ULP
Area	3.96mm x 2.4mm (Die) 3.80mm x 2.0mm (Core)
Voltage	1V/3.3V
Memory	512 KB on-chip memory 16KB ICache 16KB DCache
Frequency	200/100 MHz
Power	14mW @ 1.0V
#PEs	8x8 (17 instructions, including NOP)
Frequency	10 MHz
Memory	8 x 8KB data memory 64 x 0.25KB config memory
Power	1.1 mW @ 0.45V
Efficiency	582 TOPS/W

NUS MLIR-Based End-to-End Compiler



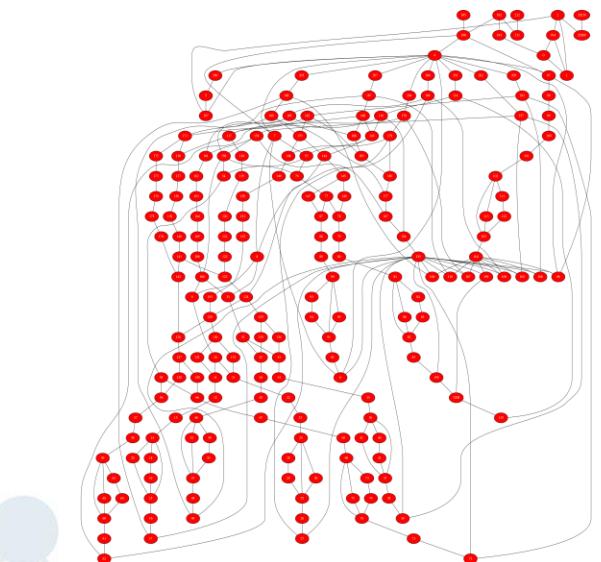
- End2End compiler that bridges the gap between ML frameworks and the Morpher toolchain by leveraging the power of MLIR (Mult-level intermediate representation).
 - Support complex applications
 - Support heterogeneous system (RISC-V & CGRA) compilation
 - Leveraging MLIR builtin dialects and passes (linalg, std, scf...)

Application mapping on CGRA

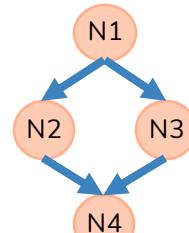
- Target: a loop kernel from applications
- Mapping the dataflow graph (DFG) of the loop body on to the CGRA
 - Placement: assigning DFG operations to ALUs
 - Routing: mapping data signals using wires and registers

```
static void
main(void)
{
    int i;
    tnt16 *src0 = gCoeffBuf;
    For (i = 0; i < 8; i++)
    {
        #ifdef CGRA_COMPILER
        /* Compute the dot product of the row vector and column vector */
        /* and store it in the result vector. */
        if (src0[i] != 0)
            psrc[i] = src0[i] * psrc[0] + psrc[1] * psrc[1] + psrc[2] * psrc[2] + psrc[3] * psrc[3] + psrc[4] * psrc[4] + psrc[5] * psrc[5] + psrc[6] * psrc[6];
        else
            psrc[i] = psrc[0];
        #else
        /* Short circuit the ID DCT if only the DC component is non-zero */
        if (src0[i] != 0)
            psrc[i] = src0[i] * psrc[0];
        else
            psrc[i] = psrc[0];
        #endif
    }
}
```

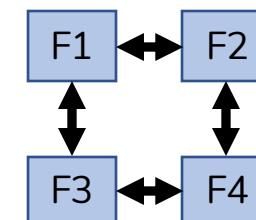
Loop kernel c code



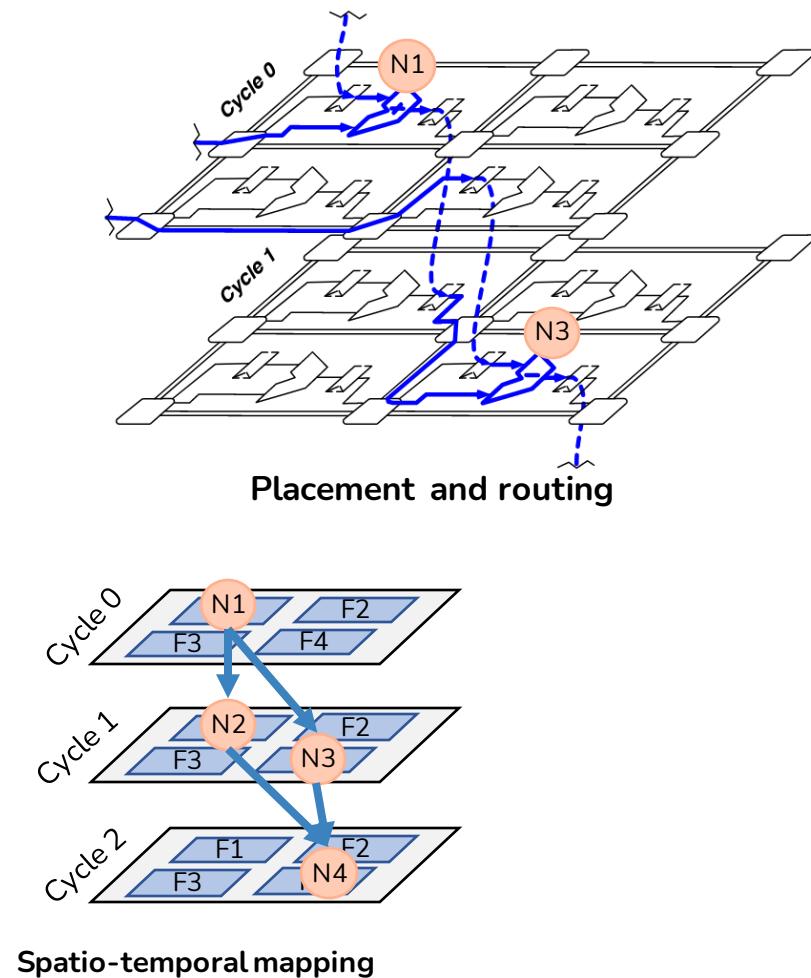
Dataflow graph (DFG) of the loop body



DFG



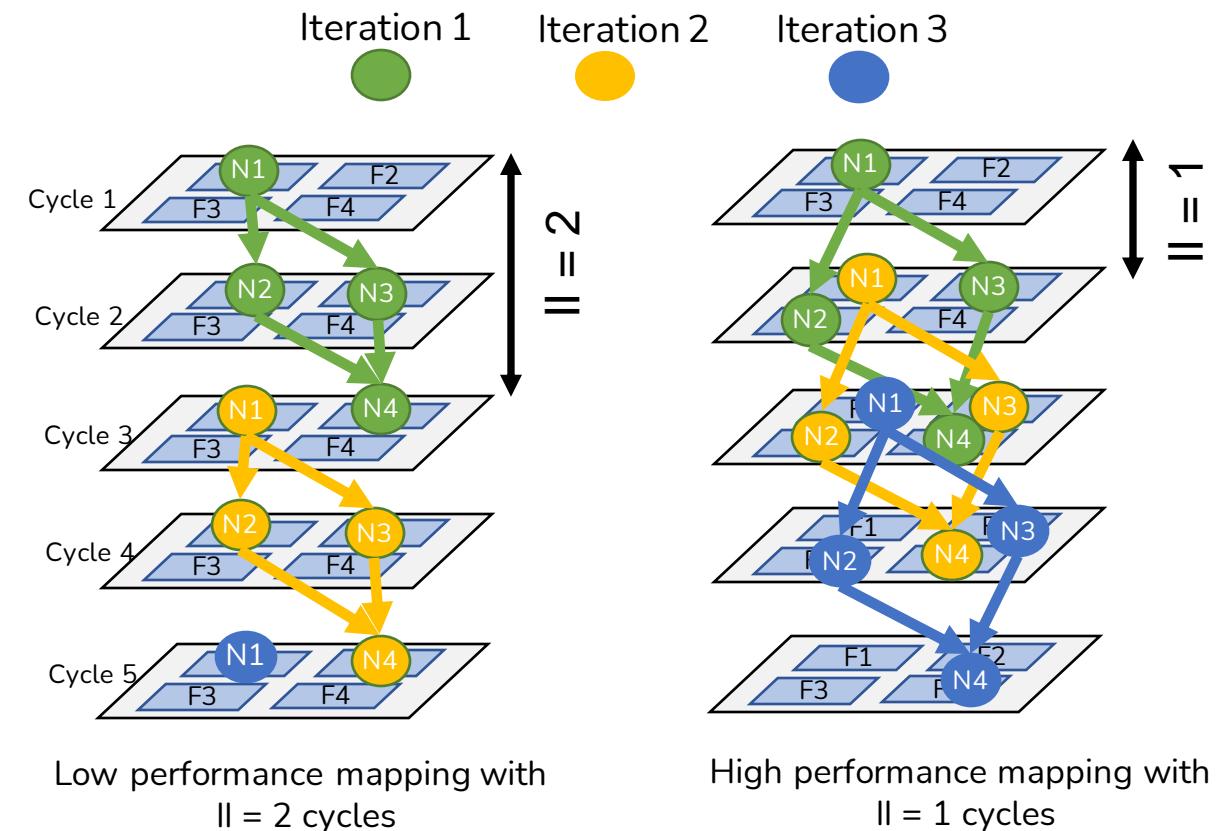
2x2 CGRA



Spatio-temporal mapping

Application mapping on CGRA

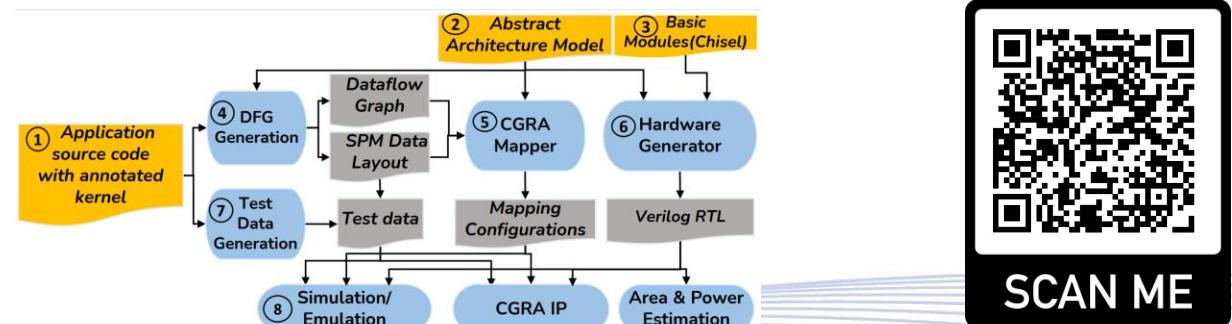
- Software pipelined schedule
- Goal: Mapping with minimum initiation interval
- Initiation interval (II) = cycle difference between initiation of consecutive iterations
- Low II -> High performance



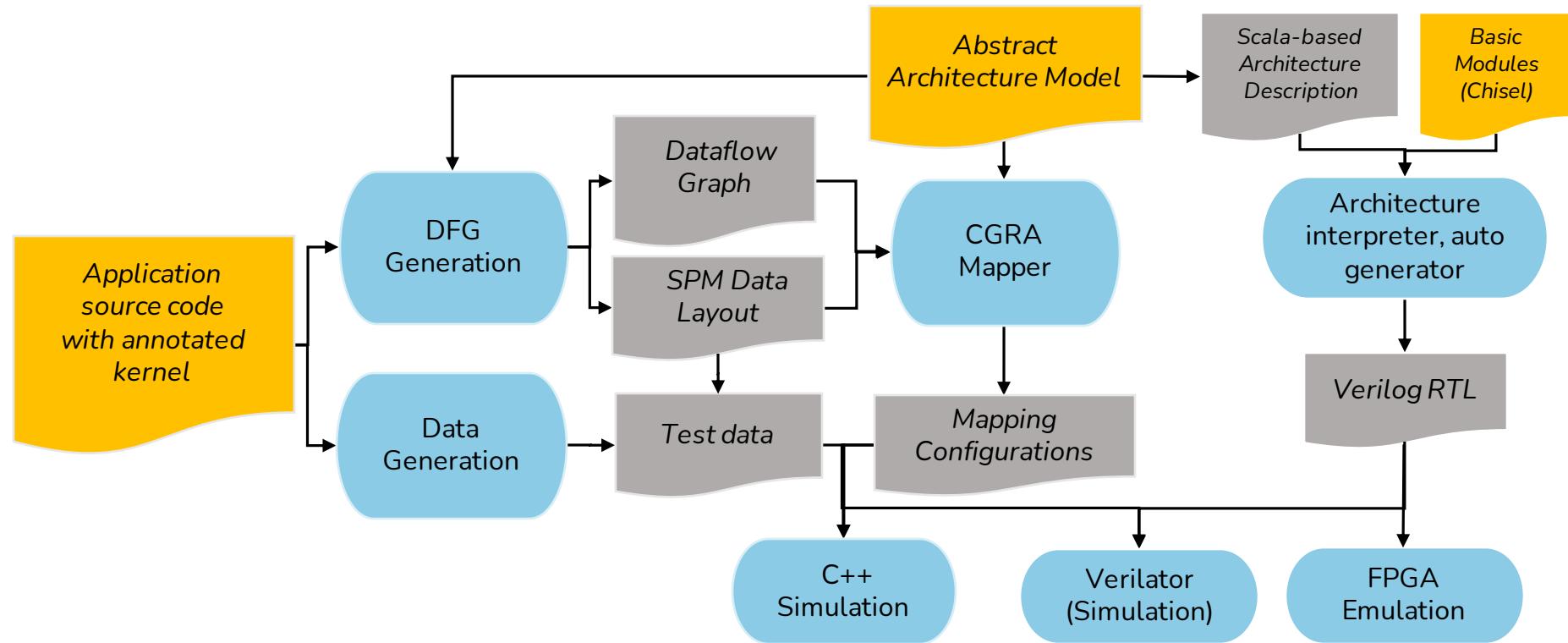
Morpher: An Integrated Compilation and Simulation Framework for CGRA

- Fully automated end-to-end CGRA compilation, architecture generation and simulation framework
 - Flexible architecture specification language
 - Efficient mapping algorithms
 - Automatic RTL generation & cycle-accurate simulation to validate the compilation results
 - Fully open-source with easily modifiable modular code base

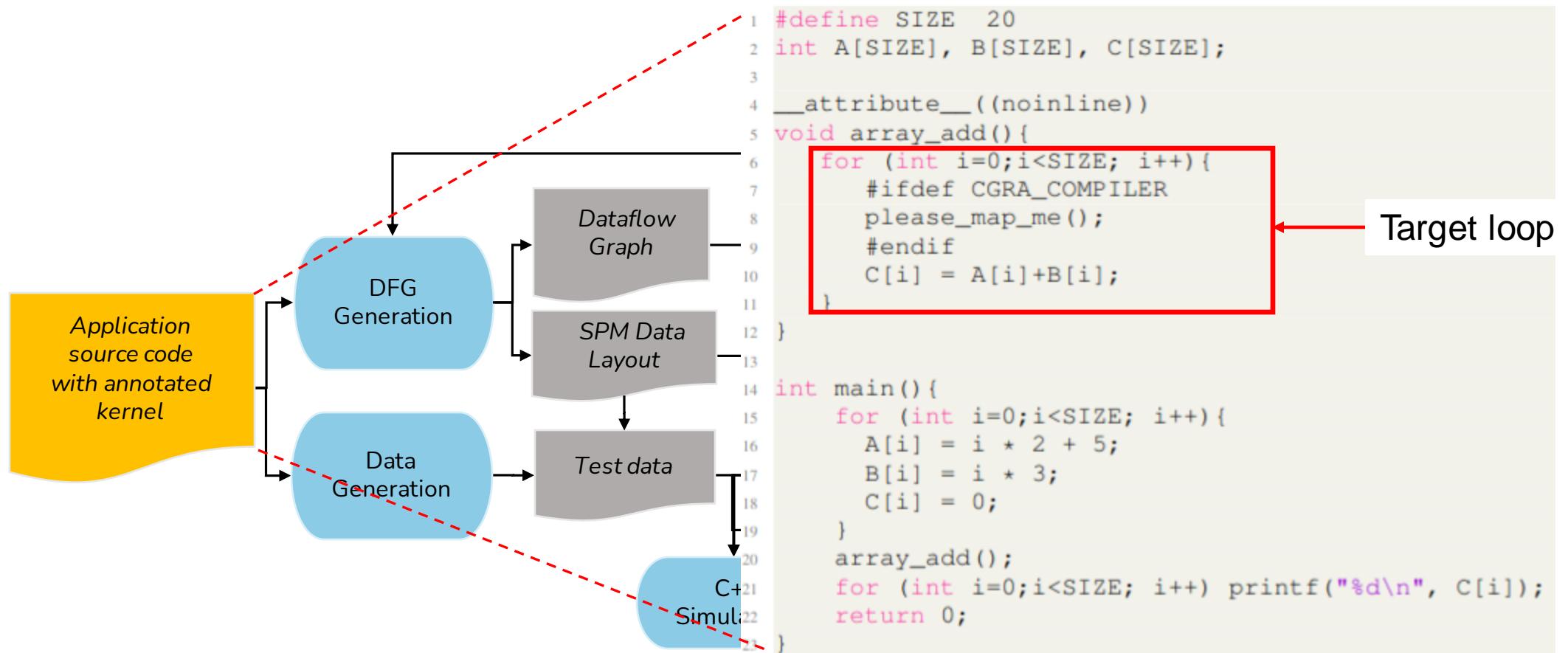
	Features	CGRA-ME	Pillars	Open-CGRA	CCF	Morpher
DFG Generation	Models control divergence	X	X	✓	✓	✓
	Recurrence edges	X	X	✓	✓	✓
Architecture Modelling	Adapt user defined architectures	✓	✓	✓	X	✓
	Multi-hop connections	X	X	X	X	✓
P&R Mapper	Different memory organizations	X	X	✓	X	✓
	Architecture adaptive mapping	✓	✓	X	X	✓
	Data layout aware mapping	X	X	X	X	✓
Simulation & validation	Recurrence aware mapping	X	X	✓	✓	✓
	Cycle accurate simulation	X	✓	✓	✓	✓
	Test data generation	X	X	X	X	✓
Hardware Generation	Validation against test data	X	X	X	X	✓
	Generate RTL	✓	✓	✓	X	✓
	Infer control paths	X	✓	✓	X	✓
	Infer multiplexers	X	X	X	X	✓



Overview



Overview



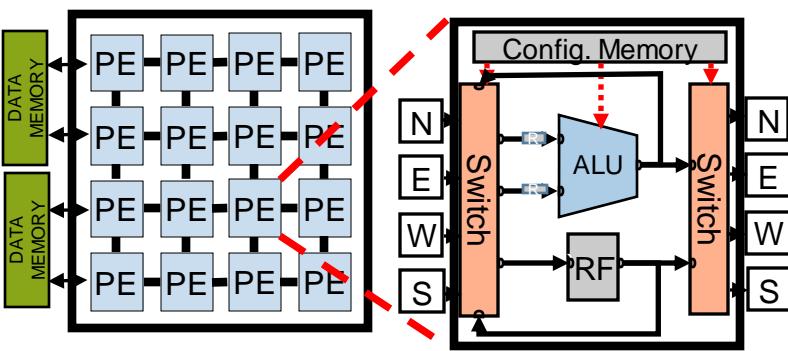
Overview

```

1 #define SIZE 20
2 int A[SIZE], B[SIZE],
3
4 __attribute__((noinline))
5 void array_add(){
6     for (int i=0;i<SIZE;
7         #ifdef CGRA_COMPILER
8             please_map_me();
9         #endif
10        C[i] = A[i]+B[i];
11    }
12 }
13
14 int main(){
15     for (int i=0;i<SIZE; i++){
16         A[i] = i * 2 + 5;
17         B[i] = i * 3;
18         C[i] = 0;
19     }
20     array_add();
21     for (int i=0;i<SIZE; i++) printf("%d\n", C[i]);
22     return 0;
23 }

```

Application source code with annotated kernel

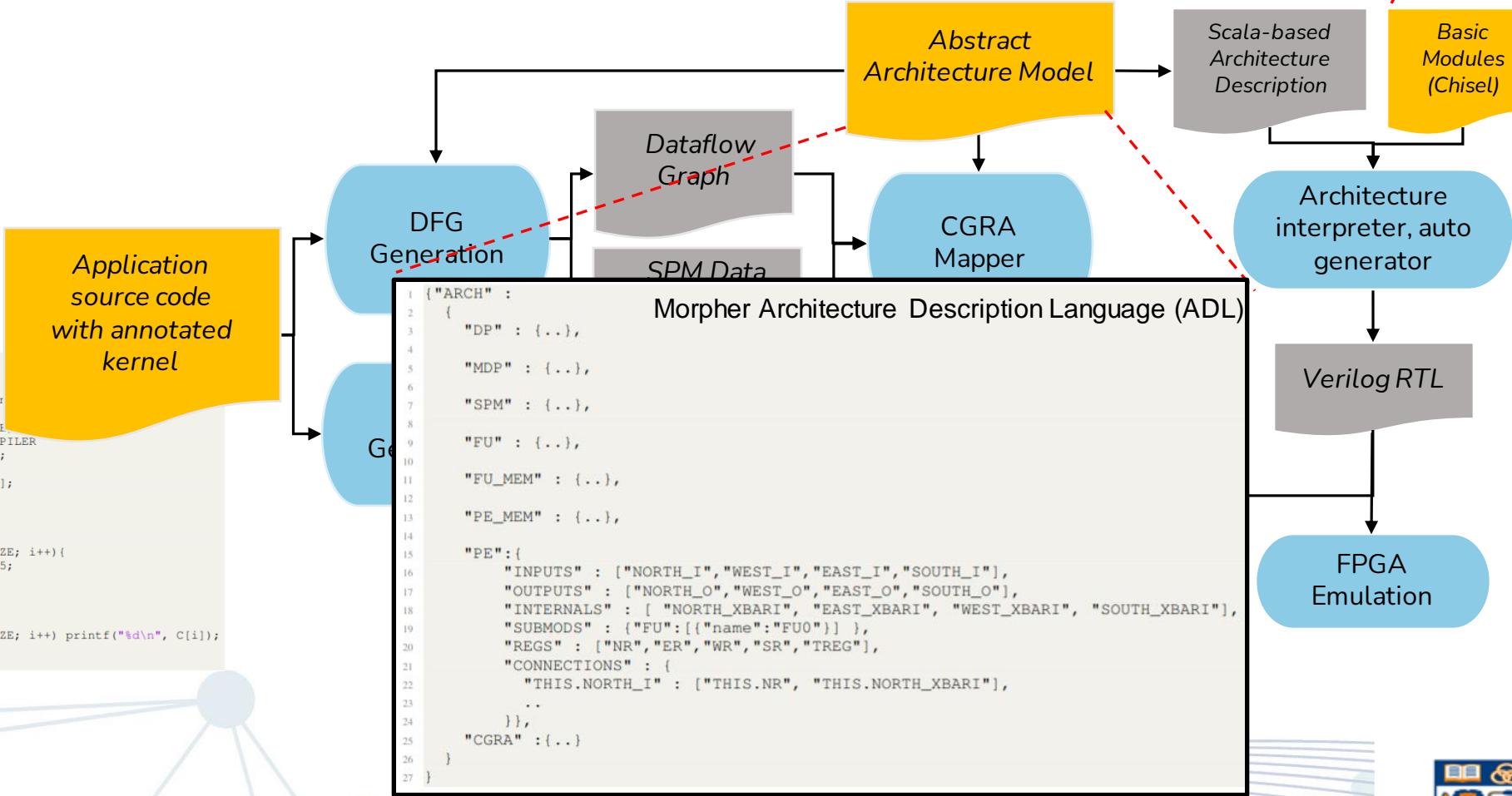


```

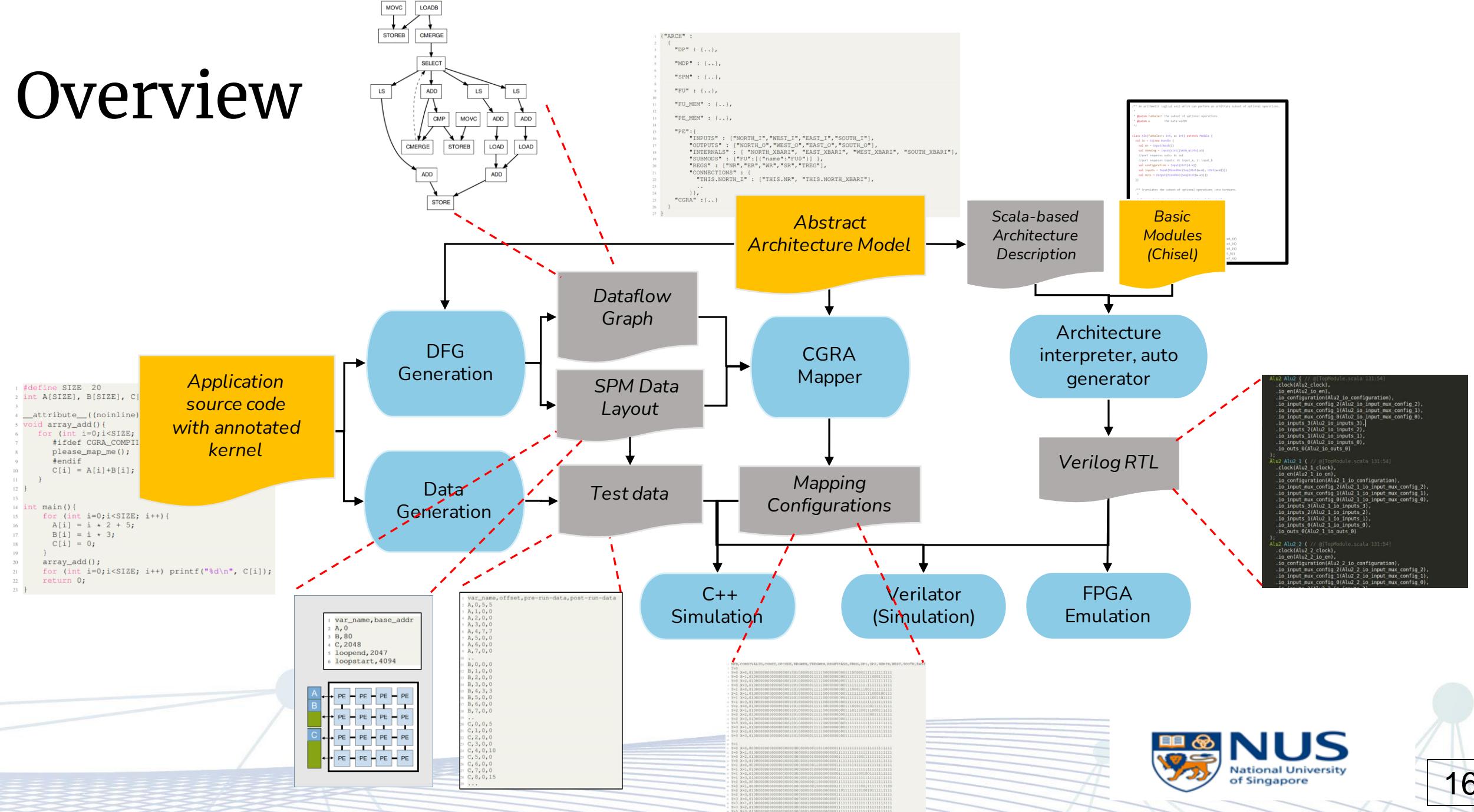
/* An arithmetic unit which can perform an arbitrary subset of optional operations.
 *
 * @param[in] val_lo = 10 new logical
 * @param[in] val_hi = 10 new logical
 * @param[in] val_wiring = InputPort((ConfigMemory.wiring))
 * @param[in] port_sequences inputs: 0: Input_A, 1: Input_B
 * @param[in] port_sequences outputs: 0: Input_A, 1: Input_B
 * @param[in] val_inputs = InputPort((ConfigMemory.inputs))
 * @param[in] val_wiring = InputPort((ConfigMemory.wiring))
 */
class ALU(funcSelect: Int, w: Int) extends module {
    val lo = IO(new Logical)
    val hi = IO(new Logical)
    val wiring = InputPort((ConfigMemory.wiring))
    //port sequences inputs: 0: Input_A, 1: Input_B
    //port sequences outputs: 0: Input_A, 1: Input_B
    val inputs = InputPort((ConfigMemory.inputs))
    val outputs = OutputPort((ConfigMemory.wiring))
    val funcSeq = new Array[InputPort[(ConfigMemory.wiring)]](2)
}

/** Translates the subset of optional operations into hardware.
 *
 * @param[in] val_wiring share the bottom log(wire) bits of "Input.h"
 */
def getFuncSeq(wiret: Int): Seq[InputPort[(ConfigMemory.wiring)]] = {
    val funcSeq = new ArrayBuffer[InputPort[(ConfigMemory.wiring)]]()
    for (i < wiret until wiret.log(wiret)) {
        if ((ConfigMemory.wiring & (1 << i)) > 0) {
            i match {
                case 0 => funcSeq.append(funcSeq.wiring(0) >> (Input.h << i))
                case 1 => funcSeq.append(funcSeq.wiring(0) >> (Input.h << i))
                case 2 => funcSeq.append(funcSeq.wiring(0) >> (Input.h << i))
                case 3 => funcSeq.append(funcSeq.wiring(0) >> (Input.h << i))
                case 4 => funcSeq.append(funcSeq.wiring(0) >> (Input.h << i))
                case 5 => funcSeq.append(funcSeq.wiring(0) >> (Input.h << i))
            }
        }
    }
    funcSeq
}

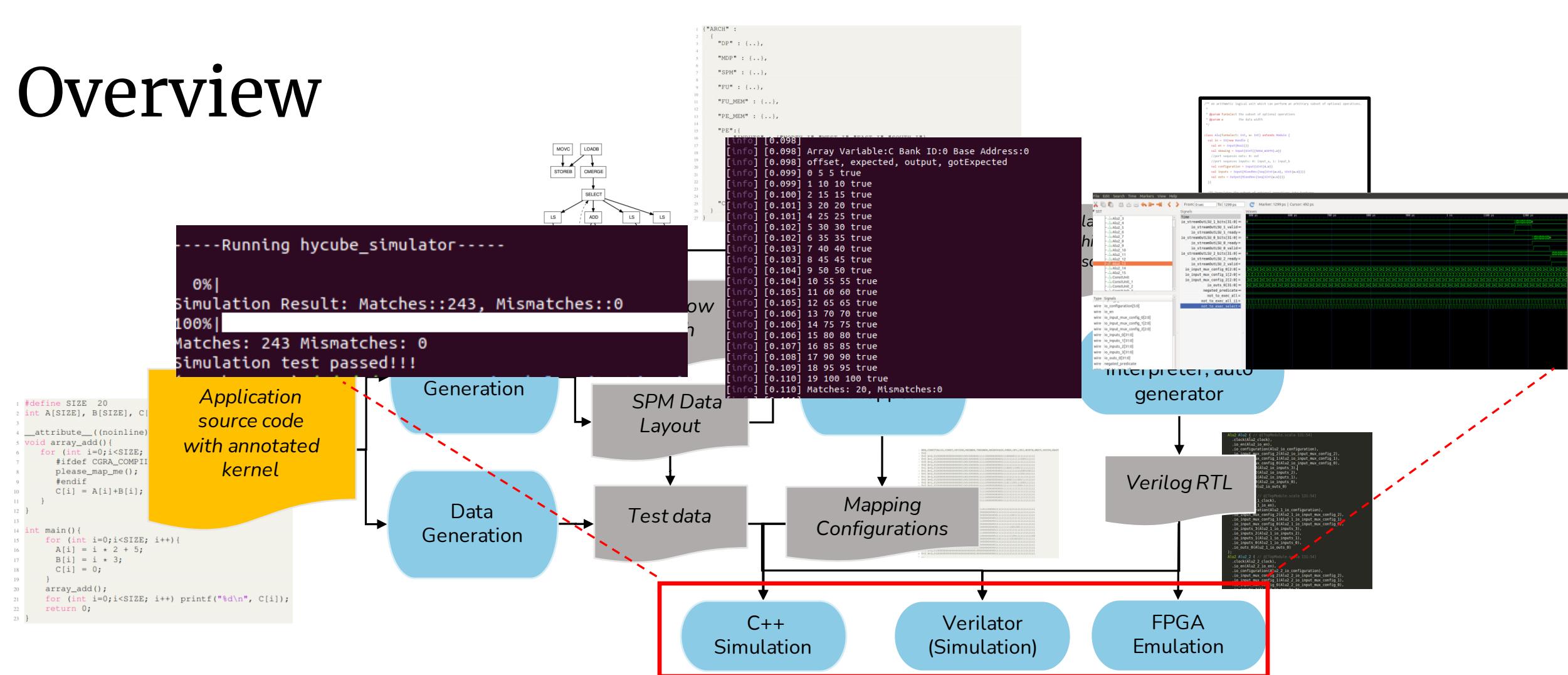
```



Overview

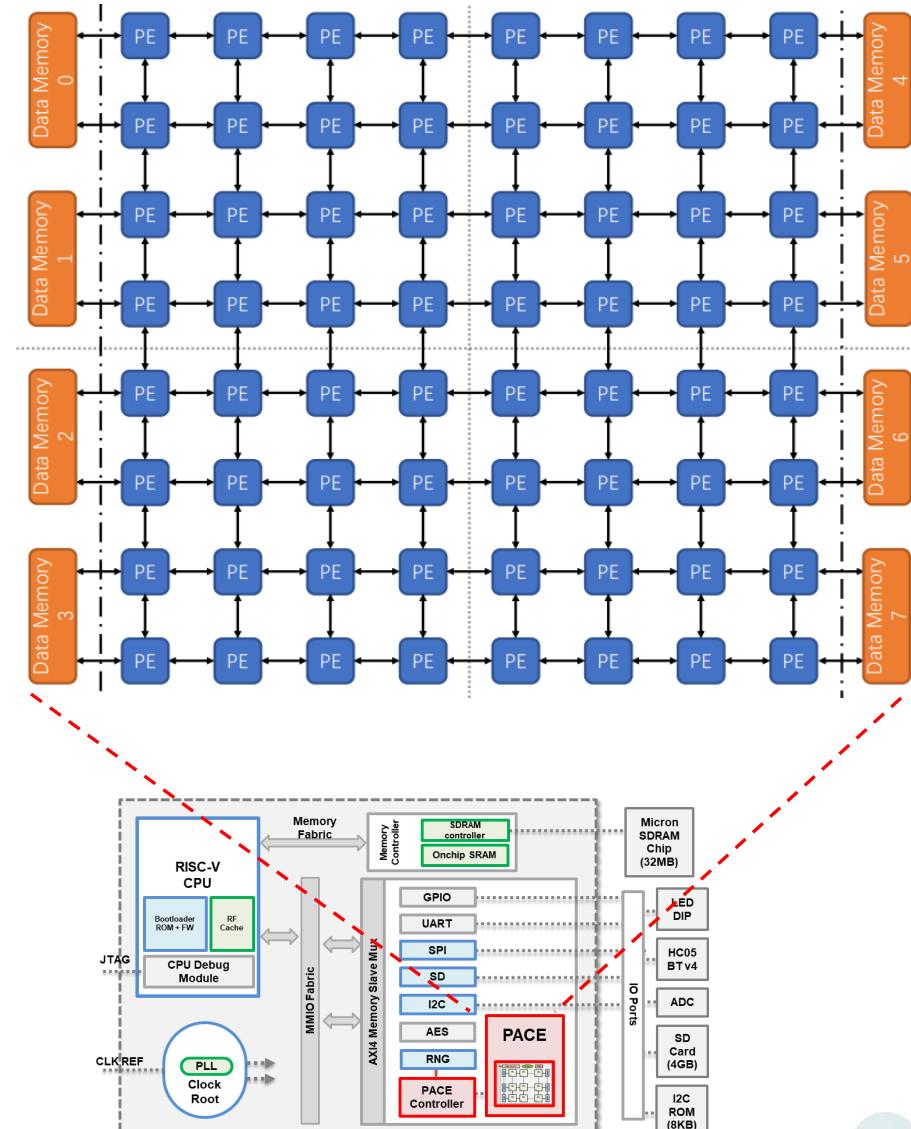


Overview



Experimental Study

- Target CGRA Design: 8x8 PE array with 8 data memories on boundary PEs
 - Logically divided into four clusters:
 - Each with a 4x4 PE array and two 8kB memory banks
- Accelerating ML Workloads:
 - Focus on GEMM and CONV Kernels
 - Kernel Dimensions & Tile Sizes:
 - GEMM: Dimension of $64 \times 64 \times 64$; Tile Size: $64 \times 16 \times 64$
 - CONV: Dimension of $64 \times 64 \times 64 \times 3 \times 3$; Tile Size: $64 \times 64 \times 1 \times 3 \times 3$

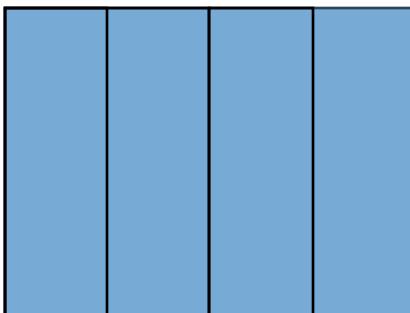


GEMM Mapping

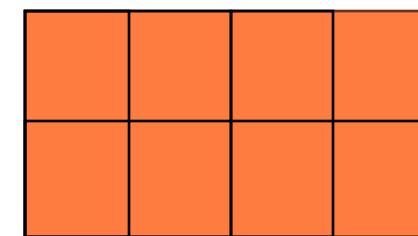
GEMM Computation



*



=



W

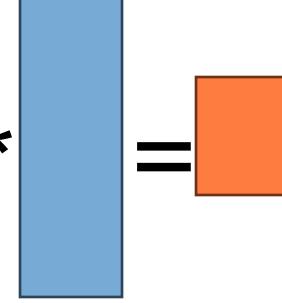
I

O

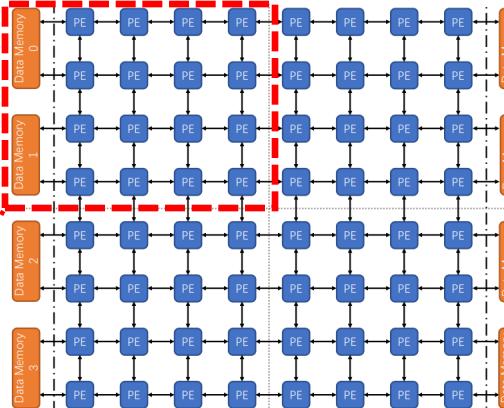
Single CGRA Cluster



*

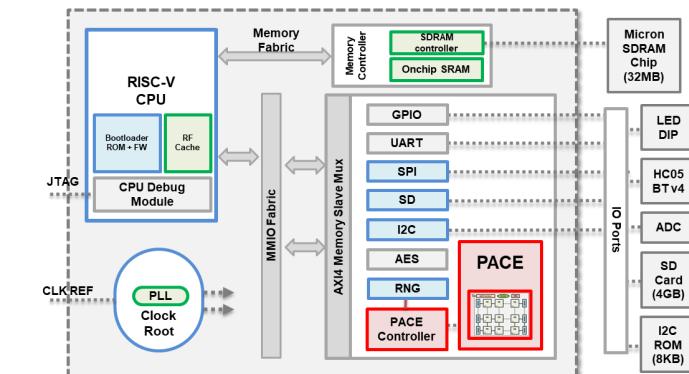


=

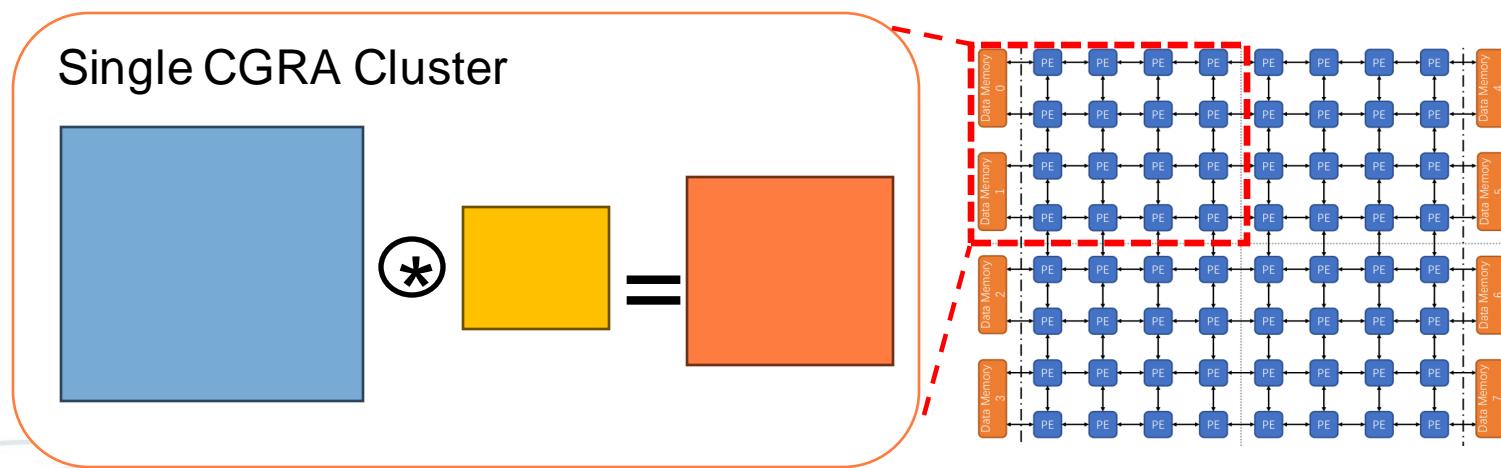
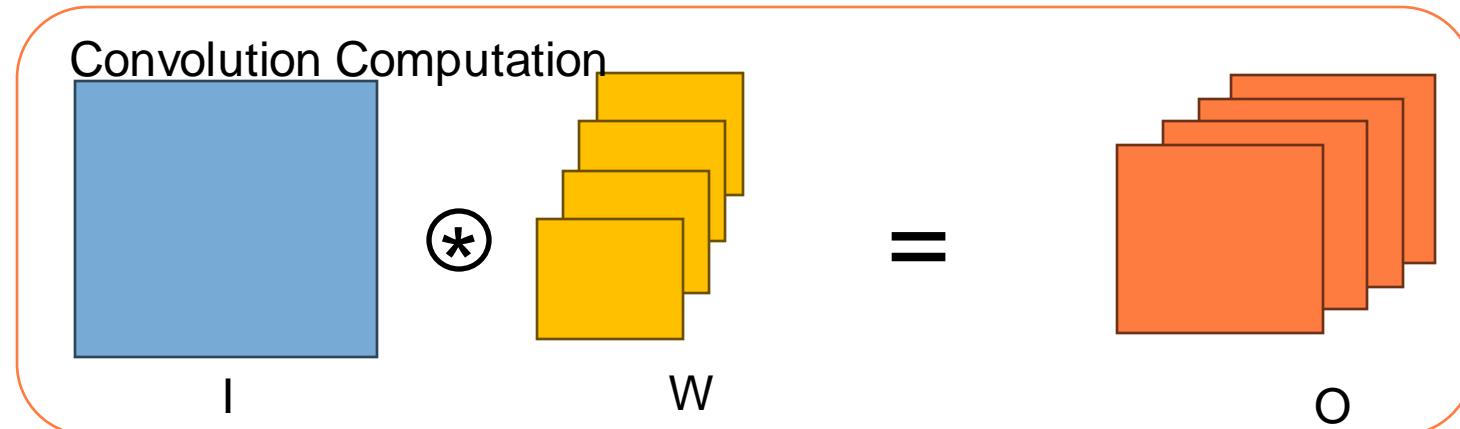


```

1 // Sequential loop: from off-chip to on-chip
2 for m in range(M/ (TI*X)) :
3     for n in range(N/ (TJ*Y)) :
4         for k in range(K/ (TK)) :
5             // Parallel loop: CGRA clusters
6             for x in range(X) :
7                 for y in range(Y) :
8                     // Single CGRA level
9                     for i in range(TI) :
10                    for j in range(TJ) :
11                        for k in range(TK) : //map this
12                            O[][] += W[][] * I[][];
    
```



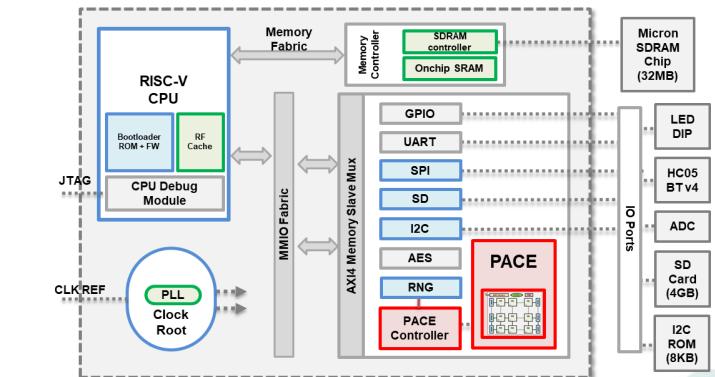
Convolution Mapping



```

1 //Sequential loop: from off-chip to on-chip
2 for i.0 in range (O1/ X*T01):
3   for j.0 in range (O2/ Y*T02):
4     for c.0 in range(Co/ TCo):
5       // Parallel loop: CGRA clusters
6       for x in range(X):
7         for y in range(Y):
8           // Single CGRA level
9           for i in range(T01):
10             for j in range(T02):
11               for c in range(TCo):
12                 temp = 0;
13                 for k1 in range(K):
14                   for k2 in range(K): // map this:
15                     temp += I[] * W[];
16               O[] = temp;

```



Micro kernel mapping on GEMM

```
1 for (i=0; i<TI; i++)  
2 for (j=0; j<TJ; j++)  
3 for (k=0; k<TK; k++) : //map this  
4 O[i][j] += W[i][k]* I[k][j];
```

Unrolling

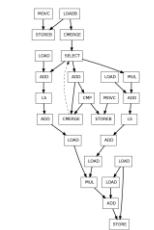
```
1 for (i=0; i<TI; i++)  
2 for (j=0; j<TJ; j++)  
3 for (k=0; k<TK; k=k+4) : //map this  
4 O[i][j] += W[i][k]* I[k][j]+ W[i][k+1]* I[k+1][j]  
5 W[i][k+2]* I[k+2][j]+W[i][k+3]* I[k+3][j];
```

Coalescing

```
1 for (n=0; i=0; j=0; k=0; n<TI*TJ*TK; n++) { : //map this  
2 O[i][j] += W[i][k]*I[k][j]+W[i][k+1]*I[k+1][j]  
3 W[i][k+2]*I[k+2][j]+W[i][k+3]*I[k+3][j]; k = k + 4;  
4 if(k+1 >= TK) {k=0; ++j;}  
5 if(j == TJ) {j=0; ++i;}}
```

DFG Nodes in
the inner Loop

26



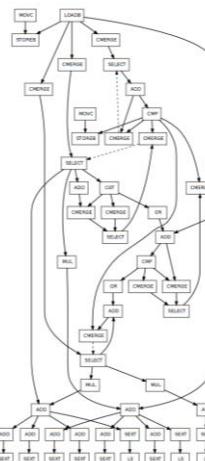
CGRA PE
Utilization

40.6%

Inner Loop
Invocations

4096 (TI x TJ)

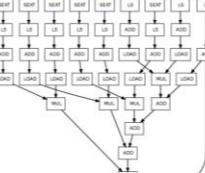
58



60.1%

4096 (TI x TJ)

79



61.7%

1

Micro kernel mapping on Convolution

```
1 for (i=0; i<TO1; i++)  
2 for (j=0; j<TO2; j++)  
3 for (c=0; c<TCo; c++) {  
4     temp = 0;  
5     for(k1 = 0; k1<K; k1++)  
6         for(k2 = 0; k2<K; k2++) { : //map this  
7             temp += I[] * W[]  
8         }  
9     O[] = temp; }
```

Unrolling

```
1 for (i=0; i<TO1; i++)  
2 for (j=0; j<TO2; j++) { //map this  
3     O[] = I[] * W[] + I[] * W[] + I[] * W[]  
4     + I[] * W[] + I[] * W[] + I[] * W[]  
5     + I[] * W[] + I[] * W[] + I[] * W[]; }
```

Coalescing

```
1 for (int ijk=0; ijk<TCo*TO1*TO2; ijk++){ : //map this  
2     O[] = I[] * W[] + I[] * W[] + I[] * W[]  
3     + I[] * W[] + I[] * W[] + I[] * W[]  
4     + I[] * W[] + I[] * W[] + I[] * W[];  
5     j = j + 1;  
6     if(j+1 > TO2) {j=0; ++i;}  
7     if(i == TO1) {i=0; ++c;}}
```

DFG Nodes in
the inner Loop

27

100

153

CGRA PE
Utilization

42.2%

52%

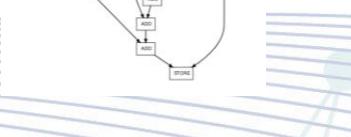
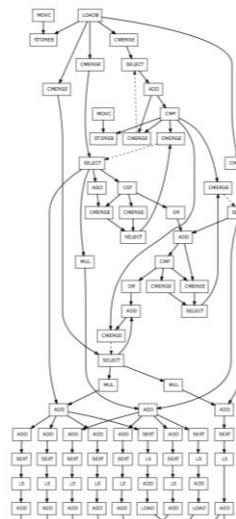
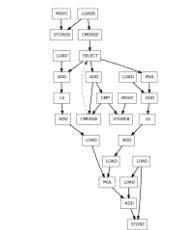
86.9%

Inner Loop
Invocations

64x64x64x3
(TO1xTO2xTcoxK)

4096 (TO1 x TO2)

1



Performance comparison

Kernel	Nodes	II (MII)	Utilization	Compute time (ms)*	Data transfer time (ms)*	Total execution time (ms)*	Speedup
GEMM	26	4 (4)	40.63%	0.56	2.13	2.69	1×
GEMM-U	58	6 (4)	60.42%	0.25	2.13	2.38	1.1×
GEMM-U-C	79	8 (8)	61.72%	0.27	0.49	0.76	3.5×
CONV	27	4 (4)	42.19%	8.32	306.38	314.7	1×
CONV-U-C-1	100	12 (7)	52.08%	1.53	12.75	14.28	22×
CONV-U-C-2	153	11 (10)	86.93%	1.26	11.19	12.45	25.2×

* At 100 MHz CGRA frequency, 50 MBps host-to-CGRA data rate (GPIO)

Conclusion

- Morpher CGRA compilation and simulation framework
 - Flexible to model modern CGRA architectures
 - Map complex workloads with a higher mapping quality at a shorter compilation time
 - Automatic RTL generation & cycle-accurate simulation to validate the compilation results
- Fully open-source
 - Modular codebase
 - Easy to modify
- Open source repository:
 - <https://github.com/ecolab-nus/morpher-v2>





Thank you!

