

# A Simulator and Compiler Framework for Agile Hardware-Software Co-design Evaluation and Exploration

*ICCAD 2020 Opensource Tools and Platforms for Agile Development of Specialized Architectures*

**Tyler Sorensen**  
**UC Santa Cruz**

Aninda Manocha, Esin Tureci, Margaret Martonosi  
Princeton University

Juan L. Aragón  
Universidad de  
Murcia



# Speaker Bio



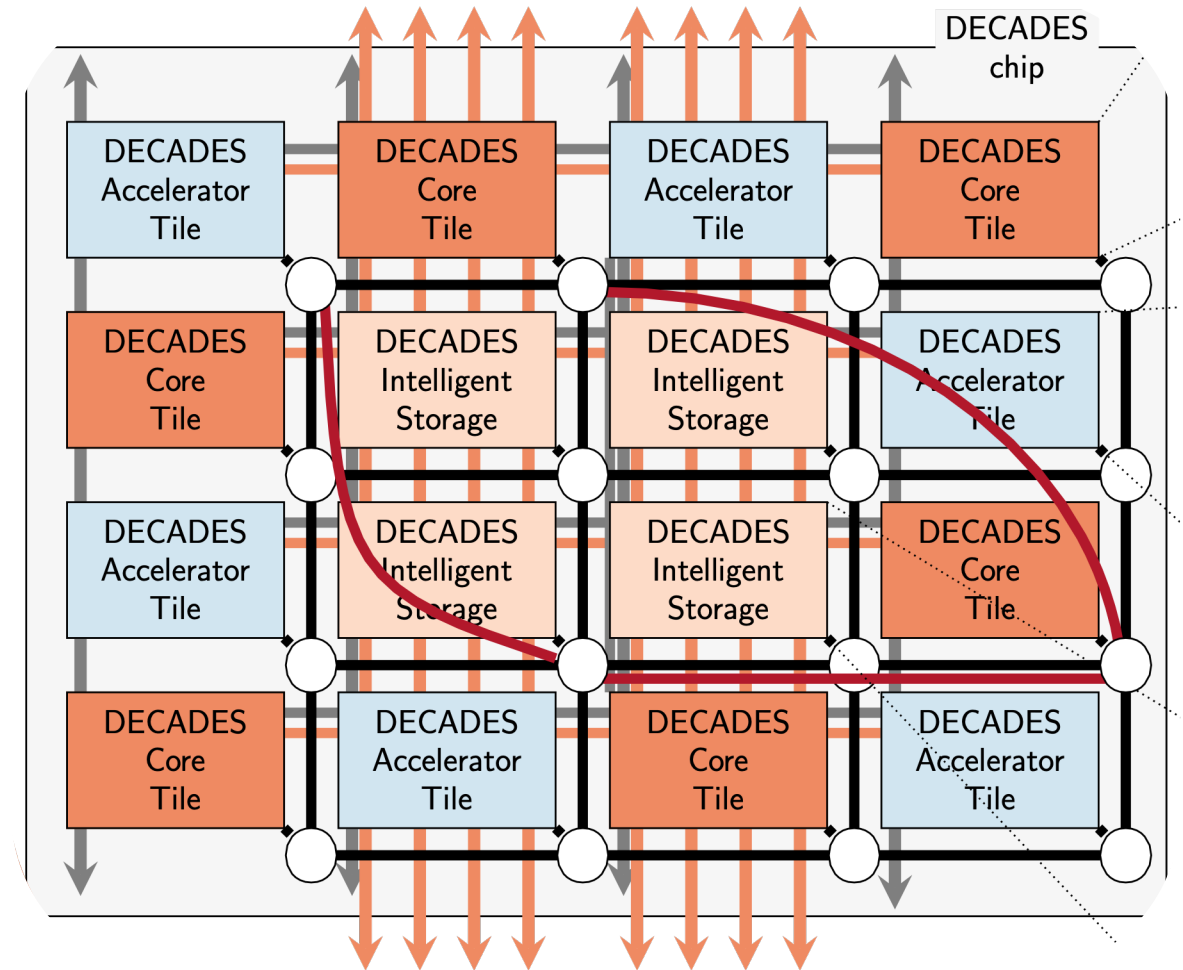
- Tyler Sorensen
- **Assistant Professor** at UC Santa Cruz Since Summer 2020
- Previously I was a **Post doc** at Princeton with Margaret Martonosi when this work was done
- Background is in Programming Languages (GPU programming models), but I was interested in peeking under the hood 😊

[https://twitter.com/Tyler UCSC](https://twitter.com/Tyler_UCSC)

<https://users.soe.ucsc.edu/~tsorensen/>

# The DECADES Project

- Part of DARPA's SDH project
- Principle Investigators:
  - David Wentzlaff (Princeton)
  - Luca P. Carloni (Columbia)
  - Margaret Martonosi



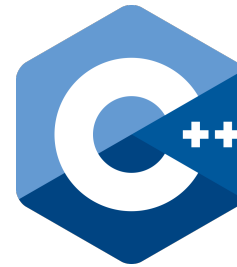
- Developing a new tiled, heterogeneous architecture with data-supply and accelerator innovations (tape out planned in near future!)

# Building a Chip is a Big Project...

**The focus of Professor Martonosi's group was:**

- Programming language support and innovations

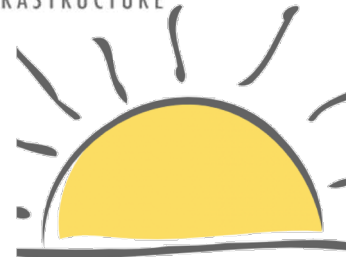
- Support for popular programming languages



- Modular and extensible design



- Early stage design space exploration



# Our Contributions: A Compiler and Simulator

- Compiler: **DEC++**

- Builds on top of LLVM, Clang frameworks.
- Kernel-centric parallel programming model (main support for C++)
- Flexible frontends, backends, and transformations

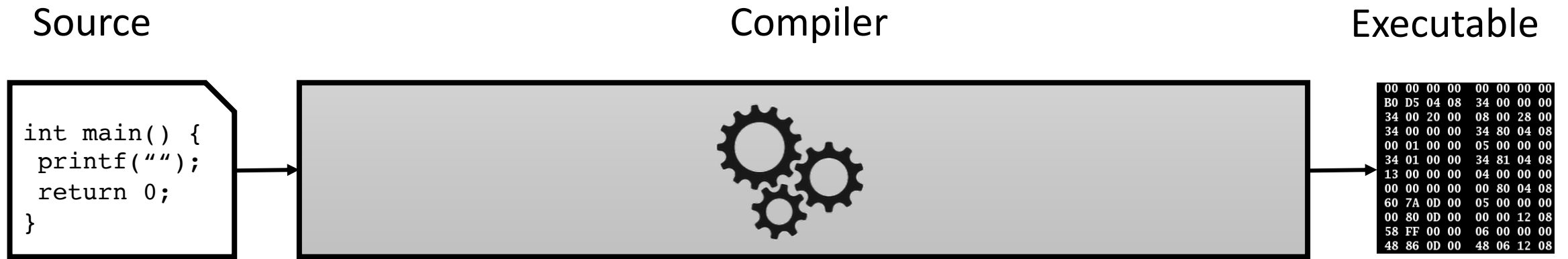
<https://github.com/PrincetonUniversity/DecadesCompiler>

- Simulator: **MosaicSim**

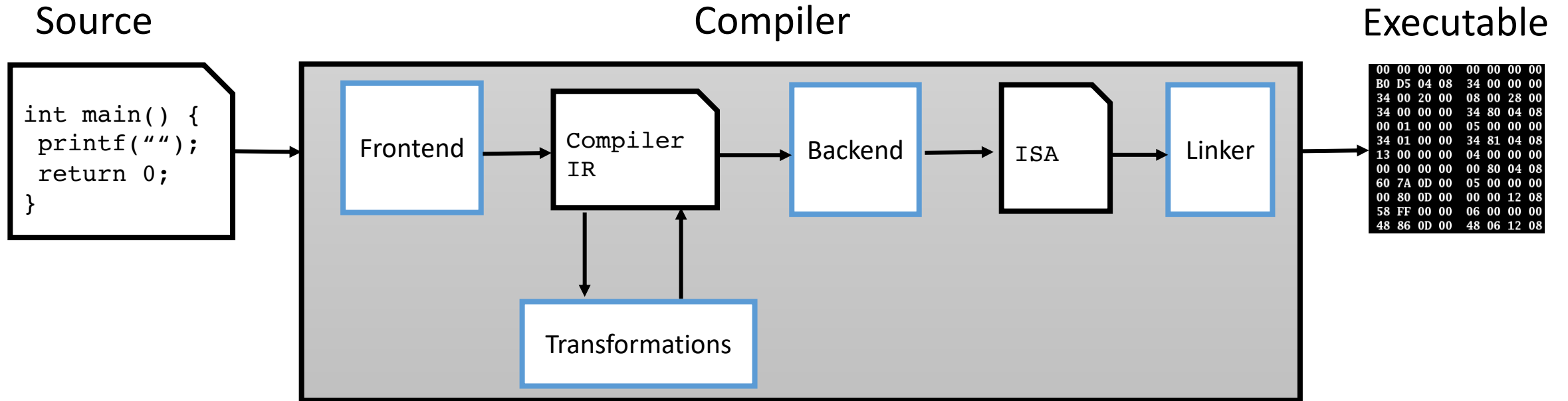
- Early-stage performance estimates (cycle-driven LLVM IR simulation)
- Tile model support heterogeneous core models (both CPUs and accelerators)
- *Best Paper Nomination at ISPASS 2020!*  
MosaicSim: A Lightweight, Modular Simulator for Heterogeneous Systems

<https://github.com/PrincetonUniversity/MosaicSim>

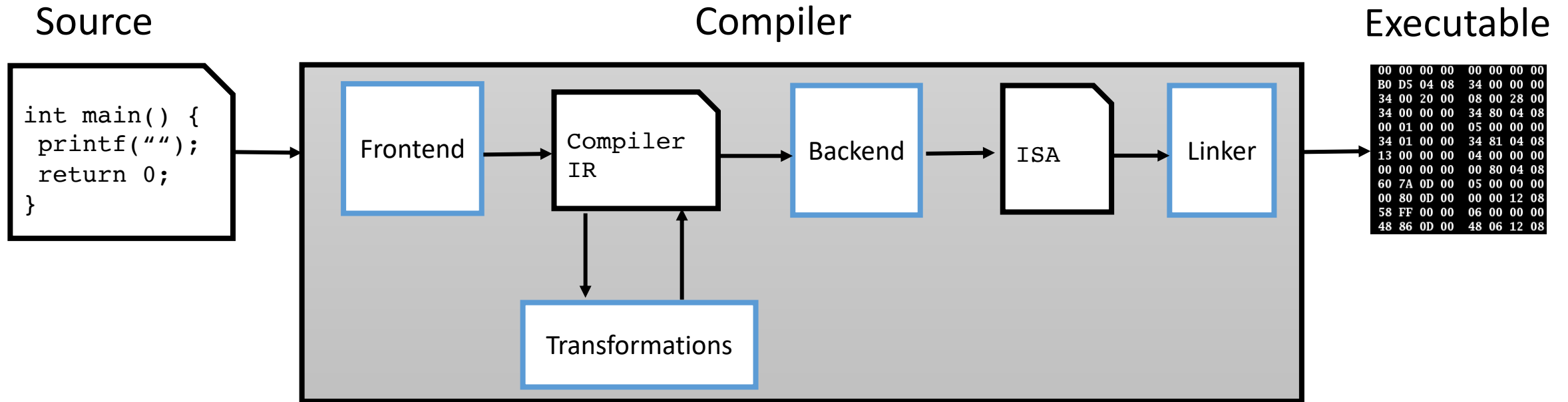
# The DEC++ Compiler



# The DEC++ Compiler



# The DEC++ Compiler



*Too complex to develop everything!  
Instead we plug into the LLVM toolflow*

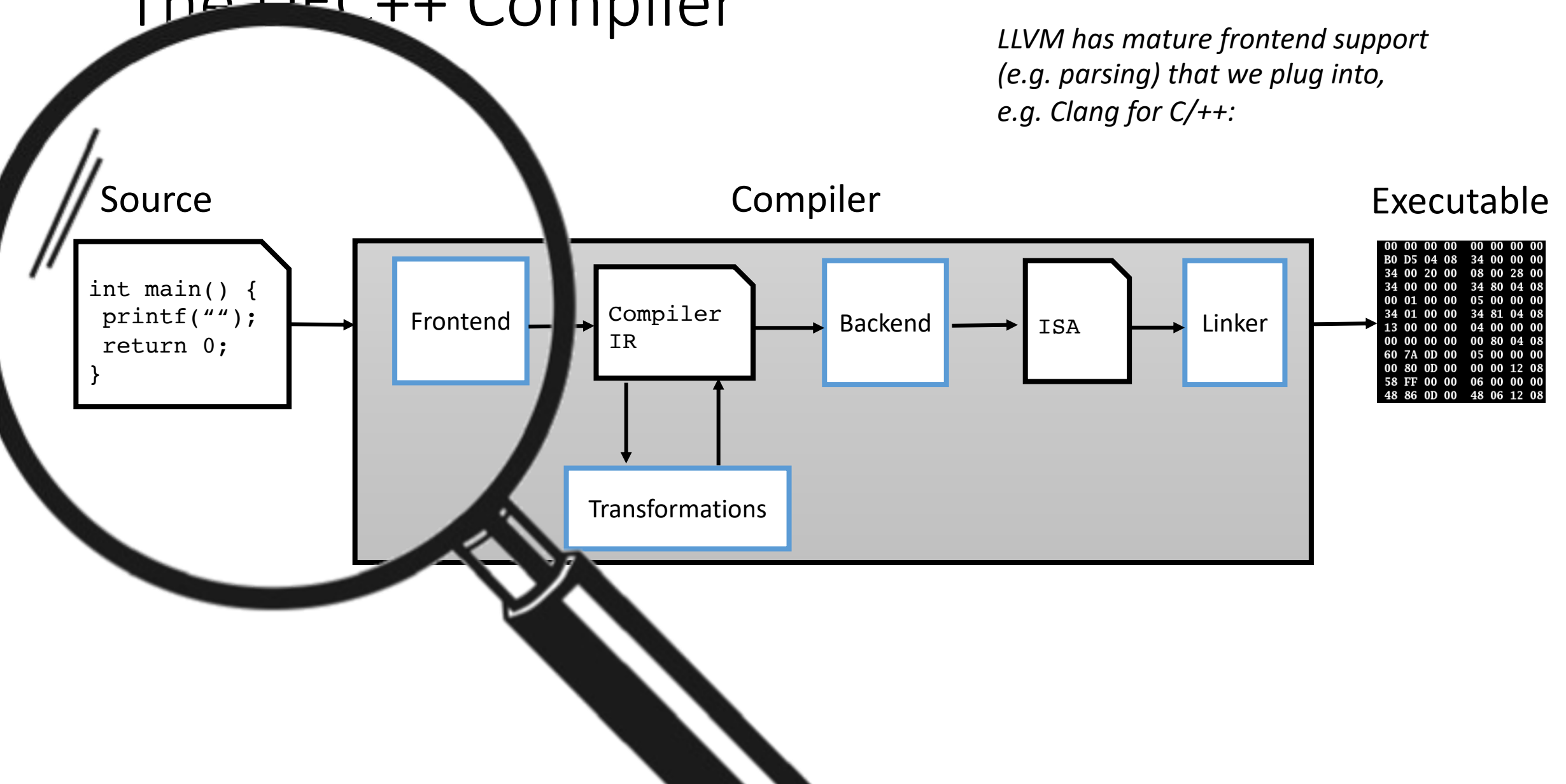


*We require no LLVM source code changes and  
simply link to public APIs or use tools directly*



# The DEC++ Compiler

*LLVM has mature frontend support  
(e.g. parsing) that we plug into,  
e.g. Clang for C/++:*



# DEC++ Front End: Programming Model

Kernel function has two required parameters

Kernel function must be identified

```
1 void _kernel_(float *a, float *b ... int thread_id, int num_threads) {  
2     for (int i = thread_id; i < SIZE; i+=num_threads) {  
3         a[i] += b[i];  
4     }  
5 }
```

Program is written in a thread-agnostic SPMD way

# DEC++ Front End: Implementation

Front end implementation must intercept kernel function call and run in SPMD parallel execution

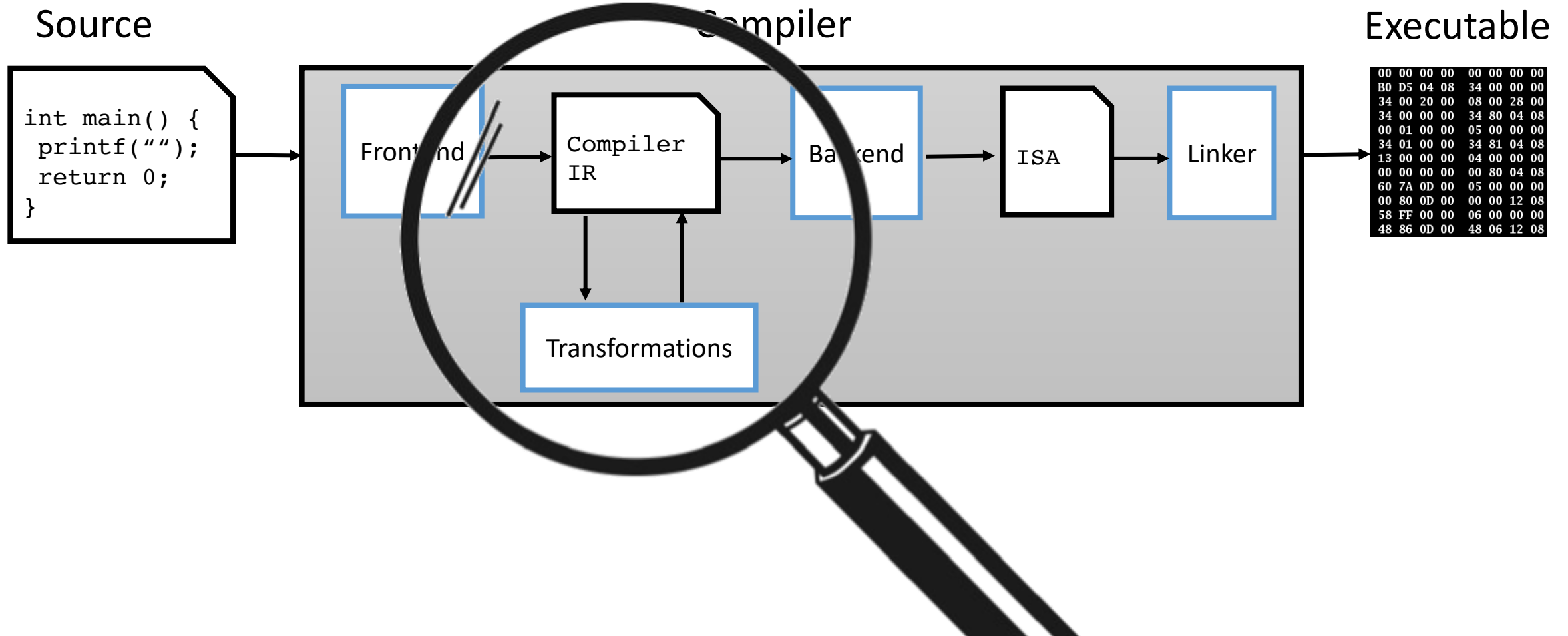
```
1 // Setting up data
2 _kernel_(a, b, ..., DEC_TID, DEC_NUMTHREADS);
3 // Cleaning up
```



*Implemented with Clang Visitor pass*

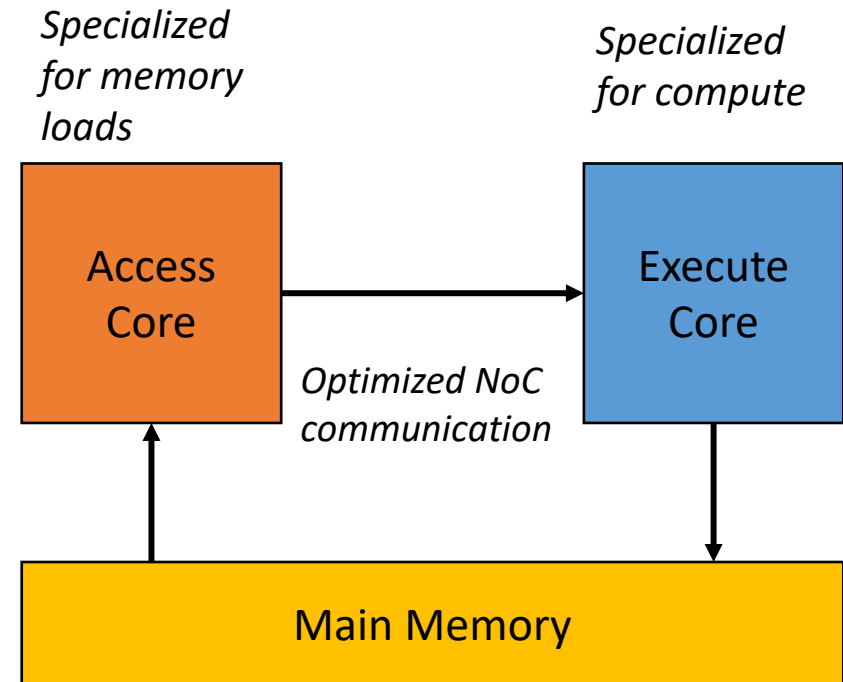
```
1 #pragma omp parallel for
2 for (int i = 0; i < NUM_THREADS; i++) {
3     _kernel_(a, b, ... i, NUM_THREADS);
4 }
```

# The DEC++ Compiler



# DEC++ Transformations

- Compiler passes over the LLVM AST performing re-writes and analysis.
- Lots of opportunity for innovation
- Example: Decoupled Access Execute (DAE)



# DEC++ Transformations: DAE Example

```
1 void _kernel_(float *a, float *b ... int thread_id, int num_threads) {  
2     for (int i = thread_id; i < SIZE; i+=num_threads) {  
3         a[i] += b[i];  
4     }  
5 }
```

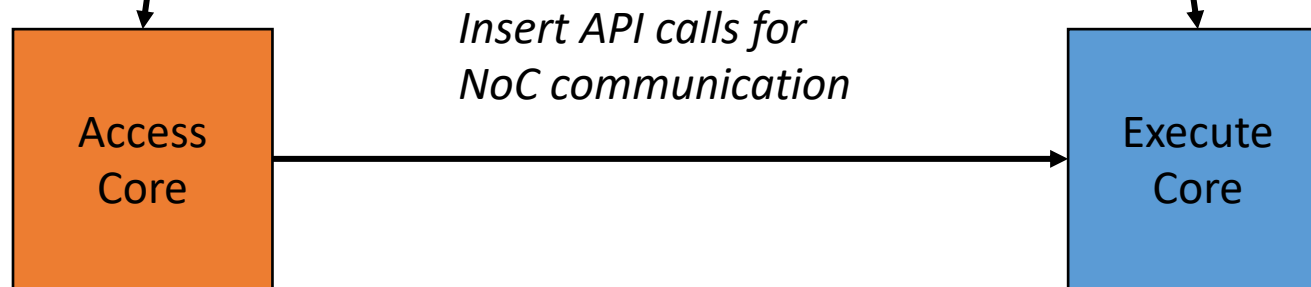
↓ *In pseudo LLVM-IR*

```
1 void _kernel_(float *a, float *b ... int thread_id, int num_threads) {  
2     for (int i = thread_id; i < SIZE; i+=num_threads) {  
3         tmp_var0 = load(a[i]);  
4         tmp_var1 = load(b[i]);  
5         tmp_var2 = tmp_var0 + tmp_var1;  
6         store(a[i], tmp_var2);  
7     }  
8 }
```

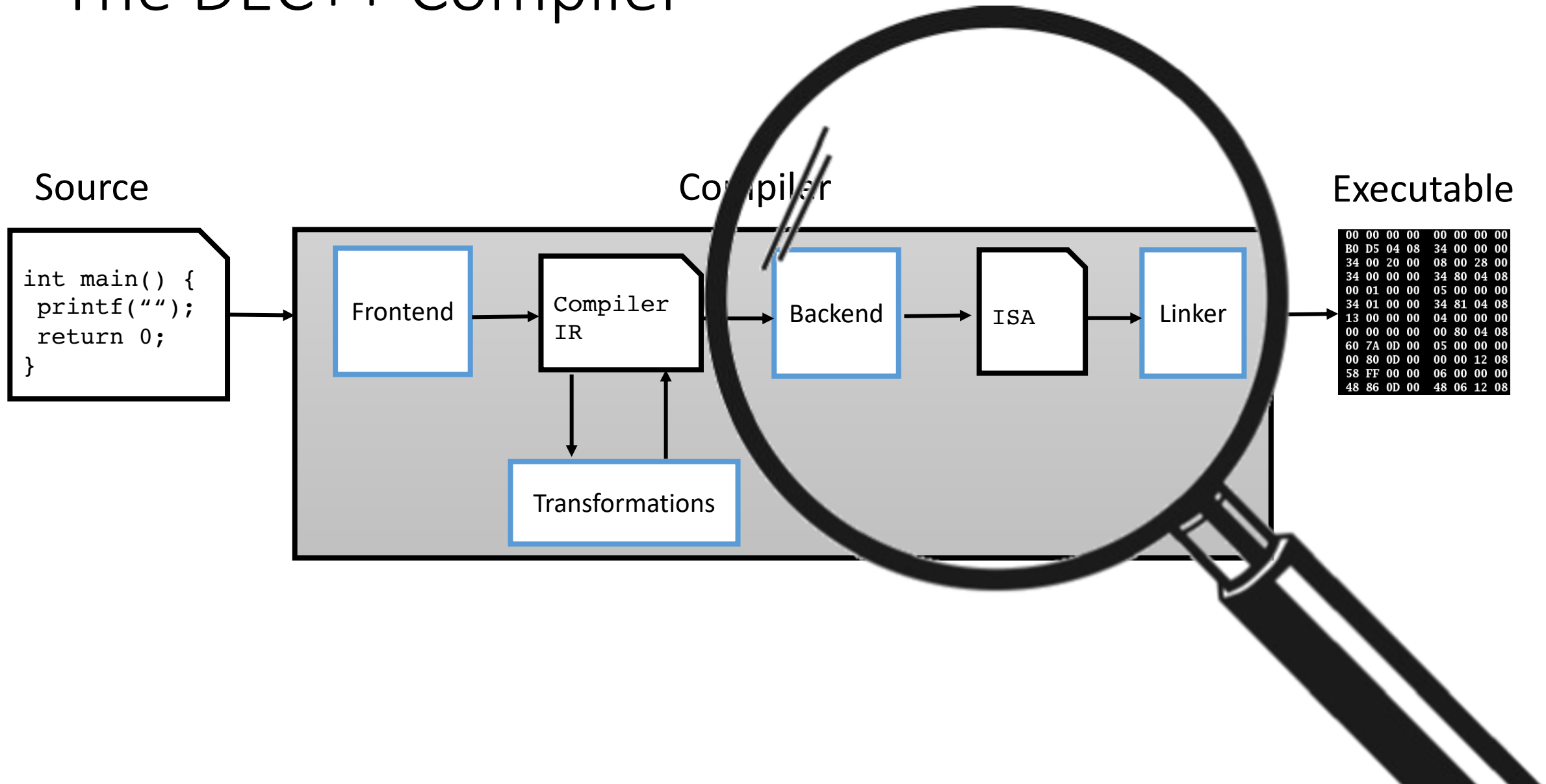
# DEC++ Transformations: DAE Example

```
1 void _kernel_(float *a, float *b ... int thread_id, int num_threads) {  
2   for (int i = thread_id; i < SIZE; i+=num_threads) {  
3     tmp_var0 = load(a[i]);  
4     tmp_var1 = load(b[i]);  
5     tmp_var2 = tmp_var0 + tmp_var1;  
6     store(a[i], tmp_var2);  
7   }  
8 }
```

*Simple code slicing  
places loads on an  
Access core and  
compute on the  
execute core*



# The DEC++ Compiler

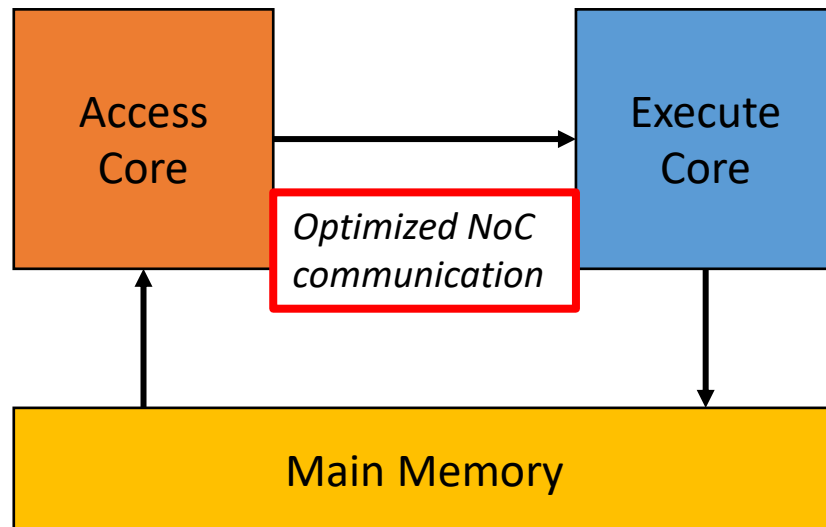




# DEC++ Backend and Linking

LLVM has backends for many architectures:

- X86: ideal for developing and debugging
- RISC-V: The ISA for the DECADES architecture



But how to deal with new architecture features?

# DEC++ Backend and Linking

- We require architecture features to be implemented behind an API with a software emulation implementation:

*Provides:  
Portable Execution  
Documentation  
Specification*

Loads (load 32 or 64 bits from the stated address and put the data into the queue):

```
void dec_load32_async(uint64_t qid, uint32_t *addr)  
void dec_load64_async(uint64_t qid, uint64_t *addr)
```

Produce and Consume from Producer to Consumer

Loads that are not asynchronous and performed on the Producer are simply loaded regularly in the Producer tile and then enqueued to the "Produce to Consumer". They appear on the Producer tile as follows:

```
void dec_produce32(uint64_t qid, uint32_t data)  
void dec_produce64(uint64_t qid, uint64_t data)
```

The mirror of this interaction is when the Execute consumes data from the queue. These instructions appear on the Execute tile as follows:

```
uint32_t data;
```

# Our Contributions: A Compiler and Simulator

- Compiler: **DEC++**

- Builds on top of LLVM, Clang frameworks.
- Kernel-centric parallel programming model (main support for C++)
- Flexible frontends, backends, and transformations

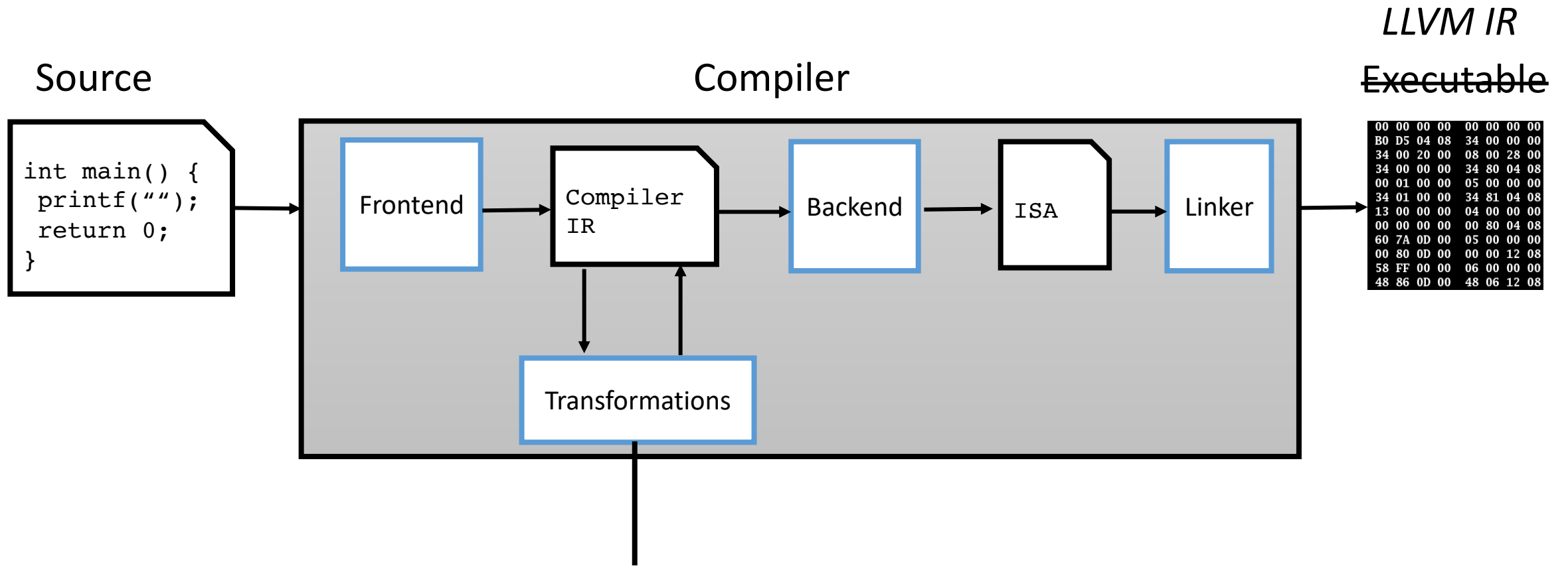
<https://github.com/PrincetonUniversity/DecadesCompiler>

- Simulator: **MosaicSim**

- Early-stage performance estimates (cycle-driven LLVM IR simulation)
- Tile model support heterogeneous core models (both CPUs and accelerators)
- *Best Paper Nomination at ISPASS 2020!*  
MosaicSim: A Lightweight, Modular Simulator for Heterogeneous Systems

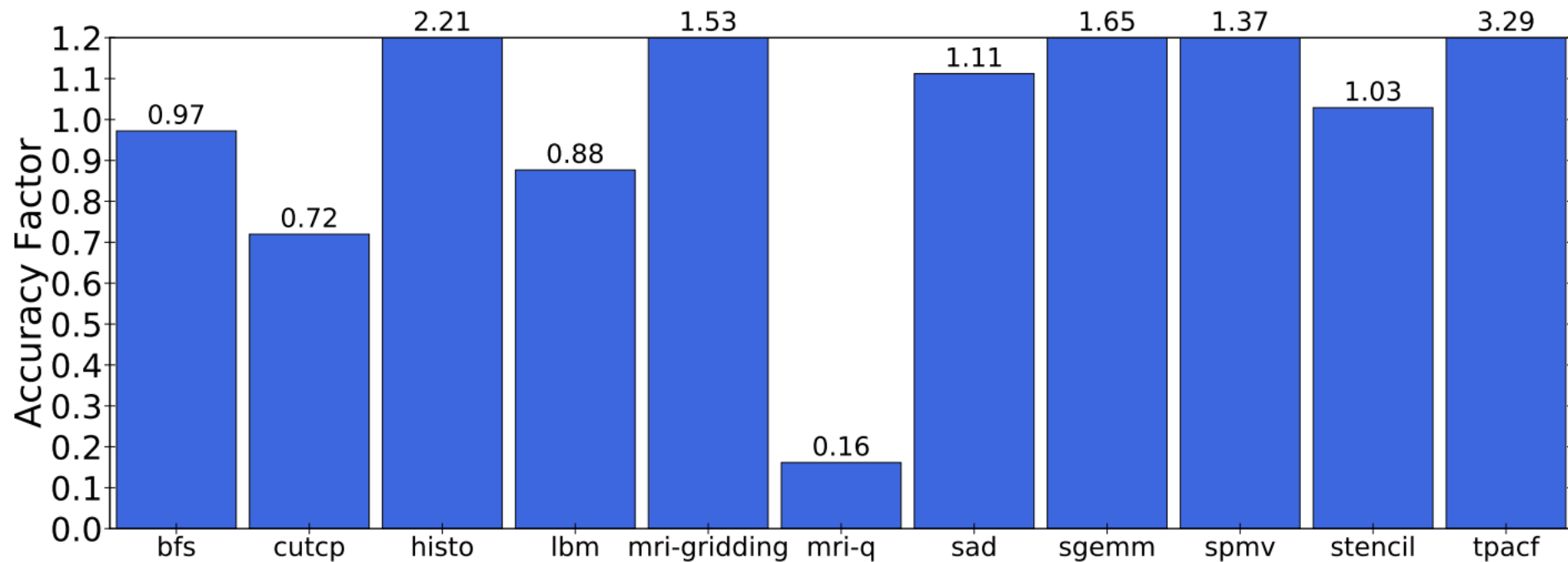
<https://github.com/PrincetonUniversity/MosaicSim>

# How DEC++ interfaces with MosaicSim



*Simulator annotates memory instructions to get a memory trace  
and generates a data-dependency graph*

# MosaicSim: How accurate is simulating LLVM IR?



*Raw cycle counts are pretty inaccurate  
but why?*

The image displays four code editors arranged in a 2x2 grid, illustrating the compilation of a simple C++ function into different instruction sets.

**Top Left Editor (C++ source #1):** Shows the original C++ code for a function named `memory_access` that takes a volatile integer array and returns 0.

```
1 int memory_access(volatile int* array) {
2     array[6] = 6;
3     return 0;
4 }
```

**Top Right Editor (x86-64 clang 10.0.1):** Shows the LLVM IR generated for the C++ code. The function is defined as `@_Z13memory_access@.l` and uses instructions like `alloca`, `store`, `call`, `load`, `getelementptr`, and `ret`.

```
1
2 define dso_local i32 @_Z13memory_access@.l(
3     %0 = alloca i32*, align 8
4     store i32* %0, i32** %2, align 8
5     call void @llvm.dbg.declare(metadata
6     %3 = load i32*, i32** %2, align 8
7     %4 = getelementptr inbounds i32,
8     store volatile i32 6, i32* %4, align 8
9     ret i32 0, !dbg !18
10 }
11
```

**Bottom Left Editor (x86-64 clang 10.0.1):** Shows the X86 assembly generated for the C++ code. The function `memory_access` is defined, and the assembly instructions include `push rbp`, `mov rbp, rsp`, `xor eax, eax`, `mov qword ptr [rbp - 8], rdi`, `mov rcx, qword ptr [rbp - 8]`, `mov dword ptr [rcx + 24], 6`, `pop rbp`, and `ret`.

```
1 memory_access(int volatile*):
2     push    rbp
3     mov     rbp, rsp
4     xor     eax, eax
5     mov     qword ptr [rbp - 8], rdi
6     mov     rcx, qword ptr [rbp - 8]
7     mov     dword ptr [rcx + 24], 6
8     pop     rbp
9     ret
```

**Bottom Right Editor (RISC-V rv64gc clang (trunk)):** Shows the RISC-V assembly generated for the C++ code. The function `memory_access` is defined, and the assembly instructions include `addi sp, sp, -32`, `sd ra, 24(sp)`, `sd s0, 16(sp)`, `addi s0, sp, 32`, `sd a0, -24(s0)`, `ld a0, -24(s0)`, `addi a1, zero, 6`, `sw a1, 24(a0)`, `mv a0, zero`, `ld s0, 16(sp)`, `ld ra, 24(sp)`, `addi sp, sp, 32`, and `ret`.

```
1 memory_access(int volatile*):
2     addi    sp, sp, -32
3     sd      ra, 24(sp)
4     sd      s0, 16(sp)
5     addi    s0, sp, 32
6     sd      a0, -24(s0)
7     ld      a0, -24(s0)
8     addi    a1, zero, 6
9     sw      a1, 24(a0)
10    mv      a0, zero
11    ld      s0, 16(sp)
12    ld      ra, 24(sp)
13    addi    sp, sp, 32
14    ret
```

*instruction mapping  
mismatches:*

- 1 C instruction maps to:*
- 3 LLVM IR instructions*
- 2 X86 instructions*
- 4 RISC-V instructions*

# MosaicSim: Scaling Trends

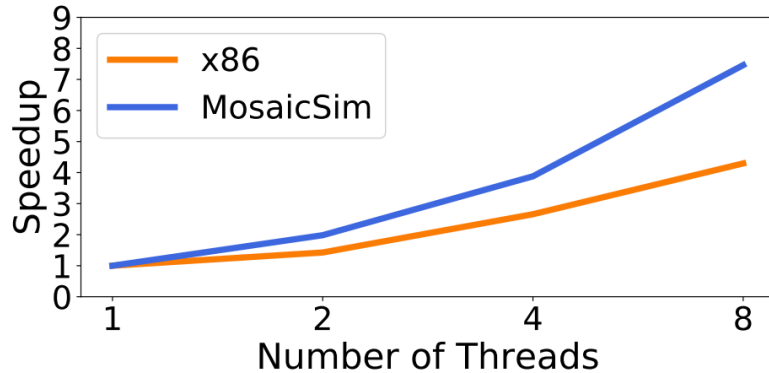


Fig. 7. BFS Scaling Trends

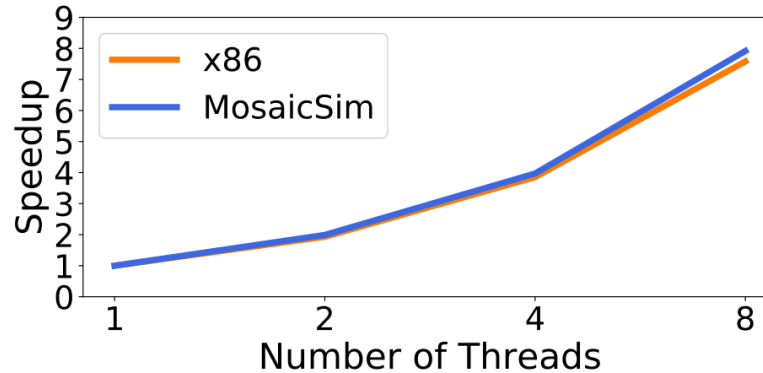


Fig. 8. SGEMM Scaling Trends

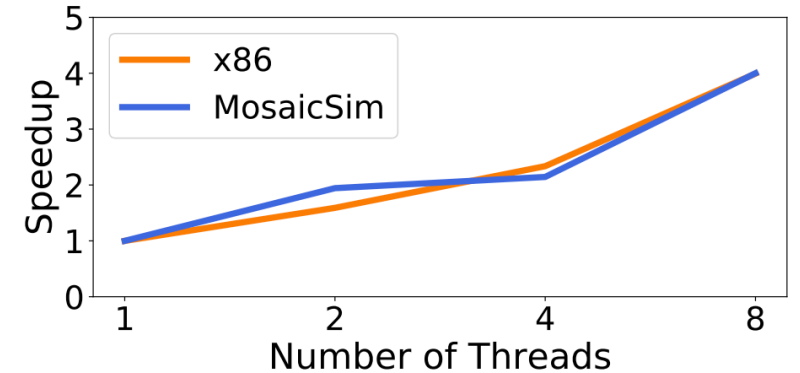


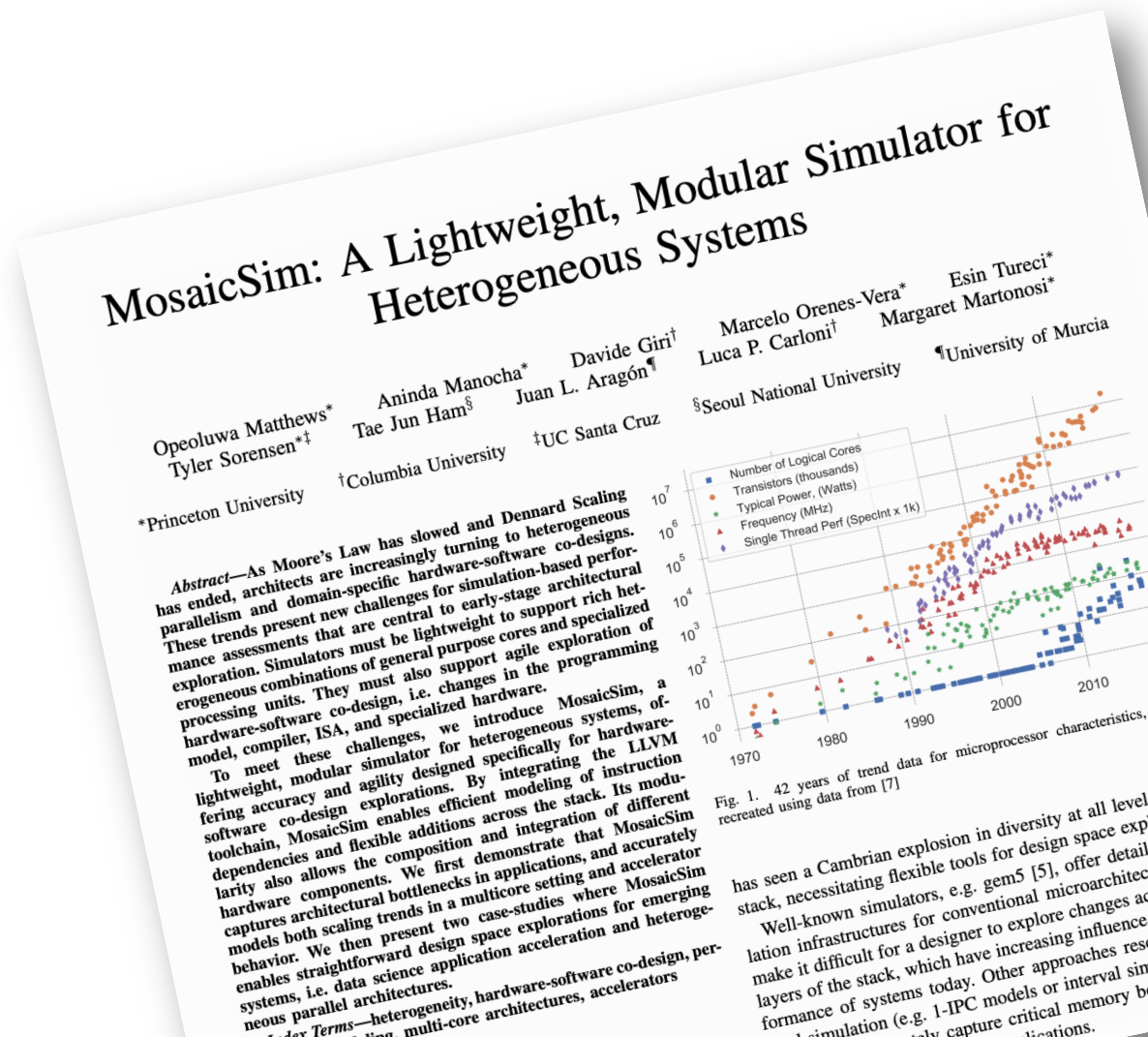
Fig. 9. SPMV Scaling Trends

*MosaicSim accurately captures trends, which is what is important for early-stage performance modeling*

# MosaicSim Extras in ISPASS paper

- modeling ASIC accelerator tiles
- how complex architecture features are efficiently modeled, e.g. RoB, LSQ
- case studies of design space exploration of applications on heterogeneous architectures

***Lead author of MosaicSim is Luwa Matthews (now at Apple)***





# Conclusion

- We present a compiler/simulator framework for hardware-software co-design
- DEC++ is built alongside LLVM, giving it flexibility in frontends and backends.
  - Straightforward to implement innovations at the IR transformation level, e.g. DAE
  - Architecture additions are provided through APIs that support native emulation
- MosaicSim provides early-stage performance estimates. Simulating LLVM is inaccurate at the cycle level, but captures trends and characterizations

*Thanks to the DECADES team and co-authors: Aninda Manocha, Esin Tureci, Marcelo Orenes-Vera, Juan L. Aragón, Margaret Martonosi*

Software:

<https://github.com/PrincetonUniversity/DecadesCompiler>

<https://github.com/PrincetonUniversity/MosaicSim>

Tyler Sorensen

[https://twitter.com/Tyler\\_UCSC](https://twitter.com/Tyler_UCSC)

<https://users.soe.ucsc.edu/~tsorensen/>