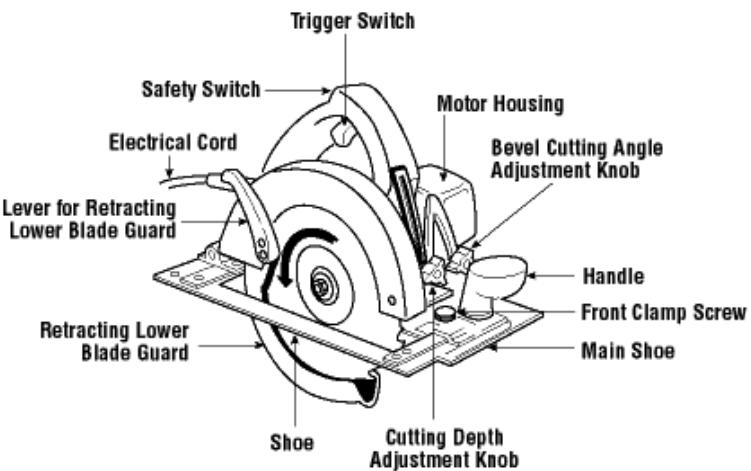
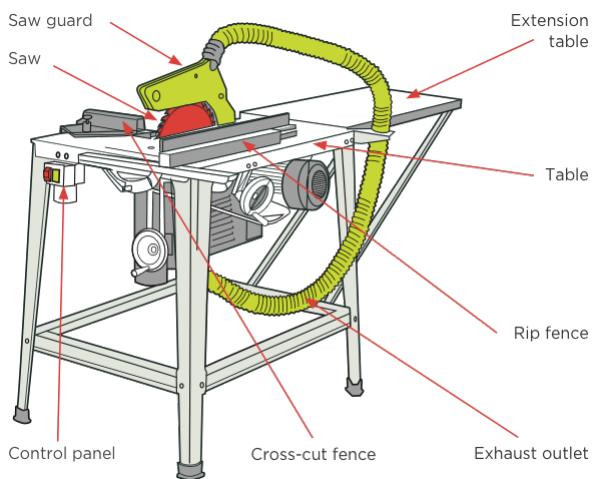


Exo: Let's Build Power Tools for Performance Engineers

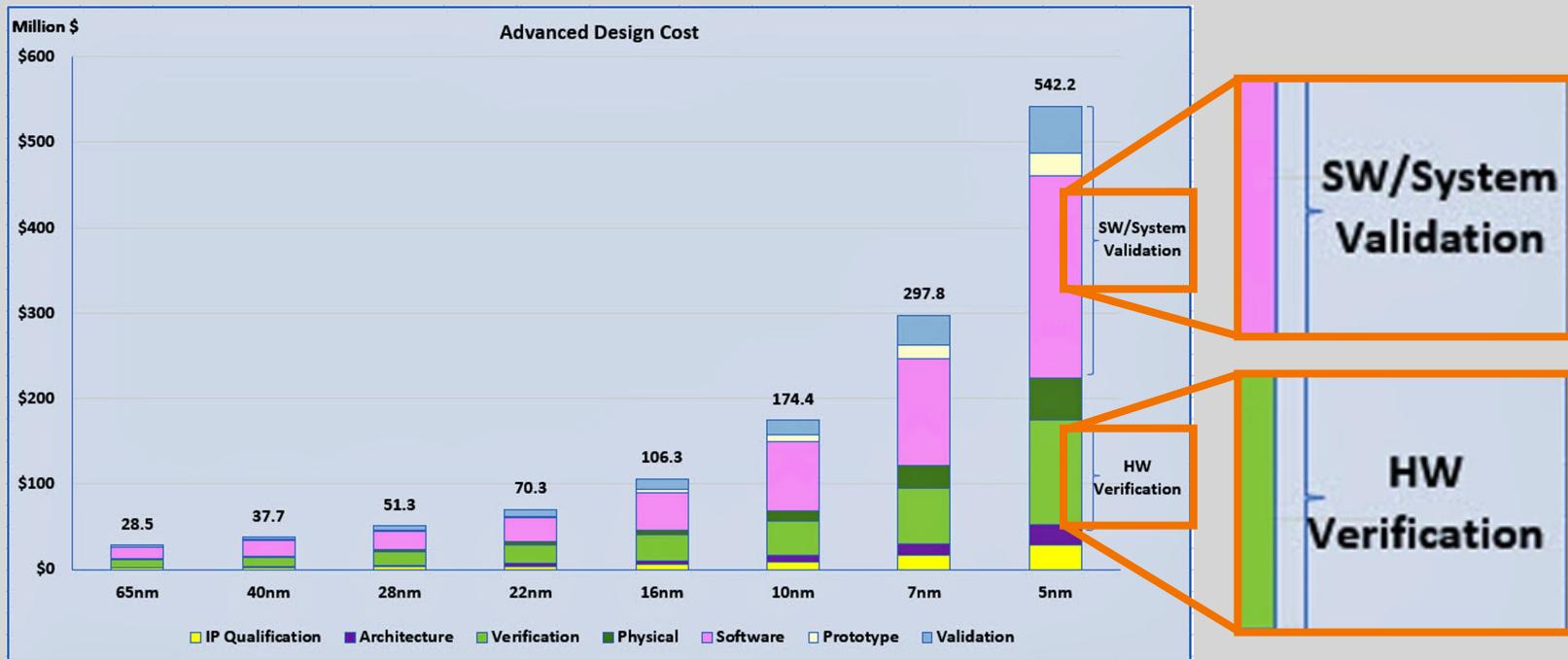




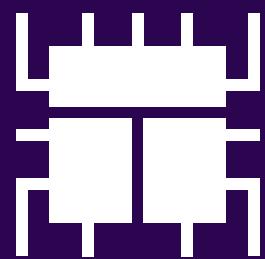




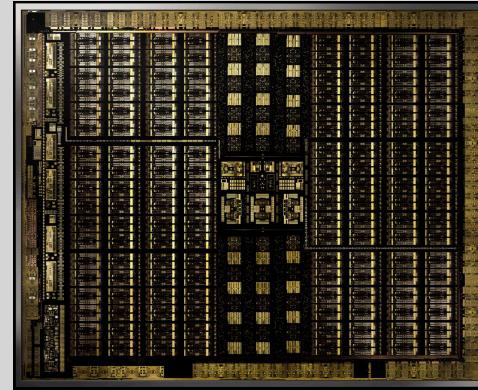
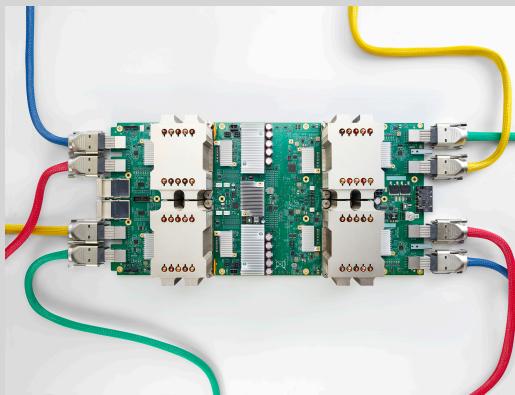
Verification & Validation Dominate HW Development Cost



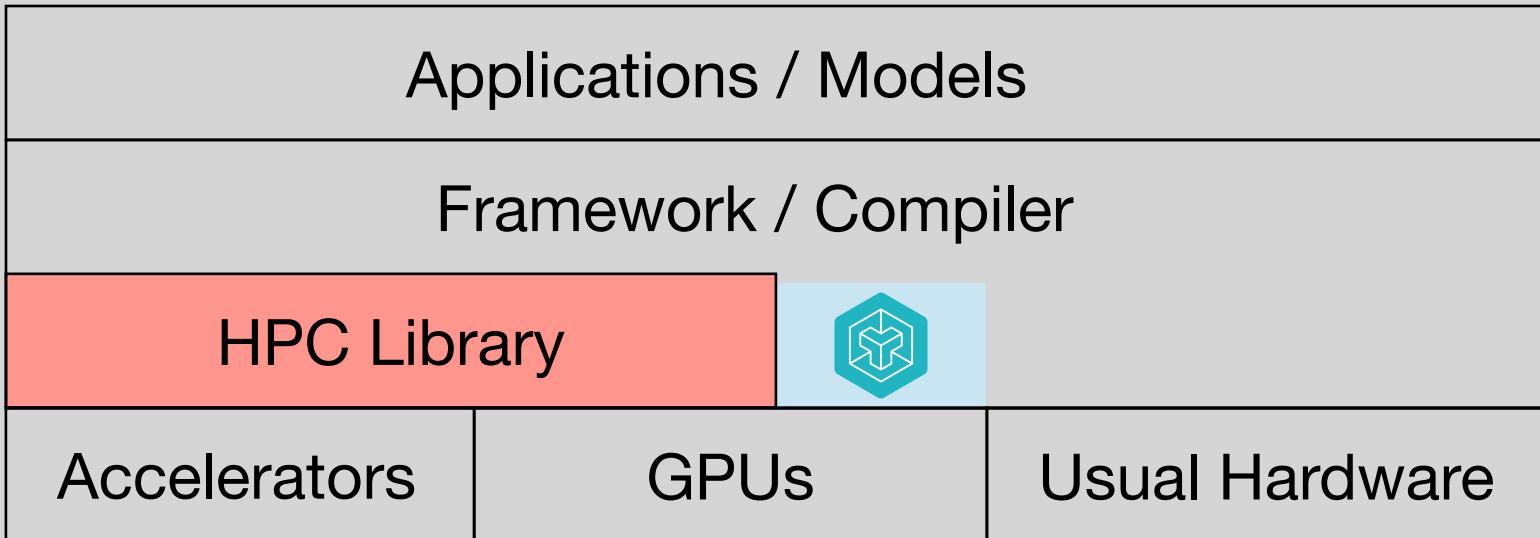
source: IBM, by way of
<https://www.eeworldonline.com/hardware-assisted-verification-technology-in-the-mainstream/>



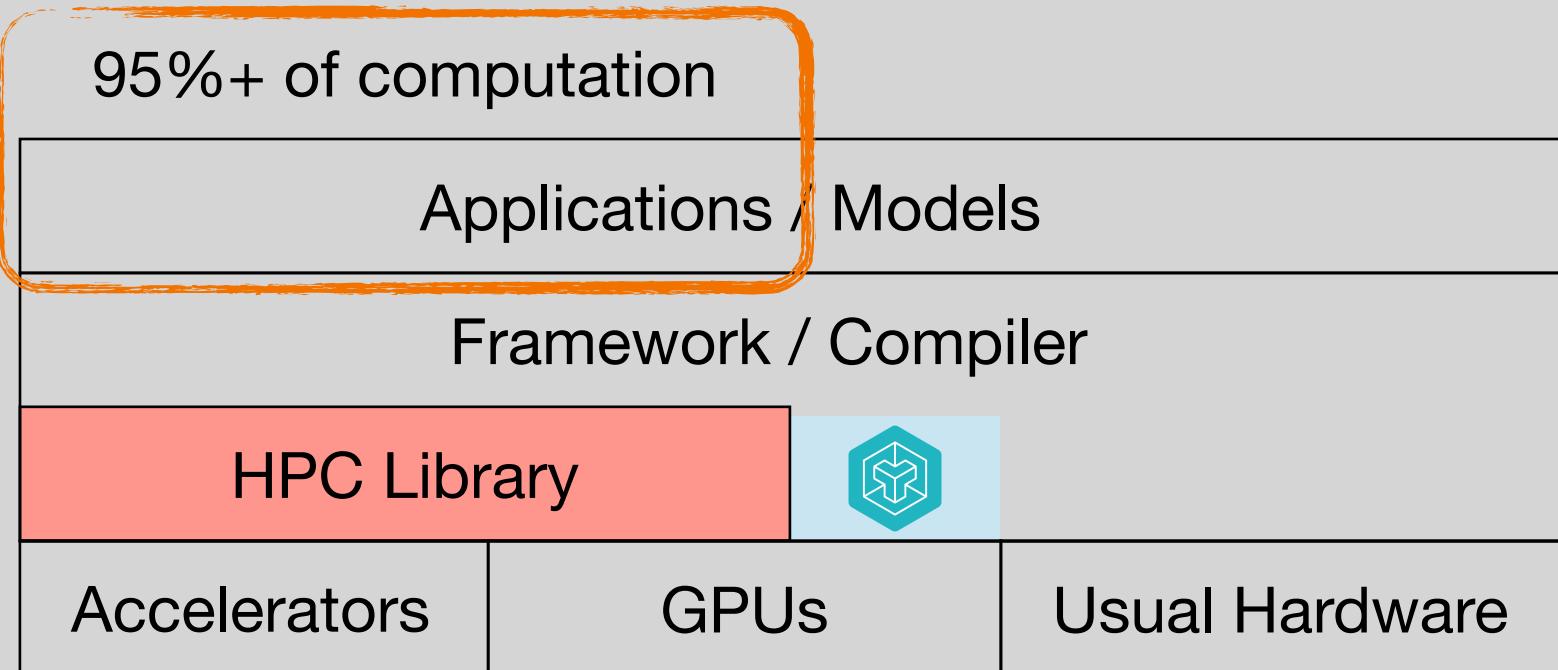
EXO



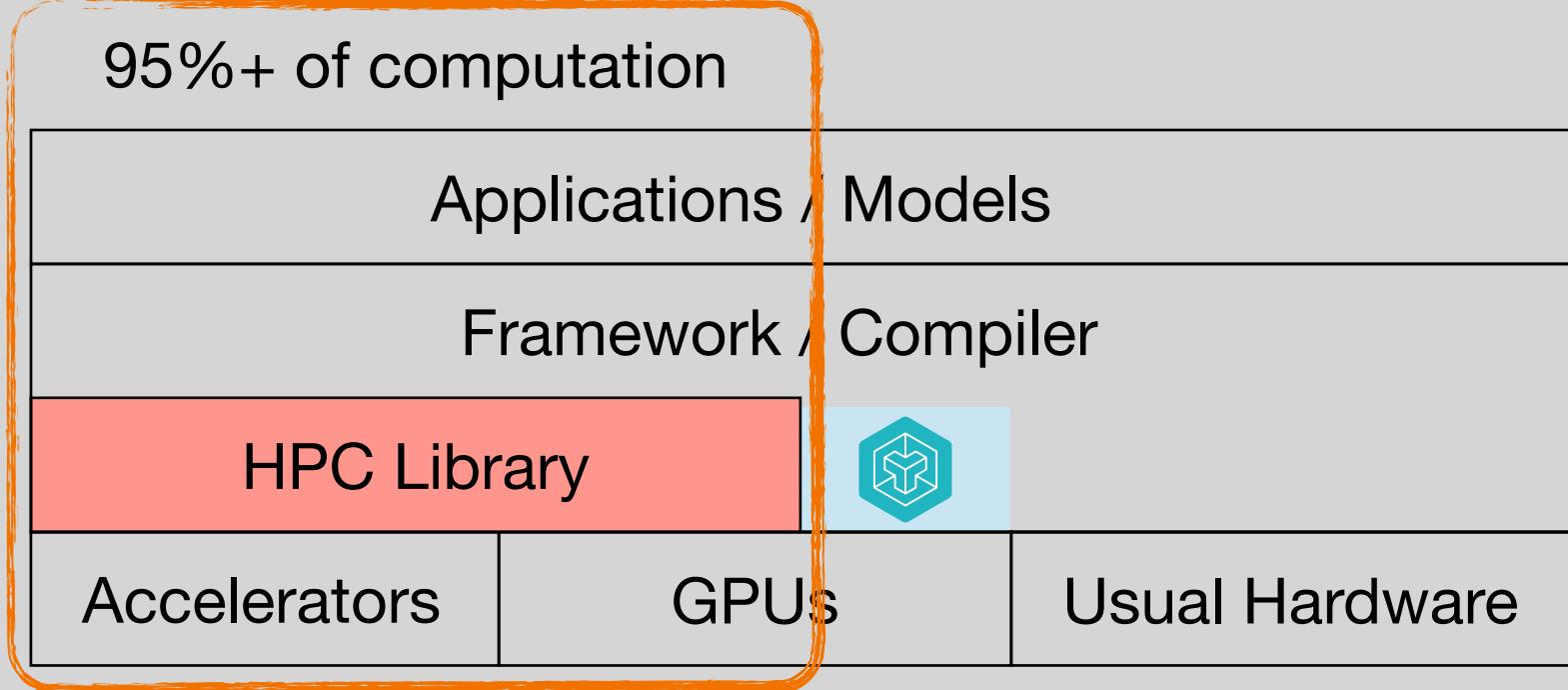
The ML Stack



The ML Stack Stovepipe



The ML Stack → Stovepipe



The ML Stack → Stovepipe

95%+ of computation

Applications / Models

Framework / Compiler

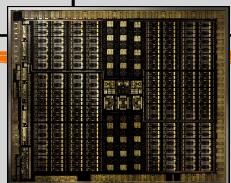
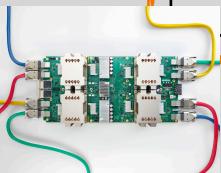
HPC Library



Accelerators

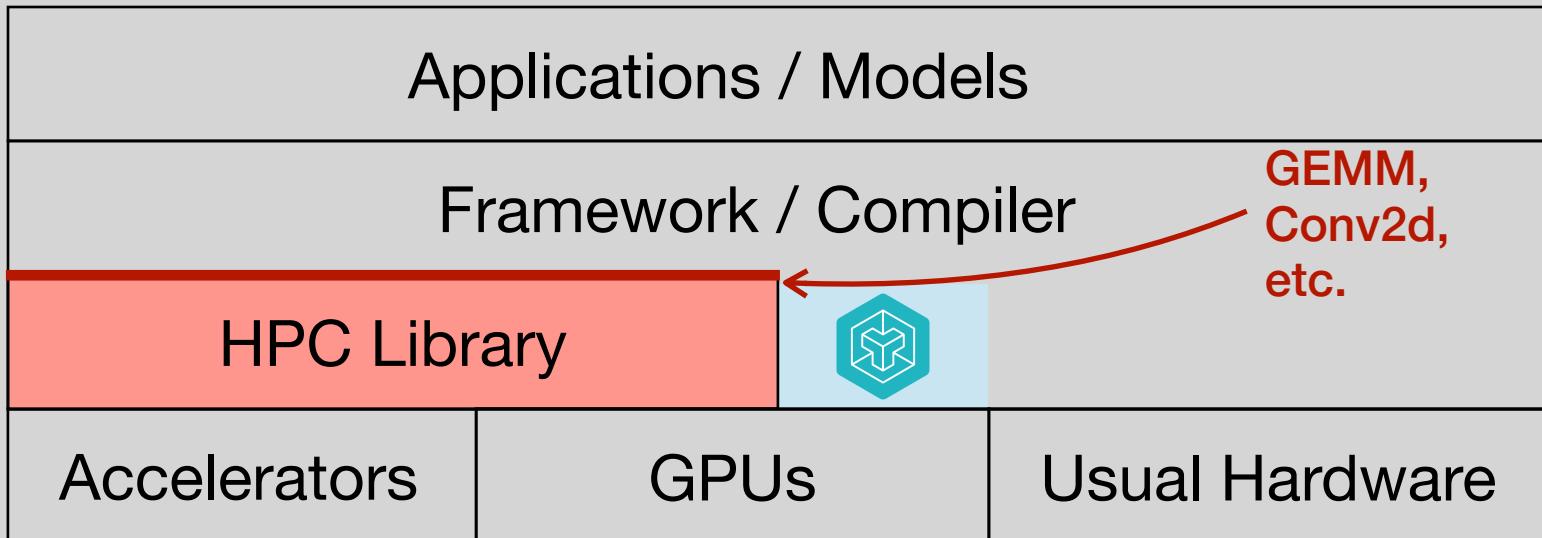
GPUs

Usual Hardware



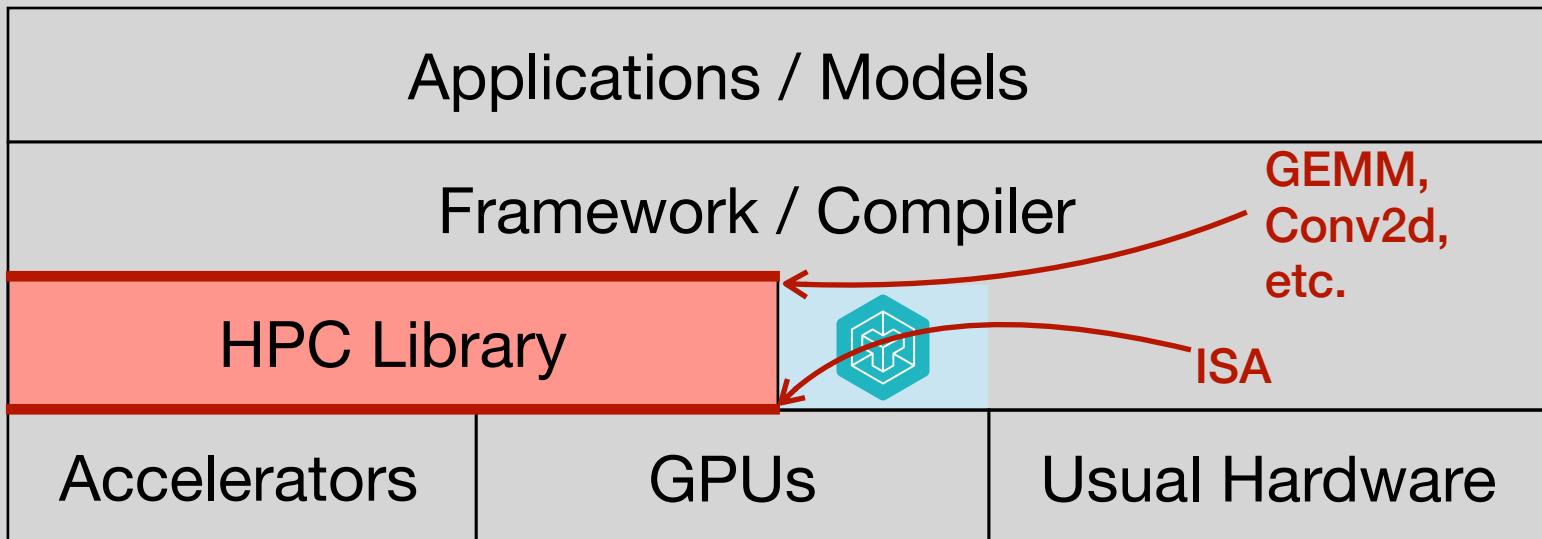
The ML Stack ~~Stack~~ Stovepipe

Consequence 1: Portability Interface



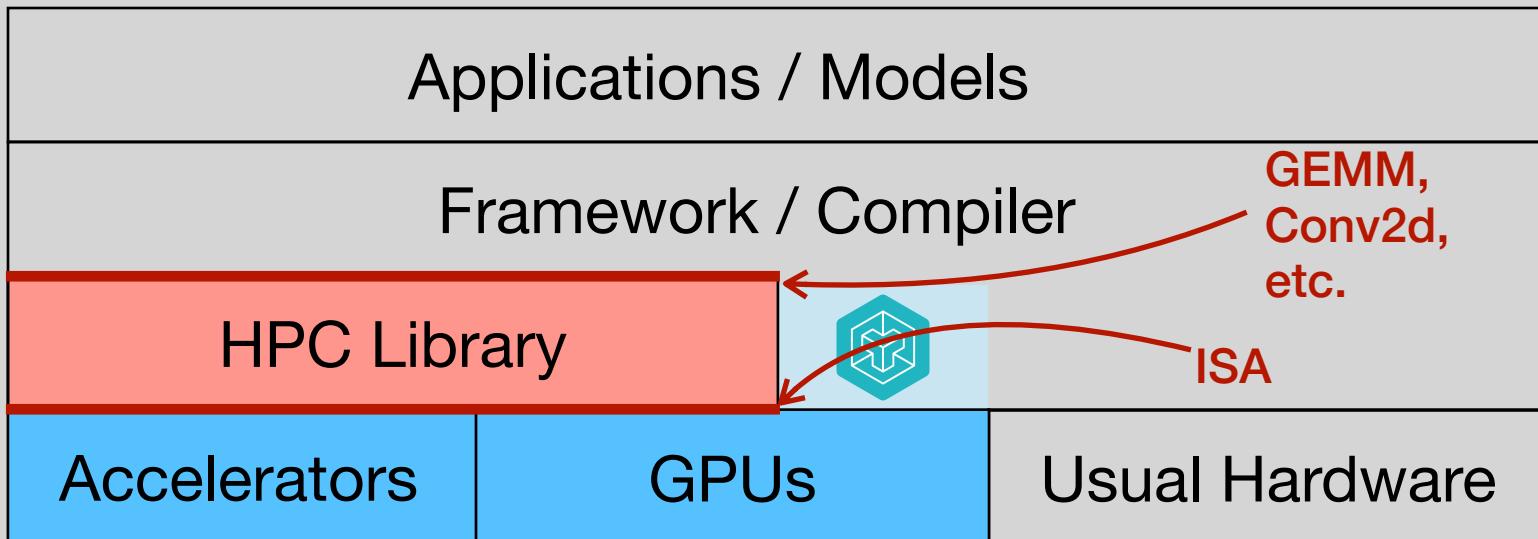
The ML Stack Stovepipe

Consequence 1: Portability Interface

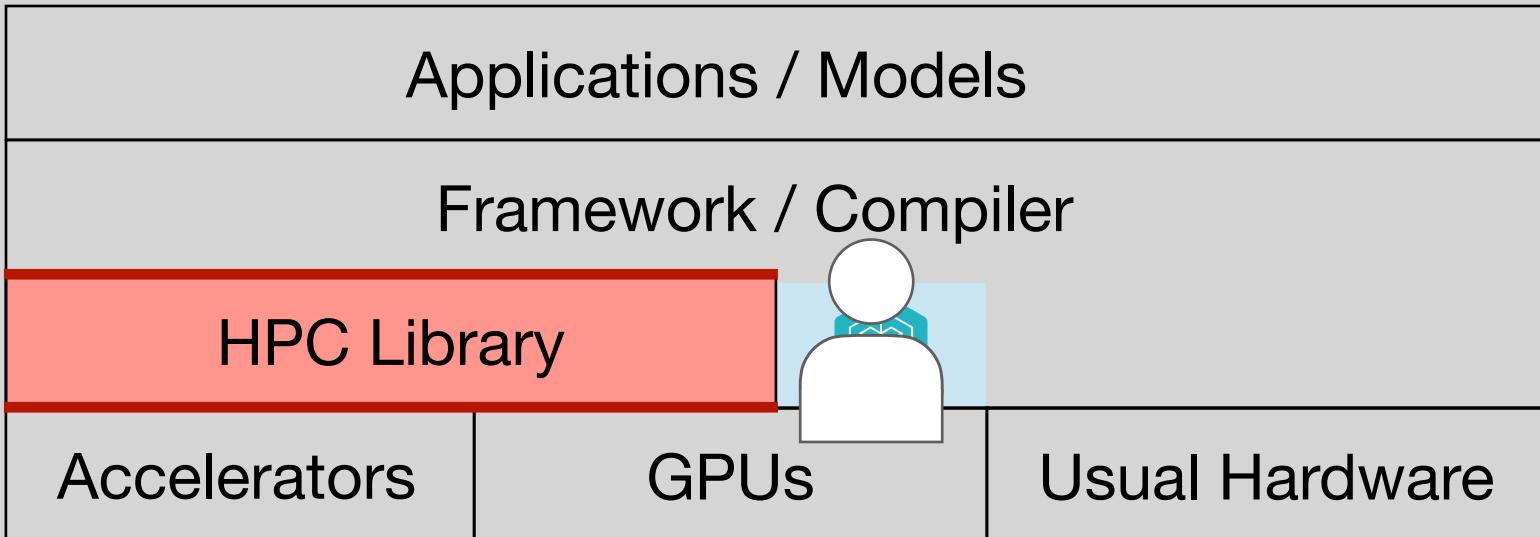


The ML Stack ~~Stack~~ Stovepipe

Consequence 2: HPC Libs are on “Critical Path”



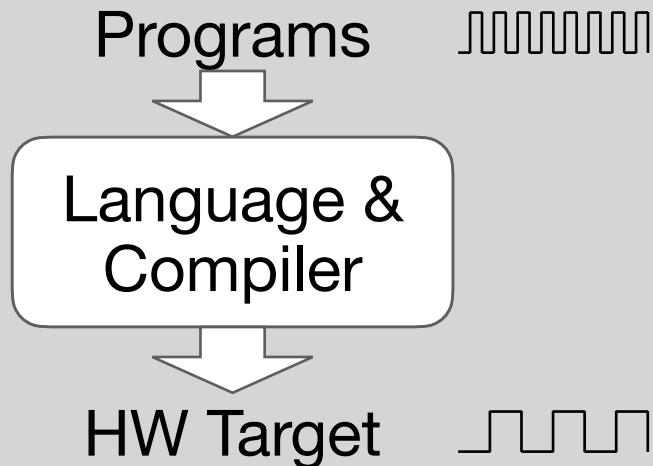
The ML Stack → Stovepipe



What's Special About Writing HPC Libraries?

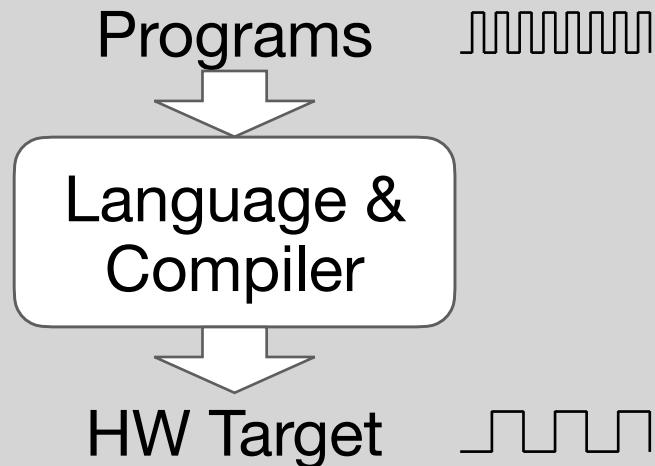
What's Special About Writing HPC Libraries?

Usual Programs

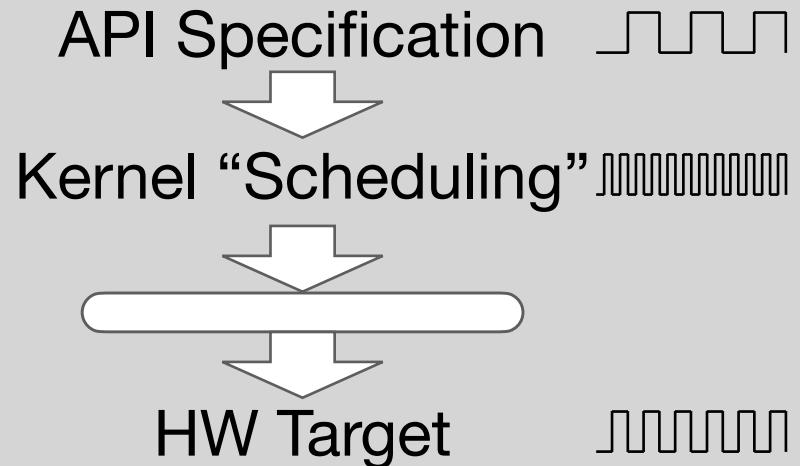


What's Special About Writing HPC Libraries?

Usual Programs

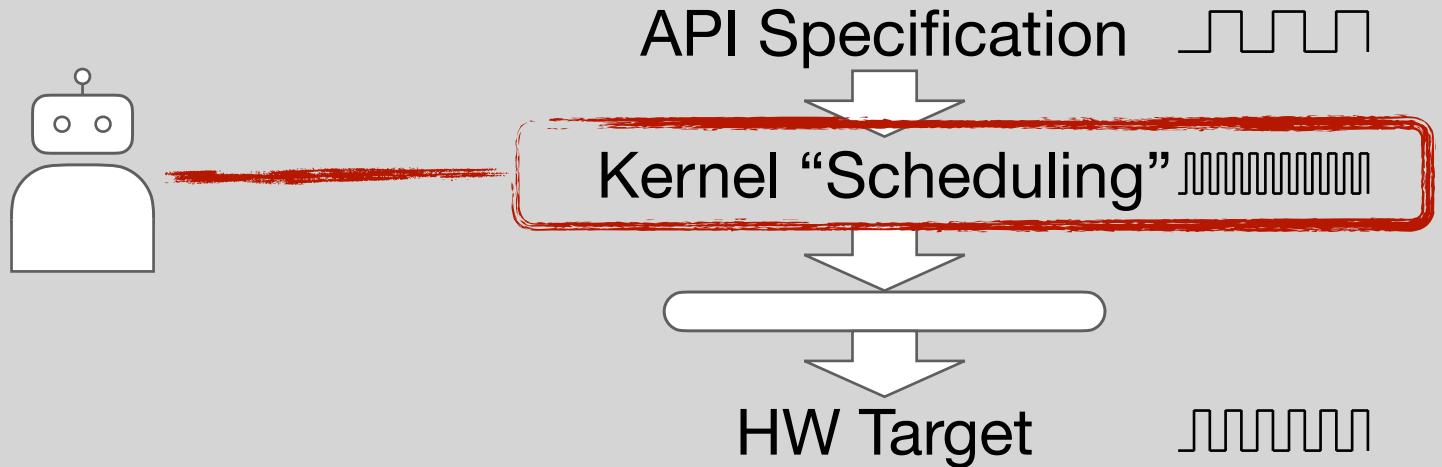


HPC Libraries

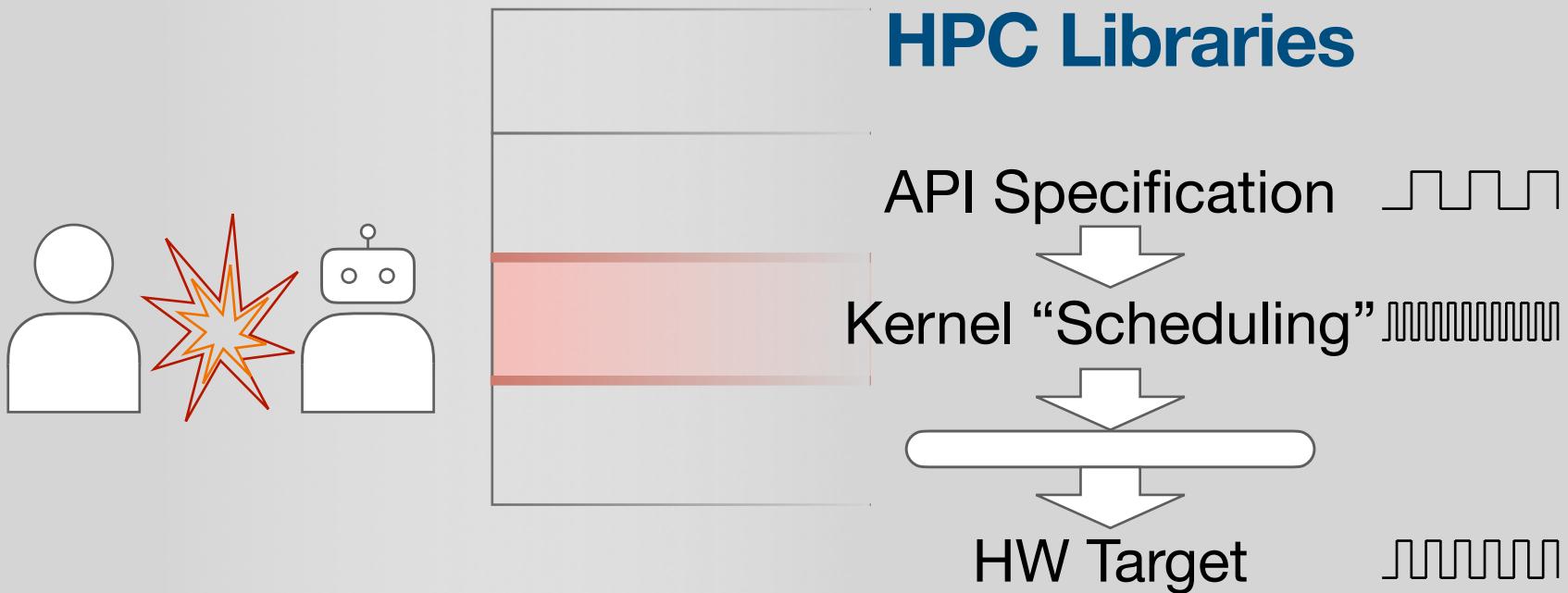


The Trouble With Automation

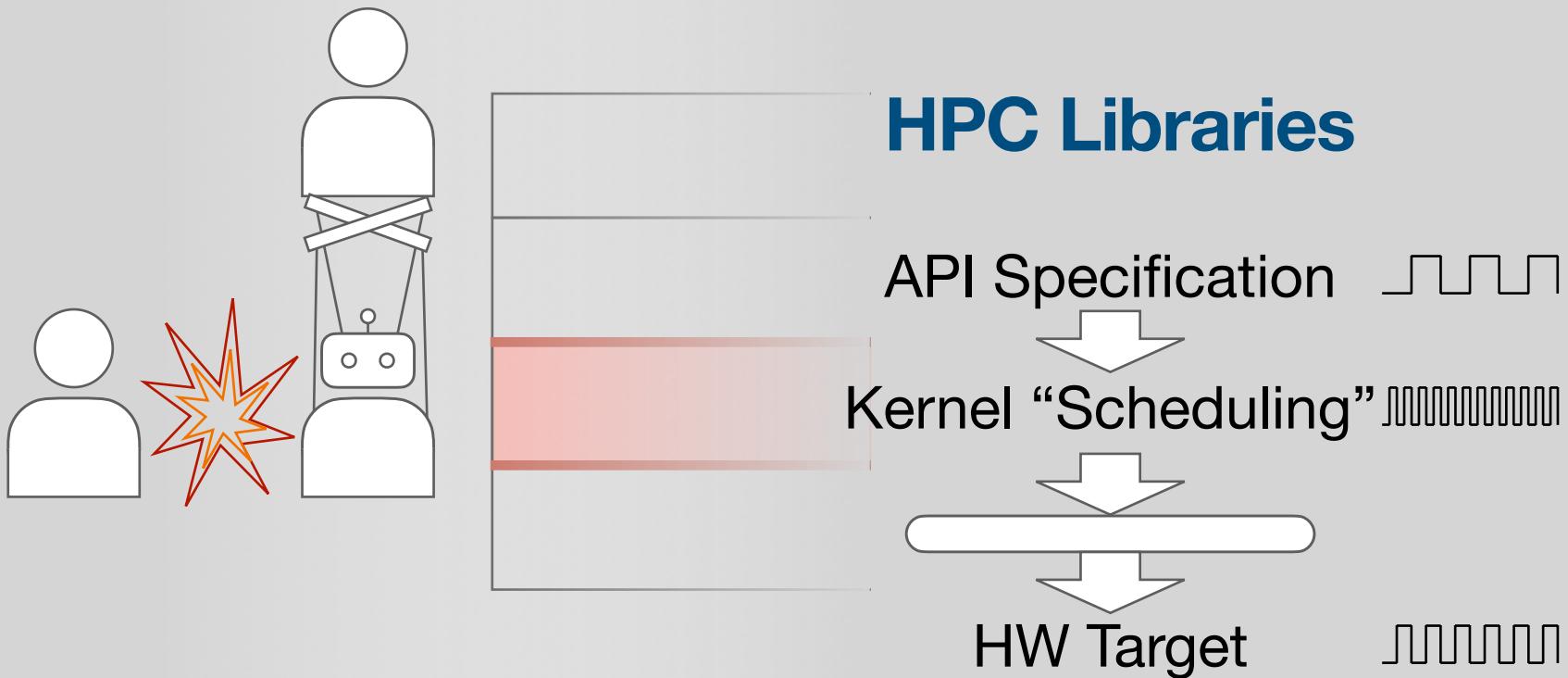
HPC Libraries



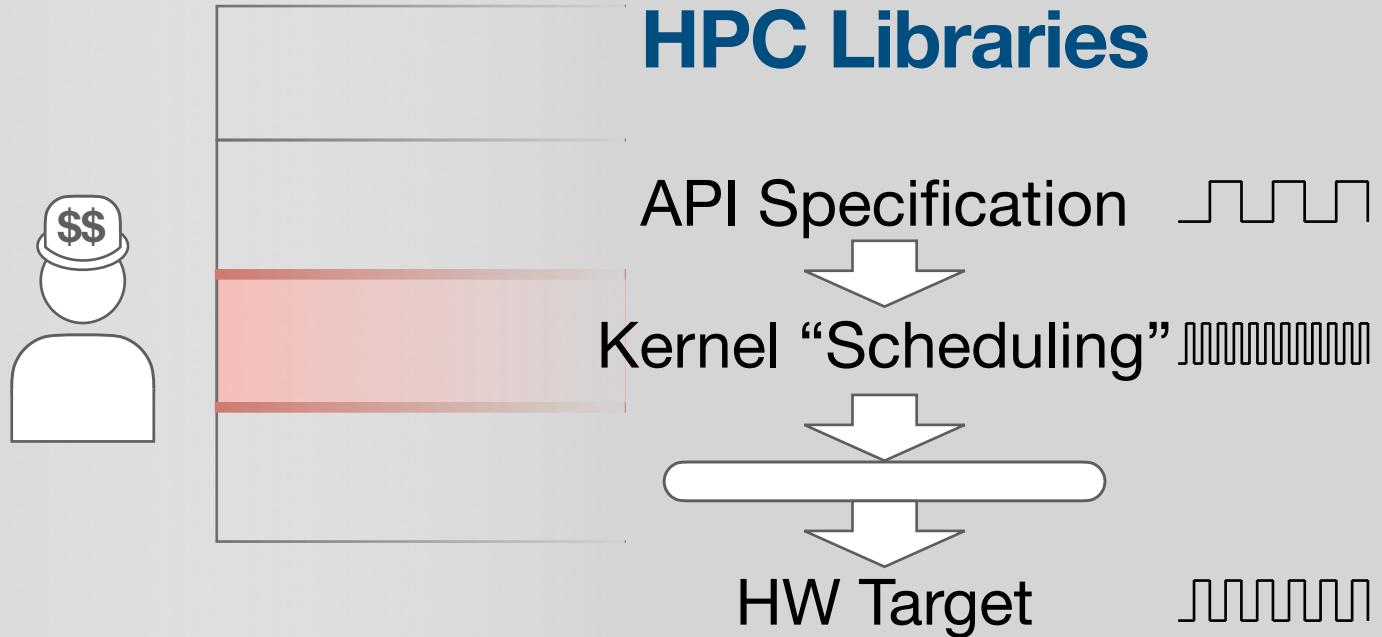
The Trouble With Automation



The Trouble With Automation

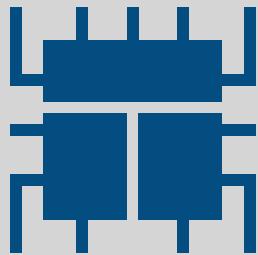


The Trouble With Automation



Exocompilation

a compiler/language design that *externalizes* parts of the compiler in order to give programmers more control



Exo

Exo code

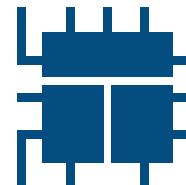
```
def gemm(N : size,      M : size,      K : size,
         A : f32[N,K], B : f32[K,M], C : f32[N,M]):
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```

“basically C”

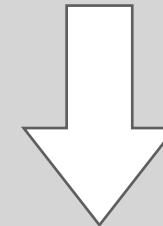


python™

```
def gemm(N : size,      M : size,      K : size,
         A : f32[N,K], B : f32[K,M], C : f32[N,M]):
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```



Exo



```
void gemm(int N, int M, int K,
          float *A, float *B, float *C) {
    for (int i=0; i < N; i++) {
        for (int j=0; j < M; j++) {
            for (int k=0; k < K; k++) {
                C[i,j] += A[i,k] * B[k,j];
            }
        }
    }
}
```



Exocompilation

a compiler/language design that *externalizes* parts of the compiler in order to give programmers more control

Two Basic Ideas in Exo

Externalize
Backends

User Scheduling

Exocompilation

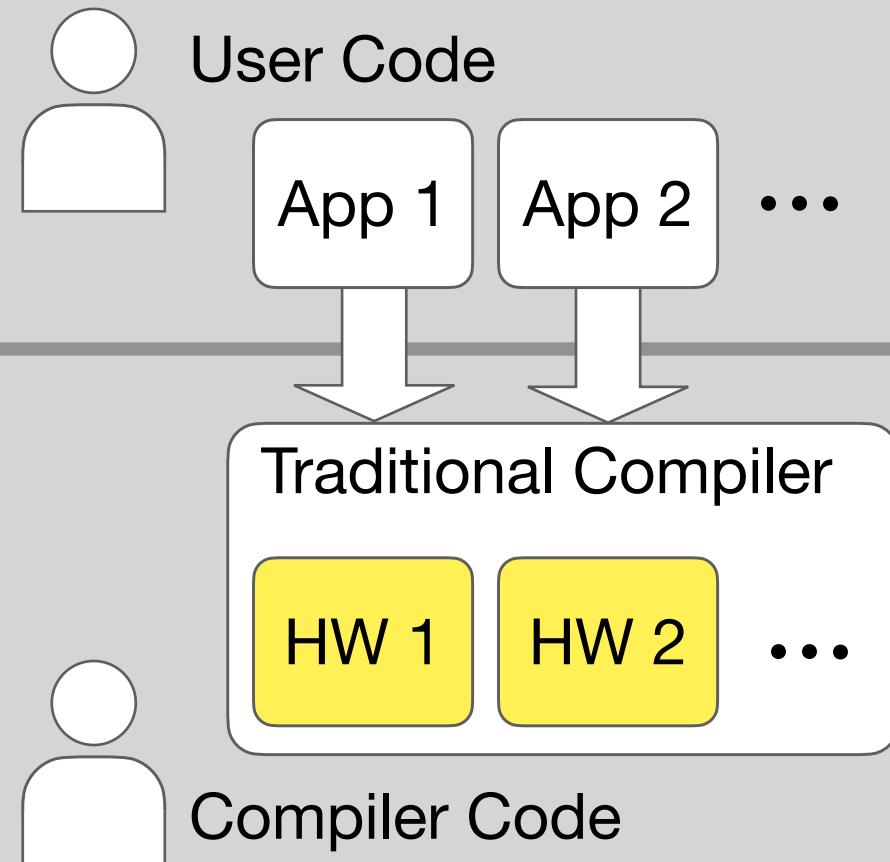
a compiler/language design that *externalizes* parts of the compiler in order to give programmers more control

Two Basic Ideas in Exo

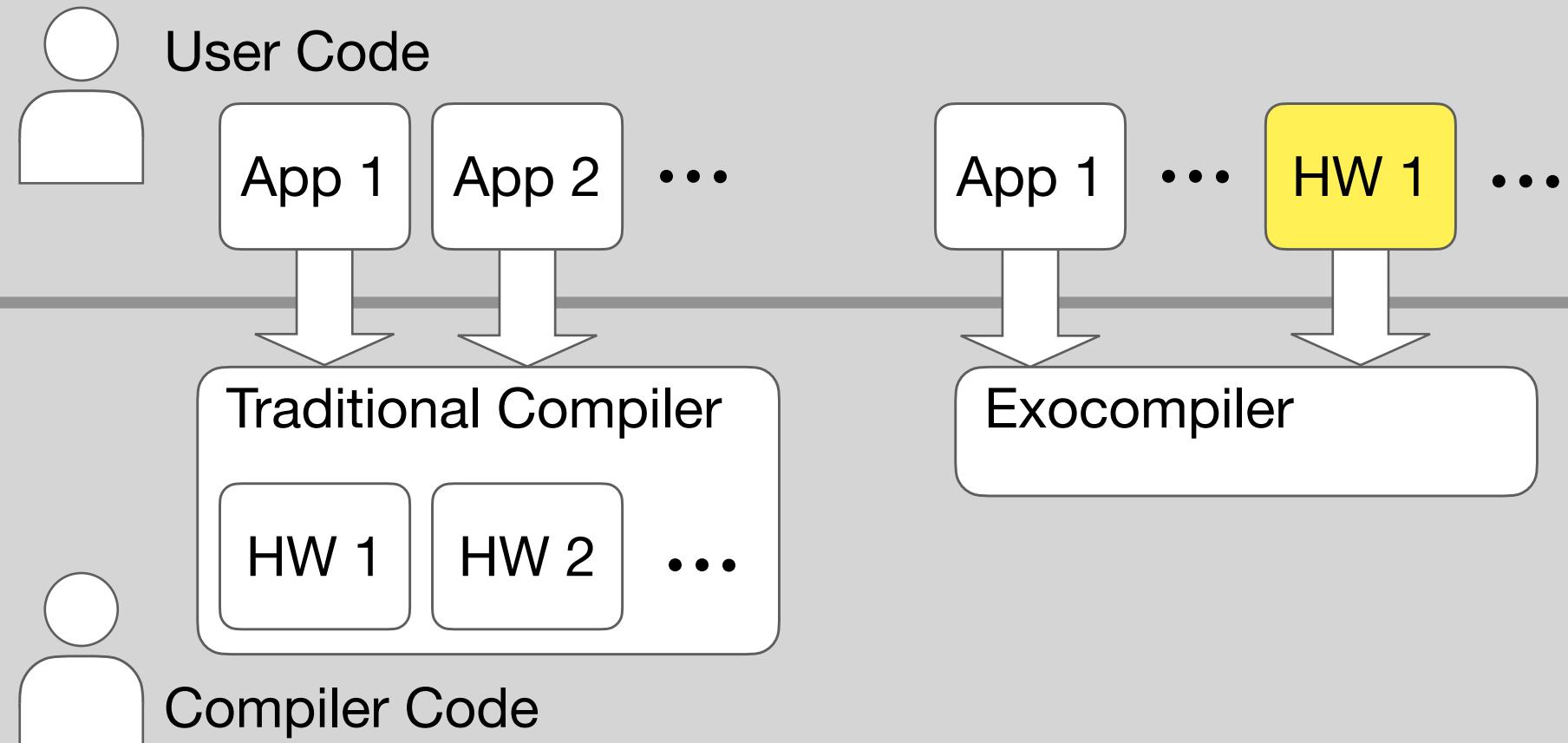
Externalize
Backends

User Scheduling

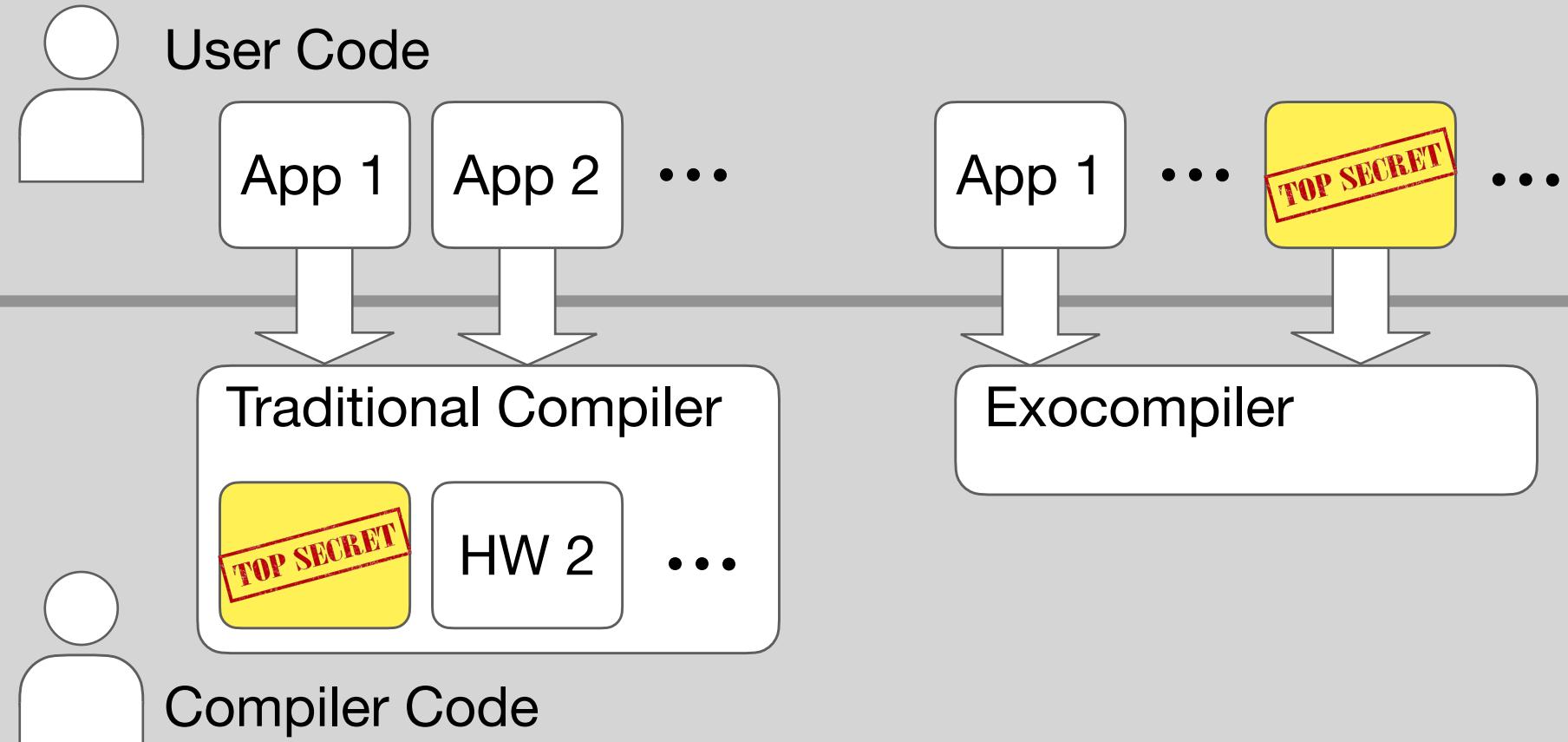
HW Backends as User-Level Libraries



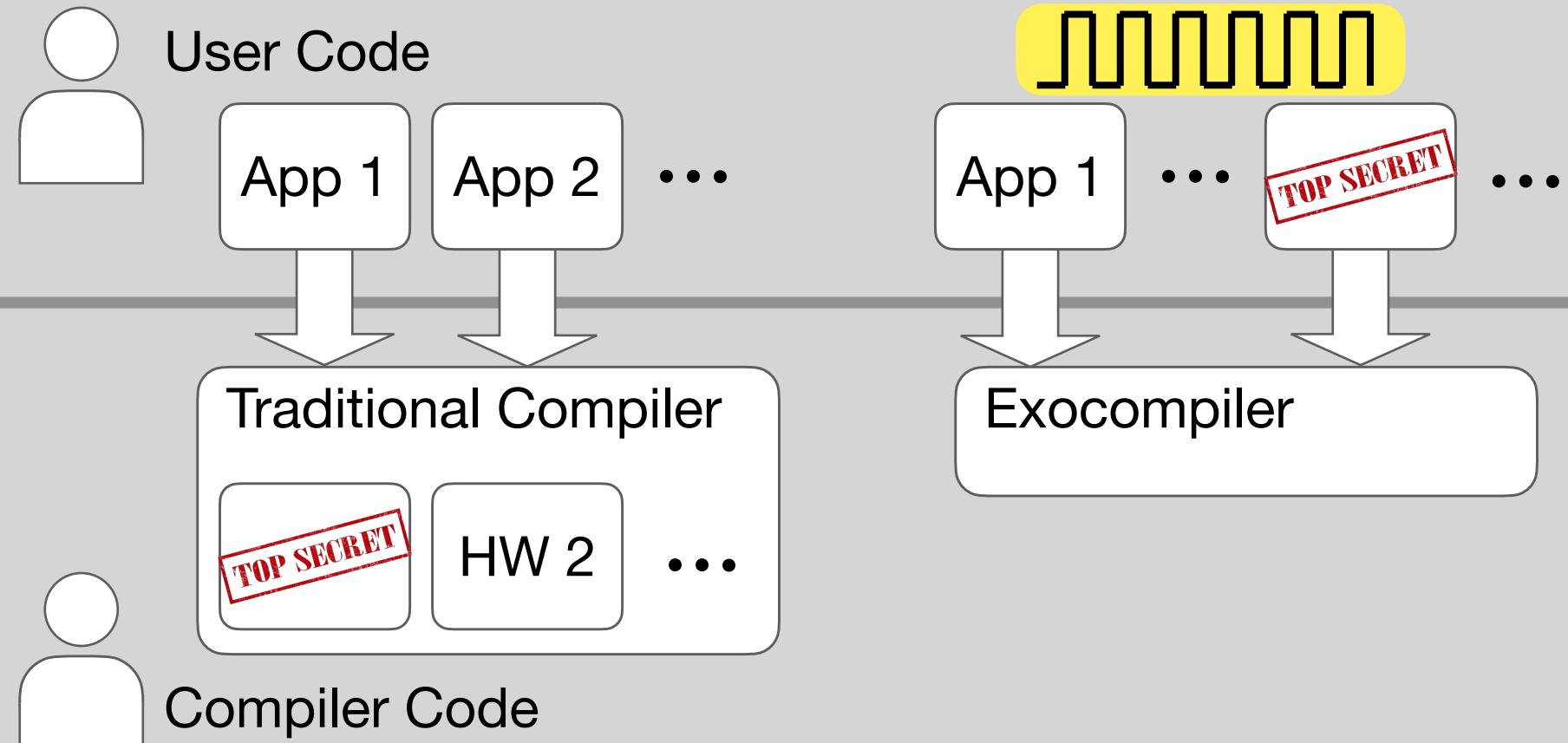
HW Backends as User-Level Libraries



HW Backends as User-Level Libraries



HW Backends as User-Level Libraries



Exocompilation

a compiler/language design that *externalizes* parts of the compiler in order to give programmers more control

Two Basic Ideas in Exo

Externalize
Backends

User Scheduling

User-Schedulable Languages

```
def gemm(N : size,      M : size,      K : size,
         A : f32[N,K], B : f32[K,M], C : f32[N,M]):
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```

User-Schedulable Languages

```
def gemm(N : size,      M : size,      K : size,
         A : f32[N,K], B : f32[K,M], C : f32[N,M]):
```

```
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```

User-Schedulable Languages

```
def gemm(N : size,      M : size,      K : size,
         A : f32[N,K], B : f32[K,M], C : f32[N,M]):
```

```
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```

User-Schedulable Languages

e.g.

CHILL [’08]

Halide [SIGGRAPH ’12]

TACO [OOPSLA ’17]

LIFT [CGO ’17]

TVM [OSDI ’18]

Tiramisu [CGO ’19]

```
def gemm(N : size,      M : size,      K : size,
        A : f32[N,K], B : f32[K,M], C : f32[N,M]): 
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```

```
def gemm(N : size,      M : size,      K : size,
        A : f32[N,K], B : f32[K,M], C : f32[N,M]): 
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```

User-Schedulable Languages

e.g.

CHILL [’08]

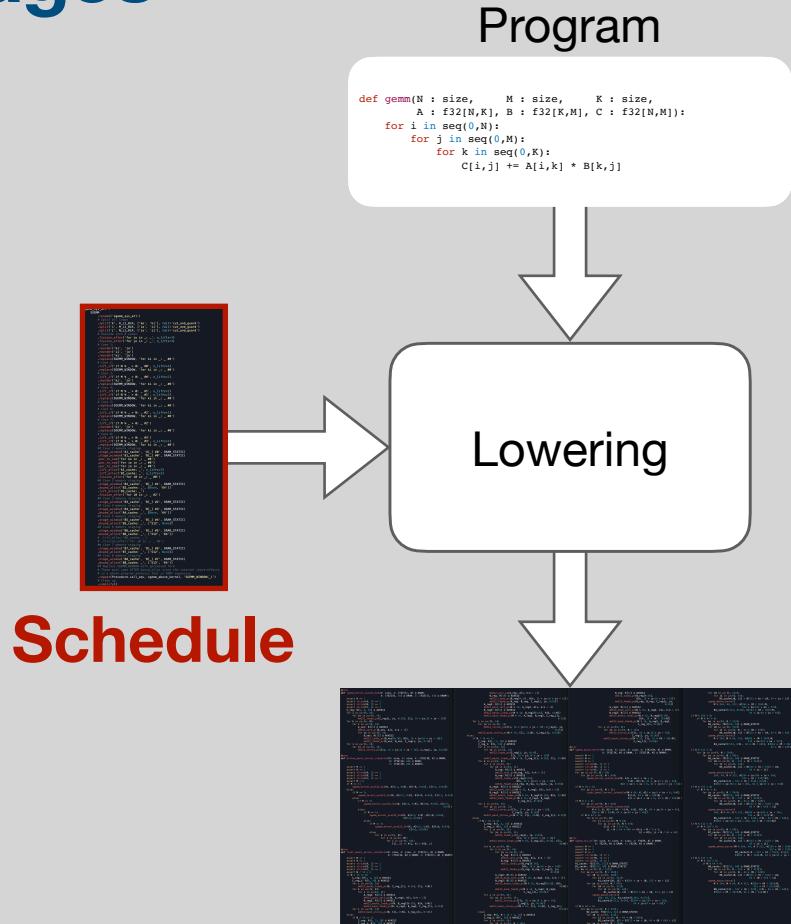
Halide [SIGGRAPH ’12]

TACO [OOPSLA ’17]

LIFT [CGO ’17]

TVM [OSDI ’18]

Tiramisu [CGO ’19]



User-Schedulable Languages

e.g.

CHILL [’08]

Halide [SIGGRAPH ’12]

TACO [OOPSLA ’17]

LIFT [CGO ’17]

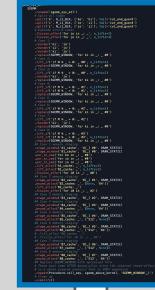
TVM [OSDI ’18]

Tiramisu [CGO ’19]

Program1

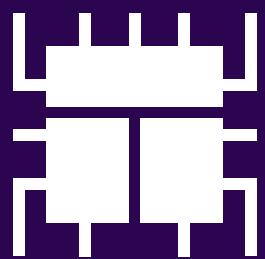
```
def gemm(N : size,
         A : f32[N,K], B : f32[K,M], C : f32[N,M]):
    for i in seq(0,N):
        for j in seq(0,M):
            for k in seq(0,K):
                C[i,j] += A[i,k] * B[k,j]
```

Schedule



```
def gemm(N : size,
         A : f32[N,K], B : f32[K,M], C : f32[N,M]):
    for i in seq(0,N):
        for k in seq(0,K):
            for j in seq(0,M):
                C[i,j] += A[i,k] * B[k,j]
```

Program2



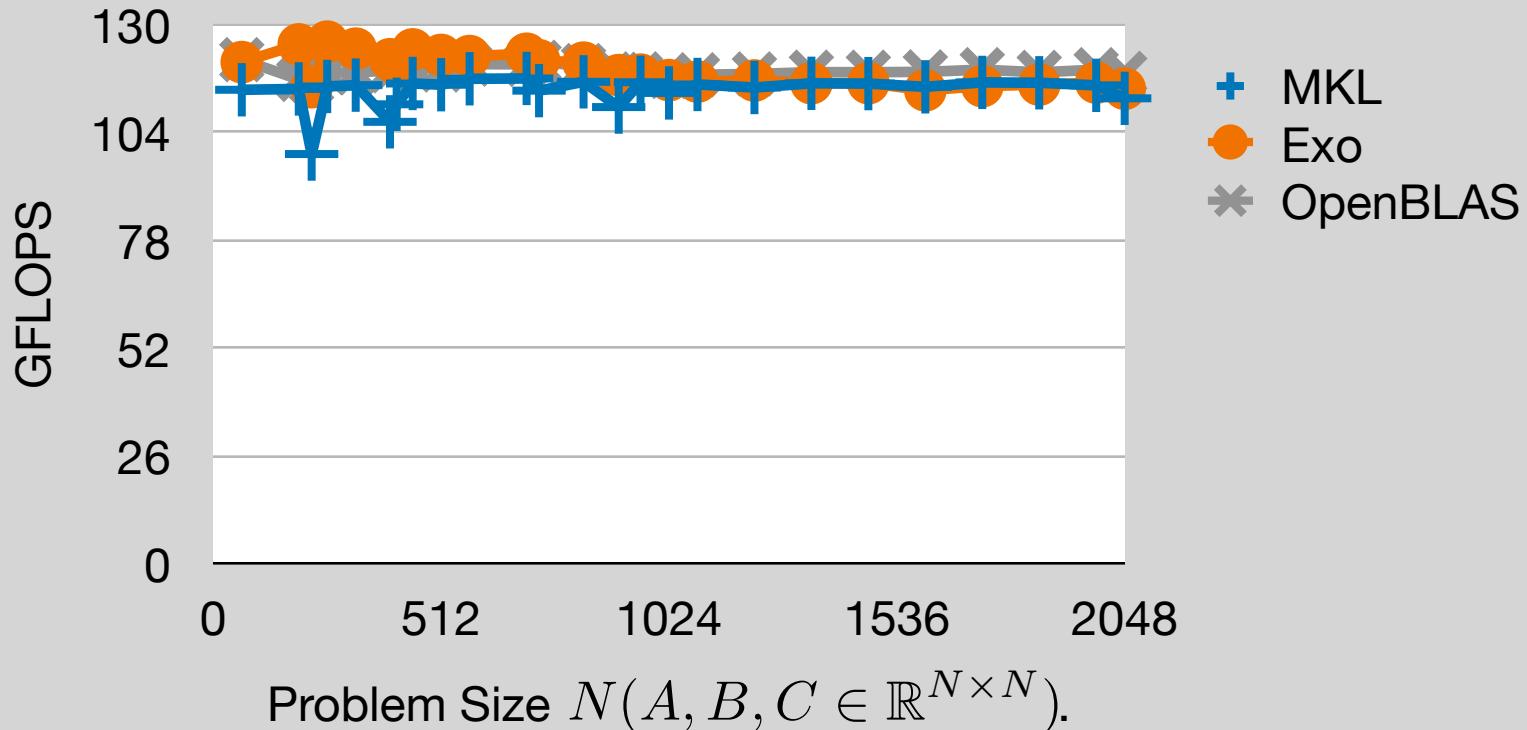
Exo Demo

Some Performance Numbers

*Results from
2021*

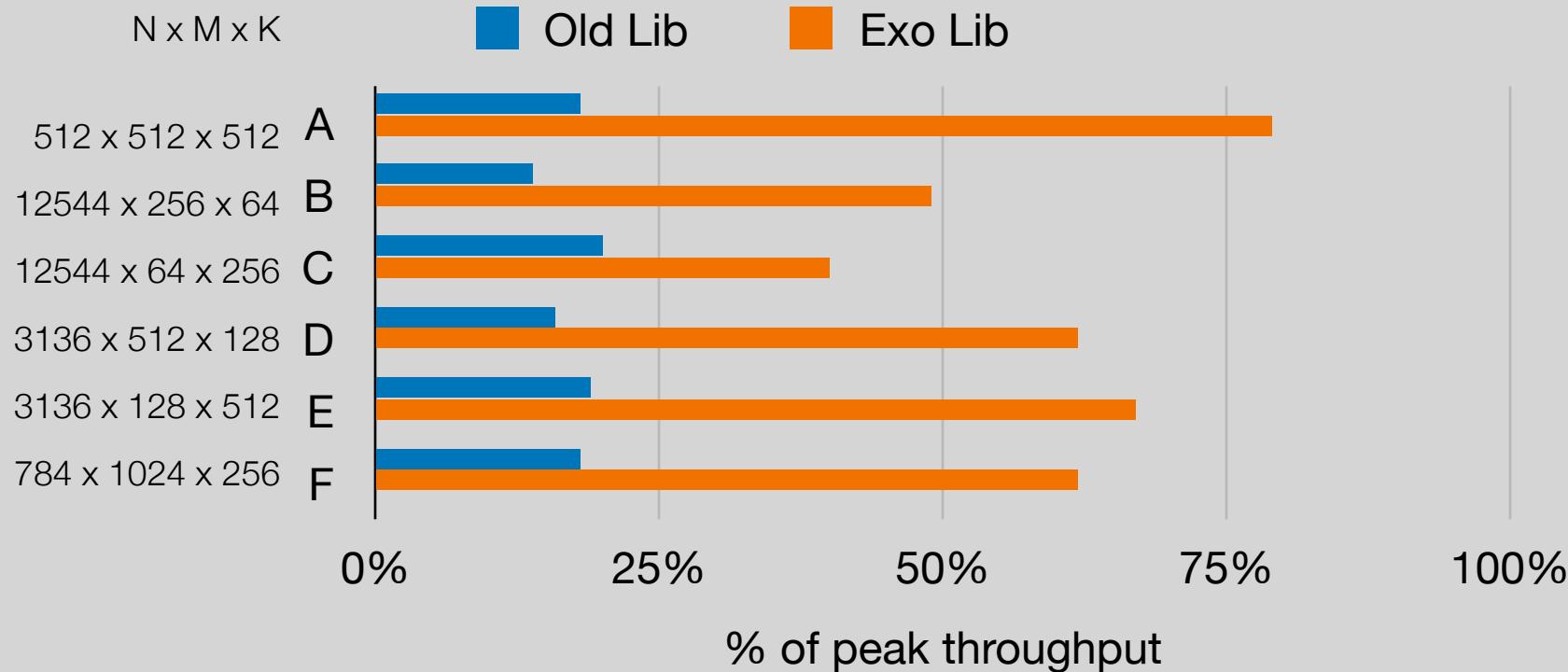
AVX-512 MatMul

Matches best possible performance



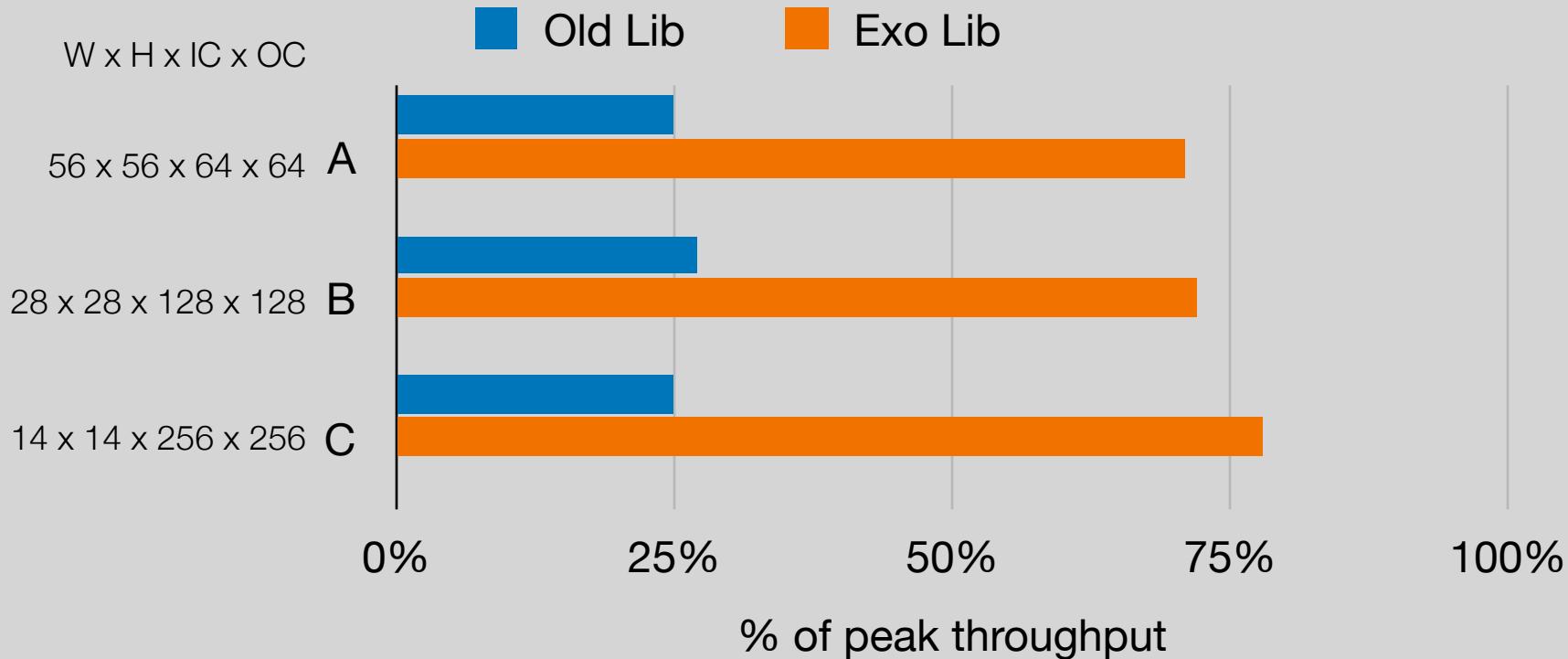
GEMMINI MatMul

2.3x - 3.7x faster than original HPC Library



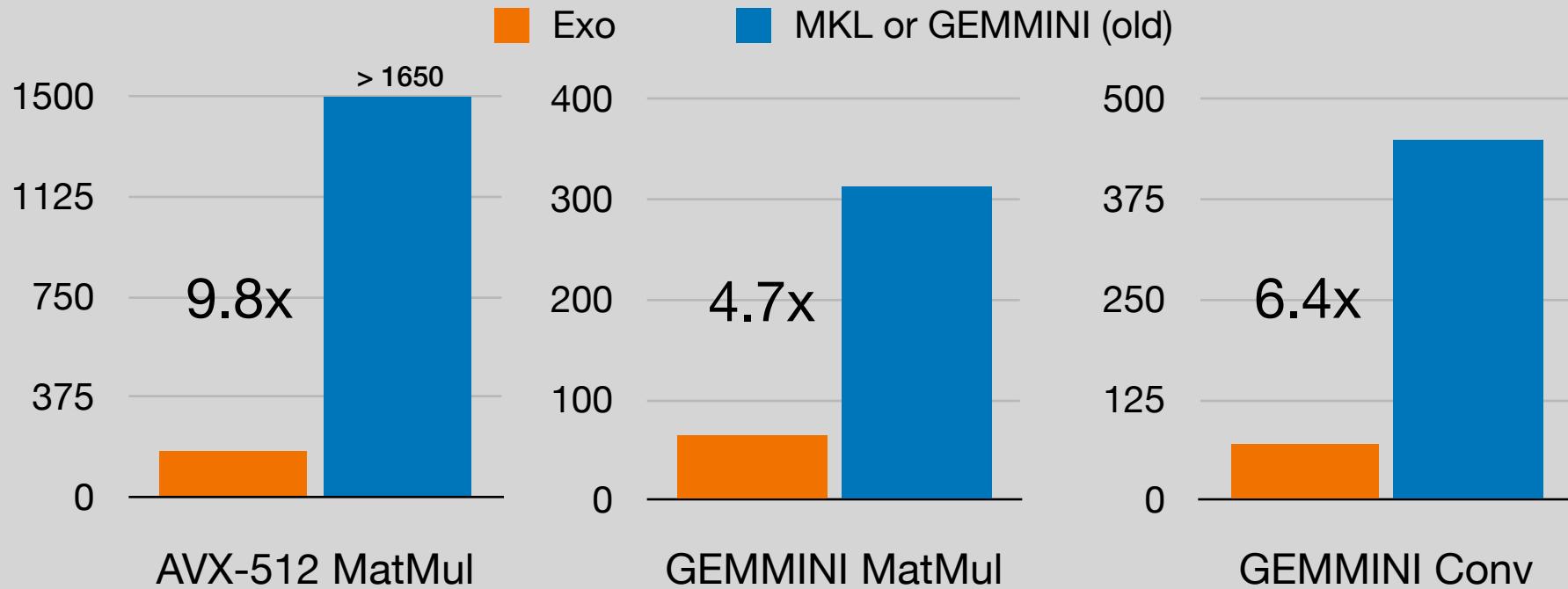
GEMMINI Conv

Similar Story for Convolution Layers



*Results from
2021*

Exo vs. Existing - Lines of Code



Some Important Limitations


$$32*j + 4*k + 8$$

Linear/Affine Indexing Expression


$$i*j + k$$

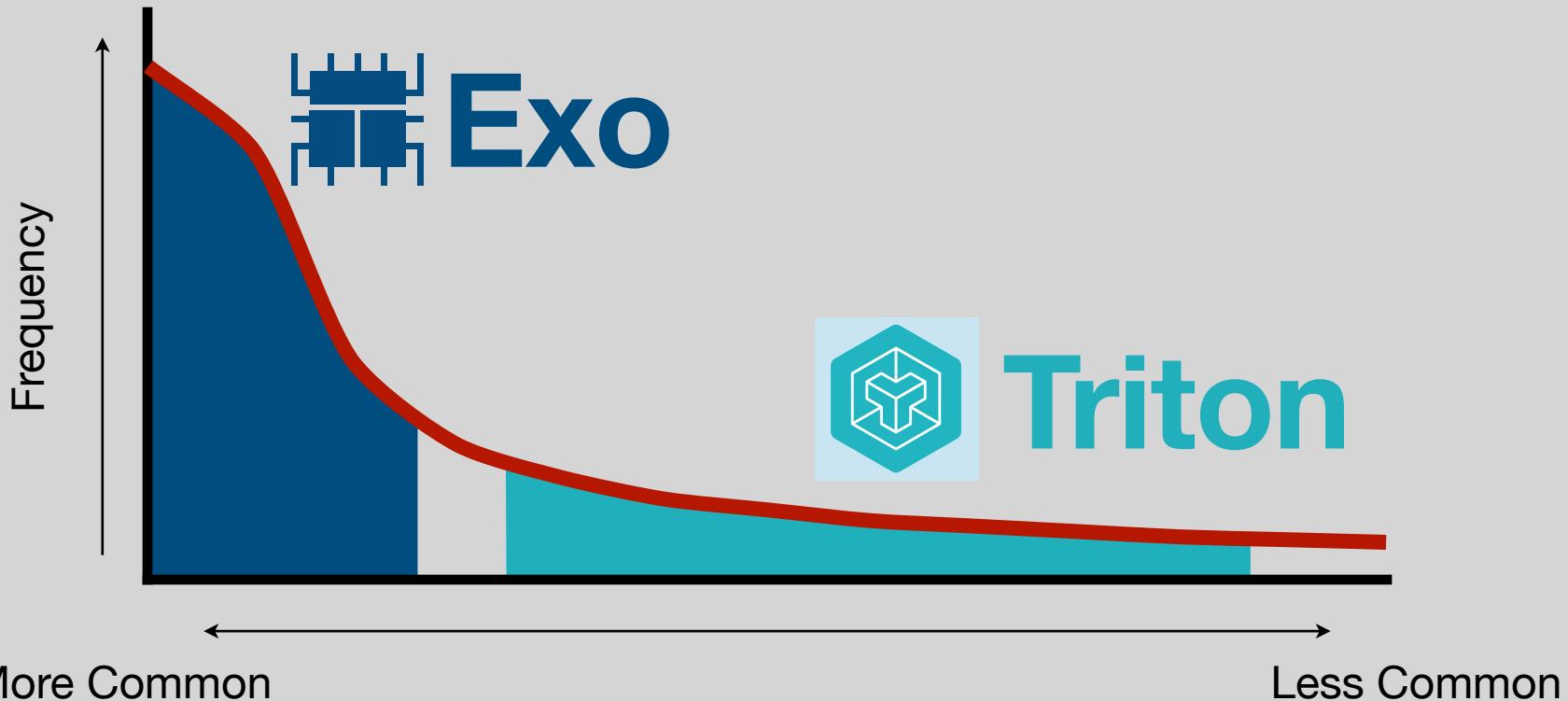
Non-Linear Indexing



Currently Unsupported



Distribution of Kernel Workload



Metaprogramming Spectrum

No Metaprogramming

Write Your Own Compiler



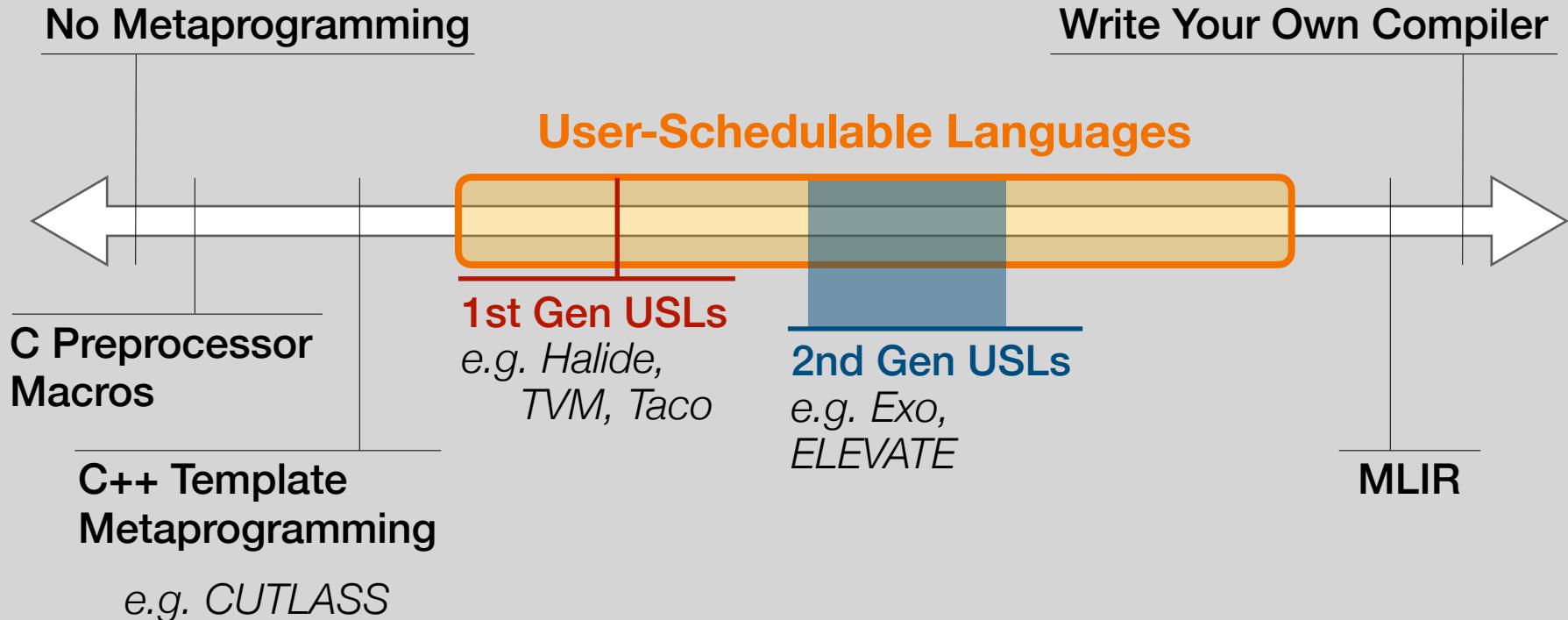
C Preprocessor
Macros

C++ Template
Metaprogramming

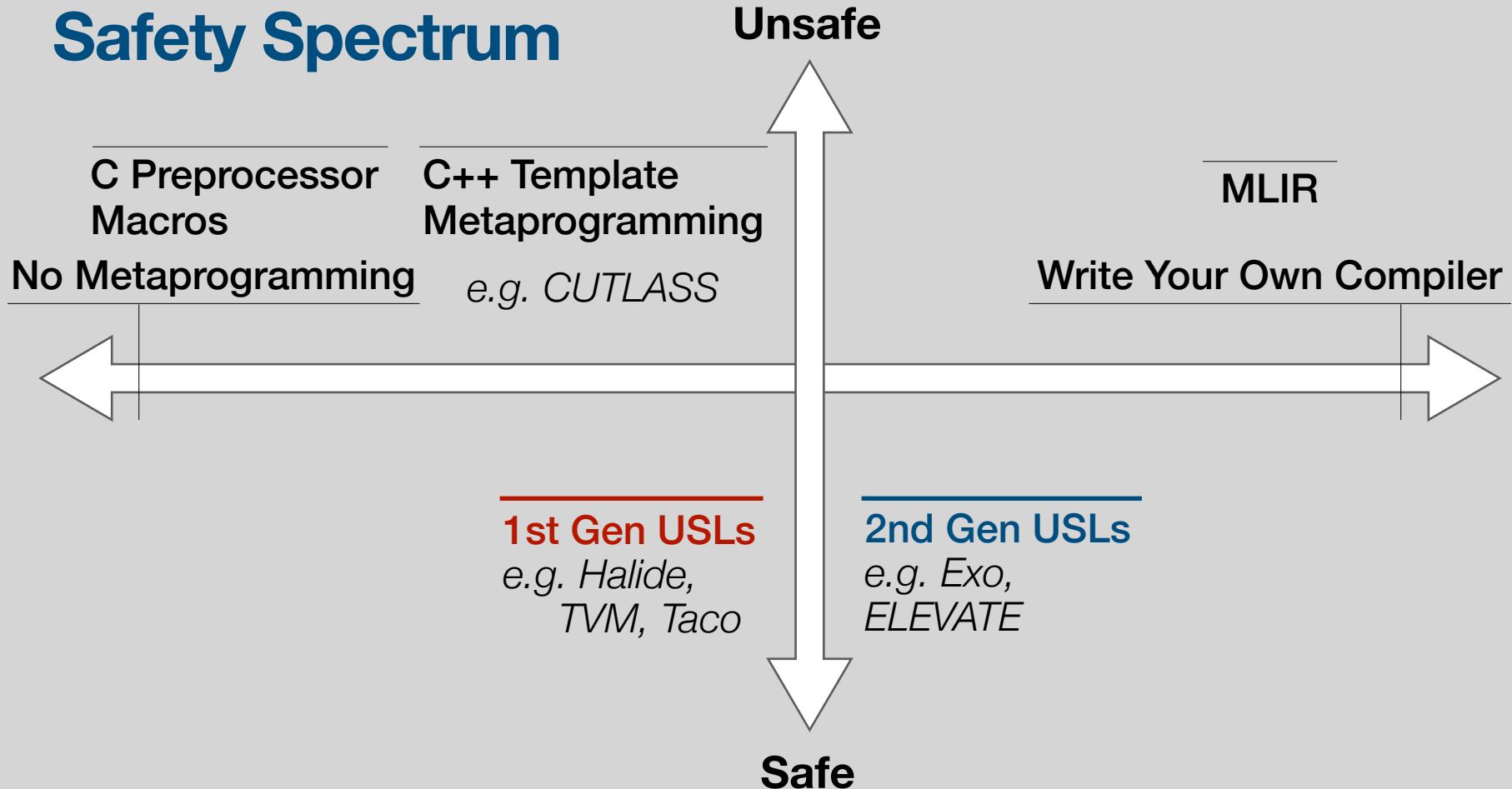
e.g. CUTLASS

MLIR

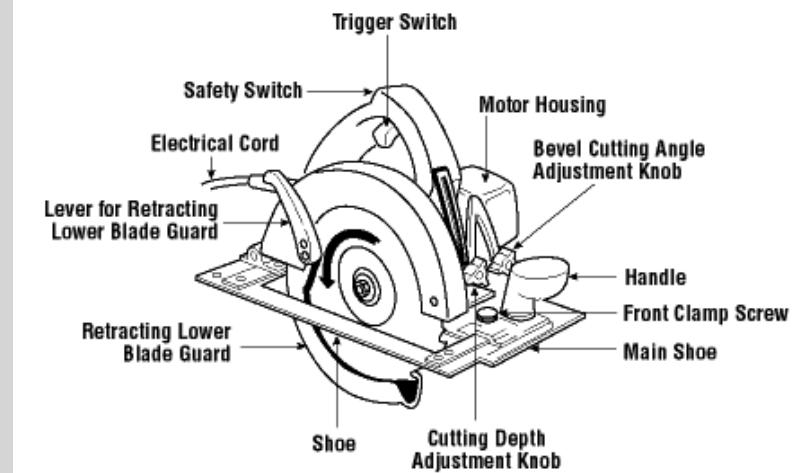
Metaprogramming Spectrum



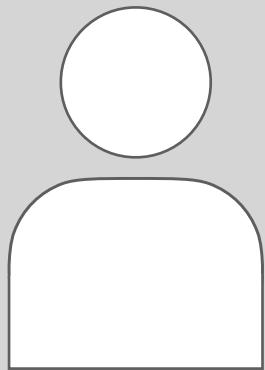
Safety Spectrum



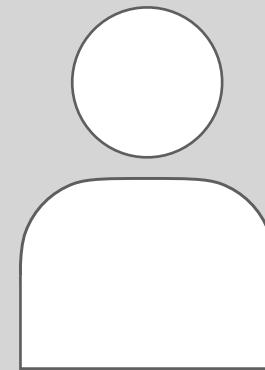
Power Tools have Safety Features



**Performance
Engineer**



Correctness



**Compiler
Programmer**

Performance
Optimization

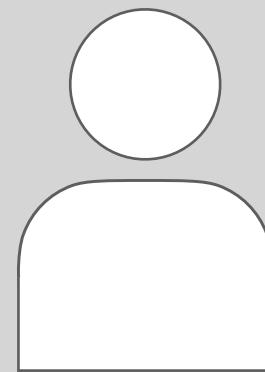
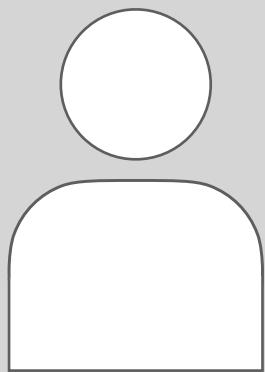
Correctness



Exo

Performance
Engineer

Compiler
Programmer





Tutorials & Documentation
are...

