# Polyhedral Rescheduling of GPU Kernels To Exploit Async Memory Movement

Ivan R. Ivanov,
William Moses, Emil Vatai, Toshio Endo, Jens Domke, Alex Zinenko

# Background

Advanced features in GPUs are underused

- Legacy code
- Code not written by GPU experts
- Untuned code
- Synthesized code

# GPU features

Async global->shared copies

# GPU features

Faster reductions

Tensor cores
TMA

# GPU features

More memory movement features

# Common theme

- A lot of async memory movement
- Specialized hardware for memory and computation (TMA and Tensor cores)

...

Need to overlap them and stress them to make **full** use of your modern GPU

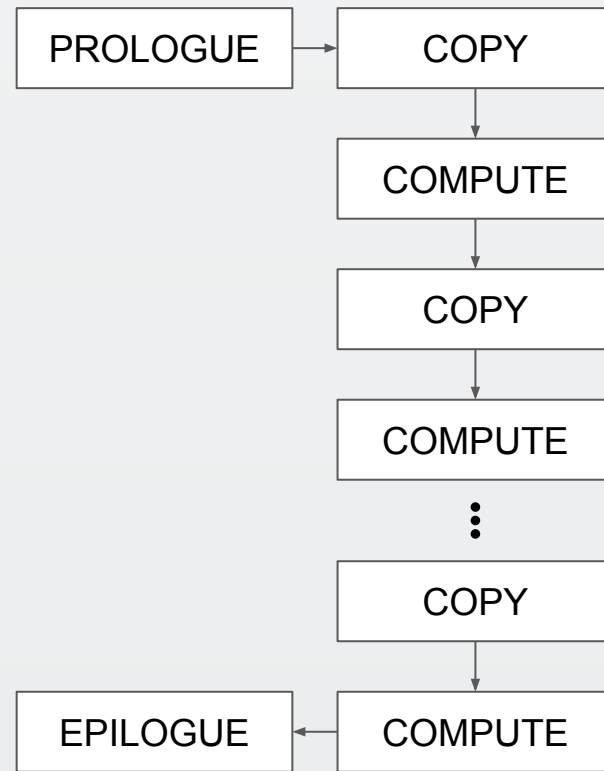# Example usage

# Blocked matmul
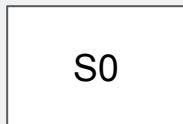
```
__global__ void MatrixMulCUDA(float *C, float *A,
                              float *B, int wA,
                              int wB) {
  int bx = blockIdx.x;
  int by = blockIdx.y;
  int tx = threadIdx.x;
  int ty = threadIdx.y;
  int aBegin = wA * BLOCK_SIZE * by;
  int aEnd   = aBegin + wA - 1;
  int aStep  = BLOCK_SIZE;
  int bBegin = BLOCK_SIZE * bx;
  int bStep  = BLOCK_SIZE * wB;
  float Csub = 0;
  for (int a = aBegin, b = bBegin;
       a <= aEnd;
       a += aStep, b += bStep) {
    __shared__ float As[BLOCK_SIZE][BLOCK_SIZE];
    __shared__ float Bs[BLOCK_SIZE][BLOCK_SIZE];
    As[ty][tx] = A[a + wA * ty + tx];
    Bs[ty][tx] = B[b + wB * ty + tx];
    __syncthreads();
    for (int k = 0; k < BLOCK_SIZE; ++k)
      Csub += As[ty][k] * Bs[k][tx];
    __syncthreads();
  }
  int c = wB * BLOCK_SIZE * by + BLOCK_SIZE * bx;
  C[c + wB * ty + tx] = Csub;
}
```
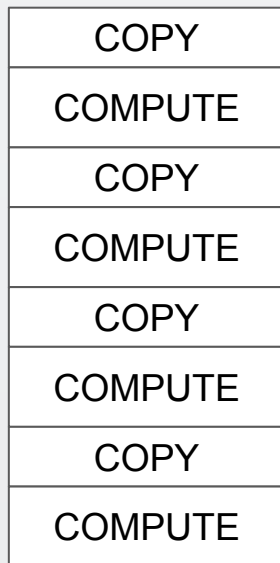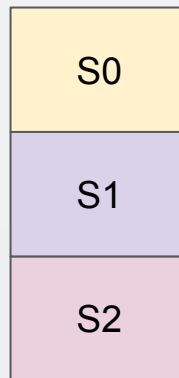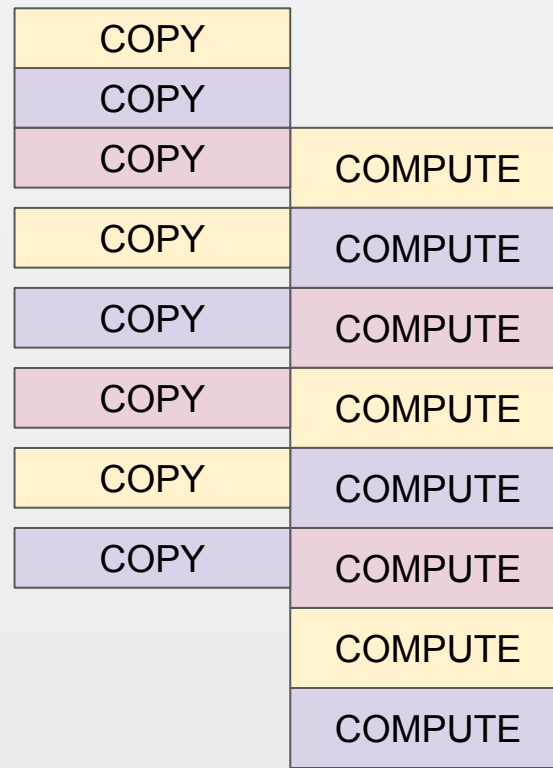
# Optimized matmul - pipelining

Memory

Computation

| S0 |
|----|

| COPY |
|------|
| COMPUTE |
| COPY |
| COMPUTE |
| COPY |
| COMPUTE |
| COPY |
| COMPUTE |

Memory

Computation

| S0 |
|----|
| S1 |
| S2 |

| COPY | |
|------|--|
| COPY | |
| COPY | COMPUTE |
| COPY | COMPUTE |
| COPY | COMPUTE |
| COPY | COMPUTE |
| COPY | COMPUTE |
| COPY | COMPUTE |
| | COMPUTE |
| | COMPUTE |

# Pipelining using nvvm intrinsics

COPY COMPUTE
COPY COMPUTE
COPY COMPUTE

```
async.cp global[i] -> shared[j]
...
commit_group
wait_group 2

compute using shared[...]
```

Can we automatically optimize existing code to use these features?

# Polyhedral model

# Polyhedral model

```
for (int i = 0; i < N; i++) {
  for (int j = 0; j < i; j++) {
    S(i, j);
    R(i, j);
  }
}
```

```
domain

{ S(i, j) : 0 <= i < N,
            0 <= j < i }
{ R(i, j) : 0 <= i < N,
            0 <= j < i }
```

```
dependencies

{ S(si, sj) -> R(ri, rj) :
    si = ri, sj = rj }
```

Everything is represented using linear algebra

# Existing work on polyhedral representation

- Extracting the polyhedral structure
  - Source code - pet
  - Intermediate representation - polly (LLVM), polygeist+polymer(MLIR)

**What about synchronisation?**

-> No prior work in importing code with synchronisation into a polyhedral representation.

It is essential to be able to handle barriers for GPU code (extremely prevalent)

# Existing work on polyhedral scheduling

- CPU/Generic: polly, isl, pluto, etc
- GPU: ppcg

Parallelism **Supported**

But everything assumes synchronous execution *within* a thread of execution
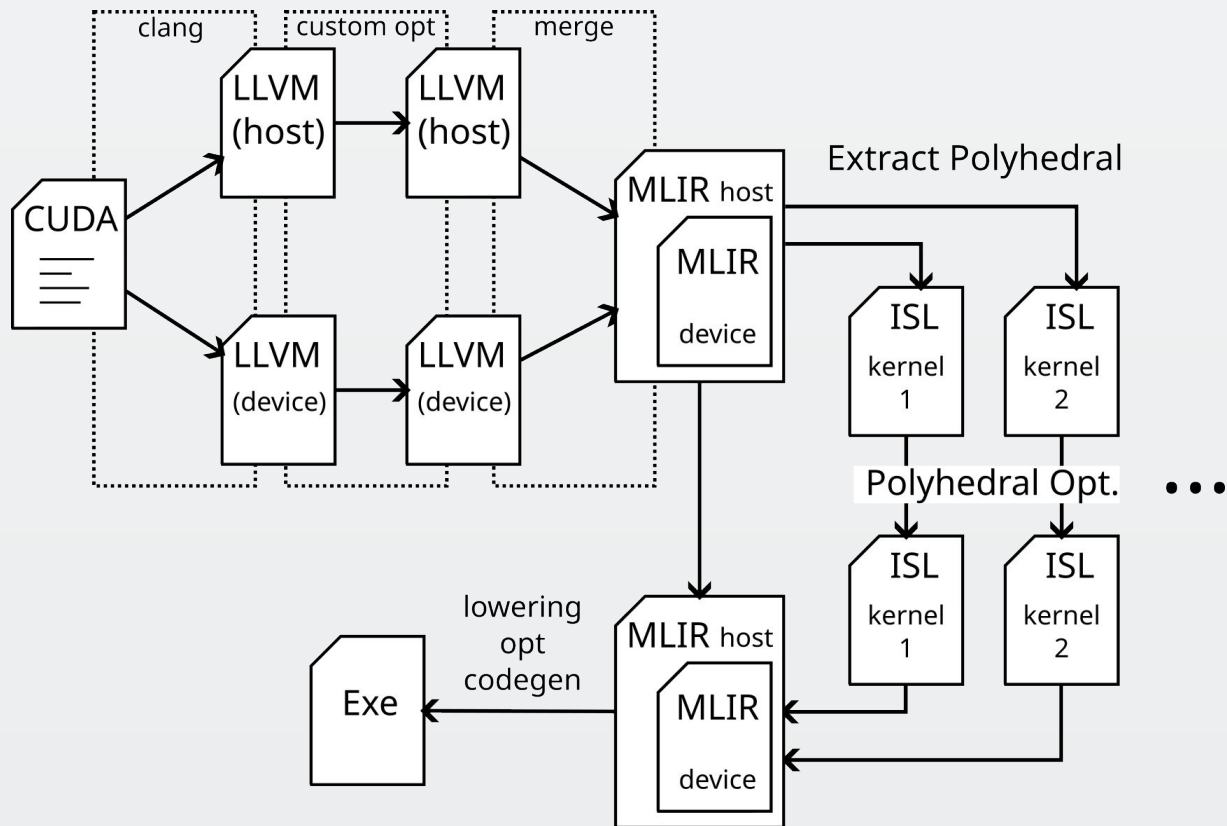
What about asynchronous execution? e.g. **nvvm.async.cp**
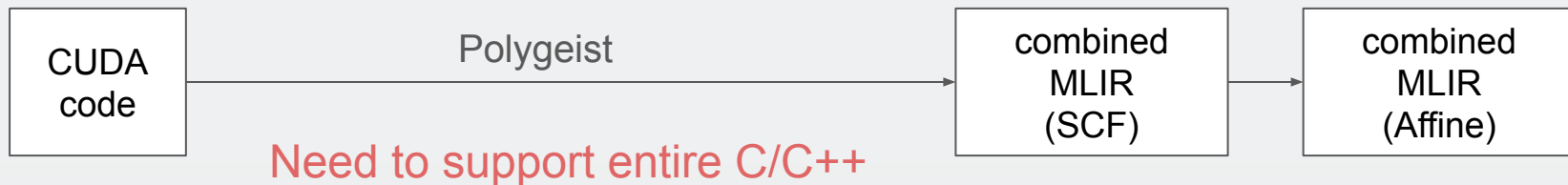
Our approach

# *Hydra*: The pipeline

We need to get the polyhedral representation of a kernel...

- The information is split
  - In the host code (the launch configuration: grid, block dims, shared mem size)
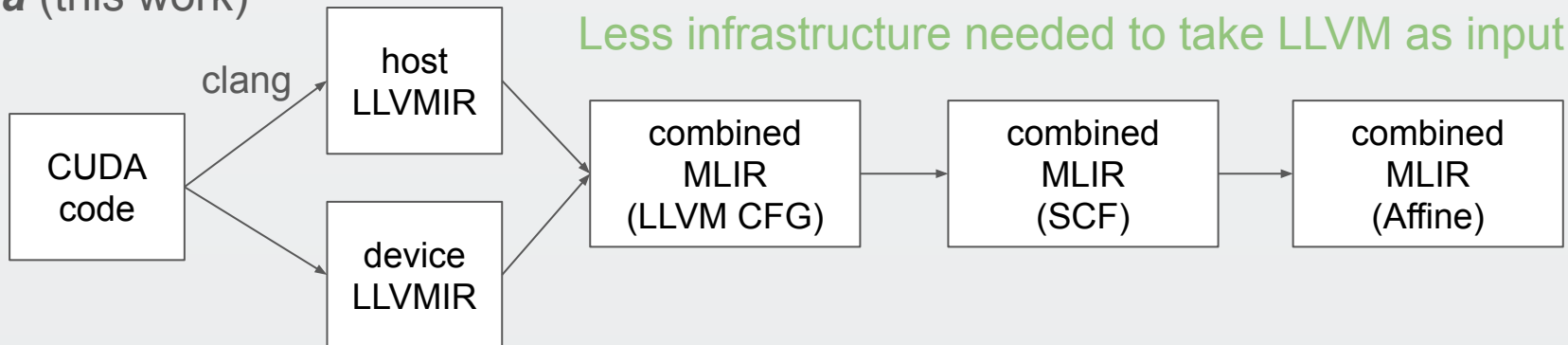  - In the device code (the actual kernel computation)
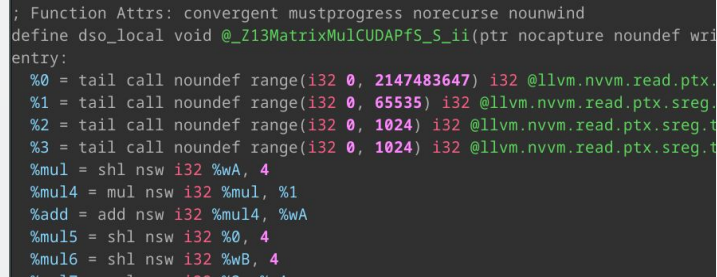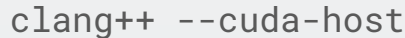
# Comparison to earlier work

# Compiling to LLVM-IR

```
#define BLOCK_SIZE 16

__global__ void MatrixMulCUDA(float *C, float *A,
                              float *B, int wA,
                              int wB) {
  int bx = blockIdx.x;
  int by = blockIdx.y;
  int tx = threadIdx.x;
  int ty = threadIdx.y;
  int aBegin = wA * BLOCK_SIZE * by;
  int aEnd   = aBegin + wA - 1;
  int aStep  = BLOCK_SIZE;
  int bBegin = BLOCK_SIZE * bx;
  int bStep  = BLOCK_SIZE * wB;
  float Csub = 0;
  for (int a = aBegin, b = bBegin;
       a <= aEnd;
       a += aStep, b += bStep) {
    __shared__ float As[BLOCK_SIZE][BLOCK_SIZE];
    __shared__ float Bs[BLOCK_SIZE][BLOCK_SIZE];
    As[ty][tx] = A[a + wA * ty + tx];
    Bs[ty][tx] = B[b + wB * ty + tx];
    __syncthreads();
    for (int k = 0; k < BLOCK_SIZE; ++k)
      Csub += As[ty][k] * Bs[k][tx];
    __syncthreads();
  }
  int c = wB * BLOCK_SIZE * by + BLOCK_SIZE * bx;
  C[c + wB * ty + tx] = Csub;
}

int MatrixMultiply(const dim3 &dimsA,
                   const dim3 &dimsB,
                   float *A, float *B, float *C) {
  // Setup execution parameters
  dim3 threads(BLOCK_SIZE, BLOCK_SIZE);
  dim3 grid(dimsB.x / threads.x, dimsA.y / threads.y);
  MatrixMulCUDA<<<grid, threads>>>(C, A, B, dimsA.x, dimsB.x);
}
```

`clang++ --cuda-device`

```
; Function Attrs: convergent mustprogress norecurse nounwind
define dso_local void @_Z13MatrixMulCUDAPfS_S_ii(ptr nocapture noundef wr
entry:
  %0 = tail call noundef range(i32 0, 2147483647) i32 @llvm.nvvm.read.ptx.
  %1 = tail call noundef range(i32 0, 65535) i32 @llvm.nvvm.read.ptx.sreg.
  %2 = tail call noundef range(i32 0, 1024) i32 @llvm.nvvm.read.ptx.sreg.t
  %3 = tail call noundef range(i32 0, 1024) i32 @llvm.nvvm.read.ptx.sreg.t
  %mul = shl nsw i32 %wA, 4
  %mul4 = mul nsw i32 %mul, %1
  %add = add nsw i32 %mul4, %wA
  %mul5 = shl nsw i32 %0, 4
  %mul6 = shl nsw i32 %wB, 4
```

`clang++ --cuda-host`

```
; Function Attrs: mustprogress uwtable
define dso_local void @_Z14MatrixMultiplyRK4dim3S1_PfS2_S2_(ptr nocapture noundef
entry:
  %0 = load i32, ptr %dimsB, align 4, !tbaa !6
  %div13 = lshr i32 %0, 4
  %y = getelementptr inbounds i8, ptr %dimsA, i64 4
  %1 = load i32, ptr %y, align 4, !tbaa !11
  %div314 = lshr i32 %1, 4
  %2 = load i32, ptr %dimsA, align 4, !tbaa !6
  tail call void @__mlir_launch_kernel__Z28__device_stub__MatrixMulCUDAPfS_S_ii(p
  ret void
}
```

Patched clang to preserve launch information better

# Merging the GPU modules

# The raising pipeline:
# finding control flow structure



Adapted from upstream numba-mlir

# Preserving loop information through LLVM transformations

**Standard pipeline**

Clang → LLVM Optimizer → Low level code

Loop unswitching, loop rotation, etc. enabled

**Our pipeline**

Clang → LLVM Pre-Optimizer → MLIR loop optimizations → LLVM Post-Optimizer → Low level code

Loop opts disabled

Loop unswitching, loop rotation, etc. enabled

# The raising pipeline:
# finding the polyhedral structure



affine loop bounds and access indices

# Representing the parallel structure of a GPU kernel



```
kernel @_Z10stencil_1dPKiPi(...) {
  A()
  barrier
  B()
  return
}
```

```
func @_Z10stencil_1dPKiPi(...) {
  affine.parallel %block = 0 to %grid_dim {
    %shared_mem = alloca(%size)
    affine.parallel %thread = 0 to %block_dim {
      A()
      affine.sync %thread
      B()
    }
  }
  return
}
```

[1]

[1] High-Performance GPU-to-CPU Transpilation and Optimization via High-Level Parallel Constructs
Authors: William S. Moses, Ivan R. Ivanov, Jens Domke, Toshio Endo, Johannes Doerfert, Oleksandr ZinenkoAuthors Info & Claims

# Propagating constants



Extremely important for subsequent analysis

# Representing synchronization

```
func @_Z10stencil_1dPKiPi(...) {
  affine.parallel %block = 0 to %grid_dim {
    %shared_mem = alloca(%size)
    affine.parallel %thread = 0 to %block_dim {
      A(%block, %thread)
      affine.sync %thread
      B(%block, %thread)
    }
  }
  return
}
```

What *is* `affine.sync` ?

# Transforming to pure polyhedral

No synchronisation allowed

# Parallel loop fission

```
parallel %i = 0 to %n
  A(%i)
  sync %i
  B(%i)
```

⇒

```
parallel %i = 0 to %n
  A(%i)
parallel %i = 0 to %n
  B(%i)
```

[1] Moses, W.S., Ivanov, I.R., Domke, J., Endo, T., Doerfert, J. and Zinenko, O., 2023, February. High-performance gpu-to-cpu transpilation and optimization via high-level parallel constructs. In Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming (pp. 119-134).
[2] Stratton, J.A., Stone, S.S. and Hwu, W.M.W., 2008, July. MCUDA: An efficient implementation of CUDA kernels for multi-core CPUs. In International Workshop on Languages and Compilers for Parallel Computing (pp. 16-30). Berlin, Heidelberg: Springer Berlin Heidelberg.

# Nested barriers?



```
func @foo(...) {
  parallel %i = 0 to 2 {
    for %j = 0 to 2 {
      A(%i, %j)
      sync %i
      B(%i, %j)
    }
  }
  return
}
```

```
func @foo(...) {
  for %j = 0 to 2
    parallel %i = 0 to 2
      A(%i, %j)
    parallel %i = 0 to 2
      B(%i, %j)
  return
}
```

[1] Moses, W.S., Ivanov, I.R., Domke, J., Endo, T., Doerfert, J. and Zinenko, O., 2023, February. High-performance gpu-to-cpu transpilation and optimization via high-level parallel constructs. In Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming (pp. 119-134).
[2] Stratton, J.A., Stone, S.S. and Hwu, W.M.W., 2008, July. MCUDA: An efficient implementation of CUDA kernels for multi-core CPUs. In International Workshop on Languages and Compilers for Parallel Computing (pp. 16-30). Berlin, Heidelberg: Springer Berlin Heidelberg.

# Let's look at matmul

```c
#define BLOCK_SIZE 16

__global__ void MatrixMulCUDA(float *C, float *A,
                              float *B, int wA,
                              int wB) {
  int bx = blockIdx.x;
  int by = blockIdx.y;
  int tx = threadIdx.x;
  int ty = threadIdx.y;
  int aBegin = wA * BLOCK_SIZE * by;
  int aEnd   = aBegin + wA - 1;
  int aStep  = BLOCK_SIZE;
  int bBegin = BLOCK_SIZE * bx;
  int bStep  = BLOCK_SIZE * wB;
  float Csub = 0;
  for (int a = aBegin, b = bBegin;
       a <= aEnd;
       a += aStep, b += bStep) {
    __shared__ float As[BLOCK_SIZE][BLOCK_SIZE];
    __shared__ float Bs[BLOCK_SIZE][BLOCK_SIZE];
    As[ty][tx] = A[a + wA * ty + tx];
    Bs[ty][tx] = B[b + wB * ty + tx];
    __syncthreads();
    for (int k = 0; k < BLOCK_SIZE; ++k)
      Csub += As[ty][k] * Bs[k][tx];
    __syncthreads();
  }
  int c = wB * BLOCK_SIZE * by + BLOCK_SIZE * bx;
  C[c + wB * ty + tx] = Csub;
}

int MatrixMultiply(const dim3 &dimsA,
                   const dim3 &dimsB,
                   float *A, float *B, float *C) {
  // Setup execution parameters
  dim3 threads(BLOCK_SIZE, BLOCK_SIZE);
  dim3 grid(dimsB.x / threads.x, dimsA.y / threads.y);
  MatrixMulCUDA<<<grid, threads>>>(C, A, B, dimsA.x, dimsB.x);
}
```

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (symbol(%12), symbol(%11), 1) {
  %As = memref.alloca() : memref<1024xi8, 3>
  %Bs = memref.alloca() : memref<1024xi8, 3>
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    %res = affine.for %i = 0 to #map()[%10] iter_args(%accum1 = %0) -> (f32) {
      %15 = affine.vector_load %A[...]
      affine.vector_store %15, %As[...]
      %16 = affine.vector_load %B[...]
      affine.vector_store %16, %Bs[...]
      "affine.barrier"(%tx, %ty, %tz)
      %added = affine.for %arg20 = 0 to 16 iter_args(%accum2 = %accum1) -> (f32) {
        %lA = affine.vector_load %As[...]
        %lB = affine.vector_load %Bs[...]
        %mul = %lA * %lB
        affine.yield %mul + %accum2: f32
      }
      "affine.barrier"(%tx, %ty, %tz)
      affine.yield %added : f32
    }
    affine.vector_store %res, %C[...]
  } {gpu.par.block}
} {gpu.par.grid}
```

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (symbol(%12), symbol(%11), 1) {
  %As = memref.alloca() : memref<1024xi8, 3>
  %Bs = memref.alloca() : memref<1024xi8, 3>
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    %res = affine.for %i = 0 to #map()[%10] iter_args(%accum1 = %0) -> (f32) {
      %15 = affine.vector_load %A[...]
      affine.vector_store %15, %As[...]
      %16 = affine.vector_load %B[...]
      affine.vector_store %16, %Bs[...]
      "affine.barrier"(%tx, %ty, %tz)
      %added = affine.for %arg20 = 0 to 16 iter_args(
        %lA = affine.vector_load %As[...]
        %lB = affine.vector_load %Bs[...]
        %mul = %lA * %lB
        affine.yield %mul + %accum2: f32
      }
      "affine.barrier"(%tx, %ty, %tz)
      affine.yield %added : f32
    }
    affine.vector_store %res, %C[...]
  } {gpu.par.block}
} {gpu.par.grid}
```

Analyzable using standard
polyhedral techniques

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (symbol(%12), symbol(%11), 1) {
  %As = memref.alloca() : memref<1024xi8, 3>
  %Bs = memref.alloca() : memref<1024xi8, 3>
  %registers = memref.alloca() : memref<16x16x1xf32, 16>
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    affine.store %0, %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>
  } {gpu.par.block}
  affine.for %arg15 = 0 to affine_map<()[s0] -> ((s0 - 1) floordiv 16 + 1)>()[%5] {
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %15 = affine.vector_load %A[...]
      affine.vector_store %15, %As[...]
      %16 = affine.vector_load %B[...]
      affine.vector_store %16, %Bs[...]
    } {gpu.par.block}
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %8 = affine.load %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
      %added = affine.for %k = 0 to 16 iter_args(%accum2 = %accum1) -> (f32) {
        %lA = affine.vector_load %As[...]
        %lB = affine.vector_load %Bs[...]
        %mul = %lA * %lB
        affine.yield %mul + %accum2: f32
      }
      affine.store %9, %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
    } {gpu.par.block}
  }
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    %8 = affine.load %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>
    affine.vector_store %8
  } {gpu.par.block}
} {gpu.par.grid}
```

# Example analysis of matmul

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (symbol(%12), symbol(%11), 1) {
  %As = memref.alloca() : memref<1024xi8, 3>
  %Bs = memref.alloca() : memref<1024xi8, 3>
  %registers = memref.alloca() : memref<16x16x1xf32, 16>
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {          INIT(%bx, %by)
    affine.store %0, %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>
  } {gpu.par.block}
  affine.for %arg15 = 0 to affine_map<()[s0] -> ((s0 - 1) floordiv 16 + 1)>()[%5] {
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %15 = affine.vector_load %A[...]
      affine.vector_store %15, %As[...]                                 COPY(%bx, %by, %k)
      %16 = affine.vector_load %B[...]
      affine.vector_store %16, %Bs[...]
    } {gpu.par.block}
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %8 = affine.load %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
      %added = affine.for %k = 0 to 16 iter_args(%accum2 = %accum1) -> (f32) {
        %lA = affine.vector_load %As[...]
        %lB = affine.vector_load %Bs[...]
        %mul = %lA * %lB                                                COMPUTE(%bx, %by, %k)
        affine.yield %mul + %accum2: f32
      }
      affine.store %9, %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
    } {gpu.par.block}
  }
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    %8 = affine.load %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>   EPILOGUE(%bx, %by)
    affine.vector_store %8
  } {gpu.par.block}
} {gpu.par.grid}
```

# Example analysis of matmul



```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (symbol(%12), symbol(%11), 1) {
  %As = memref.alloca() : memref<1024xi8, 3>
  %Bs = memref.alloca() : memref<1024xi8, 3>
  %registers = memref.alloca() : memref<16x16x1xf32, 16>
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    affine.store %0, %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>     INIT(%bx, %by)
  } {gpu.par.block}
  affine.for %arg15 = 0 to affine_map<()[s0] -> ((s0 - 1) floordiv 16 + 1)>()[%5] {
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %15 = affine.vector_load %A[...]
      affine.vector_store %15, %As[...]                                                COPY(%bx, %by, %k)
      %16 = affine.vector_load %B[...]
      affine.vector_store %16, %Bs[...]
    } {gpu.par.block}
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %8 = affine.load %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
      %added = affine.for %k = 0 to 16 iter_args(%accum2 = %accum1) -> (f32) {
        %lA = affine.vector_load %As[...]
        %lB = affine.vector_load %Bs[...]
        %mul = %lA * %lB                                                                COMPUTE(%bx, %by, %k)
        affine.yield %mul + %accum2: f32
      }
      affine.store %9, %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
    } {gpu.par.block}
  }
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    %8 = affine.load %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>       EPILOGUE(%bx, %by)
    affine.vector_store %8
  } {gpu.par.block}
} {gpu.par.grid}
```

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (%gx, %gy, 1) {
  INIT(%bx, %by)
  affine.for %k = 0 to %size {
    COPY(%bx, %by, %k)
    COMPUTE(%bx, %by, %k)
  }
  EPILOGUE(%bx, %by)
}
```

# Polyhedral scheduling

# Polyhedral model

```
for (int i = 0; i < N; i++) {
  for (int j = 0; j < i; j++) {
    S(i, j);
    R(i, j);
  }
}
```

```
domain

{ S(i, j) : 0 <= i < N,
            0 <= j < i }
{ R(i, j) : 0 <= i < N,
            0 <= j < i }
```

```
dependencies

{ S(si, sj) -> R(ri, rj) :
    si = ri, sj = rj }
```

Everything is represented using linear algebra

# Polyhedral scheduling

Integer Linear Programming Problem

**maximize** parallelism

**minimize** temporal distance between dependencies

**subject to** validity constraints (dependencies)

How do we optimize for async execution?

# Optimizing for async: Copy detection

Institute of
SCIENCE
TOKYO
RIKEN's
Programs for
Junior Scientists

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (symbol(%12), symbol(%11), 1) {
  %As = memref.alloca() : memref<1024xi8, 3>
  %Bs = memref.alloca() : memref<1024xi8, 3>
  %registers = memref.alloca() : memref<16x16x1xf32, 16>
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    affine.store %0, %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>
  } {gpu.par.block}
  affine.for %arg15 = 0 to affine_map<()[s0] -> ((s0 - 1) floordiv 16 + 1)>()[%5] {
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %15 = affine.vector_load %A[...]
      affine.vector_store %15, %As[...]                          COPY(%bx, %by, %k)
      %16 = affine.vector_load %B[...]
      affine.vector_store %16, %Bs[...]
    } {gpu.par.block}
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %8 = affine.load %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
      %added = affine.for %k = 0 to 16 iter_args(%accum2 = %accum1) -> (f32) {
        %lA = affine.vector_load %As[...]
        %lB = affine.vector_load %Bs[...]
        %mul = %lA * %lB
        affine.yield %mul + %accum2: f32
      }
      affine.store %9, %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
    } {gpu.par.block}
  }
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    %8 = affine.load %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>
    affine.vector_store %8
  } {gpu.par.block}
} {gpu.par.grid}
```
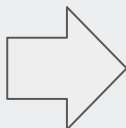
# Optimizing for async:
# Detecting dependencies on copies

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (symbol(%12), symbol(%11), 1) {
  %As = memref.alloca() : memref<1024xi8, 3>
  %Bs = memref.alloca() : memref<1024xi8, 3>
  %registers = memref.alloca() : memref<16x16x1xf32, 16>
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    affine.store %0, %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>
  } {gpu.par.block}
  affine.for %arg15 = 0 to affine_map<()[s0] -> ((s0 - 1) floordiv 16 + 1)>()[%5] {
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %15 = affine.vector_load %A[...]
      affine.vector_store %15, %As[...]
      %16 = affine.vector_load %B[...]
      affine.vector_store %16, %Bs[...]
    } {gpu.par.block}
    affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
      %8 = affine.load %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
      %added = affine.for %k = 0 to 16 iter_args(%accum2 = %accum1) ->
        %lA = affine.vector_load %As[...]
        %lB = affine.vector_load %Bs[...]
        %mul = %lA * %lB
        affine.yield %mul + %accum2: f32
      }
      affine.store %9, %registers[%arg16, %arg17, %arg18] : memref<16x16x1xf32, 16>
    } {gpu.par.block}
  }
  affine.parallel (%tx, %ty, %tz) = (0, 0, 0) to (16, 16, 1) {
    %8 = affine.load %registers[%arg15, %arg16, %arg17] : memref<16x16x1xf32, 16>
    affine.vector_store %8
  } {gpu.par.block}
} {gpu.par.grid}
```

COPY(%bx, %by, %k)

COMPUTE(%bx, %by, %k)

Async deps:

```
{ COPY(i, j, k) -> COMPUTE(i, j, k) }
```

# Optimizing for async: Optimization objective

```
affine.parallel (%bx, %by, %bz) =
    (0, 0, 0) to (%gx, %gy, 1) {
  INIT(%bx), %by)
  affine.for %k = 0 to %size {
    COPY(%bx, %by, %k)
    COMPUTE(%bx, %by, %k)
  }
  EPILOGUE(%bx, %by)
}
```
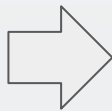
Async deps:

```
{ COPY(i, j, k) -> COMPUTE(i, j, k) }
```

Let's maximize the temporal distance of async deps

# Let's maximize the temporal distance of async deps

```
affine.parallel (%bx, %by, %bz) =
    (0, 0, 0) to (%gx, %gy, 1) {
  INIT(%bx), %by)
  affine.for %k = 0 to %size {
    COPY(%bx, %by, %k)
    COMPUTE(%bx, %by, %k)
  }
  EPILOGUE(%bx, %by)
}
```

```
affine.parallel (%bx, %by, %bz) =
    (0, 0, 0) to (%gx, %gy, 1) {
  INIT(%bx), %by)
  affine.for %k = 0 to %size {
    COPY(%bx, %by, %k)
  }
  affine.for %k = 0 to %size {
    COMPUTE(%bx, %by, %k)
  }
  EPILOGUE(%bx, %by)
}
```

We need an infinite amount of memory? Not very useful...

We need a way to constraint that

# Live range overlap constraint

Upper bound on the amount of memory used.

# Dependencies

```
affine.parallel (%bx, %by, %bz) = (0, 0, 0) to (%gx, %gy, 1) {
  INIT(%bx, %by)
  affine.for %k = 0 to %size {
    COPY(%bx, %by, %k)
    COMPUTE(%bx, %by, %k)
  }
  EPILOGUE(%bx, %by)
}
```

True dependencies (live ranges)

```
COPY(k) -> COMPUTE(k)
```
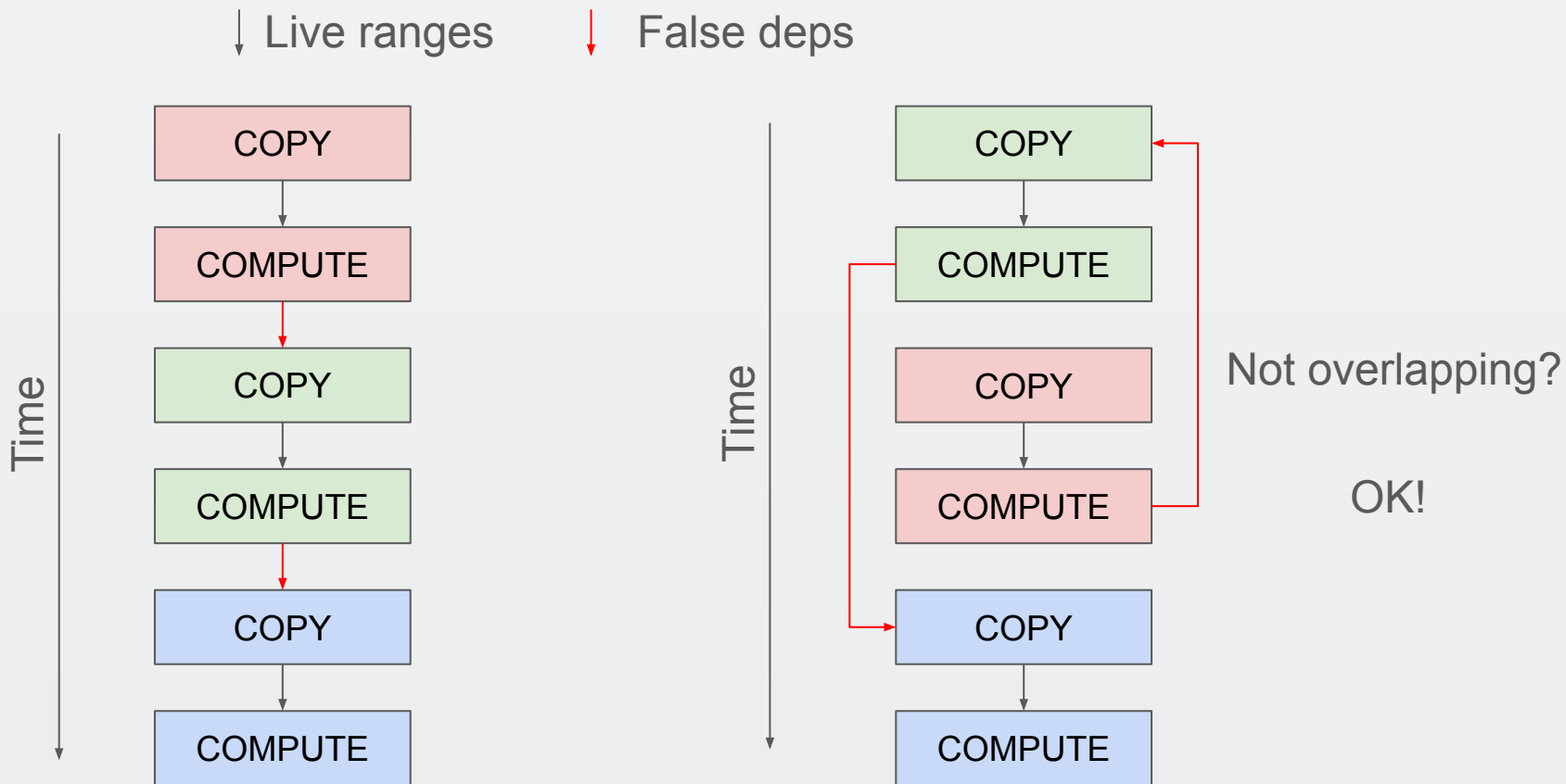
"False" dependencies

```
COMPUTE(k) -> COPY(k+1)
```

*Prior work* Live-range reordering: allow conditionally breaking the "false" dependencies.
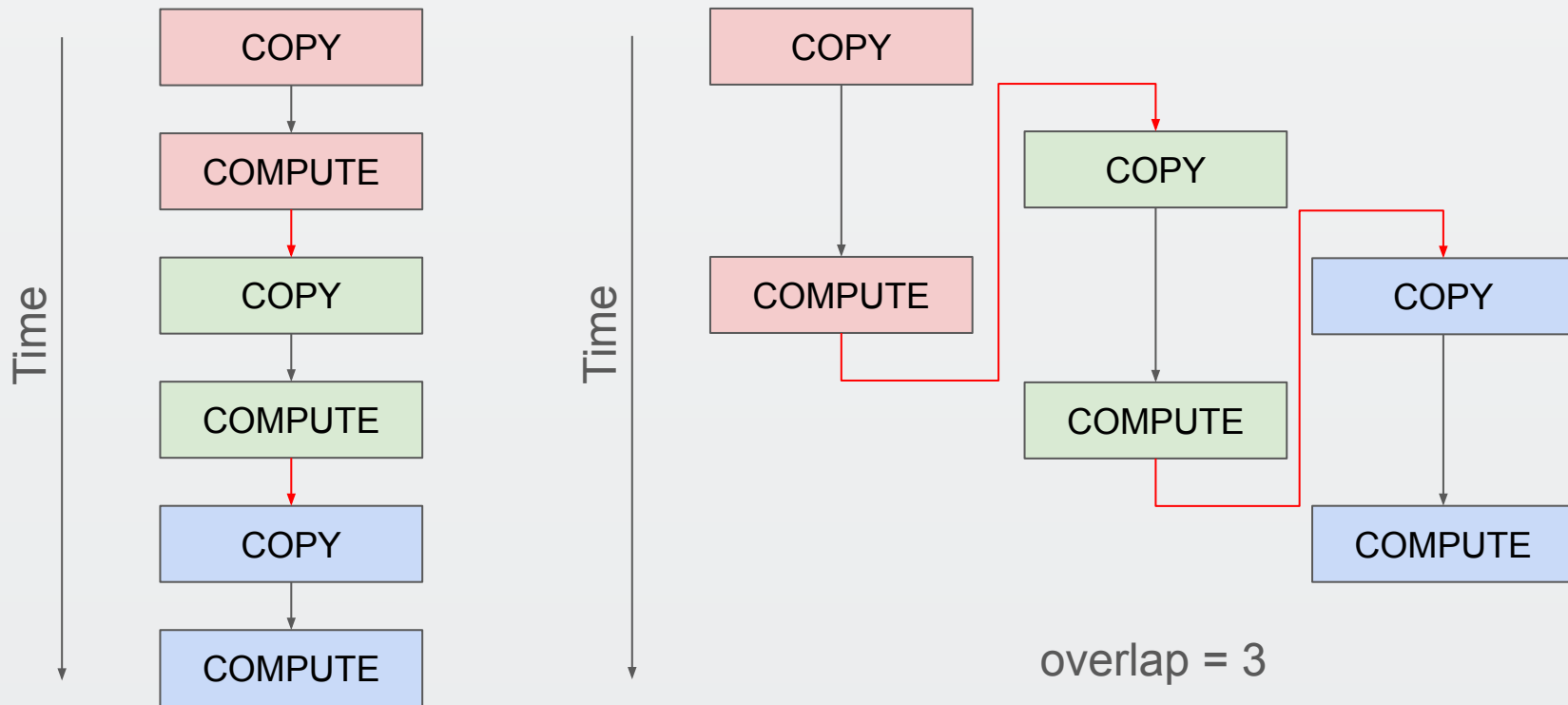
Verdoolaege, Sven, and Albert Cohen. "Live-range reordering." International Workshop on Polyhedral Compilation Techniques, Date: 2016/01/19-2016/01/19, Location: Prague, Czech Republic. 2016.

# *Prior work:* Live-range reordering

# Live-range overlapping

# Constraint in ILP

$$\sum_{a \in arrays} size(a) \times overlap(a) \leq AvailableSharedMemory$$

# Polyhedral scheduling

## ILP problem (prior work)

**maximize** parallelism
**minimize** temporal distance between dependencies


**subject to** validity constraints (dependencies)

## ILP problem (our work)

**maximize** parallelism
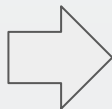**maximize** async dep distance
**minimize** temporal distance between dependencies


**subject to** validity constraints (dependencies)
**subject to** overlap constraints

# Final version

```
affine.parallel (%bx, %by, %bz) =
    (0, 0, 0) to (%gx, %gy, 1) {
  INIT(%bx), %by)
  affine.for %k = 0 to %size {
    COPY(%bx, %by, %k)
    COMPUTE(%bx, %by, %k)
  }
  EPILOGUE(%bx, %by)
}
```
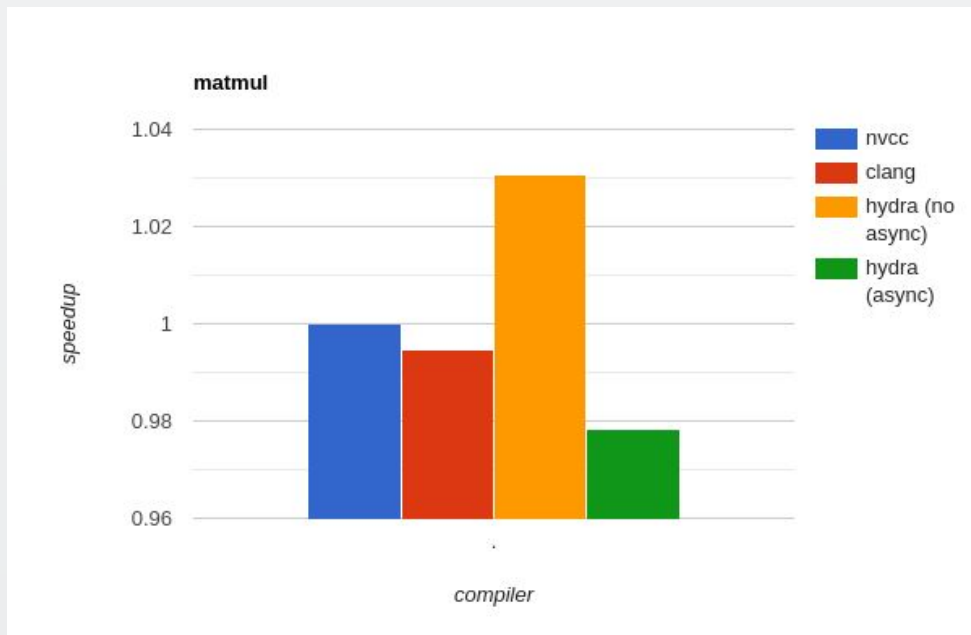
⇨

```
affine.parallel (%bx, %by, %bz) =
    (0, 0, 0) to (%gx, %gy, 1) {
  INIT(%bx), %by)
  COPY(%bx, %by, 0)
  COPY(%bx, %by, 1)
  affine.for %k = 0 to %size - 2 {
    COPY(%bx, %by, %k + 2)
    COMPUTE(%bx, %by, %k)
  }
  COMPUTE(%bx, %by, %size - 2)
  COMPUTE(%bx, %by, %size - 1)
  EPILOGUE(%bx, %by)
}
```

# Some early stage evaluation

The matrix multiplication example

# Conclusion

- LLVM -> MLIR Affine raising pipeline
- Polyhedral analysis of code with synchronisation
- Polyhedral scheduling with async

Questions?