YAKUP KORAY BUDANAZ

# SoftHier Progress Report February 10

SPCL

# Overiew of the Topics:

- **Programming Models**
- **DaCe + SoftHier**
- **Outlook For Next Weeks**

- **AI Accelerators and Programming Models**

# AI Accelerators – An Overview of Hardware Features

| Accelerator | Cache | Programmable NoC | Explicit Memory Levels | Dedicated MMU | Oversubscription to PEs | Double Precision Support |
|---|---|---|---|---|---|---|
| Nvidia A100 | ✓ | × | 2 | ✓ | ✓ | ✓ |
| Nvidia H100 | ✓ | ✓* | 2 | ✓ | ✓ | ✓ |
| Ascend 910B | × | × | 3 | ✓ | × | × |
| Graphcore IPU Mk2 | × | ✓ | 2 | × | ✓ | × |
| AMD VC2802** | × | ✓ | 2 | × | × | × |
| Microsoft Maia 100*** | × | ? | ? | ✓ | ? | × |
| Cerebras WSE-2 | × | ✓ | 2 | × | × | × |
| Intel Gaudi v3 | ✓ | × | 2 | ✓ | ✓ | × |
| AWS Trainium2 | × | × | 2 | ✓ | × | × |

Information collected from available white-papers, official tutorials and HotChips presentations.

4

# AI Accelerators

▪ TPU is not there because I found nothing on the low-level programming of TPUs.

▪ Only see how to write in Hight Level Operator (HLO) IR for Accelerated Linear Algebra (XLA) DL compiler used as backend for JAX and TensorFlow.
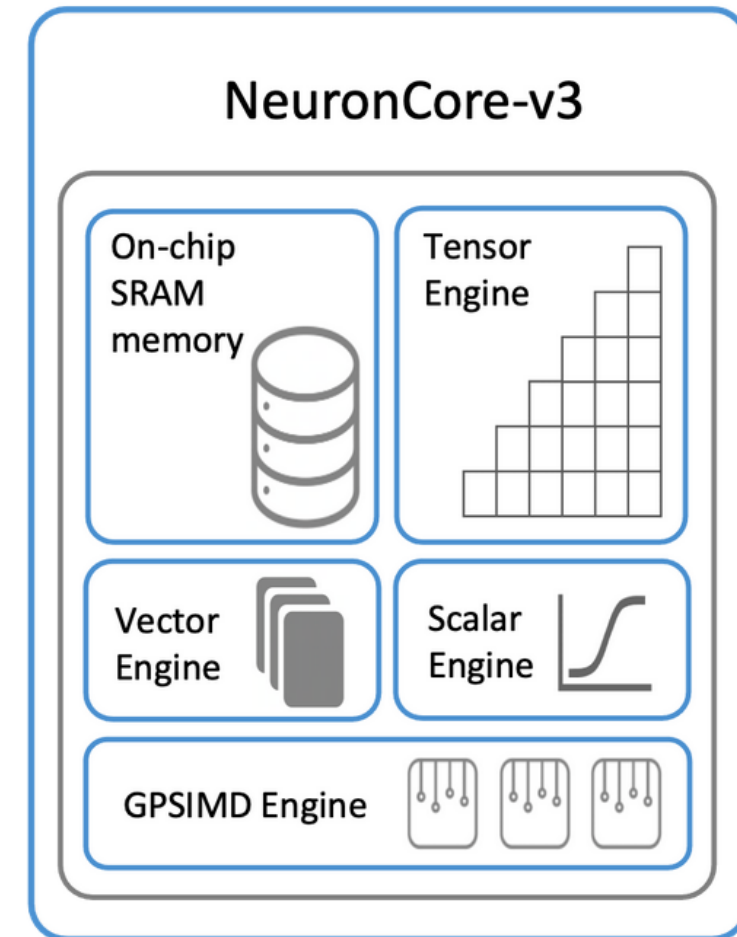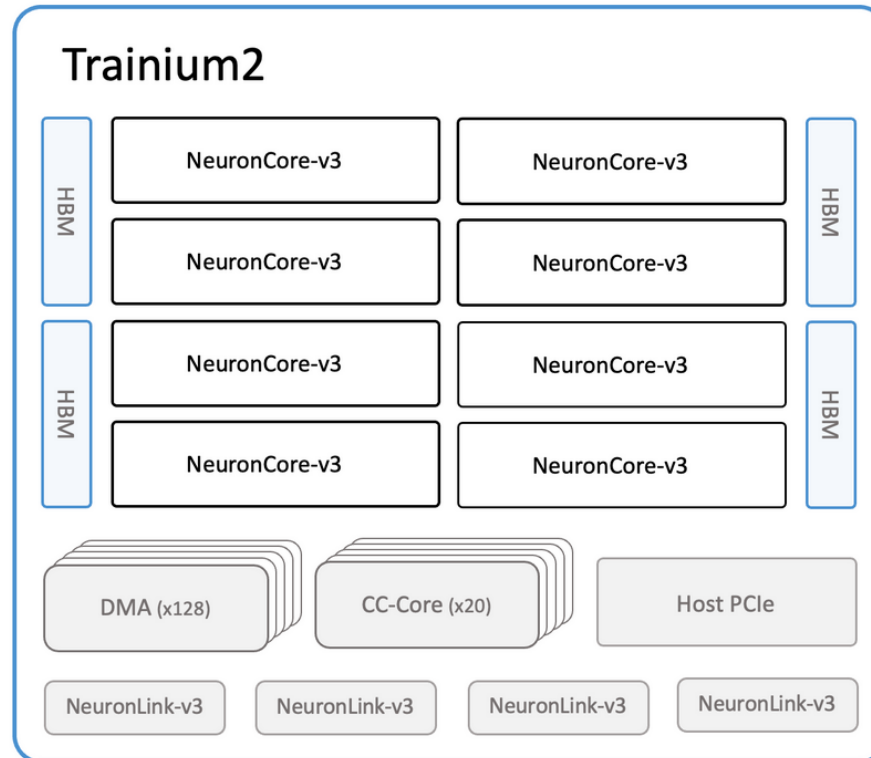
# Programming Models: High-Level IRs:

- A data-flow graph consisting of operations and input/outputs.

- Not really different a graph of consisting of ONNX operators. It is a variant of MLIR. Looks like this:

```
1    func @add_tensors(%arg0: tensor, %arg1: tensor) -> tensor {
2      %0 = "mhlo.add"(%arg0, %arg1) : (tensor, tensor) -> tensor
3      return %0 : tensor
4    }
```

**Abstracting Away Complexity – Method 1**: Provide a list of Operations on Tensors + Efficient Implementations for them
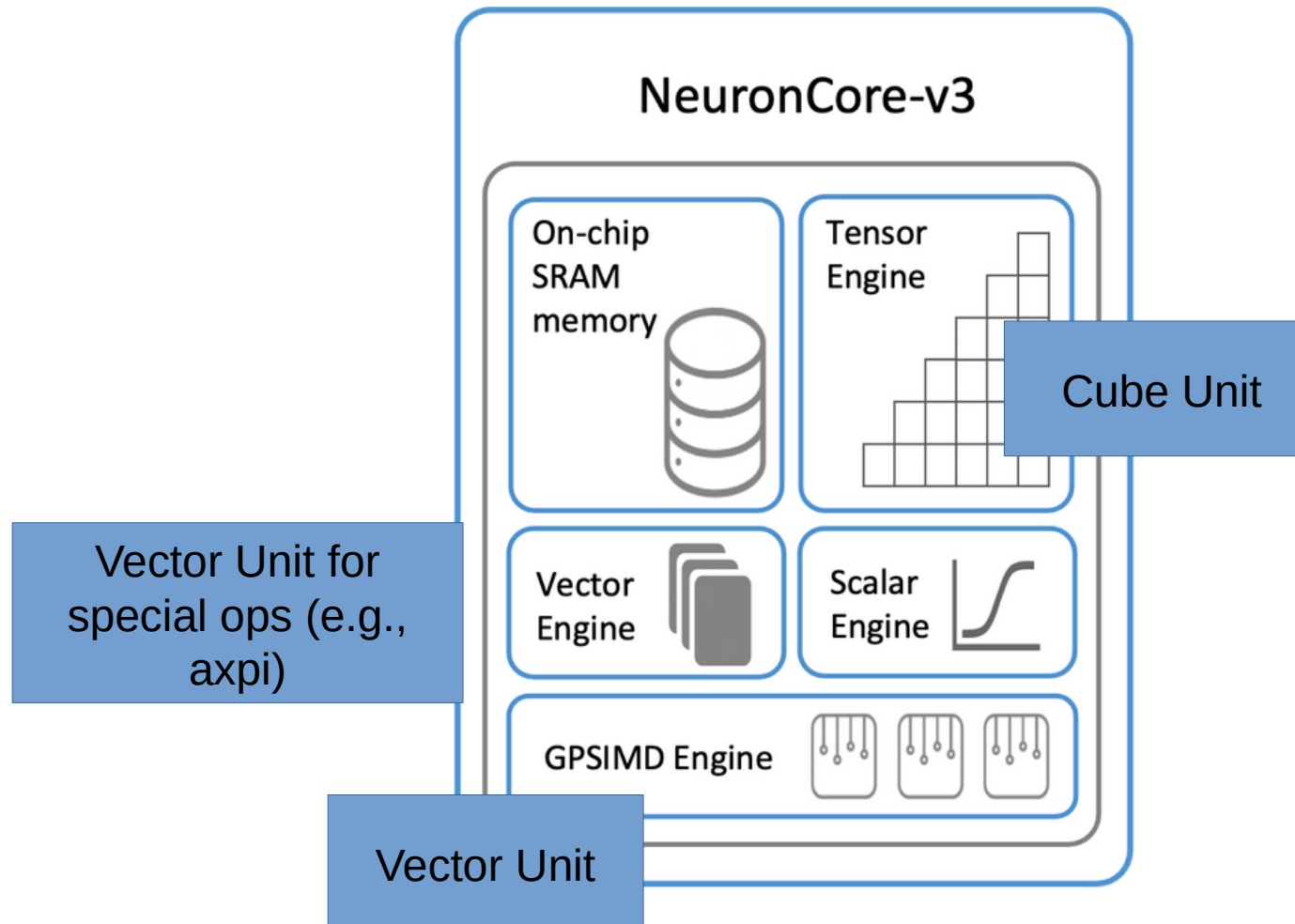
Google TPU only supports this, most AI accelerators support this
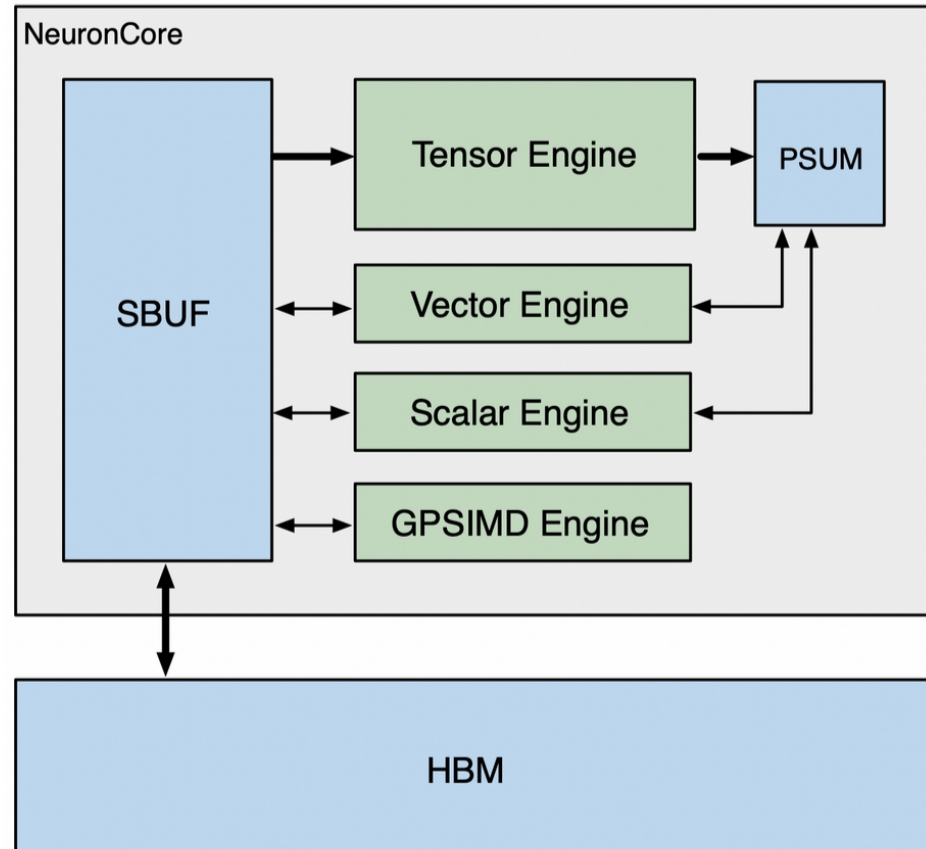
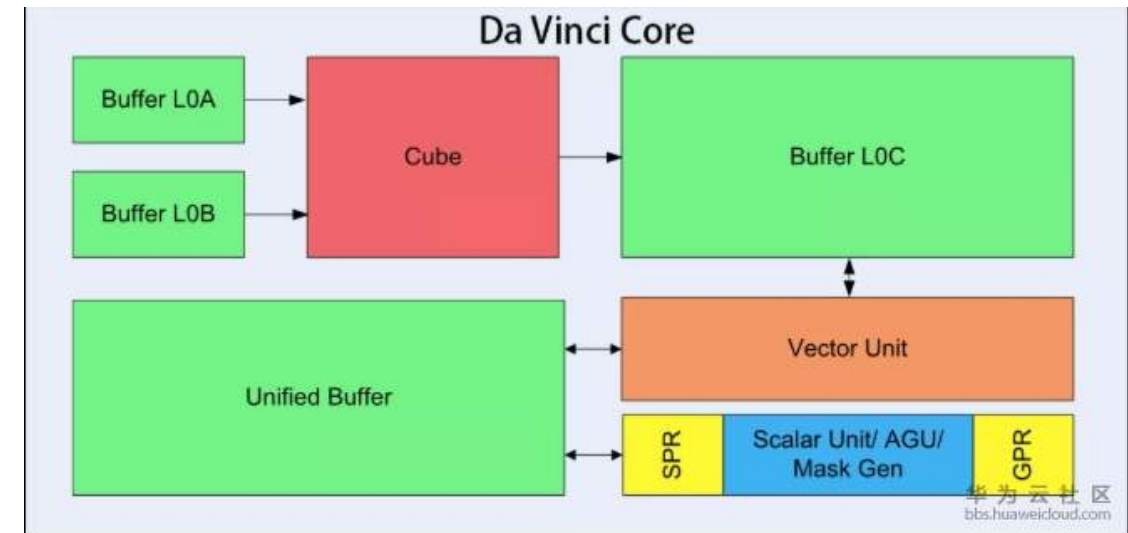# A very similar looking device - AWS Trainium2

# A very similar looking device - AWS Trainium2

# A very similar looking device - AWS Trainium2



https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/nki/programming_model.html#nki-programming-model



https://bbs.huaweicloud.com/blogs/106229

9

# A very similar looking device - AWS Trainium2

- Provides a tile-based programming interface similar to Triton named NKI

- No low-level programming interface for other computational units.

**Abstracting Away Complexity – Method 1**: Provide a list of Operations on Tensors + Efficient Implementations for them

# Programming Models

- Betting on providing efficient implementations for specific tensor operations, or betting on tile-based languages.

# Triton

- Load tiles and describe computation on tiles.

- Hopefully Triton back end + transformations will come-up with the most efficient way to move the memory

**Abstracting Away Complexity – Method 2**: Tile-Based language + a lot of memory movement optimizations
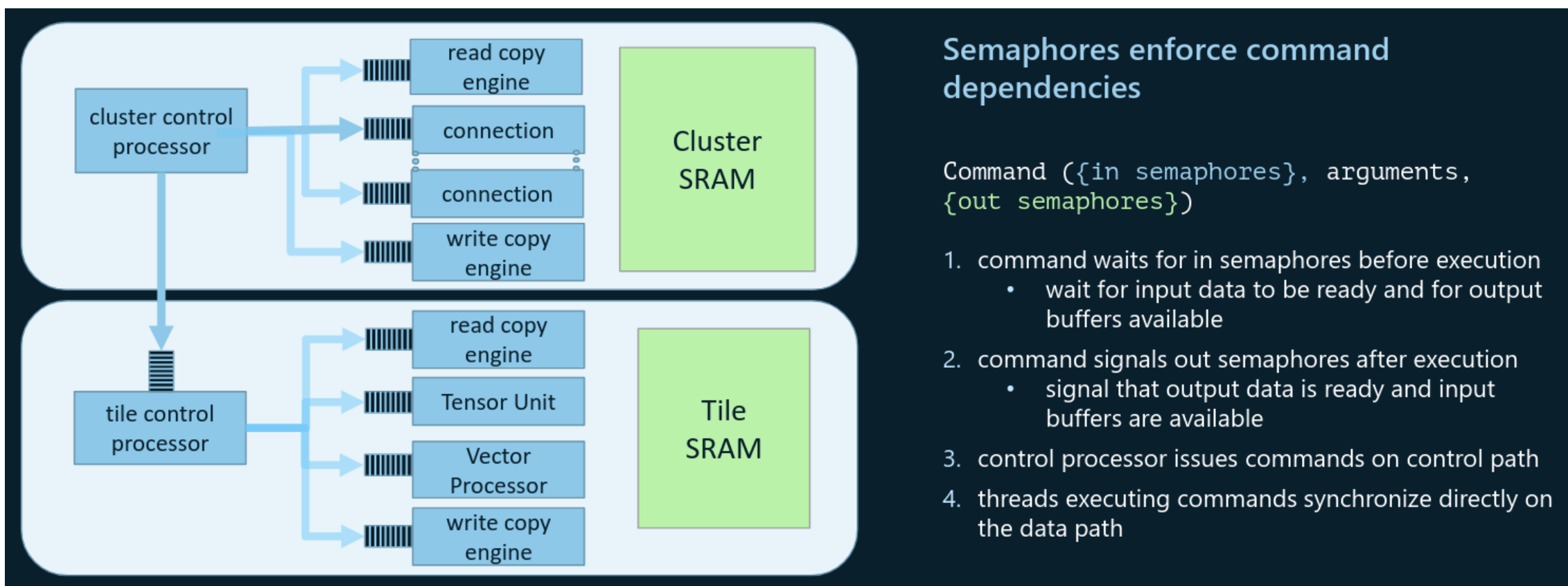
# Triton

- Tiled-based CUDA Wrapper (+embedded DSL in Python to generate Triton IR and into LLVM IR to ptx)

- Example stencil kernel.

```
1  A_bottom = triton.language.load(
2      A_ptr + (i + 2) * N + (j + 1), mask=(i < N - 2) & (j < N - 2), other=0.0
3  )
4
5  ...
6  B_new = 0.20 * (A_left + A_right + A_top + A_bottom + A_center)
7  triton.language.store(
8      B_ptr + (i + 1) * N + (j + 1), B_new, mask=(i < N - 2) & (j < N - 2)
9  )
10
```

# Microsoft Maia 100

- From the hotchips presentation. Looks like Maia will have a model similar to asynchronous tasks, dependencies will be satisfied using semaphores.

# Microsoft Maia 100

- It also plans to provide Triton as a programming model

**Abstracting Away Complexity – Method 3**: Asynchronous Tasks Based Language?

# Intel Gaudi 3

- It provides TPC language for low-level programming of the device

- Almost identical to CUDA with more intrinsics. It exposes warp (vector processor) more than the thread.

**Abstracting Away Complexity – Method 4**: Copy CUDA

# Graphcore:

- Provide poplar C++ API to create a dataflow graph consisting of tensor operations and input-output depdencies.

- Looks like SDFG with library-nodes only.

- The operation implementation uses a BSP-like model, the communication is implicitly generated using tiles. (It describes how computation is mapped to cores using a tile abstraction)

- Inter-PE communication is facilitated by the compiler (derived from tiling in form of graph partitioning).

- For more explicit control they provide assembly for their ISA.

- The bulk-synchronous parallel model of execution. This decomposes the runtime execution into three phases: local compute, global synchronisation, and data exchange.
- The graph representation of computations. The Poplar graph programming framework operates on a computational graph where vertices represent operations and edges represent the input and output data.

# Graphcore

**Abstracting Away Complexity – Method 5**: Graph-based language + Graph-partitioning

# AMD FPGAs

- Super low-level
- Super chaotic report
- Looks like assembly with some wrappers
- Inter-PE communication through an API similar to streams
- Provides a dataflow graph abstraction for operators

**Abstracting Away Complexity – Method 6**: Just Don't?

# Asynchronous Memory Pipelines

- Pipeline as a FIFO Queue following producer-consumer pattern (asynchronous).
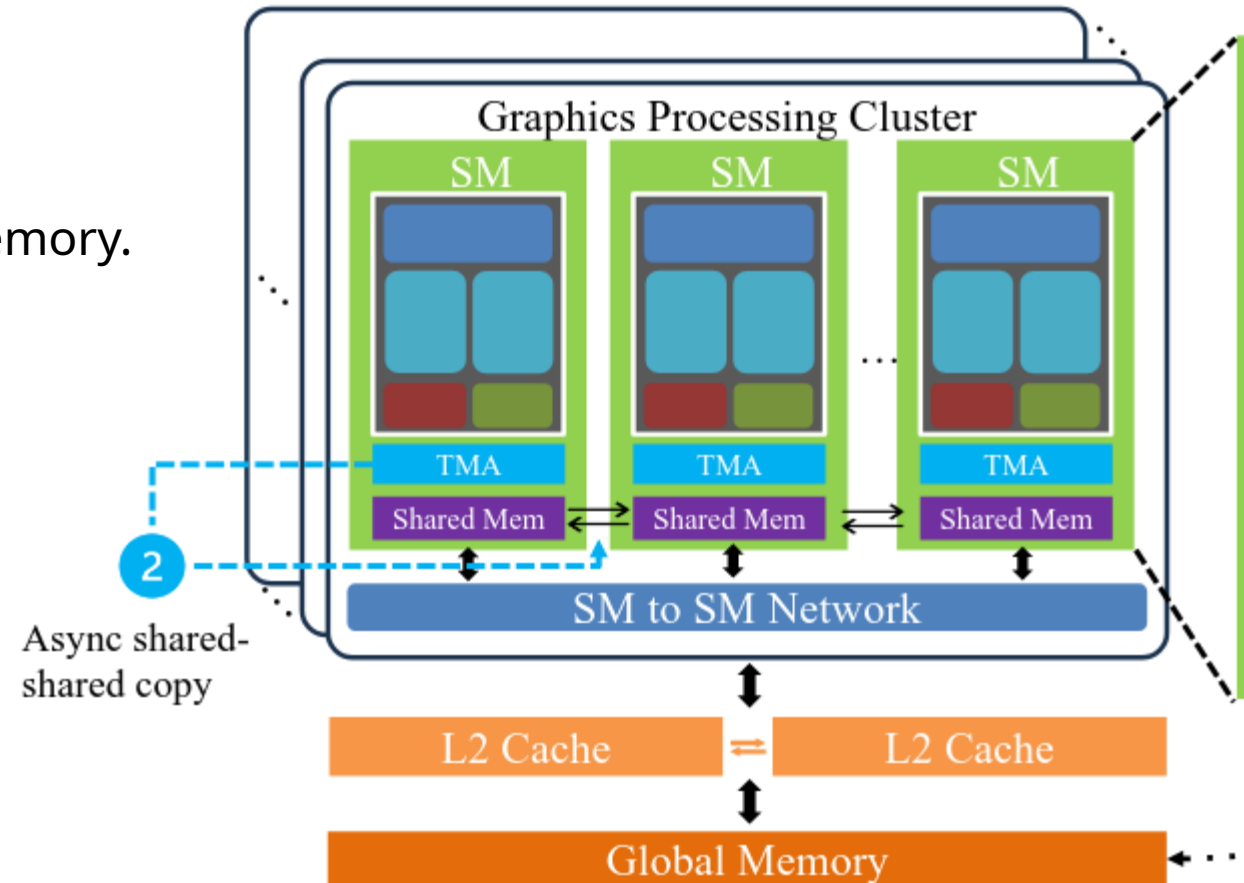
```
1      cuda::pipeline<cuda::thread_scope_thread> pipe = cuda::make_pipeline();
2
3      cuda::memcpy_async(&smem[stage][threadIdx.x], &src[idx], sizeof(int), pipe);
4          pipe.producer_commit();
5
6      for (...) {
7          cuda::pipeline_consumer_wait_prior<num_stages - 1>(pipe);
8          __syncthreads();
9
10         compute();
11
12         __syncthreads();
13         pipe.consumer_release();
14
15         pipe.producer_acquire();
16         cuda::memcpy_async(&smem[stage][threadIdx.x], &src[idx], sizeof(int), pipe);
17
18         pipe.producer_commit();
19
20         stage = (stage + 1) % num_stages;
21     }
22  }
```

https://developer.nvidia.com/blog/controlling-data-movement-to-boost-performance-on-ampere-architecture/

# Asynchronous DMA Engine:

- Tensor cores are now asynchronous too.

- One TMA per SM to asynchronously move memory.

| Arch | Precision | Programmability | Mode |
|---|---|---|---|
| Ampere | FP16,BF16, TF32,FP64, INT8,INT4,Binary | C: wmma PTX: mma, mma.sp | Sync |
| Ada | FP16,BF16,FP8, TF32,FP64, INT8,INT4,Binary | C: wmma PTX:mma, mma.sp | Sync |
| Hopper | FP16,BF16,FP8,TF32, FP64,INT8,Binary | C: wmma PTX: mma, mma.sp | Sync |
| | | PTX: wgmma, wgmma.sp | ASync |



Async shared-shared copy

https://arxiv.org/pdf/2501.12084

21

# Asynchronous DMA Engine:

- CuTe / CUTLASS provides convenient ptx wrappers for common optimizations such as memory swizzling, and calls to tensor cores / TMA.

**Abstracting Away Complexity – Method 7**: Providing Assembly Wrappers for Common Use-cases

https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html?highlight=TMA#using-tma-to-transfer-one-dimensional-arrays

# Asynchronous DMA Engine:

- TMA without high-level wrappers and only with ptx wrappers:

- Create, and copy tensor map

- Create many barriers, semaphores

- Async copies, set barriers, wait on barriers

- ...

**Abstracting Away Complexity – Method 7**: Even more convenient(?) assembly wrappers

https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html?highlight=TMA#using-tma-to-transfer-one-dimensional-arrays

# My Take on Programming Models:

- If I were to design a programming language I would only treat asynchronous tasks (defined over tiles) with input and output dependencies and construct a graph of asynchronous tasks.

- Graph IRs  are technically asynchronous. The dependencies are expressed as edges. The compiler or transformations need to find optimal memory movement pattern.
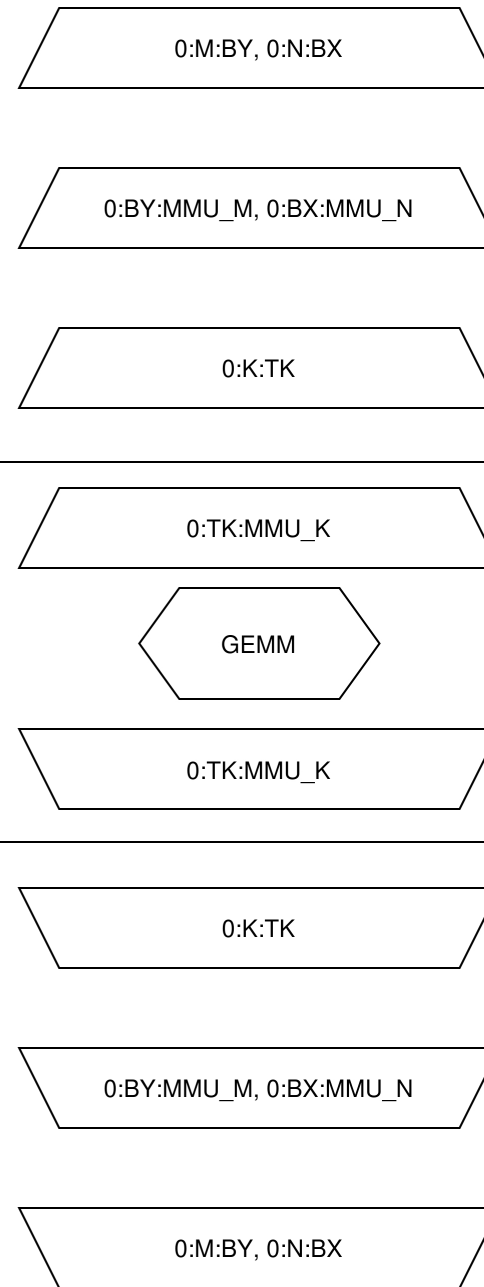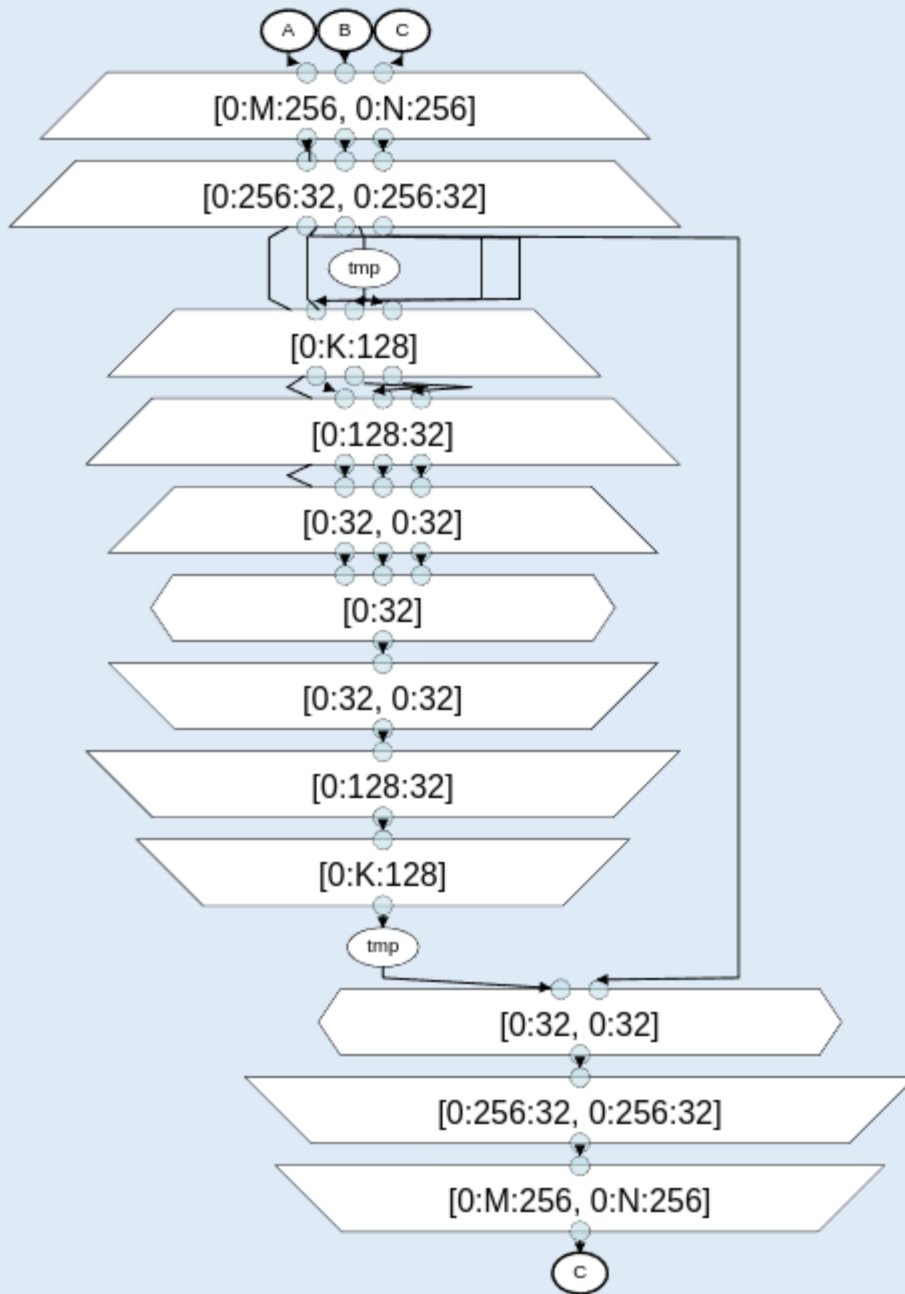
- **DaCe + SoftHier**

# Expressing NoCs in SDFGs:

- In my opinion, the best way of expressing NoC communication is using streams.

# Expressing NoCs in SDFGs:

- In my opinion, the best way of expressing NoC communication is using streams.

- Stream abstraction used in SDFGs is a FIFO queue with multiple senders and one receiver.

- How to transform SDFGs?

    - Designed a transformation that transforms a subgraph to use streams using a BSP-model-like description
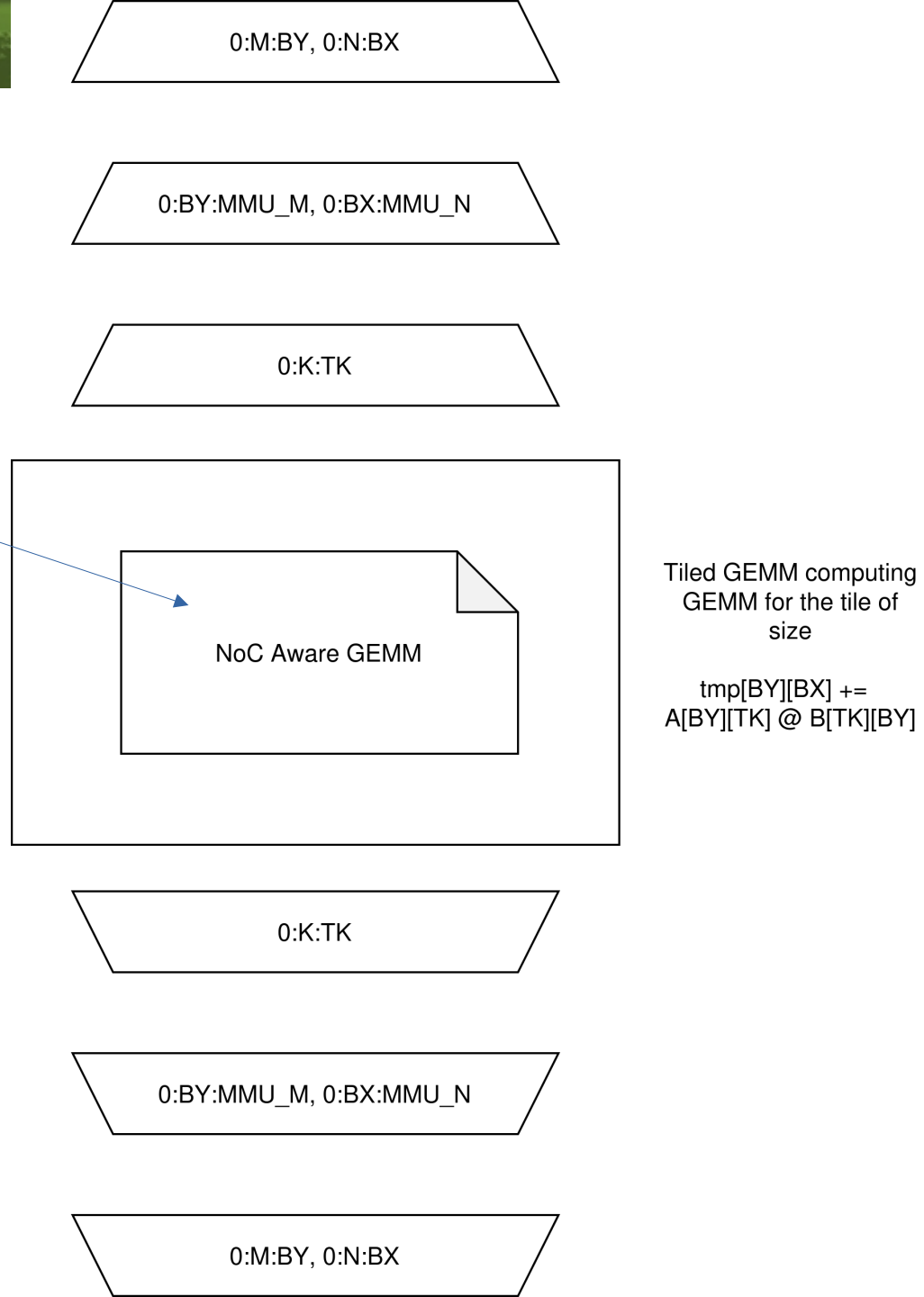
MapState



0:M:BY, 0:N:BX

0:BY:MMU_M, 0:BX:MMU_N

0:K:TK

0:TK:MMU_K

GEMM

0:TK:MMU_K

Tiled GEMM computing GEMM for the tile of size

tmp[BY][BX] += A[BY][TK] @ B[TK][BY]

0:K:TK

0:BY:MMU_M, 0:BX:MMU_N

0:M:BY, 0:N:BX

28

0:M:BY, 0:N:BX

0:BY:MMU_M, 0:BX:MMU_N

0:K:TK

Input is a BSP-model based description of the algorithm running inside.

1. Initial data dist (which core owns data)
2. Communication scheme between cores
3. Synchronization (if necessary in the architecture)

NoC Aware GEMM

Tiled GEMM computing GEMM for the tile of size

tmp[BY][BX] += A[BY][TK] @ B[TK][BY]

0:K:TK

0:BY:MMU_M, 0:BX:MMU_N

0:M:BY, 0:N:BX

1. Initial distribution:
sA and sB are 2 dimensions streams, both of size [N][N], N=NX=NY

A[i * tileSize][((i+j) % N) * tileSize] → sA[i][(i+j)%N][:][:]
B[((i+j) % N) * tileSize][j * tileSize] → sB[(i+j)%N][j][:][:]

2. Compute:
localA = sA.pop(), localB = sB.pop()
tmp += localA * localB

3. Communication

LocalA[:][:] → sA[i][(j+N-1)%N][:][:]
LocalB[:][:] → sB[(i+N-1)%N][j][:][:]

4. Synchronization
Yes|No depending on the architecture

y=0:M:BY, x=0:N:BX

0:BY:MMU_M, 0:BX:MMU_N

0:K:TK

sA          sB

0:N:1

GEMM

sA          sB

0:N:1

0:K:TK

0:BY:MMU_M, 0:BX:MMU_N

0:M:BY, 0:N:BX

Assume
NumCoreY == NumCoreX
and  BY == TK == BX

Tiled GEMM computing
GEMM for the tile of
size

tmp[BY][BX] +=
A[BY][TK] @ B[TK][BY]

No issues in part 1

2. Compute:
localA = sA.pop(), localB = sB.pop()
tmp += localA * localB

3. Communication

LocalA[:][:] → sA[i][(j+N-1)%N][:][:]
LocalB[:][:] → sB[(i+N-1)%N][j][:][:]

4. Synchronization
Yes|No depending on the architecture

y=0:M:BY, x=0:N:BX

0:BY:MMU_M, 0:BX:MMU_N

0:K:TK

Assume NumCoreY == NumCoreX and BY == TK == BX

sA          sB

For GEMM, the steps 2 and 3 can be concurrent. For Stencil, the steps 2 and 3 can't be concurrent. The user

0:N:1

GEMM

Tiled GEMM computing GEMM for the tile of size

tmp[BY][BX] +=
A[BY][TK] @ B[TK][BY]

sA          sB

0:N:1

0:K:TK

0:BY:MMU_M, 0:BX:MMU_N

0:M:BY, 0:N:BX

31

zürich

No issues in part 1

2. Compute:
localA = sA.pop(), localB = sB.pop()
tmp += localA * localB

3. Communication

LocalA[:][:] → sA[i][(j+N-1)%N][:][:]
LocalB[:][:] → sB[(i+N-1)%N][j][:][:]

No issues in part 4

y=0:M:BY, x=0:N:BX

0:BY:MMU_M, 0:BX:MMU_N

0:K:TK

sA          sB

0:N:1

GEMM

Assume
NumCoreY == NumCoreX
and  BY == TK == BX

Tiled GEMM computing
GEMM for the tile of
size

tmp[BY][BX] +=
A[BY][TK] @ B[TK][BY]

What if need to transfer only a subset?

The description should support an
update. (In stencil, for example, only
communicate the halo ranges, update
the local boundary of the local data)

sA          sB

0:N:1

0:K:TK

0:BY:MMU_M, 0:BX:MMU_N

0:M:BY, 0:N:BX

1. Initial distribution:

sA and sB are 2 dimensions streams, both of size [N][N], N=NX=NY,

A[i * tileSize][((i+j) % N) * tileSize] → sA[i][(i+j)%N][:][:]

B[((i+j) % N) * tileSize][j * tileSize] → sB[(i+j)%N][j][:][:]
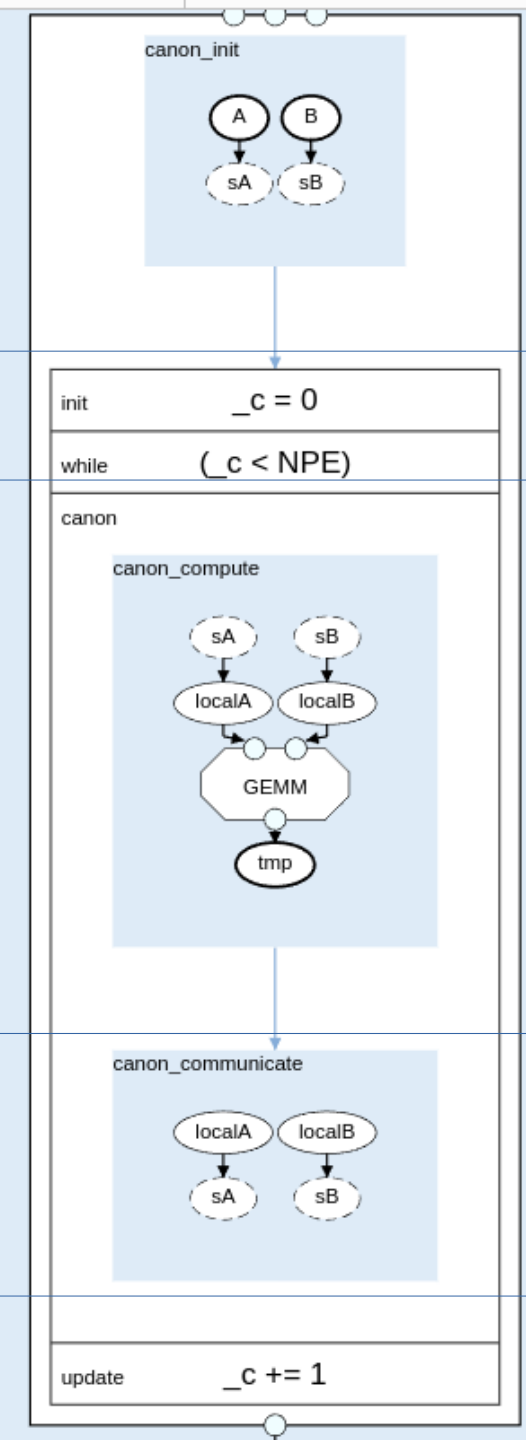
2. Loop Range

2. Compute:

localA = sA.pop(), localB = sB.pop()

tmp += localA * localB

3. Communication

LocalA[:][:] → sA[i][(j+N-1)%N][:][:]

LocalB[:][:] → sB[(i+N-1)%N][j][:][:]

4. Synchronization

Generated as a library node on a 3$^{rd}$ state or as a loop property



canon_init

A    B

sA    sB

init          _c = 0

while        ( _c < NPE)

canon

canon_compute

sA    sB

localA    localB

GEMM

tmp

canon_communicate

localA    localB

sA    sB

update        _c += 1

33

- Scalar → MMU Pass (Pattern detection + Reading Tasklet ASTs), working for couple of patterns

- Scalar → Vector Unit Pass (Pattern detection + Reading Tasklet ASTs), works for couple of patterns, need to extend patterns later as necessary

- GEMM Input Type information is added to access nodes which is used by the controller to apply the correct ExplicitMemoryMove transformation

- Passing "input purpose" to SDFGs using access nodes manually (done)

- Automating this information by extending Scalar → MMU Pass (mostly done)

- Pass to insert memory movement between different computational units (under work)

# Outlook / TODOs:

- Work further on transformations

- Extend cube unit parts of the code-gen

# My Take on Programming Models:

Sketched how it could look like. Many considerations are necessary. Message buffers, supporting out-of-order messages, low-level implementation of completions of one-sided communication.

```
1  __kernel__ vadd<N>(in A[N*1024], in B[N*1024], inout C[N*1024]){
2      decl id = PEs<N>
3      decl TS = 1024
4      a1 = reserve(localB[1024], TCM)
5      a2 = reserve(localA[1024], TCM)
6      a3 = reserve(localC[1024], TCM)
7      l1 = a1.then(load(in:A[i*TS:(i+1)*TS], out:localA[:]))
8      l2 = a2.then(load(in:B[i*TS:(i+1)*TS], out:localB[:]))
9      t1 = join(l1, l2).then(add(in:localA[:], in:localB[:], out:localC[:]))
10     s1, r1 = t1.then(send(to:(i+N-1)%N, in:localC[:], out:localC[:]))
11     c1 = r1.then(store(in:localC[:], out:C[i*TS:(i+1)*TS]))
12     return c1
13 }
14
15 ...
16
17 vadd_kernel = vadd<N>(A, B, C)
18 vadd_kernel.wait()
```