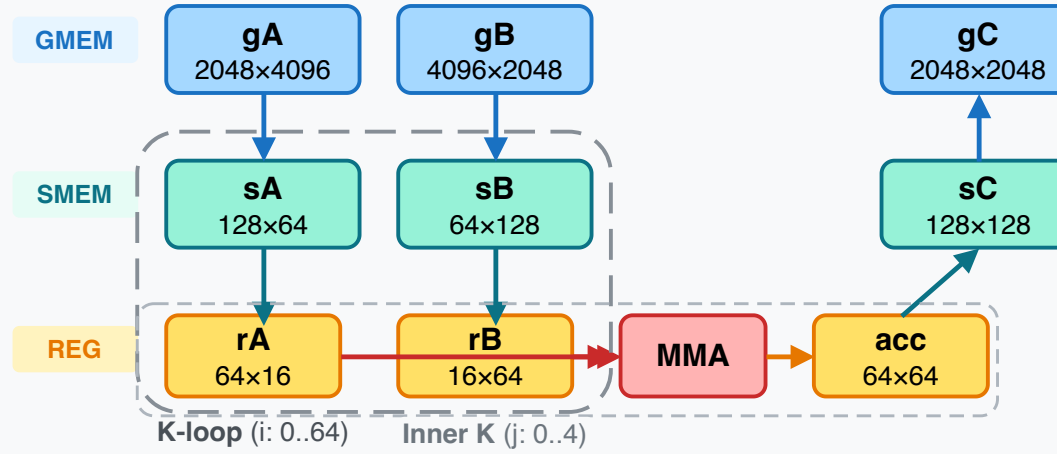# (a) AffineGraph IR (Python DSL)

```python
gA = Buffer(space=Global, dim=[2048, 4096])
sA = Buffer(space=Shared, dim=[128, 64])
rA = Buffer(space=Reg, dim=[64, 16])
acc = Buffer(space=Reg, dim=[64, 64])
i = LoopVar(name='i', domain=(0..64)) # Outer K
j = LoopVar(name='j', domain=(0..4))  # Inner K
```

```python
Map_G2S_A = AffineMap(src=gA, dst=sA, ctx=[bi], loop=[i],
          expr="sA[m,k] = gA[bi*128+m, i*64+k]")
Map_S2R_A = AffineMap(src=sA, dst=rA, ctx=[wi], loop=[j],
          expr="rA[m',k'] = sA[wi*64+m, j*16+k']")
```
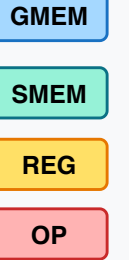
```python
RegGraph = Graph(nodes=[rA, rB, acc, GEMM_op], ...)
S2R_Block = Block(graph=RegGraph, loop=[j], maps=[...])
SharedGraph= Graph(nodes=[sA, sB, S2R_Block], ...)
G2S_Block = Block(graph=SharedGraph, loop=[i], maps=[...])
```

# (b) Lowered Dataflow Graph

GMEM
- gA 2048×4096
- gB 4096×2048
- gC 2048×2048

SMEM
- sA 128×64
- sB 64×128
- sC 128×128

REG
- rA 64×16
- rB 16×64
- MMA
- acc 64×64

K-loop (i: 0..64)    Inner K (j: 0..4)

GMEM
SMEM
REG
OP

### Map_G2S_A

$$\begin{bmatrix} m \\ k \end{bmatrix}_{sA} = \begin{bmatrix} 1 & 0 & 128 & 0 \\ 0 & 1 & 0 & 64 \end{bmatrix} \begin{bmatrix} m \\ k \\ bi \\ i \end{bmatrix}$$

### Map_S2R_A

$$\begin{bmatrix} m' \\ k' \end{bmatrix}_{rA} = \begin{bmatrix} 1 & 0 & 64 & 0 \\ 0 & 1 & 0 & 16 \end{bmatrix} \begin{bmatrix} m' \\ k' \\ wi \\ j \end{bmatrix}$$

# (c) GPU Architecture Mapping (4 Warps, 2×2 Layout)

**Global Memory (Matrix A)**

other CTAs
M=2048
128×64
CTA row 128 rows
blockIdx.x
K-loop (i: 0..64)
other CTAs
K=4096

cp.async 128B

**Shared Memory (128×64)**
**WarpLayout (2×2)**

M=128
- Warp0 64×32
- Warp1 64×32
- Warp2 64×32
- Warp3 64×32

K=64

Swizzle<2, 3, 3>

ldmatrix 16×16

**Register (Per Warp)**
**WarpTile (64×32)**

64
- 16×16  16×16
- 16×16  16×16
- 16×16  16×16
- 16×16  16×16

32

BaseTile 32 threads × 8 elem

4×2 ldmatrix

mma.sync

**Tensor Core MMA**

**mma.sync.m16n8k16**

A 16×16 × B 16×8 + C 16×8 = D 16×8

FP16 × FP16 + FP32 → FP32

**Per Warp (64×64)**
- 8 ldmatrix for A (64×32)
- 8 ldmatrix for B (32×64)
- 32 mma.sync calls
- Accumulate 64×64
**Inner K (j: 0..4)**
64-col / 16 = 4 iters

**Pipeline**

CTA → G2S → Swizzle → S2R → MMA → Acc → Store