

```

A = t.placeholder((1024, 1024))
B = t.placeholder((1024, 1024))
k = t.reduce_axis((0, 1024))
C = t.compute((1024, 1024), lambda y, x:
               t.sum(A[k, y] * B[k, x], axis=k))
s = t.create_schedule(C.op)

```

```

for y in range(1024):
    for x in range(1024):
        C[y][x] = 0
        for k in range(1024):
            C[y][x] += A[k][y] * B[k][x]

```

+ Loop Tiling

```

yo, xo, ko, yi, xi, ki = s[C].tile(y, x, k, 8, 8, 8)

```

```

for yo in range(128):
    for xo in range(128):
        C[yo*8:yo*8+8][xo*8:xo*8+8] = 0
        for ko in range(128):
            for yi in range(8):
                for xi in range(8):
                    for ki in range(8):
                        C[yo*8+yi][xo*8+xi] +=
                            A[ko*8+ki][yo*8+yi] * B[ko*8+ki][xo*8+xi]

```

+ Cache Data on Accelerator Special Buffer

```

CL = s.cache_write(C, vta.acc_buffer)
AL = s.cache_read(A, vta.inp_buffer)
# additional schedule steps omitted ...

```

+ Map to Accelerator Tensor Instructions

```

s[CL].tensorize(yi, vta.gemm8x8)

```

```

inp_buffer AL[8][8], BL[8][8]
acc_buffer CL[8][8]
for yo in range(128):
    for xo in range(128):
        VTAPushResetOp(CL)
        for ko in range(128):
            VTALoadBuffer2D(AL, A[ko*8:ko*8+8][yo*8:yo*8+8])
            VTALoadBuffer2D(BL, B[ko*8:ko*8+8][xo*8:xo*8+8])
            VTAPushGEMMOp(CL, AL, BL)
            VTASToreBuffer2D(C[yo*8:yo*8+8,xo*8:xo*8+8], CL)

```

