# Outer-product CPU Instruction for NPU-level ML Throughput

Houxiang (hj14), Husnain (mubarik3) and Nam (hn5)

# Motivation

- DNN accelerators usually accelerate **matrix-multiplications** (MM)

- Accelerators are hard to program --> poor optimization

- Extend CPU ISA to support **outer-product**
  - Compute MM as **sum of outer-products**
  - **High compute density** with low memory cost

# Motivation: MM as Sum of Outer-products

Matrix Multiplication: C = AXB

Traditional Way: $C_{ij} = \sum_k A_{ik} * B_{kj} = A_i^\top B_j$ (Inner-products)

**$C = \sum_i (A^\top)_i \otimes (B_i) = \sum_i (A^\top)_i (B_i)^\top$** (sum of outer products)

$$
\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}
\times
[b_1 \ b_2 \cdots b_n]
=
\begin{bmatrix}
a_1 b_1 & a_1 b_2 & \cdots & a_1 b_n \\
a_2 b_1 & a_2 b_2 & & a_2 b_n \\
& \vdots & \ddots & \vdots \\
a_n b_1 & a_n b_2 & \cdots & a_n b_n
\end{bmatrix}
$$

# ISA Specification (n=16, ARM based)

Usage: opa C, a, b ( C = a $\otimes$ b + C )

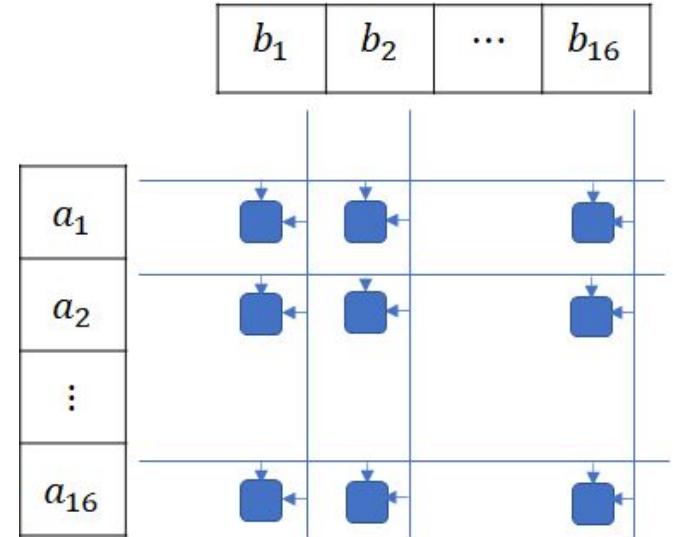a, b: vector registers of 16 INT8 values (128-bit) (v0~v16)

C: **tile register** of 16X16 INT32 accumulators (8192-bit) (m0~m3)

256 MACs per cycle -> 5 TOPs @ 2.5 GHz 4-core ( > 4 TOPs of Edge TPU)


tld/tst C, a, #n (s) : shifted load/store n-th row of the tile C from/to the vector a
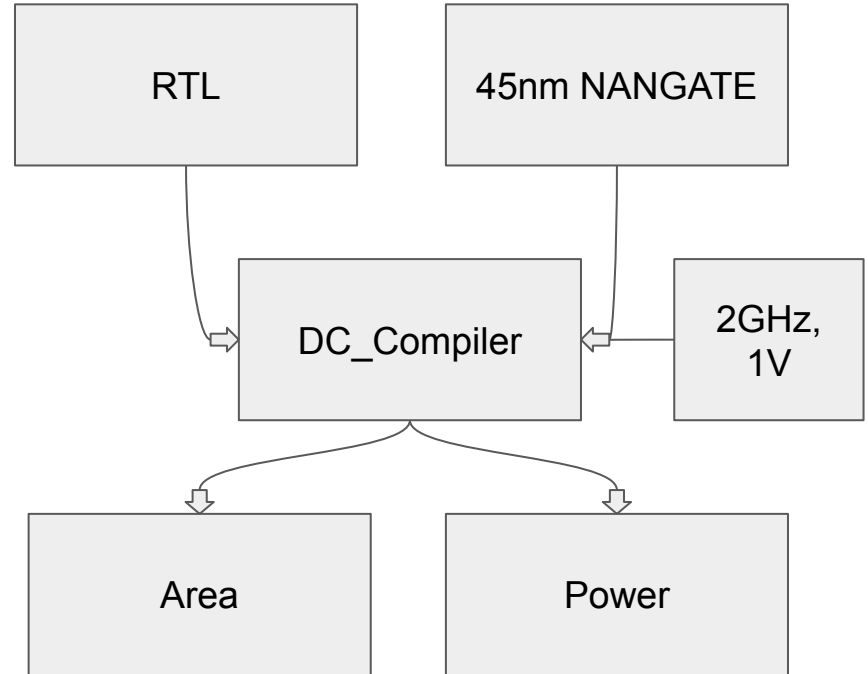
# Architecture Details (HW)

- Functional Element
  - 256 MAC units (int8)
  - Accumulated to the tile register


- Tile Accumulator Register
  - 16X16 32 bit accumulator register
  - Row load/store unit from/to 128-bit vector

# Evaluation Method - Power/Area

- Handwritten, synthesizable RTL
- Synopsys DC_Compiler
- NanGate 45nm Library
- Area and Power reports

# Current Results - Power/Area

- Technology node: 45 nm
- Cortex A72 Area: 50 mm2
- Cortex A72 Power: 4.76 W

|  | Frequency | Area mm2 | Power mW | Area overhead vs A72 | Power overhead vs A72 |
|---|---|---|---|---|---|
| macBrick | 2GHz | 0.23 | 458 | 0.45% | 9.65% |

# SW Usage

## Inner-product

```
for A_t, B_t, C_t in tile(A,B,C):
  tmp = [0, ..., 0]
  for i in 0...16:
    for j in 0...16:
      tmp = fma16(A_t[i][:], B_t[:][j])
      C_t[i][j] = C_t[i][j] + sum(tmp)
```

## Outer-product using FMA

```
for A_t, B_t, C_t in tile(A,B,C):
  for k in 0...16:
    for i in 0...16:
      tmp = [A_t[i][k], A_t[i][k], ...]
      C_t[i] = fma16(C_t[i], tmp, B_t[k][:])
```

## Ours (16X acceleration)

```
for A_t, B_t, C_t in tile(A,B,C):
  for k in 0...16:
    C_t = opa16(C_t, A_t[:][k], B_t[k][:])
```

8

# Assembly Comparison
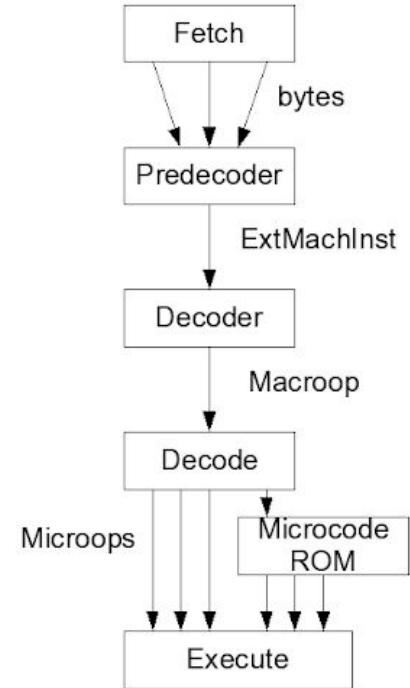
```
101     .L16:
102         ldp  x0, x1, [x6]
103         stp  x0, x1, [x12]          @Av <- A[k][i:i+32]
104         mov  x3, x12                @x3 = &Av
105         ldp  x8, x9, [x11]
106         mov  x1, x16
107         ldp  x14, x15, [x6, 16]
108         stp  x8, x9, [x30]
109         mov  x0, x16
110         stp  x14, x15, [x12, 16]
111         ldr  q3, [sp, 48]           @v3 = Bv[:16]
112         ldr  q2, [x11, 16]          @v2 = B[0][16:32]
113         str  q2, [x30, 16]          @Bv[16:] = v2
115     .L12:
116         ld1r    {v0.16b}, [x3], 1   @v0 = {Av[i++]}
117         mul  v1.16b, v3.16b, v0.16b @v1 = v3*a
118         mul  v0.16b, v0.16b, v2.16b @v0 = a*v2
119         stp  q1, q0, [x0]           @Ct[i][:] = {v1, v0}
120         add  x0, x0, 32             @i++
121         cmp  x0, x2                 @i<32
122         bne  .L12
123         add  x0, x22, x5
124         .p2align 3,,7
125     .L11:
```

```
101     .L16:
102         ldp  x0, x1, [x6]
103         stp  x0, x1, [x12]          @Av <- A[k][i:i+32]
104         mov  x3, x12                @x3 = &Av
105         ldp  x8, x9, [x11]
106         mov  x1, x16
107         ldp  x14, x15, [x6, 16]
108         stp  x8, x9, [x30]
109         mov  x0, x16
110         stp  x14, x15, [x12, 16]
111         ldr  q3, [sp, 48]           @v3 = Bv[:16]
112         ldr  q2, [x11, 16]          @v2 = B[0][16:32]
113         str  q2, [x30, 16]          @Bv[16:] = v2
114         ldp  q0, q4, [x3]           @v0 = Av[:16], v4 = Av[16:]
115         opa  m0, v0, v3
116         opa  m1, v4, v3
117         opa  m2, v0, v2
118         opa  m3, v4, v2
119         .p2align 3,,7
120     .L11:
```

Entire loop (32-iter) to 4 instructions

9

# Evaluation Method - Simulation

- Find an unallocated instruction
- Add customized OPA instruction to Gem5
- Compile the synthesized program with/without our OPA extension
- Run the simulations on Gem5

# Progress and Plan

- ## Current Progress
  - Functional element Power/Area estimation (45nm)
  - ISA definition
  - Baseline matrix multiplication assembly code

- ## Apr 30 - May 12
  - Assembly code with/w.o. custom instruction
  - Gem5 instruction extension implementation and test
  - SpeedUp compared without extension
  - Edge TPU baseline results

# Progress & Current Results

Please Put area, power (numbers we have now).