

BrainTTA: A 28.6 TOPS/W Compiler Programmable Transport Triggered NN SoC

IEEE INTERNATIONAL CONFERENCE ON COMPUTER DESIGN - 2023

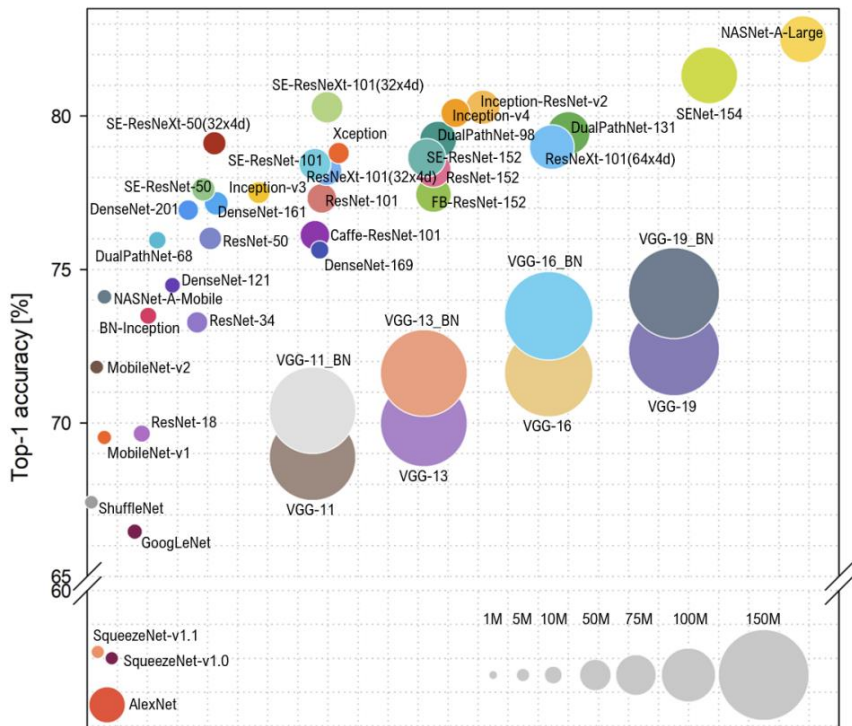
[Maarten J. Molendijk](#)^{1,2}, Floran A.M. de Putter¹, Manil Dev Gomony¹, Pekka Jääskeläinen³, Henk Corporaal¹

¹Eindhoven University of Technology, the Netherlands

²NXP Semiconductors, the Netherlands

³Tampere University, Finland

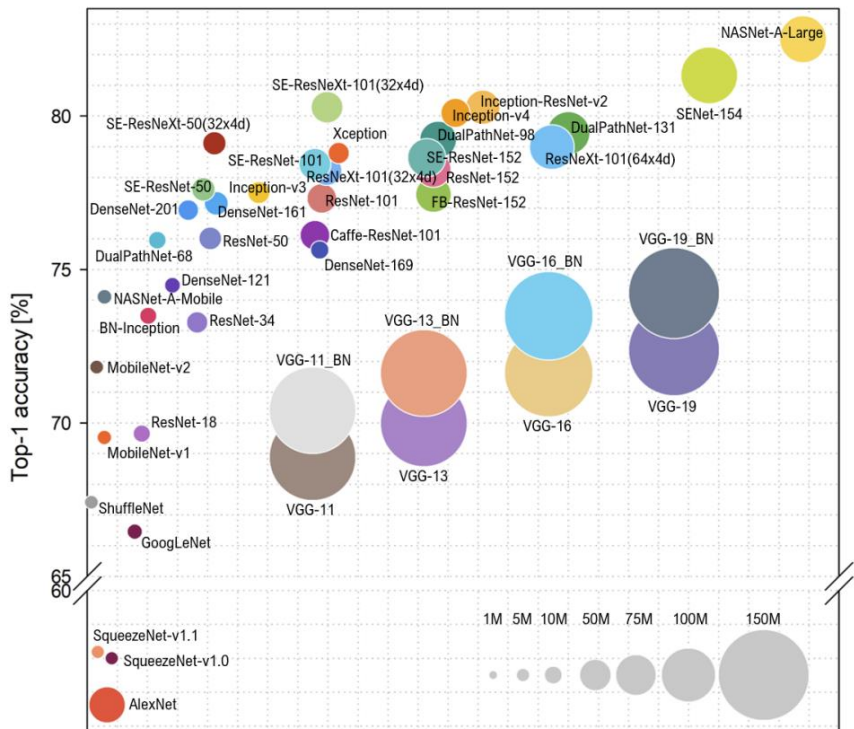
NN Architecture and Hardware Development Cost



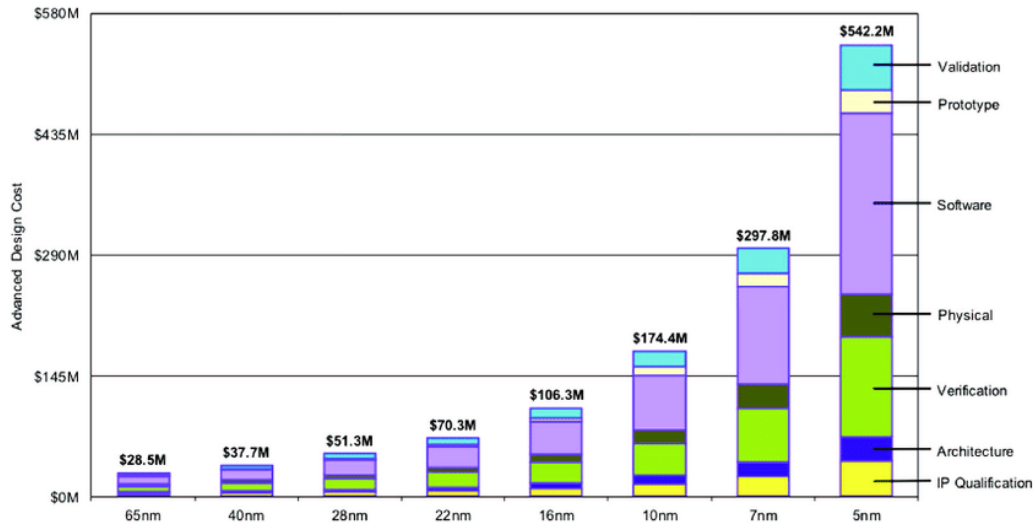
Source: [2]

- Large **variety** of NN architectures
- Rapidly evolving

NN Architecture and Hardware Development Cost

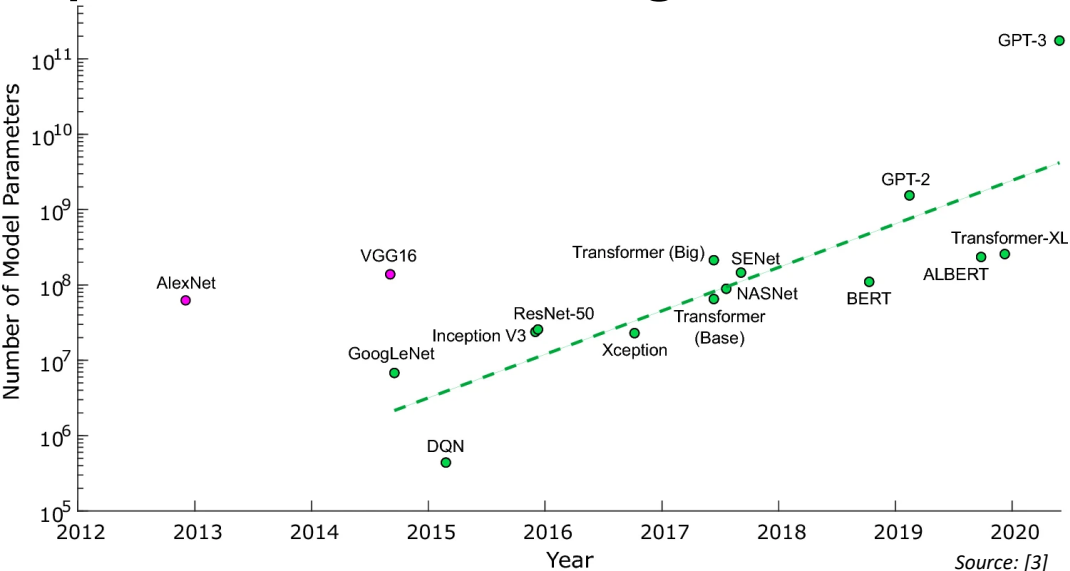


- Large **variety** of NN architectures
- Rapidly evolving



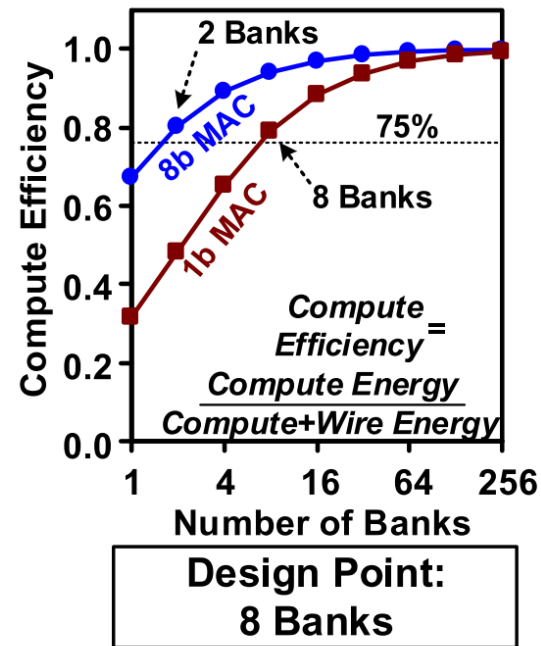
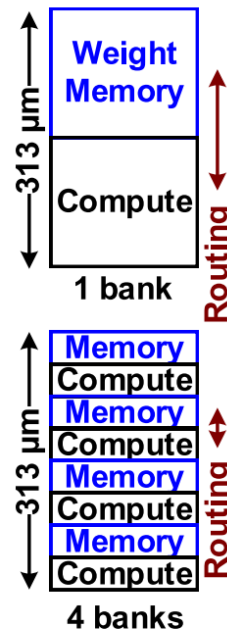
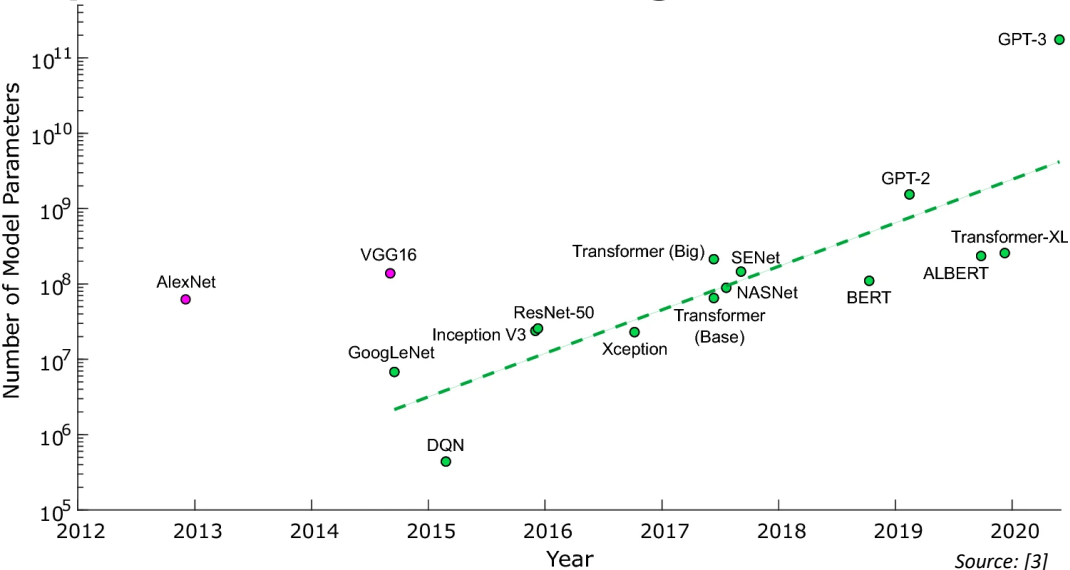
- NN memory + compute reqs
- Increased HW development cost
→ **reconfigurable solution**

Operand Precision Scaling



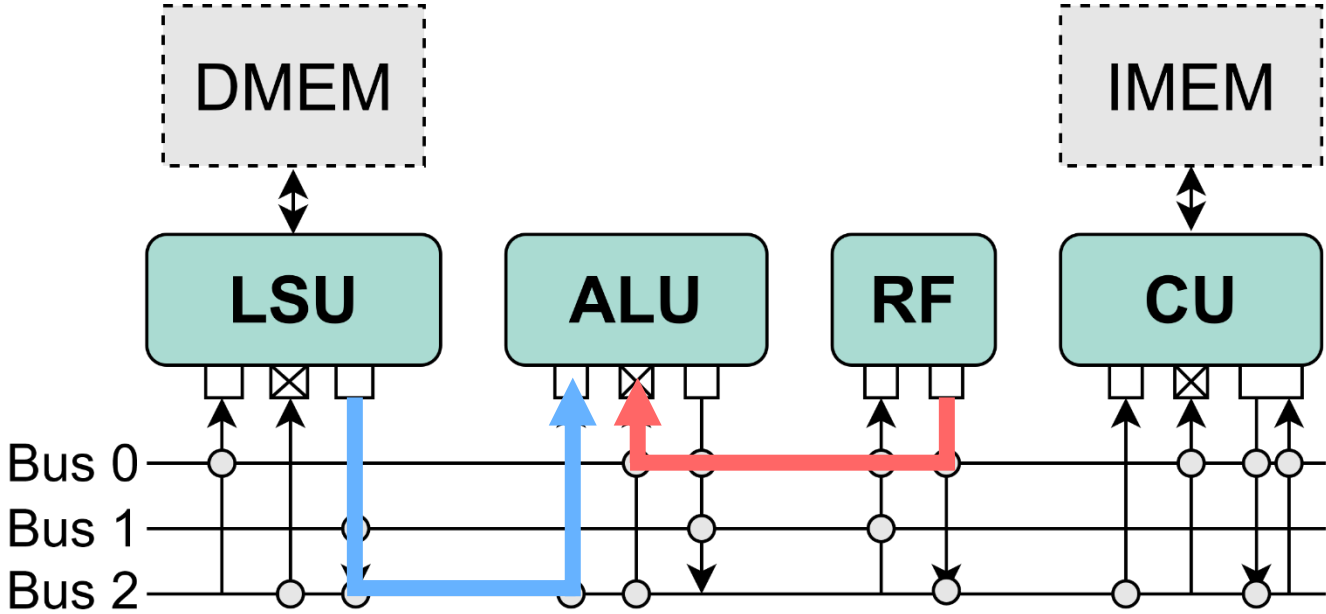
- Deployment on edge devices → **Quantization**
- Operand width ↓
 - MAC HW **superlinearly** ↓↓
 - Overhead **(sub)linearly** ↓

Operand Precision Scaling



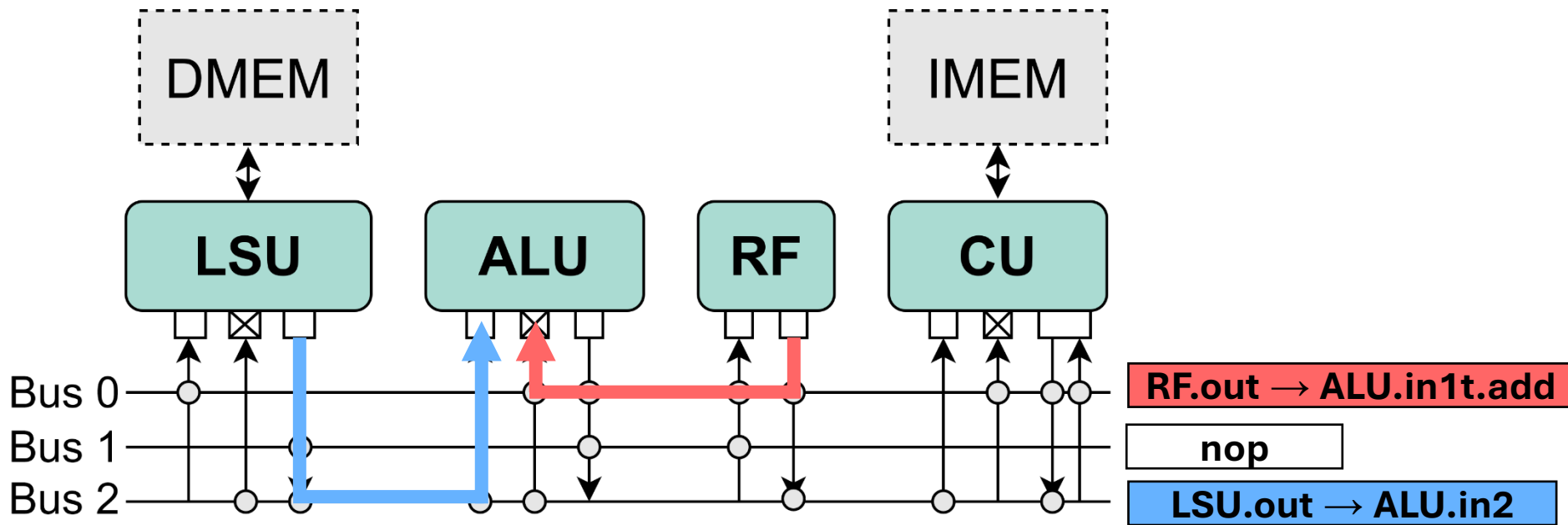
- Deployment on edge devices → **Quantization**
- Operand width ↓
 - MAC HW **superlinearly** ↓↓
 - Overhead **(sub)linearly** ↓
- Efficient data reuse + minimized data movement

Transport-Triggered Architecture



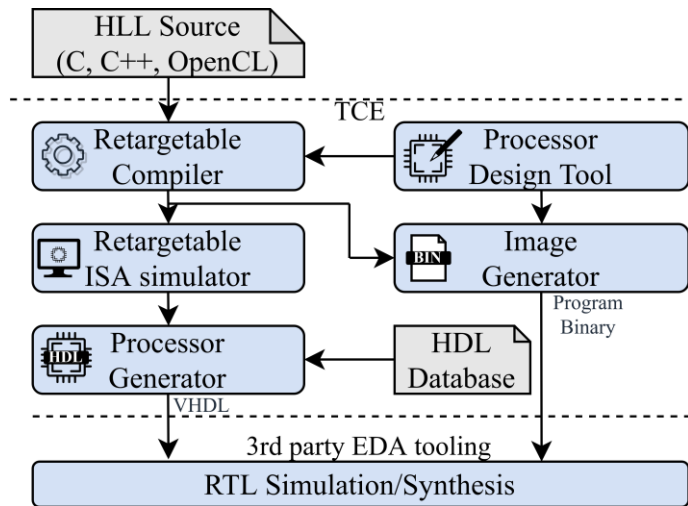
| | | |
|------------------------------|------------|--------------------------|
| RF.out → ALU.in1t.add | nop | LSU.out → ALU.in2 |
|------------------------------|------------|--------------------------|

Transport-Triggered Architecture



- + Compile-time configurable → flexible schedule
- + Exposed datapath → RF bypassing
- + Exposed datapath → Operand sharing

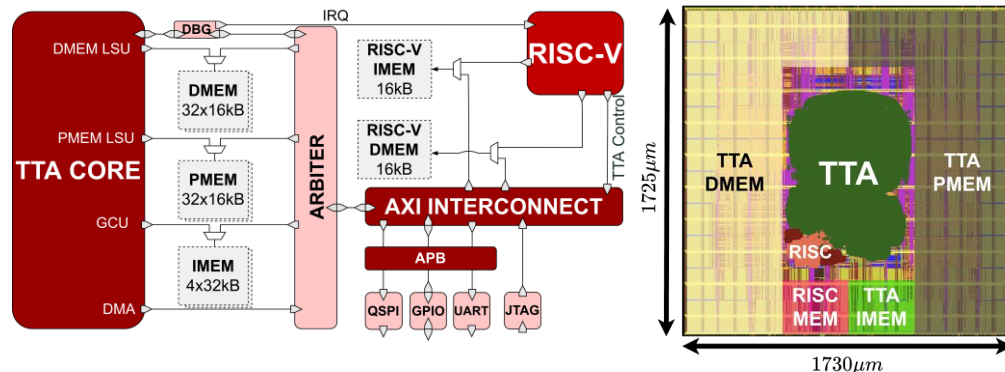
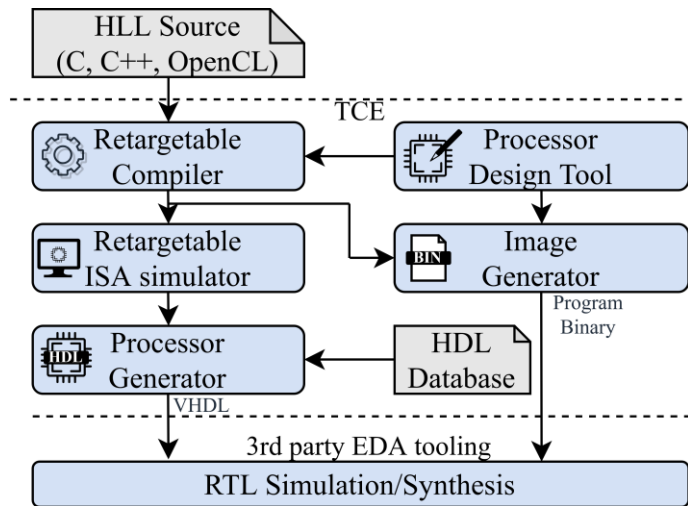
BrainTTA: Toolchain and System



Design Flow

- **OpenASIP, retargetable** [1]
- LLVM-based compiler
- ISA simulator
- HDL Database → custom units

BrainTTA: Toolchain and System



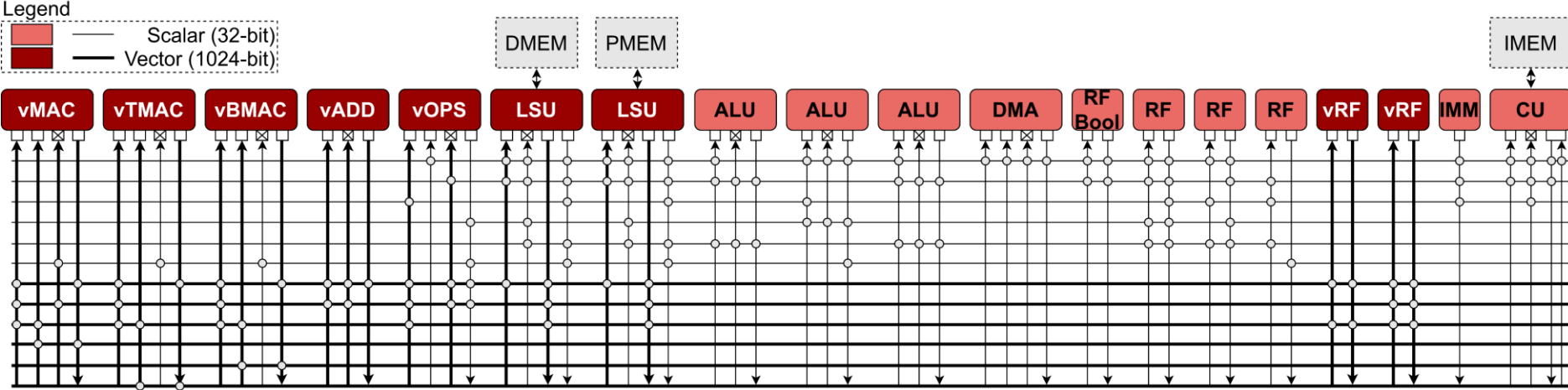
Design Flow

- **OpenASIP, retargetable** [1]
- LLVM-based compiler
- ISA simulator
- HDL Database → custom units

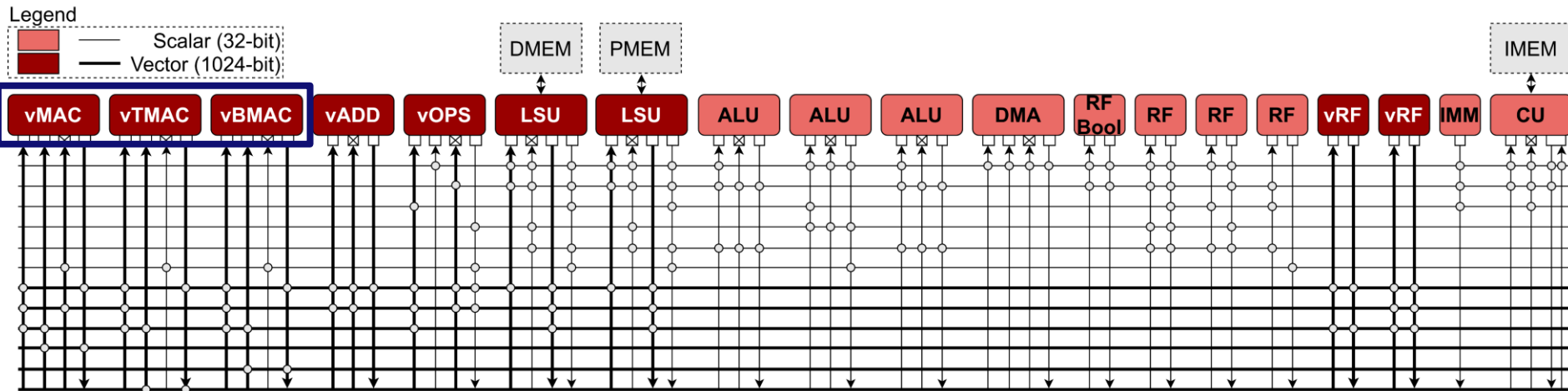
Architecture design:

- Tech: GF 22nm FDX
- RISC-V + peripherals
- DMEM/PMEM split + banked access
- IMEM + HW loopbuffer

BrainTTA: TTA Core



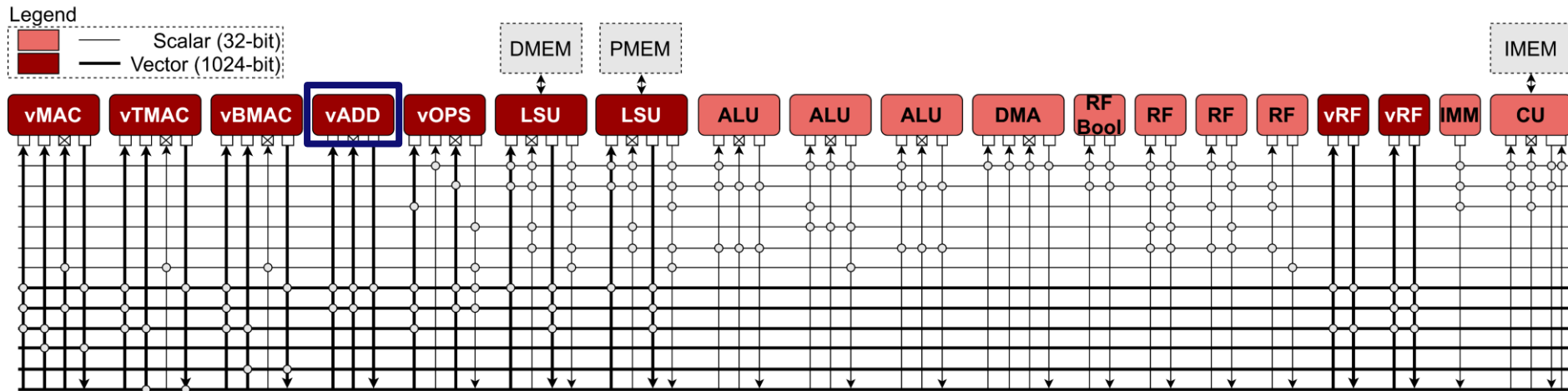
BrainTTA: TTA Core



vMAC

- 8-bit MAC
 - Scalar-Vector MAC
 - Vector-Vector MAC
- Binary MAC
- Ternary MAC

BrainTTA: TTA Core



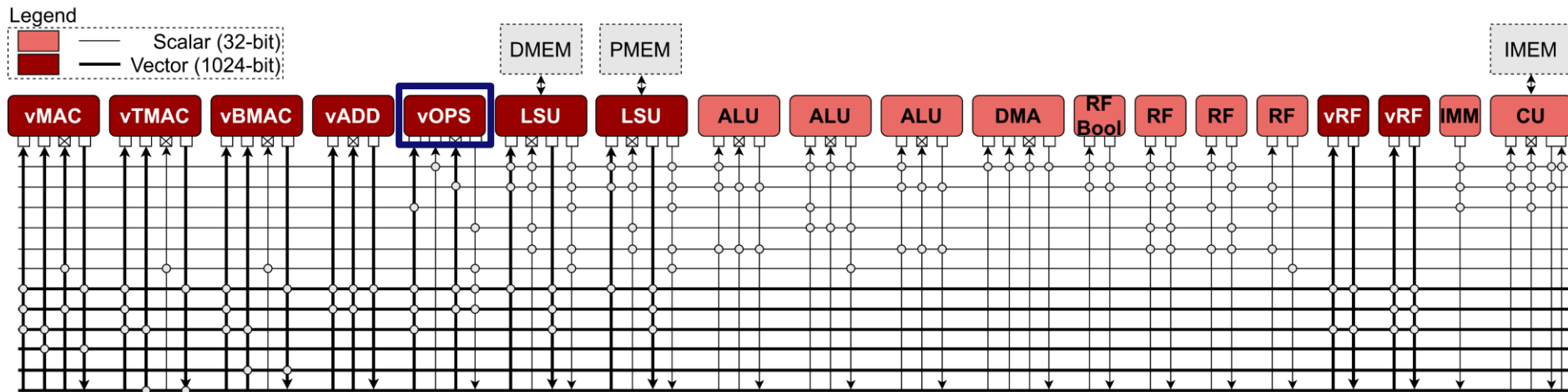
vMAC

- 8-bit MAC
 - Scalar-Vector MAC
 - Vector-Vector MAC
- Binary MAC
- Ternary MAC

vADD

- Vector-Vector addition
- Residual support

BrainTTA: TTA Core



vMAC

- 8-bit MAC
 - Scalar-Vector MAC
 - Vector-Vector MAC
- Binary MAC
- Ternary MAC

vADD

- Vector-Vector addition
- Residual support

vOPS

- Requantization
- Binarization
- Ternarization
- MaxPool
- Auxiliary ops

Application Mapping

Quantized NNs →

OS schedule

```
for h in [0, H - R + 1]:  
  for w in [0, W - S + 1]:  
    for m in [0, M]:  
      accu = bias[m]  
      for c in [0, C]:  
        for r in [0, R]:  
          for s in [0, S]:  
            accu += in[h+r][w+s][c] * w[c][r][s][m]  
          output[h][w][m] = act_function(accu)
```

Output feature map height
Output feature map width
Output channels
Input channels
Kernel height
Kernel width

Application Mapping

Quantized NNs →
OS schedule

Loop tiling
Loop interchange

```
for h in [0, H - R + 1]:  
  for w in [0, W - S + 1]:  
    for m in [0, M]:  
      accu = bias[m]  
      for c in [0, C]:  
        for r in [0, R]:  
          for s in [0, S]:  
            accu += in[h+r][w+s][c] * w[c][r][s][m]  
          output[h][w][m] = act_function(accu)
```

Output feature map height
Output feature map width
Output channels
Input channels
Kernel height
Kernel width

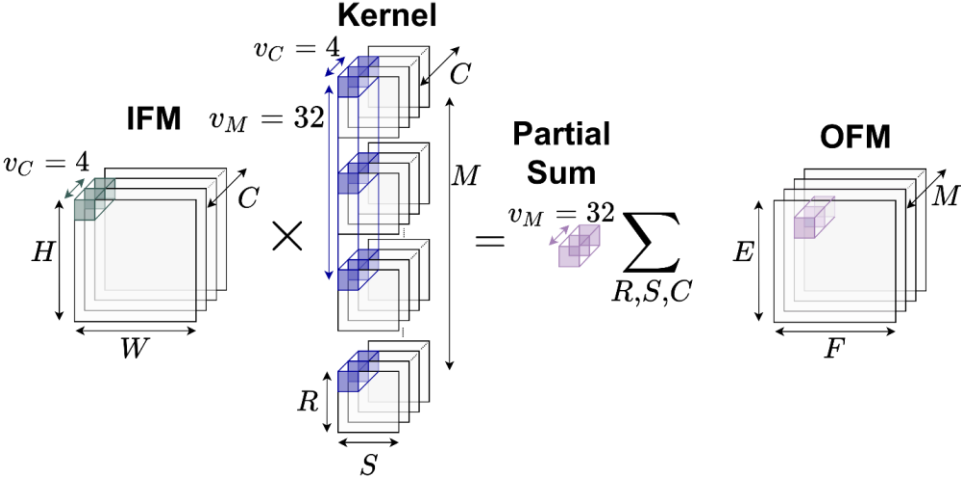
```
for h in [0, H - R + 1]:  
  for w in [0, W - S + 1]:  
    for m in [0, M/32]:  
      accu = bias[32*m]  
      for c in [0, C/4]:  
        for r in [0, R]:  
          for s in [0, S]:  
            for tm in [0, 31]:  
              for tc in [0, 3]:  
                accu += in[h+r][w+s][4*c+tc]  
                  * w[4*c+tc][r][s][32*m+tm]  
            output[h][w][m] = act_function(accu)
```

Output feature map height
Output feature map width
Output channels ($v_M = 32$)
Input channels ($v_C = 4$)
Kernel height
Kernel width

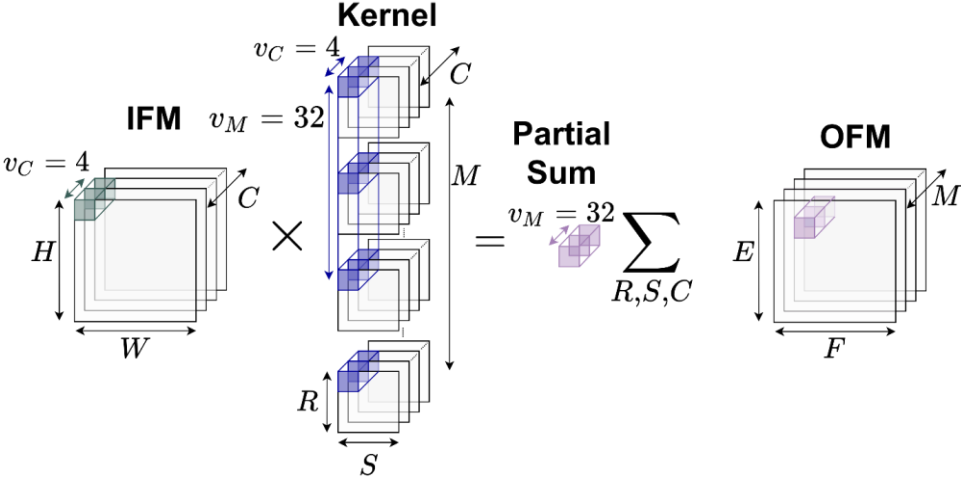
Application Mapping

```
for h in [0, H - R + 1]:           Output feature map height
  for w in [0, W - S + 1]:       Output feature map width
    for m in [0, M/32]:         Output channels ( $v_M = 32$ )
      accu = bias[32*m]
      for c in [0, C/4]:       Input channels ( $v_C = 4$ )
        for r in [0, R]:      Kernel height
          for s in [0, S]:    Kernel width
            for tm in [0, 31]:
              for tc in [0, 3]:
                accu += in[h+r][w+s][4*c+tc]
                    * w[4*c+tc][r][s][32*m+tm]
            output[h][w][m] = act_function(accu)
```

Application Mapping

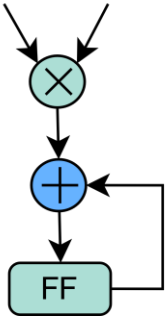


Application Mapping



$$v_C = 1$$

$$v_M = 1$$

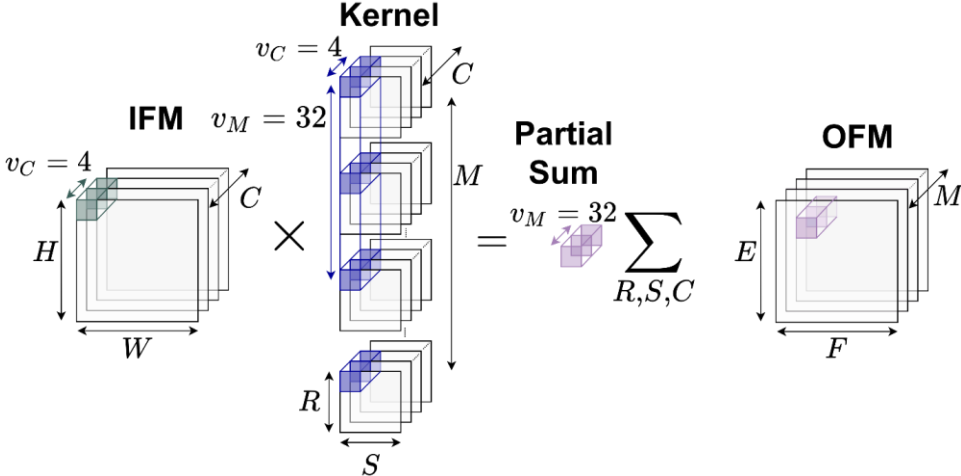
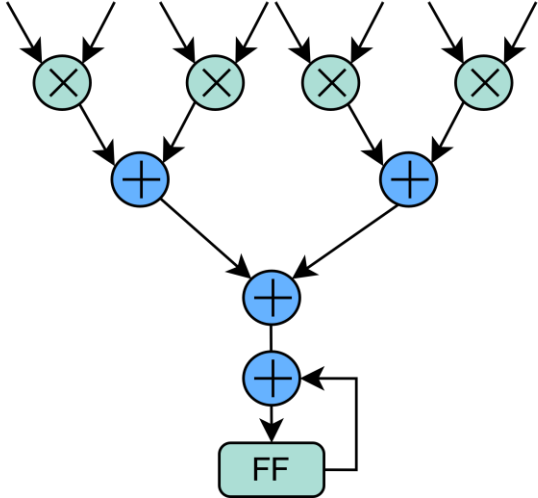


ACCU cost \gg MUL cost \rightarrow
wider reduction tree

Application Mapping

$$v_C = 4$$

$$v_M = 1$$

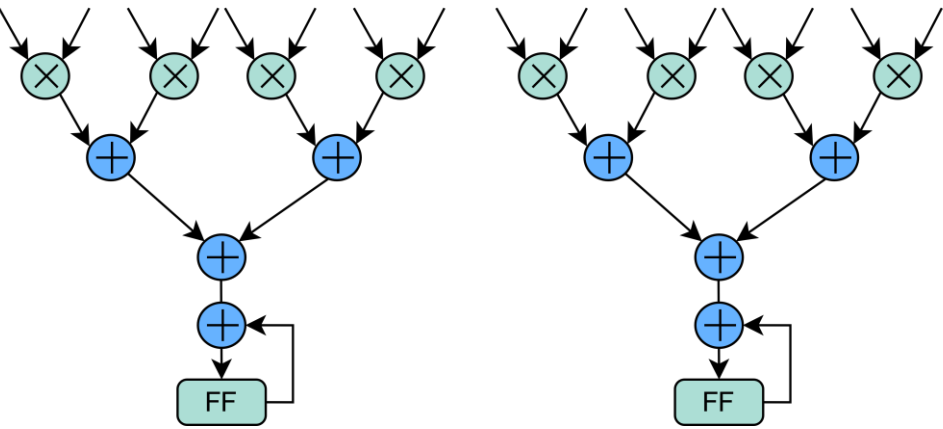
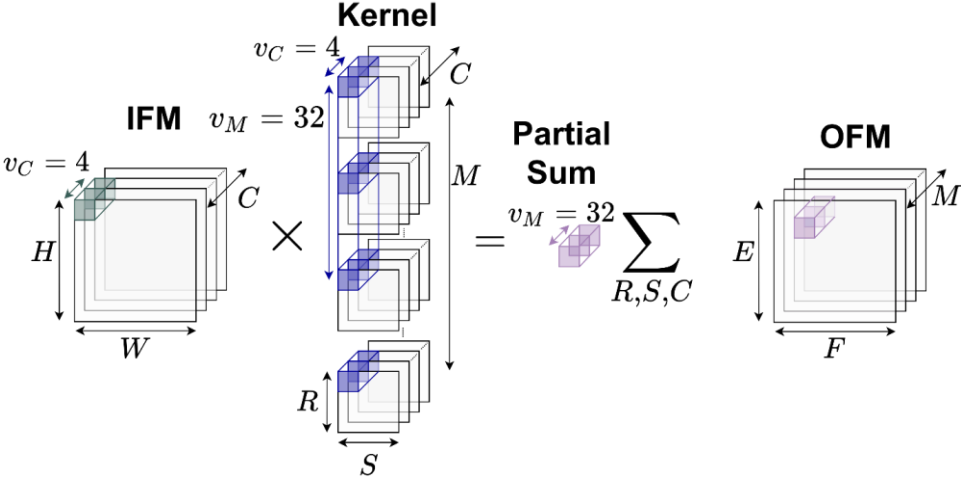


ACCU cost \gg MUL cost \rightarrow
wider reduction tree

Application Mapping

$$v_C = 4$$

$$v_M = 2$$

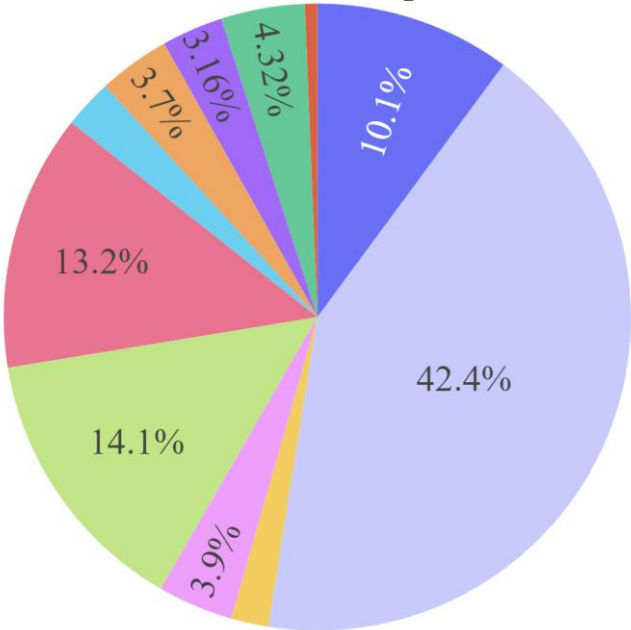


ACCU cost \gg MUL cost \rightarrow
wider reduction tree
 IFM broadcast \rightarrow **multiple reduction trees**

Post-layout Energy Consumption [GF 22nm FDX]

8-bit Convolution

E = 405 fJ/op

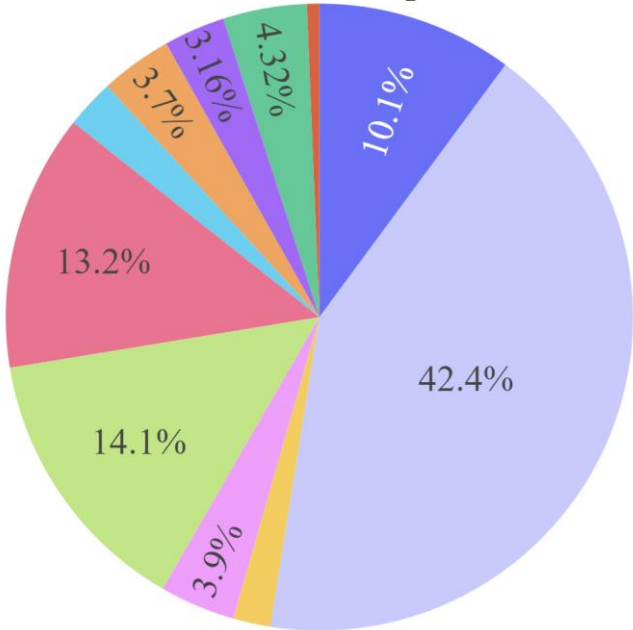


■ PMEM ■ IMEM ■ DMEM ■ Loopbuffer ■ DMEM LSU ■ PMEM LSU ■ IC ■ vRF ■ TTA other ■ RISC ■ vMAC

Post-layout Energy Consumption [GF 22nm FDX]

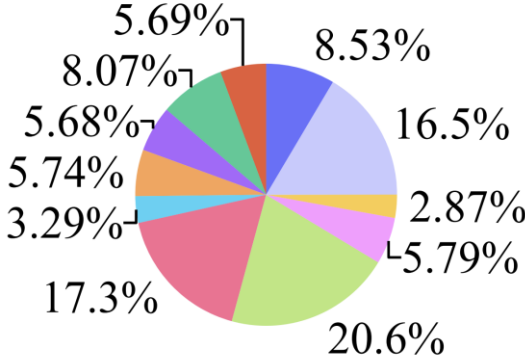
8-bit Convolution

E = 405 fJ/op



Binary Convolution

E = 35 fJ/op



■ PMEM
 ■ IMEM
 ■ DMEM
 ■ Loopbuffer
 ■ DMEM LSU
 ■ PMEM LSU
 ■ IC
 ■ vRF
 ■ TTA other
 ■ RISC
 ■ vMAC

Comparison to SotA

| | Eyeriss v2 | XNE | Samurai | XPULPNN | Dustin | This work |
|----------------------------------|-----------------------------|-----------------|-----------------|---|--|---|
| Tech | 65nm | 22nm | 28nm | 22nm | 65nm | 22nm |
| Programmability | Configurable | ASM | ASM | Compiler | Compiler | Compiler |
| Energy efficiency | 252 GOPS/W(8b) ¹ | 8.7 TOPS/W (1b) | 1.3 TOPS/W (8b) | 2.2 TOPS/W (8b) 6.1 TOPS/W (2b) | 606 GOPS/W (8b) 2304 GOPS/W (2b) | 2.5 TOPS/W (8b) 14.9 TOPS/W (T) 28.6 TOPS/W (1b) |
| Memory Cap. [kB] | 246 | 520 | 464 | 640 | 80 | 1024 ² |
| Area Eff. [GOPS/m ²] | 5.5 (8b) ³ | 28.9 (1b) | 0.6 (8b) | 21.7 (8b) 70.7 (2b) | 0.9 (8b) 3.46 (2b) | 25.8 (8b) 103.0 (T) 206.0 (1b) |

*after technology scaling

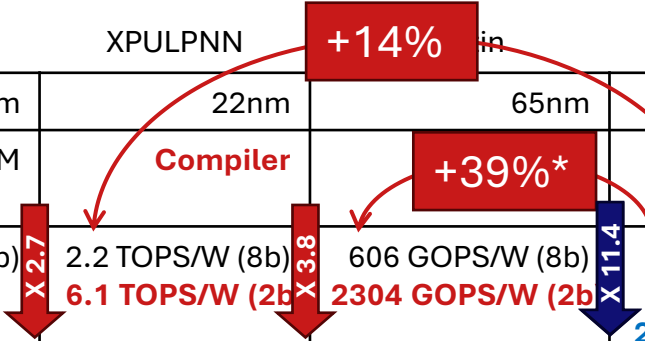
¹Average on AlexNet.

²Excluding instruction memory.

³Area estimated using gatecount diff. EyerissV1, EyerissV2.

Comparison to SotA

| | Eyeriss v2 | XNE | Samurai | XPULPNN | BrainTTA | This work |
|----------------------------------|-----------------------------|-----------------|-----------------|------------------------------------|-------------------------------------|--|
| Tech | 65nm | 22nm | 28nm | 22nm | 65nm | 22nm |
| Programmability | Configurable | ASM | ASM | Compiler | Compiler | Compiler |
| Energy efficiency | 252 GOPS/W(8b) ¹ | 8.7 TOPS/W (1b) | 1.3 TOPS/W (8b) | 2.2 TOPS/W (8b) 6.1 TOPS/W (2b) | 606 GOPS/W (8b) 2304 GOPS/W (2b) | 2.5 TOPS/W (8b) 14.9 TOPS/W (T) 28.6 TOPS/W (1b) |
| Memory Cap. [kB] | 246 | 520 | 464 | 640 | 80 | 1024 ² |
| Area Eff. [GOPS/m ²] | 5.5 (8b) ³ | 28.9 (1b) | 0.6 (8b) | 21.7 (8b) 70.7 (2b) | 0.9 (8b) 3.46 (2b) | 25.8 (8b) 103.0 (T) 206.0 (1b) |

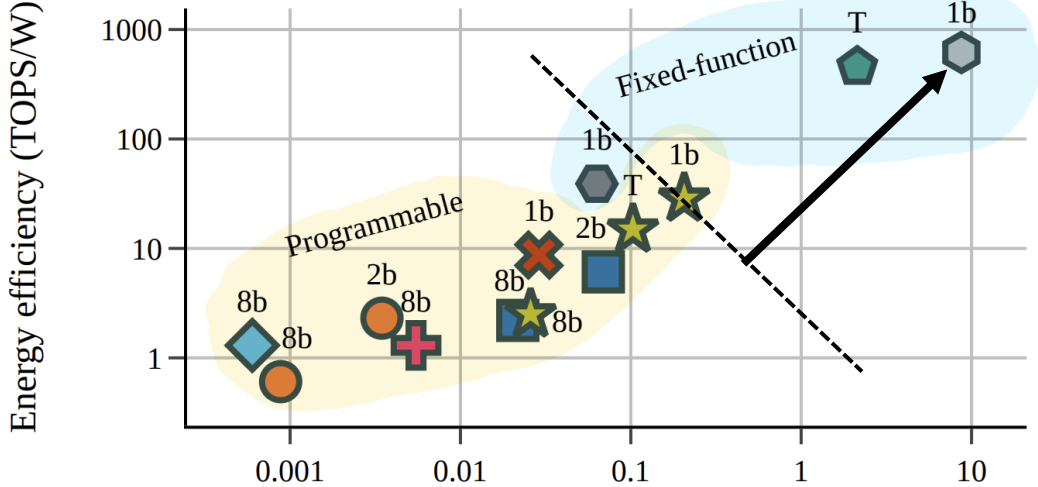


*after technology scaling

¹Average on AlexNet.
²Excluding instruction memory.
³Area estimated using gatecount diff. EyerissV1, EyerissV2.

Programmable vs. Fixed-Function Trade-off

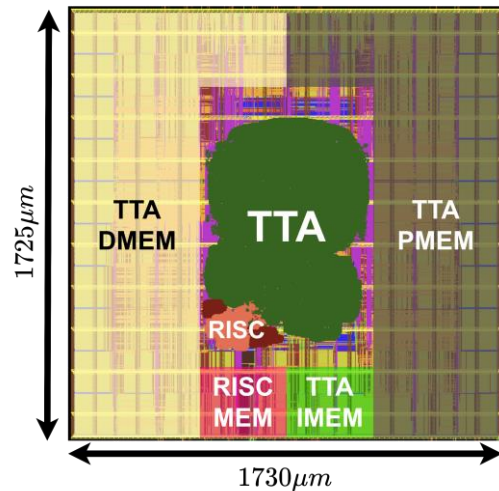
- Programmable architectures
- Fixed-function
 - Spatial loop unrolling
 - Fixed FM/weight buffering



| | | | |
|---|-----------------------------|---|---|
| ● | Dustin, 65nm, 0.8V, [18] | + | Eyeriss V2, 65nm, [21] |
| ■ | XPULPNN, 22nm, 0.6V, [17] | ⬠ | CUTIE, 22nm, 0.65V, [23] |
| ◆ | SamuraiAI, 28nm, 0.45V [20] | ⬠ | Knag et al, 10nm, 0.37V, [14] |
| ⊗ | XNE, 22nm, 0.6V [19] | ⬠ | ChewBaccaNN, 22nm, 0.4V, [22] (3x3 kernel) |
| ★ | BrainTTA, 22nm, 0.5V | | |

BrainTTA - Conclusion

- Efficient and flexible NN inference engine:
 - Mixed-precision
 - Compile-time reconfigurable
 - Eff: **2.47** / **14.9** / **28.6** [TOPS/W]
 - Throughput: **77** / **307** / **614** [GOPPS]
- Superlinear energy eff. scaling
 - 8-bit → ternary: **x5.96**
 - 8-bit → binary: **x11.4**



Q&A