Hardware Accelerators for Deep Neural Networks

Energy Efficiency of Hardware Acceleration

Olivier Sentieys Univ. Rennes, Inria, IRISA

olivier.sentieys@irisa.fr

https://gitlab.inria.fr/sentieys/dnn_acc









What will you learn in this course?

- Principle and design of DNN Accelerators
 - Energy efficiency of hardware accelerators
 - Speeding up the GEMM kernel
 - Designing hardware accelerators for GEMM/CONV
 - Available accelerators for DNNs
 - Computing at the right precision

https://gitlab.inria.fr/sentieys/dnn_acc



Complexity Issues of Deep Neural Networks



- Two main tasks
 - training determine set of network parameters to solve a *task* (minimize a loss on a *training set*)
 - inference given an input, compute (forward propagate) using the trained network



Computing Demand of Al

• is higher than what computer architectures can bring



4

Evolution of the Number of Parameters

• is much higher than available (on-chip) memory capacity



Memory Bottleneck



Data movement

- move input data & model from memory to compute units
- send partial results back to memory

Computations

- vector/matrix manipulations
- done on CPU, GPU, or custom accelerators (e.g., FPGA, ASIC)

Evolution of Bandwidth

• is much slower than FLOPS



On the Computer Architecture Side

The Hardware Lottery

Sarah Hooker, The Hardware Lottery, Communications of The ACM, 2021

Silicon Technology Evolution

- Now several billions or transistors!
 - Apple M1: 33 B.Tr, 5nm, 2.5cm2

ALL CONTRACTOR OF CONTRACTOR O

- Amazing compute progress
 - 12 orders of magnitude performance i
 - A supercomputer in every body's pocket

HPE Frontier

1.6 Exa FLOPS

The Many Walls of Computer Architecture



48 Years of Microprocessor Trend Data

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2019 by K. Rupp

Energy Cost in a Chip

• Fetching operands costs more than computing



Energy Efficiency Energy Performance Efficiency e.g., Tera op/s (TOPS) e.g., TOPS/Watt Joules Operation Operations Second Power =Х

- Power budget is fixed
- How to increase energy efficiency while maintaining performance?

Improving Energy Efficiency

- Technology?
 - What can advanced technology nodes bring?
 - Dark Silicon Era
- Accelerate
 - Energy advantages of specialized hardware
- Approximate
 - Playing with precision and number representations to reduce energy
- Key message of this course: specialized hardware that computes at the right (lowest) precision

Outline

- Part I: the clear need for specialized hardware
 - Energy efficiency of hardware accelerators
- Part II: accelerating GEMM
 - DNN kernels with a focus on convolution
 - Speeding up the GEMM kernel
 - Designing hardware accelerators for GEMM/CONV
 - Available accelerators for DNNs
- Part III: computing at the right precision
 - A bit of arithmetic
 - Inference and training with low precision

Energy Cost in a Processor D-cache 6% Operations I-cache 23% – 32-bit addition: 0.05pJ - 16-bit multiply: 0.25pJ Datapath - 64-bit FPU: 20pJ/op Fetch/ 38% Decode 19% Reg. File Instructions 14% - fetch, decode, read two

operands from RF,

execute, write back

General Purpose Processor 91 pJ/instruction

Achieving Higher Performance

- Pushing clock frequency...
- Branch/value prediction
- Cache memory
- In-core parallelism
 - Superscalar
 - Out of order execution
 - VLIW+good compilers



• Multiple cores on a single chip

Pushing for Hardware Acceleration!

What is a Hardware Accelerator?

- Specialized hardware for a given set of kernels
- With limiting programmability
- Computes just right!

• e.g., Matrix Multiplication





An example: Bitcoin Mining

Туре	Model	Mhash/s	Mhash/J	Power (W)
GPP	Intel Xeon X5355 (dual)	22.76	0.09	120
GPP	ARMCortex-A9	0.57	1.14	1.5
GPP	Intel Core i7 3930k	66.6	0.51	130
GPU	AMD 7970x3	2050	2.41	850
GPU	Nvidia GTX460	158	0.66	240
ASIC	AntMiner S1	180.000	500	360
ASIC	AntMiner S5	1.155.000	1957	590
FPGA	Bitcoin Dominator X5000	100	14.7	6.8
FPGA	Butterflylabs Mini Rig	25.200	20.16	1250





Apple Silicon M2 Max

- 5 nm (TSMC 2G), 40 billion transistors
- 8 performance cores
 - 38 Int MOPS, 56 MFLOPS
 - NEON vector processor
- 4 power-efficiency cores
- Unified memory (32-96 GB LPDDR5-6400) next to the (400 GB/s bandwidth) for GPU&CPU
- 89W (peak CPU+GPU), CPU 36W (peak)
- High-performance media engine
- 16 TOPS Neural Engine
 - 10,000 times the GPU speed for ML tasks
 - "power-efficient" (but no reported power figures)



Accelerators for ML



CPU	GPU	FPGA	TPU
Threads SIMD	Massive Threads SIMD HBM	LUTs DSP BRAM	MM Unit BRAM

Key Takeaways

- Energy efficiency requires deeply specialized hardware
 - which also may come with pain from the programmer/designer
- Basic tasks of DNNs are easy to accelerate
 - this course is mainly focused on matrix multiplication
- Number representations and precisions are key techniques
 - also memory access since execution is often memorybound

TARAN Team at a Glance

Accelerators

Approximate

Computing

Resilient

Computing

Embracing Technologies



Lannion



∮ ≸IRISA



Domain-Specific Computers in the post Moore's law era

- ~40 people, Rennes and Lannion campuses
- Our focus: hardware specialization and acceleration
 - Energy Efficiency of hardware accelerators
 - Domain-specific architectures, languages and compilers
- Automatically create hardware that is resilient, secured, and computes just right
 - From Sensors to Clouds

Outline

- Part I: the clear need for specialized hardware
 Energy efficiency of hardware accelerators
- Part II: accelerating GEMM
 - DNN kernels with a focus on convolution
 - Speeding up the GEMM kernel
 - Designing hardware accelerators for GEMM/CONV
 - Available accelerators for DNNs
- Part III: computing at the right precision
 - A bit of arithmetic
 - Inference and training with low precision

Speeding Up GEMM Efficient Processing of Matrix Multplication

Focus on Convolution Neural Networks (CNN)

Input Activations



W

- H: Height of Input Activation
- W: Width of Input Activation
- R: Height of Weight

Н

- S: Width of Weight
- T: Height of Output Activation
- U: Width of Output Activation



Example is with:

Weights

stride=1

of rows/columns traversed per step

Output Activations

padding=0

of zero rows/columns added



 $A = a \times 1 + b \times 2 + c \times 3 + f \times 4 + g \times 5 + h \times 6 + k \times 7 + l \times 8 + m \times 9$

 $(R \times S)$ Multiply and Accumulate (MAC) operations



$B = b \times 1 + c \times 2 + d \times 3 + g \times 4 + h \times 5 + i \times 6 + l \times 7 + m \times 8 + n \times 9$

 $(R \times S)$ Multiply and Accumulate (MAC) operations

 $T = \frac{H - R}{stride} + 1$

 $U = \frac{W - S}{stride} + 1$

2D Convolution



 $(T \times U \times R \times S)$ MAC operations in total

A lot of potential data reuse for memory accesses

С

а

f

k

р

u



Convolution Loop Nest

```
for (n=0; n<N; n++) {
                                             // for each Batch
    for (k=0; k<K; k++) { // for each Output Channel</pre>
      for (t=0; t<T; t++) {
                                             // OA Height
        for (u=0; u<U; u++) {</pre>
                                             // OA Width
          OA[n][k][t][u]= 0;
          for (r=0; r<R; r++) { // W Height
            for (s=0; s<S; s++) { // W Width
              for (c=0; c<C; c++) { // for each Input Channel
                h = t * stride - pad + r;
CONV
                w = u * stride - pad + s;
kernel
                OA[n][k][t][u] += IA[n][c][h][w] * W[k][c][r][s];
                                                        d
          Activation(OA[n][k][t][u]);
                                                                      K
                                                      h
                                                       i
                                                                             D
                                                 Н
                                                     I m n o
                                                    q r
                                                   р
                                                       st
                                                                R
                                                      W
                                                                                 32
```

Opportunities for Data Reuse

Convolution Kernel

Input Activations (feature maps)





Output Activations



N: Batch size

W



• Input reuse

- different filters
 are applied to the
 same input
- each input is reused K times

C: # of Input Channels K: # of Output Channels N: Batch size 35



Filter (weight) reuse

- when processing a batch of size N, all inputs are applied to the same filter
- each filter weight is reused N times

C: # of Input Channels K: # of Output Channels N: Batch size 36


Conv. reuse

- filters slide across different positions of the same input
- each weight is reused \approx T.U times
- each input is reused \approx R.S times

C: # of Input Channels K: # of Output Channels N: Batch size 37

Other Kernels

- Fully-Connected Layer
 H=W=R=S=T=U=1
- Depth-Wise Convolution
 K=1
- Pooling Layer
 - [MAX, AVG], pooling stride and kernel size
- BatchNorm Layer
 - provides zero-mean, unit-variance activations
- Activations
 - ReLU, L-ReLU, sigmoid, tanh, clipping







Traversal Order

Caching

- CPU caches are orders of magnitude faster, but much smaller, so using them correctly is critical
 - Automatically managed by the CPU.
 - Every time we fetch data from the main memory, the CPU automatically loads it and its neighboring memory into the cache, hoping to utilize locality of reference.



Caching

- In our case:
 - once we access
 A[m, k], the next
 element in the
 row, A[m, k+1] is
 already cached
 - but we get a cache miss for each data from matrix B fetched B[k, n]

for (m=0; m<M; m++)
for (n=0; n<N; n++)
for (k=0; k<K; k++)
C[m][n] += A[m][k] * B[k][n]</pre>
Baseline



М

.....

K





.....





Traversal Order

44

```
for (m=0; m<M; m++) Baseline Reordered
for (k=0; k<K; k++)
for (n=0; n<N; n++)
C[i][j] += A[i][k] * B[k][j]</pre>
```

- Reordering the loops
 from m, n, k to m, k, n
 - Improve data locality (better cache usage)



46



K

m

Traversal Order

N



- Reordering the loops
 from m, n, k to m, k, n
 - Improve data locality (better cache usage)





- Tiling
 - Looping on smaller submatrices (tiles of size $T \times T$)
 - small enough to fit in the cache



- Tiling
 - Looping on smaller submatrices (tiles of size $T \times T$)
 - small enough to fit in the cache



• Tiling

- Looping on smaller submatrices (tiles of size $T \times T$)

```
for (m=0; m<M/T; m++)
  for (n=0; n<N/T; n++)
    for (k=0; k<K/T; k++)
```

```
for (mt=0; mt<T; mt++) // Tile-level mult.
  for (nt=0; nt<T; nt++)</pre>
    for (kt=0; kt<T; kt++)</pre>
```

- // Tile row index // Tile column index
 - // Tile inner index

note: for this code M, N and K must be divisible by T

```
C[m*T+mt][n*T+nt] += A[m*T+mt][k*T+kt]
                   * B[k*T+kt][n*T+nt]
```

Tiling



Results on CPU (Apple Silicon M2 Pro)

- More of Tiling
 - Tile inner loops can be vectorized and unrolled
 - Tiles can run in parallel (multithreading)

Try it yourself

git clone <u>https://gitlab.inria.fr/sentieys/dnn_acc.git</u>
cd GEMM

look at the C code: baseline.c, baseline_reordered.c, opti-l1.c,
opti-l2.c

make all

How to add SIMD and vectorization?

Designing Hardware Accelerators for GEMM/CONV

Building the Accelerator

Pipeline (Dataflow) Architecture



Li: Layer i



Sequential Architecture

Building the Accelerator

• Sequential Architecture





Building the Accelerator Are

• Dataflow Architecture





Exploiting Data Reuse in PE Array

- Temporal Architecture
 - SIMD (CPU), SIMT (GPU)
 - Classical Memory Hierarchy



- Spatial Architecture
 - Dataflow accelerators



Exploiting Data Reuse

- Why reuse is important?
 - Relative energy costs
 - Memory access is the bottleneck





DRAM: 200x SRAM: 6x PE: 2x MAC: 1x Energy (external) (cache, buffer)

Exploiting Data Reuse

- Temporal Reuse
 - e.g., memory hierarchy
 - the same data is used more than once over time by the same PE



- Temporal and Spatial Reuse
 - Memory hierarchy and multiple PEs

- Spatial Architecture
 - e.g., systolic, multicast
 - the same data is used by more than one PE at different spatial locations of the hw



Back to the GEMM Kernel



- Parallelizing most inner loop
 - (1) Adder tree
 - Typical width: 8-64
 - e.g. NVDLA, NVIDIA Tensor cores, FINN

```
for (m=0; m<M; m++) {
  for (n=0; n<N; n++) {
    C[m][n] = 0;
    parallel_for (k=0; k<K; k++) {
        C[m][n] += A[m][k] * B[k][n];
    }
    for: temporal execution order
    parallel_for: parallel execution</pre>
```



- Parallelizing most inner loop
 - (2) Systolic Multiply-And-Accumulate (MAC)
 - Typical width: 8-256
 - e.g. Gemmini, Google TPU

```
for (m=0; m<M; m++) {
  for (n=0; n<N; n++) {
    C[m][n] = 0;
    parallel_for (k=0; k<K; k++) {
        C[m][n] += A[m][k] * B[k][n];
    }</pre>
```



- Parallelizing second inner loop
 - (3) Multicasting a (sub)line of weights
 - Typical width: 8-16
 - e.g. NVDLA, NVIDIA Tensor cores

```
for (m=0; m<M; m++) {
    parallel_for (n=0; n<N; n++) {
        C[m][n] = 0;
        for (k=0; k<K; k++) {
            C[m][n] += A[m][k] * B[k][n];
        }
        for: temporal execution order</pre>
```

parallel_for: parallel execution



[Credit: Sophia Shao, Hardware for Machine Learning, Course@UC Berkeley]

64

- Parallelizing second inner loop
 - (4) Systolic multicast
 - Typical width: 8-256
 - e.g. Gemmini, Google TPU

```
for (m=0; m<M; m++) {
  parallel for (n=0; n<N; n++) {
     C[m][n] = 0;
                                                           IA-REG
                                                                    IA-REG
     for (k=0; k<K; k++) {
         C[m][n] += A[m][k] * B[k][n];
                                                       [Credit: Sophia Shao, Hardware for
                                                       Machine Learning, Course@UC Berkeley]
```

-REG

₹

65

Parallelizing second and most inner loops

```
- (1)+(3) Adder Tree + Multicast
```

– e.g. NVDLA, FINN

```
for (m=0; m<M; m++) {
    parallel_for (n=0; n<N; n++) {
        C[m][n] = 0;
        parallel_for (k=0; k<K; k++) {
            C[m][n] += A[m][k] * B[k][n];
        }
    }
}</pre>
```



• Parallelizing second and most inner loops

- (2)+(4) Systolic MAC + Systolic Multicast

– e.g. Gemmini, Google TPU

```
for (m=0; m<M; m++) {
    parallel_for (n=0; n<N; n++) {
        C[m][n] = 0;
        parallel_for (k=0; k<K; k++) {
            C[m][n] += A[m][k] * B[k][n];
        }
    }
}</pre>
```



[Credit: Sophia Shao, Hardware for Machine Learning, Course@UC Berkeley]

67

Systolic Arrays

Replace single Processing Element (PE) with an array of regular PEs



- Orchestrate data flow for high throughput with less memory access than classical architectures
- Each PE may have (small) local instruction and data memory
- Analogy with the heart → blood → (many)cells → heart

Systolic Arrays are New?

SYSTOLIC ARCHITECTURES FOR CONNECTED SPEECH RECOGNITION (*)

Francois CHAROT Patrice FRISON Patrice QUINTON

IRISA, Campus de Beaulieu, 35042 RENNES-CEDEX FRANCE

July 1984

ABSTRACT

Systolic arrays for two connected speech recognition methods are presented. The first method is based on the dynamic time warping algorithm which is applied directly on acoustic feature patterns. The second method is the probabilistic matching algorithm which requires that the input sentence be preprocessed by a phonetic analyzer. It is shown that both methods may be implemented on either a two-dimensional or a linear systolic array. Advantages of each of these implementations are discussed. The architecture of a 12000 transistors programmable NMOS prototype IC which can be used as the basic processor of

(*) This work has been partly funded by a grant of the CNET in Lannion (French Telecommunication Research Center).

Systolic Array Matrix Multiplication



a _{0,0}	a _{0,1}	a _{0,2}	a _{0,3}
a _{1,0}	a _{1,1}	a _{1,0}	a _{1,3}
a _{2,0}	a _{2,1}	a _{2,0}	a _{2,3}
a _{3,0}	a _{3,1}	a _{3,0}	a _{3,3}
















- Cycle 7: all PEs contain $C_{i,i}$ results
- $C_{i,j}$ values can be shifted to last column/row through the PEs at each cycle



• Weight (B) Stationary

First phase:
 load weights





• Cycle 1

3.3



 $a_{0,0} \, \mathbf{x} \, b_{0,0}$

• Cycle 2



$a_{0,1} x b_{0,0}$	$a_{0,0} x b_{0,1}$
$a_{1,0} \times b_{1,0}$	
$a_{0,0} \land D_{0,0}$	

• Cycle 3



• Cycle 4



 $c_{0,0} = a_{3,0} \times b_{3,0+} \\ a_{2,0} \times b_{2,0+} \\ a_{1,0} \times b_{1,0+} \\ a_{0,0} \times b_{0,0}$

• Cycle 5

- New inputs can start to be broadcasted



• Cycle 6



• Cycle 7



• Cycle 10: all $C_{i,j}$ results have been outputted



Outline

- Part I: the clear need for specialized hardware
 - Energy efficiency of hardware accelerators
- Part II: accelerating GEMM
 - DNN kernels with a focus on convolution
 - Speeding up the GEMM kernel
 - Designing hardware accelerators for GEMM/CONV
 - Available accelerators for DNNs
- Part III: computing at the right precision
 - A bit of arithmetic
 - Inference and training with low precision

Accelerators for ML



CPU	GPU	FPGA	TPU
Threads SIMD	Massive Threads SIMD HBM	LUTs DSP BRAM	MM Unit BRAM

- Accelerators: GPU, TPU
- Open-source accelerators: NVDLA, Gemmini
- FPGA: overlays, dataflow (e.g., FINN)



• NVIDIA Volta GV100 (2017, 14nm)



- NVIDIA Hopper GH100 (2022, 4nm)
 - 144 SMs



• Volta GV100: SM's details

	SM															
	L1 Instruction Cache															
	L0 Instruction Cache								_							
	Warp Scheduler (32 thread/clk)							Warp Scheduler (32 thread/clk)								
	Dispatch Unit (32 thread/clk)							Dispatch Unit (32 thread/clk)								
	Register File (16,384 x 32-bit)						Register File (16,384 x 32-bit)									
FP64	INT	INT	FP32	FP32	F	H	HHH.		FP64	INT	INT	FP32	FP32	H		
FP64	INT	INT	FP32	FP32	-				FP64	INT	INT	FP32	FP32			
FP64	INT	INT	FP32	FP32	Ħ				FP64	INT	INT	FP32	FP32			
FP64	INT	INT	FP32	FP32	TEN	SOP	TENSOR		FP64	INT	INT	FP32	FP32	TEN	0.P	TENSOR
FP64	INT	INT	FP32	FP32	CC	DRE	CORE		FP64	INT	INT	FP32	FP32	CO	RE	CORE
FP64	INT	INT	FP32	FP32	Ħ				FP64	INT	INT	FP32	FP32			
FP64	INT	INT	FP32	FP32	Ħ				FP64	INT	INT	FP32	FP32			
FP64	INT	INT	FP32	FP32	-				FP64	INT	INT	FP32	FP32			
LD/ LD/ ST ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	SFU		LD/ LD/ ST ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	SFU
	L0 Instruction Cache						Warp Scheduler (32 thread/clk)									
	War	m Sch	odulo	r (32 f	hread	(clk)				Wai	m Sch	olulo	r (32 t	hread/	clk)	
	War Di:	r <mark>p Sch</mark> spatcl	edule 1 Unit	r (32 t (32 th	hread read/o	/clk) :lk)				Wa Di	r <mark>p Sch</mark> spatcl	iedule h Unit	r (32 t (32 th	hread/ read/c	cik) ik)	
	War Di: Reg	rp Sch spatcl lister	edule 1 Unit File ('	r (32 t (32 th 16,38+	hread read/o 4 x 32	/clk) clk) 2-bit)				Wa Di Reg	rp Sch spatcl jister	edule h Unit File ('	r (32 t (32 th 16,384	hread/ read/c 4 x 32	clk) lk) -bit)	
FP64	War Di: Reg	rp Sch spatcl lister INT	File (' FP32	r (32 t (32 th 16,384 FP32	hread read/c 4 x 32	/clk) :lk) 2-bit)			FP64	Wai Di Reg	rp Sch spatc jister INT	File (* File (*	r (32 t (32 th 16,384 FP32	hread/ read/c 4 x 32	clk) lk) -bit)	
FP64 FP64	War Di: Reg INT INT	rp Sch spatcl jister INT INT	File (* FP32 FP32	r (32 t (32 th 16,384 FP32 FP32	hread read/ 4 x 32	/clk) slk) 2-bit)			FP64 FP64	Wai Di Reg INT	rp Sch spatcl lister INT INT	File (* File (* FP32 FP32	r (32 t (32 th 16,384 FP32 FP32	hread/ read/c 4 x 32	clk) lk) -bit)	
FP64 FP64 FP64	War Di Reg INT INT	rp Sch spatcl jister INT INT	File (' FP32 FP32 FP32	r (32 t (32 th 16,38 FP32 FP32 FP32	hread read/ 4 x 32	/clk) slk) 2-bit)			FP64 FP64 FP64	Wan Di Reg INT INT	rp Sch spatc lister INT INT	edule h Unit File (* FP32 FP32 FP32	r (32 t (32 th 16,384 FP32 FP32 FP32	hread/ read/c 4 x 32	clk) lk) -bit)	
FP64 FP64 FP64 FP64	Wan Di: Reg INT INT INT	rp Sch spatcl jister INT INT INT	File (* FP32 FP32 FP32 FP32 FP32	r (32 t (32 th 16,38 FP32 FP32 FP32 FP32	hread read/d 4 x 32 TEN	/clk) slk) 2-bit)	TENSOR		FP64 FP64 FP64 FP64	War Di Reg INT INT INT	rp Sch spatc lister INT INT INT	FP32 FP32 FP32 FP32 FP32 FP32	r (32 t (32 th 16,384 FP32 FP32 FP32 FP32	hread/ read/c 4 x 32 4 x 32	clk) lk) -bit) SOR	TENSOR
FP64 FP64 FP64 FP64 FP64	War Di Reg INT INT INT INT	rp Sch spatcl jister INT INT INT INT	edule 1 Unit File (* FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,38 FP32 FP32 FP32 FP32 FP32	tead/ t x 32	/clk) clk) 2-bit) ISOR DRE	TENSOR		FP64 FP64 FP64 FP64 FP64	Wal Di Reg INT INT INT INT	INT INT INT INT INT INT	FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 t (32 th 16,38 FP32 FP32 FP32 FP32 FP32	tead/ read/c 4 x 32	cik) ik) -bit) SOR RE	TENSOR
FP64 FP64 FP64 FP64 FP64 FP64	War Di Reg INT INT INT INT INT	INT INT INT INT INT INT INT	File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tread/ 4 x 32	/clk) 2lk) 2-bit) ISOR DRE	TENSOR		FP64 FP64 FP64 FP64 FP64 FP64	Wall Di Reg INT INT INT INT INT	INT INT INT INT INT INT INT	File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 t (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32	hread/c 4 x 32 TEN: CO	clk) lk) -bit) SOR	TENSOR
FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wan Di Reg INT INT INT INT INT INT	rp Sch spatcl jister INT INT INT INT INT	edule File (' FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32	hread/4 4 x 32 TEN CC	/clk) :lk) ?-bit) !SOR DRE	TENSOR		FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wai Di Reg INT INT INT INT INT	PP Sch spatcl sister INT INT INT INT INT INT	File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tread/c 4 x 32	cik) ik) -bit) SOR RE	TENSOR
FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	War Di: Reg INT INT INT INT INT INT	rp Sch spatcl ister INT INT INT INT INT INT INT	edule Junit File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	 (32 th (32 th	hread/ read/ 4 x 32	/clk) clk) 2-bit) ISOR DRE	TENSOR		FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wai Di Reg INT INT INT INT INT INT	P) Sef spatcl ister INT INT INT INT INT INT INT	FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	hread/ read/c 4 x 32	cik) ik) -bit) SOR RE	TENSOR
FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Water Discrete Sector INT INT INT INT INT INT INT INT INT INT	P) Schr spatcl ister INT INT INT INT INT INT INT INT INT INT	edule Junit File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	 (32 th (32 th	tread/ tread/ 4 x 32 TEN CC	ICIK) ISIK) 2-bit) ISOR DRE	TENSOR CORE		FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Waa Di Di NT NT NT NT NT NT NT NT NT	P) Self spatcl ister INT INT INT INT INT INT INT INT INT	FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tread/c tailed tai tai tai tai tai tai tai tai tai tai	clk) lk) -bit) SOR RE	TENSOR CORE
FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Water Distribution Regg INT INT INT INT INT INT INT INT LDT	P) Scholar spatch ister INT INT INT INT INT INT INT INT INT	edule Junit File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	(32 th (32 th (32 th (32 th (32 th (32 th (32 th (32 th)) (32 th)) (32 th) (32	tread/or 4 x 32 TEN CC	rcik) cik) 2:bit) SSOR DRE	TENSOR CORE SFU	che /	FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Waa Di Reg INT INT INT INT INT INT INT INT INT INT	P Soft spatcl ister INT INT INT INT INT INT INT INT	edule h Unit File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tread/c tax 32 TEN: CO	clk) lk) -bit) SOR RE LD/ ST	TENSOR CORE
FP64 FP64 FP64 FP64 FP64 FP64 FP64 ED1 LD1 LD1 LD1	Wan Di- Reg INT INT INT INT INT INT INT INT INT	P) Scholar spatch jister INT INT INT INT INT INT INT INT INT	edule Junit File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	(32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tread/ tread/ t x 32 TEN CC	rcik) cik)	TENSOR CORE SFU	- - - - - - - - - - - - - - - - - - -	FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Waa Di Di NT NT NT NT NT NT NT NT NT NT NT NT NT	P Sch spatc ister INT INT INT INT INT INT INT INT INT	File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	hread/c read/c 4 x 32 TEN: CO	clk) lk) -bit) SOR RE LD/ ST	TENSOR



NVIDIA Tensor Cores

- Mixed-Precision Matrix Math
 - 4x4 matrices operations in one cycle



Tensor Core Throughputs

Multiply-Accumulates per clock per SM

(multiply by 2x for ops counts)

	FP32	FP16	INT8	INT4	INT1
Volta	64	512			
Turing	64	512	1,024	2,048	8,192

NVIDIA GPU

• Highest Peak Performance

Peak performance (TFLOPS)								
Device	Year	fp64	fp32	tfloat32	fp16	bfloat16	fp8	
P100	2016	5	9		19	140		
V100	2017-2019	8	16	•	125			
A100	2020-2021	19	19	156	312	312		
H100	2022	48	48	400	800	800	1600	

• A new chip in less than two years now...



• GPU have high control overhead

Half-precision Fused Multiply-Add Half-precision 4-way Dot Product H.-p. Matrix Multiply and Acc.

Operation	Energy**	Overhead*
HFMA	1.5pJ	2000%
HDP4A	6.0pJ	500%
HMMA	110pJ	27%

- But also have very efficient tensor cores
- Last generation: Hopper GH100 GPU
 - 144 SMs with 128 FP32 cores, 64 FP64 cores, 64 INT32 cores, and four Tensor Cores per SM
 - Support for FP8 format
 - 4nm, 700 Watts...

TPU: Tensor Processing Unit (Google)

- TPUv1
 - 2016, 28nm, 700MHz
 - 8GB DDR3, 28MB on-chip mem
 - 75W, 23 TOPS
- TPUv3
 - 2018, 16nm, 940 MHz
 - 32GB HBM, 32 MB
 - 450W, 92 TOPS





TPU

- Coarse-grained matrix multiply and data read/write instructions
- Compute intensive
 - 64K MACs per cycle
- Memory intensive
 - 4 MB of on-chip Accumulator Memory
 - 24 MB of on-chip Unified Buffer (activation memory)
 - Two 2133MHz DDR3 DRAM, 8 GiB of off-chip (weight DRAM)
- >25X as many MACs vs GPU



TPU

- Matrix Multiply Unit: 256x256 (65,536) 8-bit MAC as a systolic array
 - Peak: 92 TOPS
 - 2x65,536x700MHz
- Datapath:
 - Parallel-K: systolic accumulation
 - Parallel-N: systolic multicast
- Memory
 - Custom systolic registers
 - Dedicated accumulation and weight buffers
 - Double-buffered, weightstationary dataflow



Accelerators for ML



CPU	GPU	FPGA	TPU
Threads SIMD	Massive Threads SIMD HBM	LUTs DSP BRAM	MM Unit BRAM

- Accelerators: GPU, TPU
- Open-source accelerators: NVDLA, Gemmini
- FPGA: overlays, dataflow (e.g., FINN)

NVDLA

- NVIDIA Deep Learning Accelerator
 - 8-16 bit datapath, weight compression



Gemmini

https://github.com/ucb-bar/gemmini

- UC Berkeley
 - open-source (Chisel), systolic array



Gemmini

• Weight-Stationary or Output-Stationary dataflow



Accelerators for ML



CPU	GPU	FPGA	TPU
Threads SIMD	Massive Threads SIMD HBM	LUTs DSP BRAM	MM Unit BRAM

- Accelerators: GPU, TPU
- Open-source accelerators: NVDLA, Gemmini
- FPGA: overlays, dataflow (e.g., FINN)

FPGA

- Xilinx Alveo U55C card
 - PCIe[®] Gen3x16 or dual Gen4x8
- Ethernet 2 x 100 Gb/s
- XCU55 UltraScale+ FPGA
 - 16 GB High-Bandwidth Memory (HBM2), 460GB/s bandwidth
 - 1.3M CLBs
 - 270+70M BRAM
 - 9K DSP blocks, 4 INT8 MAC/DSP

D -

- 28 TOPS INT8 (peak, 800MHz)





27 x 18 Multiplier

Pre-adde

HLS vs. Overlays

- HLS synthesizes C-like high-level design and performs code transformations and synthesis optimizations
- FPGA overlay is a coarse-grained design abstraction layer over fine-grained FPGA resources



HLS vs. Overlays

- HLS synthesizes C-like high-level design and performs code transformations and synthesis optimizations
 - FINN (Xilinx): SIMD + parallel Matrix Vector Activate Unit
 - HLS4ML: "true" HLS style + reuse
 - HLS Libraries: e.g., Vitis Accelerated Blas Library, GEMM_HLS, HLS_LIB, AC_ML
- FPGA overlay is a coarse-grained design abstraction layer over fine-grained FPGA resources
 - Xilinx Deep Learning Processing Unit (DPU) + Vitis Al software stack + runtime (XRT)
 - VTA+TVM

107

TVM+VTA

- End-to-end hardware-software DL system stack
 - hardware design
 - sequential architecture
 - drivers, JIT runtime
 - optimizing compiler stack based on TVM



FINN

- Experimental framework from Xilinx Research Labs
 - Provides an HLS library of standard NN layers
 - Design space exploration of QNN accelerators on FPGAs
 - Generates HLS code that supports a wide range of precisions
- Inference only, focused on Quantized NNs
 - Pipelined dataflow architecture
 - AXI controlled dataflow structure for low latency




Outline

- Part I: the clear need for specialized hardware
 - Energy efficiency of hardware accelerators
- Part II: accelerating GEMM
 - DNN kernels with a focus on convolution
 - Speeding up the GEMM kernel
 - Designing hardware accelerators for GEMM/CONV
 - Available accelerators for DNNs
- Part III: computing at the right precision
 - Number representation and precision
 - Inference and training with low precision

Computing just right! (with the right accuracy and precision)

Number Representations and Precision

• Energy, delay, and area vary a lot between numeric formats and word-length

	Addition	Multiplication	
8-bit integer	0.03pJ / 36µm²	0.2pJ / 282µm²	
32-bit float	0.9pJ / 4184µm²	3.7pJ / 7700µm ²	



High-precision computations often lead to inefficiency

Resilience of ANN?

Aoccdrnig to a rscheearch at Cmabrigde Uinervtisy, it deosn't mttaer in waht oredr the Itteers in a wrod are, the olny iprmoatnt tihng is taht the frist and Isat Itteer be at the rghit pclae. And we spnet hlaf our Ifie Iarennig how to splel wrods. Amzanig, no! [O. Temam, ISCA10]

- Our biological neurons are tolerant to computing errors and noisy inputs
- Quantization of parameters and computations provides benefits in throughput, energy, storage

Number Representations

• Floating-Point (FIP)

 $x = (-1)^s \times m \times 2^{e-127}$

s: sign, m: mantissa, e: exponent



- Easy to use
- High dynamic range
- Hardware cost and power

• Fixed-Point (FxP)

$$x = p \times K$$

p: integer, *K*=2⁻ⁿ: fixed scale factor

- Integer arithmetic
- Efficient operators
 - Speed, power, cost
- Harder to use...



Integer part: *m* bits Fractional part: *n* bits

What can be Customized?

- Floating-Point (FIP)
- Precision
 - Exponent (E) and Mantissa(M) bit-width
 - *e* and *m* both impact accuracy
- Play with exponent bias
- Sub-normal numbers or not?
- 0, ∞, NaN?
- Rounding modes
 - stochastic, to nearest, truncation, to $0/\infty$

- Fixed-Point (FxP)
- Precision
 - Integer (m) and fractional (n) bit-width
 - *n* impacts accuracy
 - *m* impacts dynamic range
 - Wordlength (W=m+n+1)
- Rounding modes
 - stochastic, to nearest, truncation, to $0/\infty$
- Saturation modes
 - wrap, max/min

Energy Gains of Low Precision

• Multiplier (float)

• Adder (float)





Very Low-Precision float Multiplication

• Example: 7 bits (7,5,1)



How does this apply to DNNs?

Low-Precision Inference

• Not only Weights, but also Activations, Per-Layer Quantization, etc.



4-bit activations and10-bit weights keepsaccuracy near (98.4%)32-bit float reference

Even Worse for Training...

• Carbon footprint of DNN training

Analyzing the carbon footprint of current natural-language processing models shows an alarming trend: **training one huge model for machine translation emits the same amount of CO2 as five cars in their lifetimes (fuel included)**

[Strubell et al., ACL 2019]

- Many more operations than inference
- More pressure on memory access
- Much more difficult to accelerate

Need for a Significant Reduction of the Carbon Footprint of Neural Network Training Hardware

Mixed Precision DNN Training

- GEMMs and weight updates are performed in custom precision
- GEMMs are performed on GPU or FPGA kernels
- Different arithmetic configurations for FWD, BWD and WU operations



Mixed Precision DNN Training

- We can reach baseline accuracy using lower precision
- Too low of a precision could lead to divergence

• A Performance-Precision equilibrium must be found



Putting It All Together

• MPtorch-FPGA: a Custom Mixed Precision Framework for FPGA-based DNN Training



Results

- Optimize <N,M,C,F> configuration for a given DNN model
- Explore model accuracy for different arithmetic configurations
 - FP8 multiplier, FP12 (E6M5) with Stochastic Rounding
 - with less than 2% accuracy loss w.r.t FP32

	Multiplier	Accumulator	LeNet5	ResNet20	VGG16 ‡	ResNet50
		E6M5-RZ	97.10	10.00	10.00	10.00
		E6M5-RO	98.00	10.00	10.00	10.00
	E5M2-NR	E6M5-RN	98.61	10.00	10.00	10.00
FP8/FP12 SR		E6M5-SR	99.00	90.55	88.99	80.88
		E5M10-RN	99.05	91.24	89.81	82.97
FP32 baseline	E8M23-RN	E8M23-RN	99.18	91.91	90.67	82.92
	FXP4, 4-RN		99.06	10.00	10.00	10.00
	FXP4, 4-SR		99.14	10.00	10.00	10.00
	FXP4, 4-RZ	FXP8, 8	98.85	10.00	10.00	10.00
	FXP4, 4-RO		10.00	10.00	10.00	10.00
		Datasets: †MNI	ST. [‡] CIFAI	R10, ^Δ Image	woof	

What's next?

- Efficiency of hardware specialization
 - Domain-specific architectures and languages
- Computing just right
 - @design-time or @run- time
- Hardware-aware optimizations are mandatory
 - Deep knowledge of the hardware is required to propose energy-efficient DNN models

Pruning

Network Pruning



• Does pruning always translate into energy savings?

 Solution: structured pruning



Backup Slides

Key Metrics

- Accuracy
 - Evaluate using the appropriate DNN model and dataset
- Programmability
 - Support multiple applications
- Throughput / Latency
 - GOPS, frame rate, delay
- Energy / Power
 - Energy per operation and memory access
- DRAM Bandwidth
- Area Cost (memory size, # of cores)
- Number of FLOPs, MACS, and Weights are not a good proxy for energy and latency















[Eyeriss tutorial]



• Small configuration

	INT8 MACs (# instances) Conv. Buffer (KB) Area (mm2) Memer BW (GB/	Aroa	Memory	ResNet50			
(# instances)		BW (GB/s)	Perf (frames/s)	Power (mW)	Power Eff. (DL TOPS/W)		
2048	512	3.3	20	269	388	5.4	
1024	256	1.8	15	153	185	6.3	
512	256	1.4	10	93	107	6.8	
256	256	1.0	5	46	64	5.6	
128	256	0.84	2	20	41	3.8	
64	128	0.55	1	7.3	28	2.0	

https://en.wikichip.org/wiki/nvidia/microarchitectures/nvdla

TSMC 16 nm process at 1 GHz 129



• Large configuration

Configuration				ResNet50			
INT16/FP16	512 MACs	Data Type	Data Internal Type RAM Size	Perf	Power (mW)	Power Eff. (DL TOPS/W)	
INT8	1024 MACs	.,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,		(frames/s)			
Conv Buffer	256 KB	INT8	none	165	267	4.8	
Area	2.4 mm2	FP16	none	59	276	1.6	
DRAM BW	15 GB/s	INT8	2M	230	348	5.1	
TCM R/W BW	25/25 GB/s	FP16	2M	115	475	1.9	

https://en.wikichip.org/wiki/nvidia/microarchitectures/nvdla

TSMC 16 nm process at 1 GHz 130

FPGA: DSP

• DSP Double Data Rate (DDR) technique



Deep Learning Processing Unit (DPU)

- DPU: programmable engine optimized for DNNs
- Includes
 - high-performance scheduler
 - hybrid computing array module
 - instruction fetch unit module
 - global memory pool module
- Uses a specialized instruction set
- Sequential architecture model



APU - Application Processing Unit 132

Vitis Al Stack

- Vitis Al
 - Al Compiler
 - Al Quantizer
 - Al Optimizer
 - Al Profiler
 - Al Library
 - Xilinx Runtime
 Library (XRT)



User Application

Vitis Al Stack



xDNN processing engine

- **instruction memory**
 - execution controller
 - element-wise processing units
 - systolic array







xfDNN Inference Toolbox

Graph Compiler



 Python tools to quickly compile networks from common
 Frameworks – Caffe, MxNet and Tensorflow

Network Optimization



 Automatic network optimizations for lower latency by fusing layers and buffering on-chip memory



 Maintains 32bit accuracy at 8 bit within 2%

Mixed Low Precision



Fixed low-precision quantization already showed competitive results.

Next generation: **Variable** precision of activation/weights among layers









*accuracy drop less than 1%

BW	2	3	4	5	6	7	8
wgt	0	3	4	6	0	0	3
act	0	0	0	2	5	10	5

BW	2	3	4	5	6	7	8
wgt	0	0	3	22	17	10	2
act	0	0	0	16	41	13	3

BW	2	3	4	5	6	7	8
wgt	0	0	0	15	84	38	13
act	0	0	0	0	6	84	99

Preliminary experiments on popular networks. (vgg-16,resNet-50,inceptionv4)

VTA: Versatile Tensor Accelerator

- VTA: an open, generic, and customizable DL accelerator
 - Sequential architecture model
- Complete TVMbased compiler stack



TVM: Open Deep Learning Compiler Stack

- Apache TVM is an open source machine learning (ML) compiler framework for CPUs, GPUs, and accelerators
 - optimize and run computations efficiently on any hardware backend



FINN Matrix Vector Activate Unit

• Two versions: block and streaming



- PEs based on SIMD units
- #PE and #SIMD/PE are configurable

Verilog

wrapper

weight memory

FINN Matrix Vector Activate Unit



PE = 4 # SIMD = 4

PE = 4 # SIMD = 2

4 clock cycles

HLS sources: MVAU, MAC

HLS4ML

- hls4ml converts DNN models into a full HLS project
- Supports many kinds of models (MLP, CNN, RNN) and different formats (Keras/Tensorflow, Pytorch, ONNX)

 Lacks good optimization (resource usage), models needs to be quantized aggressively

Limited board,
 FPGA, recent HLS
 compiler support



Integer Addition


Integer Multiplication

 Many implementation depending on compression scheme and final adder



J.-M. Muller et al., Handbook of floating-point arithmetic, Springer, 2009.

Fixed-Point Representation

Representation

 (w,m,n)
 w=m+n



• Examples

97 = 01100001₂ in (8,2,6) 1+1/2+1/64 = 1.515625

97 = 01100001_2 in (8,-2,10) $2^{-4}+2^{-5}+2^{-10}$ 0.0947265625

Floating-Point Multiplication



Floating-Point Multiplication

- Representation (W,E,M)
 - Exponent e on E bits
 - Mantissa *m on M* bits

- Floating-point hardware is doing the job for you!
 - FIP operators are therefore more complex than FxP





Floating-Point Addition



- swap to have $e_y \le e_x$
- determine if effective subtraction
- calculate $e_x e_y$
- shift mantissa m_y accordingly
- cancellation may occur only if $e_x e_y < 2$
- normalization (LZC/shift)
- rounding

 $e_x - e_y \ge 2$ (far path)

$$e_x - e_y < 2$$
 (close path)





Floating-Point Addition

- Representation (W,E,M)
 - Exponent e on E bits
 - Mantissa *m on M* bits

$$s_z = s_x \text{ XOR } s_y$$
$$x + y = (-1)^{s_x} \cdot \left(|x| + (-1)^{s_z} \cdot |y| \right)$$

Arithmetic Support in Latest Chips?

- Hopper GH100 GPU from Nvidia
 - FP8 support in tensor cores provides up to 4x speedup





Allocate 1 bit to either range or precision Support for multiple accumulator and output types