Advanced Architectures

Master Informatics Eng.

2020/21 A.J.Proença

SIMD extensions at Intel, AMD & ARM

(some slides are borrowed)

AJProença, Advanced Architectures, MiEI, UMinho, 2020/21

公



3

Vector architectures

SIMD Parallelism

- Read sets of data elements (gather from memory) into "vector registers"
- Operate on those registers
- Store/scatter the results back into memory

SIMD extensions to scalar architectures

- Have limitations, compared to vector architectures (before AVX...)
- Number of data operands encoded into op code
- No sophisticated addressing modes (strided, scatter-gather)
- No mask registers

• Graphics Processor Units (GPUs) (next set of slides)



SIMD Implementations

- Intel implementations:
 - MMX (1996)
 - Eight 8-bit integer ops or four 16-bit integer ops
 - Streaming SIMD Extensions (SSE) (1999)
 - Eight 16-bit integer ops
 - Four 32-bit integer/fp ops or two 64-bit integer/fp ops
 - Advanced Vector eXtensions (AVX) (2010...)
 - Eight 32-bit fp or four 64-bit fp ops (integers only in AVX-2)
 - 512-bits wide in AVX-512 (and also in Larrabee & Phi-KNC)
 - Operands <u>must / should be in consecutive and</u> <u>aligned</u> memory locations
- AMD Zen/Epyc (Opteron follow-up): up to AVX-2
- ARMv8 (64-bit) architecture: NEON & SVE



Registers for vector processing in Intel 64



Intel evolution to the AVX-512



Intel SIMD ISA evolution



The AVX-512 across Intel devices

公

Microarchitecture	Support Level													
	F	CD	ER	PF	BW	DQ	٧L	IFMA	VBMI	4FMAPS	VNNI	4VNNIW	VPOPCNTDQ	BF16
Knights Landing	~	~	~	~	×	×	×	×	×	×	×	×	×	×
Knights Mill	~	~	~	~	×	×	×	×	×	~	~	~	~	×
Skylake (server)	~	~	x	x	~	~	~	×	×	×	×	×	×	×
Cascade Lake	~	~	x	x	~	~	~	×	×	×	~	×	×	×
Cannon Lake	~	~	×	x	~	~	~	~	~	×		×	×	×
Ice Lake (server)	~	~	×	x	~	~	•	~	~	~	~	×	×	×
Cooper Lake	~	~	×	x	~	~	~	×	×	×	~	×	×	~

AVX512F - <u>AVX-512 Foundation</u> AVX512CD - <u>AVX-512 Conflict Detection</u> AVX512BW - <u>AVX-512 Byte and Word</u> AVX512DQ - <u>AVX-512 Doubleword and Quadword</u> AVX512VL - <u>AVX-512 Vector Length</u> AVX512IFMA - <u>AVX-512 Integer Fused Multiply-Add</u> AVX512_VNNI - <u>AVX-512 Vector Neural Network Instructions</u> AVX512_BF16 - AVX-512 BFloat16 Instructions

"I hope AVX512 dies a painful death, and that Intel starts fixing real problems..." Linus Torvalds

The ISA chaos at Intel SIMD extensions...

Nehalem (2009), Westmere (2010): Intel Xeon	Sandy Bridge (2012): Intel Xeon Processor	Haswell (2014): Intel Xeon Processor	Knights Corner (2012): Intel Xeon Phi	Knights Landing (2016): Intel Xeon Phi	Skylake (2017): Intel Xeon Scalable Processor Family
Processors (legacy)	E3/E5 family	E3 v3/E5 v3/E7 v3 Family	Coprocessor x100 Family	Processor x200 Family	AVX-512VL
					AVX-512DQ
	Ivy Bridge (2013): Intel Xeon	Broadwell (2015): Intel Xeon		512-bit	AVX-512BW
	Processor E3 v2/E5 v2/E7 v2	Processor E3 v4/E5 v4/E7 v4		AVX-512ER	512-bit
	Family	Family		AVX-512PF	
				AVX-512CD	AVX-512CD
			512-bit	AVX-512F	AVX-512F
		256-bit	IMCI		
	256-bit	AVX2		AVX2	AVX2
128-bit	AVX	AVX		AVX	AVX
SSE*	SSE*	SSE*		SSE*	SSE*
		- primary instruction	set	- legacy instru	ction set

The AVX-512 across Intel microarchitectures



Intel AVX-512 Instruction Types					
AVX-512-F	AVX-512 Foundation Instructions				
AVX-512-VL	Vector Length Orthogonality : ability to operate on sub-512 vector sizes				
AVX-512-BW	512-bit Byte/Word support				
AVX-512-DQ	Additional D/Q/SP/DP instructions (converts, transcendental support, etc.)				
AVX-512-CD	Conflict Detect : used in vectorizing loops with potential address conflicts				

AJProença, Advanced Architectures, MiEI, UMinho, 2020/21



Intel Advanced Matrix Extension (AMX)

AMX instructions in Accelerator 1



NEON: the SIMD extension in ARM

公

What is NEON?

- NEON[™] is a wide SIMD data processing architecture
 - Extension of the ARM® instruction set
 - 32 registers, 64-bit wide
 (AArch64: 32 registers, 128-bit wide)
- NEON Instructions perform "Packed SIMD" processing
 - Registers are considered as <u>vectors</u> of <u>elements</u> of the same <u>data type</u>
 - Data types can be: signed/unsigned 8-bit, 16-bit, 32-bit, 64-bit, single precision float (AArch64: Double precision float)
 - Instructions perform the same <u>operation</u> in all lanes





NEON & FP registers: the move to ARMv8 (64-bit)

\sim

3

Changes in AArch64

- More registers
 - AArch32: 16x128-bit"Q-regs" Q0-Q15
 - AArch64: 32x128-bit "V-regs" V0-V31
- 'Dual view' no longer packed
 - 32 registers of each type: S, D, V
 - Clearer mapping of overlap
- Asm language changes
 - No 'v' in mnemonics
 - Width specifier moved to register description
 - I 28-bit Q-regs renamed V-regs



Vector & FP register sizes in ARMv8 (64-bit)

公



Figure 4-10 Arrangement of floating-point values

16-bit floating-point is supported, but only as a format to be converted from or to. It is not AJProença, Adv supported for data processing operations.

Note

NEON FP registers in ARMv8 (64-bit)

~~



Figure 7-1 Divisions of the V register

ARMv8 - A Scalable Vector Extension (SVE)



The 1st implementation of SVE: Fujitsu A64FX



Beyond vector extensions

\sim

- Vector/SIMD-extended architectures are hybrid approaches
 - mix (super)scalar + vector op capabilities on a single device
 - highly pipelined approach to reduce memory access penalty
 - tightly-closed access to shared memory: lower latency
- Evolution of vector/SIMD-extended architectures
 - computing accelerators optimized for number crunching (GPU)
 - add support for matrix <u>multiply + accumulate</u> operations; why?
 - most <u>scientific, engineering, AI & finance</u> applications use matrix computations, namely the dot product: multiply and accumulate the elements in a row of a matrix by the elements in a column from another matrix
 - manufacturers typically call these extensions **Tensor Processing Unit** (TPU)
 - support for half-precision FP & 8-bit integer; why?
 - machine learning using neural nets is becoming very popular; to compute the model parameter during training phase, intensive matrix products are used and with very low precision (is adequate!)