# LLVM-CTH @ ISC 2020

Will Lovett – Principal Technical Product Owner

arm

will.lovett@arm.com

🥑 @hpc\_will

26<sup>th</sup> June 2020

# <sup>\*</sup> Supercomputer Fugaku: Fastest Supercomputer in the World

\* \* \* \* \* \* \* \* \* \* \* \* \* \* \* \* \*

	+ +	<b>500</b> The List. <b>CERTIFICATE</b>	
	+ +	is ranked	
		No. 1	
		among the World's TOP500 Supercomputers	
+	+ +	with 415.53 Pflop/s Linpack Performance	+
		Conference on June 22nd, 2020.	
	+ +		
	+ +	he the Jack Dongame Mouth Much Much	
		Erich Strohmaier Jack Dongarra Horst Simon Martin Meuer NERSC/Berkeley Lab University of Tennessee NERSC/Berkeley Lab Prometeus	
	© 2020 Arm Limited (or its affiliate	) Erich Strohmaier Uack Dongarra Horst Simon Martin Meuer NERSC/Berkeley Lab Laliversity of Tennassee NERSC/Berkeley Lab Promitieus + + +	









### Important concepts and further reading



5 © 2020 Arm Limited (or its affiliates)

### arm

### **Scalable Vectors and Predication**

- 32 scalable vector registers (Z0-Z31):
  - 128-2048 bits vector length is decided by implementation
  - Bottom 128 bits are overlaying with the Floating-point & NEON vector register bank (V0-V31)
  - Supporting 8-bit, 16-bit, 32-bit, and 64-bit element sizes



2048 bit – 128 bit

- Flags that represent whether each vector lane is active
- 1/8<sup>th</sup> size of vector registers, so can represent the lanes of 8-bit elements
- For vector elements > 8-bit, only the least-significant bit is used



arm

Vx (NEON)

Zx (SVE)

128 bit

128 bit





# Teach LLVM IR to understand scalable vectors

How long is a piece of (scalable) string?

- 1. Add the concept of scalable vectors
- 2. Fix up VectorType::getNumElements() et al
- 3. Stack frame layout with scalable sizes
- 4. Representing shuffles of scalable vectors





# LLVM IR 1/4: Add the concept of scalable vectors

### The problem

- Traditionally, vector size was known at compile-time
  - eg. <4  $\,$  x  $\,$  132> is a 128-bit vector of 4 32-bit integer values
- With SVE, vector size is now unknown at compile-time
  - Some unknown multiple of 128-bits, up to 2048 bits

- LLVM IR chooses to assume that this value is unknown but constant
  - If you break that assumption, you've broken the <u>SVE calling convention</u> and you're on your own
- Extend vector type to include vscale, which represents this unknown constant
  - eg. <vscale x 4 x i32> is a vector with a multiple of 4 32-bit integer values
  - On a 128-bit implementation of SVE, vscale = 1, and there will be 4 elements
  - On a 512-bit implementation of SVE, vscale = 4, and there will be 16 elements

# LLVM IR 2/4: Fix up VectorType::getNumElements() et al

#### The problem

unsigned getNumElements() const {return VectorTypeBits.NumElements;}

• This, and other similar accessors are used everywhere

- **Split** VectorType into FixedVectorType and ScalableVectorType
- Methodically go through all 787 calls to getNumElements ()
- Convert them to use FixedVectorType or to correctly support scalable vectors
- Kudos to Christopher Tetreault from Qualcomm!



# LLVM IR 3/4: Stack frame layout with scalable sizes

#### The problem

- Traditionally, all values on the stack have a known size in bytes
- So the offset from the stack base of every stack element is a compile-time constant
- We now need to support spilling/filling scalable vectors

- Add a scalable area to each function's stack frame
- This area only contains scalable values
- Address calculation becomes

```
• pointer + offset + offset2 * vscale
```



# LLVM IR 4/4: Representing shuffles of scalable vectors

### The problem

- LLVM needs to represent permute the elements within a vector
- Eg a zip of the even-numbered elements of two <4 x i32> vectors looks like this:
   <result> = shufflevector <4 x i32> %v1, <4 x i32> %v2, <4 x i32> <i32 0, i32 4, i32 1, i32 5> ;
- But how do we construct an index into a vector when we don't know its length?
- Other uses of shufflevector have similar issues
  - splat, reverse, concat, split\_lo/hi, zip\_lo/hi, unzip\_even/odd, splice

- Undecided upstream!
  - Current discussion is around using named shuffle patterns and/or pattern-generating IR intrinsics
- Downstream, we have added additional constants:
  - We use constant stepvector and vscale to derive all the access patterns for shufflevector







### Teach LLVM vectorizer to use SVE features

Leveraging new instructions for fun and profit

- 1. Per-lane predication
- 2. Scalar tail removal
- 3. Multiple exits and unknown tripcounts
- 4. Loops with gather/scatters
- 5. Reductions
- 6. Complex number operations



# SVE features 1/6: Per-lane predication

Allows efficient conditional execution within loops

#### 15 © 2020 Arm Limited (or its affiliates)

#### Today

- Supported in LLVM today using select to filter out inactive lanes in results
- This is incompatible with FP operations that have side-effects, such as raising an exception on FP divide-by-zero

#### Future

- First-class predication support in LLVM IR is a work is in progress in the community
- Will make it easier to support fully compliant IEEE FP exception handling



# SVE features 2/6: Scalar tail removal

Introducing WHILELO and friends

- SVE has instructions such as whilelo which test the content of a predicate vector and set condition codes
- These condition codes can then control conditional branches
- This allows us to eliminate scalar tail blocks



# SVE features 2/6: Scalar tail removal

### SVE example

Setup	<pre>foo: // <setup snipped=""></setup></pre>		<pre>void foo (double* a,</pre>				
	.LBB0 2:		<pre>unsigned n) {</pre>				
Vector loop	<pre>ld1d { z0.d }, p1/z, fadd z0.d, p0/m, z0 st1d { z0.d }, p1, incd x8 whilelo p1.d, x8, x9 b.first .LBB0 2</pre>	[x0, x8, lsl #3] d, #1.0 x0, x8, lsl #3]	<pre>for (unsigned i=0; i<n; +="1.0;" a[i]="" i++)="" pre="" }<=""></n;></pre>				
	.LBB0_3:	Advantages					
	rec	<ul> <li>Scalar tail is removed</li> </ul>					
		<ul> <li>Great for code density/icache</li> </ul>					

#### Disadvantages

• Adds a whilelo -> b.first dependency

arm

# SVE features 3/6: Multiple exits and unknown tripcounts

Allows vectorization with early exits

```
void strcpy(char *a, char *b)
```

```
{
    int i=0;
    while (a[i] != 0)
        b[i] = a[i];
}
```

- Loops with unknown tripcounts are traditionally hard to vectorize
- Vector loads or stores past the end of either array could trigger a page fault which could never occur in a scalar implementation
- First-faulting loads and stores in SVE allow a safe vector implementation
- This is unsupported by LLVM's vectorizer today
- Easy to implement for library functions using the ACLE



# SVE features 4/6: gather/scatters and structured accesses

Allows indirect memory accesses within vectorized loops

- SVE supports vector loads from a vector of (calculated) addresses
- Supported in LLVM today
- Cost modeling is overly simplistic
- InterleavedAccessPass replaces gathers/scatters with interleaved accesses, which are commonly faster
- Work is needed to update this pass to handle indexing of scalable vector types



# SVE features 5/6: Horizontal reductions

Allows efficient reduction loops

```
double sum(double* a, int length) {
  double res = 0;
  for (int i = 0; i < length; ++i)
    res += a[i];
  return res;
}</pre>
```

Two options, depending on FP requirements:

- 1. Strict ordering
  - Scalar partial result is updated in each vector loop iteration
  - Significant cost to cross-lane operation within the vector body

#### 2. Loose ordering

- -fassociative-math -fno-signed-zeros -fno-trapping-math
- Allows fast vertical reductions within loop body
- Relatively fast recursive pairwise reduction of final vector

Both options are fully supported in LLVM



# SVE features 6/6: Complex number operations

Allows efficient operations on complex numbers

- SVE supports operations on complex numbers, operating on pairs of (real, imaginary) lanes
- Can be identified using SLP-style vectorization with fixed width vectors
- This is more difficult with scalable vectors
- Some work done downstream (not yet productised)



# SVE features 7\*/6: Vector-length Specific SVE support

#### Purpose

- Allow LLVM to assume a specific SVE vector width
- Behaviour with a different vector width is undefined

#### Usage:

• -msve-vector-bits=<length>

#### ACLE

- Allows LLVM to transparently cast between SVE ACLE types and GNU C/C++ vectors
- This is necessary for storing values within structs and classes
- Likely to be used in some HPC frameworks such as Grid and Kokkos SIMD

#### Autovec

- Allows LLVM vector IR to codegen to SVE
- Likely to be used by some downstream LLVM compilers

22 © 2020 Arm Limited (or its affiliates)

\* This one isn't really a feature of SVE!





### Status and roadmap in LLVM

#### LLVM 9 (September 2019)

• SVE/SVE2 assembly and disassembly

#### LLVM 11 (September 2020)

- SVE/SVE2 intrinsic support
- Stabilisation, stack unwinding, debuginfo
- Armv8.6-a support (bfloat16, matmul)
- Vector-length specific SVE codegen

#### LLVM 12 (March 2021)

- ACLE code quality improvements (let us know!)
- Vector-length agnostic SVE vectorization of a few simple loops

#### LLVM 13 (September 2021)

• Full support for Vector-length agnostic SVE vectorization

24 © 2020 Arm Limited (or its affiliates)

### arm

### **Contributions welcome!**

- Huge thanks to the Arm partners and the wider LLVM community, especially
  - Qualcomm
  - HPE/Cray
  - Fujitsu
  - Linaro
  - Huawei
- Get involved!
  - <u>SVE/SVE2 bi-weekly sync up signup sheet</u>
  - Meeting notes from recent calls
  - Phabricator patches
  - List of identified work items



<sup>+</sup> Thảnk Yỏu
, Danke
Merci
, , 谢谢
ありがとう
+ Gracias
Kiitos
감사합니다
धन्यवाद
شکرًا
ধ্বন্যবা     দ

תודה

+ + + + +

A legale reference and the second sec

© 20ื่20 Arm Limited (or its affiliates)

a	r'n	$\mathbf{n}^{+}$				<sup>+</sup> The Arm trademarks featured in this presentation are registered trademarks or trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere. All rights reserved. All other marks featured may be trademarks of their respective owners.						
								+	+	+	+	
				www.arm.com/company/policies/trademarks								