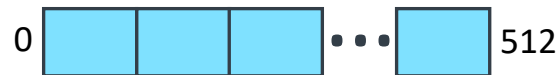**arm**

# Arm SVE Fundamentals

# Arm's Scalable Vector Extension (SVE)
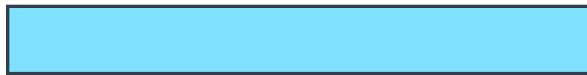
An ISA feature which Si partners can implement at length – 128 to 2048 bits

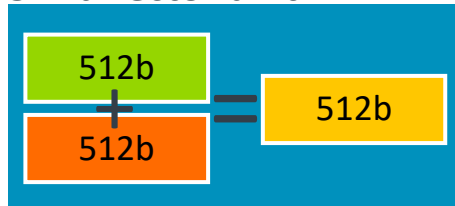## How SVE works

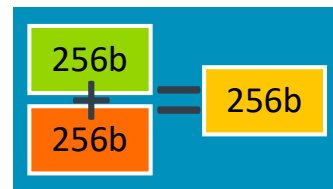The hardware sets the vector length

0 $\cdots$ 512

In software, vectors have no length

The *exact same* binary code runs on hardware with different vector lengths
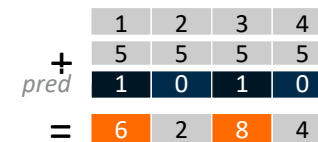
A + B = C

**512b vector unit**

512b
+
512b
=
512b

**256b vector unit**

256b
+
256b
=
256b

## SVE improves auto-vectorization

Gather-load and scatter-store

| 1 | 2 | 3 | 4 |
| 5 | 5 | 5 | 5 |

*pred* | 1 | 0 | 1 | 0 |

= | 6 | 2 | 8 | 4 |

Per-lane predication

```
for (i = 0; i < n; ++i)
```

| *INDEX i* | n-2 | n-1 | n | n+1 |
| *WHILELT n* | 1 | 1 | 0 | 0 |

Predicate-driven loop control and management

| 1 | 2 |

+ | 1 | 2 | 0 | 0 |
*pred* | 1 | 1 | 0 | 0 |

Vector partitioning and software-managed speculation

1 + 2 + 3 + 4 =
1 + 2 + 3 + 4
=         =
3    +    7    =

Extended floating-point horizontal reductions

arm

# VLA

**V**ector **L**ength **A**gnostic
programming model

Write once

Compile once

Vectorize more loops

arm

# SVE vs Traditional ISA

How do we compute data which has ten chunks of 4-bytes?

## Aarch64 (scalar)
❑ Ten iterations over a 4-byte register

## NEON (128-bit vector engine)
❑ Two iterations over a 16-byte register + two iterations of a drain loop over a 4-byte register

## SVE (128-bit VLA vector engine)
❑ Three iterations over a 16-byte **VLA register** with an adjustable **predicate**
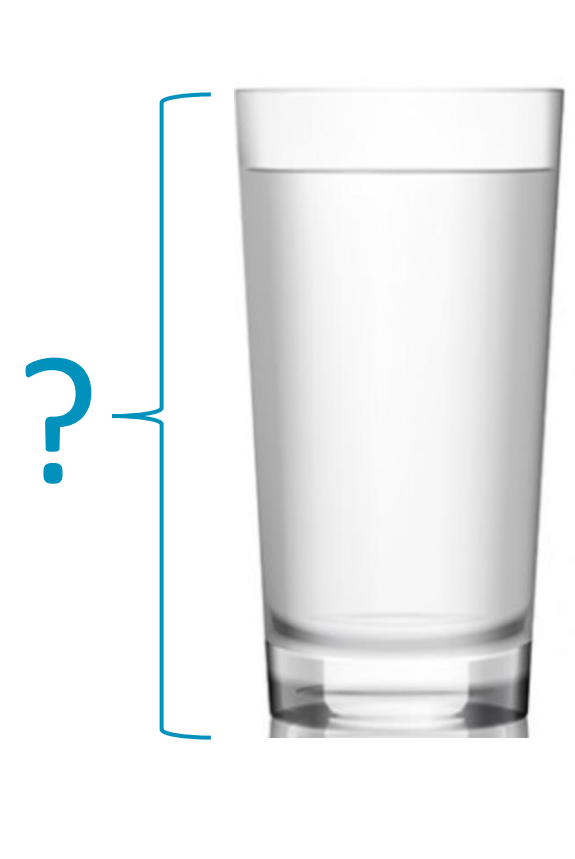
arm

# How big can an SVE vector be?

Any multiple of 128 bits up to 2048 bits, and it can be dynamically reduced.

$$(A) \quad VL = LEN \times 128$$

$$(B) \quad VL <= 2048$$

$VL$ is *implementation dependent*,

can be *reduced by the OS/Hypervisor*.

?

**arm**

# How can you program when the vector length is unknown?

SVE provides features to enable VLA programming from the assembly level and up

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **+** | 5 | 5 | 5 | 5 |
| *pred* | 1 | 0 | 1 | 0 |
| **=** | 6 | 2 | 8 | 4 |

## Per-lane predication

Operations work on individual lanes under control of a predicate register.

```
for (i = 0; i < n; ++i)
```

| *INDEX i* | n-2 | n-1 | n | n+1 |
|---|---|---|---|---|
| *WHILELT n* | 1 | 1 | 0 | 0 |

## Predicate-driven loop control and management

Eliminate scalar loop heads and tails by processing partial vectors.

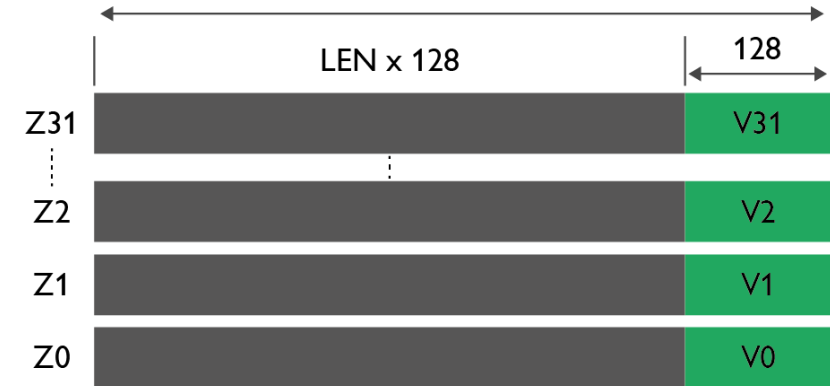| | 1 | 2 | | |
|---|---|---|---|---|
| **+** | 1 | 2 | 0 | 0 |
| *pred* | 1 | 1 | 0 | 0 |

## Vector partitioning & software-managed speculation

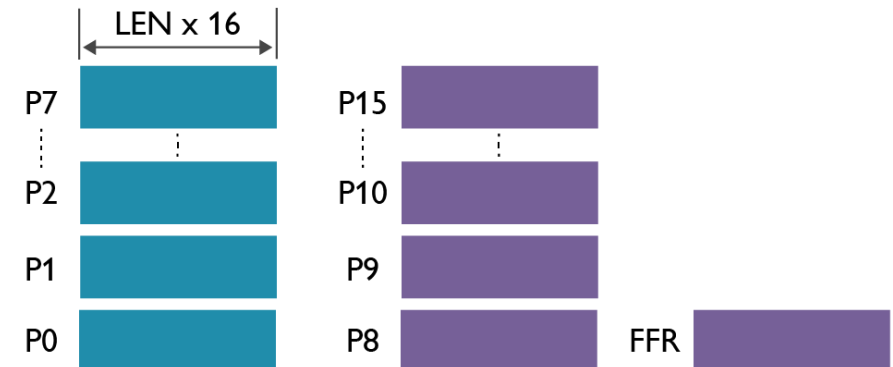First Faulting Load instructions allow memory accesses to cross into invalid pages.

**arm**

# SVE Registers

- **Scalable vector registers**

  - `Z0`-`Z31` extending NEON's 128-bit `V0`-`V31`.

  - Packed DP, SP & HP floating-point elements.

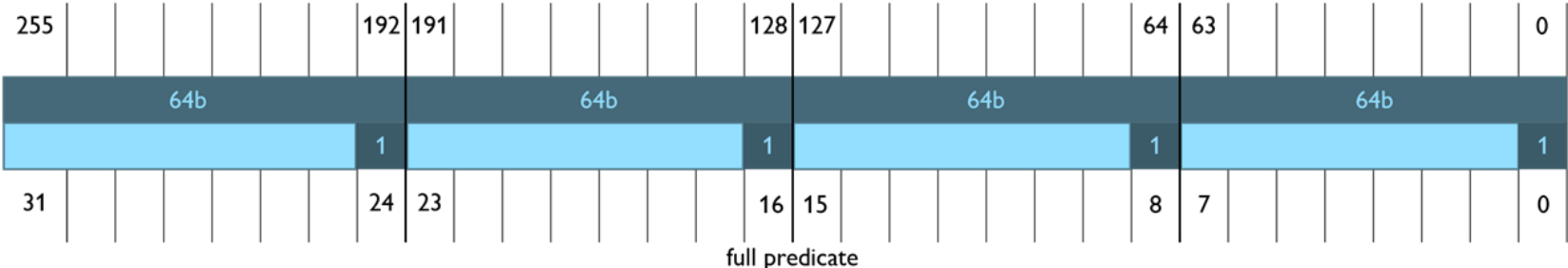  - Packed 64, 32, 16 & 8-bit integer elements.



- **Scalable predicate registers**

  - `P0`-`P7`   governing predicates for load/store/arithmetic.

  - `P8`-`P15`  additional predicates for loop management.

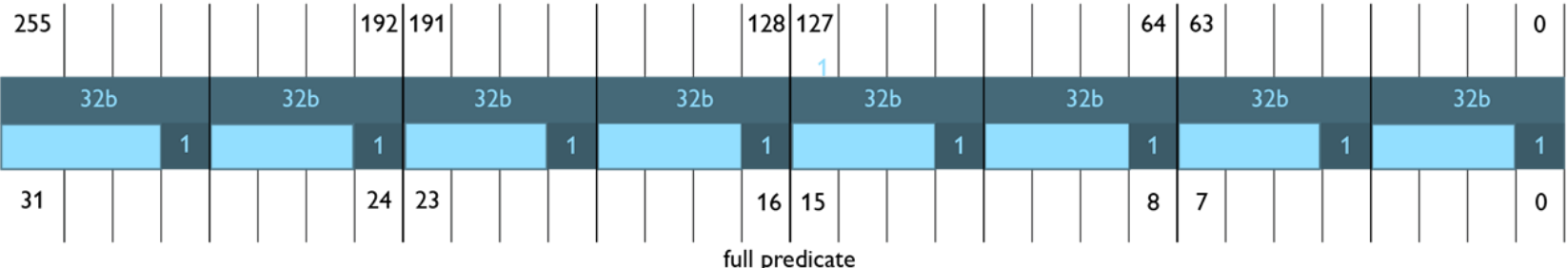  - `FFR`    first fault register for software speculation.

arm

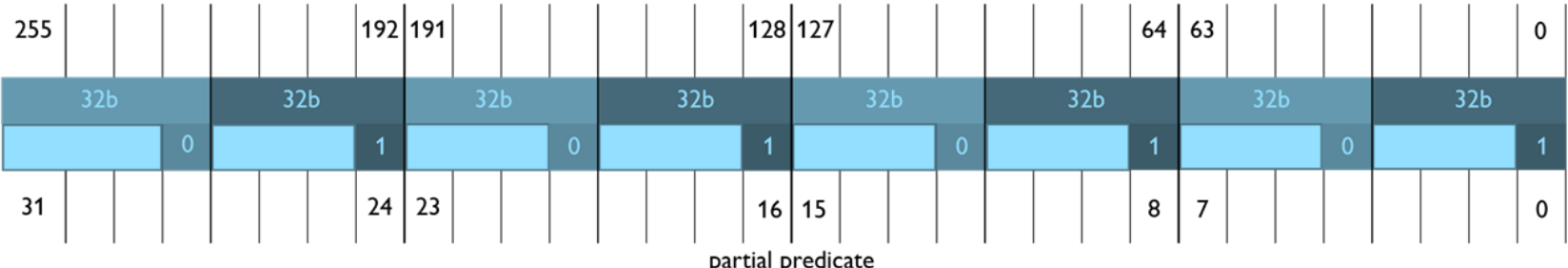# SVE vector & predicate register organization



256-bit vector, 64-bit elements

256-bit vector, packed 32-bit elements

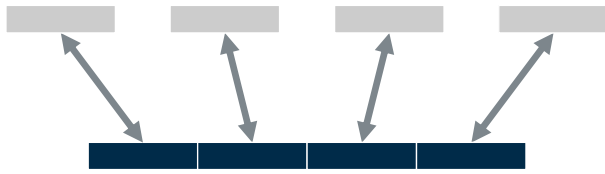256-bit vector, unpacked 32-bit elements

# VLA Programming Approaches

Don't panic!

- Compilers:
  - Auto-vectorization: GCC, Arm Compiler for HPC, Cray, Fujitsu
  - Compiler directives, e.g. OpenMP
    - #pragma omp parallel for **simd**
    - #pragma vector always

- Libraries:
  - Arm Performance Library (ArmPL)
  - Cray LibSci
  - Fujitsu SSL II

- Intrinsics (ACLE):
  - Arm C Language Extensions for SVE
  - Arm Scalable Vector Extensions and Application to Machine Learning

- Assembly:
  - Full ISA Specification: The Scalable Vector Extension for Armv8-A
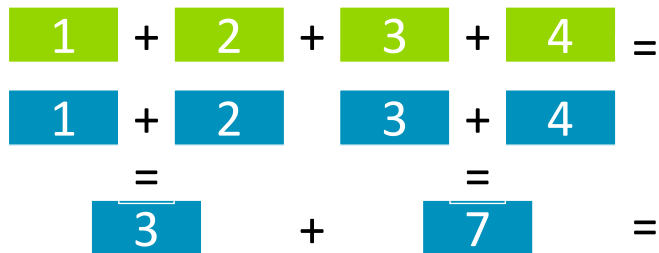
arm

# SVE supports vectorization in complex code

Right from the start, SVE was engineered to handle codes that usually won't vectorize

### Gather-load and scatter-store

Loads a single register from several non-contiguous memory locations.

$1 + 2 + 3 + 4 =$

$1 + 2 \quad 3 + 4$

$= \quad =$

$3 \quad + \quad 7 \quad =$

### Extended floating-point horizontal reductions

In-order and tree-based reductions trade-off performance and repeatability.

arm

# Portability

Is it really possible to run a vectorized application anywhere?

**Write once**: can my code **_compile_** for machines with different VL?

👍 Code that is auto-vectorized by the compiler

👍 Hand-written assembly

👍 Hand-written C intrinsics

**Compile once**: Can I take my executable and **_run it_** on machines with different VL?

👍 Self contained programs with no external dependencies

❓ But what about programs that depend on external libraries? … (spoiler: 👍 )

**arm**

# Auto-vectorize external calls: `libm` example.

`float sinf(float);`

## NEON

- Neon has 128-bit and 64-bit register split.

- The library has to provide at least 2 symbols, because it doesn't know where the auto-vec code comes from:

  - `_ZGVnN2v_sinf`
  - `_ZGVnN4v_sinf`

## SVE

- Does `libm` need to provide a symbol for each VL?

  - `_ZGVsM4v_sinf`
  - `_ZGVsM6v_sin`
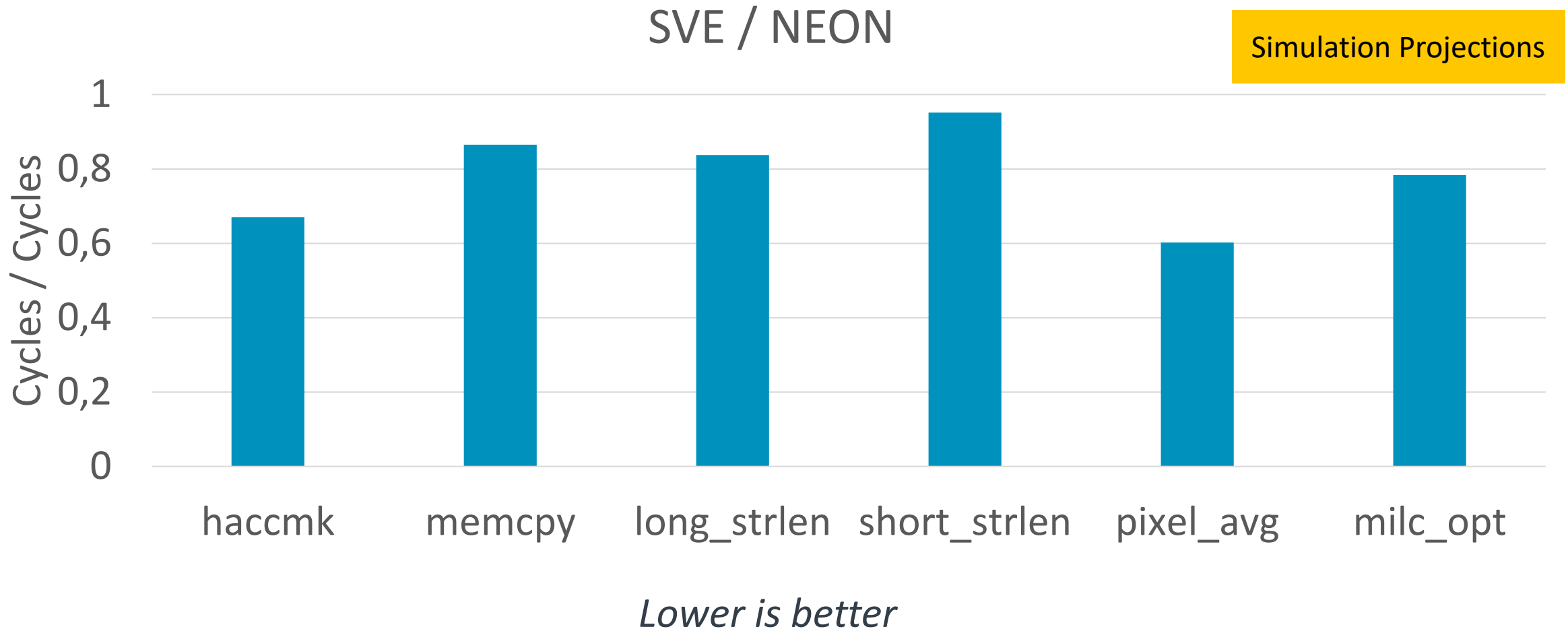  - `_ZGVsM8v_sinf`
  - `_ZGVsM10v_sinf`
  - …

- One symbol!

  `_ZGVsMxv_sinf`

arm

# Open source support

- **Arm actively posting SVE open source patches upstream**
  - Beginning with first public announcement of SVE at HotChips 2016

- **Available upstream**
  - GNU Binutils-2.28:       released Feb 2017,  includes SVE assembler & disassembler
  - GCC 8:                        Full assembly, disassembly and basic auto-vectorization
  - LLVM 7:                      Full assembly, disassembly
  - QEMU 3:                     User space SVE emulation
  - GDB 8.2                      HPC use cases fully included

- **Under upstream review**
  - LLVM:                        Since Nov 2016, as presented at LLVM conference
  - Linux kernel:             Since Mar 2017, LWN article on SVE support

**arm**

# SVE: More Powerful Vectorization on V1

SVE vectorizes more codes and makes better use of the vector units



SVE / NEON

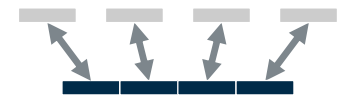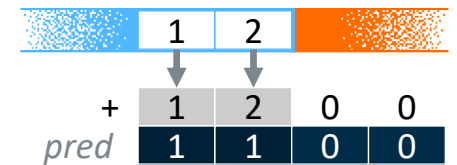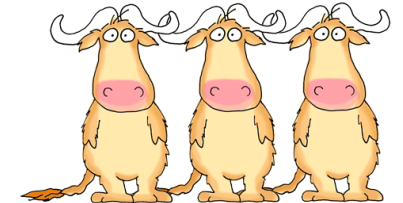Simulation Projections

*Lower is better*

arm

# Quick Recap

- SVE enables Vector Length Agnostic (VLA) programming

- VLA enables portability, scalability, and optimization

- Predicates control which operations affect which vector lanes
  - Predicates are not bitmasks
  - You can think of them as dynamically resizing the vector registers

- The actual vector length is set by the CPU architect
  - Any multiple of 128 bits up to 2048 bits
  - May be dynamically reduced by the OS or hypervisor

- SVE was designed for HPC and can vectorize complex structures

- Many open source and commercial tools currently support SVE

```
for (i = 0; i < n; ++i)
```

| INDEX i | n-2 | n-1 | n | n+1 |
|---|---|---|---|---|
| WHILELT n | 1 | 1 | 0 | 0 |

|  | 1 | 2 |  |  |
|---|---|---|---|---|
| + | 1 | 2 | 0 | 0 |
| pred | 1 | 1 | 0 | 0 |

GNU! GNU! GNU!

arm

Hands On: HACC

# 05_Apps/01_HACC

See README.md for details

- Computationally intensive part of an N-body cosmology code.

- Application performance is dominated by a long chain of floating point instructions

- Performance scales well with vector length

- FOM: Wall clock time spent in the application loop reported in seconds

```
------------------------------------------------
./hacc_arm_neon.exe 1000
Maximum OpenMP Threads: 48
Iterations: 1000
Gravity Short-Range-Force Kernel (5th Order): 12823.6 -444.108 -645.349: 5.19743 s
------------------------------------------------
./hacc_arm_sve.exe 1000
Maximum OpenMP Threads: 48
Iterations: 1000
Gravity Short-Range-Force Kernel (5th Order): 12823.6 -444.108 -645.349: 1.55594 s
------------------------------------------------
```

NEON128 on A64FX

SVE512 on A64FX

arm

# arm

Thank You
Danke
Merci
谢谢
ありがとう
Gracias
Kiitos
감사합니다
धन्यवाद
شكرًا
תודה