# Vectorization

Silvio Stanzani , Raphael Cóbe , Rogério Iope

Núcleo de Computação Científica – UNESP

**ERAD-SP 2016 - Mackenzie**
**05 August 2016**

# Agenda

- Hybrid Parallel Architectures**;**

- Memory System and Vector Processing Units;

- Intel Architectures;

- Profiling;

- Optimizing Memory Access;

- Auto Vectorization;

- Guided Vectorization;

- Examples.

# Agenda

- **Hybrid Parallel Architectures;**

- Memory System and Vector Processing Units;

- Intel Architectures;

- Profiling;

- Optimizing Memory Access;

- Auto Vectorization;

- Guided Vectorization;

- Examples.

# Hybrid Parallel Architectures

- Heterogeneous computational systems:
  - Multicore processors;
  - Multi-level memory sub-system;
  - Input and Output sub-system;
- Multi-level parallelism:
  - Processing core;
  - Chip multiprocessor;
  - Computing node;
  - Computing cluster;
- Hybrid Parallel architectures
  - Coprocessors and accelerators;

# Hybrid Parallel Architectures

- ## Heterogeneous computational systems:
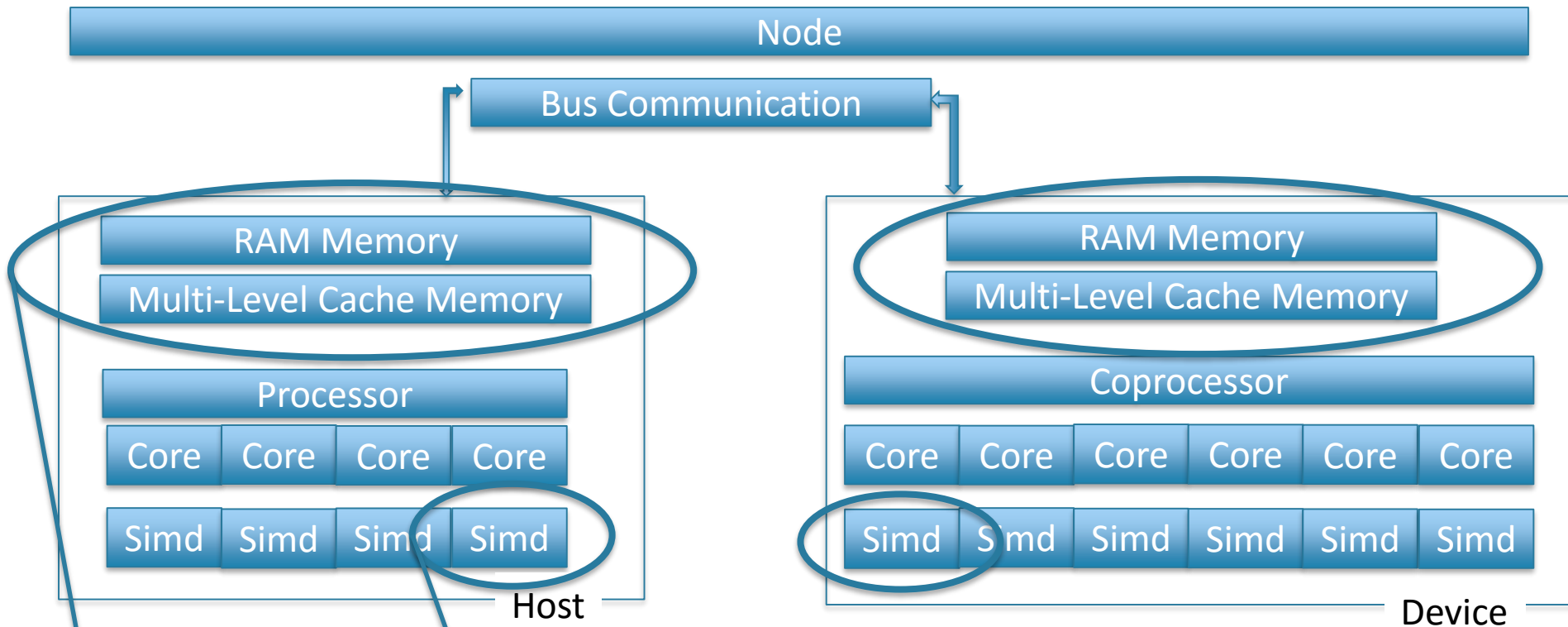
  - ### Scalar and Vector Instructions

**Vector Instructions (SIMD)**

**Scalar Instructions**

| A7 | A6 | A5 | A4 | A3 | A2 | A1 | A0 |
|----|----|----|----|----|----|----|----|

+

| B7 | B6 | B5 | B4 | B3 | B2 | B1 | B0 |
|----|----|----|----|----|----|----|----|

=

| A7+B7 | A6+B6 | A5+B5 | A4+B4 | A3+B3 | A2+B2 | A1+B1 | A0+B0 |
|-------|-------|-------|-------|-------|-------|-------|-------|

A

+

B

=

A+B

  - ### Multi-level memory

    - ❏ RAM Memory;
    - ❏ Multi-level Cache.

| Processor 1 | | | Processor 2 | | |
|---|---|---|---|---|---|
| Core 1 | Core 2 | Core N | Core 1 | Core 2 | Core N |
| L1 | L1 | L1 | L1 | L1 | L1 |
| L2 | L2 | L2 | L2 | L2 | L2 |
| L3 | | | L3 | | |
| RAM | | | | | |

# Hybrid Parallel Architectures

Node

Bus Communication

**Host**

RAM Memory

Multi-Level Cache Memory

Processor

| Core | Core | Core | Core |

| Simd | Simd | Simd | Simd |

**Device**

RAM Memory

Multi-Level Cache Memory

Coprocessor

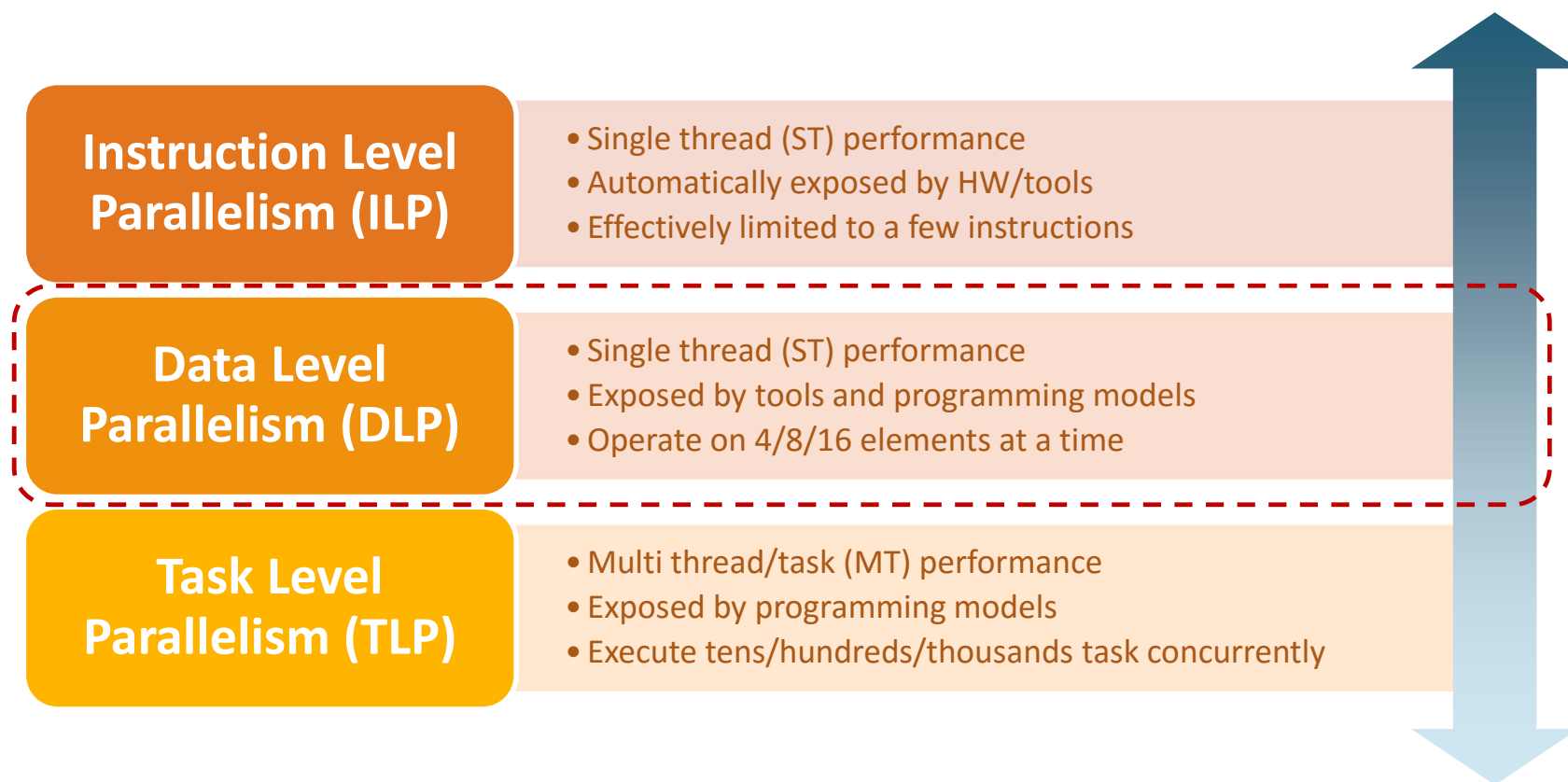| Core | Core | Core | Core | Core | Core |

| Simd | Simd | Simd | Simd | Simd | Simd |

- SIMD Instructions
- Optimized Memory Access → Vectorization!

# Don't use a single thread or vector lane

# Exploiting the parallel universe

**Instruction Level Parallelism (ILP)**
- Single thread (ST) performance
- Automatically exposed by HW/tools
- Effectively limited to a few instructions

**Data Level Parallelism (DLP)**
- Single thread (ST) performance
- Exposed by tools and programming models
- Operate on 4/8/16 elements at a time

**Task Level Parallelism (TLP)**
- Multi thread/task (MT) performance
- Exposed by programming models
- Execute tens/hundreds/thousands task concurrently

## Programmers' responsibility to expose DLP/TLP

# Agenda

- Hybrid Parallel Architectures;
- **Memory System and Vector Processing Units;**
- Intel Architectures;
- Profiling;
- Optimizing Memory Access;
- Auto Vectorization;
- Guided Vectorization;
- Examples.

# Memory System

- CPU Register: internal Processor Memory. Stores data or instruction to be executed;

- Cache: stores segments of programs currently being executed in the CPU and temporary data frequently needed in the present calculations;

- Main memory: only program and data currently needed by the processor resides in main memory;

- Auxiliary memory: devices that provides backup storage.

# Memory Hierarchy

Fast

Larger in Size

CPU Registers

Cache Memory

Main Memory

Auxiliary Memory

# Cache Memory

- Cache Memory is employed in computer systems to compensate for the difference in speed between main memory access time and processor logic.

- Operating System controls the load of Data to Cache; such load can be guided by the developer

# Cache Memory

- The Performance of cache memory is frequently measured in terms of hit ratio.

    - When the CPU refers to memory and finds the word in cache, it is said to produce a **hit**.

    - If the word is not found in cache, it is in main memory and it counts as a **miss**

# Locality

- Temporal locality: if an item was referenced, it will be referenced again <u>soon</u> (e.g. cyclical execution in loops);

- Spatial locality: if an item was referenced, items <u>close</u> to it will be referenced too (the very nature of every program – serial stream of instructions)

CPU

*Increasing distance from CPU in access time*

L1

L2

...

Main system memory

*Size of memory at each level*

# Vectorization

- Scalar Code computes this one-element at a time.

- Vector (or SIMD) Code computes more than one element at a time. SIMD stands for **S**ingle **I**nstruction **M**ultiple **D**ata.

```
float *A, *B, *C;
for(i=0;i<n;i++){
  A[i] = B[i] + C[i];
}
```

- Scalar



- SIMD

# Vectorization

- Vectorization
  - Loading data into cache accordingly;
  - Store elements on SIMD registers or vectors;
  - Apply the same operation to a set of Data at the same time;
  - Iterations need to be independent;
  - Usually on inner loops.

Scalar loop

```
for (int i = 0; i < N; i++)
    c[i] = a[i] + b[i];
```

SIMD loop (4 elements)

```
for (int i = 0; i < N; i += 4)
    c[i:4] = a[i:4] + b[i:4];
```

|      | x | y | z | w |
| a[i:4] | 1.0 | 2.0 | 3.0 | 4.0 |
|      | + | + | + | + |
| b[i:4] | 5.0 | 6.0 | 7.0 | 8.0 |
| c[i:4] | 6.0 | 8.0 | 10.0 | 12.0 |

# Past, present, and future of Intel SIMD types



Current Intel® Xeon® processors

64-bit SIMD

128-bit SIMD

512-bit SIMD

| | Exponential & Reciprocal Instructions (ERI) |
| MultiMedia eXtensions (MMX) | Prefetch Instructions (PFI) |
| Advanced Vector eXtensions (AVX) | Foundation instructions (FI) |
| | Conflict Detection Instructions (CDI) |
| AVX2 | Byte & Word Instructions (BWI) |
| | Double-/Quad-word Instructions (DQI) |
| AVX-512 | Vector Length Extensions (VLE) |

256-bit SIMD

Future Intel® Xeon Phi™ coprocessors (including Knights Landing)

Future Intel® Xeon® processors

Initial Many Core Instructions (IMCI)

Current Intel® Xeon Phi™ coprocessors (Knights Corner)

512-bit SIMD

# Intel® AVX2/IMCI/AVX-512 differences

| | Intel® Initial Many Core Instructions<br>**IMCI** | Intel® Advanced Vector Extensions 2<br>**AVX2** | Intel® Advanced Vector Extensions 512<br>**AVX-512** |
|---|---|---|---|
| **Introduction** | 2012 | 2013 | 2015 |
| **Products** | Knights Corner | Haswell, Broadwell | Knights Landing, future Intel® Xeon® and Xeon® Phi™ products |
| **Register file** | SP/DP/int32/int64 data types<br>32 x 512-bit SIMD registers<br>8 x 16-bit mask registers | SP/DP/int32/int64 data types<br>16 x 256-bit SIMD registers<br>No mask registers (instr. blending) | SP/DP/int32/int64 data types<br>32 x 512-bit SIMD registers<br>8 x (up to) 64-bit mask |
| **ISA features** | Not compatible with AVX*/SSE*<br>No unaligned data support<br>Embedded broadcast/cvt/swizzle<br>MVEX encoding | Fully compatible with AVX/SSE*<br>Unaligned data support (penalty)<br><br>VEX encoding | Fully compatible with AVX*/SSE*<br>Unaligned data support (penalty)<br>Embedded broadcast/rounding<br>EVEX encoding |
| **Instruction features** | Fused multiply-and-add (FMA)<br>Partial gather/scatter<br>Transcendental support | Fused multiply-and-add (FMA)<br>Full gather | Fused multiply-and-add (FMA)<br>Full gather/scatter<br>Transcendental support (ERI only)<br>Conflict detection instructions<br>PFI/BWI/DQI/VLE (if applies) |

# AVX-512

- AVX512-F: similar to the core feature set of the AVX2 instruction set, with the difference of wider registers, and more double precision and integer support;

- AVX512-CD ("Conflict Detection"):
- AVX512-ER ("Exponential and Reciprocal" )
- AVX512-PF ("prefetch")

# Agenda

- Hybrid Parallel Architectures;
- Memory System and Vector Processing Units;
- **Intel Architectures;**
- Profiling;
- Optimizing Memory Access;
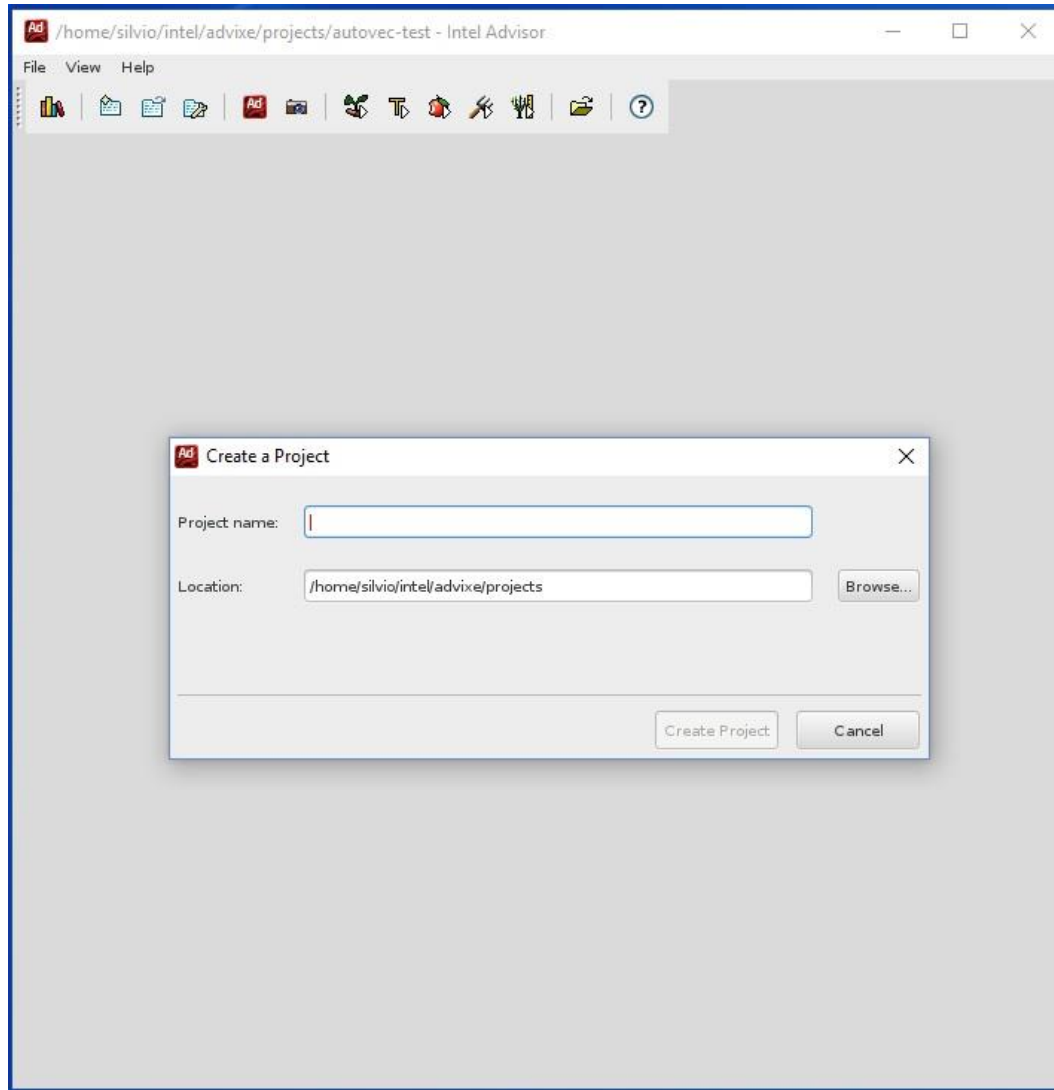- Auto Vectorization;
- Guided Vectorization;
- Examples.

# Intel® Xeon Phi™ Architecture Overview



**8 memory controllers
16 Channel GDDR5 MC
PCIe GEN2**

**Cores: 61 core s, at 1.1 GHz
in-order, support 4 threads
512 bit Vector Processing Unit
32 native registers**

**Distributed tag
directory to uniquely
map physical
addresses**

**High-speed bi-directional
ring interconnect
Fully Coherent L2 Cache**

**Reliability Features
Parity on L1 Cache, ECC on memory
CRC on memory IO, CAP on memory**

# Knights Landing (KNL)

Over 3 TF DP peak
Full Xeon ISA compatibility through AVX-512
~3x single-thread vs. compared to Knights Corner

Up to 72 cores
2D mesh architecture

Up to 16GB high-bandwidth on-package memory (MCDRAM)
Exposed as NUMA node
~500 GB/s sustained BW

2x 512b VPU per core
(Vector Processing Units)

6 channels
DDR4
Up to
384GB

**Up to 72 cores**

MCDRAM MCDRAM MCDRAM MCDRAM

DDR4 DDR4 DDR4

DDR4 DDR4 DDR4

MCDRAM MCDRAM MCDRAM MCDRAM

Wellsburg PCH

DMI

Common with Grantley PCH

HFI

Connector

Twinax Cable
Twinax Cable

Micro-Coax Cable (IFP)

Micro-Coax Cable (IFP)

PCIe Gen3 x36

2 ports
**Intel Omni-Path Fabric** On-package
50 GB/s bi-directional

## Tile

2 VPU | HUB | 2 VPU

Core | 1MB L2 | Core

Based on Intel® Atom Silvermont processor with many HPC enhancements
Deep out-of-order buffers
Gather/scatter in hardware
Improved branch predition
4 threads/core
High cache bandwidth
& more

23

# Knights Landing (KNL)

- KNL Knights Landing has 72 cores;

- Each one has an L1 cache;

- Pairs of cores are organized into tiles with a slice of the L2 cache symmetrically shared between the two cores;

- All caches are kept coherent;

- Two VPUs (Vector Processing Units) per core;

# Cluster modes

## One single space address

**Hemisphere:**
the tiles are divided into two parts called hemisphere

**Quadrant:**
tiles are divided into two parts called hemisphere or into four parts called qudrants

Node 0

| Tile | Tile | Tile | Tile | Tile | Tile |
|------|------|------|------|------|------|
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |

Node 0

| Tile | Tile | Tile | Tile | Tile | Tile |
|------|------|------|------|------|------|
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |

# Cluster modes

## Cache data are isolated in each sub numa domain

**SNC-2**:
the tiles are
divided into two
Numa Nodes

Node 1

Node 0

| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |

**SNC-4**:
the tiles are
divided into two
Numa Nodes

Node 1

Node 0

Node 3

Node 2

| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |
| Tile | Tile | Tile | Tile | Tile | Tile |

# Knights Landing Integrated On-Package Memory

## Multi-Channel DRAM (High-bandwidth memory)

# Integrated On-Package Memory Usage Models

Model configurable at boot tiem and software exposed through NUMA



Split Options:
25/75% or 50/50%

| Cache Model | Flat Model | Hybrid Model |
| --- | --- | --- |
| Hardware automatically manages the MCDRAM as a "L3 cache" between CPU and ext DDR memory | Manually manage how the app uses the integrated on-package memory and external DDR for peak perf | Harness the benefits of both Cache and Flat models by segmenting the integrated on-package memory |
| ▪ App and/or data set is very large and will not fit into MCDRAM<br>▪ Unknown or unstructured memory access behavior | ▪ App or portion of an app or data set that can be, or is needed to be "locked" into MCDRAM so it doesn't get flushed out | ▪ Need to "lock" in a relatively small portion of an app or data set via the Flat model<br>▪ Remaining MCDRAM can then be configured as Cache |

# Agenda

- Hybrid Parallel Architectures;

- Memory System and Vector Processing Units;

- Intel Architectures;

- **Profiling;**

- Optimizing Memory Access;

- Auto Vectorization;

- Guided Vectorization;

- Examples.

# Intel Advisor

- Evaluate multi-threading parallelization

- Intel® Advisor XE
  - ❑ Performance modeling using several frameworks for multi-threading in processors and co-processors:
    - o OpenMP, Intel® Cilk ™ Plus, Intel® Threading Bulding Blocks
    - o C, C++, Fortran (OpenMP only) and C# (Microsoft TPL)

  - ❑ Identify parallel opportunities
    - o Detailed information about vectorization;
    - o Check loop dependencies;

  - ❑ Scalability prediction: amount of threads/performance gains
  - ❑ Correctness (deadlocks, race conditions)

# Intel Advisor

- Survey Target;
  - Vectorization of loops: detailed information about vectorization;
  - Total Time: elapsed time on each loop considering the time involved in internal loops;
  - Self Time: elapsed time on each loop not considering the time involved in internal loops;
- Find Trip Counts;
  - Analysis to identify how many time particular loops run;
- Check Dependencies;
  - Analysis it there are many loop-carried dependencies;
- Check Memory Access Patterns.
  - Analysis to identify how your code is iterating with memory.

# Advisor – New Project

# Advisor - Analysis



File   View   Help

e000

**Vectorization Workflow**          **Threading Workflow**

OFF  Batch mode

**1. Survey Target**
▶ Collect                  → Collect Profiling Information

**1.1 Find Trip Counts**
▶ Collect                  → Obtain the amount of times a marked loop is executed

**Mark Loops for Deeper Analysis**
Select loops in the Survey Report for Dependencies
and/or Memory Access Patterns analysis.
-- There are no marked loops --

**2.1 Check Dependencies**
▶ Collect                  → Verify if a marked loop has dependencies between iterations
-- Nothing to analyze --

**2.2 Check Memory Access Patterns**
▶ Collect                  → Obtain the stride distribution of marked loops
-- Nothing to analyze --

**Summary of predicted parallel behavior**
Summary   Survey Report   Refinement Reports   Annotation Report

# Advisor – Survey Target

# Advisor – Survey Target

## Instruction Set Analysis

| Traits | Data Types | Number of Vector Registers | Vector Widths | Instruction Sets |
|---|---|---|---|---|
| Extracts; Type Conversions | Float32; Float64; Int32; UIn... | 15 | 128/256 | AVX; AVX2 |
| Extracts; Inserts; Type Conversions | Float32; Float64 | 14; 15 | 128; 256 | AVX |
| Extracts; Inserts; Type Conversions | Float32; Float64; Int32; UIn... | 16 | 256 | AVX; AVX2 |

## Advanced

| Transformations | Unroll Factor | Vectorization Details | Optimization Details | Location |
|---|---|---|---|---|
| | | | | vec.c:50 |
| Fused; Unrolled | 4 | | LOOP WAS DISTRIBUTED, CHUNK 1; LOOP WAS DISTRIBUTED, C... | vec.c:63 |
| | | | | vec.c:55 |

# Advisor – Memory Access Patterns

# Agenda

- Hybrid Parallel Architectures;
- Memory System and Vector Processing Units;
- Intel Architectures;
- Profiling;
- **Optimizing Memory Access;**
- Auto Vectorization;
- Guided Vectorization;
- Examples.

# Stride (array elements)

- Stride:
  - Step size between consecutive access of array elements;

- Strided access with stride k means touching every kth memory element
  - Unit Stride :
    - ❑ Sequential access (0, 1, 2, 3, 4, 5, 6, ...)
  - Non-unit stride
    - ❑ Constant Stride =
      - ○ 2 is (0, 2, 4, 6, 8, ...)
    - ❑ k is (0, k, 2k, 3k, 4k, ...)
    - ❑ Random Access;

- Strides > 1 commonly found in multidimensional data
  - Row accesses (stride=N) & diagonal accesses (stride=N+1)
  - Scientific computing (e.g., matrix multiplication)

# Padding

- Data structures may have members with different sizes.
- To maintain proper alignment the translator normally inserts additional unnamed data members so that each member is properly aligned.
- Example:

```
struct stu_a {
    int i;
    char c;
};
```

- Actual size 4+1 (5)

| i | i | i | i | c | i | i | i | c | i | i | i | c |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| 8 bytes | 8 bytes |

- After Padding size 4+4 (5)

| i | i | i | i | c | c | c | c | i | i | i | i | c |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| 8 bytes | 8 bytes |

- ...

# Padding

- Vectorization more efficient with unit strides
  - Non-unit strides will generate gather/scatter
  - Unit strides also better for data locality

Demo: padd.c

Icc padd.c –o padd

./padd

# Data layout

- AoS vs SoA (Array of Structures vs Structure of Arrays)
  - Layout your data as Structure of Arrays (SoA)

```
// Array of Structures (AoS)
struct coordinate {
    float x, y, z;
} crd[N];
    …
for (int i = 0; i < N; i++)
    … = … f(crd[i].x, crd[i],y,
crd[i].z);
```

Consecutive elements in memory

| x0 | y0 | z0 | x1 | y1 | z1 | … | x(n-1) | y(n-1) | z(n-1) |

```
// Structure of Arrays (SoA)
struct coordinate {
    float x[N], y[N], z[N];
} crd;
    …
for (int i = 0; i < N; i++)
    … = … f(crd.x[i], crd.y[i],
crd.z[i]);
```

Consecutive elements in memory

| x0 x1 … x(n-1) | y0 y1 … y(n-1) | z0 z1 … z(n-1) |

# Data Alignment

| How to… | Syntax | Semantics |
|---|---|---|
| …align data | `void* _mm_malloc(int size, int n)`<br>`void* _mm_free(int size)` | Allocate memory on heap aligned to *n* byte boundary. |
| | `int posix_memalign`<br>`    (void **p, size_t n, size_t size)` | |
| | `__declspec(align(n)) array` | Alignment for variable declarations. |
| …tell the compiler about it | `#pragma vector aligned` | Vectorize assuming all array data accessed are aligned (may cause fault otherwise). |
| | `__assume_aligned(array, n)` | Compiler may assume array is aligned to *n* byte boundary. |

# Loop Splitting

- Loop Splitting
  - Set of techniques to breaking the loop into multiple loops which have the same body, but iterate over different contiguous portions of the index range.
    - ❑ Body
    - ❑ Peel Loop: beginning of loop
    - ❑ Remainder Loop:  end of loop

- Loop Unrolling
  - Execute a set of iterations as a single iteration;

# Vectorization with multi-version loops

**Peel loop**
Alignment purposes
Might be vectorized

**Main loop**
Vectorized
Unrolled by x2 or x4

**Remainder loop**
Remainder iterations
Might be vectorized

```
LOOP BEGIN at gas_dyn2.f90(2330,26)
<Peeled>
   remark #15389: vectorization support: reference AMAC1U has unaligned
access
   remark #15381: vectorization support: unaligned access used inside loop
body
   remark #15301: PEEL LOOP WAS VECTORIZED
LOOP END

LOOP BEGIN at gas_dyn2.f90(2330,26)
   remark #25084: Preprocess Loopnests: Moving Out Store
   remark #15388: vectorization support: reference AMAC1U has aligned access
   remark #15399: vectorization support: unroll factor set to 2
   remark #15300: LOOP WAS VECTORIZED
   remark #15475: --- begin vector loop cost summary ---
   remark #15476: scalar loop cost: 8
   remark #15477: vector loop cost: 0.620
   remark #15478: estimated potential speedup: 15.890
   remark #15479: lightweight vector operations: 5
   remark #15488: --- end vector loop cost summary ---
   remark #25018: Total number of lines prefetched=4
   remark #25019: Number of spatial prefetches=4, dist=8
   remark #25021: Number of initial-value prefetches=6
LOOP END

LOOP BEGIN at gas_dyn2.f90(2330,26)
<Remainder>
   remark #15388: vectorization support: reference AMAC1U has aligned access
   remark #15388: vectorization support: reference AMAC1U has aligned access
   remark #15301: REMAINDER LOOP WAS VECTORIZED
LOOP END
```

# Agenda

- Hybrid Parallel Architectures;
- Memory System and Vector Processing Units;
- Intel Architectures;
- Profiling;
- Optimizing Memory Access;
- **Auto Vectorization;**
- Guided Vectorization;
- Examples.

# Vectorization on Intel® compilers

**Auto Vectorization**

- Compiler knobs

**Guided Vectorization**

- Compiler hints/pragmas
- Array notation
- Elemental Functions

**Low level Vectorization**

- C/C++ vector classes
- Intrinsics/Assembly

Easy of use

Fine control

# Auto vectorization

- Relies on the compiler for vectorization
  - No source code changes
  - Enabled with `-vec` compiler knob (default in `-O2` and `-O3` modes)
- Compiler smart enough to apply loop transformations
  - It will allow to vectorize more loops

| Option | Description |
|---|---|
| `-O0` | Disables all optimizations. |
| `-O1` | Enables optimizations for speed which are know to not cause code size increase. |
| `-O2/-O` (default) | Enables intra-file interprocedural optimizations for speed, including:<br>• **Vectorization**<br>• **Loop unrolling** |
| `-O3` | **Performs O2 optimizations and enables more aggressive loop transformations such as:**<br>• **Loop fusion**<br>• **Block unroll-and-jam**<br>• **Collapsing IF statements**<br>This option is recommended for applications that have loops that heavily use floating-point calculations and process large data sets. However, it might incur in slower code, numerical stability issues, and compilation time increase. |

# Vectorization: target architecture options

| Option | Description |
|---|---|
| `-mmic` | Builds an application that runs natively on Intel® MIC Architecture. |
| `-xfeature`<br>`-xHost` | Tells the compiler which processor features it may target, referring to which instruction sets and optimizations it may generate (not available for Intel® Xeon Phi™ architecture). Values for *feature* are:<br>• `COMMON-AVX512` (includes AVX512 FI and CDI instructions)<br>• **`MIC-AVX512` (includes AVX512 FI, CDI, PFI, and ERI instructions)**<br>• **`CORE-AVX512` (includes AVX512 FI, CDI, BWI, DQI, and VLE instructions)**<br>• **`CORE-AVX2`**<br>• `CORE-AVX-I` (including RDRND instruction)<br>• `AVX`<br>• `SSE4.2`, `SSE4.1`<br>• `ATOM_SSE4.2`, `ATOM_SSSE3` (including MOVBE instruction)<br>• `SSSE3`, `SSE3`, `SSE2`<br>**When using `-xHost`, the compiler will generate instructions for the highest instruction set available on the compilation host processor.** |
| `-axfeature` | Tells the compiler to generate multiple, feature-specific auto-dispatch code paths for Intel® processors if there is a performance benefit. Values for *feature* are the same described for `-xfeature` option. Multiple features/paths possible, e.g.: `-axSSE2,AVX`. It also generates a baseline code path for the default case. |

# Auto vectorization: not all loops will vectorize

- Data dependencies between iterations
  - Proven Read-after-Write data (i.e., loop carried) dependencies
  - Assumed data dependencies
    - ❏ Aggressive optimizations

RaW dependency
```
for (int i = 0; i < N; i++)
    a[i] = a[i-1] + b[i];
```

- Vectorization won't be efficient
  - Compiler estimates how better the vectorized version will be
  - Affected by data alignment, data layout, etc.

Inefficient vectorization
```
for (int i = 0; i < N; i++)
    a[c[i]] = b[d[i]];
```

- Unsupported loop structure
  - While-loop, for-loop with unknown number of iterations
  - Complex loops, unsupported data types, etc.
  - (Some) function calls within loop bodies

Function call within loop body
```
for (int i = 0; i < N; i++)
    a[i] = foo(b[i]);
```

# Validating vectorization

- ## Generate compiler report about optimizations

  `-qopt-report[=n]`   Generate report (level [1..6], default 2)

```
LOOP BEGIN at gas_dyn2.f90(193,11) inlined into gas_dyn2.f90(4326,31)
    remark #15300: LOOP WAS VECTORIZED
    remark #15448: unmasked aligned unit stride loads: 1
    remark #15450: unmasked unaligned unit stride loads: 1
    remark #15475: --- begin vector loop cost summary ---
    remark #15476: scalar loop cost: 53
    remark #15477: vector loop cost: 14.870
    remark #15478: estimated potential speedup: 2.520
    remark #15479: lightweight vector operations: 19
    remark #15481: heavy-overhead vector operations: 1
    remark #15488: --- end vector loop cost summary ---
    remark #25456: Number of Array Refs Scalar Replaced In Loop: 1
    remark #25015: Estimate of max trip count of loop=4
LOOP END
```
<div align="right">Vectorized loop</div>

```
LOOP BEGIN at gas_dyn2.f90(2346,15)
    remark #15344: loop was not vectorized: vector dependence prevents vectorization
    remark #15346: vector dependence: assumed OUTPUT dependence between IOLD line 376 and IOLD line 354
    remark #25015: Estimate of max trip count of loop=3000001
LOOP END
```
<div align="right">Non-vectorized loop</div>

# Agenda

- Hybrid Parallel Architectures;
- Memory System and Vector Processing Units;
- Intel Architectures;
- Profiling;
- Optimizing Memory Access;
- Auto Vectorization;
- **Guided Vectorization;**
- Examples.

# Intel® compiler directives for vectorization

| Directive | Clause | Description |
| --- | --- | --- |
| `ivdep` | | Instructs the compiler to ignore assumed vector dependencies. |
| `vector` | `always` | Force vectorization even when it might be not efficient. |
| | `[un]aligned` | Use [un]aligned data movement instructions for all array vector references. |
| | `[non]temporal(`*`var1`*`[,…])` | Do or do not generate non-temporal (streaming) stores for the given array variables. On Intel® MIC architecture, generates a cache-line-evict instruction when the store is known to be aligned. |
| | `[no]vecreminder` | Do (not) vectorize the remainder loop when the main loop is vectorized. |
| | `[no]mask_readwrite` | Enables/disables memory speculation causing the generation of [non-]masked loads and stores within conditions. |

# Intel® compiler directives for vectorization

| Directive | Clause | Description |
|---|---|---|
| `ivdep` | | Instructs the compiler to ignore assumed vector dependencies. |
| `simd` | `vectorlength(n1[,…])` `vectorlengthfor(dtype)` | Assume safe vectorization for the given vector length values or data type. |
| | `private(var1[,…])` `firstprivate(var1[,…])` `lastprivate(var1[,…])` | Which variables are private to each iteration; *firstprivate*, initial value is broadcasted to all private instances; *lastprivate*, last value is copied out from the last instance. |
| | `linear(var1:step1[,…])` | Letting know the compiler that *var1* is incremented by *step1* on every iteration of the original loop. |
| | `reduction(oper:var1[,…])` | Which variables are reduction variables with a given operator. |
| | `[no]assert` | Warning or error when vectorization fails. |
| | `[no]vecremainder` | Do (not) vectorize the remainder loop when the mail loop is vectorized. |

# Guided vectorization: disambiguation hints

- ## Assume function arguments won't be aliased
  - C/C++: Compile with `-fargument-noalias`

- ## C99 "restrict" keyword for pointers
  - Compile with `-restrict` otherwise

```
void v_add(float *restrict c,
           float *restrict a,
           float *restrict b)
{
    for (int i = 0; i < N; i++)
        c[i] = a[i] + b[i];
}
```

```
void v_add(float *c, float *a, float *b)
{
for (int i = 0; i < N; i++)
        c[i] = a[i] + b[i];
}
```

# Guided vectorization:

- `#pragma simd or #pragma ivdep`
  - Force loop vectorization ignoring **all** dependencies
    - ❑ Additional clauses for specify reductions, etc.

```
void v_add(float *c, float *a, float *b)
{
#pragma simd
    for (int i = 0; i < N; i++)
        c[i] = a[i] + b[i];
}
```
SIMD loop

```
__declspec(vector)
void v_add(float c, float a, float b)
{
    c = a + b;
}

    …
for (int i = 0; i < N; i++)
    v_add(C[i], A[i], B[i]);
```
SIMD function

# Agenda

- Hybrid Parallel Architectures;
- Memory System and Vector Processing Units;
- Intel Architectures;
- Profiling;
- Optimizing Memory Access;
- Auto Vectorization;
- Guided Vectorization;
- **Examples.**

# Matrix Multiplication - Serial

```
void multiply(int msize, int tidx, int numt, TYPE a[][NUM], TYPE
b[][NUM], TYPE c[][NUM], TYPE t[][NUM])
{

int i,j,k;
    for(i=0; i<msize; i++) {
        for(k=0; k<msize; k++) {
            for(j=0; j<msize; j++) {
                c[i][j] = c[i][j] + a[i][k] * b[k][j];
            }
        }
    }
}
```

# Matrix Multiplication



| Function Call Sites and Loops | ♦ | Vector Issues | Self Time▼ | Total Time | Type | Why No Vectorization? |
|---|---|---|---|---|---|---|
| ⌄ ⏱ [loop in multiply3 at multiply.c:228] | | 💡 2 Assume … | 0.170s ■ | 0.170s ■ | Scalar | ▣ vector dependence prevents vectorization |
| ⌄ ⏱ [loop in __libc_csu_init] | | | 0.000s [ | 0.000s [ | Scalar | |
| ⌄ ⏱ [loop in _INTERNAL_16_offload_host_cpp_ad92… | | | 0.000s [ | 0.000s [ | Scalar | |
| ⌄ ⏱ [loop in func@0x5b810] | | 💡 2 Data typ.. | 0.000s [ | 0.000s [ | Scalar | |
| ⌄ ⏱ [loop in main at matrix.c:144] | | 💡 2 Data typ.. | 0.000s [ | 0.000s [ | Scalar | ▣ inner loop was already vectorized |
| ⌄ ⏱ [loop in multiply3 at multiply.c:227] | | 💡 2 Assume … | 0.000s [ | 0.170s ■ | Scalar | ▣ vector dependence prevents vectorization |
| ⌄ ⏱ [loop in multiply3 at multiply.c:226] | | 💡 2 Assume … | 0.000s [ | 0.170s ■ | Scalar | ▣ vector dependence prevents vectorization |
| ⊞ ⏱ [loop in main at matrix.c:144] | | 💡 1 Data typ.. | 0.000s [ | 0.000s [ | Vectorized (B… | |

| Source | Top Down | Loop Analytics | Loop Assembly | 💡 Recommendations | ▣ Compiler Diagnostic Details |
|---|---|---|---|---|---|

**File: multiply.c:228 multiply3**

| Line | Source |
|---|---|
| 218 | `void multiply3(int msize, int tidx, int numt, TYPE a[][NUM], TYPE b[][NUM], TYPE c[][NUM], TYPE t[][NUM])` |
| 219 | `{` |
| 220 | |
| 221 | `//#pragma omp target device(0) map(a[0:NUM][0:NUM]) \` |
| 222 | `//map(b[0:NUM][0:NUM]) map(c[0:NUM][0:NUM])` |
| 223 | `//{` |
| 224 | `  int i,j,k;` |
| 225 | `//   #pragma omp parallel for collapse (2) //num threads(60)` |
| 226 | `    for(i=0; i<msize; i++) {` |
| | ⏱ [loop in multiply3 at multiply.c:226]<br>  Scalar loop. Not vectorized: vector dependence prevents vectorization<br>  No loop transformations applied |
| 227 | `      for(k=0; k<msize; k++) {` |
| | ⏱ [loop in multiply3 at multiply.c:227]<br>  Scalar loop. Not vectorized: vector dependence prevents vectorization<br>  Remainder loop |
| 228 | `        for(j=0; j<msize; j++) {` |
| | ⏱ [loop in multiply3 at multiply.c:228]<br>  Scalar loop. Not vectorized: vector dependence prevents vectorization<br>  Loop was unrolled by 2 |
| 229 | `          c[i][j] = c[i][j] + a[i][k] * b[k][j];` |
| 230 | `        }` |
| 231 | `      }` |
| 232 | `    }` |
| 233 | `//}` |
| 234 | `}` |
| 235 | |

# Matrix Multiplication

- Check dependency analysis shows that it is safe to enforce the vectorization of this loop

# Matrix Multiplication - vectorized

```
void multiply(int msize, int tidx, int numt, TYPE a[][NUM], TYPE
b[][NUM], TYPE c[][NUM], TYPE t[][NUM])
{

int i,j,k;
    for(i=0; i<msize; i++) {
        for(k=0; k<msize; k++) {
            #pragma simd
            for(j=0; j<msize; j++) {
                c[i][j] = c[i][j] + a[i][k] * b[k][j];
            }
        }
    }
}
```

# Matrix Multiplication

# Example

- Particle Binning Problem[1]

- Optimizations:
  - Automatic Vectorization
  - Data Alignment

[1] http://colfaxresearch.com/optimization-techniques-for-the-intel-mic-architecture-part-2-of-3-strip-mining-for-vectorization/

# Particle Binning - Serial

```
for (int i = 0; i < inputData.numDataPoints; i++) {

    // Transforming from cylindrical to Cartesian coordinates:
    const FTYPE x = inputData.r[i]*COS(inputData.phi[i]);
    const FTYPE y = inputData.r[i]*SIN(inputData.phi[i]);

    // Calculating the bin numbers for these coordinates:
    const int iX = int((x - xMin)*binsPerUnitX);
    const int iY = int((y - yMin)*binsPerUnitY);

}
```

# Particle Binning - Vectorized

```
for (int ii = 0; ii < inputData.numDataPoints; ii += STRIP_WIDTH) {

    int iX[STRIP_WIDTH];
    int iY[STRIP_WIDTH];

    const FTYPE* r   = &(inputData.r[ii]);
    const FTYPE* phi = &(inputData.phi[ii]);

    // Vector loop

    for (int c = 0; c < STRIP_WIDTH; c++) {
        // Transforming from cylindrical to Cartesian coordinates:
        const FTYPE x = r[c]*COS(phi[c]);
        const FTYPE y = r[c]*SIN(phi[c]);

        // Calculating the bin numbers for these coordinates:
        iX[c] = int((x - xMin)*binsPerUnitX);
        iY[c] = int((y - yMin)*binsPerUnitY);
    }

}
```

# Particle Binning - Vectorized

# Particle Binning - Data Alignment

```
for (int ii = 0; ii < inputData.numDataPoints; ii += STRIP_WIDTH) {

    int iX[STRIP_WIDTH] __attribute__((aligned(64)));
    int iY[STRIP_WIDTH] __attribute__((aligned(64)));

    const FTYPE* r   = &(inputData.r[ii]);
    const FTYPE* phi = &(inputData.phi[ii]);

    // Vector loop
#pragma vector aligned
    for (int c = 0; c < STRIP_WIDTH; c++) {
        // Transforming from cylindrical to Cartesian coordinates:
        const FTYPE x = r[c]*COS(phi[c]);
        const FTYPE y = r[c]*SIN(phi[c]);

        // Calculating the bin numbers for these coordinates:
        iX[c] = int((x - xMin)*binsPerUnitX);
        iY[c] = int((y - yMin)*binsPerUnitY);
    }

}
```

# Diffusion - Serial

# Diffusion - Vectorized

Potential inefficient memory access;

| Function Call Sites and Loops▲ | 🌡 | Vector Issues | Self Time | Total Time | Type | Why No Vectorization? | Vectorized Loops | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Vec... | Efficiency | Gai... |
| ⊻⏱ [loop in diffusion_openmpv at diffusion_ompvect.c:106] | | 💡 1 Potential.. | 0.640s▎ | 5.050s■ | Scalar | 🔒 inner loop was already vectorized | | | |
| ⊞⏱ [loop in diffusion_openmpv at diffusion_ompvect.c:106] | | 💡 1 Potential.. | 0.000s[ | 0.060s[ | Scalar Versions | 🔒 1 inner loop was already vectorized | | | |
| ⊻⏱ [loop in diffusion_openmpv at diffusion_ompvect.c:108] | | | 4.410s■ | 4.410s■ | Vectorized (B... | | AVX2 | ~62% | 4.93x |
| ⊞⏱ [loop in diffusion_openmpv at diffusion_ompvect.c:108] | | 💡 1 Possible i.. | 0.060s[ | 0.060s[ | Vectorized Ve... | | AVX2 | ~50% | 4.00x |
| ⊻⏱ [loop in diffusion_openmpv at diffusion_ompvect.c:108] | | 💡 1 Possible i.. | 0.030s[ | 0.030s[ | Vectorized (B... | | AVX2 | ~50% | 4.00x |
| ⊻⏱ [loop in diffusion_openmpv at diffusion_ompvect.c:108] | | 💡 1 Possible i.. | 0.030s[ | 0.030s[ | Vectorized (B... | | AVX2 | ~50% | 4.00x |
| ⊻⏱ [loop in init at diffusion_ompvect.c:71] | | 💡 1 Data typ.. | 0.000s[ | 0.060s[ | Scalar | 🔒 inner loop was already vectorized | | | |

| Source | Top Down | Loop Analytics | Loop Assembly | 💡 Recommendations | 🔒 Compiler Diagnostic Details |
|---|---|---|---|---|---|

**File: diffusion_ompvect.c:108 diffusion_openmpv**

| Line | Source | Tota |
|---|---|---|
| 97 | `void` | |
| 98 | `diffusion openmpv(REAL *restrict f1, REAL *restrict f2, int nx, int ny, int nz,` | |
| 99 | `REAL ce, REAL cw, REAL cn, REAL cs, REAL ct,` | |
| 100 | `REAL cb, REAL cc, REAL dt, int count) {` | |
| 101 | `REAL *f1 t = f1;` | |
| 102 | `REAL *f2 t = f2;` | |
| 103 | | |
| 104 | ⊞ `for (int i = 0; i < count; ++i) {` | |
| 105 | ⊞ `for (int z = 0; z < nz; z++) {` | |
| 106 | ⊞ `for (int y = 0; y < ny; y++) {` | |
| 107 | `#pragma simd` | |
| 108 | ⊟ `for (int x = 0; x < nx; x++) {` | |

🔴 [loop in diffusion_openmpv at diffusion_ompvect.c:108]
Vectorized AVX; FMA loop processes Float32 data type(s) and includes FMA
No loop transformations applied
🔴 [loop in diffusion_openmpv at diffusion_ompvect.c:108]
Vectorized AVX; AVX2; AVX2GATHER; FMA loop processes Float32; Int32 data type(s) and includes FMA; Gathers
No loop transformations applied
🔴 [loop in diffusion_openmpv at diffusion_ompvect.c:108]
Vectorized AVX; AVX2; AVX2GATHER; FMA loop processes Float32; Int32 data type(s) and includes FMA; Gathers
No loop transformations applied
🔴 [loop in diffusion_openmpv at diffusion_ompvect.c:108]
Vectorized AVX; AVX2; AVX2GATHER; FMA loop processes Float32; Int32 data type(s) and includes FMA; Gathers
No loop transformations applied
🔴 [loop in diffusion_openmpv at diffusion_ompvect.c:108]
Vectorized AVX; AVX2; AVX2GATHER; FMA loop processes Float32; Int32 data type(s) and includes FMA; Gathers
No loop transformations applied

# Diffusion - alignment

# Interpolation

```
__declspec(vector)
int FindPosition(double x) {
    return (int)(log(exp(x*steps)));
}


__declspec(vector)
double Interpolate(double x, const point*
vals)
{
    int ind = FindPosition(x);

    ...


    return res;
}
```

```
int main ( int argc , char argv [] )
{
  . . .
  for ( i=0; i <ARRAY_SIZE;++ i ) {
      dst[i] = Interpolate( src[i], vals ) ;
  }
  . . .
}
```

# Vectorization report - Interpolate

Begin optimization report for: Interpolate.._simdsimd3__H2n_v1_s1.P(double, const point *)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(74,48) ]
========================================================================

Begin optimization report for: Interpolate.._simdsimd3__H2m_v1_s1.P(double, const point *)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(74,48) ]
========================================================================

Begin optimization report for: Interpolate.._simdsimd3__L4n_v1_s1.V(double, const point *)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(74,48) ]
remark #15415: vectorization support: gather was generated for the variable pnt:  indirect access, 64bit indexed   [ main.c(78,26) ]
remark #15415: vectorization support: gather was generated for the variable pnt:  indirect access, 64bit indexed   [ main.c(78,36) ]
========================================================================

Begin optimization report for: Interpolate.._simdsimd3__L4m_v1_s1.V(double, const point *)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(74,48) ]
remark #15415: vectorization support: gather was generated for the variable pnt: masked, indirect access, 64bit indexed   [ main.c(78,26) ]
remark #15415: vectorization support: gather was generated for the variable pnt: masked, indirect access, 64bit indexed   [ main.c(78,36) ]

egin optimization report for: FindPosition.._simdsimd3__H2n_v1.P(double)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(70,28) ]
========================================================================

Begin optimization report for: FindPosition.._simdsimd3__H2m_v1.P(double)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(70,28) ]
========================================================================

Begin optimization report for: FindPosition.._simdsimd3__L4n_v1.V(double)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(70,28) ]
========================================================================

Begin optimization report for: FindPosition.._simdsimd3__L4m_v1.V(double)

   Report from: Vector optimizations [vec]

remark #15301: FUNCTION WAS VECTORIZED   [ main.c(70,28) ]
========================================================================

# Lattice Boltzmann

# Questions?

Silvio Stanzani , Raphael Cóbe , Rogério Iope
Núcleo de Computação Científica – UNESP

silvio@ncc.unesp.br
rmcobe@ncc.unesp.br
rogerio@ncc.unesp.br