

使用 Stencil 评估 Intel AVX2 Vgather 指令

林新华^{1,2} 秦 强¹ 李 硕³ 文敏华¹ 松岗聪²

(上海交通大学高性能计算中心 上海 200240)¹ (东京工业大学学术国际情报中心 东京 152-8550)²
(Intel 公司软件与服务部门 波特兰 999039)³

摘 要 为了更好地在向量化时读取离散的数据, Intel 在 Haswell CPU 提供了 AVX2vgather 指令。由于 Stencil 在设置边界条件时使用了条件判断, 因此编译器生成了 vgather 指令, 并降低了 Stencil 在 Haswell 上的性能。提出使用 peel 优化或 intrinsic load 的方法来避免 vgather 指令的生成, 并把该方法应用到 3 个 Stencil 基准算例、长程 Stencil 程序 3DFD 以及混合 Stencil 应用 3DEW 上。这些 Stencil 在 Haswell 上的性能都获得了 1.22X 至 3.88X 不等的提升。通过研究指令的实现, 发现 vgather 指令会被解码成多个微操作(μ ops), 并为每个要读入的元素生成一个 μ ops。由于 vgather 指令解码时会产生较高的开销, 导致 vgather 指令成为 Stencil 在 Haswell 上的性能瓶颈。了解 AVX2 vgather 指令的实现以及掌握避免生成 vgather 指令的优化方法, 对在 Haswell 上调优具有良好空间局部性应用的性能有一定的参考价值。

关键词 AVX2 vgather 指令, Stencil, 性能评估
中图法分类号 TP391 **文献标识码** A **DOI** 10.11896/j.issn.1002-137X.2017.01.004

Evaluating Intel AVX2 Vgather Instructions with Stencils

LIN Xin-hua^{1,2} QIN Qiang¹ LI Shuo³ WEN Min-hua¹ MATSUOKA Satoshi²

(Center for High Performance Computing, Shanghai Jiao Tong University, Shanghai 200240, China)¹
(Global Scientific Information and Computing Center, Tokyo Institute of Technology, Tokyo 152-8550, Japan)²
(Software and Services Group, Intel Corporation, Portland 999039, USA)³

Abstract Intel provided AVX2 vgather instruction on Haswell CPU to better support reading discontinued data in vectorization. We found the compiler generates vgather instructions, which slow down the performance of Stencil on Haswell, because the branches exist in defining boundary condition of Stencils. We proposed to utilize peel optimization or intrinsic load to avoid these vgather instructions. We applied these optimizations to three Stencil benchmarks, a long-range Stencil 3DFD, and a hybrid Stencil application, and archived the speedup from 1.22X to 3.88X on Haswell. By analyzing the implementation of the instruction, we found the vgather instructions are decoded into multiple micro-operations (μ ops), and the instructions generate one μ ops for each element to be gathered. Due to the high overhead of decoder, the vgather instructions become the performance bottleneck of Stencils on Haswell. It is believed that the understanding of the implementation of AVX2 vgather instructions and adopting the optimizations to avoid the vgather instructions are quite helpful for performance tuning the applications with good spatial locality on Haswell.

Keywords AVX2 vgather, Stencil, Performance evaluation

1 简介

为了更好地在向量化时读取离散的数据, Intel 陆续在不同平台上提供了硬件支持的 vgather 指令: 2013 年上半年发布的 Knight Corner(缩写为 KNC)上的 IMCI(Initial Many Core Instructions)vgather 指令; 2013 年 6 月发布的 Haswell(缩写为 HSW)CPU 上的 AVX(Advanced Vector Extension)2vgather 指令; 且 2014 年发布的 Broadwell CPU 改进了该指令的实现; 2016 年发布的 Knight Landing(KNL)以及 SkylakeXeon CPU 都会支持 AVX-512vgather 指令。

有一些文献表明 KNC 上的 IMCI vgather 指令会成为数据离散的应用^[1]和 Stencil^[2]的性能瓶颈。因此本文拟回答以下 2 个问题:

- 1) AVX2vgather 指令是否也会降低 Stencil 在 HSW 上的性能? 编译器在何种情况下会生成 vgather 指令? 应该如何避免?
 - 2) 为什么 AVX2 vgather 指令会降低 Stencil 的性能?
- 为了回答清楚这 2 个问题, 首先分析了 Stencil 的源代码, 发现在 Stencil 中设置边界条件时需要元素地址进行条件判断。编译器因为在编译时无法确定这些元素的地址, 就

到稿日期: 2015-12-14 返修日期: 2016-02-01 本文受国家重点研发计划(2014AA01A302, 2016YFB0201800), 日本学术振兴会 RONPAKU Fellowship 资助。
林新华(1979—), 男, 讲师, 高级工程师, 主要研究方向为体系结构与代码优化; 秦 强(1992—), 男, 硕士生, 主要研究方向为体系结构; 李 硕(1969—), 男, 工程师, 主要研究方向为高性能计算在量化金融中的应用; 文敏华(1988—), 男, 工程师, E-mail: wenminhua@sjtu.edu.cn (通信作者); 松岗聪(1963—), 男, 教授, 主要研究方向为高性能计算。

生成了 vgather 指令。提出使用 peel 优化或 intrinsic load 的方法来避免生成 vgather 指令。对 1D3P, 2D5P, 3D7P 这 3 个 Stencil 基准算例、长程 Stencil 程序 3DFD 以及混合 Stencil 应用 3DEW 进行了性能验证,发现消除了 vgather 指令之后,所有 Stencil 都获得了 1.22X 至 3.88X 不等的性能提升。其次,使用 Intel Architecture Code Analysis(IACA)分析 vgather 指令的实现,发现 vgather 指令会为每个要读取的元素生成一个 μops 。据此初步分析了 vgather 指令对单核和多核性能的影响。

概括来说,本文有 3 个创新点:

1) 首先发现 AVX2 vgather 指令会降低 Stencil 在 HSW 上的性能。

2) 发现由于设定元素地址时进行了条件判断,编译器会生成 AVX2 vgather 指令,因此首先提出使用 peel 优化或 intrinsic load 的方法来避免生成 vgather 指令。

3) 发现 vgather 指令会被解码成多个 μops ,并为每个要读入的元素生成一个 μops ,并首先用 IACA 予以证明。

本文第 2 节介绍相关工作;第 3 节介绍 vgather 指令;第 4 节说明实验配置;第 5 节分析了编译器生成 vgather 指令的原因,提出了相应的优化方法,并评估了不同 Stencil 优化前后的性能;第 6 节深入分析了 AVX2 vgather 指令的实现,并评估了 vgather 指令对单核和多核性能的影响;最后总结全文并展望下一步工作。

2 相关工作

有一些文献讨论了 KNC 上 IMCI vgather 指令的性能瓶颈问题。George Hager 等人^[1]将不规则数据访问的医学图像算法 FDK 进行性能优化,并用性能建模的方法分析了 vgather 成为 FDK 在 KNC 上性能瓶颈的原因。本文作者^[2]使用 semi-empirical 的性能建模方法,分析 vgather 指令在 KNC 上对 Stencil 性能的影响。实验结果和模型推测一致,vgather 引起的流水线停滞会占到 Stencil 全部运行时间的 20% 左右。本文与这些文献不同,是针对 AVX2 vgather 指令进行的分析。

3 vgather 指令

Intel 从 AVX2 开始提供了专门的 vgather 指令。之前的 SSE 和 AVX 没有专门的 gather 指令,需要使用标量 load 逐个读入每个元素,然后用 shuffles (比如 pinsrd, extractps, vinsertf128 等指令)将这些元素插入到向量寄存器中。

• 指令格式与寻址方式

表 1 中 vgatherdpq 的 v 表示 AVX 操作, d 表示元素索引的大小,这里是 doubleword(4Byte)。p 是指 packed,表示这是量化的操作。p 出现在第 2 个位置时(即 vpgatherdp),表示操作数为整数;出现在倒数第 2 个位置时(即 vgatherdpq),表示操作数为浮点数。q 表示要元素的大小,这里是 quadword(8 Byte)。元素大小和索引大小有多种组合,对应的 vgather 指令如图 1 所示。

表 1 vgather 的指令格式

| vgatherdpq ymm0, [esi+ymm1*4], ymm1 | | | |
|-------------------------------------|--------|------------|------------|
| 元素大小 | | | |
| 索引大小 | 4 Byte | 8 Byte | |
| | 4 Byte | vgatherdps | vgatherdpd |
| 索引大小 | 8 Byte | vgatherqps | vgatherqpd |

图 1 vgather 指令的不同形式

表 1 中的指令使用了 Intel 汇编格式,因此 ymm0 是目的地寄存器, $[\text{esi} + \text{ymm1} * 4 + 8]$ 是源地址, ymm1 是 mask 寄存器。ymm 是 256bit 的向量寄存器, HSW 由于没有专门的 mask 寄存器,因此就用 ymm 寄存器来代替,但这可能会造成 mask 寄存器空间的浪费(通用向量寄存器有 256bit,而用作 mask 则最多只需要 8bit)以及通用向量寄存器的紧张。如表 1 中 vgatherdpd 读入 4 个双精度元素,mask 寄存器只需要 4bit 就可以了。ymm1 寄存器将 256bit 分成 4 个 lane,每个 lane 的长度为 64bit,正好对应 1 个双精度元素。若某 lane 中的数值为 0,则对应的元素不必读入或者已经读入;反之,则读入对应的元素,并将该 lane 中的数值置 0。

• 指令操作

以表 1 中的指令为例,vgather 读入元素分成 3 步。第一步是生成地址。针对每一个要读入的元素,将指令中的相对地址发送到 AGU 中,生成绝对地址。如图 2 所示,ymm1 中要读取的第 2 个元素的 index 是 3,即通用寄存器 esi 所指到的数组的第 3 个元素;第二步,检测 mask 寄存器 ymm1 中每个元素对应的数值,如果为 0,则不读入元素;反之,则进入第三步。比如,如果 ymm1 中第 3 个元素对应的 mask 数值为 0,就不需要读入元素,因此 vgather 只需读入 3 个元素。注意到地址生成是在检测 mask 寄存器之前,这就意味着即使 mask 寄存器对应的数值为 0,vgather 指令还是会为那个元素生成绝对地址。第三步是读入该元素到目的地寄存器的相应位置中。注意,这 3 个元素的读入是乱序的,不一定是从高位到低位或是从低位到高位;而且 Haswell 有 2 个读入端口(端口 2 和 3),因此能同时读入 2 个元素。

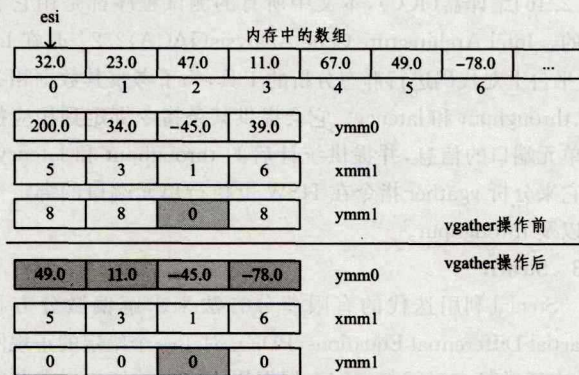


图 2 vgatherdpq ymm0, $[\text{esi} + \text{ymm1} * 4]$, ymm1 的操作示意图

• 指令解码

x86 指令集属于 CISC,其每条指令的长度都不同。为了简化执行单元(execution unit)的设计,解码器中将 x86 指令动态翻译成类似 RISC 的指令,称为微操作(micro-operations,即 μops),但 μops 并不等价于 RISC 指令,而是更接近于解码后的 RISC 指令。x86 的解码过程需要好几个时间周期(即 cycle),因此会成为处理器 Front End 性能瓶颈的来源之一。

如图 3 所示,解码可分为 2 个阶段:预解码阶段和解码阶段。预解码阶段中,指令长度解码器(Instruction length Decoder, ILD)将原始的二进制流划分成有效的 x86 指令流,然后发送给指令队列(Instructions Queue, IQ)。解码阶段则将从 IQ 传来的 x86 指令流动态翻译成对应的 μops 。与 Nehalem 和 Sandy Bridge 架构相同,Haswell 有 3 个简单解码器和 1 个复杂解码器。简单解码器负责解码那些可以翻译成 1 条 μops 的指令,而复杂解码器则负责动态解码那些最多可以

翻译为 4 个 μops 的指令。需要翻译成 4 个以上 μops 的指令,如 `vgather`,则首先会将其发送至复杂指令解码器,然后停止正常的解码流水线,再发送给 MSROM(micro-sequencer)单元。MSROM 单元有一个 sequencer 回路和一个 ROM 数组,它会输出一个预先翻译好的 μops 程序。

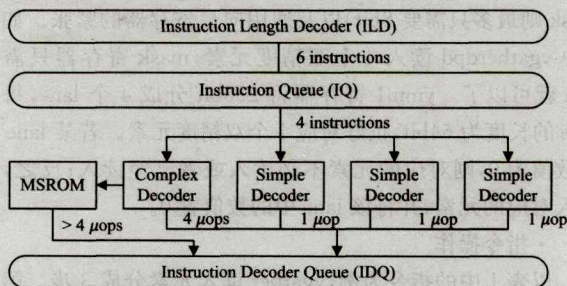


图3 Haswell Decoder

4 实验配置

4.1 硬件

硬件环境为 Intel Haswell E5-2693, 14 核, 主频为 2.3GHz。L1 和 L2 是每个核私有的, 而 L3 是所有核共享的。L1D 和 L1I 的大小均为 32kB, 8 路, latency 为 4 cycle。L2 cache 的大小为 256kB, 也是 8 路, latency 为 11cycle。L3 cache 的大小为 35MB。内存峰值带宽为 60GB/s, Stream 实测带宽为 53GB/s。

4.2 软件

主要使用了以下 2 个软件: Intel C/C++ Compiler 15.2.164 编译器 (ICC), 本文中所有的测试程序都是用它编译的。Intel Architecture Code Analysis (IACA) 2.2^[6] 是在 Intel 平台上对代码进行静态分析的工具, 为了考察其数据相关性、throughput 和 latency。它会提供某条指令绑定到相应执行单元端口的信息, 并提供统计后的 throughput 和 latency。用它来分析 `vgather` 指令在 HSW 上执行单元端口的绑定情况以及 throughput。

4.3 Stencil

Stencil 利用迭代的有限差分方法来求解偏微分方程 (Partial Differential Equations, PDE), 对于一个给定的正规网格, 它通常是在时间与空间上根据周边格点数值加权求和的过程。Stencil 需要读入离散的数据, 但又具有一定的空间局部性, 且其性能往往受限于内存带宽。Stencil 被广泛应用于计算流体力学、数值天气模拟、石油勘探等多个与国民经济息息相关的领域。

Stencil 有 2 种不同的元素更新方法: Jacobi 和 Gauss-Seidel。Jacobi 方法是读和写在不同的网格上, 而 Gauss-Seidel 方法是读和写在同一网格上。本文中 3 个 Stencil 基准算例、1 个长程 Stencil 程序和 1 个混合 Stencil 应用都采用了 Jacobi 方法。

• Stencil 基准算例

Stencil 代码按维度可以分成 1D3P, 2D5P 和 3D7P 这 3 种形式, 如图 4 所示。1D3P 更新中间的点, 需要自己以及左右 2 个点的数据; 类似地, 2D5P 需要自己以及上下左右 4 个点的数据; 3D7P 需要自己以及上下前后左右 6 个点的数据。从图 4 中还可以看出, 1D3P 在内存中具有良好的空间局部性; 2D3P 的 X 方向与 1D3P 相同, 而 Y 方向则有了 NX 的间距; 类似地, 3D7P 中 Z 方向则有了 NX * NY 的间距。

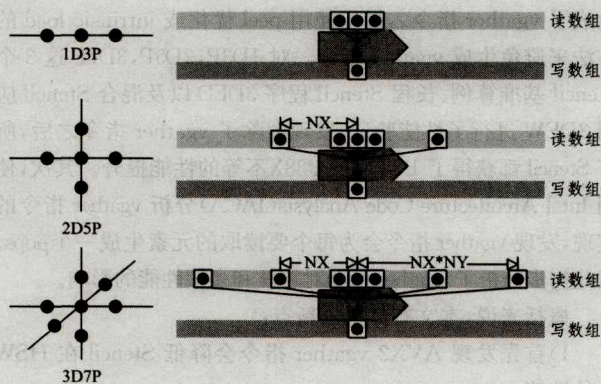


图4 Stencil 基准算例示意图

• 长程 Stencil 程序: 3DFD

3DFD^[7] 的全称为各向同性介质中三维有限差分程序 (3-Dimensional Finite Difference Code with an Isotropic)。该程序利用有限差分方法求解各向同性介质中噪音波动方程。更新 3DFD 的一个点, 需要自己以及上下前后左右 24 个点, 一共 25 个点, 即 3D25P。

• 混合 Stencil 应用: 3DEW

3DEW 的全称为三维纵横波分离的弹性波方程模拟 (3D pure P and S wave elastic wave equation modeling)。该应用由中国石油东方地球物理公司自主研发^[8], 采用波场延拓技术来模拟弹性波在各向同性的弹性介质中传播。3DEW 分别模拟纵波 (P 波) 和横波 (S 波), 以便更好地研究这两种波在弹性介质中的传递, 同时利用高阶有限差分方法分析弹性波的传递情况。3DEW 是一个混合 Stencil, 波被分解为 u, v, w 3 个方向, 空间被分解为 x, y, z 3 个方向。更新一个点, 首先需要计算 1 维上波的传播, 比如 u 在 x 方向上的传播, 这是一个 1D10P Stencil, 根据组合, 这样的 Stencil 一共有 9 个。其次需要计算 2 维平面上波的传播, 这是一个 2D100P Stencil, 这样的 Stencil 根据排列组合一共有 6 个。因此 3DEW 每更新一个点, 就需要进行 9 个 1D10P 加上 6 个 2D100P 的 Stencil 计算。

5 Stencil 代码在 Haswell 上的性能评估

5.1 Stencil 代码优化方法

对 3 个 Stencil 基准算例以及 3DFD 和 3DEW 分别进行了各种优化, 如表 2 所列。版本 A 是最基本的串行实现。版本 B 是在 Stencil 最外层循环前添加 `#pragma omp for` 进行多线程并行。版本 C 是在最内层循环前添加 `#pragma simd` 进行向量化。版本 D 和版本 E 都是对 cache 空间局部性进行的优化。Tiling 是将数据分块以增加 cache 中的空间局部性。Tiling 是比较常用的优化方式, 但其缺点是分块的大小需要根据应用以及所运行的硬件进行调整才能达到最佳性能。版本 F 和 G 将在 5.2 节中介绍。已实现了的用 \checkmark 表示, 有些因为技术原因没来得及实现的用 \times 表示。

表2 不同 Stencil 的各种优化版本

| 优化版本 | 1D3P | 2D5P | 3D7P | 3DFD | 3DEW |
|----------------------|--------------|--------------|--------------|--------------|--------------|
| A: baseline | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| B: A+OpenMP | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| C: B+Vector | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| D: C+Tiling | \times | \checkmark | \checkmark | \checkmark | \checkmark |
| E: C+Cache Oblivious | \checkmark | \checkmark | \checkmark | \times | \times |
| F: D+Peel | \checkmark | \checkmark | \checkmark | \times | \checkmark |
| G: D+Intrinsic | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |

5.2 编译器生成 vgather 指令的原因及对策

检查不同 Stencil 的版本 C, D, E 生成的汇编代码,发现编译器都生成了 vgather 指令。我们认为这是由于编译器在进行向量化时,遇到设置边界条件所使用的条件判断语句引起的。以 3D7P 为例,需要为 6 个方向都设置边界条件。比如 w 是当前元素西边的一个元素,通常情况是当前元素的索引减 1,即 $c-1$ 。但当前元素是最左边的那个元素时, w 就应该还是这个元素自己,即 c 。类似地,用这种条件判断语句来设定其他 5 个方向的边界条件。对于编译器而言, w 的取值到底是 c 还是 $c-1$ 在编译时并不确定,要等到运行时知道 x 的数值时才能确定,这种访问就属于间接读入,因此也就会导致 ICC 在 HSW 上生成 vgather 指令。

有 2 种方法可以避免编译器生成 vgather,其核心思想相同,即避免在定义数组元素时使用条件判断语句,分别是 Peel 优化和 Intrinsic load。限于篇幅,这 2 种优化的完整示例代码可以从 Github^[11]上下载。

• Peel 优化

比较大的 Stencil 遇到边界而需要更换读入元素以避免越界的情况是非常少的,因此可以把边界条件的判断从内层循环剥离出去,这是由于绝大多数的内层循环是不需要边界条件的。以 3D7P 为例, x 方向上只有第一个元素 0 和最后一个元素 $nx-1$ 需要进行边界条件的判断,因此将这 2 个元素作为头尾剥离出来单独处理,把不需要进行边界判断的元素 1 到 $nx-2$ 仍留在内层循环中。然后由于内层循环中的元素不会遇到边界,因此可以直接使用 $e=c-1$ 和 $w=c+1$ 消除对东边元素和西边元素的索引计算。这种方式既能避免生成 vgather 指令,又能提高向量化的效率。建议在 HSW 上对 Stencil 尽量采用这种优化,即表 2 中的版本 F。

• Intrinsic load

使用 intrinsic 编程和使用其他 C/C++ 函数库类似,都需要包含正确的头文件,然后调用相应的 intrinsic 函数。编译器通常会将 Intrinsic 函数逐一翻译成汇编指令,不同之处在于使用汇编语言需要程序员自己管理寄存器,而使用 intrinsic 则还是由编译器来管理寄存器。

表 2 中版本 G 就是采用这种优化。以 3D7P 为例,首先,定义一组 `__m256d` 的数据,256 表示总长为 256bit, d 表示每个元素为双精度浮点数(64bit),因此这个数据包含了 4 个元素。其次,使用 `_mm256_load_pd()` 读入 4 个连续地址的双精度元素到 `__m256d` 中,且这 4 个元素需要 32 byte 对齐。使用 `_mm256_permute4x64_pd()` 设定边界条件。使用 `_mm256_set1_pd` 将某一个常数广播到其余 3 个元素上。使用 `_mm256_fmadd_pd()` 执行 $(a * b) + f2_t$ 的 FMA 操作。使用 `_mm256_store_pd()` 存储更新后的元素。最后,由于 `_mm256_load_pd()` 一次读入 4 个元素,因此内层循环的步长应该加 4,而不再是加 1。

5.3 优化结果

通过检查汇编代码,发现版本 C, D, E 中都存在 vgather 指令,而版本 A, B, F 和 G 则都没有。这说明在内存循环前添加 `#pragma simd` 进行向量化,由于边界条件设定中存在条件判断语句,编译器会生成 vgather 指令。同时, Tiling 和 cache oblivious 这些改进空间局部性的优化方法无法避免编译器生成 vgather 指令,只有使用 Peel 优化和 Intrinsic 优化才可以。

5 个 Stencil 不同优化版本的性能如图 5 所示。由此发现

AVX2 vgather 指令会降低 Stencil 在 HSW 上的性能。在消除 vgather 指令之后,各 Stencil 的版本 F 和 G 上均有 1.22X 至 3.88X 不等的性能提升。

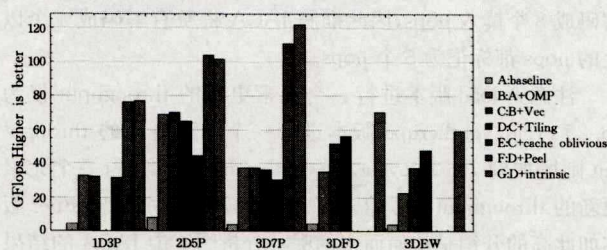


图 5 不同 Stencil 的优化版本在 Haswell 上的性能

6 vgather 指令实现分析

为了进一步探究 vgather 指令影响 Stencil 性能的原因,以 3D7P 为例,从 μops 层面分析 vgather 指令的实现。

6.1 使用 IACA 分析 vgather 的 μops

使用 IACA 检测 intrinsic load 以及 intrinsic vgather 的双精度和单精度版本。首先简单介绍如何解释 IACA 的结果。如图 6 所示, `vgatherqpd` 版本进行一个元素更新的 throughput 为 70.00 cycle。IACA 明确提示了 throughput 的瓶颈在 Gather 操作。 `vgatheqpd` 指令在端口 2 和端口 3 上都执行了 2.0 cycle 的 AGU 和 2.0 cycle 的读入 μops 。类似地, `vmadd132pd` 指令被解码为 2 个 μops 。在端口 0 和端口 1 上分别有 0.1 和 0.9 cycle,这是因为 HSW 在这 2 个端口上分别有一个 FMA 单元,而 IACA 显示的是统计结果,即假设有 100 条 FMA 指令,其中 10 条在端口 0 上执行,90 条在端口 1 上执行,那么统计结果就是端口 0 上有 0.1 个 cycle,端口 1 上有 0.9 个 cycle。另外,由于这条 FMA 指令还要访存,因此与 `vmovupd` 指令相似,它在端口 2 上执行了 1.0 cycle 的读入 μops 以及 1.0 cycle 的 AGU。最后的 `vmovpd` 指令是 store 操作,使用端口 2 和端口 3 的 AGU 生成地址的同时,在端口 4 执行 1.0 cycle 的写回 μop 。注意到虽然端口 7 是专门为 store 准备的 AGU,但这是一个精简版本的 AGU,不能处理地址中带 index 的 store 操作。 `vgatherdps` 版本进行一个元素更新的 throughput 更是高达了 105.00cycle,而且 IACA 也明确提示了 throughput 的瓶颈在 Gather 操作。而 load 版本进行一个元素更新的 throughput 仅为 16.25 cycle。限于篇幅,此处不能一一列出完整的 IACA 结果。完整的结果可以从 Github 上^[9]下载。

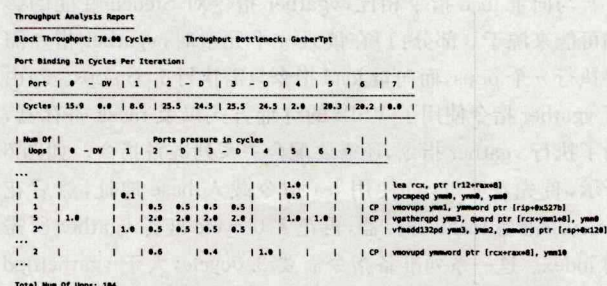


图 6 intrinsic vgather 双精度版本的 IACA 结果摘要

因为 `vgatherqpd` 指令需要读取 4 个双精度元素,而 IACA 显示 `vgatheqpd` 指令在端口 2 和端口 3 上都执行了 2.0 cycle 的 AGU 和 2.0 cycle 的读入 μops ;而 `vgatherdps` 指令需要读入 8 个单精度元素,而 IACA 显示 `vgatherdps` 指令在端口 2 和端口 3 上都执行了 4.0 cycle 的 AGU 和 4.0 cycle 的

读入 μops 。因此推测 `vgather` 指令会被解码成多个读入 μops ，每一个 μops 对应读入一个元素。IACA 提示 `vgather` 指令被解码成了 5 个 μops ，但实际上 `vgatherdps` 至少应该被解码成 8 个读入 μops ，因此推测 IACA 将所有解码成 4 个以上的 μops 都标记为 5 个 μops 。

注意到 `load` 版本进行一个元素更新的 throughput 仅为 16.25 cycle，`vgatherqpd` 版本进行一个元素更新的 throughput 则增加到 70.00 cycle，而 `vgatherdps` 版本进行一个元素更新的 throughput 更是高达了 105.00 cycle。那么 `vgather` 版本如此高的开销是从何而来的呢？分析图 6 中 IACA 的结果不难看出，`vgather` 版本在执行单元（即 Back-end）上花费的时间最高也就是 25.5 cycle，多出的 $70 - 25.5 = 44.5$ cycle 很可能是 Front-end 带来的。由于每条 `vgather` 指令至少会被解码成 4 个 μops 以上，因此推测这是由于使用 MSROM 进行解码所带来的高开销。

6.2 Vgather 对 Stencil 性能的影响

• 单核性能

为了了解清楚 `vgather` 指令对 stencil 性能的影响，对比了 `load` 版本和 `vgather` 版本在 HSW 单核上的性能。如图 7 所示，纵轴是性能 GFlops，横轴是 3D7P 立方体一个方向上点的数量，因此其所占的存储空间可由式(1)获得：

$$\text{Data_size} = \text{Cubic_size}^3 \times 2 \times 8 \quad (1)$$

其中，2 是因为使用了 Jacobi 方法更新元素，读和写各有一个数组；8 是因为双精度浮点大小为 8 byte。HSW 的 L1D 为 32kB，L2 为 256kB，L3 为 35MB，因此当 cubic size 小于 12 时，数据在 L1D 中；在 13~25 之间时，数据在 L2 中；26~129 之间时，数据在 L3 中；130 以上时，数据在内存中。当数据在 L1D 或 L2 中时，由于数据太小，不足以填满流水线，因此 `load` 版本和 `vgather` 版本都不能获得最高性能；当数据在 L3 时，`load` 版本的最高性能接近 `vgather` 版本性能的 3 倍；当数据在内存中时，`load` 版本的性能依然为 `vgather` 性能的 1.5 倍左右。

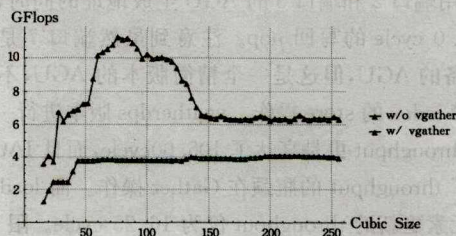


图 7 vgather 指令对单核性能的影响

与向量 `load` 指令相比，`vgather` 指令对 Stencil 性能的影响可能来源于 3 部分：1) 在读入 n 个元素时，`vgather` 指令需要执行 n 个 μops ，而向量 `load` 指令只需执行 1 个 μops 。2) 由于 `vgather` 指令使用了 VSIB 的寻址方式以及 mask 寄存器，为了执行 `vgather` 指令，还需要配合一系列准备指令。如图 6 所示，首先花 1.0 cycle 使用 `lea` 指令载入 base 地址，然后花 1.0 cycle 准备 mask 寄存器，再花 1.0 cycle 准备 `vgather` 所需的 index。这一系列准备指令需要 3.0 cycle，大于 `vgatherqpd` 本身的 2.0 cycle。而向量 `load` 指令则不需要这些准备指令。3) 如 6.1 节中分析的，`vgather` 指令由于需要被解码为多个 μops ，因此就由解码器中的 MSROM 进行解码，从而产生了较高的 Front-end 开销。

• 多核性能

`load` 版本和 `vgather` 版本在 HSW 多核上的性能对比如图 8 所示，由于 Stencil 代码受限于内存带宽，因此 `load` 版本在

9 核时达到了带宽饱和点 (Saturation Point)。而 `vgather` 版本的性能由于没有达到饱和点，因此在 14 个核上的性能达到了线性加速，但其 14 核的性能还不如 `load` 版本 4 核的性能。

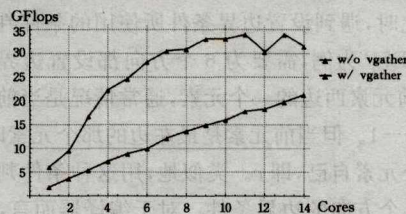


图 8 vgather 指令对多核性能的影响

结束语 本文通过评估 AVX2 `vgather` 指令在 Intel Haswell 上的性能，得到了以下 3 个结论：

- 1) AVX2 `vgather` 指令会降低 Stencil 的性能，这个结论对其他具有良好空间局部性的应用也适用。
- 2) `vgather` 指令会为每个要读入的元素生成一个 μops 。
- 3) 由于 `vgather` 指令需要由 MSROM 进行解码，因此会产生较高的 Front-end 开销。

据此，对 Stencil 这类具有良好空间局部性的应用在 Haswell 上的编程建议是要尽量避免让编译器生成 `vgather` 指令。有 2 种优化方式：①不要在指定元素地址时进行条件判断，推荐对 Stencil 使用 Peel 优化；②使用 intrinsic `load` 读入元素。

本文的工作有一些局限性，比如没有建立一个性能模型来细致分析 `vgather` 指令对 Stencil 性能的影响，这是下一步工作的主要方向之一。此外，Haswell 的下一代 Broadwell 对 `vgather` 指令的实现进行了一些改进，使用 Gather Index Table (GIT) 后大约能减少最多 60% 的 μop ，因此下一步计划对 GIT 进行评估。

致谢 感谢上海交通大学高性能计算中心 π 集群提供测试环境。

参考文献

- [1] HOFMANN J, TREIBIG J, HAGER G, et al. Performance Engineering for a Medical Imaging Application on the Intel Xeon Phi Accelerator[C]// 27th International Conference on Architecture of Computing Systems (ARCS2014). VDE, 2014.
- [2] PENNYCOOK S J, HUGHES C J, Smelyanskiy M, et al. Exploring SIMD for Molecular Dynamics, Using Intel Xeon Processors and Intel Xeon Phi Coprocessors[C]// IPDPS'13. IEEE, 2013; 1085-1097.
- [3] HOFMANN J, TREIBIG J, HAGER G, et al. Comparing the performance of different x86 SIMD instruction sets for a medical imaging application on modern multi- and manycorechips[C]// Proceedings of the 2014 Workshop on Programming models for SIMD/Vector Processing (WPMVP'14). New York, 2014.
- [4] KUSSWURM D. Modern X86 Assembly Language Programming 32bit, 64bit, SSE, and AVX[M]. Apress, 2014.
- [5] IACA[OL]. <https://software.intel.com/en-us/articles/intel-architecture-code-analyzer>.
- [6] 3DFD [OL]. <https://software.intel.com/en-us/articles/eight-optimizations-for-3-dimensional-finite-difference-3dfd-code-with-an-isotropic-iso>.
- [7] ZHANG C W J, TIAN Z P. and s-wave separated elastic wave equation numerical modeling using 2d staggered-grid[C]// SEG/San Antonio 2007 Annual Meeting, 2007.
- [8] AVX2-vgather 的部分源代码以及 IACA 结果[OL]. <https://github.com/jameslinsjtu/AVX2-vgather>.