Analysis of Intel's Haswell Microarchitecture Using The ECM Model and Microbenchmarks

J. Hofmann¹, D. Fey¹, J. Eitzinger², G. Hager² and G. Wellein²

 Computer Architecture, University Erlangen-Nuremberg
 Erlangen Regional Computing Center (RRZE), University Erlangen-Nuremberg johannes.hofmann@fau.de

Abstract. This paper presents an in-depth analysis of Intel's Haswell microarchitecture for streaming loop kernels. Among the new features examined is the dual-ring Uncore design, Cluster-on-Die mode, Uncore Frequency Scaling, core improvements as new and improved execution units, as well as improvements throughout the memory hierarchy. The Execution-Cache-Memory diagnostic performance model is used together with a generic set of microbenchmarks to quantify the efficiency of the microarchitecture. The set of microbenchmarks is chosen such that it can serve as a blueprint for other streaming loop kernels.

Keywords: Intel Haswell, Architecture Analysis, ECM Model, Performance Modeling

1 Introduction and Related Work

In accord with Intel's tick-tock model, where a tick corresponds to a shrinking of the process technology of an existing microarchitecture and a tock corresponds to a new microarchitecture, Haswell is a tock and thus represents a new microarchitecture. This means major changes to the preceding Ivy Bridge release have been made that justify a thorough analysis of the new architecture. This paper demonstrates how the Execution-Cache-Memory (ECM) diagnostic performance model [3.11.2.10] can be used as a tool to evaluate and quantify the efficiency of a microarchitecture. The ECM model is a resource-centric model that allows to quantify the runtime of a given loop kernel on a specific architecture. It requires detailed architectural specifications and an instruction throughput prediction as input. It assumes Perfect instruction level parallelism for instruction execution as well as bandwidth-bound data transfers. As a consequence the model yields a practical upper limit for single core performance. The only empirically determined input for the model is that of sustained memory bandwidth, which can be different for each benchmark. The model quantifies different runtime contributions from instruction execution and data transfers within the complete memory hierarchy as well as potential overlap between contributions. Runtime contributions are divided into two different categories: $T_{\rm nOL}$, i.e. cycles in which the core executes instructions that forbid simultaneous transfer of data between the L1 and L2 caches; and $T_{\rm OL}$, i.e. cycles that do not contain non-overlapping instructions, thus allowing for simultaneous instruction execution and data transfers between L1 and L2 caches. Note that one improvement to the original ECM model

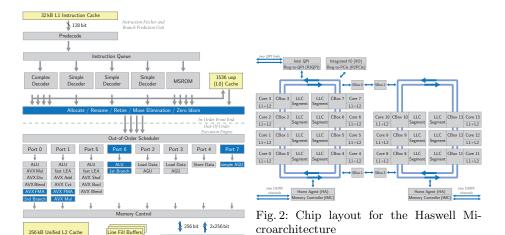


Fig. 1: Core design for the Haswell Microarchitecture

is that apart from load instructions, store instructions are now also considered non-overlapping. Instruction times as well as data transfer times, e.g. $T_{\rm L1L2}$ for the time required to transfer data between L1 and L2 caches, can be summarized in shorthand notation: $\{T_{\rm OL} \mid T_{\rm nOL} \mid T_{\rm L1L2} \mid T_{\rm L2L3} \mid T_{\rm L3Mem}\}$. The in-core execution time $T_{\rm core}$ is the maximum of either overlapping or non-overlapping instructions. Predictions for cache/memory levels is given by $\max(T_{\rm OL}, T_{\rm nOL} + T_{\rm data})$ with $T_{\rm data}$ the sum of the individual contributions up to the cache/memory level under consideration, e.g. for the L3 cache $T_{\rm data} = T_{\rm L1L2} + T_{\rm L2L3}$. A similar shorthand notation exists for the model's prediction: $\{T_{\rm core} \mid T_{\rm L2} \mid T_{\rm L3} \mid T_{\rm Mem}\}$. For details on the ECM model refer to the previously provided references.

Related work covers in-detail analysis of architectural features using microbenchmarks, e.g., [1,7,9]. We are not aware of any work though using an analytic model to quantify the efficiency of a microarchitecture.

Section 2 presents major improvements in Intel Haswell. In Section 3 we introduce a comprehensive set of microbenchmarks that serves as a blueprint for streaming loop kernels. To evaluate the hardware, obtained measurements are correlated with the performance predictions in Section 4.

2 Haswell Microarchitecture

2.1 Core Design

Fig. 1 shows a simplified core design of the Haswell microarchitecture with selected changes to previous microarchitectures highlighted in blue. Due to lack of space we focus on new features relevant for streaming loop kernels.

The width of all three data paths between the L1 cache and processor registers has been doubled in size from $16\,\mathrm{B}$ to $32\,\mathrm{B}$. This means that two Advanced Vector Extensions (AVX) loads and one store ($32\,\mathrm{B}$ in size) can now retire in a

single clock cycle as opposed to two clock cycles required on previous architectures. The data path between the L1 and L2 caches has been widened from $32\,\mathrm{B}$ to $64\,\mathrm{B}$.

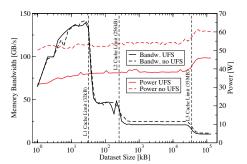
While the core is still limited to retire four μ ops per cycle, the number of issue ports has been increased from six to eight. The newly introduced port 6 contains the primary branch unit; a secondary unit has been added to port 0. In previous designs only a single branch unit was available and located on port 5. By moving it to a dedicated port, port 5—which is the only port that can perform AVX shuffle operations—is freed up. Adding a secondary branch unit benefits branch-intensive codes. The other new port is port 7, which houses a so-called simple Address Generation Unit (AGU). This unit was made necessary by the increase in register-L1 bandwidth. Using AVX on Sandy Bridge and Ivy Bridge, two AGUs were sufficient, because each load or store required two cycles to complete, not making it necessary to compute three new addresses every cycle, but only every second cycle. With Haswell this has changed, because potentially a maximum of three load/store operations can now retire in a single cycle, making a third AGU necessary. Unfortunately, this simple AGU can not perform the necessary addressing operations required for streaming kernels on its own (see Section 4.3 for more details).

Apart from adding additional ports, Intel also extended existing ones with new functionality. Instructions introduced by the Fused Multiply-Add (FMA) Instruction Set Architecture (ISA) extension are handled by two new, AVX-capable units on ports 0 and 1. Haswell is the first architecture to feature the AVX2 ISA extension and introduces a second AVX multiplication unit on port 1 while there is still just one low-latency add unit.

2.2 Package Layout

Figure 2 shows the layout of a 14-core Haswell processor package. Apart from the processor cores, the package consists of what Intel refers to as the Uncore. Attached to each core and its private L1 and L2 caches, there is a Last-Level Cache (LLC) segment, that can hold 2.5 MB of data. The physical proximity of core and cache segment does however not imply that data used by a core is stored exclusively or even preferably in its LLC segment. Data is placed in all LLC segments according to a proprietary hash function that is supposed to provide uniform distribution of data and prevent hotspots for a wide range of data access patterns. An added benefit of this design is that single-threaded applications can make use of the full accumulated LLC capacity.

The cores and LLC segments are connected to a bidirectional ring interconnect that can transfer one Cache Line (CL) (64 B in size) every two cycles in each direction. In order to reduce latency, the cores are arranged to form two rings, which are connected via two queues. To each ring belongs a Home Agent (HA) which is responsible for cache snooping operations and reordering of memory requests to optimize memory performance. Attached to each HA is a Memory Controller (MC), each featuring two 8 byte-wide DDR4 memory channels. Also accessible via the ring interconnect are the on-die PCIe and QPI facilities.



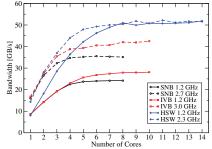


Fig. 3: Impact of UFS on Bandwidth and Power Usage.

Fig. 4: Stream Triad Bandwidth as Function of Frequency.

Haswell introduces an on-die Fully Integrated Voltage Regulators (FIVR). This FIVR draws significantly less power than on previous microarchitectures, because it enabled for faster switching of power-saving states. It also allows a more fine-grained control of CPU states: instead of globally setting the CPU frequencies for all cores within a package, Haswell can now set core frequencies and sleep states individually.

2.3 Uncore Frequency Scaling

In the new Haswell microarchitecture, Intel reverted from Sandy and Ivy Bridges' unified clock domain for core and the Uncore to the Nehalem design of having two separate clock domains [4,5]. Haswell introduced a feature called Uncore Frequency Scaling (UFS), in which the Uncore frequency is dynamically scaled based on the number of stall cycles in the CPU cores. Despite reintroducing higher latencies, the separate clock domain for the Uncore offers a significant potential for power saving, especially for serial codes. Fig. 3 shows the measured sustained bandwidth (left y-axis) for the Schönauer vector triad (cf. Table 1) using a single core along with the power consumption (right y-axis) for varying dataset sizes. As expected the performance is not influenced by whether UFS is active or not when data resides in a core's private caches; however, power requirements are reduced by about 30%! Although we observe a difference in performance as soon as the LLC is involved, the performance impact is limited. The bandwidth drops from 24 to 21 GB/s (about 13%) in the LLC, but power usage is reduced from 55 W to 40 W (about 27%).

2.4 Memory

Intel's previous microarchitectures show a strong correlation between CPU frequency and sustained memory bandwidth. Fig. 4 shows the measured chip bandwidth for the Stream Triad (cf. Table 1)—adjusted by a factor of 1.3 to account for the write-allocate when storing—on the Sandy Bridge, Ivy Bridge, and Haswell microarchitectures. For each system, the bandwidth was measured using the lowest possible frequency (1.2 GHz) and the advertised nominal clock

speed. While we find differences around 33% in the maximum achievable sustained memory bandwidth depending on the CPU frequency for Sandy and Ivy Bridge, on Haswell we can observe a frequency-independent sustained bandwidth of 52.3 GB/s. On Haswell the CPU frequency can be lowered, thereby decreasing power consumption, while the memory bandwidth stays constant. Further research regarding the Stream Triad with a working set size of 10 GB has shown that Haswell offers an improvement of 23% respectively 12% over the Sandy and Ivy Bridge architectures when it comes to energy consumption and 55% respectively 35% in terms of Energy-Delay Product [3].

2.5 Cluster on Die

In Cluster on Die (CoD) mode, cores get equally separated into two ccNUMA memory domains. This means that instead of distributing requests between both memory controllers each core is assigned a dedicated memory controller. To keep latencies low, the strategy is to make a core access main memory through the memory controller attached to its ring. However, with memory domains being equal in size, the asymmetric core count on the two physical rings makes exceptions necessary. In the design shown in Fig. 2 the 14 cores are divided into two memory domains of 7 cores each. Using microbenchmarks and likwid-perfctr [12] to access performance counters in order to measure the number of memory accesses for each individual memory channel, we find that cores 0–6 access main memory through the memory channels associated with the memory controller on the left ring, and cores 7–13 those associated with the memory controller on ring 1. Thus, only core number 7 has to take a detour across rings to access data from main memory. Note that with CoD active the LLC also is divided. As each domain contains seven LLC segments (2.5 MB each), the total amount of LLC for each domain is only 17.5 MB instead of 35 MB.

CoD mode is intended for NUMA-optimized codes and serves two purposes: First latency is decreased by reducing the number of endpoints in the memory domain. Instead of 14 LLC segments, data will be distributed in only 7 segments inside each memory domain, thereby decreasing the mean hop count. Also, the requirement to pass through the two buffers connecting the rings is eliminated for all but one LLC segment. Also, bandwidth is increased by reducing the probability of ring collisions by lowering participant count from 14 to 7.

3 Microbenchmarks

A set of microbenchmarks chosen to provide a good coverage of relevant data access patterns was used to evaluate the Haswell microarchitecture and is summarized in Table 1. For each benchmark, the table lists the number of load and store streams—the former being divided into explicit and Read for Ownership (RFO) streams. RFO refers to implicit loads that occur whenever a store miss in the current cache triggers a write-allocate. On Intel architectures all cache levels use a write-allocate strategy on store misses. The table also includes the predictions of the ECM model and the actually measured runtimes in cycles

Table 1: Overview of microbenchmarks: Loop Body, Memory Streams, ECM prediction and Measurement in c/CL, and Model Error.

		Load Streams	Write	ECM Prediction	Measurement	Model Error
Benchmark	Description	Explicit / RFO	Streams	L1/L2/L3/Mem	L1/L2/L3/Mem	L1/L2/L3/Mem
ddot	s+=A[i]*B[i]	2 / 0	0	{2 4 8 17.1}	2.1] 4.7] 9.6] 19.4	$5\% \; \rceil \; 17\% \; \rceil \; 20\% \; \rceil \; 13\%$
load	s+=A[i]	1 / 0	0	$\{2 \mid 2 \mid 4 \mid 8.5\}$	$2 \mid 2.3 \mid 5 \mid 10.5$	$0\% \mid 15\% \mid 25\% \mid 23\%$
store	A[i]=s	0 / 1	1	{2 4 8 20.5}	2] 6] 8.2] 17.7	$0\% \; \rceil \; 33\% \; \rceil \; 3\% \; \rceil \; 16\%$
update	A[i]=s*A[i]	1 / 0	1	$\{2 \mid 4 \mid 8 \mid 20.5\}$	2.1] 6.5] 8.3] 17.6	$5\% \mid 38\% \mid 4\% \mid 16\%$
copy	A[i]=B[i]	1 / 1	1	$\{2\; \; 5\; \; 11\; \; 27.8\}$	2.1 8 13 27	5%] $38%$] $15%$] $3%$
STREAM triad A[i]=B[i]+s*C[i]		2 / 1	1	$\{3 \; \rceil \; 7 \; \rceil \; 15 \; \rceil \; 36.7\}$	$3.1 \mid 10 \mid 17.5 \mid 37$	3%] $30%$] $14%$] $1%$
Schönauer tri	ad A[i]=B[i]+C[i]*D[i]	3 / 1	1	{4] 9] 19] 45.5}	4.1] 11.9] 21.9] 46.8	3%] 24%] 13%] 3%

along with a quantification of the model's error. In the following, we discuss the ECM model for each of the kernels and show how to arrive at the prediction shown in the table.

The sustained bandwidths used to derive the L3-memory cycles per CL inputs can be different for each benchmark, which is why for each individual kernel the sustained bandwidth is determined using a benchmark with the exact data access pattern that is modeled; note that for our measurements CoD mode was active and the measured bandwidth corresponds to that of a single memory domain.

3.1 Dot Product and Load

The dot product benchmark ddot makes use of the new FMA instructions introduced in the FMA3 ISA extension implemented in the Haswell microarchitecture. $T_{\rm nOL}$ is two clock cycles, because the core has to load two CLs (A and B) from L1 to registers using four AVX loads (which can be processed in two clock cycles, because each individual AVX load can be retired in a single clock cycle and there are two load ports). Processing data from the CLs using two AVX FMA instructions only takes one clock cycle, because both issue ports 0 and 1 feature AVX FMA units. A total of two CLs has to be transfered between the adjacent cache levels. At 64 B/c this means 2 c to transfer the CLs from L2 to L1. Transferring the CLs from L3 to L2 takes 4 c at 32 B/c. The empirically determined sustained (memory domain) bandwidth is 32.4 GB/s. At 2.3 GHz, this corresponds to a bandwidth of about 64 B/CL · 2.3 GHz/32.4 GB/s ≈ 4.5 c/CL or 9.1 c for two CLs. The ECM model input is thus $\{1 \mid 2 \mid 2 \mid 4 \mid 9.1\}$ c and the corresponding prediction is $\{2 \mid 4 \mid 8 \mid 17.1\}$ c.

For the load kernel the two AVX loads to get the CL containing A from L1 can be retired in a single cycle, yielding $T_{\rm nOL}=1\,\rm c$. With only a single AVX add unit available on port 1, processing the data takes $T_{\rm OL}=2\,\rm c$. Because only a single CL has to be transferred between adjacent cache levels and the measured bandwidth corresponds exactly to that of the ddot kernel, the time required is exactly half of that needed for the ddot benchmark. The ECM model input for this benchmark is $\{2 \parallel 1 \mid 1 \mid 2 \mid 4.5\}$ c, yielding a prediction of $\{2 \mid 2 \mid 4 \mid 8.5\}$ c.

3.2 Store, Update, and Copy

For the *store* kernel, two AVX stores are required per CL. With only a single store unit available, $T_{\rm nOL}=2\,{\rm c}$; as there are no other instructions such as arithmetic

operations, $T_{\rm nOL}$ is zero. When examining CL transfers along the cache hierarchy, we have to bear in mind that a store-miss will trigger a write-allocate, resulting in two CL transfers for each CL update: one to write-allocate the CL which data gets written to and one to evict the modified CL once the cache becomes full. This results in a transfer time of 2 c to move the data between the L1 and L2 cache and a transfer time of 4 c for L2 and L3. The sustained bandwidth of 23.6 GB/s (corresponding to approximately $6.2 \, \text{c/CL}$) for a kernel involving evictions is significantly worse than that of the previous load-only kernels. The resulting ECM input and prediction are $\{0 \, \| \, 2 \, \| \, 2 \, \| \, 4 \, \| \, 1 \, 2.5 \}$ c respectively $\{2 \, \| \, 4 \, \| \, 8 \, \| \, 20.5 \}$ c.

For the *update* kernel, two AVX stores and two AVX loads are required. Limited by a single store port, $T_{\rm nOL}=2\,\rm c.$ The multiplications take $T_{\rm OL}=2\,\rm c.^3$ The number of CL transfers is identical to that of the *store* kernel, the only difference being that the CL load is caused by explicit loads and not a write-allocate. With a memory bandwidth almost identical to that of the *store* kernel, the time to transfer a CL between L3 and memory again is approximately $6.2\,\rm c/CL$, yielding an ECM input of $\{2\,\|\,2\,\|\,2\,\|\,4\,\|\,12.5\}$ c and a prediction that is identical to that of the *store* kernel.

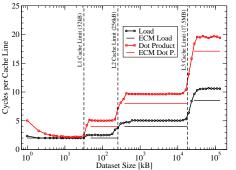
The copy kernel has to perform two AVX loads and two AVX stores to copy one CL. The single store port is the bottleneck, yielding $T_{\rm nOL}=2\,{\rm c}$; absent arithmetic instructions $T_{\rm nOL}$ is zero. Three CLs have to be transferred between adjacent cache levels: load B, write-allocate and evict A. This results in a requirement of 3 c for L1–L2 transfers and 6 c for L2–L3 transfers. With a sustained memory bandwidth of 26.3 GB/s the time to transfer one CL between main memory and LLC is approximately 5.6 c/CL or 16.8 c for three CLs. This results in the following input for the ECM model $\{0\,\|\,2\,\|\,3\,\|\,6\,\|\,16.8\}$ c, which in turn yields a prediction of $\{2\,\|\,5\,\|\,11\,\|\,27.8\}$ c.

3.3 Stream Triad and Schönauer Triad

For the STREAM Triad [6], the AGUs prove to be the bottleneck: it is impossible to retire two AVX loads and an AVX store that use indexed addressing in the same cycle, because there are only two full AGUs available supporting this addressing mode. The resulting $T_{\rm nOL}$ thus is not 2 but 3 c to issue four AVX loads (two each for CLs containing B and C) and two AVX stores (two for CL A). Both FMAs can be retired in one cycle, because two AVX FMA units are available, yielding $T_{\rm OL}=1$ c. Traffic between adjacent cache levels is 4 CLs: load CLs containing B and C, write-allocate and evict the CL containing A. The measured sustained bandwidth of 27.1 GB/s corresponds to approximately 5.4 c/CL—or about 21.7 c for all four CLs. The input parameters for the ECM model are thus $\{1 \mid 3 \mid 4 \mid 8 \mid 21.7\}$ c leading to the follow prediction: $\{3 \mid 7 \mid 15 \mid 36.7\}$ c.

For the Schönauer Triad [8], again the AGUs are the bottleneck. Six AVX loads (CLs B, C, and D) and two AVX stores (CL A) have to be performed;

³ Normally, with two AVX mul ports available, $T_{\rm OL}$ should be 1 c. However, the frontend can only retire 4 μ ops/c; this, along with the fact that stores count as 2 μ ops, means that if both multiplications were paired with the first store, there would not be enough full AGUs to retire the second store and the remaining AVX load instructions in the same cycle.



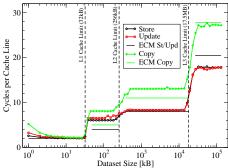


Fig. 5: ECM predictions and measurement Fig. 6: ECM predictions and measurement results for load and dot product kernels.

results for store, update, and copy kernels.

these eight instructions have to share two AGUs, resulting in $T_{\text{nOL}}=4\,\text{c}$. The two AVX FMAs can be performed in a single cycle, yielding $T_{\rm OL}=1\,{\rm c}$. Data transfers between adjacent caches correspond to five CLs: B, C, and D require loading while CL A needs to be write-allocated and evicted. For the L1 cache, this results in a transfer time of 5 c. The L2 cache transfer time is 10 cycles. The measured sustained memory bandwidth of 27.8 GB/s corresponds to about 5.3 c/CL or 26.5 c for all five CLs. The resulting ECM input parameters are thus $\{1 \| 4 | 5 | 10 | 26.5\}$ c and the resulting prediction is $\{4 | 9 | 19 | 45.5\}$ c.

4 Results

The results presented in this section were obtained using hand-written assembly kernels that were benchmarked using the likwid-bench tool [12]. No software prefetching was used in the code; the results therefore show the ability of the hardware prefetchers to hide data access latencies. The machine used for benchmarking was a standard two-socket server using Xeon E5-2695 v3 chips, featuring 14 cores each. Each core comes with its own 32 kB private L1 and 256 kB private L2 caches; the shared LLC is 35 MB in size. Each chip features four DDR4-2166 memory channels, adding up to a theoretical memory bandwidth of 69.3 GB/s per socket or 138.6 GB/s for the full node. For all benchmarks, the clock frequency was fixed at the nominal frequency of 2.3 GHz, CoD was activated, and UFS was disabled.

4.1 Dot Product and Load

Fig. 5 illustrates ECM predictions and measurement results for both the load and ddot benchmarks. While core execution time for both benchmarks is two cycles as predicted by the model, dot performance is slightly lower than predicted with data coming from the L2 cache. The worse than expected L2 cache performance has been a general problem with Haswell. In contrast to Haswell, Sandy and Ivy Bridge delivered the advertised bandwidth of 32 B/c [10]. On Haswell, in none of the cases the measured L2 bandwidth could live up to the advertised 64 B/c.

For the load kernel, the performance in L2 is almost identical to that with data residing in the L1 cache: this is because the CL can theoretically be transfered from L2 to L1 a single cycle at $64\,\mathrm{B/c}$, which is exactly the amount of slack that is the difference between $T_{\mathrm{OL}}=2\,\mathrm{c}$ and $T_{\mathrm{nOL}}=1\,\mathrm{c}$. In practise, however, we observe a small penalty of $0.3\,\mathrm{c/CL}$, so again, we do can observe the specified bandwidth of $64\,\mathrm{B/c}$.

As soon as the working set becomes too large for the core-local L2 cache, the ECM prediction is slightly off. For kernels with a low number of cycles per CL an empirically determined penalty for transferring data from off-core locations was found to be one cycle per load stream and cache-level, e.g. 2 c for the *ddot* benchmark with data residing in L3 and 4 c with data from memory. In all likelihood, this can be attributed to latencies introduced when data is passing between different clock domains (e.g. core, cbox, mbox) that cannot be entirely hidden for kernels with a very low core cycle count.

4.2 Store, Update, and Copy

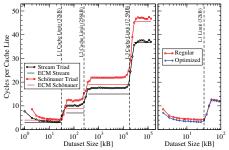
Fig. 6 shows ECM predictions and measurements for the *store*, *update*, and *copy* kernels. With data in L1 cache, measurements for all three benchmarks match the prediction. In the L2 cache, measured performance is off about one cycle per stream: two cycles for the *store* and *update* benchmarks, and four cycles for the *copy* benchmark. This means that it takes the data exactly twice as long to be transfered than what would be the case assuming a bandwidth of 64 B/c.

Measurements in L3 for the *store* and *update* kernel fit the prediction. This suggests either that either overlap between transfers is happening or some other undocumented optimization is taking place, as we would normally expect the poor L2 performance to trickle down to L3 and memory measurements (as is the case for the *copy* kernel). Suspicion about overlap or undocumented improvements is substantiated by better than expected in-memory performance.

4.3 Stream Triad and Schönauer Triad

Fig. 7 shows model predictions and measurements for both the Stream and Schönauer Triads. The measurement fits the model's prediction for data in the L1 cache. We observe the same penalty for data in the L2 cache. This time, the penalty also propagates: measurement and prediction for data in L3 is still off. The match of measurement and prediction for the in-memory case suggests either overlap of transfers or other unknown optimization as was the case before for the *store*, *update*, and *copy* kernels.

In addition, Fig. 7 shows measurement results for the naive Schönauer Triad as it is currently generated by compilers (e.g. the Intel C Compiler 15.0.1) and an optimized version that makes use of the newly introduced simple AGU on port 7. Typically, address calculations in loop-unrolled streaming kernels require two steps: scaling and offset computation. Both AGUs on ports 2 and 3 support this addressing mode called "base plus index plus offset." The new simple AGU can only perform offset computations. However, it is possible to make use of this AGU by using one of the "fast LEA" units (which can perform only indexed



```
lea rbx, [r8+rax*8]
vmovapd ymm0, [rsi+rax*8]
vmovapd ymm1, [rsi+rax*8+32]
vmovapd ymm8, [rdx+rax*8]
vmovapd ymm9, [rdx+rax*8+32]
vfmadd231pd ymm0,ymm8,[rcx+rax*8]
vfmadd231pd ymm1,ymm9,[rcx+rax*8+32]
vmovapd [rbx], ymm0
vmovapd [rbx+32], ymm1
```

(left) and comparison of naive and opti- kernel. mized Schönauer Triad (right).

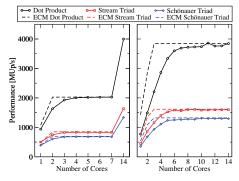
Listing 1.1: Shortened two-way unrolled, Fig. 7: ECM predictions and measurement hand-optimized code for Schönauer Triad. results for Stream and Schönauer Triads Eight-way unrolling used in real benchmark

and no offset addressing) to pre-compute an intermediary address. This precomputed address is fed to the simple AGU, which can then perform the still outstanding offset addition. Using all three AGUs, it is possible to complete the eight addressing operations in three instead of four cycles. The assembly code for this optimized version is shown in Listing 1.1.

Multi-Core Scaling and Cluster-on-Die Mode

When using the ECM model to estimate multi-core performance, single-core performance is scaled until a bottleneck is hit—which currently on Intel CPUs is main memory bandwidth. Fig. 8 shows ECM predictions along with actual measurements for the ddot, Stream Triad, and Schönauer Triad kernels using both CoD and non-CoD modes. The L3-Memory CL transfer time used for each prediction is based on the sustained bandwidth of the CoD respectively non-CoD mode. While the measurement fits the prediction in CoD mode, we find a non-negligible discrepancy in non-CoD mode. This demonstrates how the ECM model can be used to uncover the source of performance deviations. In non-CoD mode, the kernel execution time is no longer just made up of in-core execution and bandwidth-limited data transfers as predicted by the model. Although we can only speculate, we attribute the penalty cycles encountered in non-CoD mode to higher latencies caused by longer ways to the memory controllers: due to equal distribution of memory requests, on average every second request has to go the "long way" across rings, which is not the case in CoD mode.

The measurements indicate that peak performance for both modes is nearly identical, e.g. for ddot performance saturates slightly below 4000 MUp/s for non-CoD mode while CoD saturates slightly above the 4000 mark. Although the plots indicate the bandwidth saturation point is reached earlier in CoD mode, this conclusion is deceiving. While it only takes four cores to saturate the memory bandwidth of an memory domain, a single domain is only using two memory controllers; thus, saturating chip bandwidth requires 2×4 threads to saturate both memory domains, the same amount of cores it takes to achieve the sustained bandwidth in non-CoD mode.



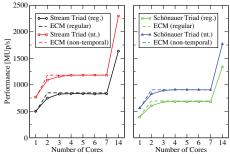


Fig. 8: Core-Scaling using CoD mode (left) and non-CoD mode (right).

Fig. 9: Performance using regular vs. nontemporal stores for Stream (left) and Schönauer Triads (right).

4.5 Non-Temporal Stores

For streaming kernels with dataset sizes that do not fit into the LLC it is imperative to use non-temporal stores in order to achieve the best performance. Not only is the total amount of data to be transfered from memory reduced by getting rid of RFO stream(s), but in addition, data does not have to travel through the whole cache hierarchy. On Haswell, non-temporal stores are sent to the L1 cache by the core, just like regular stores; they do however not update any entries in the L1 cache but are relayed to core-private line fill buffers, from which data is transfered directly to main memory.

Fig. 9 shows the performance gain offered by non-temporal stores. The left part shows the Stream Triad, which using regular stores is made up of two explicit load streams for arrays B and C plus a store and an implicit RFO stream for array A. Looking at transfered data volumes, we expect an performance increase by a factor of $1.33\times$, because using non-temporal stores gets rid of the RFO stream, thereby reducing streams count from four to three. However, the measured speedup is higher: 1181 vs. $831\,\mathrm{MUp/s}$ $(1.42\times)$ using a single memory domain respectively 2298 vs $1636\,\mathrm{MUp/s}$ $(1.40\times)$ when using a full chip. A possible explanation for this higher than anticipated speedup is that we have observed the efficiency of the memory subsystem degrade with an increasing number of streams. Vice verse, we could conclude that the efficiency increases by getting rid of the RFO stream.

A similar behavior is observed for the Schönauer Triad. Data volume analysis suggests a performance increase of $1.25 \times$ (4 streams instead of 5). However, the measured performance using non-temporal stores is 905 vs. $681\,\mathrm{GUp/s}$ ($1.33 \times$) using one memory domain resp. 1770 vs. $1339\,\mathrm{MUp/s}$ ($1.32 \times$) using a full chip.

5 Conclusion

This paper investigated new architectural features of the Intel Haswell microarchitecture with regard to the execution of streaming loop kernels. It demonstrated how to employ the ECM model together with microbenchmarking to

quantify the efficiency of architectural features. On the example of a comprehensive set of streaming loop kernels deviations from official specifications as well as the overall efficiency was evaluated. Besides incremental improvements and core related things Haswell addresses two main areas: Energy efficiency and to provide low latency data access within the chip while increasing the core count. Sustained main memory bandwidth is no longer impaired by the selection of low clock frequencies, enabling power savings of more than 20% respectively 10% over the previous Sandy respectively Ivy Bridge architectures. Uncore Frequency Scaling can further improve power savings by more than 20% for single-core workloads at no cost for data in core-private caches respectively a small performance penalty with data off-core. The new Cluster-on-Die mode offers performance improvements for single-threaded and parallel memory-bound codes, and has major benefits with regard to latency penalties.

References

- 1. Cache Coherence Protocol and Memory Performance of the Intel Haswell-EP Architecture. IEEE (2015)
- Hager, G., Treibig, J., Habich, J., Wellein, G.: Exploring performance and power properties of modern multicore chips via simple machine models. Concurrency Computat.: Pract. Exper. (2013), DOI: 10.1002/cpe.3180
- 3. Hofmann, J., Treibig, J., Fey, D.: Execution-cache-memory performance model: Introduction and validation (2015)
- 4. Intel Corporation: Intel Xeon Processor E52600/4600 Product Family Technical Overview,
 https://software.intel.com/en-us/articles/intel-xeon-processor-e5-26004600-product-family-technic
- 5. Intel Corporation: Intel Technology Journal, Vol. 14, Issue 3 (2010)
- McCalpin, J.D.: Memory bandwidth and machine balance in current high performance computers. IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter pp. 19–25 (Dec 1995)
- Molka, D., Hackenberg, D., Schöne, R.: Main memory and cache performance of intel sandy bridge and amd bulldozer. In: Proceedings of the Workshop on Memory Systems Performance and Correctness. pp. 4:1–4:10. MSPC '14, ACM (2014)
- 8. Schönauer, W.: Scientific Supercomputing: Architecture and Use of Shared and Distributed Memory Parallel Computers. Self-edition (2000)
- Schöne, R., Hackenberg, D., Molka, D.: Memory performance at reduced cpu clock speeds: An analysis of current x86_64 processors. In: Proceedings of the 2012 USENIX Conference on Power-Aware Computing and Systems. pp. 9–9. Hot-Power'12, USENIX Association (2012)
- Stengel, H., Treibig, J., Hager, G., Wellein, G.: Quantifying performance bottlenecks of stencil computations using the Execution-Cache-Memory model. In: Proceedings of the 29th ACM International Conference on Supercomputing. ICS '15, ACM, New York, NY, USA (2015), http://doi.acm.org/10.1145/2751205.2751240
- Treibig, J., Hager, G.: Introducing a performance model for bandwidth-limited loop kernels. In: Wyrzykowski, R., Dongarra, J., Karczewski, K., Wasniewski, J. (eds.) Parallel Processing and Applied Mathematics. Lecture Notes in Computer Science, vol. 6067, pp. 615–624. Springer Berlin / Heidelberg (2010)

12. Treibig, J., Hager, G., Wellein, G.: likwid-bench: An extensible microbenchmarking platform for x86 multicore compute nodes. In: Parallel Tools Workshop. pp. 27–36 (2011)