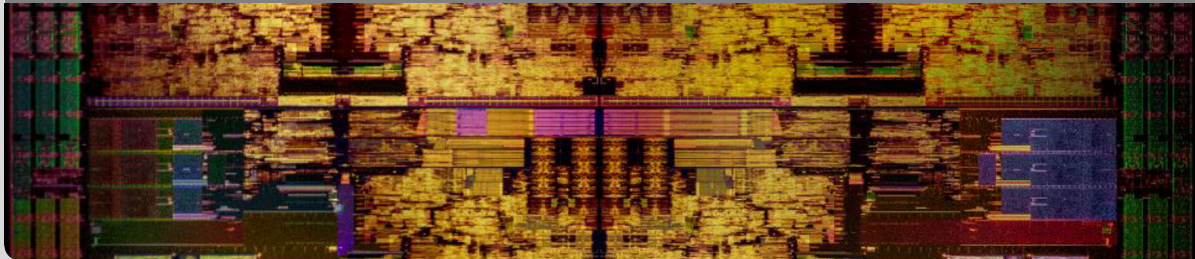


Analysis and Optimization of Dynamic Voltage and Frequency Scaling for AVX Workloads Using a Software-Based Reimplementation

Yussuf Khalil | 25 September 2019

OPERATING SYSTEMS GROUP



Motivation

- AVX: modern vector processing extension for x86 processors
- Intel CPUs reduce clock frequency when executing AVX code
- Lowered frequencies retained after last AVX instruction
 - ⇒ negative impact on performance in heterogeneous workloads

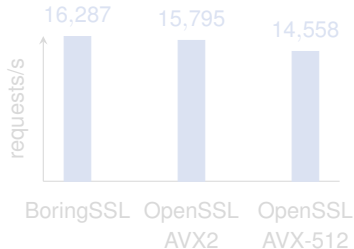


Figure: nginx throughput with different ChaCha20-Poly1305 implementations¹

¹Vlad Krasnov. *On the dangers of Intel's frequency scaling*. 2017. URL: <https://blog.cloudflare.com/on-the-dangers-of-intels-frequency-scaling/>.

Motivation

- AVX: modern vector processing extension for x86 processors
- Intel CPUs reduce clock frequency when executing AVX code
- Lowered frequencies retained after last AVX instruction
 - ⇒ negative impact on performance in heterogeneous workloads

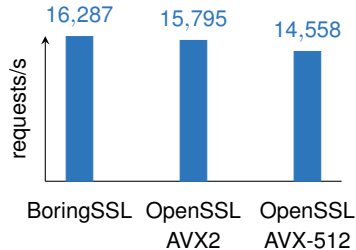


Figure: nginx throughput with different ChaCha20-Poly1305 implementations¹

¹Vlad Krasnov. *On the dangers of Intel's frequency scaling*. 2017. URL:
<https://blog.cloudflare.com/on-the-dangers-of-intels-frequency-scaling/>.

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- Frequency switches are costly
 - ⇒ wait before raising frequency after AVX execution
- What if we could predict when AVX code will be executed again?
- Could we raise the frequency earlier and improve performance during scalar phases?
- Similar approaches are used for device power management

Try different reclocking algorithms!

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- Frequency switches are costly
 - ⇒ wait before raising frequency after AVX execution
- What if we could predict when AVX code will be executed again?
- Could we raise the frequency earlier and improve performance during scalar phases?
- Similar approaches are used for device power management

Try different reclocking algorithms!

- Frequency switches are costly
 - ⇒ wait before raising frequency after AVX execution
- What if we could predict when AVX code will be executed again?
- Could we raise the frequency earlier and improve performance during scalar phases?
- Similar approaches are used for device power management

Try different reclocking algorithms!

Yussuf Khalil

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- AVX reclocking is controlled solely by CPU itself
 - We can not modify the hardware
 - But we can (mostly) disable automatic reclocking
 - ⇒ reimplement algorithm in software
 - Allows to derive and test alternative algorithms
 - Needed to measure overhead of approach
- Intel only provides vague description of algorithm²
 - Downclocking after $\leq 500 \mu\text{s}$
 - Upclocking after 2 ms
 - Two AVX frequency levels (“turbo licenses”)
- ⇒ need to analyze what processors really do to create reimplementation

²Intel 64 and IA-32 Architectures Optimization Reference Manual. Intel Corporation. Apr. 2019.

Yussuf Khalil

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- AVX reclocking is controlled solely by CPU itself
 - We can not modify the hardware
 - But we can (mostly) disable automatic reclocking
 - ⇒ reimplement algorithm in software
 - Allows to derive and test alternative algorithms
 - Needed to measure overhead of approach
- Intel only provides vague description of algorithm²
 - Downclocking after $\leq 500 \mu\text{s}$
 - Upclocking after 2 ms
 - Two AVX frequency levels (“turbo licenses”)
 - ⇒ need to analyze what processors really do to create reimplementation

² *Intel 64 and IA-32 Architectures Optimization Reference Manual*. Intel Corporation. Apr. 2019.

Methodology

We want to find out

- when exactly a core will reduce or raise its frequency,
- how long it needs to do that,
- and whether Intel's description is correct and complete.

How can we measure that?

- Use performance counters
- Run synthetic code snippets designed to trigger specific behavior

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

Methodology

We want to find out

- when exactly a core will reduce or raise its frequency,
- how long it needs to do that,
- and whether Intel's description is correct and complete.

How can we measure that?

- Use performance counters
- Run synthetic code snippets designed to trigger specific behavior

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

Yussuf Khalil

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

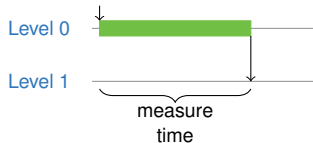
Conclusion

- Linux kernel module
 - configures Performance Monitoring Unit (PMU)
 - handles interrupts
- User-space program
 - instructs kernel component
 - conducts measurements

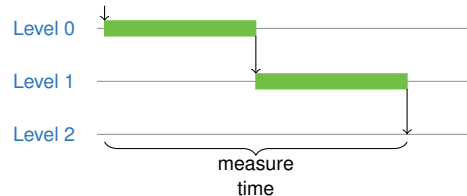
Measurement Modes

Downclocking

■ AVX



(a) Level 1



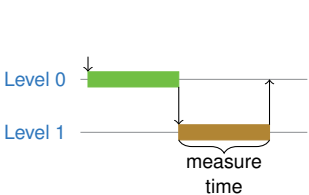
(b) Level 2

Measurement Modes

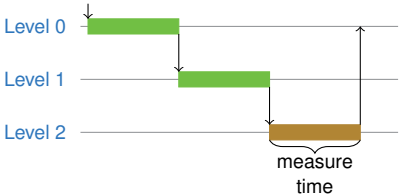
Upclocking

■ AVX

■ Scalar



(a) Level 1



(b) Level 2

Yussuf Khalil

Motivation

Idea

Analysis

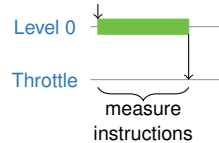
Reimplementation

Evaluation

Future Work

Conclusion

■ AVX

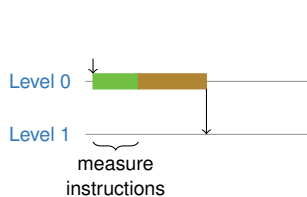


Measurement Modes

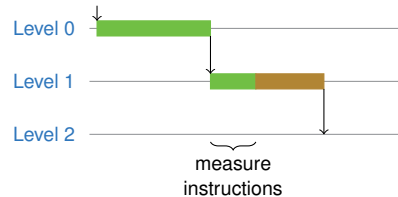
Trigger Instructions

■ AVX

■ Scalar



(a) Level 1



(b) Level 2

Findings

Coarse behavior (Intel Core i9-7940X):

- 1 Throttle out-of-order engine after first AVX instruction
- 2 Switch to frequency level 1, needs $\approx 25 \mu\text{s}$
- 3 Switch to frequency level 2 (only heavy AVX-512 instructions), $\approx 27 \mu\text{s}$
- 4 Return to baseline frequency $\frac{2}{3}$ ms after last AVX instruction

In contrast, Intel claims

- downclocking takes “up to $500 \mu\text{s}$ ”
- upclocking after 2 ms

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

Findings

Coarse behavior (Intel Core i9-7940X):

- 1 Throttle out-of-order engine after first AVX instruction
- 2 Switch to frequency level 1, needs $\approx 25 \mu\text{s}$
- 3 Switch to frequency level 2 (only heavy AVX-512 instructions), $\approx 27 \mu\text{s}$
- 4 Return to baseline frequency $\frac{2}{3}$ ms after last AVX instruction

In contrast, Intel claims

- downclocking takes “up to $500 \mu\text{s}$ ”
- upclocking after 2 ms

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

Yussuf Khalil

Motivation

Idea

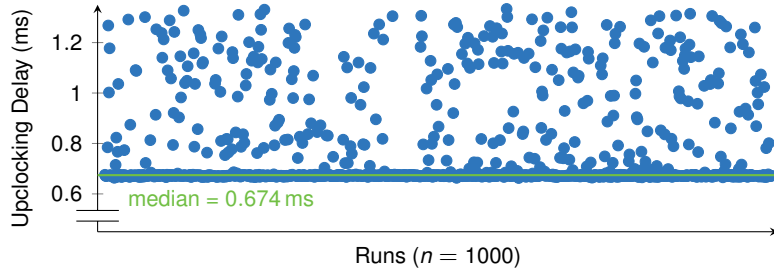
Analysis

Reimplementation

Evaluation

Future Work

Conclusion



- Heavy AVX-512 instructions *sometimes* need up to $\frac{4}{3}$ ms to return from level 1
- We do not have an explanation for this

Yussuf Khalil

Motivation

Idea

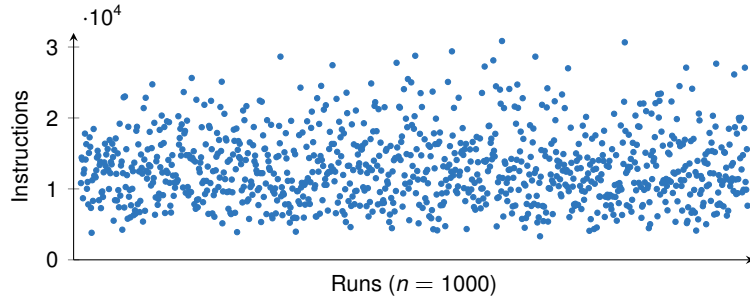
Analysis

Reimplementation

Evaluation

Future Work

Conclusion



- 512-bit instructions always trigger frequency reductions after first instruction
- 256-bit ones exhibit high variance in the amount of required operations (shown above)

Design

Motivation

Idea

Analysis

Reimplementation

Evaluation

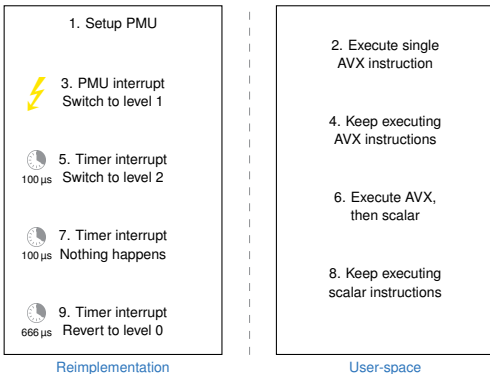
Future Work

Conclusion

- Modify `intel_pstate` Linux driver
 - `intel_pstate` manages core frequency
- Introduce virtual AVX frequency levels similar to hardware levels
- Use performance counters to measure executed AVX instructions
 - **Only available for floating-point**

Algorithm

1. Interrupt after first 512-bit instruction
2. Switch to level 1 frequency
3. Measure instruction throughput in 100 μ s intervals
4. If higher than one per cycle: switch to level 2
5. Raise frequency after 666 μ s with no 512-bit instructions



Oracle Mechanism

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- Want to find theoretical limit for optimized algorithms
 - ⇒ implement oracle mechanism
- System call API to tell reimplementation when to switch frequency levels
- Synthetic workload with scalar and vector phases
 - ⇒ scalar phases of about $\frac{2}{3}$ ms expose worst case
 - hardware's worst case is best case for oracle
 - Count iterations in each phase

Reimplementation

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- Need to evaluate correctness and quality of reimplementation
 - ⇒ use analysis framework, compare with hardware results
- Analysis framework and reimplementation both use performance monitoring
 - ⇒ modify analysis kernel module and reimplementation to work together
- Measure performance overhead

Yussuf Khalil

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- Downclocking to level 1 works as expected
- Reaching level 2 takes $\approx 104 \mu\text{s}$
 - unlike the hardware, we can measure throughput only over time
- Upclocking after $666 \mu\text{s}$ as designed
 - however, up to $100 \mu\text{s}$ more possible due to limited measurement accuracy

Overhead and Oracle Performance

AVX phase (200 ms)

Motivation

Idea

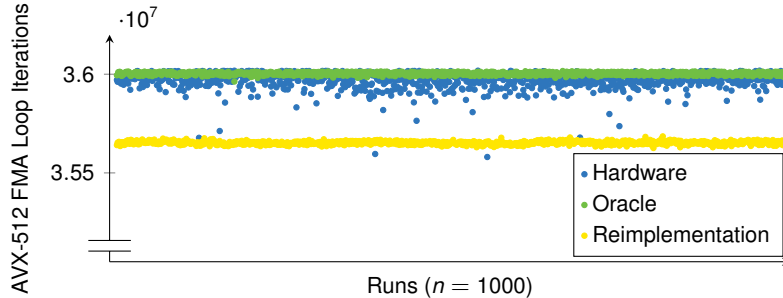
Analysis

Reimplementation

Evaluation

Future Work

Conclusion



- Reimplementation is $\approx 1\%$ slower in AVX phases
- Oracle yields same performance as hardware

Overhead and Oracle Performance

Scalar phase (666 μ s)

Motivation

Idea

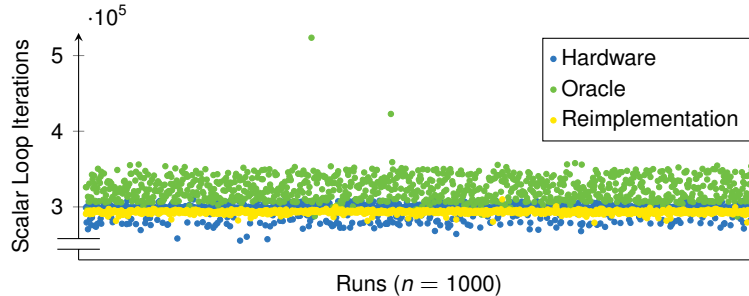
Analysis

Reimplementation

Evaluation

Future Work

Conclusion



- Reimplementation $\approx 2\%$ slower
- Oracle gives $\approx 8\%$ median performance increase (99th percentile: 15%)

Future Work

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- Implement alternative reclocking algorithms
 - Take ideas from power management
- More analysis
 - Other processors (Haswell, Icelake)
 - SMT
 - Instruction mixtures
- Reverse engineer AVX frequency offset configuration

Conclusion

Motivation

Idea

Analysis

Reimplementation

Evaluation

Future Work

Conclusion

- Performance of workloads with scalar and vector phases suffers from AVX reclocking
- We want to find improvements to the reclocking behavior
- Intel's claims about the algorithm are not accurate
- Improvements for heterogeneous workloads are theoretically possible
 - up to 15 % performance increase
- Possibilities for software reimplementation are limited
 - Still a potentially promising approach
 - Improvements likely possible