

Team 58 Presentation

Richard Wohlbold

Aljaz Medic

Aleksa Micanovic

Aleks Stepancic



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Typical Organization I

- **Algorithm that you consider (maybe 2 slides)**
 - State problem that it solves (input:..., output: ...)
 - If possible visualize how it works or show high-level pseudocode
 - State asymptotic runtime
- **Cost analysis (cost measure, exact count)**
- **Baseline implementation (briefly explain), maybe show already performance plot and extract percentage of peak**
- **Optimizations you performed**
 - Briefly discuss major optimizations/code versions
 - Maybe explain the most interesting in a bit greater detail
 - Any analysis (e.g., profiling) you performed is interesting – show the result
 - If too much, explain only some things and just state the rest

Typical Organization II

■ Experimental results

- Very brief: Experimental setup (platform, compiler)
- Performance plot over a range of sizes with different code versions
- Make sure you also push input size to the limit in the experiments
- Extract overall speedup

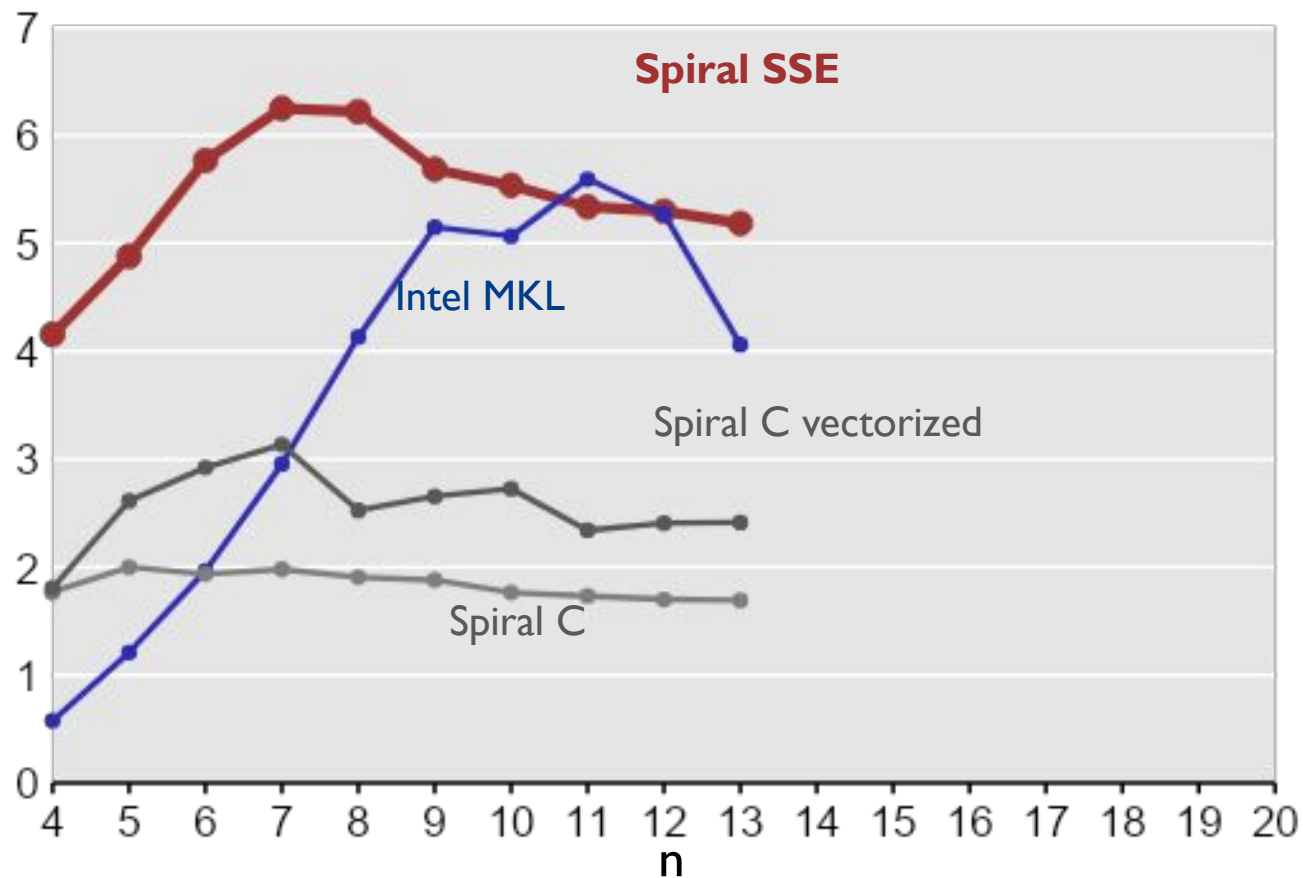
■ Every project is different – so adapt as needed

■ Focus on the most interesting things, don't explain everything that will be in the final report.

Try to Make Nice Plots

DFT 2^n (single precision) on Pentium 4, 2.53 GHz

[Gflop/s]



Prime Implicants for Quine-McCluskey

- The algorithm determines all prime implicants of a boolean function, given the initial set of minterms
- **Input:** number of bits and list of minterms
(e.g. num_bits=7 0010011, 100101, ...)
- **Output:** list of all prime implicants

Prime Implicants for Quine-McCluskey

- Every implicant has a bit that describes it in a 3^n array
- There are 3 such arrays: **implicants**, **merged**, **primes**
- Each is divided into $n+1$ sections made up of chunks
 - For $0 \leq k \leq n$, the k th section contains :
 - *all implicants with k*
 - *$\binom{n}{k}$ many chunks (all implicants have dashes in the same place)*
 - *chunks of size 2^{n-k} bits*
- Chunks within the same section have a combinatorial ordering, based on the location of the dashes
- A single merge step takes one chunk of input and produces multiple chunks of output

$n = 4$ (not drawn to scale)

$k = 0$	$k = 1$	$k = 2$	$k = 3$	$k = 4$
$\binom{4}{0} \cdot 2^4$ = 16 implicants/bits	$\binom{4}{1} \cdot 2^3$ = 32 implicants/bits	$\binom{4}{2} \cdot 2^2$ = 24 implicants/bits	$\binom{4}{3} \cdot 2^1$ = 8 implicants/bits	$\binom{4}{4} \cdot 2^0$ = 1 impl./bits

$n = 4$ (not drawn to scale)

$k = 0$	$k = 1$	$k = 2$	$k = 3$	$k = 4$
$\binom{4}{0} \cdot 2^4$ = 16 implicants/bits	$\binom{4}{1} \cdot 2^3$ = 32 implicants/bits	$\binom{4}{2} \cdot 2^2$ = 24 implicants/bits	$\binom{4}{3} \cdot 2^1$ = 8 implicants/bits	$\binom{4}{4} \cdot 2^0$ = 1 impl./bits

$$k = 1$$
 $k = 2$

implicants	***-	**-*	*-**	-***	**--	*--*	---*	*---	--**	----
------------	------	------	------	------	------	------	------	------	------	------

[illegible][illegible]

$n = 4$ (not drawn to scale)

$k = 0$	$k = 1$	$k = 2$	$k = 3$	$k = 4$
$\underbrace{\hspace{10em}}_{\substack{\binom{4}{0} \cdot 2^4 \\ = 16 \text{ implicants/bits}}} \underbrace{\hspace{10em}}_{\substack{\binom{4}{1} \cdot 2^3 \\ = 32 \text{ implicants/bits}}} \underbrace{\hspace{10em}}_{\substack{\binom{4}{2} \cdot 2^2 \\ = 24 \text{ implicants/bits}}} \underbrace{\hspace{10em}}_{\substack{\binom{4}{3} \cdot 2^1 \\ = 8 \text{ implicants/bits}}} \underbrace{\hspace{10em}}_{\substack{\binom{4}{4} \cdot 2^0 \\ = 1 \text{ impl./bits}}}$				

[illegible]

$n = 4$ (not drawn to scale)

$k = 0$	$k = 1$	$k = 2$	$k = 3$	$k = 4$
$\binom{4}{0} \cdot 2^4$ = 16 implicants/bits	$\binom{4}{1} \cdot 2^3$ = 32 implicants/bits	$\binom{4}{2} \cdot 2^2$ = 24 implicants/bits	$\binom{4}{3} \cdot 2^1$ = 8 implicants/bits	$\binom{4}{4} \cdot 2^0$ = 1 impl./bits

[illegible]

Dummy Baseline “Dense Byte”

- Two phase approach
 1. **Generation (Merge):** AND pair of implicants that differ by one bit
 2. **Scan:** Find implicants that are not merged
- **Drawback:**

“Dense byte” stores **one implicant per byte**

Optimization #1 “64bit vectorization”

Improvement:

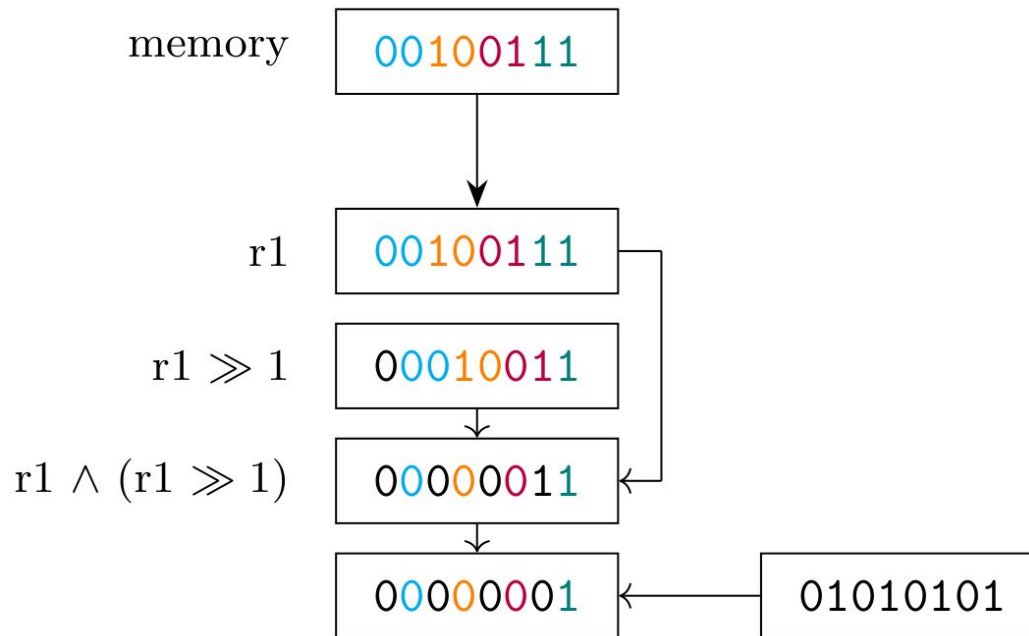
Store 64 implicants in **64 bit registers**

Block sizes ≥ 64 , merging **EASY**

Problem:

Block sizes < 64 , need to find approach to extract the merged bits within the register.

$n \geq 3$, block size: 1 bit, register width: 8 bits



How to reorganize bits into lower half, i.e. turn 01010101 into 00001111 ?

```
unsigned __int64 _pext_u64 (unsigned __int64 a, unsigned __int64 mask)
```

Synopsis

```
unsigned __int64 _pext_u64 (unsigned __int64 a, unsigned __int64 mask)
#include <immintrin.h>
Instruction: pext r64, r64, r64
CPUID Flags: BMI2
```

Description

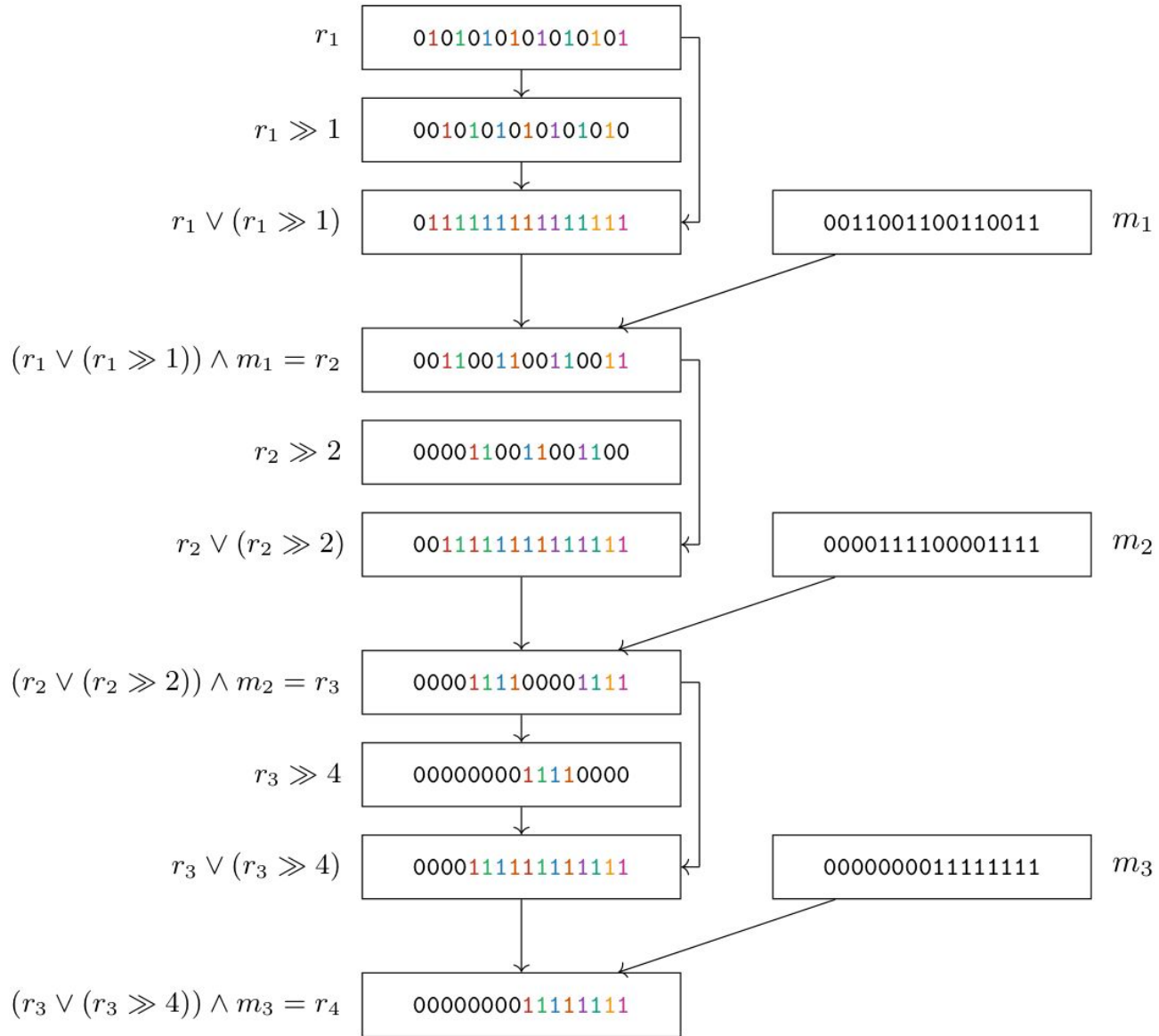
Extract bits from unsigned 64-bit integer `a` at the corresponding bit locations specified by `mask` to contiguous low bits in `dst`; the remaining upper bits in `dst` are set to zero.

Operation

```
tmp := a
dst := 0
m := 0
k := 0
DO WHILE m < 64
    IF mask[m] == 1
        dst[k] := tmp[m]
        k := k + 1
    FI
    m := m + 1
OD
```

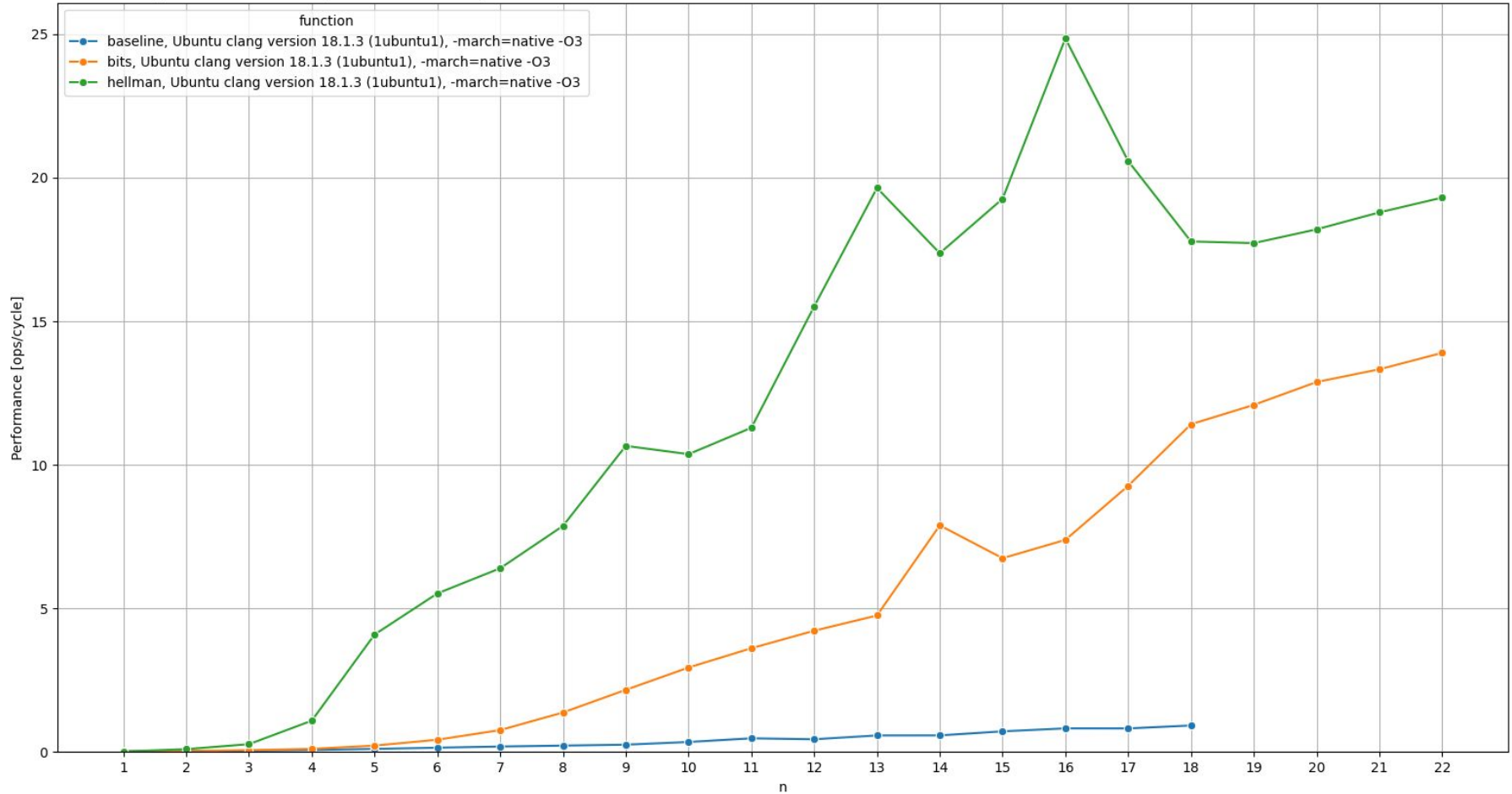
Latency and Throughput

Architecture	Latency	Throughput (CPI)
Alderlake	3	1
Icelake Intel Core	3	1
Icelake Xeon	3	1
Sapphire Rapids	3	1
Skylake	3	1



PLOT #1 (BYTE, BITS, HELMAN)

Performance of all implementations for different number of bits n
CPU: AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics



Optimization #2 (single pass)

- **Additional space (3^n bits)** for a merge bitmap
- **Additional time and cache misses** for second prime-marking traversal through **implicants, merge, prime** bitmaps
- **Improvement:** in the first pass, populate **prime** bitmap

```
merged = implicants[i] AND implicants[i + (1 << k)]
```

```
primes[i] = primes[i] AND NOT merged
```

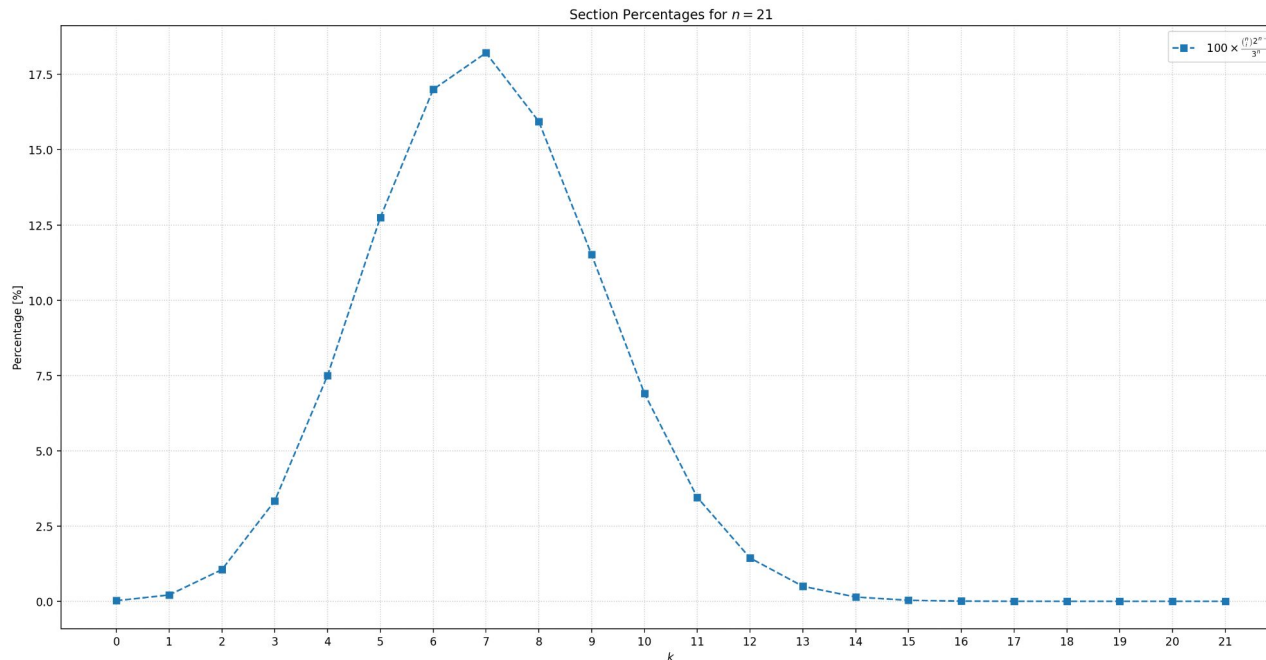

Cost Analysis

$$\text{\#Bit-Ops} = n3^n$$

$$Q(n) \geq \frac{2 \cdot 3^n}{8} \text{ bytes}$$

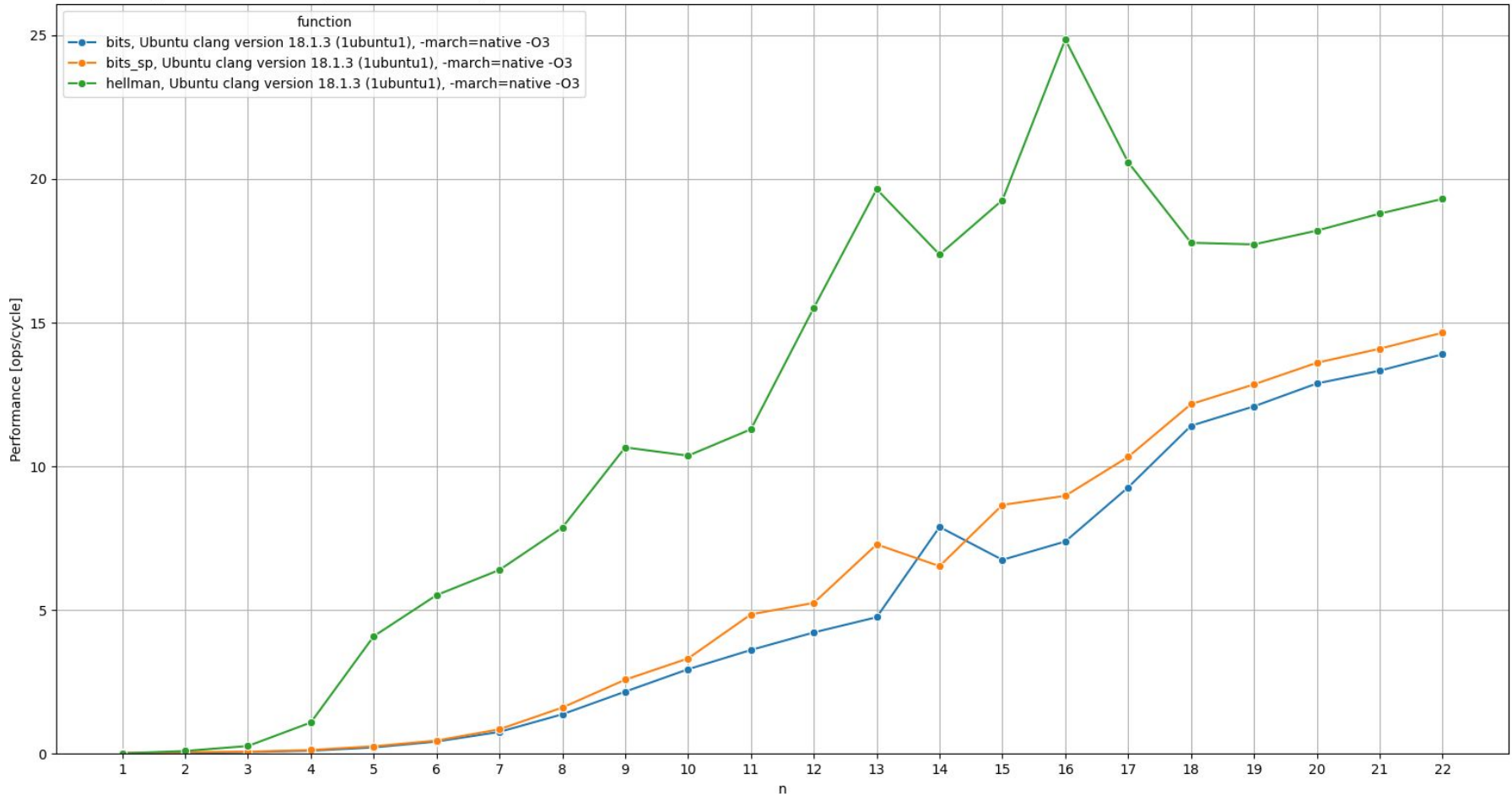
$$I(n) \leq n \cdot 4 \cdot 3^n$$

Biggest section for $k = \left\lfloor \frac{n+1}{3} \right\rfloor$



PLOT #2 (BITS, BITS_SP, HELLMAN)

Performance of all implementations for different number of bits n
CPU: AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics

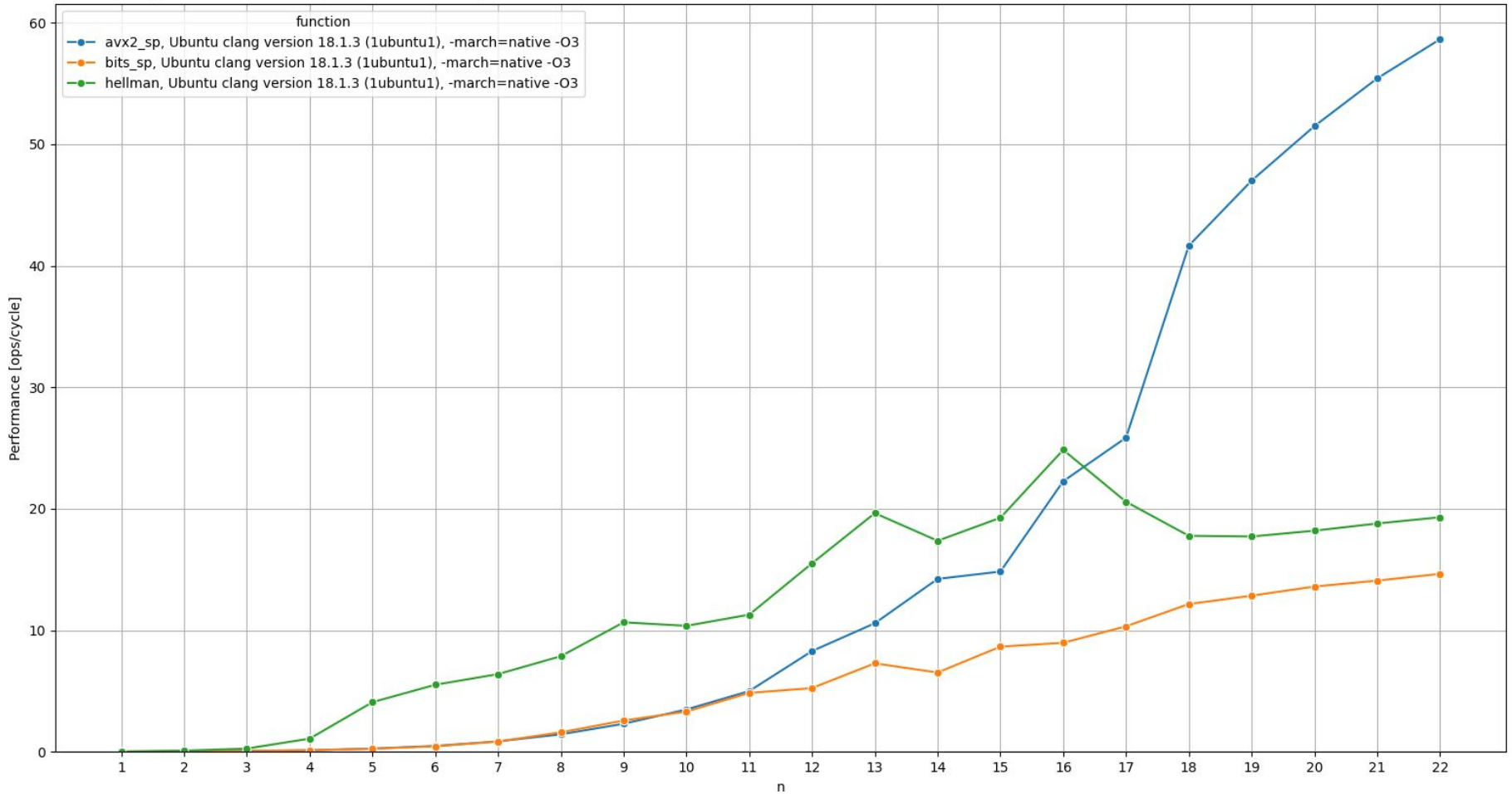


Optimization #3 (AVX2 Vectorization)

- Using 64-bit registers
- **Improvement:** use 256-bit registers

PLOT #3 (bits_sp, avx2_sp, Hellman)

Performance of all implementations for different number of bits n
CPU: AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics



N=22 $\text{Speedup}_{\text{Hellman}} = 3.9x$ $\text{Speedup}_{\text{bits_sp}} = 3.1x$

Optimization #4 (SSA+ILP, Unroll)

- **Convolutd IF branches for inter-register merge**
 - No ILP
- **Improvement 1:** write straightforward code - add ILP

In **merge step** register will be **loaded multiple times** for each block length

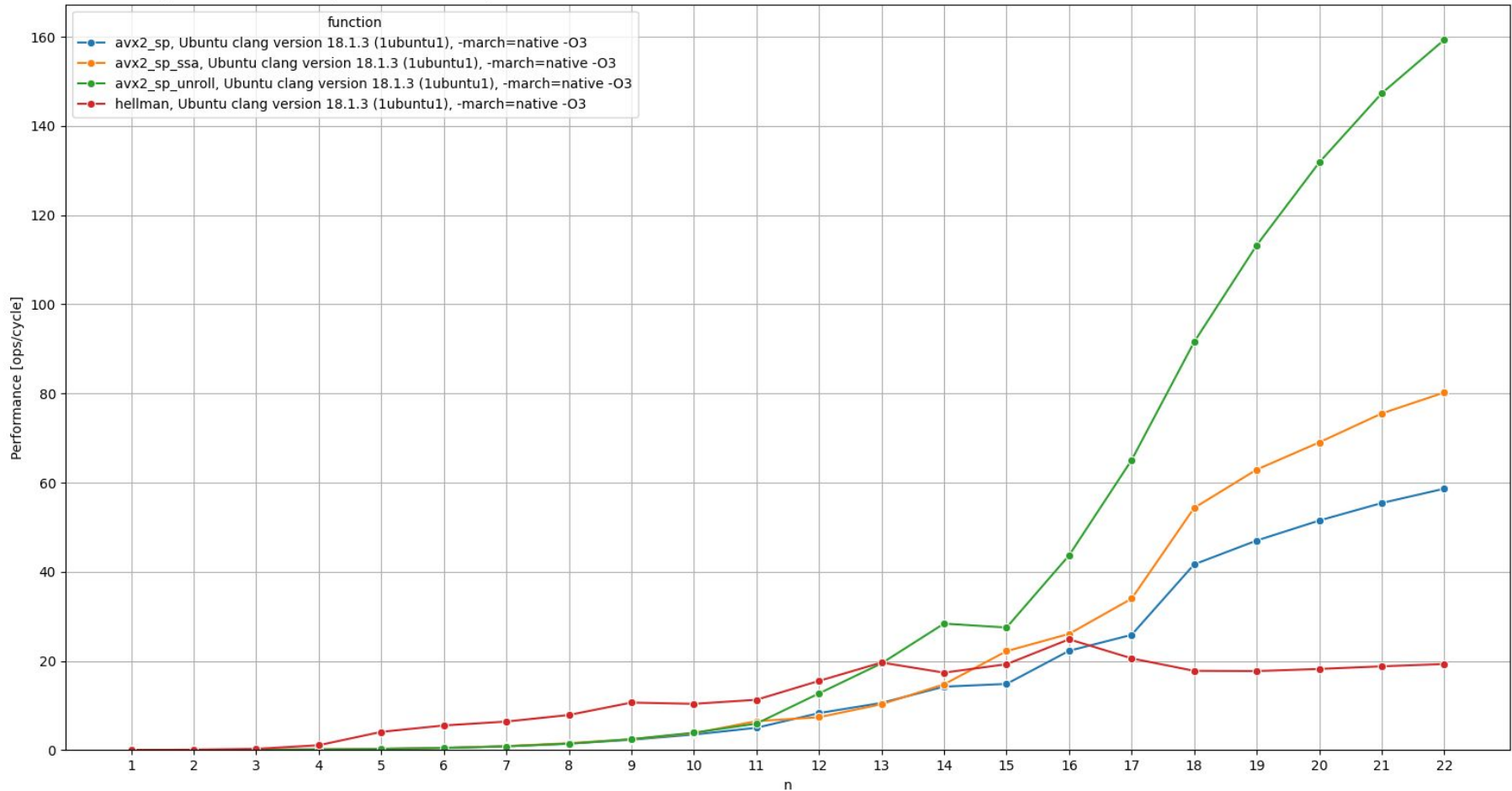
$$\#register_loads = k \frac{2^k}{register_size} = k2^{k-8}$$

- **Suboptimal traversal order of inter-register merge in a chunk**
 - Bad temporal locality - multiple reads of same input bits
 - For large chunks can evict the whole cache
- **Improvement 2:** Inter-register merge processes all **possible block lengths** for a register before moving further + unroll loop for 4 inter-register steps

$$\#register_loads = (k - 7)2^{k-8}$$

PLOT #4 (avx2_sp, avx2_sp_ssa, avx2_sp_unroll, Hellman)

Performance of all implementations for different number of bits n
CPU: AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics



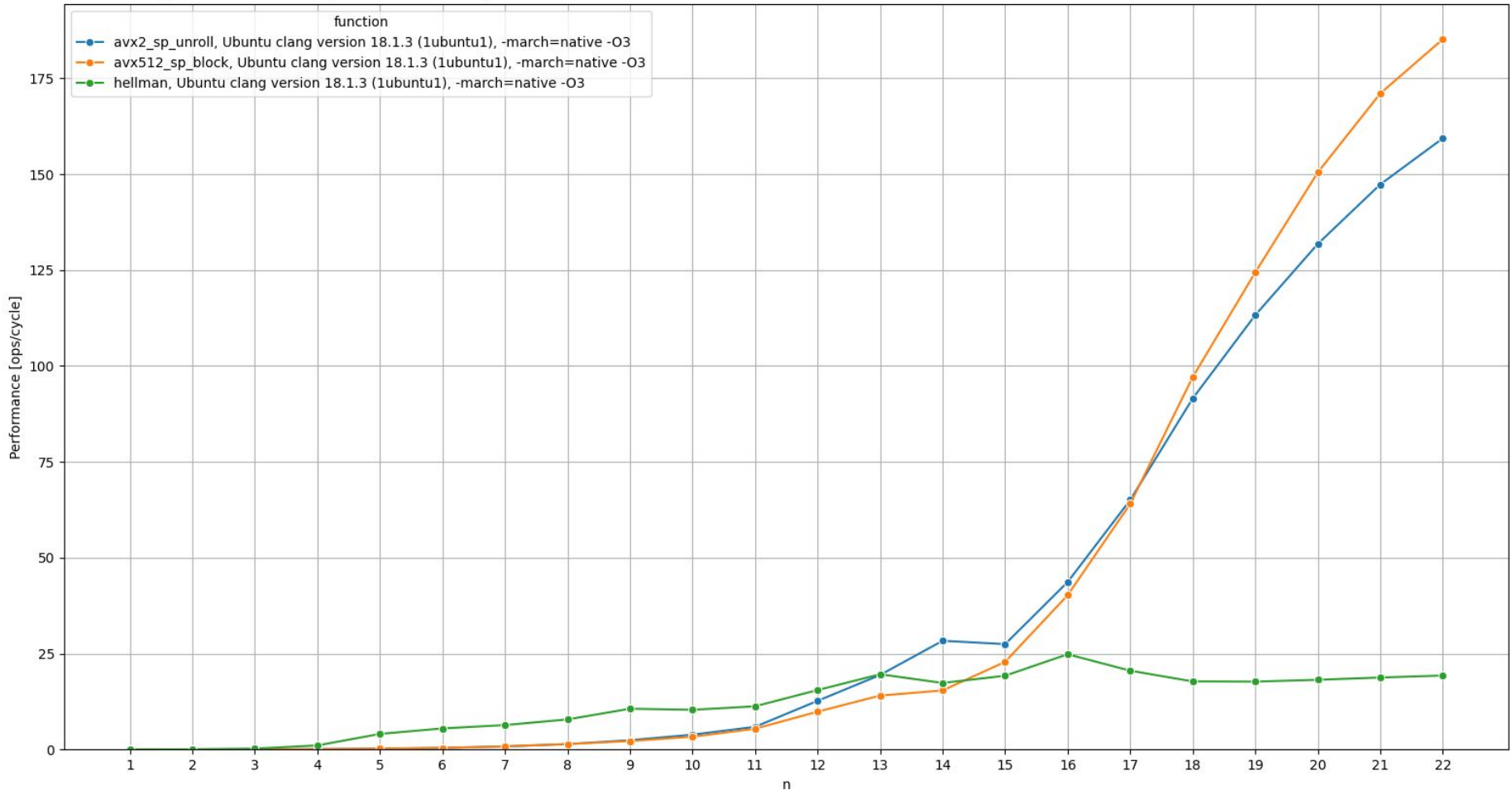
N=22 $\text{Speedup}_{\text{Hellman}} = 8.0x$ $\text{Speedup}_{\text{avx2_sp}} = 2.7x$

Optimization #5 (AVX512)

- Using 256-bit registers
- **Improvement:** use 512-bit registers

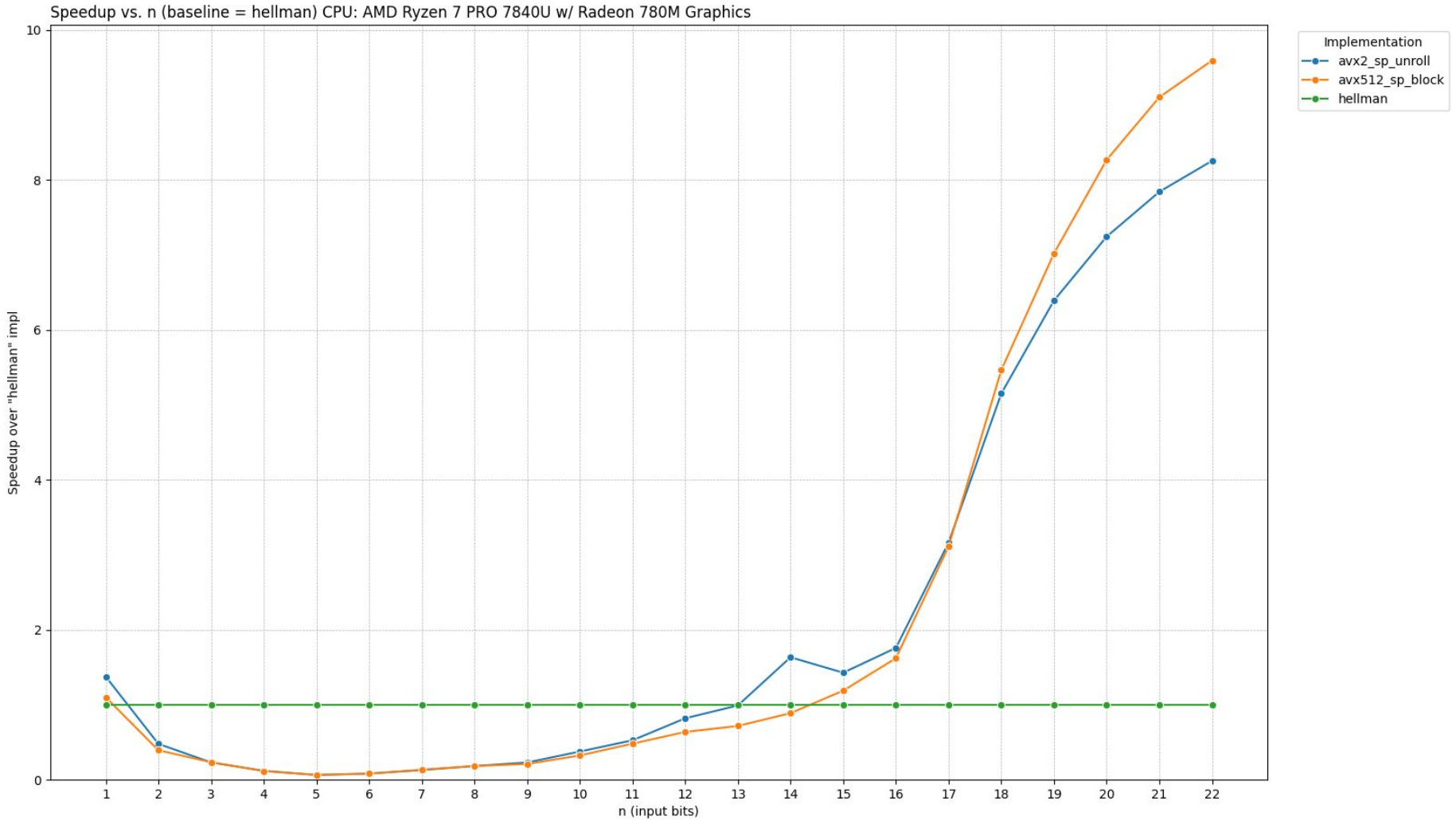
PLOT #5 (avx2_sp_unroll, avx512, Hellman)

Performance of all implementations for different number of bits n
CPU: AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics



N=22 $\text{Speedup}_{\text{Hellman}} = 9.3x$ $\text{Speedup}_{\text{avx2_sp_unroll}} = 1.1x$

SPEEDUP (avx2_sp_unroll & avx512 vs Hellman)



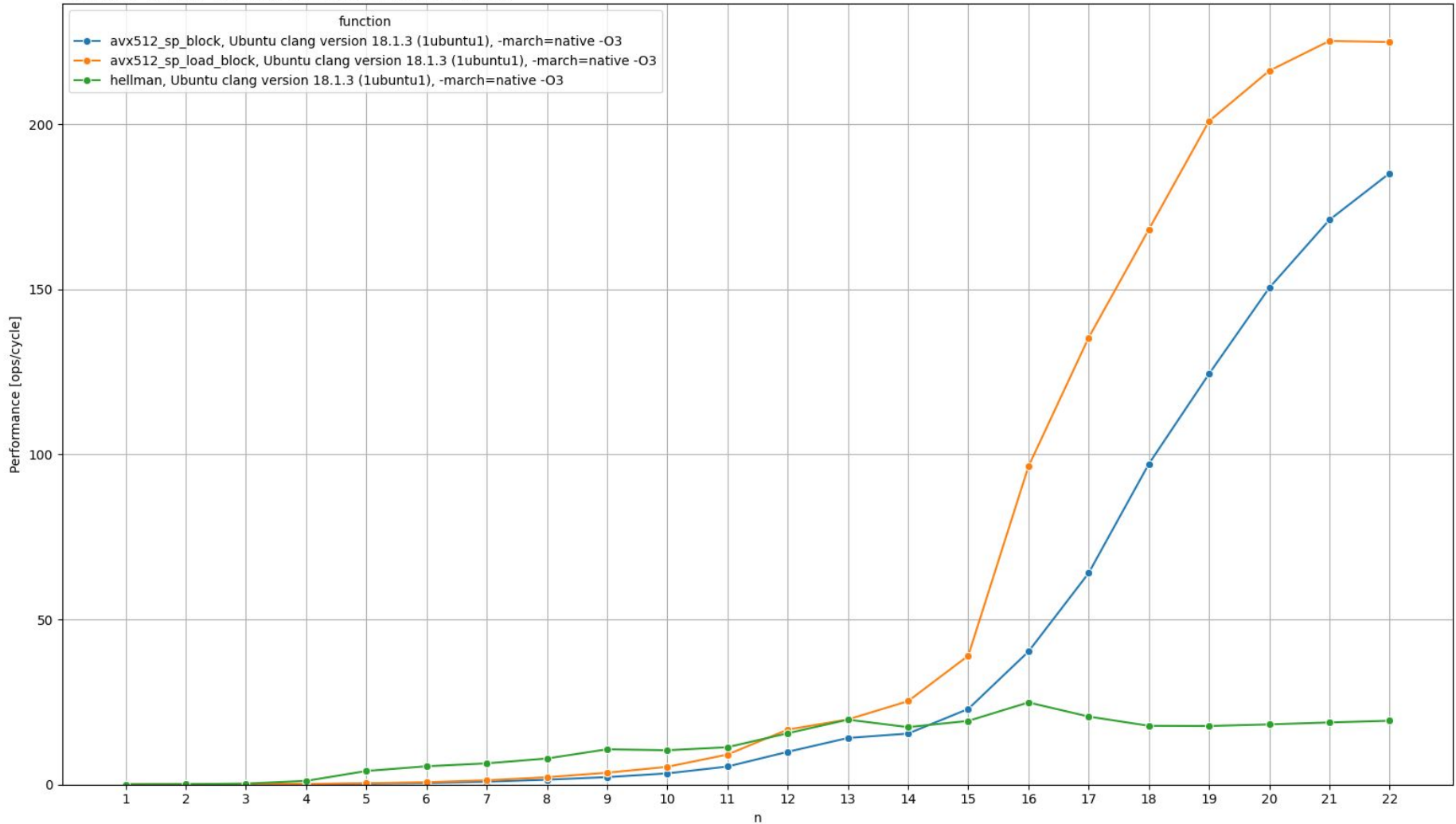
Final Optimization (loading traversal order)

- The traversal order of implicants is being calculated, which adds cycles to the execution
- **Improvement:** precompute the traversal order
- Load **additional** $10 * 2^n$ bytes
 - For $n = 20, 21, 22$ **negligible compared to bitmaps (< 1% of storage)**

```
// one linear pass
for (size_t i = 0; i < op_count; i++) {
    MergeOp *op = &ops[i];
    MERGE_FUNCTION(
        implicants,
        primes,
        op->in_idx,
        op->out_idx,
        op->rem_bits,
        op->first_diff
    );
}
```

FINAL PLOT #5 (avx512, avx512_load, Hellman)

Performance of all implementations for different number of bits n
CPU: AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics



SPEEDUP (avx512 & avx512_load vs Hellman)

