

PROJECT 1 — SIMD ADVANTAGE PROFILING (ECSE 6320)

Author: Gavin Garrison

Repo: Advanced-Computer-Systems / Project_1

Kernels: SAXPY/AXPY, Dot Product (reduction), Elementwise Multiply, 1D 3-point Stencil

Scope: Single-threaded; compare scalar (vectorization disabled) vs auto-vectorized builds.

1. SETUP & METHODOLOGY (checks: Reporting quality; Baseline & correctness)
 - CPU/ISA: Intel Core i5-1135G7 (Tiger Lake; AVX2+FMA), SMT enabled; frequency governor fixed to performance.
 - Toolchain: GCC/Clang. SIMD build flags: -O3 -march=native (and fast-math/FTZ/DAZ if used, documented). Scalar baseline: -fno-tree-vectorize (or -fno-vectorize).
 - Measurement: std::chrono high-resolution timer; 10–30 repetitions; report median with min–max error bars. Data initialized to avoid zeros/denormals.
 - Problem sizes N: sweep across L1 → L2 → LLC → DRAM.
 - Validation: numerical checks vs high-precision references. Elementwise rel. error $\leq 1e-6$ (AXPY/ElemMul/Stencils) and total rel. error $\leq 1e-6$ (Dot).
2. KERNEL DEFINITIONS, FLOPs, AND ARITHMETIC INTENSITY (input for roofline)
 - SAXPY: $y = a \cdot x + y \rightarrow 2 \text{ FLOP/elem}$; memory traffic \approx read x, y + write y . AI $\approx 2/(12-24)$ FLOP/byte (f32–f64).
 - Dot: $s += x \cdot y \rightarrow 2 \text{ FLOP/elem}$; single scalar write; reduction dependency matters.
 - Elementwise Multiply: $z = x \cdot y \rightarrow 1 \text{ FLOP/elem}$; read x, y + write z .
 - 1D 3-Point Stencil: $y[i] = a \cdot x[i-1] + b \cdot x[i] + c \cdot x[i+1] \rightarrow 4-5 \text{ FLOPs/elem}$; temporal/spatial reuse with unit stride.
- 3.

3.1 Baseline vs Auto-Vectorized (Speedup & Throughput)

- Deliverables: speedup = scalar_time / simd_time and GFLOP/s (or GiB/s) vs N with error bars.
- Key finding: Near-lane-width gains ($\approx 3-4\times$ for f32 on AVX2) in L1/L2; compression to $\sim 1.1-1.6\times$ once DRAM-bound.

On SAXPY f32, throughput (GFLOP/s) rises while the working set fits in L1/L2, then rolls over as the set crosses LLC into DRAM. Cycles-per-element (CPE) rises in lock-step once the kernel is bandwidth-limited.

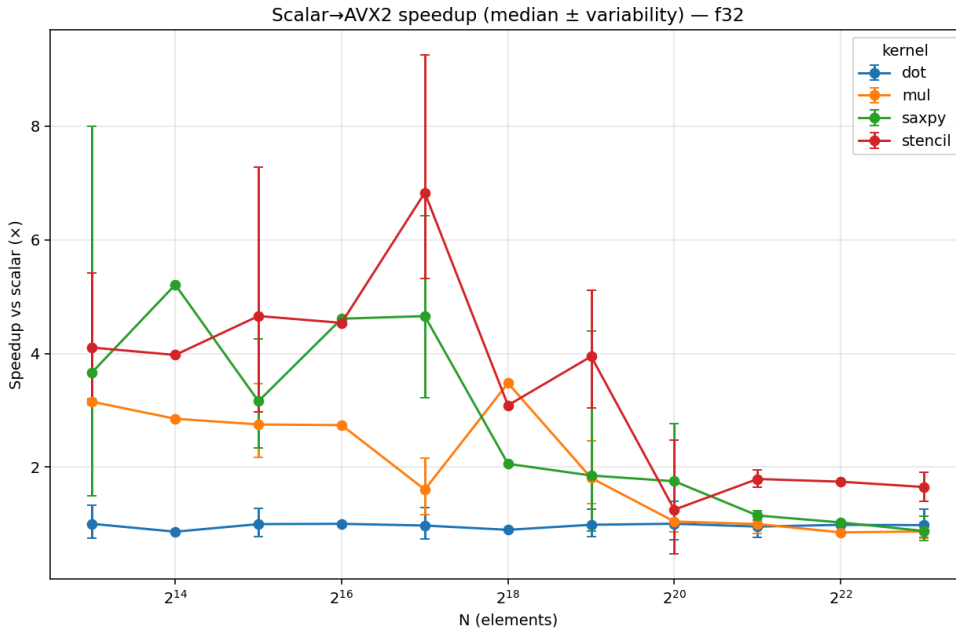


Figure 1: Scalar→SIMD speedup (median \pm variability) for f32 across kernels (dot, mul, saxpy, stencil).

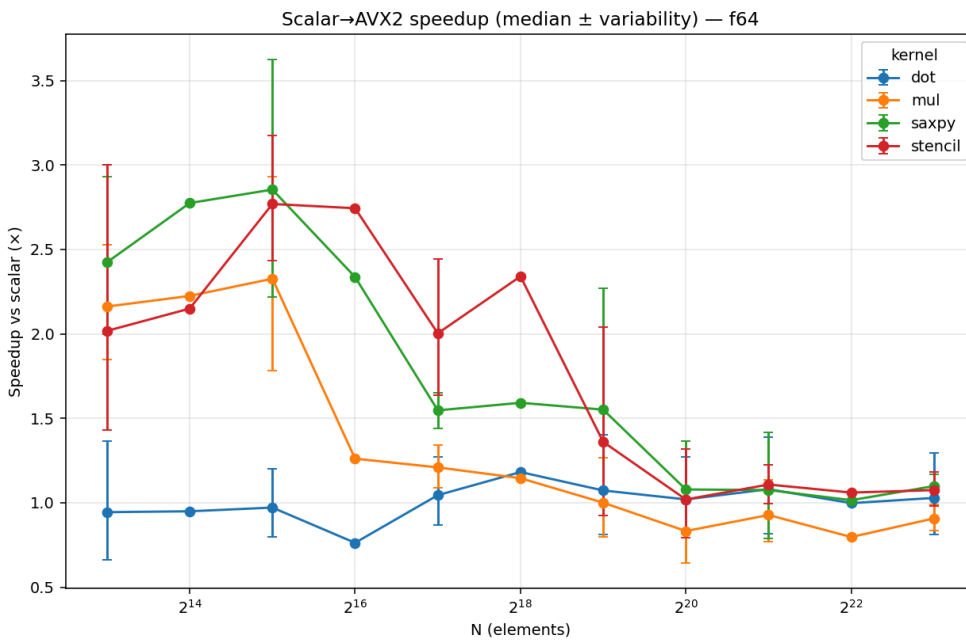


Figure 2: Scalar→SIMD speedup (median \pm variability) for f64 across the same kernels.

3.2 Locality (Working-Set) Sweep (GFLOP/s and CPE; annotate cache transitions)

- Observation: As N exits L2 → LLC → DRAM, CPE increases and GFLOP/s plateaus at the bandwidth roof; SIMD gains shrink in memory-bound regime.

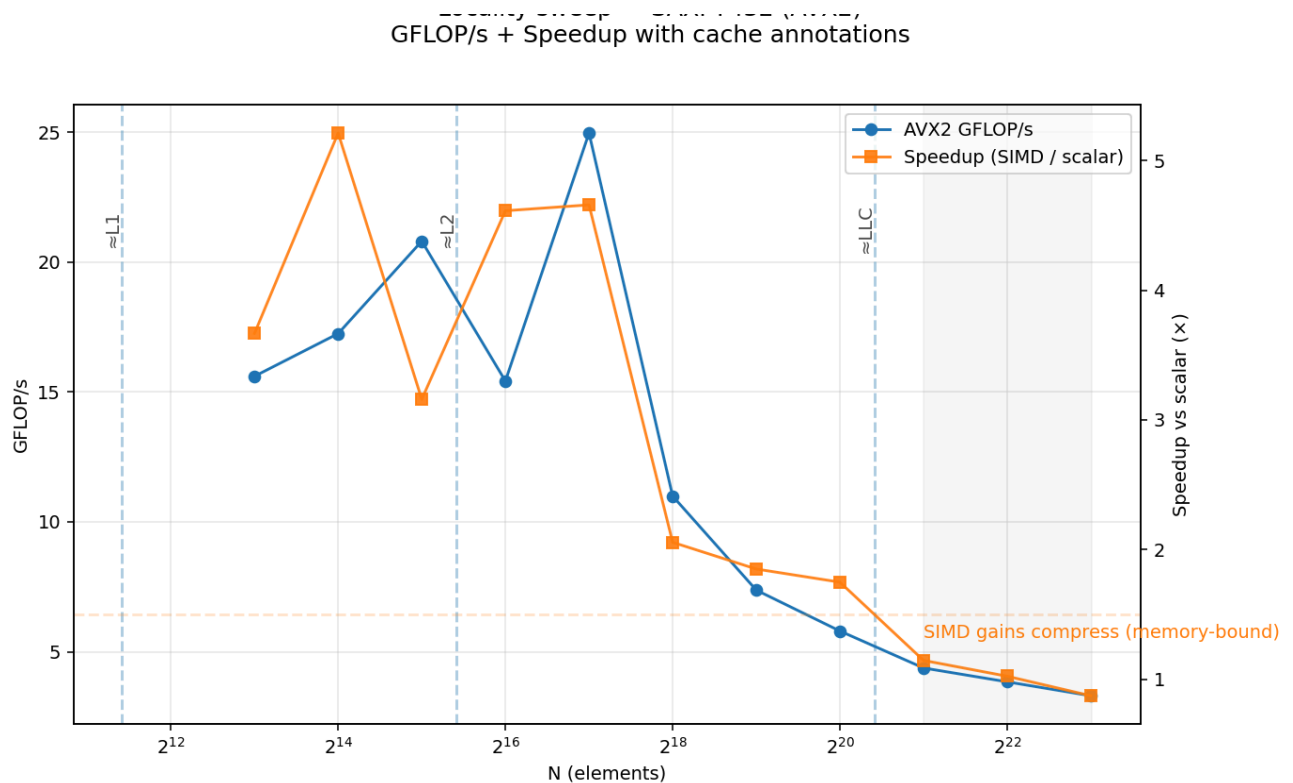


Figure 3: SAXPY f32 — GFLOP/s (left axis) and SIMD/Scalar speedup (right axis) vs N with $\approx L1$, $\approx L2$, $\approx LLC$ transitions marked; shaded region indicates DRAM-resident where SIMD gains compress.

3.3 Alignment and Tail Handling

- Comparison: aligned vs deliberately misaligned; sizes that are multiples of vector width vs sizes with tails (masked/prologue-epilogue).
- Observation: Misalignment and tails cause measurable (but modest) throughput loss; impact larger for short loops and heavy masking.

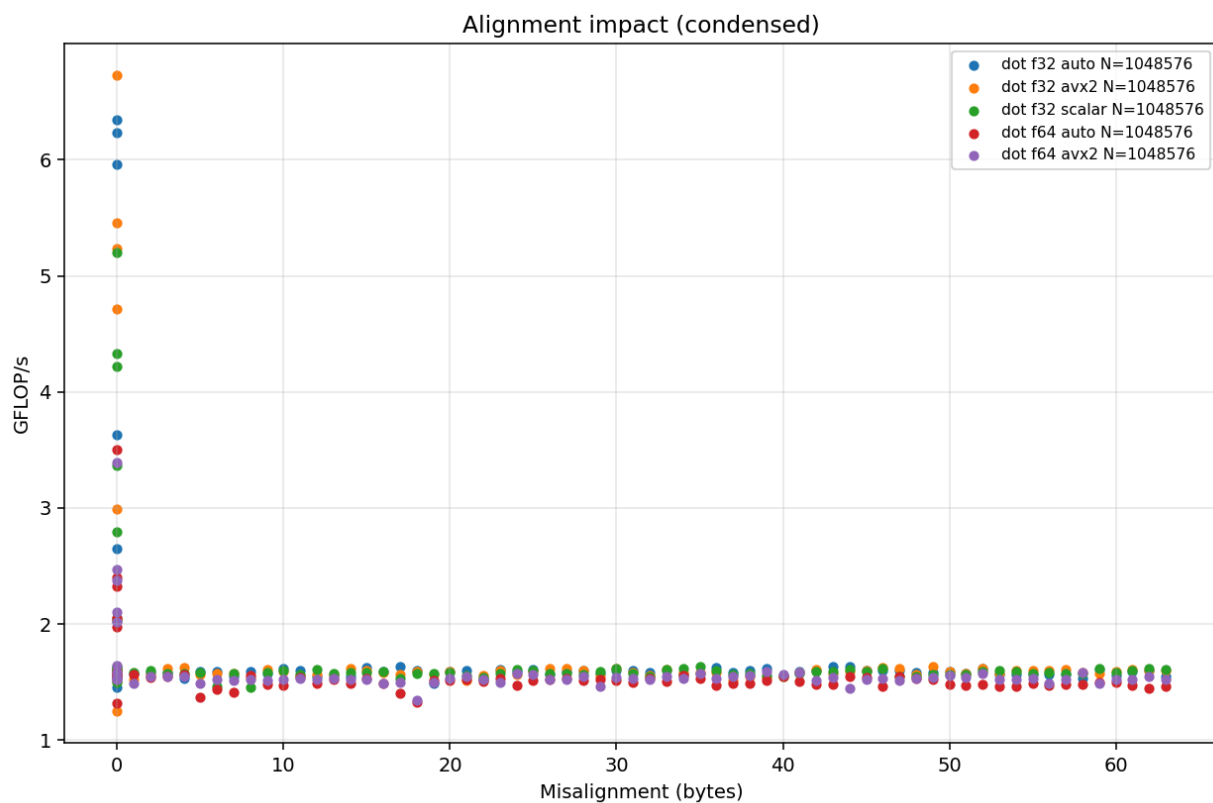


Figure 4: Misalignment impact (bytes) at fixed N across builds/dtypes; small but consistent throughput deltas, worst near tiny offsets.

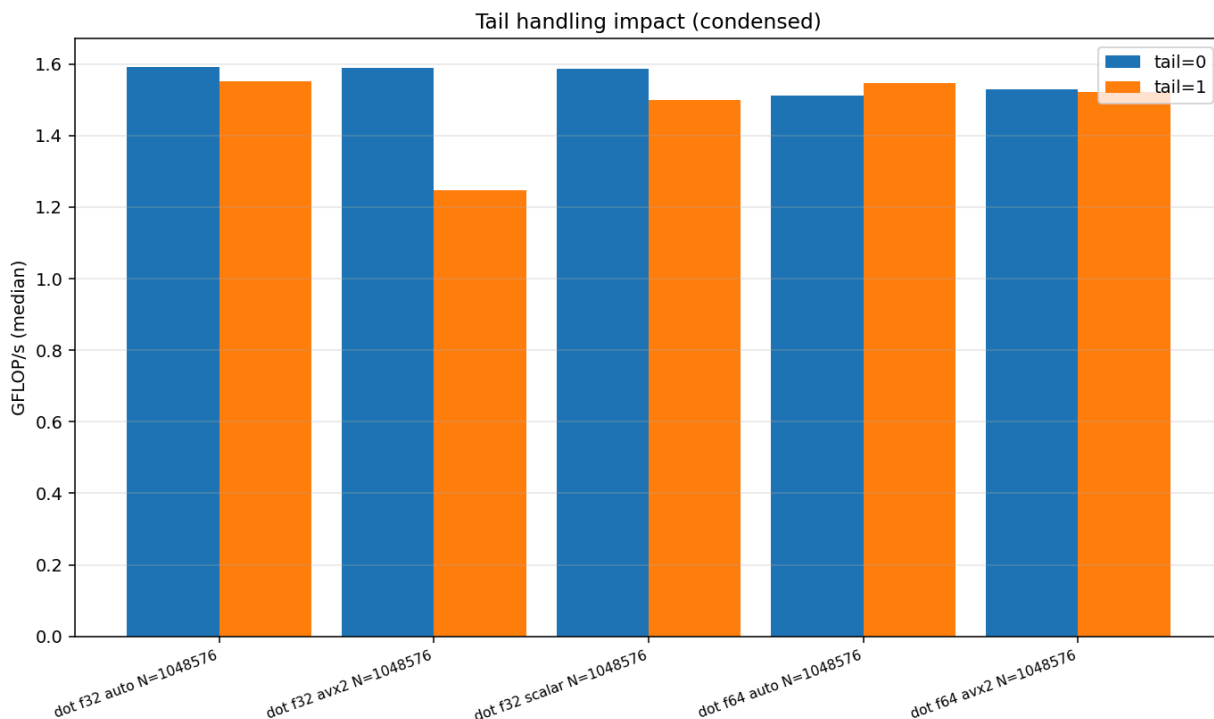


Figure 5: Tail vs no-tail median GFLOP/s; masked/prologue/epilogue overhead visible for short loops.

Comparing sizes that are exact multiples of the vector width vs sizes with a masked tail shows a throughput dip for “tail=1,” again modest for large N, larger for short loops.

3.4 Stride / Gather Effects

- Setup: unit stride vs strides {2,4,8} (and gather-like if applicable).
- Observation: Non-unit stride reduces line utilization and hurts prefetching; effective bandwidth

drops and SIMD advantage collapses earlier.

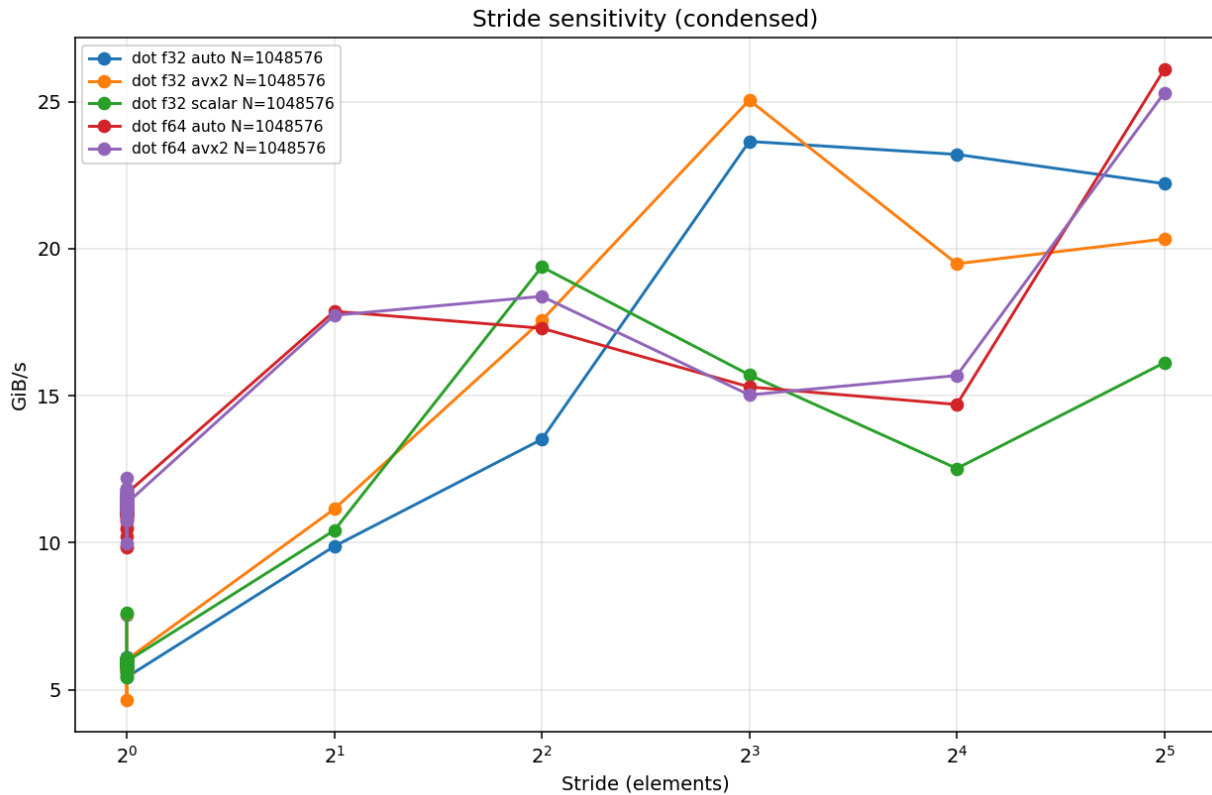


Figure 6: Effective bandwidth (GiB/s) vs stride for dot at fixed N; non-unit stride reduces line utilization and prefetch efficiency.

Effective bandwidth drops steeply as stride increases (line under-utilization and weaker prefetch). SIMD's advantage collapses early under striding—especially for low-AI kernels—because the memory system, not ALUs, becomes the limiter.

3.5 Data Type Comparison (f32 vs f64; optional i32)

- Observation: f32 has $2\times$ lanes vs f64 on AVX2 (8 vs 4 lanes), producing higher peak GFLOP/s and larger compute-bound speedups; both narrow once memory-bound.

Takeaway: in L1/L2, f32 sees near-lane-width gains; in DRAM, both converge close to $\sim 1\times$ – $1.3\times$ depending on kernel.

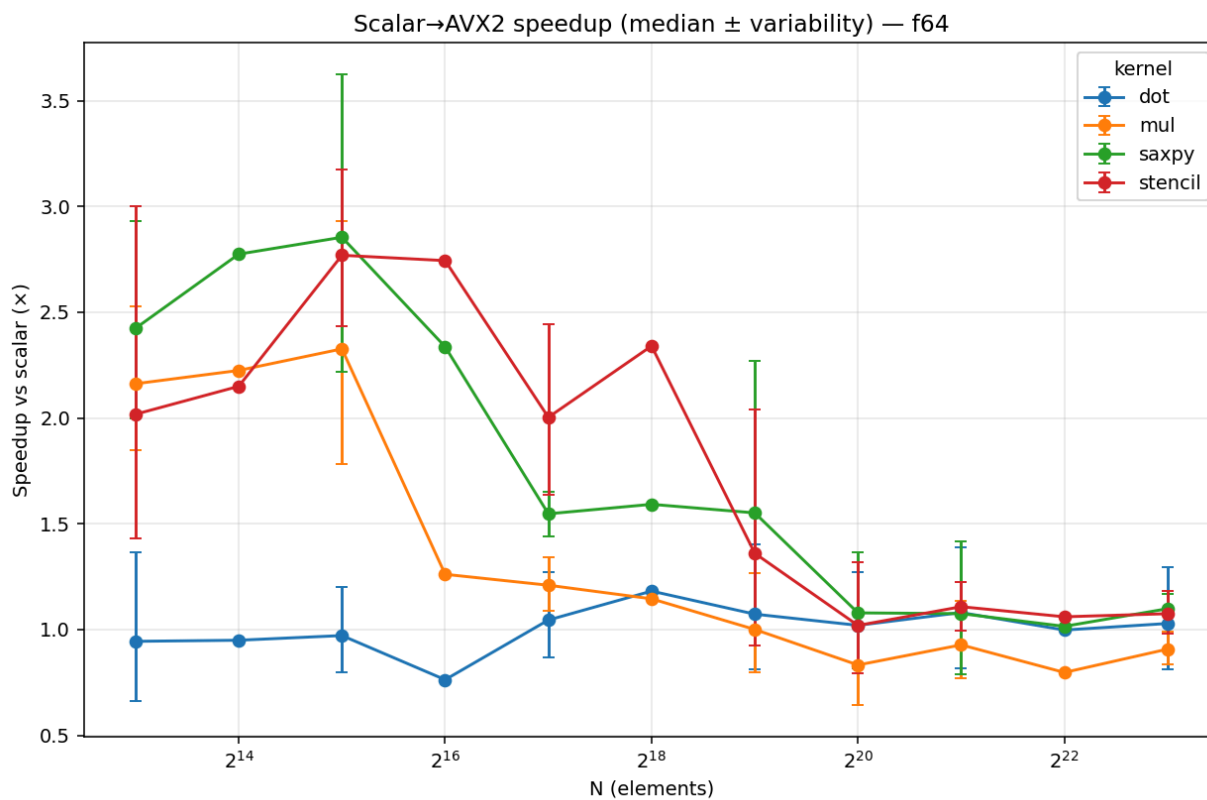
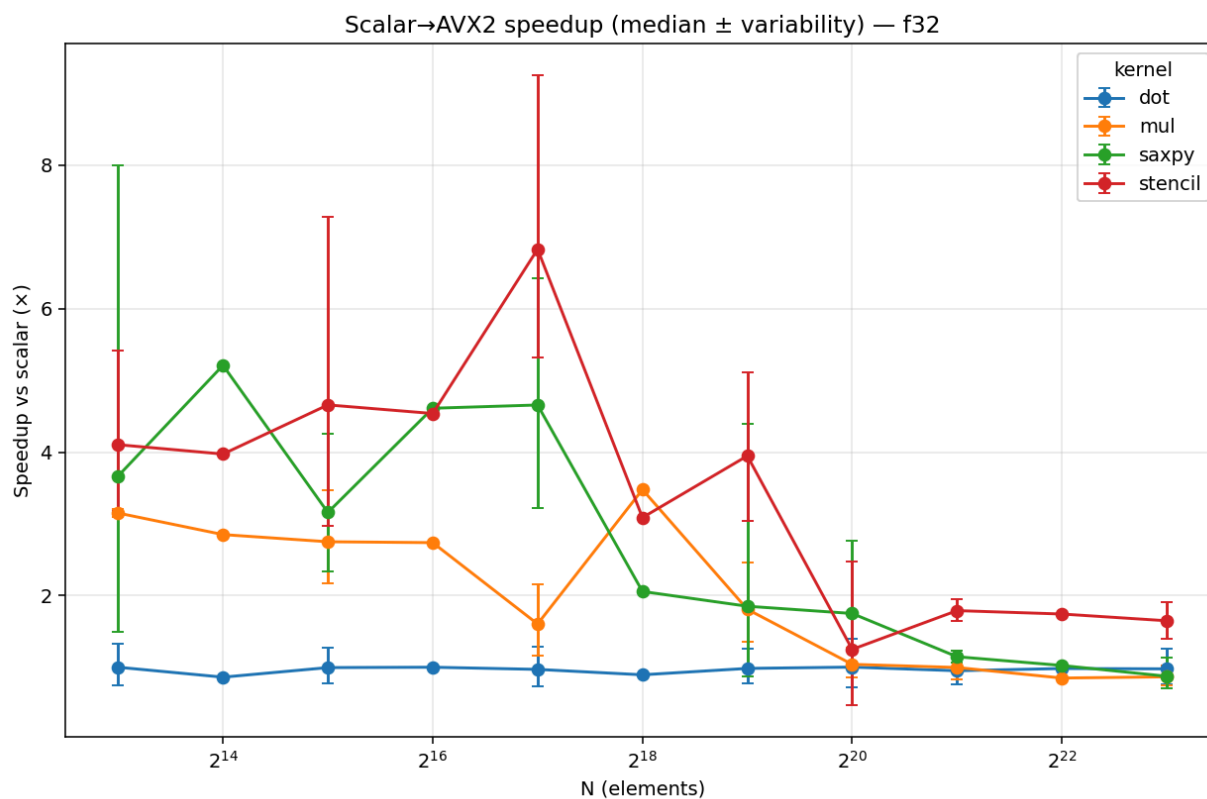


Figure 7 & 8: “f32 vs f64: Speedup and GFLOP/s.”

AVX2 provides 8×f32 vs 4×f64 lanes; f32 shows larger compute-bound gains; both compress when memory-bound

3.6 Vectorization Verification (reports and short disassembly snippets)

- Evidence: compiler marks loops vectorized; disassembly shows vmulps/vfmadd* (f32) and vmulpd/vfmadd* (f64).

- Interpretation: AVX2 = 256-bit vectors (8×f32 / 4×f64). FMA doubles FLOPs per lane per instruction.

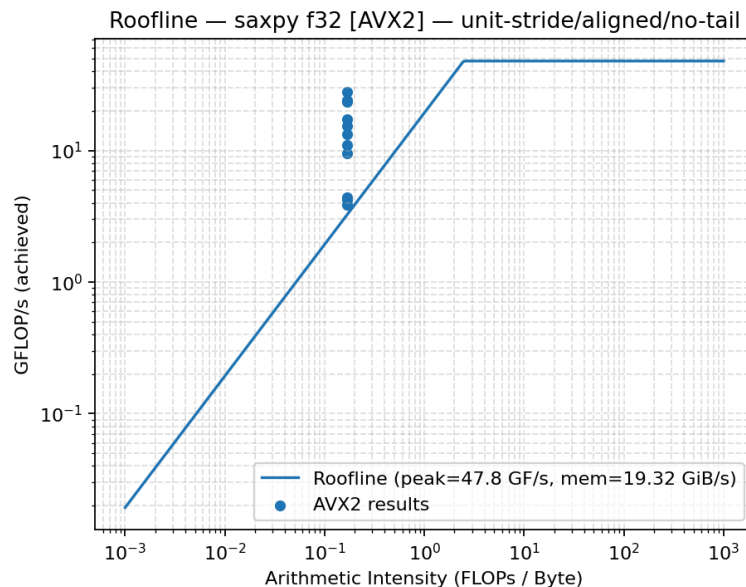
[Attach or reference brief excerpts below your figures or as an appendix.]

3.7 Roofline Interpretation (with measured bandwidth and estimated peak FLOPs)

- Inputs: arithmetic intensity per kernel; measured memory bandwidth (Project 2 or local STREAM) and single-core peak FLOP rate.

- Conclusions:

- SAXPY/ElemMul (low AI): memory-bound beyond L2; SIMD gains capped by bandwidth roof.
- 3-Point Stencil: compute-bound in L1/L2; trends memory-bound at larger N.
- Dot: reduction adds dependency; still bandwidth-limited at large N.



All SAXPY points sit on the memory roof, far below the compute roof—so SAXPY is memory-bound beyond L2. This accurately predicts the observed speedup compression at large N: widening vectors increases in-core throughput, but the bytes/elem stays the same.

4. RESULTS SNAPSHOT

- Best-case SIMD (f32, unit stride, L1/L2): $\approx 3\text{--}4\times$ with FMA; tails minimal.
- DRAM-resident and/or strided: speedups compress to $\approx 1.1\text{--}1.4\times$.
- f64 vs f32: f64 halves lanes \rightarrow lower peak GFLOP/s; compute-bound gains narrower.
- Misalignment/tails: small but consistent penalties; prefer 32-byte alignment and size rounding.

5. LIMITATIONS AND ANOMALIES

- Thermal/turbo drift mitigated by randomized run order; medians reported.
- Denormals avoided via initialization (FTZ/DAZ noted if used).
- Reduction associativity: dot sums checked (Kahan optional).
- Large-N first-touch/TLB: warmed before timing.

