

Fast numerics, small changes

The minimal playbook

Marco Barbone

Flatiron Institute

October 17, 2025



<https://github.com/DiamonDinoia/fwam2025>

~~How to shorten coffee breaks~~

The minimal playbook

Marco Barbone
Flatiron Institute

October 17, 2025

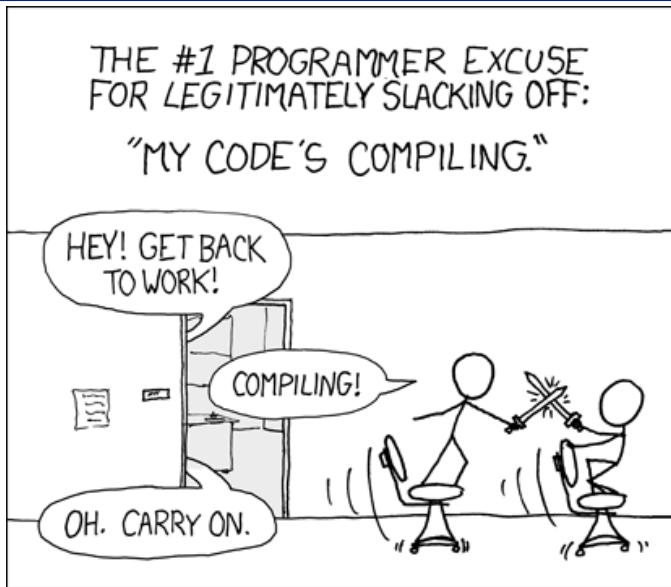


THE #1 PROGRAMMER EXCUSE
FOR LEGITIMATELY SLACKING OFF:
"MY CODE'S COMPILING."

HEY! GET BACK
TO WORK!

COMPILING!

OH. CARRY ON.



Single Instruction, Multiple Data: one instruction applies to multiple elements in parallel.

Scalar vs SIMD

Scalar

```
for i in 0:4:  
    ci = ai + bi
```

Vector

a ₀	a ₁	a ₂	a ₃
b ₀	b ₁	b ₂	b ₃
c ₀	c ₁	c ₂	c ₃

+
=

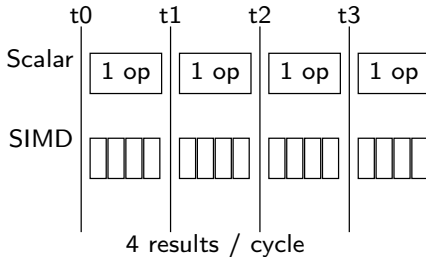
Vector width

$$128b = 2 \times 64$$

$$256b = 4 \times 64$$

$$512b = 8 \times 64$$

Throughput intuition



- **Independent operations** Each iteration or element must be computed without depending on previous results. *Example:* elementwise addition, scaling, normalization.
- **Uniform control flow** All lanes execute the same instruction sequence. Branching (*if/else*) should be minimal or identical across elements.
- **Aligned and contiguous memory access** Data should be stored in consecutive memory addresses for efficient loading/storing. Use arrays or structs-of-arrays layout.
- **Sufficient data parallelism** The problem size should be larger than or equal to the SIMD width (e.g., at least 4, 8, or 16 elements depending on 128b/256b/512b).
- **Identical operation type** All lanes perform the same arithmetic or logic operation (e.g., all do addition or multiplication).
- **No data hazards** No read-after-write (RAW), write-after-read (WAR), or write-after-write (WAW) conflicts within the vectorized section.

- **High arithmetic intensity** Computation per memory access is large enough to hide load/store overhead.
- **Large data sets** Loops with many iterations or large arrays where setup cost is amortized.
- **Regular memory access patterns** Contiguous or strided access that fits SIMD load/store instructions efficiently.
- **Tight inner loops** Performance-critical sections executed millions of times—SIMD pays off most here.
- **Uniform operations across data** Same operation repeated (e.g., vector addition, dot products, FFT kernels).
- **Hardware support available** Wider vector units (e.g., AVX2/AVX-512) and multiple execution ports increase benefit.

Goal: Transform scalar loops into SIMD code

We will use C++ and xsimd.

Step 0: Choose the right algorithm

Chebyshev polynomials of the first kind, $T_n(x)$

Recurrence relations for $T_n(x)$ – $O(3N)$

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2x T_n(x) - T_{n-1}(x)$$

Sequential dependency between iterations but lowest algorithmic complexity

Barycentric interpolation – $O(5N+1)$

$$p(x) = \frac{\sum_{j=0}^n \frac{w_j f_j}{x - x_j}}{\sum_{j=0}^n \frac{w_j}{x - x_j}}, \quad x \neq x_j$$

Iterations are independent but higher arithmetic cost.

Step 1: define SIMD width

```
1 using batch = xsimd::batch<double>;  
2 constexpr int64_t simd_width = batch::size; // lanes
```

1. Calculate the SIMD loop boundary

- Determine how many elements can be processed in full SIMD-width blocks.
- Define **n_simd** is the largest multiple of the vector width \leq 'N'.

2. Split the loop body

- first loop for 'i < n_simd'

NOTE: Increment by simd_width

- Scalar cleanup for remaining elements is basically the old loop with only the index changed from n_simd to N

Step 2: Split the loop

1. Calculate the SIMD loop boundary

- Determine how many elements can be processed in full SIMD-width blocks.

```
1 const auto n_simd = (N / simd_width) * simd_width;
```

- Ensures that 'n_simd' is the largest multiple of the vector width $\leq N$.

2. Split the loop body

- Vectorized section for 'i < n_simd'
- Scalar cleanup for remaining elements is basically the old loop with only the index changed

```
1 // vectorized loop
2 for (int i = 0; i < n_simd; i += simd_width) {
3 }
4 // scalar remainder
5 for (int i = n_simd; i < N; ++i) {
6     // old body unchanged
7 }
```

Step 2: Split the loop

```
1  const auto n_simd = N & (-simd_width); // round down to multiple of simd_width
2  std::size_t i = 0;
3  for (; i < n_simd; i += simd_width) {
4  // we will populate this
5  }
6
7  double num = 0;
8  double den = 0;
9
10 for (; i < N; ++i) {
11     double diff = pt - x[i];
12     double q = w[i] / diff;
13     num += q * fvals[i];
14     den += q;
15 }
16 return num / den;
```

Step 3: loading data in register

```
1 // it takes a pointer in input!
2 // it returns a simd register
3 auto bdata = xsimd::load_unaligned(data* ptr);
4
5 // alternatively xsimd::batch can be initialized with a variable:
6 batch bnum(0); // this sets all the elements of the register to 0
7
8 // HINT: in c++, std constainers usually have a data() function that returns a pointer
9 // i.e. x.data()
```

Step 3: loading data in register

```
1 // it takes a pointer in input!
2 // it returns a simd register
3 auto bdata = xsimd::load_unaligned(data* ptr);
4
5 // alternatively xsimd::batch can be initialized with a variable:
6 batch bnum(0); // this sets all the elements of the register to 0
7
8 // HINT: in c++, std constainers usually have a data() function that returns a pointer
9 // i.e. x.data()
10 // load_unaligned(x[i]) will not work!
```

Step 3: loading data in register

```
1 const auto bx = batch::load_unaligned(x.data() + i);  
2 const auto bw = batch::load_unaligned(w.data() + i);  
3 const auto bf = batch::load_unaligned(fvals.data() + i);
```


Step 4: vectorizing the loop

`xsimd::batch` implements operator overloads, arithmetic operators `+`, `-`, `*`, and `/` and functions `xsimd::fma`, `xsimd::sqrt`, `xsimd::abs`, etc.) can be used directly on SIMD batches just like scalar values.

Now you have everything needed to vectorize the loop!

Step 5: Understanding partial SIMD results

```
1  const batch bpt(pt);
2  batch bnum(0);
3  batch bden(0);
4
5  std::size_t i = 0;
6  for (; i < n_simd; i += simd_width) {
7      const auto bx = batch::load_unaligned(x.data() + i);
8      const auto bw = batch::load_unaligned(w.data() + i);
9      const auto bf = batch::load_unaligned(fvals.data() + i);
10
11     const auto bdiff = bpt - bx;
12     const auto bq = bw / bdiff;
13     bnum += bq * bf;
14     // bnum = xsimd::fma(bq, bf, bnum);
15     bden += bq;
16 }
```

By vectorizing the loop, each SIMD register now holds **partial reductions** over different strided subsets of the data. For example:

$$b_{\text{num}}[0] = \sum_{i=0,8,16,\dots} q_i f_i, \quad b_{\text{num}}[1] = \sum_{i=1,9,17,\dots} q_i f_i, \quad \text{and so on,}$$

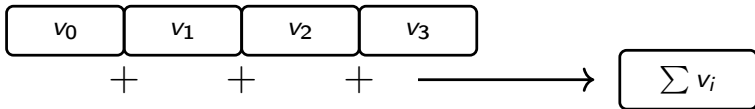
$$b_{\text{den}}[0] = \sum_{i=0,8,16,\dots} q_i, \quad b_{\text{den}}[1] = \sum_{i=1,9,17,\dots} q_i, \quad \text{and so on.}$$

Step 5: Horizontal sum

The **horizontal reduction** aggregates partial SIMD results into a scalar. Each lane in the batch contributes to the final sum.

$$\text{reduce_add}([v_0, v_1, \dots, v_{W-1}]) = \sum_{k=0}^{W-1} v_k$$

```
1 // Collapse SIMD register into a scalar value
2 const auto res = xsimd::reduce_add(data);
```



This step finalizes the SIMD computation by combining all per-lane partial results. Use other reductions like `reduce_max`, `reduce_min`, or `reduce_mul` for different aggregation operations.

Step 5: Horizontal sum

```
1 double num = xsimd::reduce_add(bnum);  
2 double den = xsimd::reduce_add(bden);
```

```
1 double operator()(const double pt) const {
2     for (int i = 0; i < N; ++i) {
3         if (pt == x[i]) { return fvals[i]; }
4     }
5     // // shorthand for the xsimd type
6     using batch = xsimd::batch<double>;
7     // simd width since it is architecture/compile flags dependent
8     constexpr std::int64_t simd_width = batch::size;
9
10    const auto n_simd = N & (-simd_width); // round down to multiple of simd_width
11
12    const batch bpt(pt);
13    batch bnum(0);
14    batch bden(0);
15    ...
16 }
```

```
1 double operator()(const double pt) const {
2     ...
3     std::size_t i = 0;
4     for (; i < n_simd; i += simd_width) {
5         const auto bx = batch::load_unaligned(x.data() + i);
6         const auto bw = batch::load_unaligned(w.data() + i);
7         const auto bf = batch::load_unaligned(fvals.data() + i);
8
9         const auto bdiff = bpt - bx;
10        const auto bq = bw / bdiff;
11        bnum += bq * bf;
12        // bnum = xsimd::fma(bq, bf, bnum);
13        bden += bq;
14    }
15
16    double num = xsimd::reduce_add(bnum);
17    double den = xsimd::reduce_add(bden);
18    ...
19 }
```

```
1 double operator()(const double pt) const {  
2     ...  
3     for (; i < N; ++i) {  
4         double diff = pt - x[i];  
5         double q = w[i] / diff;  
6         num += q * fvals[i];  
7         den += q;  
8     }  
9     return num / den;  
10 }
```