

Level 1 BLAS

Deval Deliwala

Table of Contents

- Level 1 BLAS
 - Optimizing the Dot Product
 - Naive implementation
 - Optimized Implementation
 - Optimizing Vector Addition
 - Naive Implementation
 - Optimized Implementation
 - LLVM is very impressive
 - Analyzing Assembly
 - What to look for in the assembly
 - Naive saxyo autovectorizes
 - SIMD saxpy lowers to the same NEON structure
 - Why both implementations are equally fast

BLAS is divided into three levels:

- Level 1 - vector-vector routines.
- Level 2 - matrix-vector routines.
- Level 3 - matrix-matrix routines.

Each level is progressively more difficult to optimize than the previous. With modern CPUs, Level 1 and 2 are mostly memory-bound. And Level 3 is mostly compute-bound.

Consequently, optimizing Level 1 is straightforward and leans heavily on LLVM and portable-simd. portable-simd is a natively SIMD interface in the Rust standard library that maps cleanly onto modern vector structures across architectures.

Previously, I designed a clean API for my BLAS implementation in Rust. It contains `VectorRef` and `VectorMut` types that internally handle vector buffers, strides, and offsets cleanly. The separation of `Ref`/`Mut` types also intuitively allow function calls to be impossible to confuse:

- `VectorRef` means "this routine may only `read` from it."
- `VectorMut` means "this routine may `write` into it."

Optimizing the Dot Product

The `sdot` routine calculates the dot product of two vectors:

$$\vec{x} \cdot \vec{y} = \sum_{i=0}^{n-1} x_i y_i,$$

where n is the length of \vec{x} and \vec{y} .

This routine does not mutate or overwrite any vector. It only outputs the calculated f32 product. So it uses `VectorRef`s.

Naive implementation

Contiguous memory means the vector is tightly packed. The next element is at the next index.

For example, consider $x = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{pmatrix}$. It could be represented as follows:

- contiguous:
 - let $x = \text{vec}([0, 1];$
 - increment $\text{incx} = 1$
- not contiguous:
 - let $x = \text{vec}([0, 0, 1];$
 - increment $\text{incx} = 2$, accessing every other element.

When memory is contiguous, it can be brought to the CPU together in cachelines. This yields a much faster execution time. Consequently I have a fast path for when vectors are contiguous ($\text{inc} = 1$) and a slow path otherwise.

```
use crate::types::VectorRef;

// Computes the dot product of two [VectorRefs].
#[inline]
pub fn sdot(
    x: VectorRef<_, f32>,
    y: VectorRef<_, f32>,
) -> f32 {
    let xn = x.n();
    let yn = y.n();

    if xn != yn {
        panic!("x < y vector dimensions must match!");
    }

    // empty vector
    if xn == 0 {
        return 0.0;
    }

    let mut acc_sum = 0.0;

    // fast path
    if let (Some(xs), Some(ys)) = (x.contiguous_slice(), y.contiguous_slice()) {
        for (xk, yk) in xs.iter().zip(ys.iter()) {
            acc_sum += xk * yk;
        }
    }

    // slow path
    let inx = x.stride();
    let incy = y.stride();

    let ix = x.offset();
    let iy = y.offset();

    let xs = x.as_slice();
    let ys = y.as_slice();
    let xs_it = xs[ix..].iter().step_by(inx).take(xn);
    let ys_it = ys[iy..].iter().step_by(incy).take(yn);

    for (xk, yk) in xs_it.zip(ys_it) {
        acc_sum += xk * yk
    }

    acc_sum
}
```

When vectors x and y contain 1024 elements, this routine runs in 750 nanoseconds on average, which is already extremely fast. On modern CPUs, LLVM can often vectorize patterns like

$$acc_sum += xk + yk$$

into SIMD instructions automatically when the access pattern is simple and contiguous. The work that has gone into making modern compilers like LLVM intelligent enough to do this is incredible.

I'll show in the [Assembly](#) section at the end for SAXPY, these instructions compile down to SIMD instructions for multiplying and adding.

Optimized Implementation

However, I can make it faster by writing the SIMD myself. I use `portable-simd`, which [compile to the best available SIMD instructions](#) for all modern computer architectures. On AArch64 (including Apple M4) architectures, the SIMD interface is called "NEON".

The algorithm is the same. And SIMD only really works with BLAS when vectors are contiguous, so the slow path stays exactly the same.

The rough procedure for working with SIMD in the fast path goes as follows:

```
// define chunk size at compile time
const LANES: usize = <some value>;

// decompose vector into chunks of size LANES
// and the leftover tail of length < LANES
let (chunks, tail) = vector_as_chunks(<LANES>);

for chunk in chunks
    // convert to SIMD vector
    let simd_chunk = Simd::from_array(chunk);
    |<do some stuff>
end

// leftover tail scalar path
for value in tail
    |<do some stuff>
end
```

I hope my code is readable enough to understand this:

```
/// Level 1 [~?DOT~](https://www.netlib.org/lapack/explore-html/d1/dcc/group_dot.html)
/// routine in single precision.
/// [
///   (sum_{i=0}^{n-1} x_i) * y_i
/// ]
/// # Author
/// # Deval Deliwala

use std::simd::Simd;
use std::simd::num::SimdFloat;
use crate::types::VectorRef;
use crate::debug_assert_n_eq;
```

```
/// Takes the dot product over logical elements in [VectorRef]
/// `x` and `y` are SIMD vectors.
/// Arguments:
///   * `x` : [VectorRef] - over [f32]
///   * `y` : [VectorRef] - over [f32]
/// Returns:
///   - [f32] dot product.
#[inline]
pub fn sdot(
    x: VectorRef<_, f32>,
    y: VectorRef<_, f32>,
) -> f32 {
    // ensures x and y have same length `n`
    debug_assert_n_eq!(x, y);

    let n = x.n();
    if n == 0 {
        return 0.0;
    }

    // fast path
    if let (Some(xs), Some(ys)) = (x.contiguous_slice(), y.contiguous_slice()) {
        const LANES: usize = 32;
        let a = Simd::from_f32(LANES); // splat(alpha);
        let (xv, xt) = xs.as_chunks::<LANES>();
        let (yv, yt) = ys.as_chunks::<LANES>();

        acc += xv * yv;
    }

    // slow path
    let mut acc = 0.0;
    let inx = x.stride();
    let incy = y.stride();
    let ix = x.offset();
    let iy = y.offset();

    let xs = x.as_slice();
    let ys = y.as_slice();
    let xs_it = xs[ix..].iter().step_by(inx).take(n);
    let ys_it = ys[iy..].iter().step_by(incy).take(n);

    for (xv, yv) in xs_it.zip(ys_it) {
        acc += xv * yv;
    }

    acc
}
```

This is the only tricky part is learning the `portable-simd` syntax and tuning the LANES vector length. I have an Apple M4 Pro, whose NEON vector registers are 128-bit wide, i.e. 4 f32s per vector operation. This means SIMD can apply the same arithmetic to four f32s in parallel.

Specifically,

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} * \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix} = \begin{pmatrix} x_1 y_1 \\ x_2 y_2 \\ x_3 y_3 \\ x_4 y_4 \end{pmatrix}$$

Applies the same instruction across four lanes at once. Based on this, setting `LANES = 4` would be reasonable. This would separate x and y into chunks of 4 f32s at a time, which is perfect for Apple M4's 128-bit registers.

However, within the hot loop, there is still per-iteration overhead: loop control, bounds/tail handling, and moving chunks within SIMD values. Increasing to `LANES = 32` batches more work per iteration, so the loop runs 8x fewer iterations than `LANES = 4`. This is because 32 f32s get processed per iteration, instead of just 4.

For example, let \vec{x} and \vec{y} have length 1024.

$$if \text{LANES} = 4 \rightarrow 256 \text{ chunks of } x, y$$

$$\text{if } \text{LANES} = 32 \rightarrow 32 \text{ chunks of } x, y$$

$$\rightarrow 8 \times \text{less iterations}$$

Despite vector registers only storing 4 f32s at a time, processing 8 registers (`LANES = 32`) at a time is more efficient than matching native register width.

When vectors x and y contain 1024 elements, this routine runs in

$$\text{LANES} = 4: 416ns$$

$$\text{LANES} = 16: 166ns$$

$$\text{LANES} = 32: 125ns.$$

These are all extremely fast. But setting `LANES = 32` is fastest and 10.328x faster than the naive scalar loop implementation.

The `saxy` routine performs

$$y \leftarrow \alpha x + y,$$

A * ph * $\text{XP}^T \text{Y}$, and y gets overwritten with the solution. Hence, we use `VectorRef` for x and a mutable `VectorMut` for y .

The procedure is similar to the dot product. However, the SIMD-optimized improvement is very different.

Naive Implementation

For this comparison, only three instruction patterns matter:

- NEON vector loads and stores: `ldr q1, q2, [x12]`, `str s1, s2, [y12]`
- SIMD arithmetic on four f32s at once: `fmul.4s, fadd.4s`
- Scalar remainder loops: `ldr s1, s2, f12`, `fadd.4s, fadd.4s`

When these appear in a tight loop with no intervening calls or branches to panic paths, the loop is fully vectorized and bounds checks have been hoisted out.

Naive saxy autovectorizes

The core of the naive implementation is:

```
for (xk, yk) in xs.iter().zip(ys.iter_mut()) {
    *yk += alpha * *xk;
}
```

LLVM recognizes this as a SAXPY pattern and emits a NEON loop in the contiguous fast path. The following block from this naive `saxy` implementation is the hot loop:

```
LBB0_31:
    ldp q1, q2, [x12], #32
    ldp q3, q4, [x12], #64
    ldp a4, v1, v2, v8@0
    ldp a5, q6, [x13], #32
    ldp a6, q7, q16, [x13]
    fadd.4s v1, v2, v6, v2
    fadd.4s v3, v7, v3
    fadd.4s v4, v8, v4
    fadd.4s v5, v9, v5
    fadd.4s v6, v10, v6
    fadd.4s v7, v11, v7
    fadd.4s v8, v12, v8
    fadd.4s v9, v13, v9
    fadd.4s v10, v14, v10
    fadd.4s v11, v15, v11
    fadd.4s v12, v16, v12
    fadd.4s v13, v17, v13
    fadd.4s v14, v18, v14
    fadd.4s v15, v19, v15
    fadd.4s v16, v20, v16
    fadd.4s v17, v21, v17
    fadd.4s v18, v22, v18
    fadd.4s v19, v23, v19
    fadd.4s v20, v24, v20
    fadd.4s v21, v25, v21
    fadd.4s v22, v26, v22
    fadd.4s v23, v27, v23
    fadd.4s v24, v28, v24
    fadd.4s v25, v29, v25
    fadd.4s v26, v30, v26
    fadd.4s v27, v31, v27
    fadd.4s v28, v32, v28
    fadd.4s v29, v33, v29
    fadd.4s v30, v34, v30
    fadd.4s v31, v35, v31
    fadd.4s v32, v36, v32
    fadd.4s v33, v37, v33
    fadd.4s v34, v38, v34
    fadd.4s v35, v39, v35
    fadd.4s v36, v40, v36
    fadd.4s v37, v41, v37
    fadd.4s v38, v42, v38
    fadd.4s v39, v43, v39
    fadd.4s v40, v44, v40
    fadd.4s v41, v45, v41
    fadd.4s v42, v46, v42
    fadd.4s v43, v47, v43
    fadd.4s v44, v48, v44
    fadd.4s v45, v49, v45
    fadd.4s v46, v50, v46
    fadd.4s v47, v51, v47
    fadd.4s v48, v52, v48
    fadd.4s v49, v53, v49
    fadd.4s v50, v54, v50
    fadd.4s v51, v55, v51
    fadd.4s v52, v56, v52
    fadd.4s v53, v57, v53
    fadd.4s v54, v58, v54
    fadd.4s v55, v59, v55
    fadd.4s v56, v60, v56
    fadd.4s v57, v61, v57
    fadd.4s v58, v62, v58
    fadd.4s v59, v63, v59
    fadd.4s v60, v64, v60
    fadd.4s v61, v65, v61
    fadd.4s v62, v66, v62
    fadd.4s v63, v67, v63
    fadd.4s v64, v68, v64
    fadd.4s v65, v69, v65
    fadd.4s v66, v70, v66
    fadd.4s v67, v71, v67
    fadd.4s v68, v72, v68
    fadd.4s v69, v73, v69
    fadd.4s v70, v74, v70
    fadd.4s v71, v75, v71
    fadd.4s v72, v76, v72
    fadd.4s v73, v77, v73
    fadd.4s v74, v78, v74
    fadd.4s v75, v79, v75
    fadd.4s v76, v80, v76
    fadd.4s v77, v81, v77
    fadd.4s v78, v82, v78
    fadd.4s v79, v83, v79
    fadd.4s v80, v84, v80
    fadd.4s v81, v85, v81
    fadd.4s v82, v86, v82
    fadd.4s v83, v87, v83
    fadd.4s v84, v88, v84
    fadd.4s v85, v89, v85
    fadd.4s v86, v90, v86
    fadd.4s v87, v91, v87
    fadd.4s v88, v92, v88
    fadd.4s v89, v93, v89
    fadd.4s v90, v94, v90
    fadd.4s v91, v95, v91
    fadd.4s v92, v96, v92
    fadd.4s v93, v97, v93
    fadd.4s v94, v98, v94
    fadd.4s v95, v99, v95
    fadd.4s v96, v100, v96
    fadd.4s v97, v101, v97
    fadd.4s v98, v102, v98
    fadd.4s v99, v103, v99
    fadd.4s v100, v104, v100
    fadd.4s v101, v105, v101
    fadd.4s v102, v106, v102
    fadd.4s v103, v107, v103
    fadd.4s v104, v108, v104
    fadd.4s v105, v109, v105
    fadd.4s v106, v110, v106
    fadd.4s v107, v111, v107
    fadd.4s v108, v112, v108
    fadd.4s v109, v113, v109
    fadd.4s v110, v114, v110
    fadd.4s v111, v115, v111
    fadd.4s v112, v116, v112
    fadd.4s v113, v117, v113
    fadd.4s v114, v118, v114
    fadd.4s v115, v119, v115
    fadd.4s v116, v120, v116
    fadd.4s v117, v121, v117
    fadd.4s v118, v122, v118
    fadd.4s v119, v123, v119
    fadd.4s v120, v124, v120
    fadd.4s v121, v125, v121
    fadd.4s v122, v126, v122
    fadd.4s v123, v127, v123
    fadd.4s v124, v128, v124
    fadd.4s v125, v129, v125
    fadd.4s v126, v130, v126
    fadd.4s v127, v131, v127
    fadd.4s v128, v132, v128
    fadd.4s v129, v133, v129
    fadd.4s v130, v134, v130
    fadd.4s v131, v135, v131
    fadd.4s v132, v136, v132
    fadd.4s v133, v137, v133
    fadd.4s v134, v138, v134
    fadd.4s v135, v139, v135
    fadd.4s v136, v140, v136
    fadd.4s v137, v141, v137
    fadd.4s v138, v142, v138
    fadd.4s v139, v143, v139
    fadd.4s v140, v144, v140
    fadd.4s v141, v145, v141
    fadd.4s v142, v146, v142
    fadd.4s v143, v147, v143
    fadd.4s v144, v148, v
```

