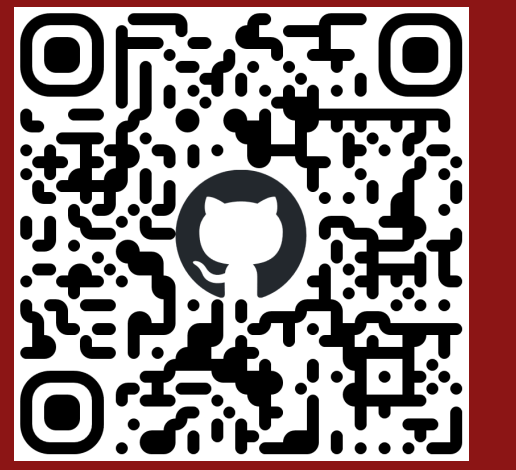


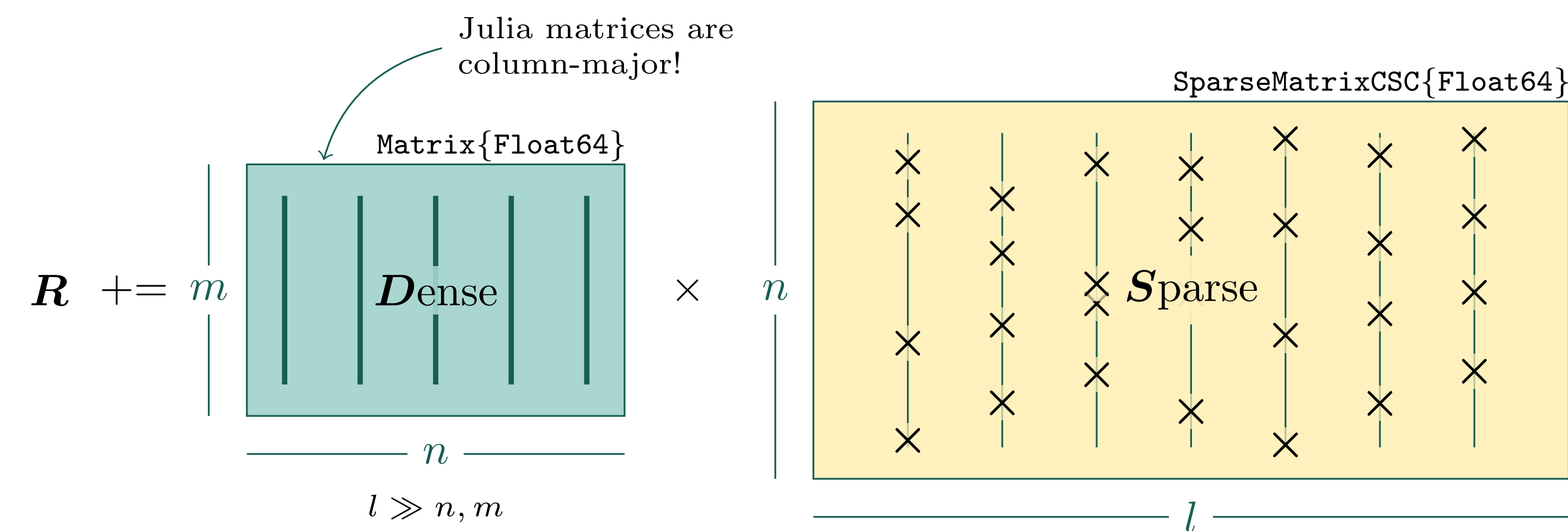
# ThreadedDenseSparseMul.jl

## Outperforming MKLSparse in 6 LOC.



Romeo Valentin (romeov@stanford.edu), Mykel J. Kochenderfer  
Stanford Intelligent Systems Lab

We need...

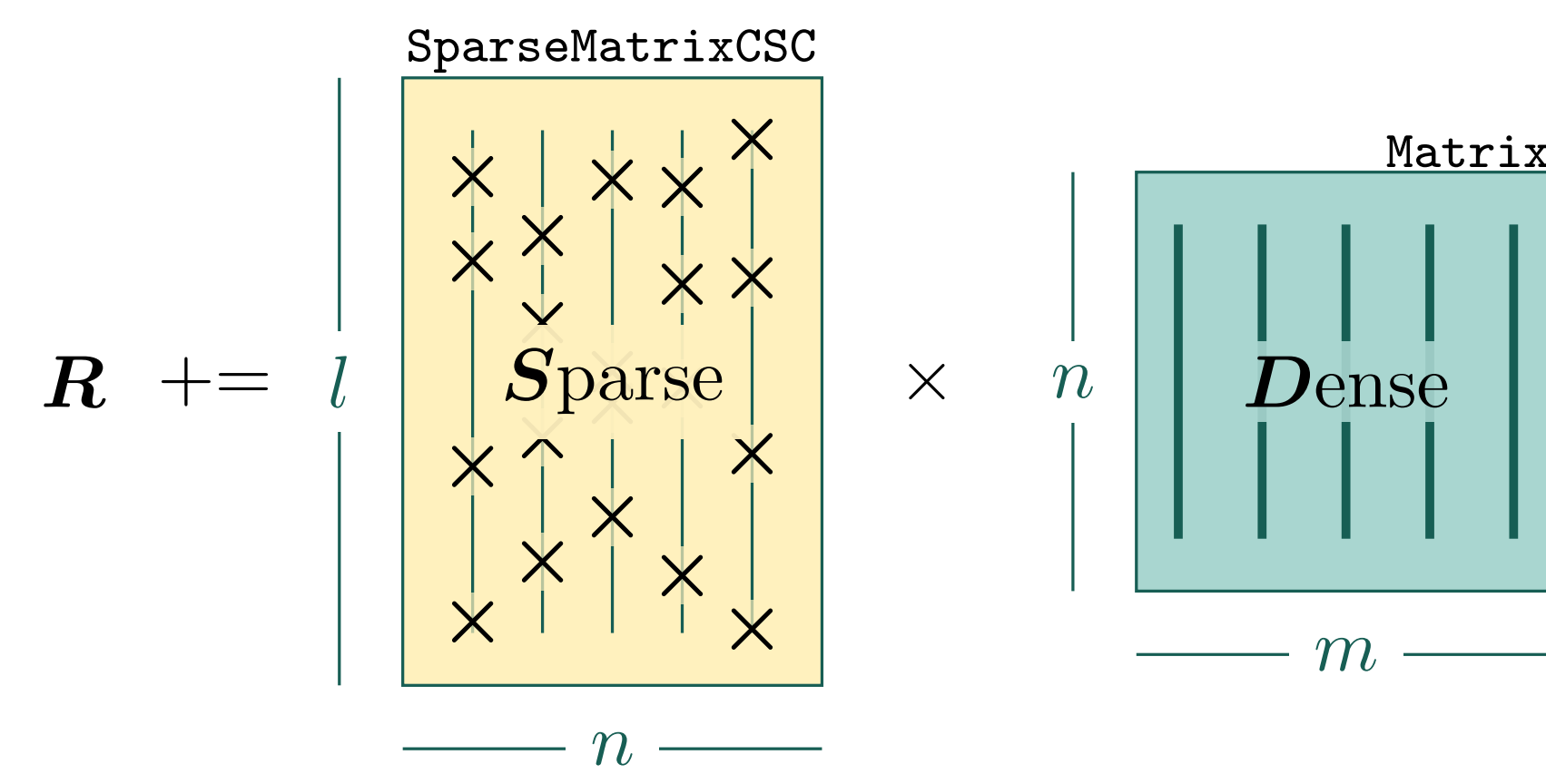


i.e., threaded dense times sparse matmul.

Our options:

- ThreadedDenseSparseMul.jl (this poster)

...but everyone else does



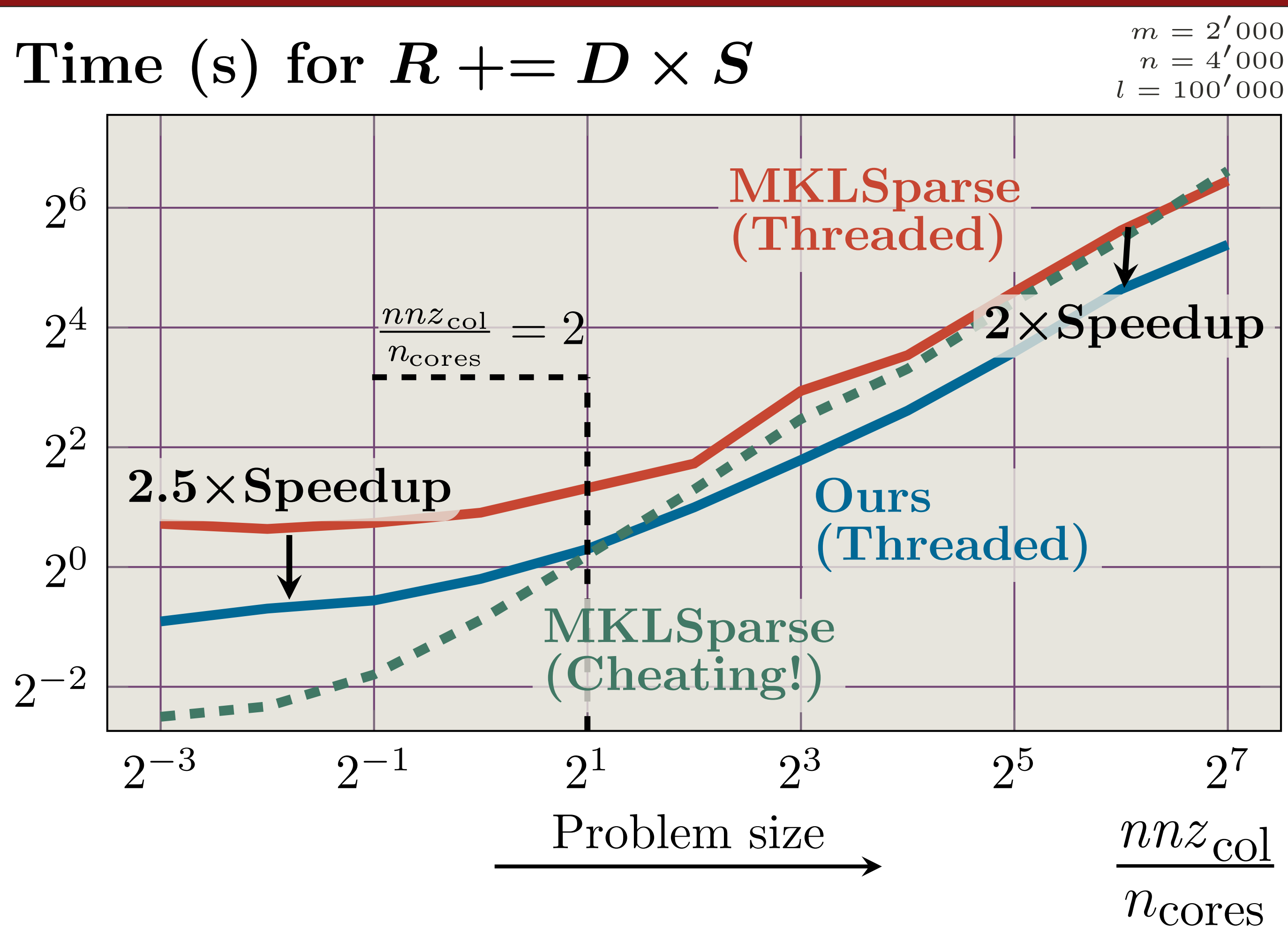
i.e., threaded sparse times dense matmul.

Our options:

- MKLSparse.jl, ThreadedSparseCSR.jl, ThreadedSparseArrays.jl

## Outperforming MKLSparse: Benchmark results.

Time (s) for  $R += D \times S$



**Ours** outperforms **MKLSparse** with speedup 2×!  
For large  $nnz_{col}$ , **Ours** even outperforms **MKLSparse (Cheating!)** which computes sparse×dense directly.

## Why can we win this?

Better memory layout. MKLSparse transposes (twice).  
Let's count memory movements of  $mul!(R', S', D')'$ .

- Transpose:  $\approx 2 \times \Theta(m \cdot l)$
- Multiply:  $\Theta(nnz_{col} \cdot m \cdot l) / n_{cores}$

$$\Rightarrow \left(2 + \frac{nnz_{col}}{n_{cores}}\right) \times \Theta(m \cdot l)$$

$\Rightarrow$  For  $(nnz_{col} > 2 \times n_{cores})$ , transposing dominates!

## Takeaways

- Good memory layout lets us outperform MKL.
- Don't be afraid to "hand-roll" tight loops.
- Know when you are memory or compute bound.
- Benchmark, benchmark, benchmark!

## The Secret Sauce: Chris Elrod's Polyester.jl

Hand-rolling loops has never been so easy! Also, memory layout is actually perfect for dense×sparse!

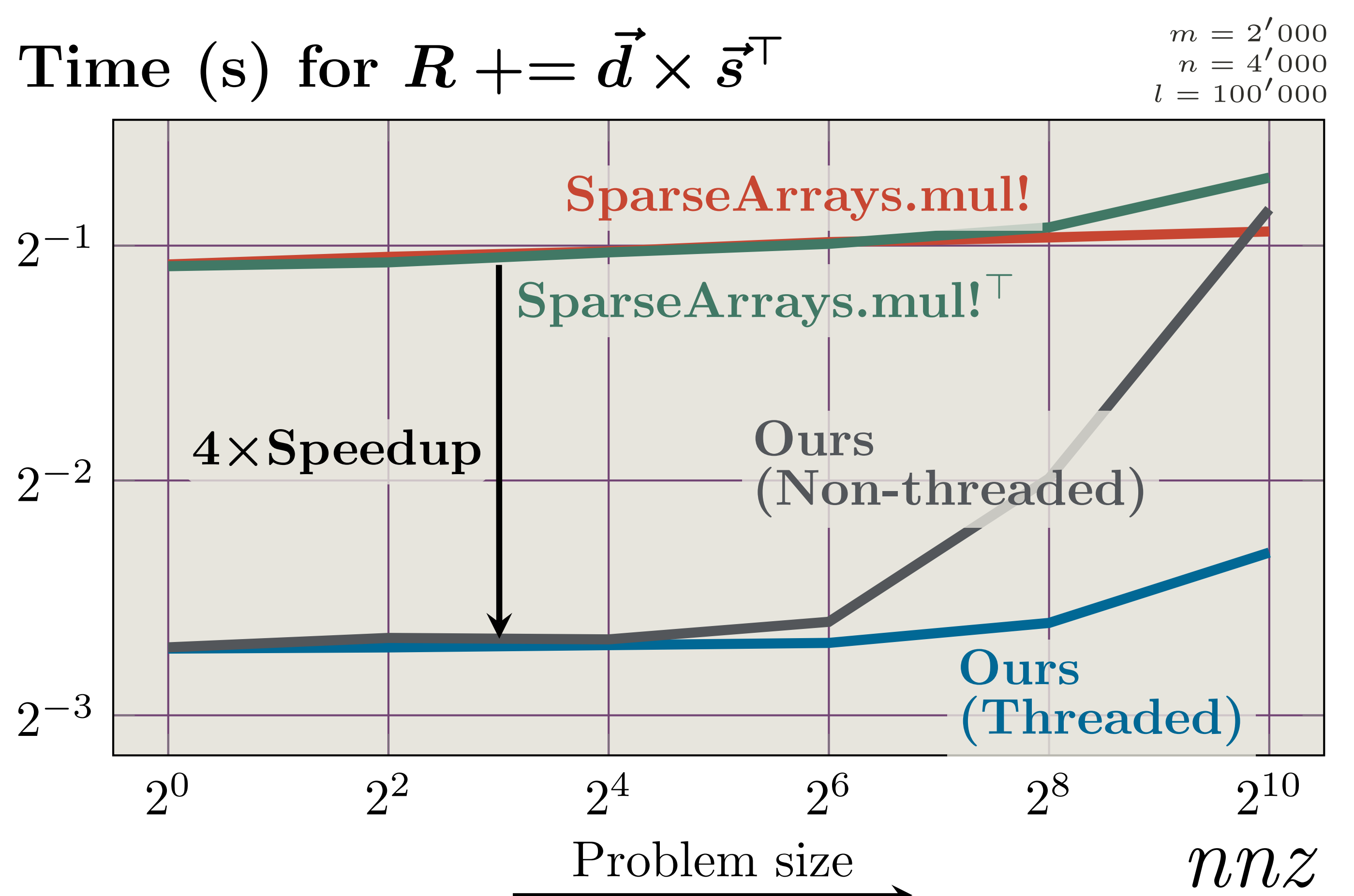
```
using Polyester
function mul!(C::Matrix{T}, A::Matrix{T}, B::SparseMatrixCSC{T},
              α::Number, β::Number)
1  @inbounds begin
2  |   C .= β
3  |   @batch for j in axes(B, 2)
4  |       C[:, j] += A * (α.*B[:, j])
5  |   end
6  |   return C
| end
end
```

## Things get even better: outer product.

MKLSparse doesn't have an outer product at all!

And we can outperform SparseArrays.

Time (s) for  $R += \vec{d} \times \vec{s}^T$



**Ours** beats **SparseArrays.mul! / mul!^T** with speedup 4× in 5 LOC (single-threaded, no Polyester.jl), or 9 LOC (multi-threaded, w/ Polyester.jl).

All benchmarks collected on mobile workstation with 11th Gen Intel(R) Core(TM) i7-11800H @ 2.30GHz, 32GB RAM, using the Chairmarks.jl library.