

# An optimized Baum-Welch algorithm

Cheuk Yu Chan

Franz Knobel

Josua Cantieni

Ramon Witschi



Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

# Baum-Welch algorithm

- Computes **Probability distribution** over a model
- The hidden states of a **hidden Markov model**
- Using **expectation-maximization**

We decide/are given

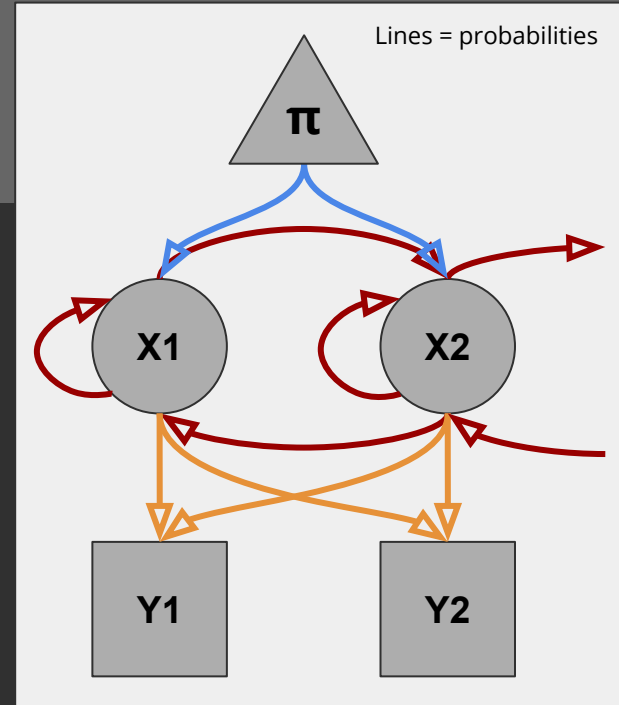
**N:** # hidden states

**M:** # observations

**K:** # observation sequences

**T:** # time steps

**Max\_iterations:** # iterations to perform



# Task of the Baum-Welch algorithm

Given observation sequence

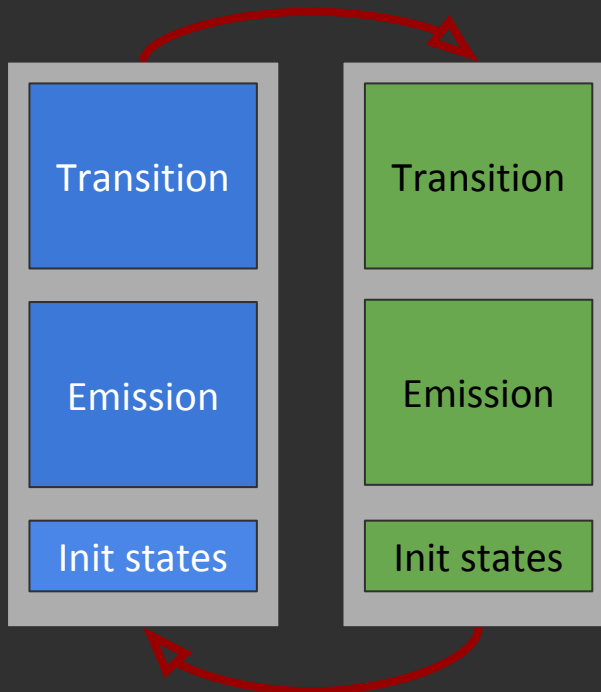
O1 O2 O3 O4 O5 O6 ... .. OM

Random init

Transition

Emission

Init states



Local Maxima

Transition

Emission

Init states



# Cost Analysis

*Asymptotic runtime per iteration:*

$$O(KTN * (M+N)) \approx O(N^2)$$

*Flops per iteration:*

$$9*TKNN - 5*KNN + NN + 8*TKN + 2*KNM + 3*KN + 2*TK + NM + K + N \text{ Flops}$$

*Memory usage:*

$$(\text{max\_iterations} + KTNN + KNN + NN + KTN + NM + KN + 2*KT + N)*8 + 144 \text{ Bytes}$$

**Featured weight:** NN

# Baseline Implementation + Performance

## *Implementation:*

- Implemented code from the referenced tutorial slides
- Struct with pointer to arrays for convenience

## *Assumptions:*

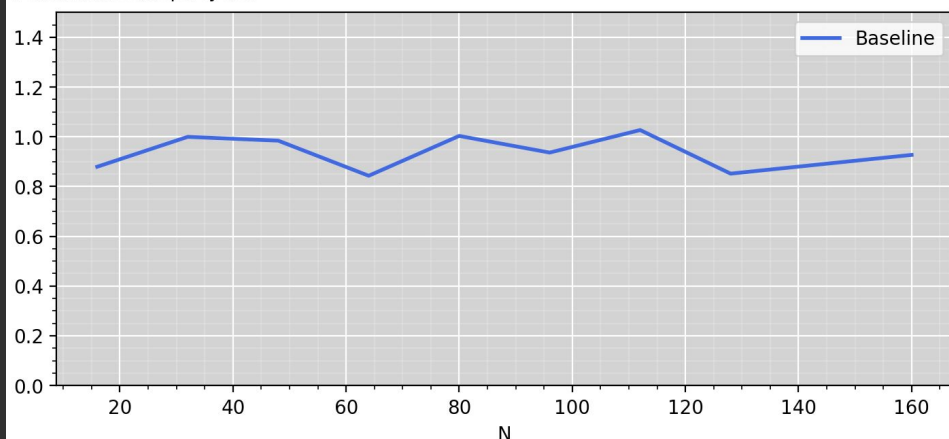
- Divisibility of N, M and K by 16 and T by 32

## *Verification:*

1. Compare baseline against other implementations known to be correct
2. Compare our optimizations against our baseline

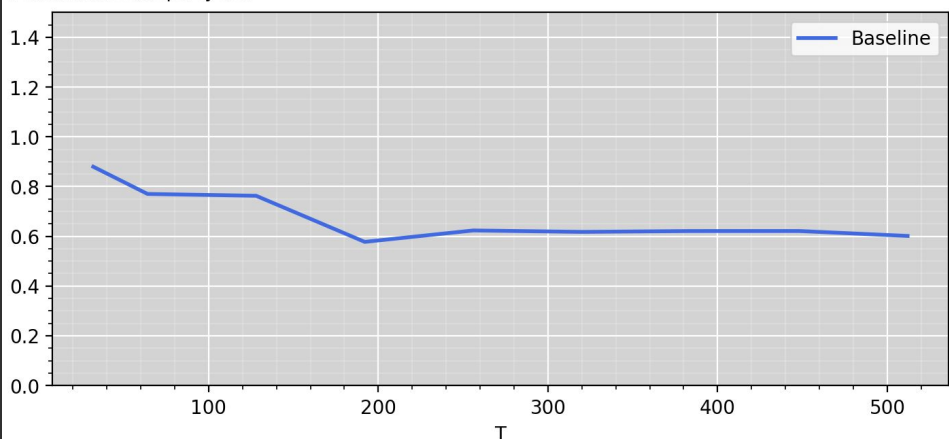
**Performance Plot - (M=16, T=32, K=16)**  
**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**

Performance [flops/cycle]



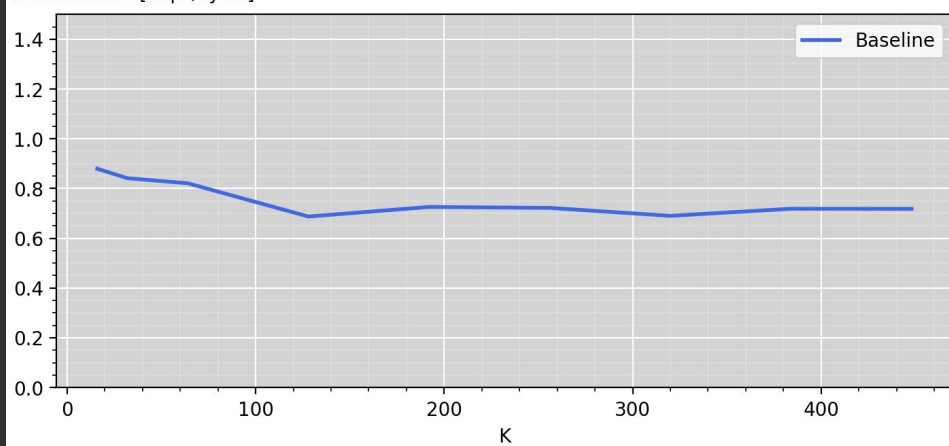
**Performance Plot - (M=16, K=16, N=16)**  
**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**

Performance [flops/cycle]



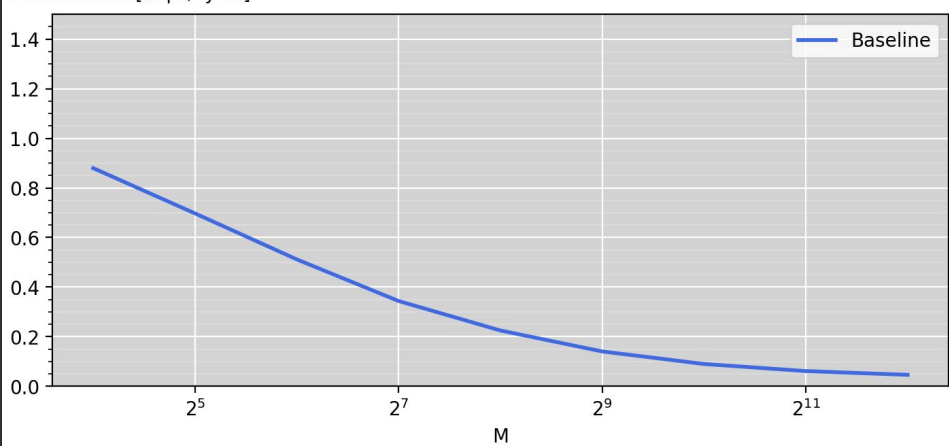
**Performance Plot - (M=16, T=32, N=16)**  
**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**

Performance [flops/cycle]



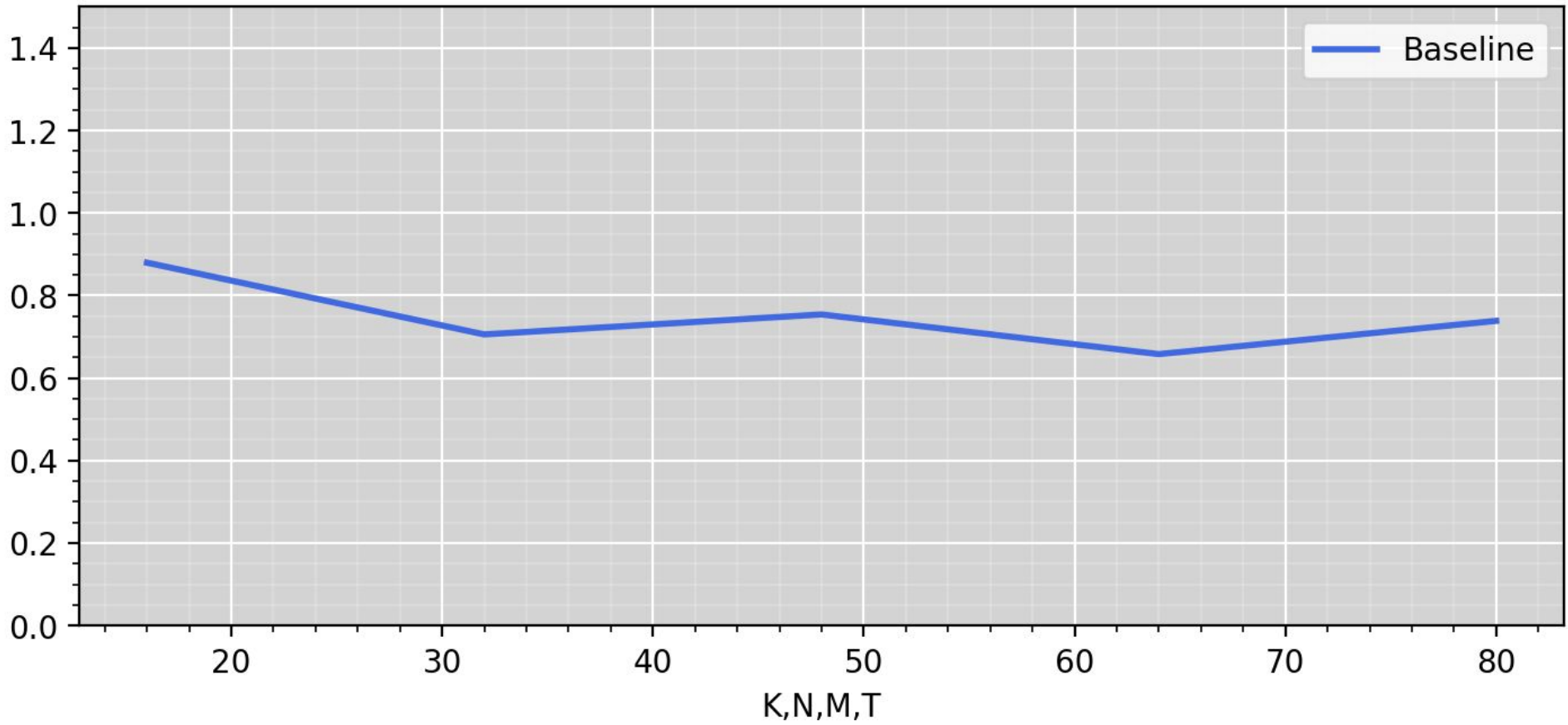
**Performance Plot - (T=32, K=16, N=16)**  
**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**

Performance [flops/cycle]



**Performance Plot - (T has a step size of 32)**  
**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**

Performance [flops/cycle]

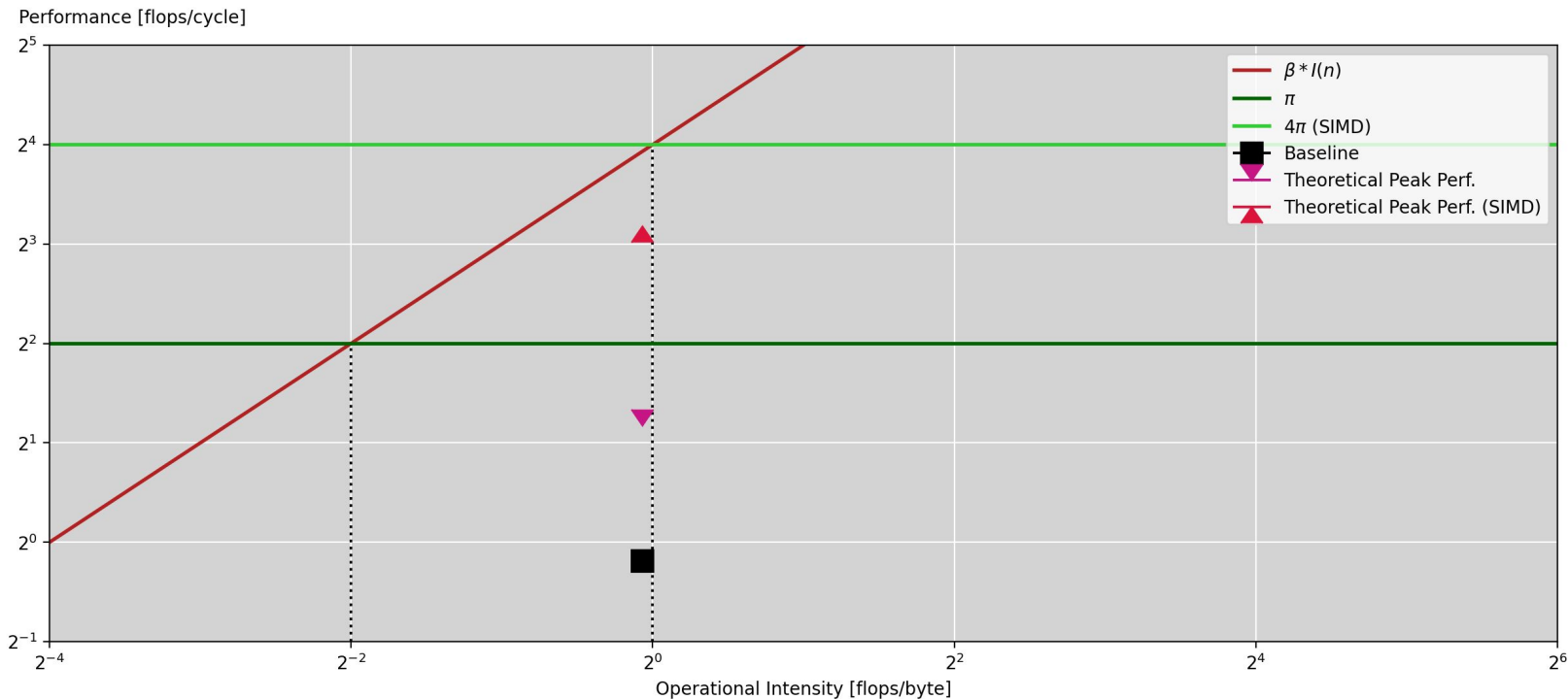


# Baseline Roofline Plot

Haswell i7-4710MQ, 2.5GHz, 12GB RAM

**Roofline Plot - (K=16,N=16,M=16,T=32)**

**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**





# Optimizations we performed

- Precomputing (taking code out of the loop)
- Reordering of the code (combining loops, transformed some terms)
- Scalar replacement
- Unrolling
- Blocking
- Vectorization
- Using alternative input (transformed a matrix)

# Optimizations we performed - reordering

## *Observation:*

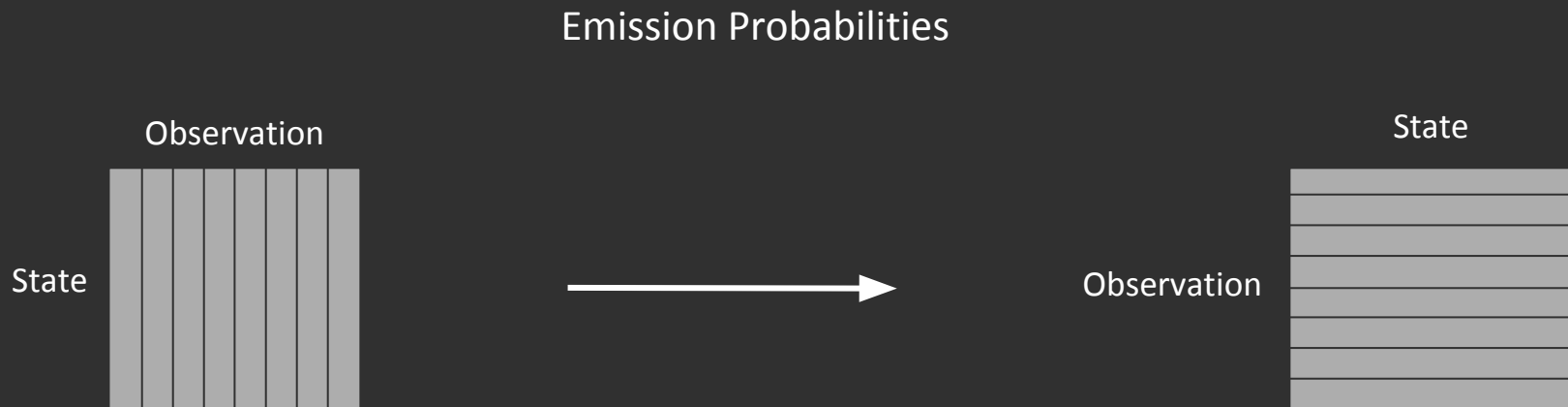
- Data dependency (e.g. forwards-step before backwards-step)
- Once all data is computed, dependency can be computed (gamma, sigma)

*Conclusion:* Merging of loops possible

## *Result:*

- + Efficient loop usage
- + Less Division operations
- + Less calls to log

# Optimizations we performed - alternative input



An implementation can request the input matrix for the emission probabilities to be transposed. This improves locality.

# Optimizations we performed - vectorization

- Application of knowledge gained in exercises  
(e.g. *horizontal addition*)
- Awareness of data dependency during unrolling  
(e.g. *usage of  $T-1$  or  $T+1$  in calculus prevents unrolling  $T$* )

# Experiments

*CPU:* Haswell i7-4710MQ, 2.5GHz, 12GB RAM

- L1: 32 KB (8-way associative, instruction+data cache)
- L3: 6 MB 12-way set associative shared cache
- AVX2 supported

*Compiler:* gcc 9.2.0 and clang 7.01

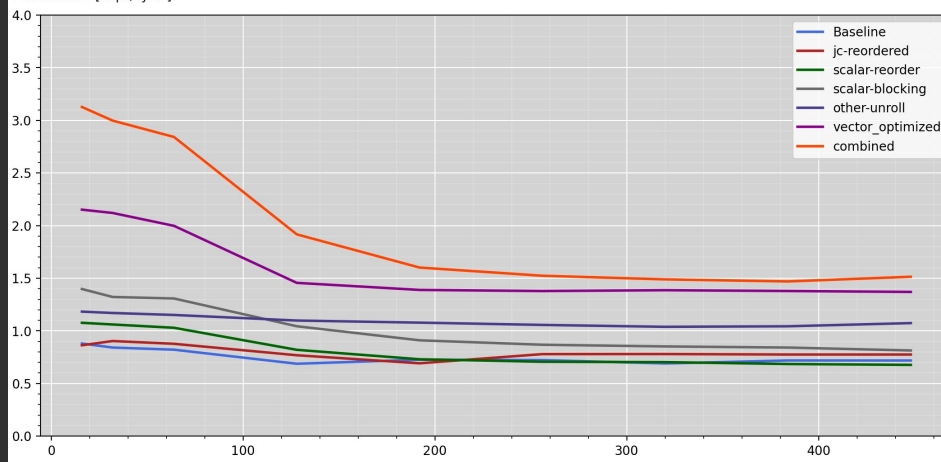
*Flags:*

- -O3 -ffast-math -mavx2 -mfma -march=native
- -O3 -fno-vectorize -ffast-math -mavx2 -mfma -march=native
- -O3 -funroll-loops -mavx2 -mfma -march=native

Performance Plot - (N=16, T=32, M=16)

gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native

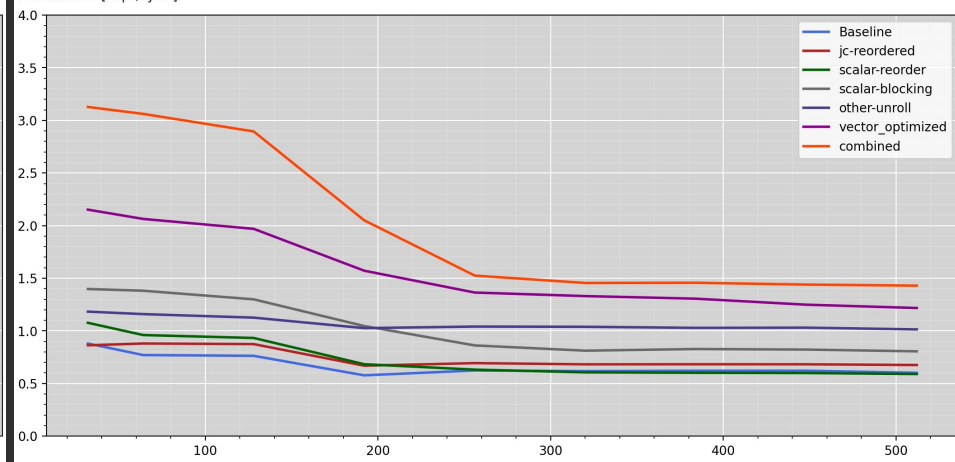
Performance [flops/cycle]



Performance Plot - (K=16, N=16, M=16)

gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native

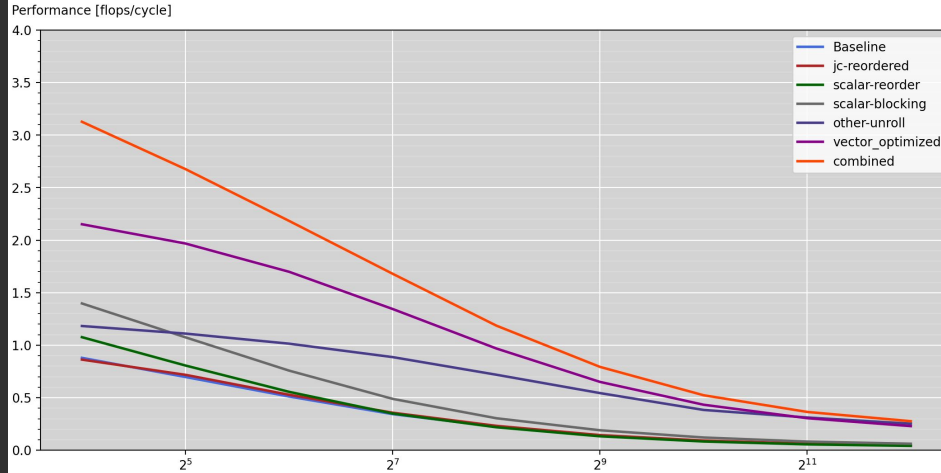
Performance [flops/cycle]



Performance Plot - (K=16, N=16, T=32)

gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native

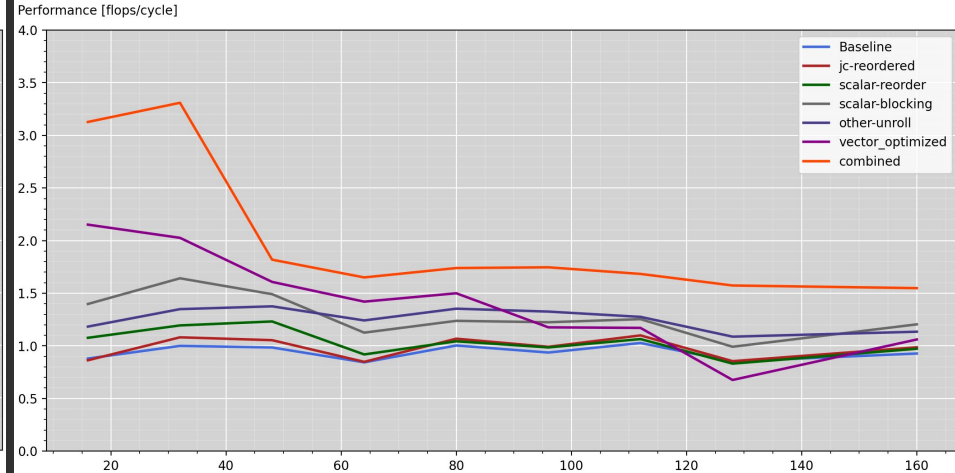
Performance [flops/cycle]



Performance Plot - (K=16, T=32, M=16)

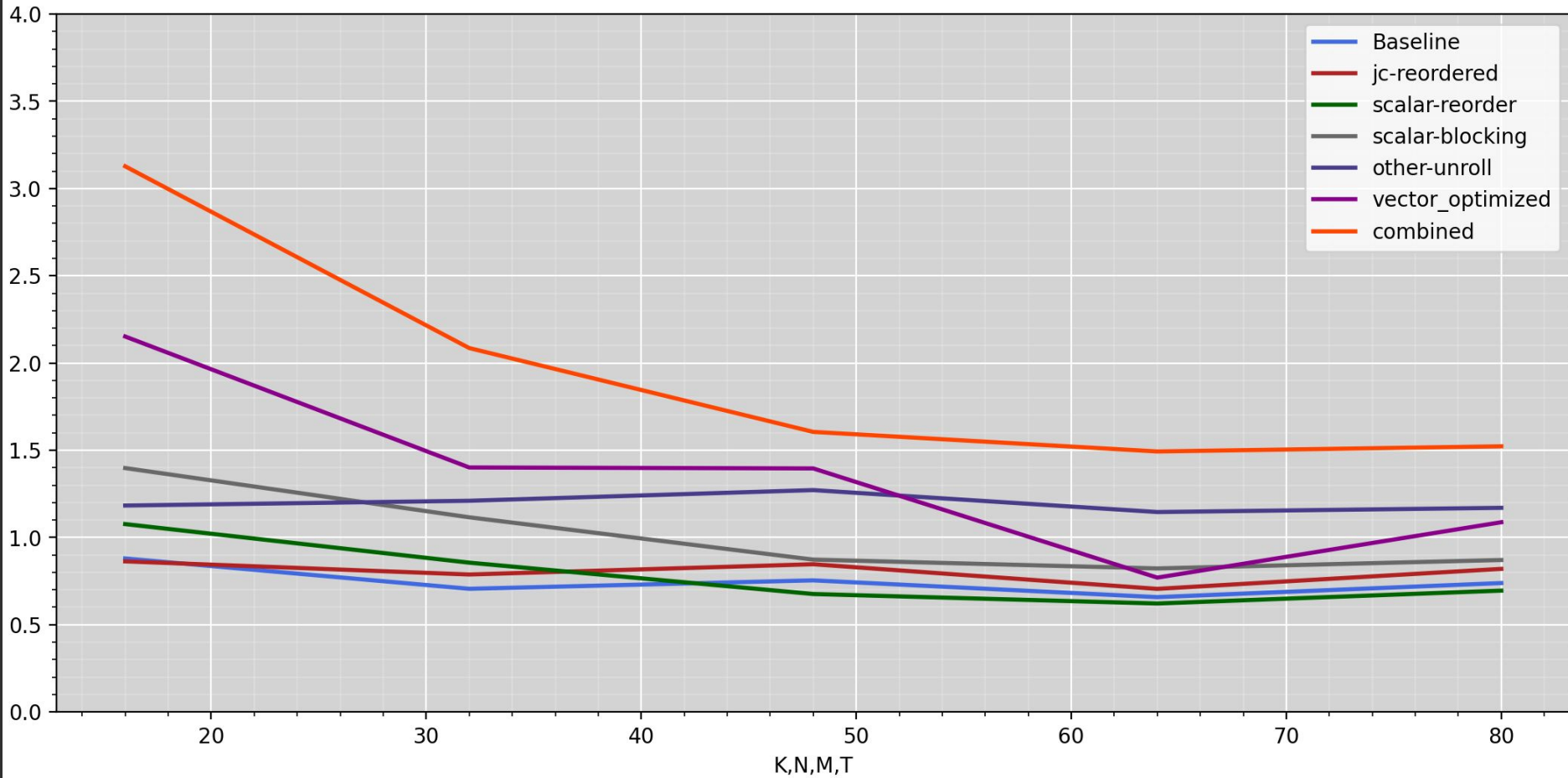
gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native

Performance [flops/cycle]



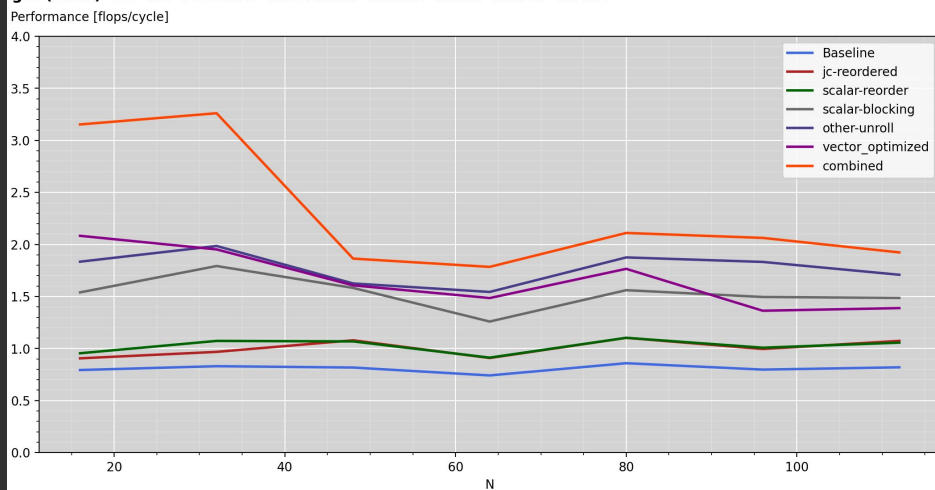
**Performance Plot - (N=16, T=32, M=16)**  
**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**

Performance [flops/cycle]

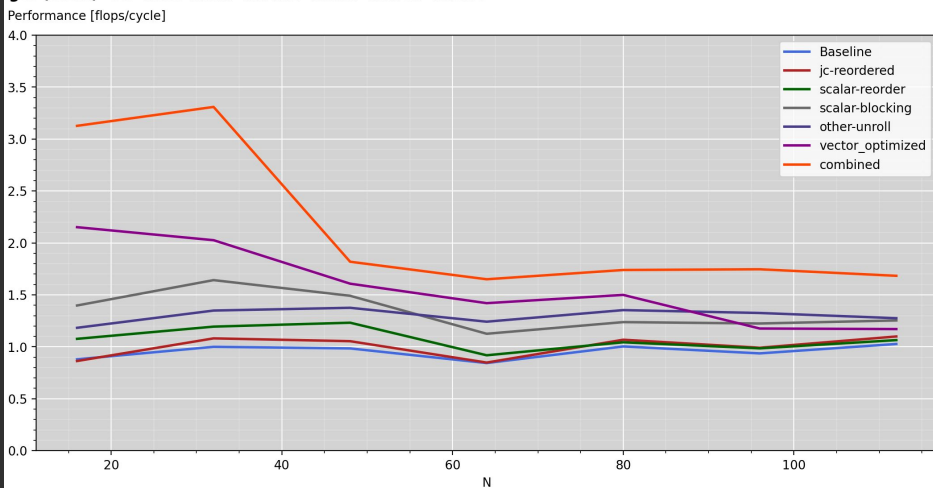


# Different Flags (gcc)

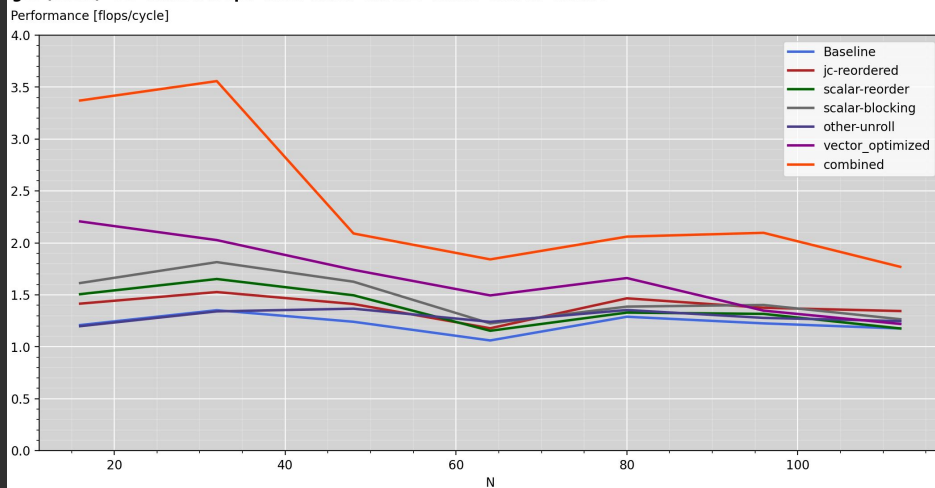
Performance Plot - (M=16, K=16, T=32)  
gcc (9.2.0) -O3 -fno-vectorize -ffast-math -mavx2 -mfma -march=native



Performance Plot - (M=16, K=16, T=32)  
gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native



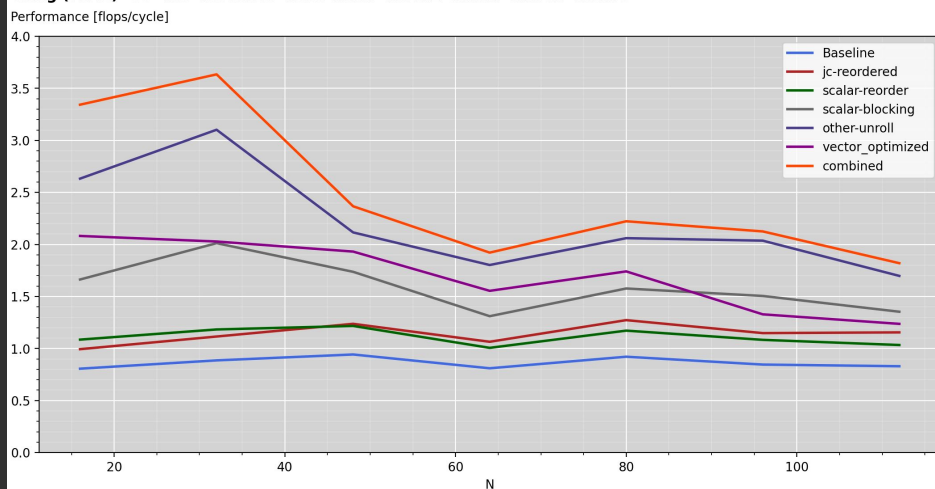
Performance Plot - (M=16, K=16, T=32)  
gcc (9.2.0) -O3 -funroll-loops -ffast-math -mavx2 -mfma -march=native



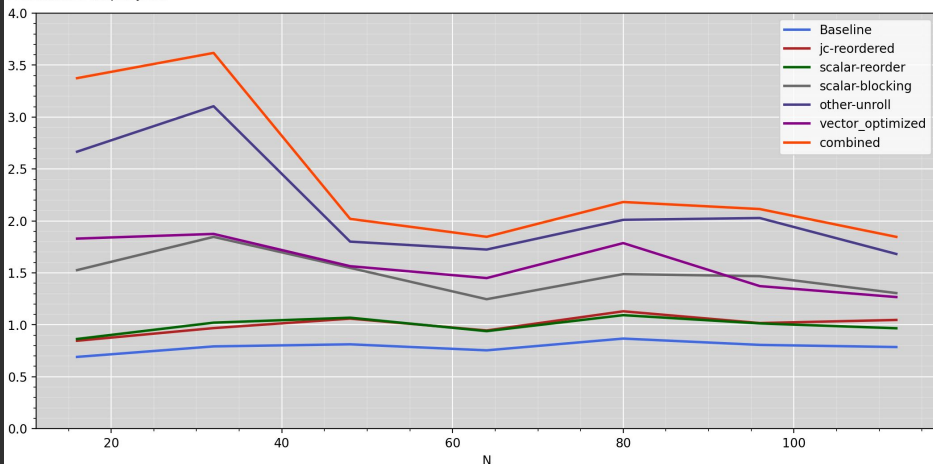


# Different Flags (clang)

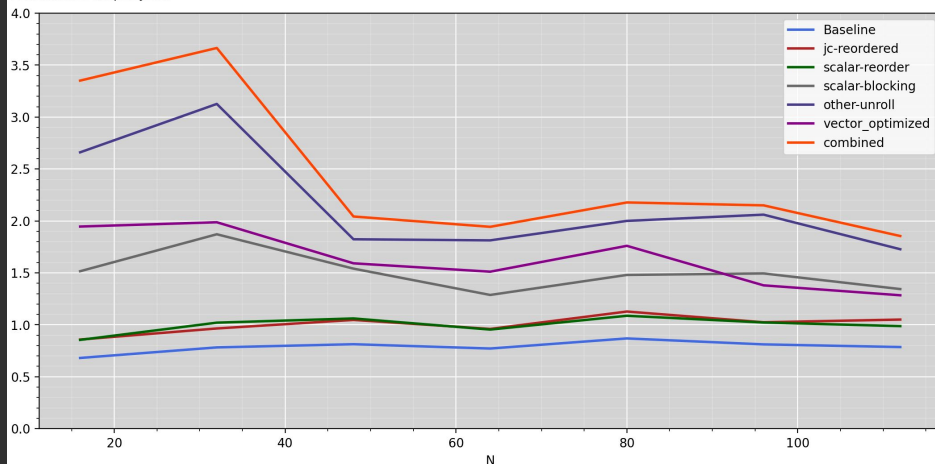
**Performance Plot - (M=16, K=16, T=32)**  
**clang (7.0.1) -O3 -fno-vectorize -ffast-math -mavx2 -mfma -march=native**



**Performance Plot - (M=16, K=16, T=32)**  
**clang (7.0.1) -O3 -ffast-math -mavx2 -mfma -march=native**



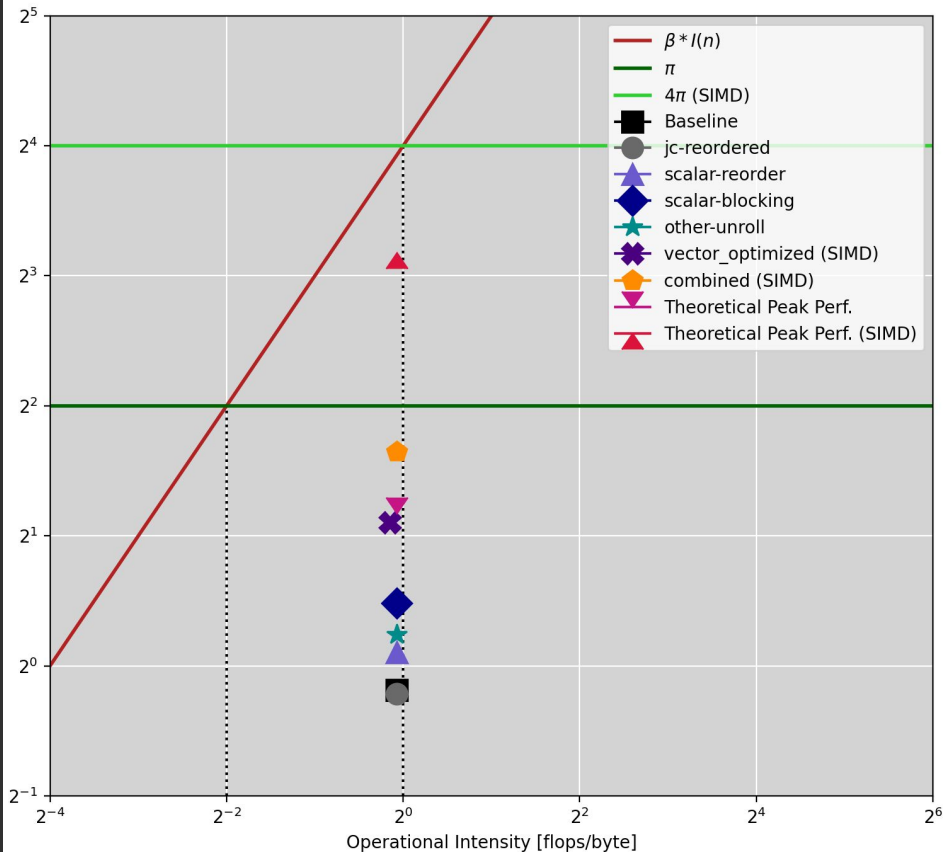
**Performance Plot - (M=16, K=16, T=32)**  
**clang (7.0.1) -O3 -funroll-loops -ffast-math -mavx2 -mfma -march=native**



# Roofline Plot - (K=16,N=16,M=16,T=32)

gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native

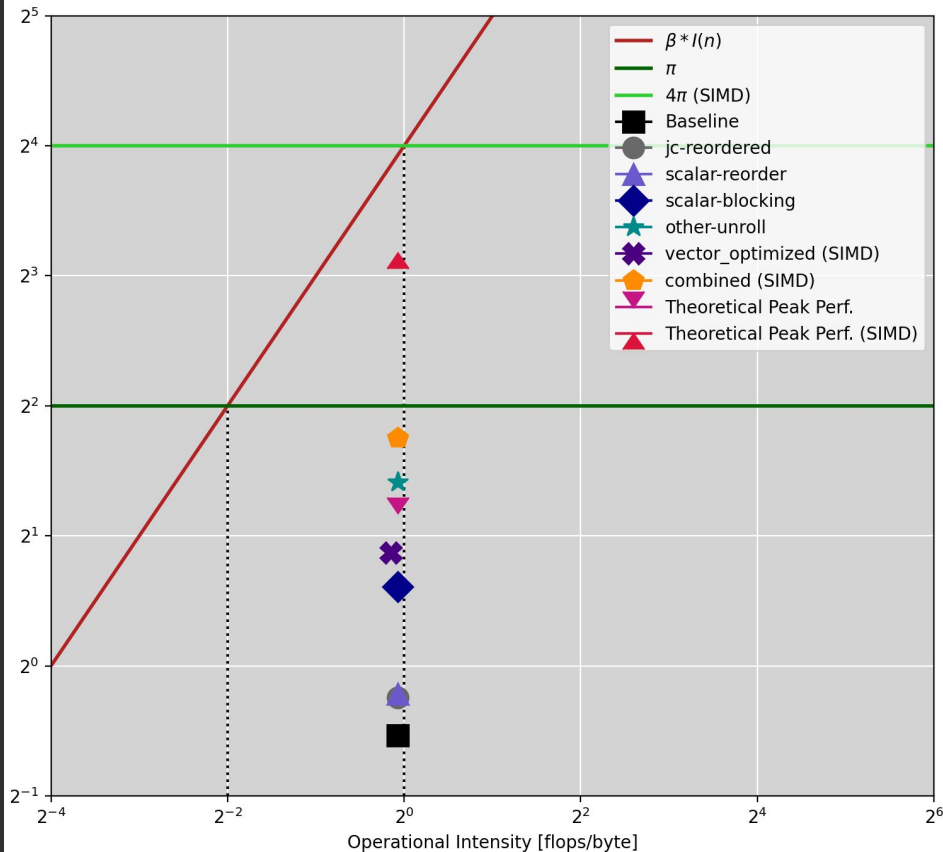
Performance [flops/cycle]



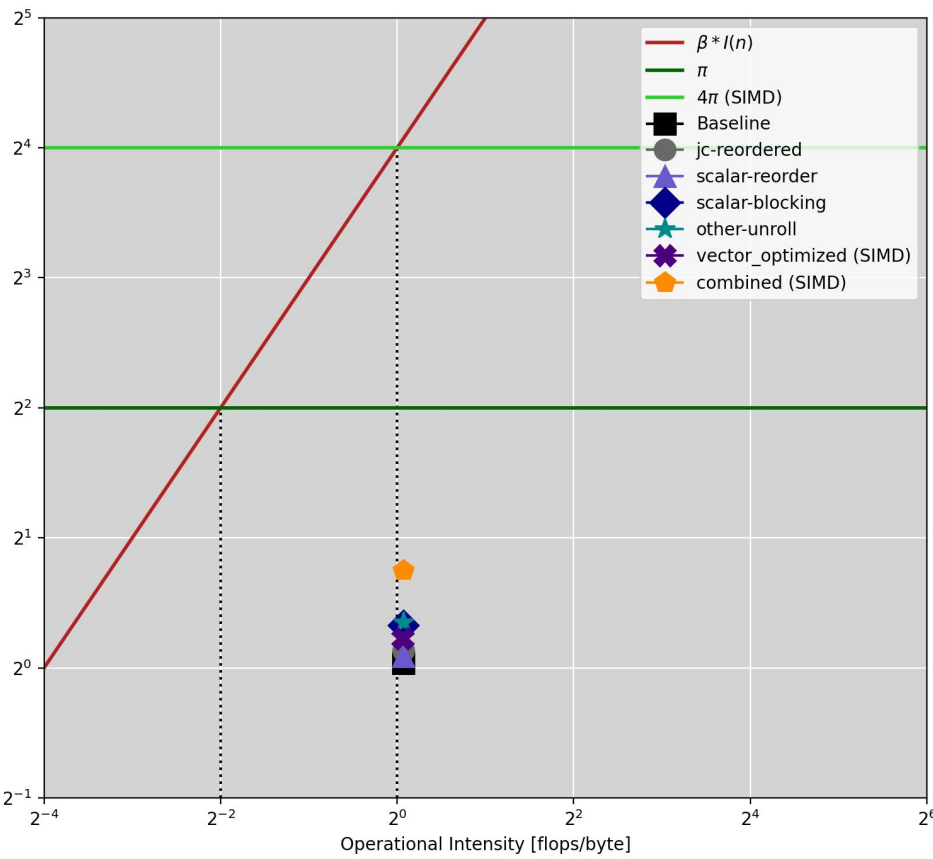
# Roofline Plot - (K=16,N=16,M=16,T=32)

clang (7.0.1) -O3 -ffast-math -mavx2 -mfma -march=native

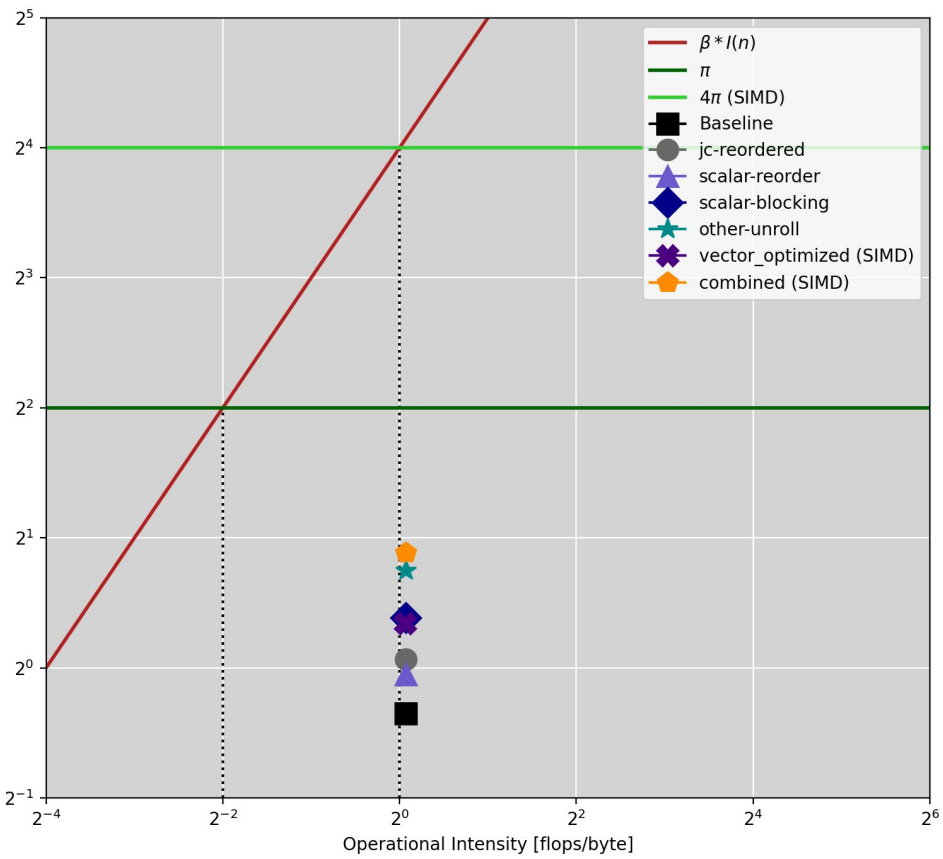
Performance [flops/cycle]



**Roofline Plot - (K=16,N=112,M=16,T=32)**  
**gcc (9.2.0) -O3 -ffast-math -mavx2 -mfma -march=native**  
Performance [flops/cycle]



**Roofline Plot - (K=16,N=112,M=16,T=32)**  
**clang (7.0.1) -O3 -ffast-math -mavx2 -mfma -march=native**  
Performance [flops/cycle]



# Conclusion

- Combination makes the difference.
- Dependant on the compiler and optimization flags, optimizations can differ or even worsen the performance
- Plots as analysis for further optimizations