

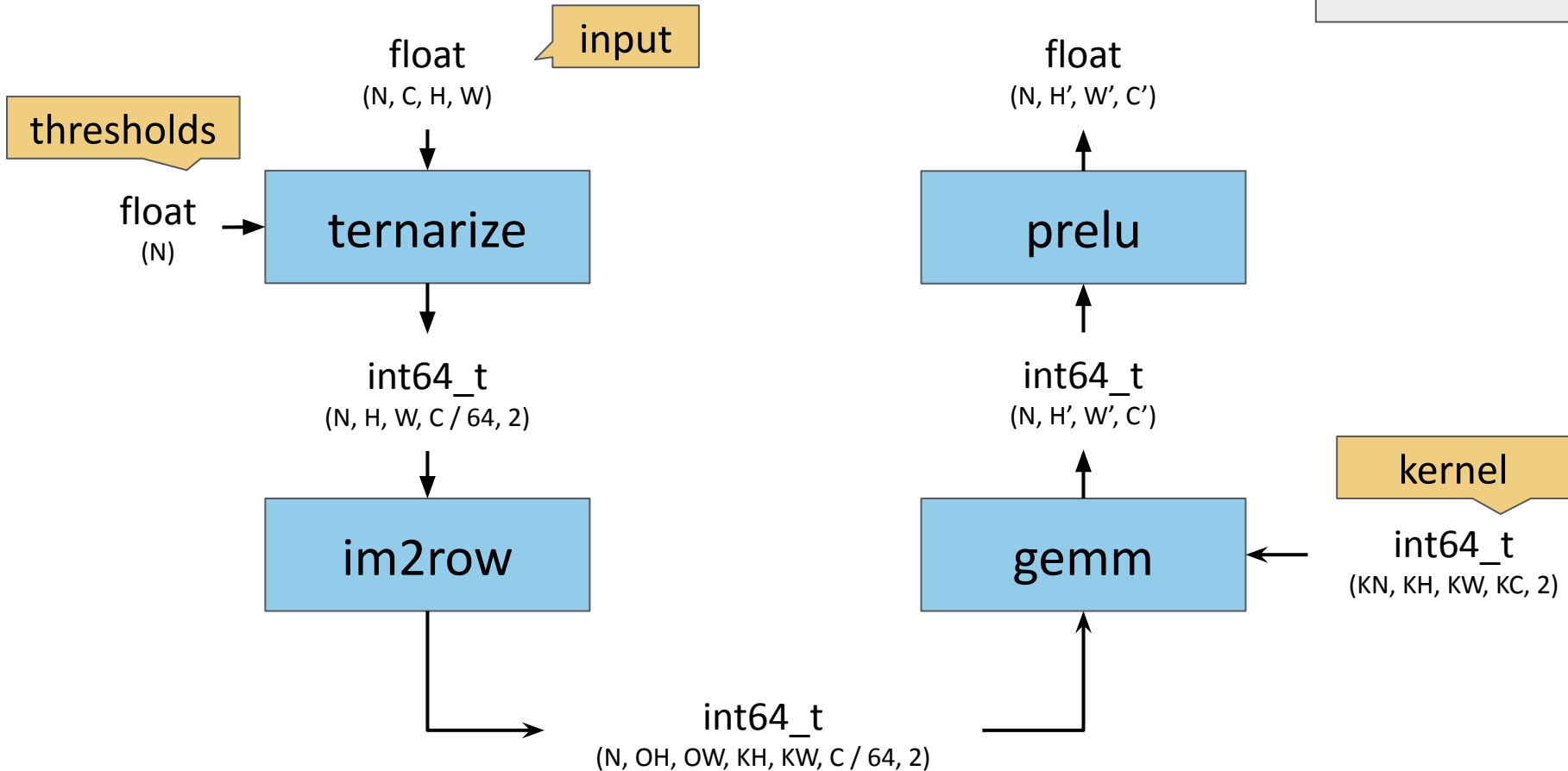
Bitwise Convolution for Ternary Neural Networks

Felix Möller, Daniel Nezamabadi, Rudy Peterson, Luca Tagliavini
Supervisor: Shien Zhu

Advanced Systems Lab - Final Presentation, ETH Zurich
June 7th, 2024

Overview: Algorithm

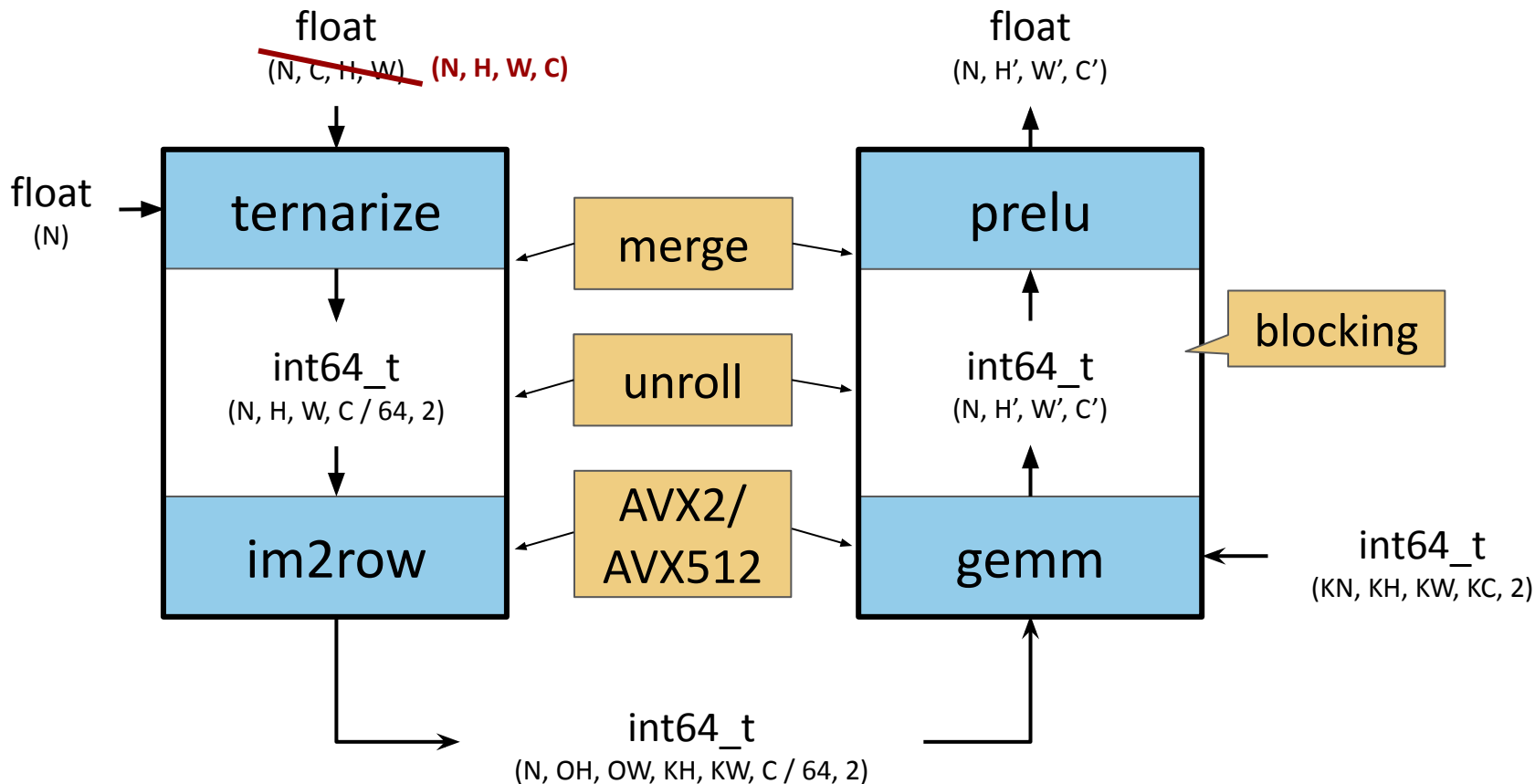
std::vector



Overview: Our Work

Tensor

~~std::vector~~



Using Tensors

```
class Tensor3D {  
    T *data;  
    const size_t dim1;  
    const size_t dim2;  
    const size_t dim3;  
}
```

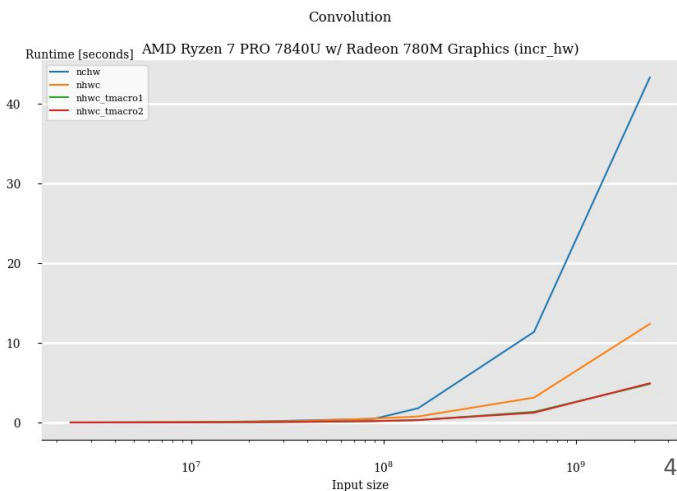
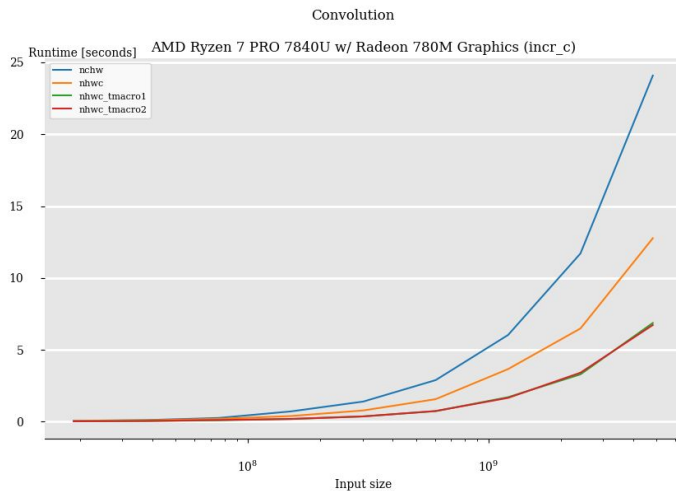
```
#define tensor3d_set(data, dim2, dim3, value, i, j, k) \  
    ((data)[((i) * ((dim2) * (dim3))) + ((j) * (dim3)) + (k)] = (value))
```

```
#define tensor3d_set(data, dim2, dim3, value, i, j, k) \  
    ((data)[((i) * (dim2) + (j)) * (dim3) + (k)] = (value))
```

nhwc_tmacro1

nhwc_tmacro2

gcc with O3
and disabled
vectorization

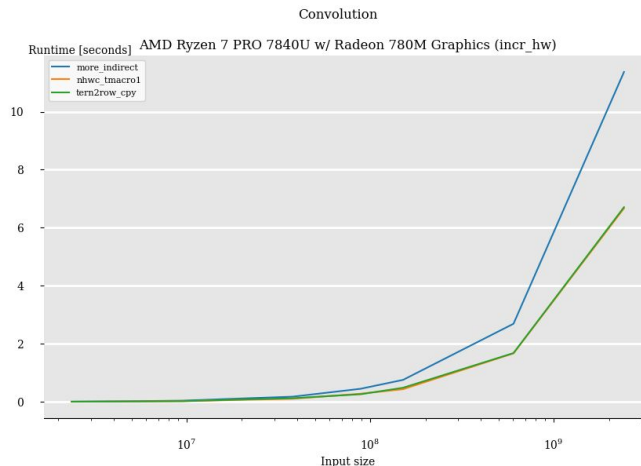
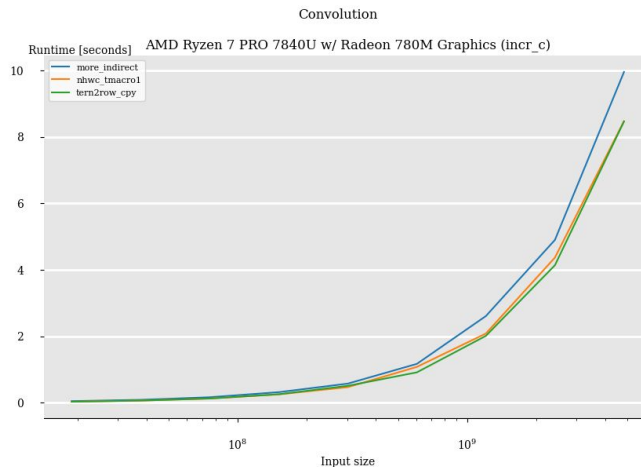


Merging im2row

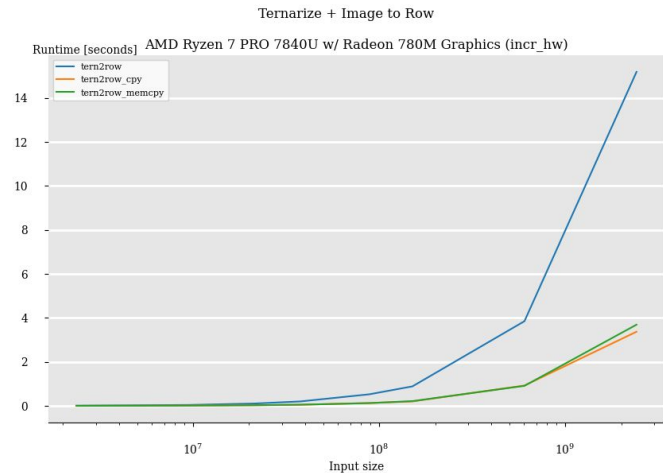
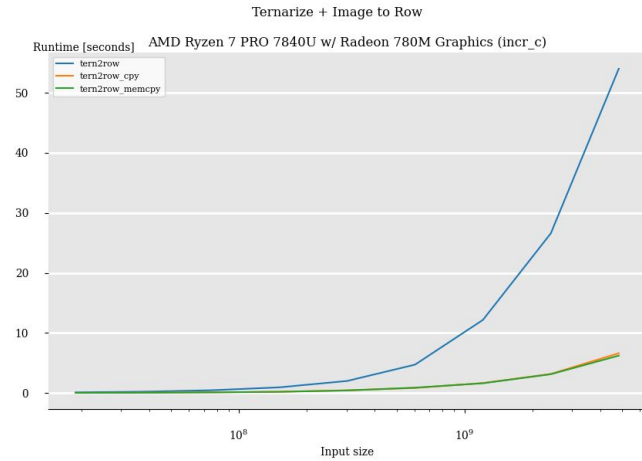
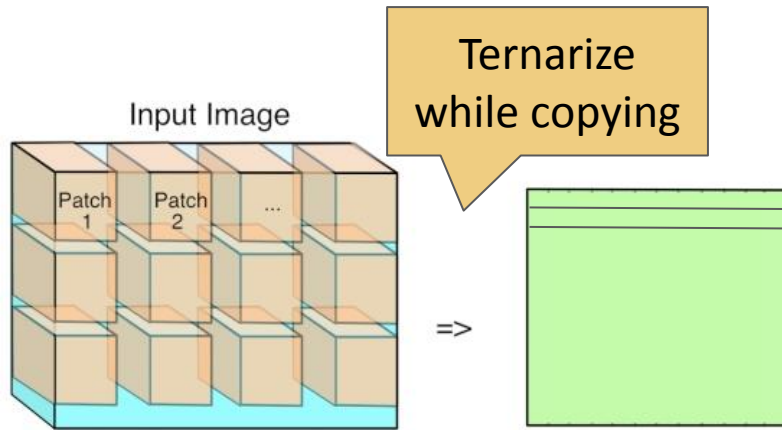
Don't merge im2row

Merge ternarize and
im2row

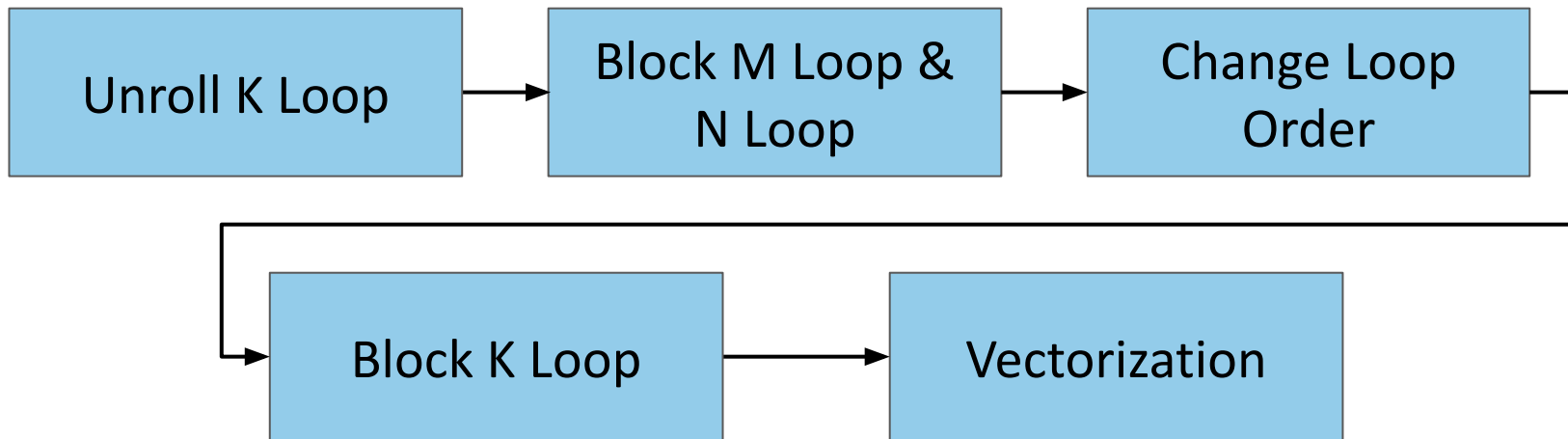
Merge im2row and gemm
(indirect)



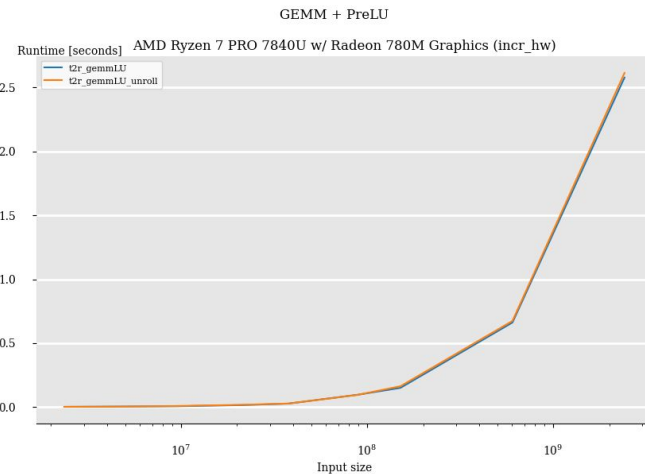
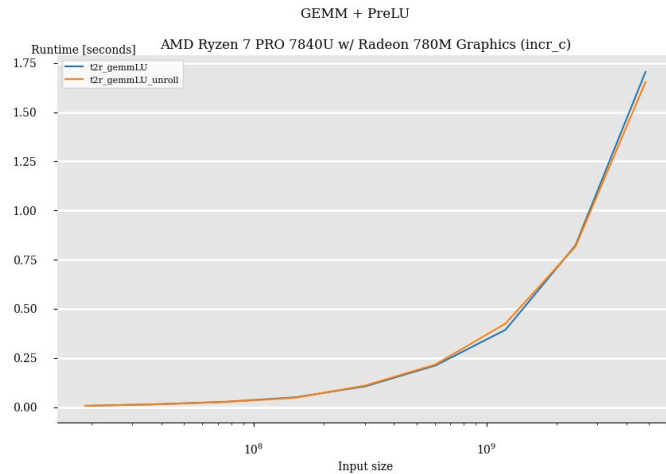
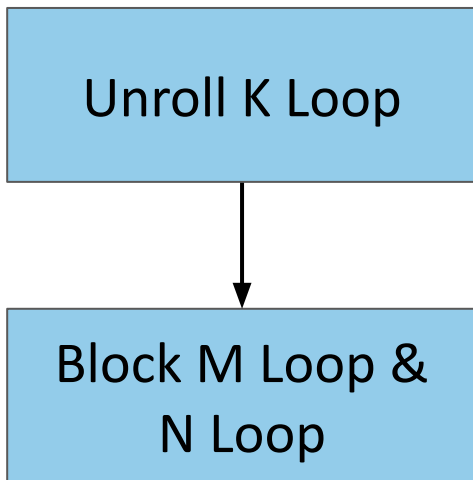
ternarize + im2row



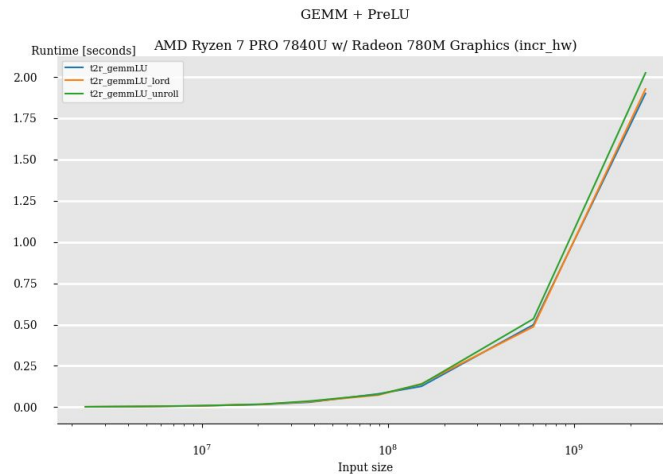
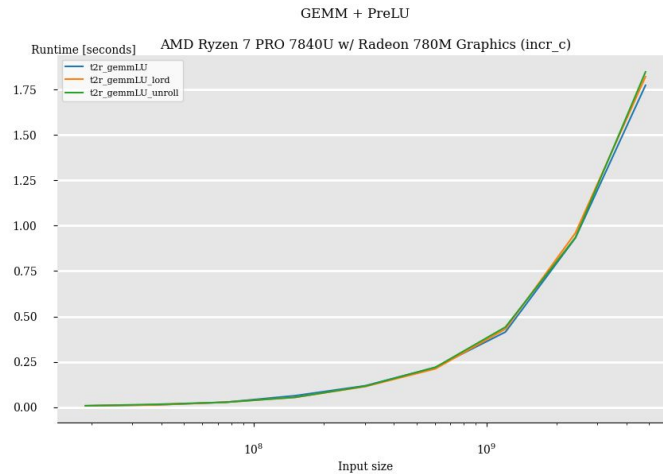
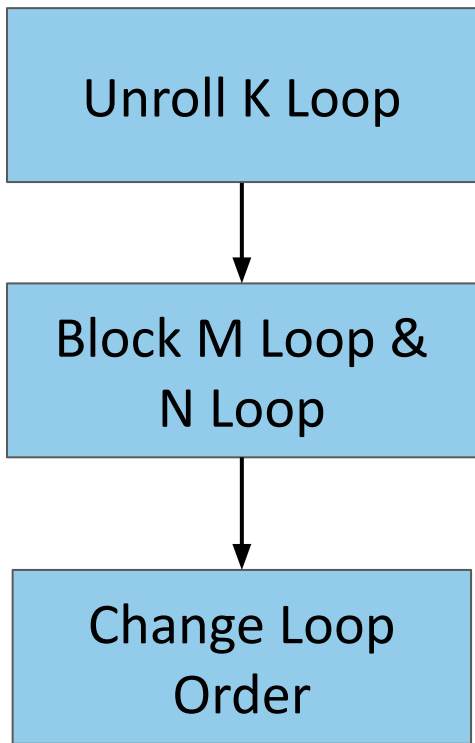
GEMM + PReLU - Optimization Plan



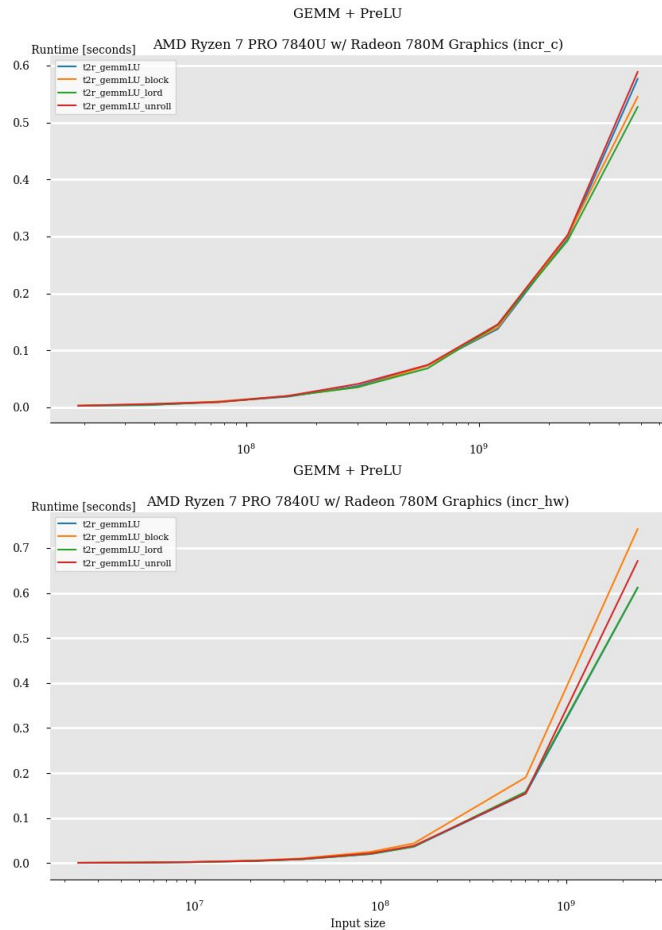
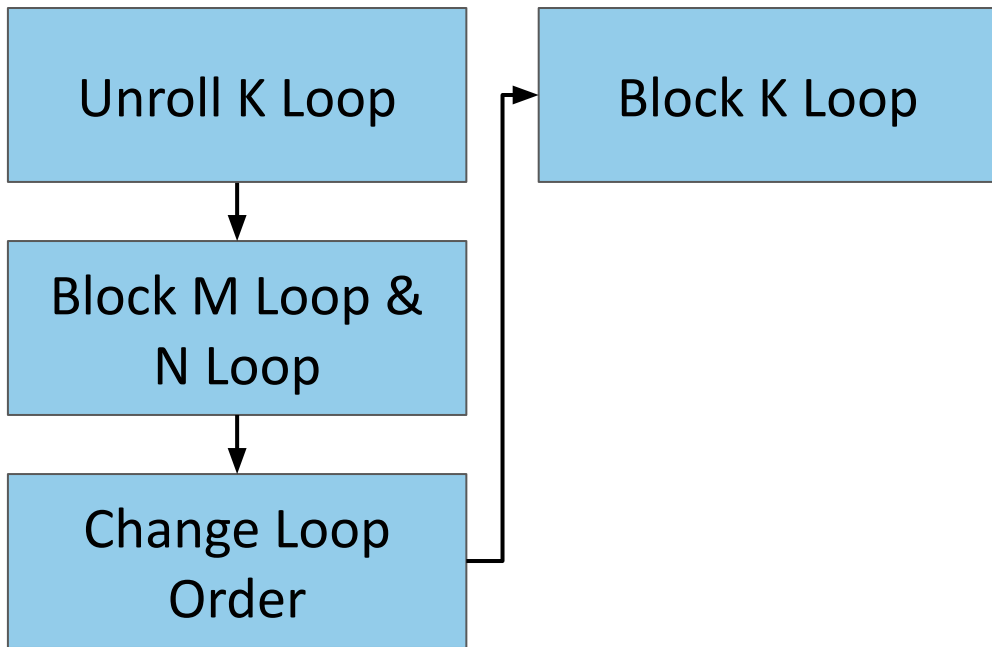
GEMM + PReLU - Unrolling over K & Blocking over M and N



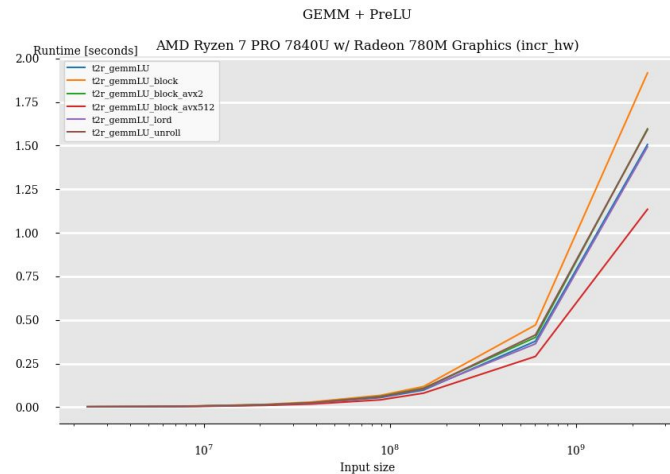
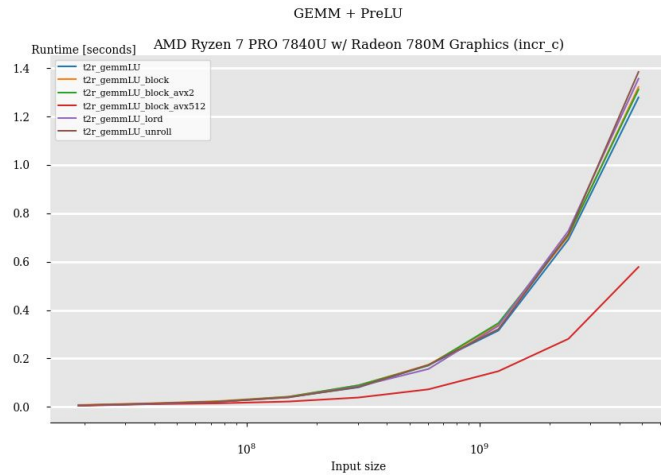
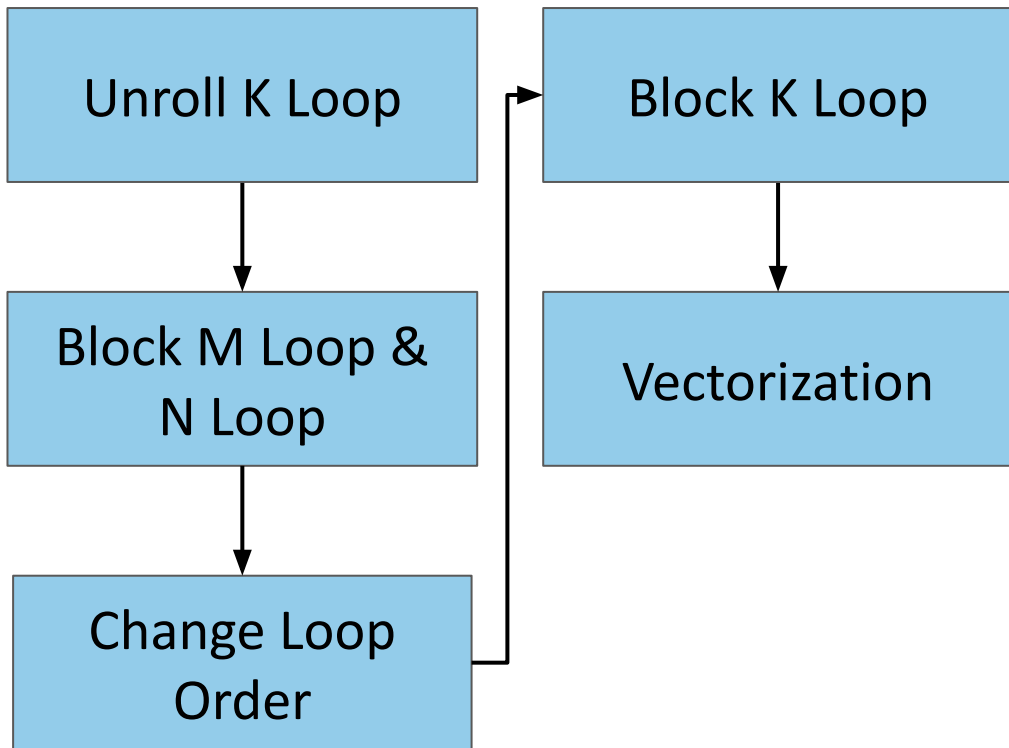
GEMM + PReLU - Changing the Loop Order



GEMM + PReLU - Blocking over K



GEMM + PReLU - Vectorization with AVX2 and AVX512

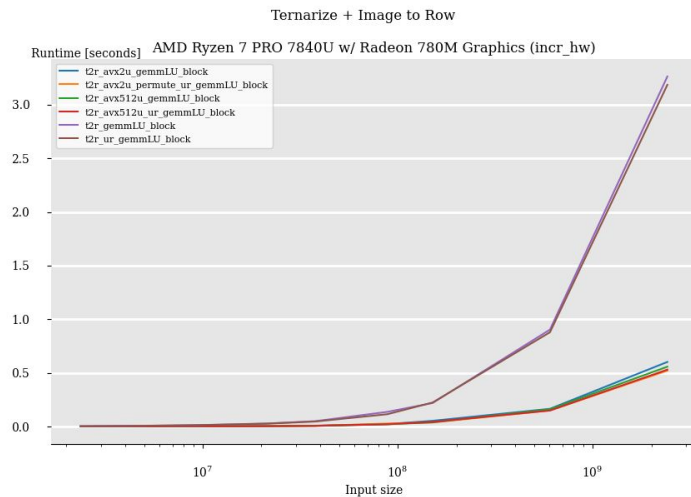
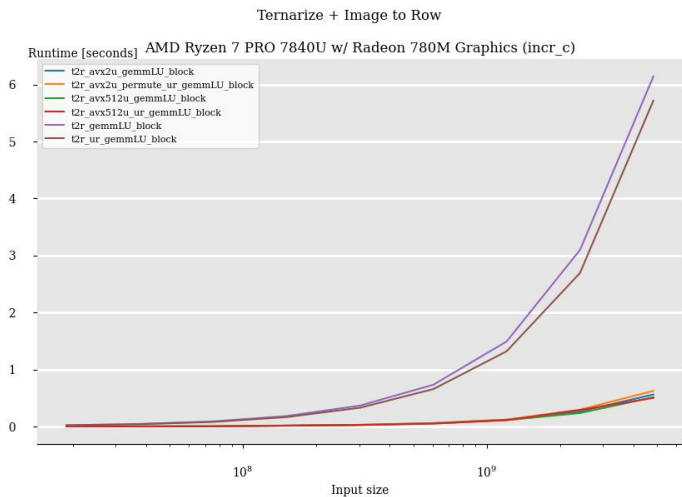


Optimizing Ternarize + im2row

Loop Unrolling

AVX2

AVX512

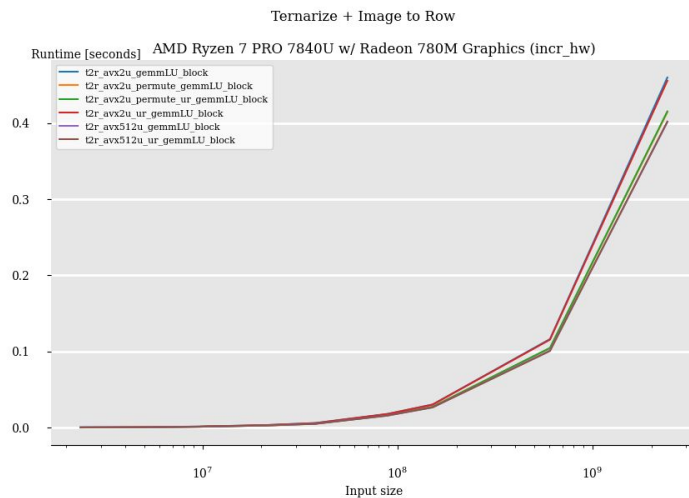
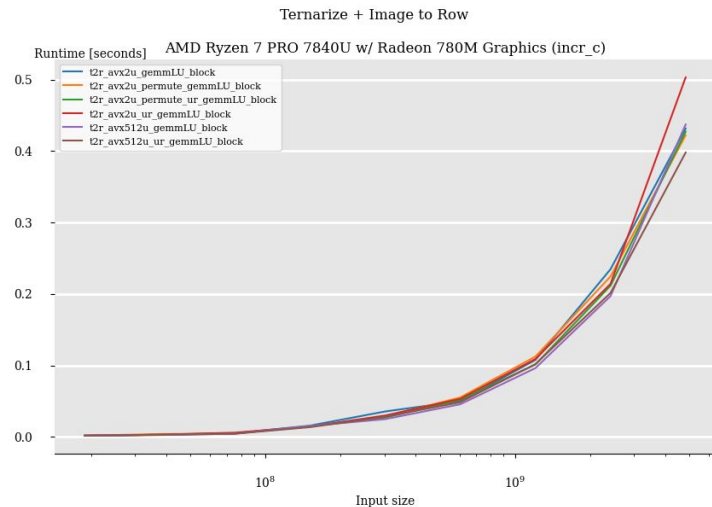


Optimizing Ternarize + im2row

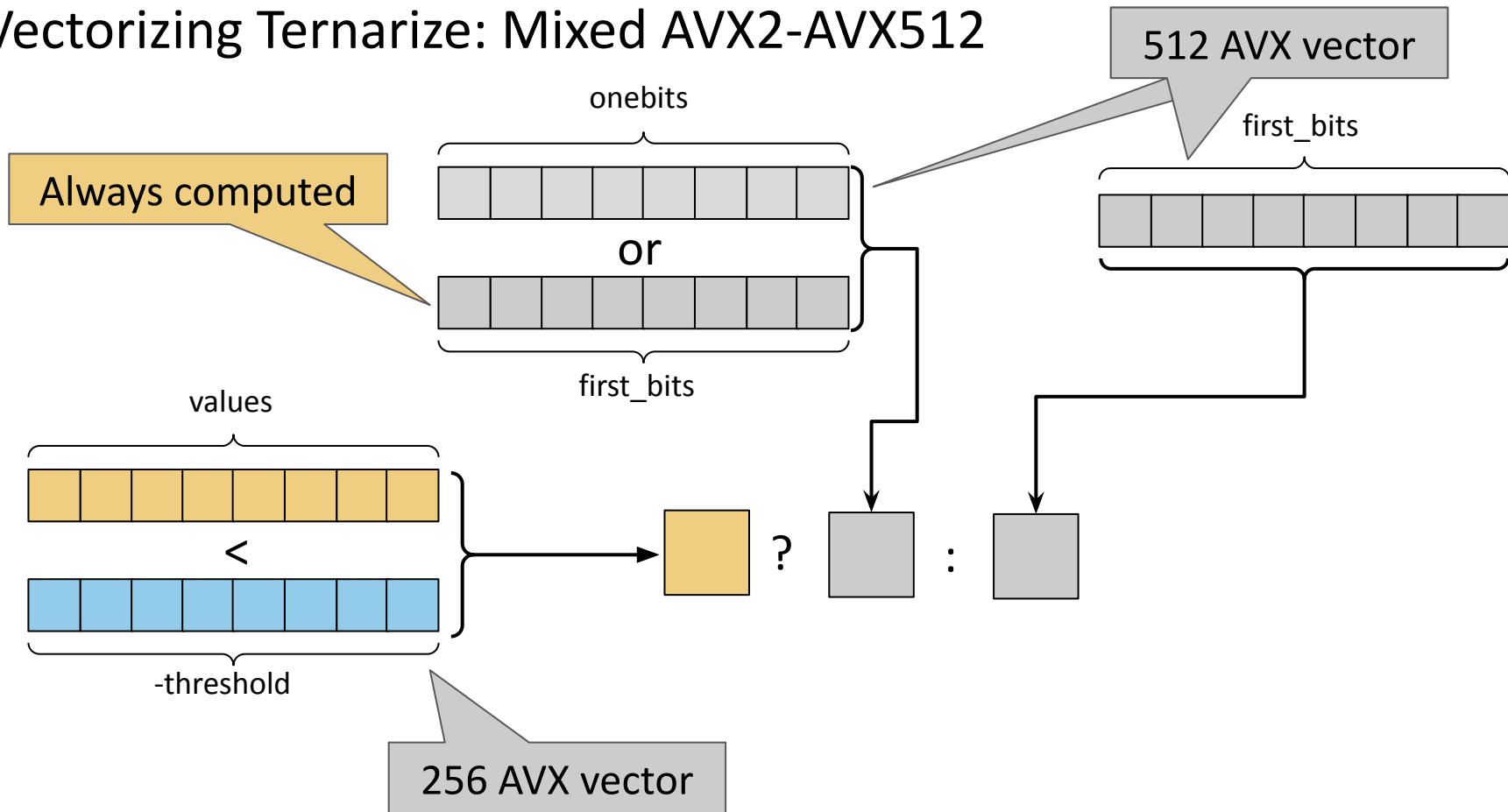
Loop Unrolling

AVX2

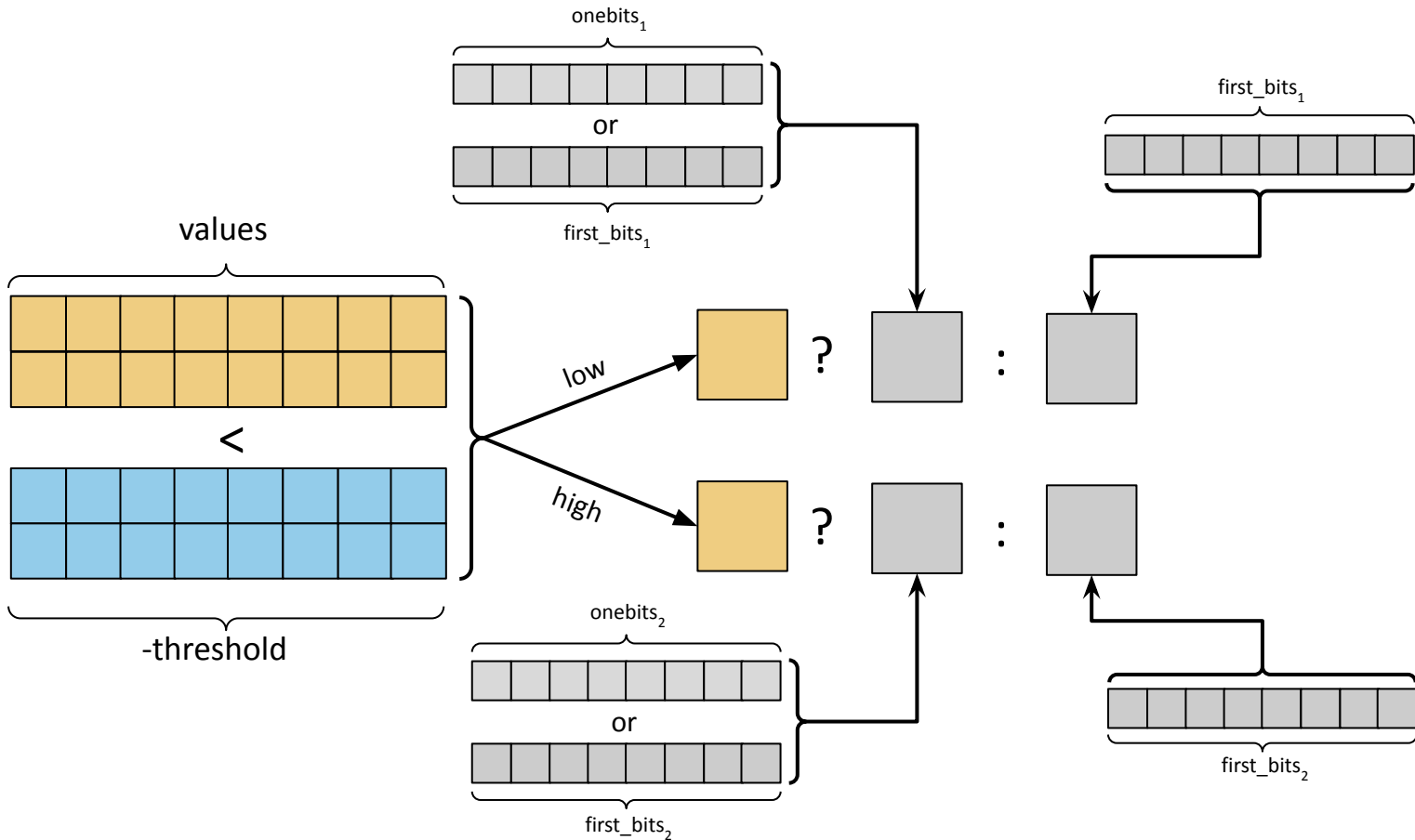
AVX512



Vectorizing Ternarize: Mixed AVX2-AVX512

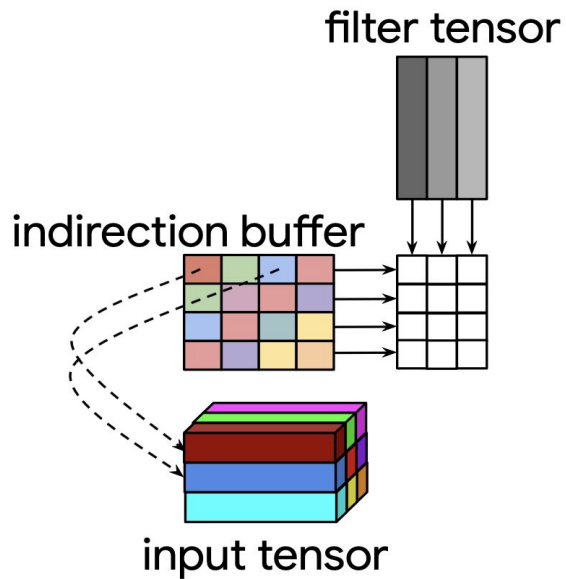


Vectorizing Ternarize: AVX512 only

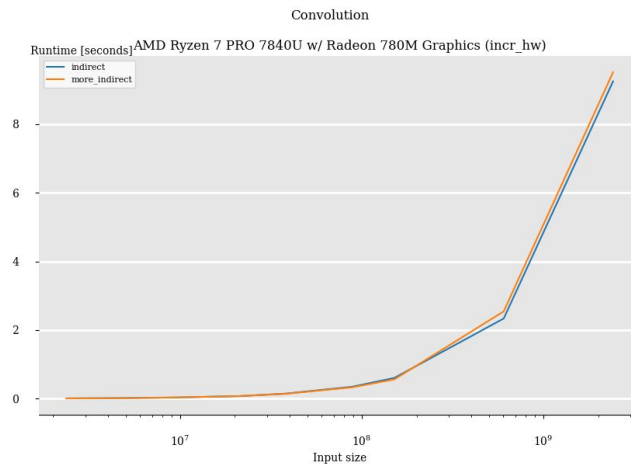
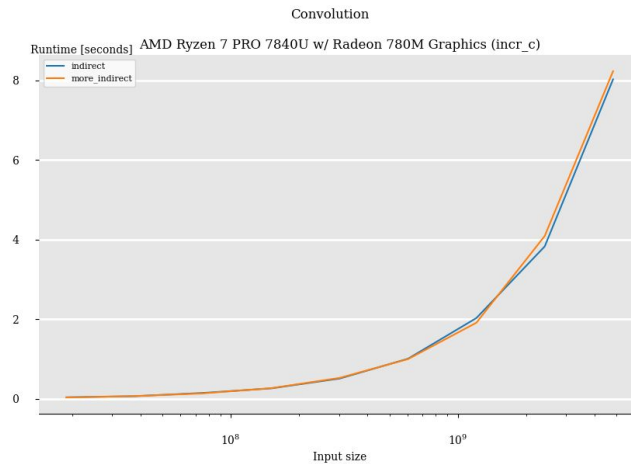


im2row + GEMM

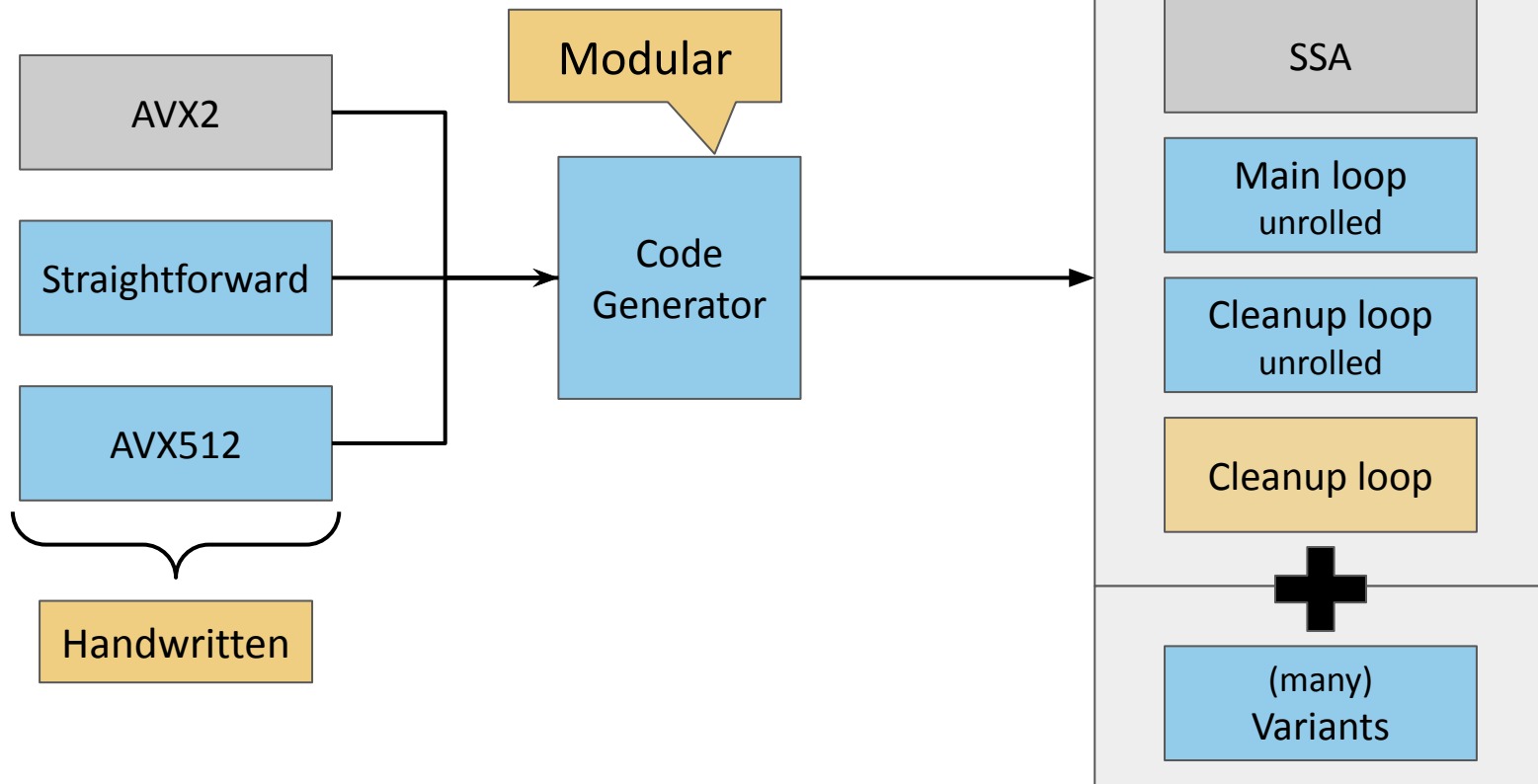
Indirect Convolution



Source: M.Dukhan, The Indirect Convolution Algorithm, *arxiv* 2019



Code Generation



Code Generation - Variant: Autotuning + libpopcnt

TNN GEMM

```
for m in range(0, M):  
    for n in range(0, N):  
  
        cntp1 = cntp2 = 0  
        for k in range(0, K, 2):  
            ...  
            cntp1 += popcnt(p2)  
            cntp2 += popcnt(p1 & p2)  
  
        output[m,n] = ..
```

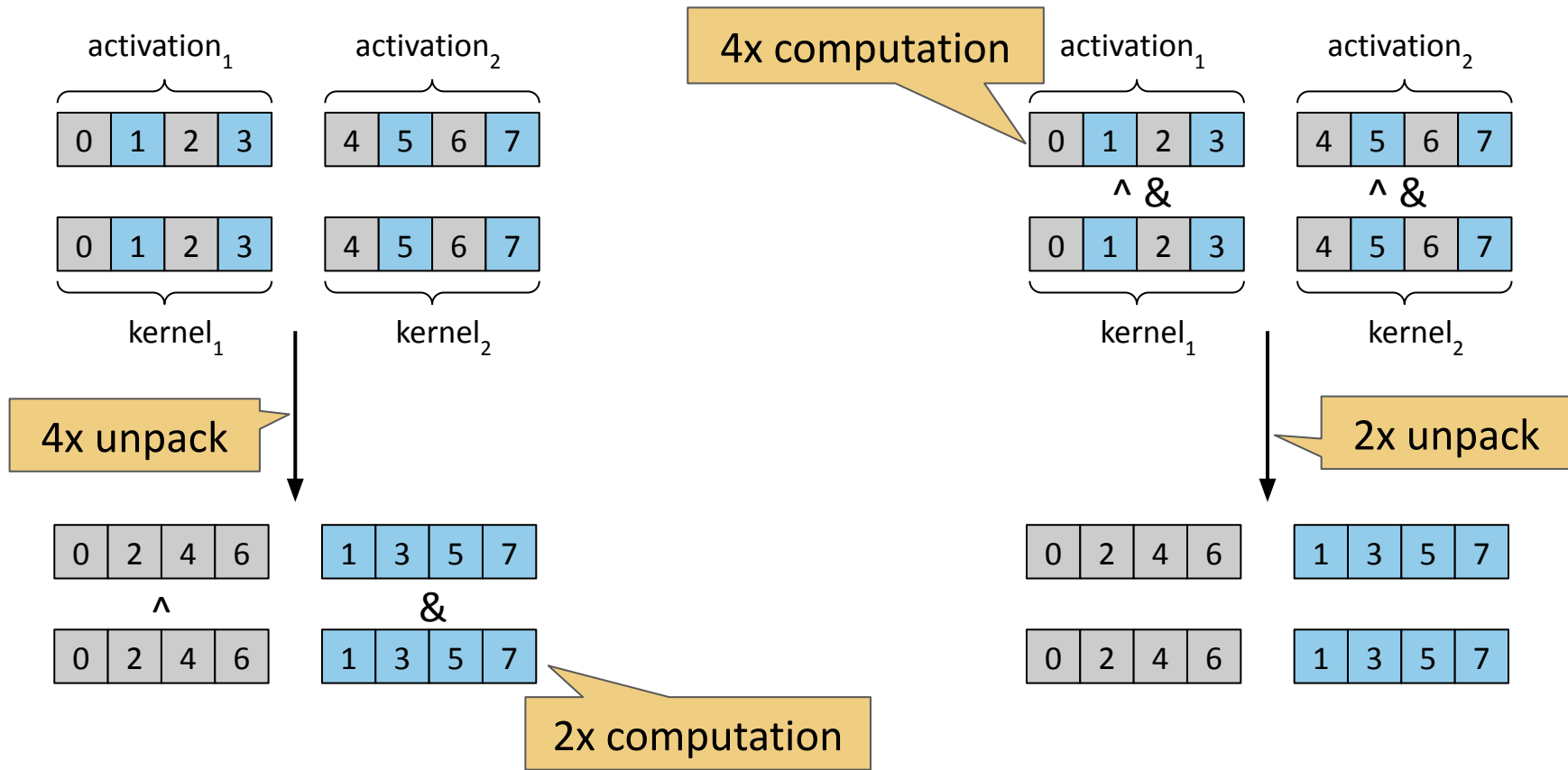
Autotuning on M and N

popcount on a **big** vector

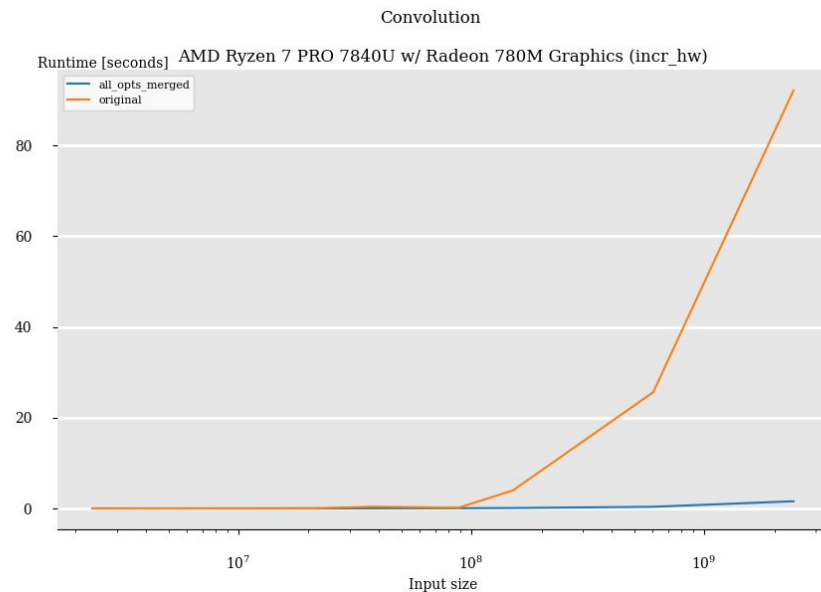
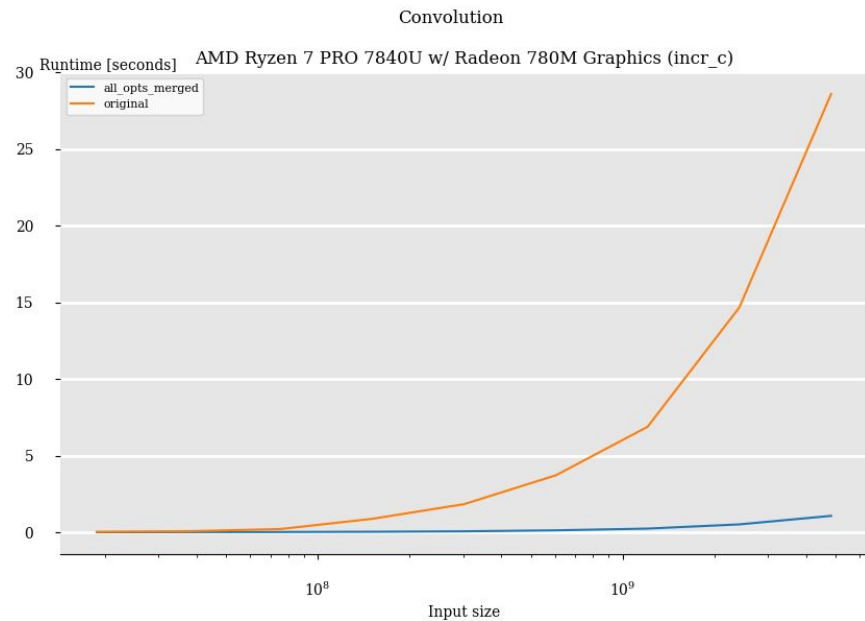
TNN GEMM

```
for m in range(0, M):  
    for n in range(0, N):  
  
        vcntp1[], vcntp2[]  
        for k in range(0, K, 2):  
            ...  
            vcntp1[k/2] = p2  
            vcntp2[k/2+1] = p1 & p2  
  
        cntp1 = libpopcnt(vcntp1)  
        ..
```

Code Generation - Variant: Less Unpacking



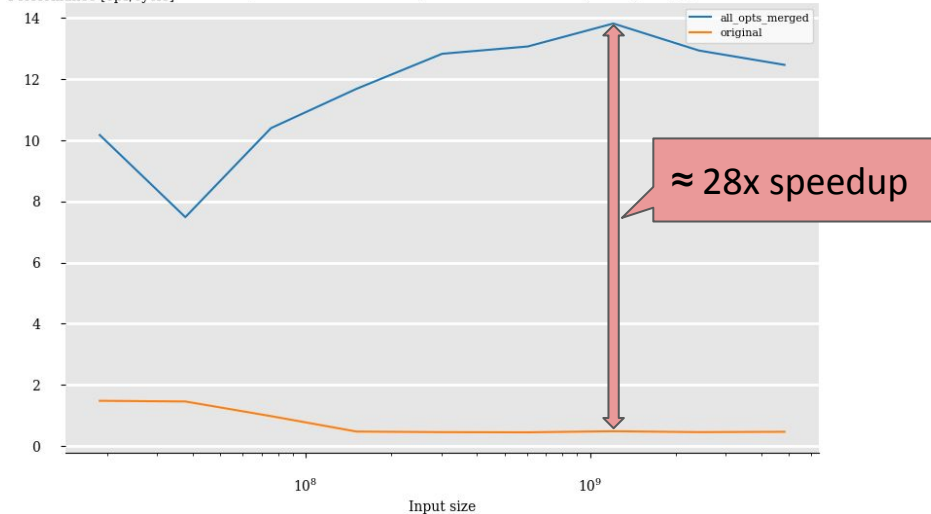
Comparing to the Original



Comparing to the Original

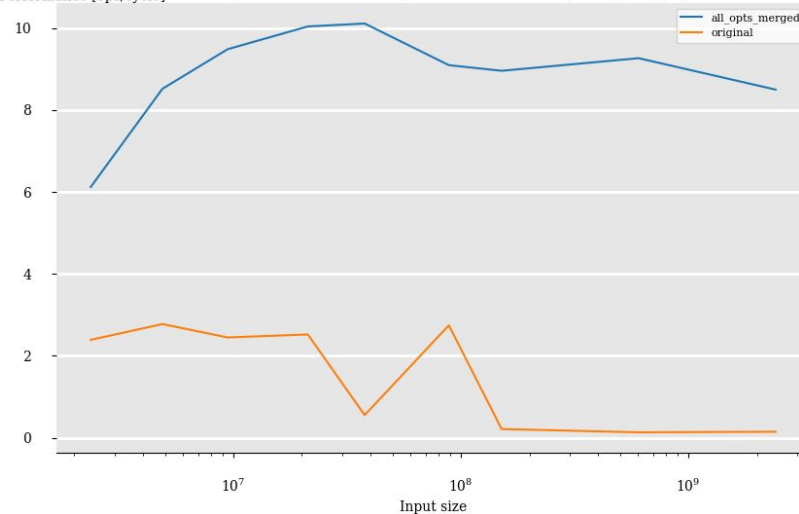
Convolution

Performance [ops/cycle] AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)



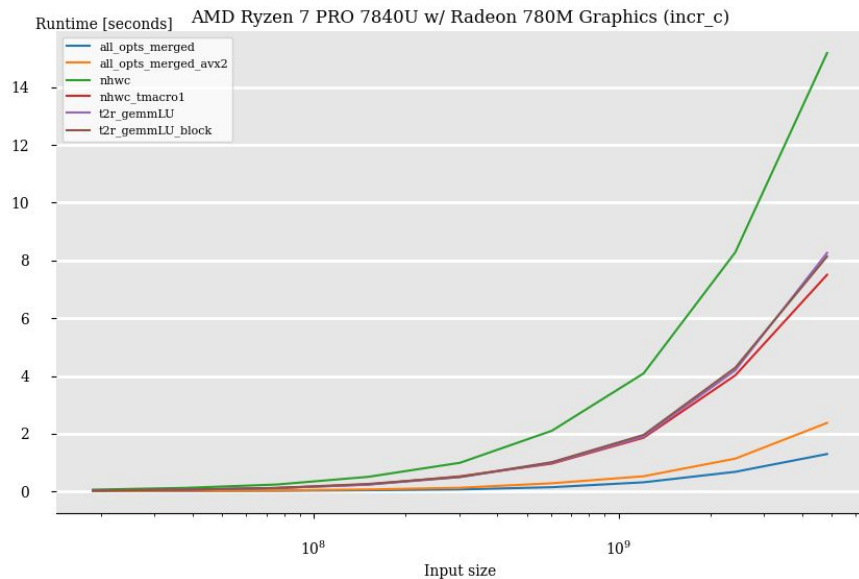
Convolution

Performance [ops/cycle] AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)

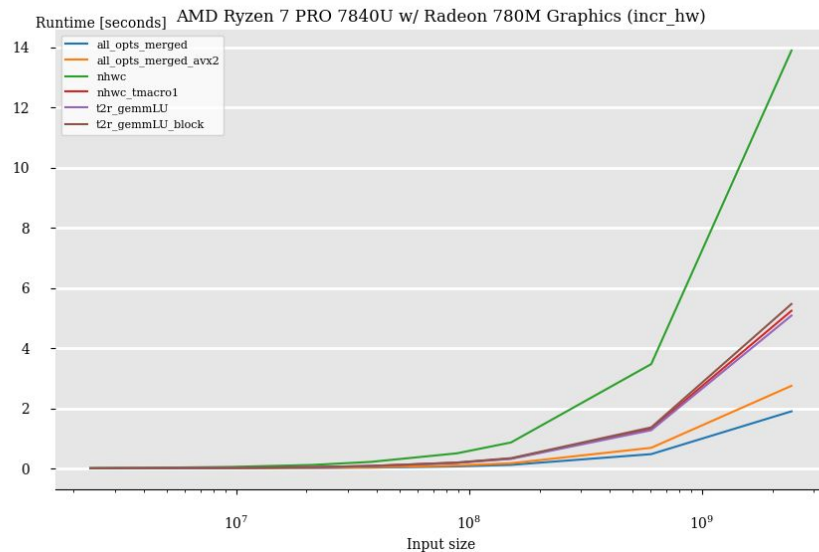


Result: Optimization Steps

Convolution



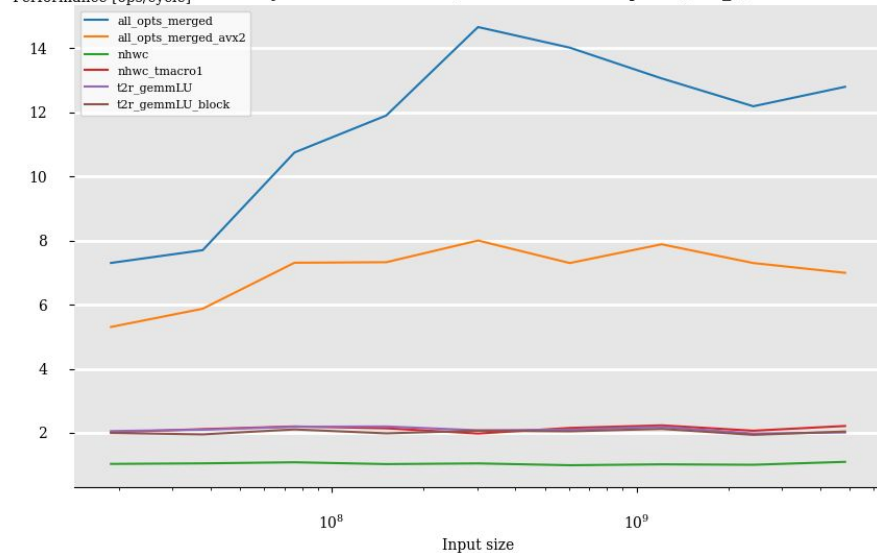
Convolution



Result: Optimization Steps

Convolution

Performance [ops/cycle] AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)



Convolution

Performance [ops/cycle] AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)

