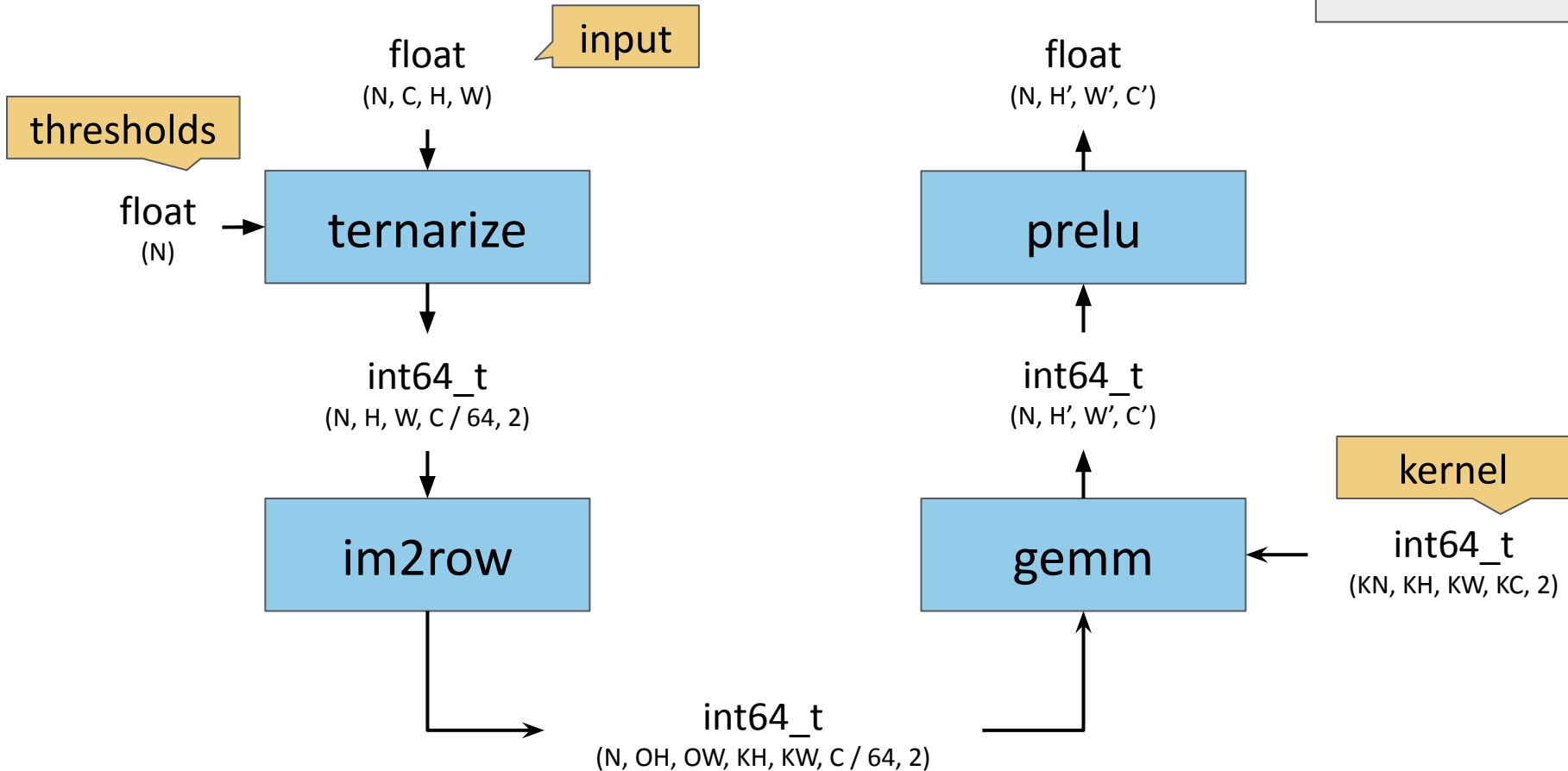# Bitwise Convolution for Ternary Neural Networks

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

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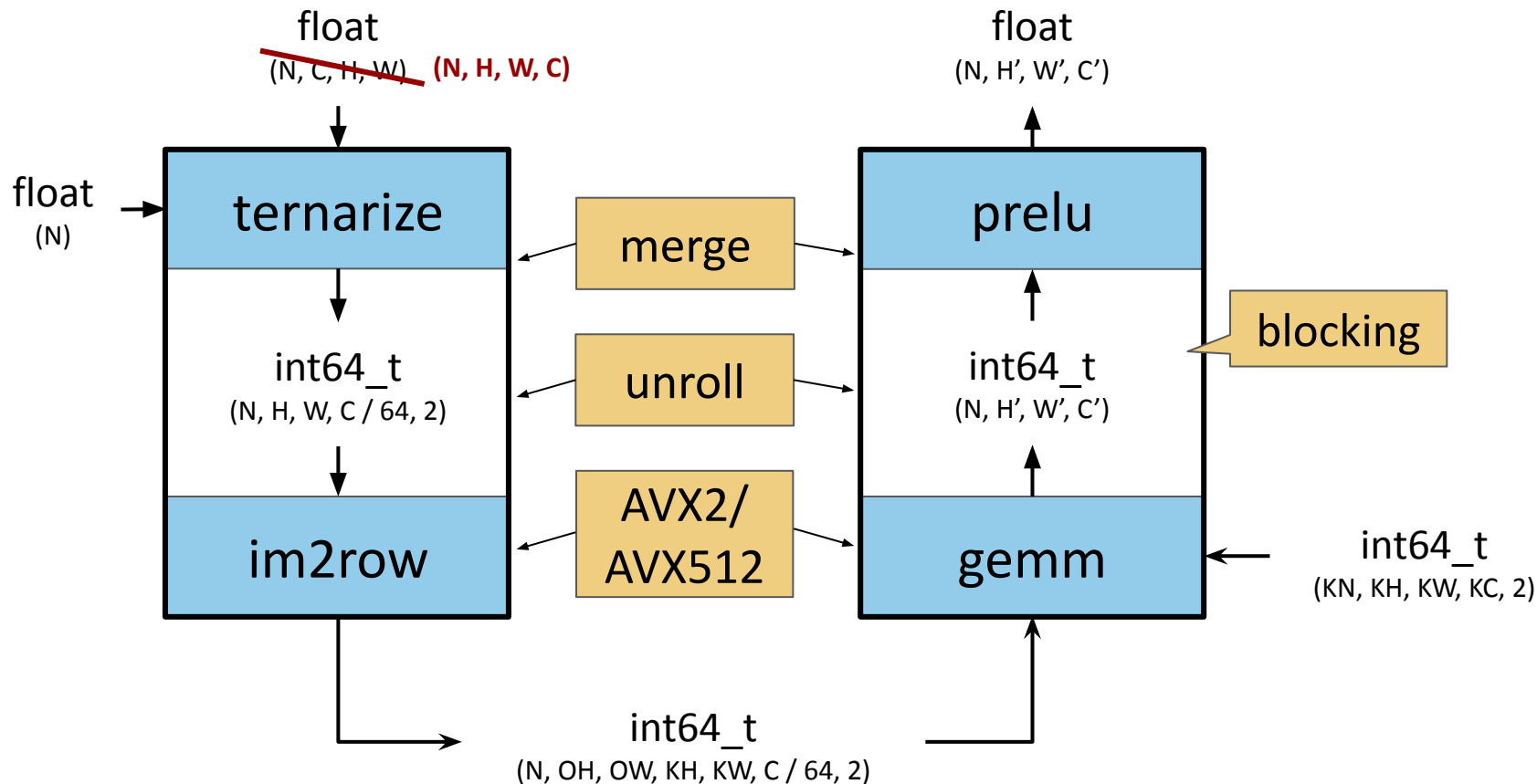June 7th, 2024

# Overview: Algorithm

std::vector

float
(N, C, H, W)

input

thresholds

float
(N)

→ **ternarize**

↓

int64_t
(N, H, W, C / 64, 2)

↓

**im2row**

↓

int64_t
(N, OH, OW, KH, KW, C / 64, 2)

float
(N, H', W', C')

↑

**prelu**

↑

int64_t
(N, H', W', C')

↑

**gemm** ← kernel

int64_t
(KN, KH, KW, KC, 2)

# Overview: Our Work

**Tensor** ~~std::vector~~

~~float~~
~~(N, C, H, W)~~ **(N, H, W, C)**

float
(N, H', W', C')

float
(N)

| ternarize | | prelu |

int64_t
(N, H, W, C / 64, 2)

int64_t
(N, H', W', C')

merge

unroll

blocking

| im2row | AVX2/AVX512 | gemm |

int64_t
(KN, KH, KW, KC, 2)

int64_t
(N, OH, OW, KH, KW, C / 64, 2)

# Using Tensors

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

nhwc_tmacro1

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

nhwc_tmacro2

```
#define tensor3d_set(data, dim2, dim3, value, i, j, k)          \
  ((data)[((i) * (dim2) + (j)) * (dim3) + (k)] = (value))
```
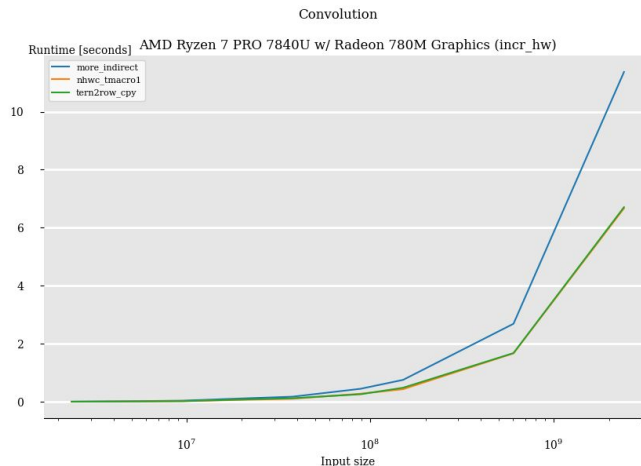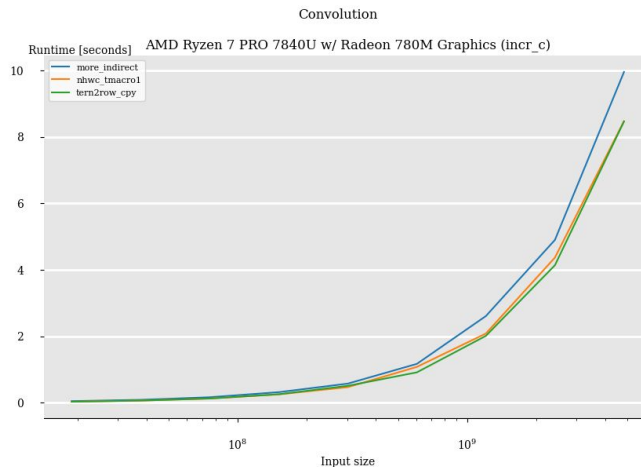
gcc with O3 and disabled vectorization



Convolution
Runtime [seconds]  AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)



Convolution
Runtime [seconds]  AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)
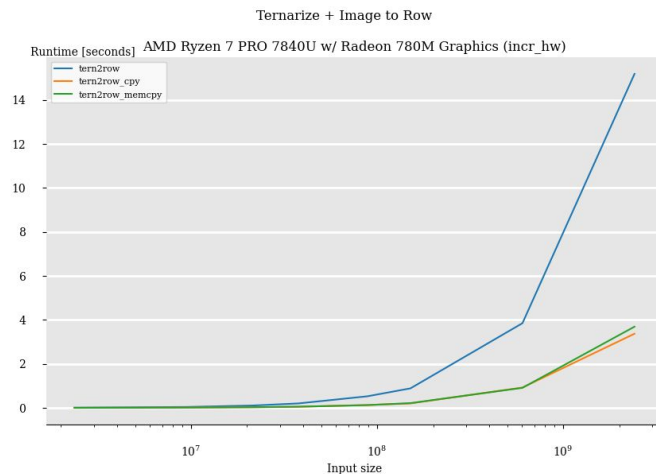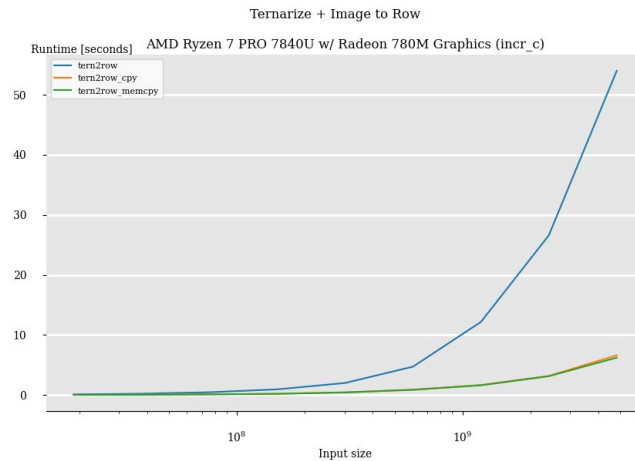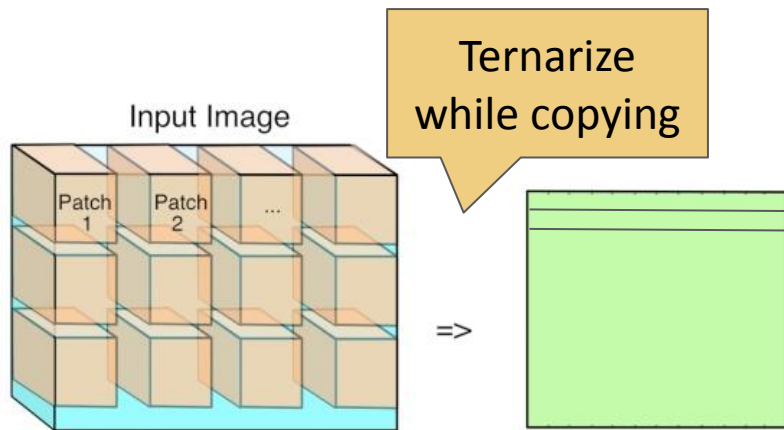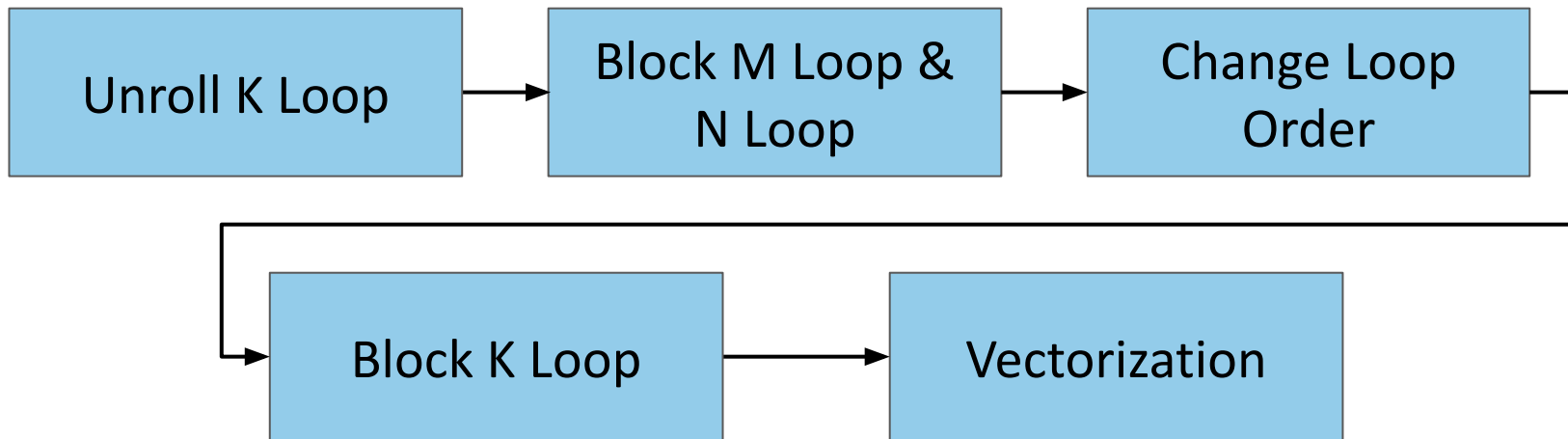
4

# Merging im2row

Don't merge im2row

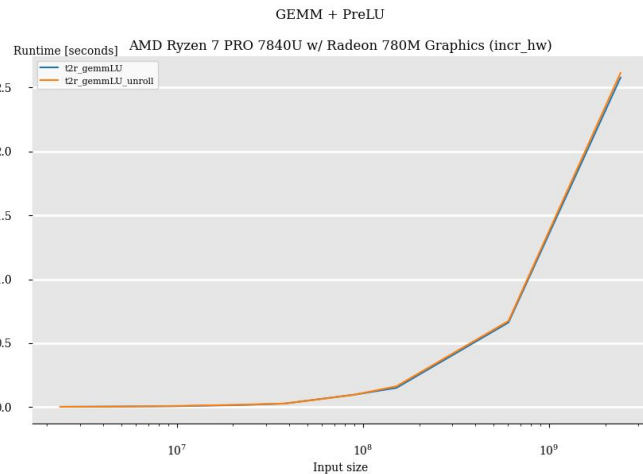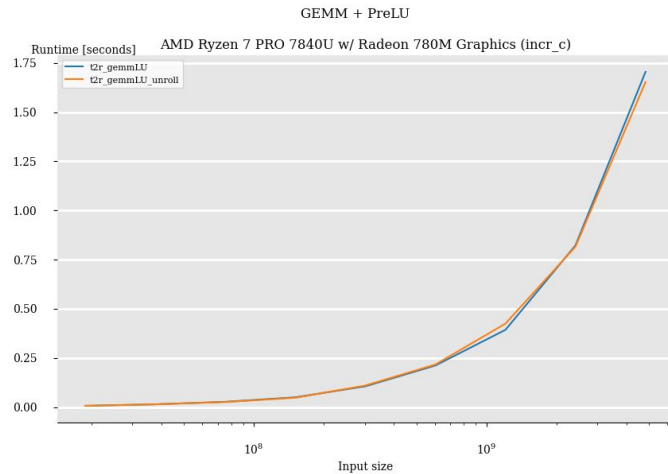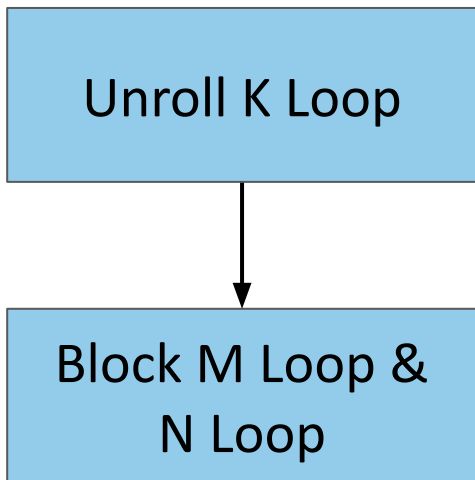Merge ternarize and im2row

Merge im2row and gemm (indirect)



Convolution
Runtime [seconds]    AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)

Legend:
more_indirect
nhwc_tmacro1
tern2row_cpy

Input size



Convolution
Runtime [seconds]    AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)

Legend:
more_indirect
nhwc_tmacro1
tern2row_cpy

Input size

# ternarize + im2row



Ternarize while copying

Input Image

Patch 1  Patch 2  ...

=>

Ternarize + Image to Row

Runtime [seconds]   AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)

- tern2row
- tern2row_cpy
- tern2row_memcpy

Input size

Ternarize + Image to Row

Runtime [seconds]   AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)

- tern2row
- tern2row_cpy
- tern2row_memcpy

Input size

Image adapted from https://leonardoaraujosantos.gitbook.io/artificial-inteligence/machine_learning/deep_learning/convolution_layer/making_faster

# GEMM + PReLU - Optimization Plan

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

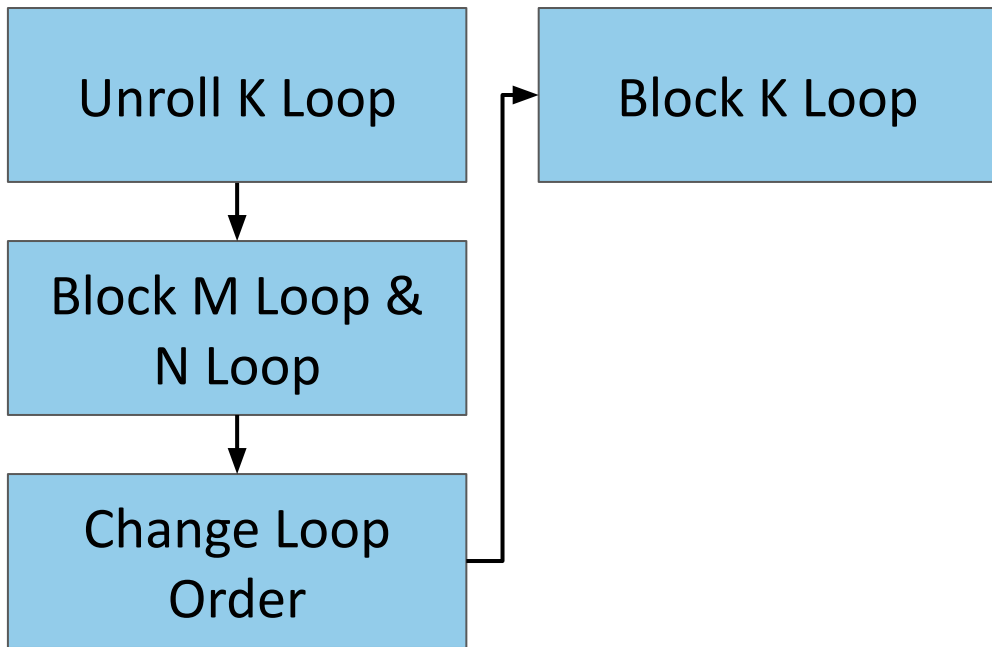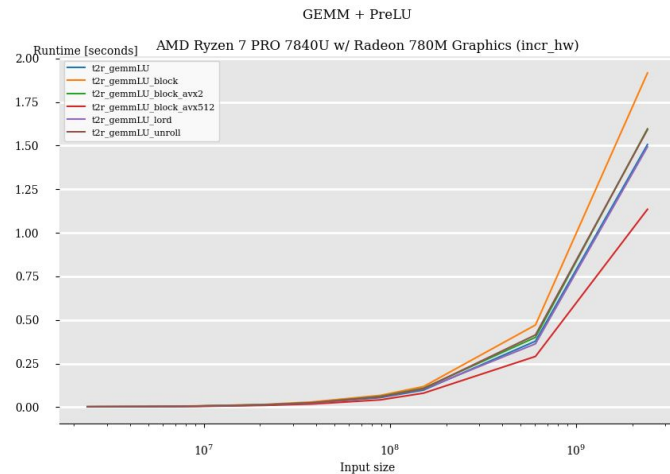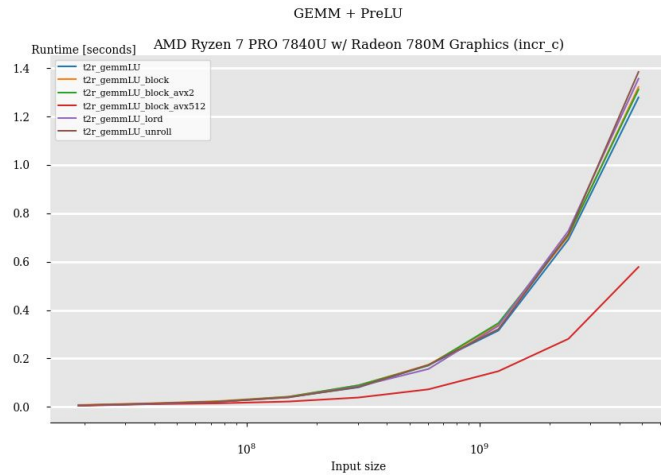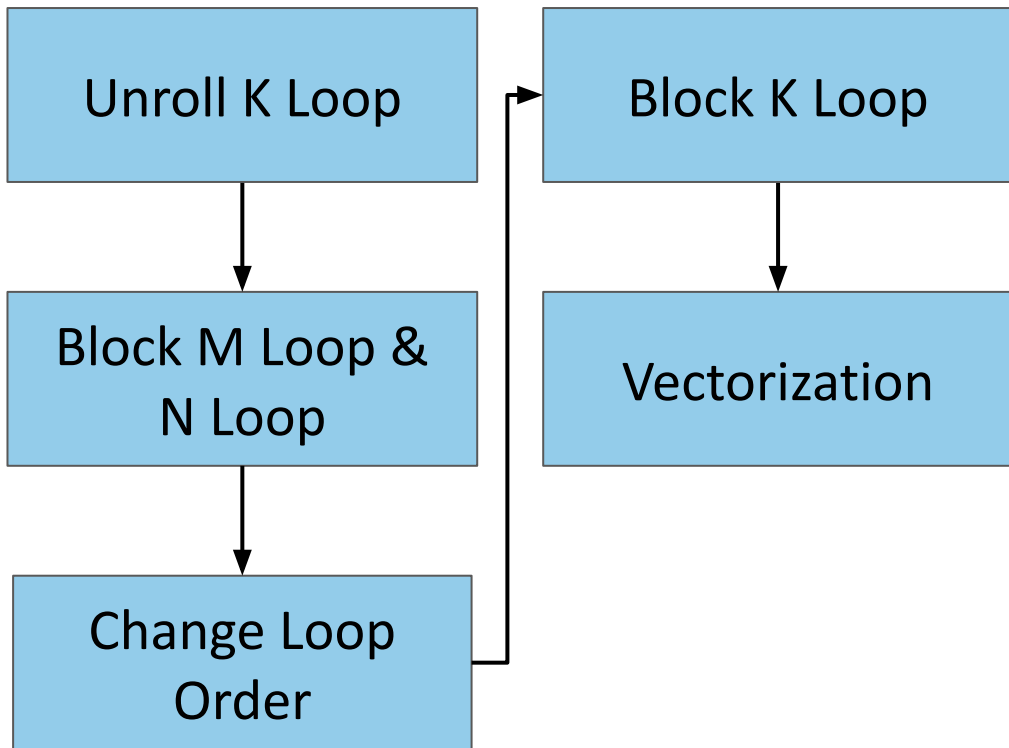Unroll K Loop

Block M Loop & N Loop

# GEMM + PReLU - Changing the Loop Order

# GEMM + PReLU - Blocking over K



10

# GEMM + PReLU - Vectorization with AVX2 and AVX512

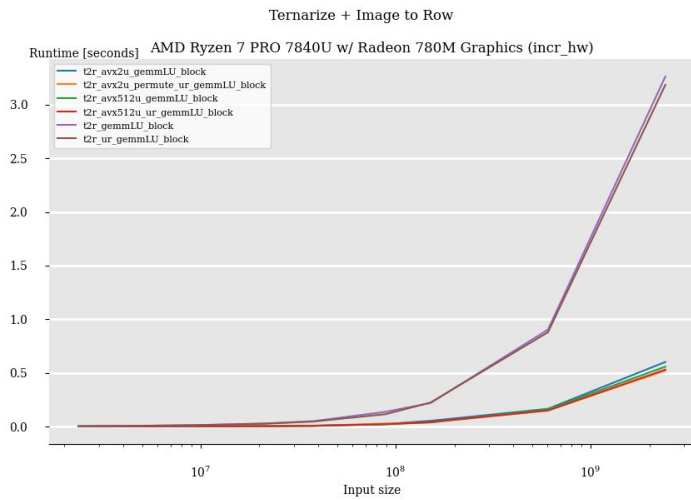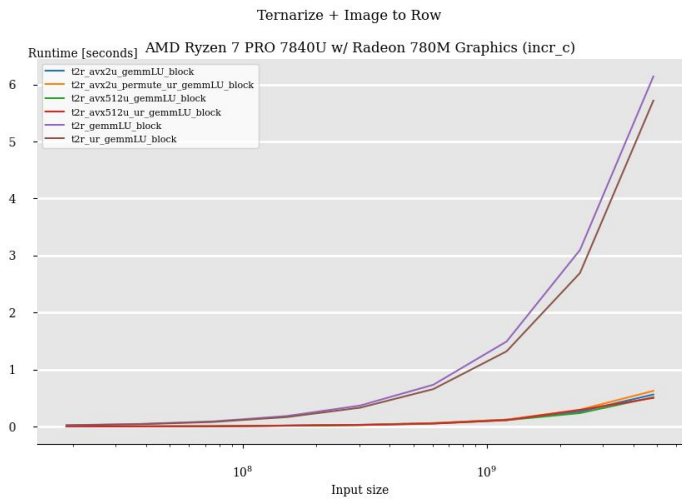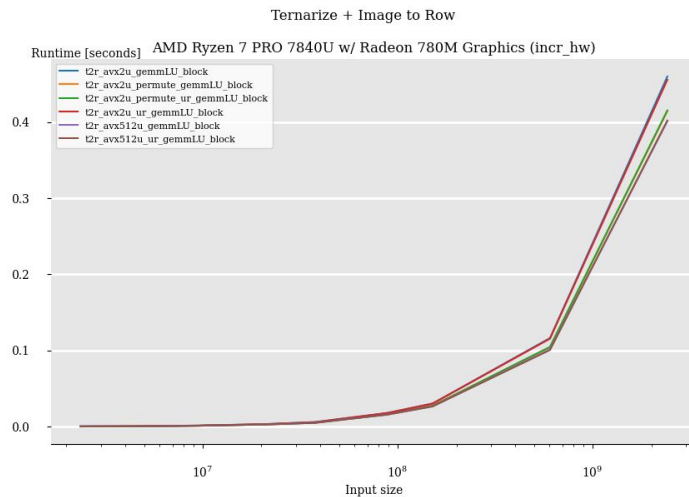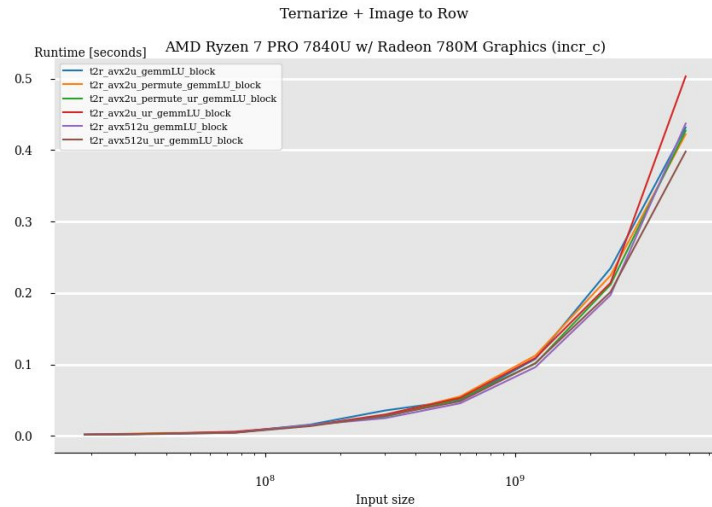# Optimizing Ternarize + im2row

| Loop Unrolling | AVX2 | AVX512 |
|---|---|---|

# Optimizing Ternarize + im2row

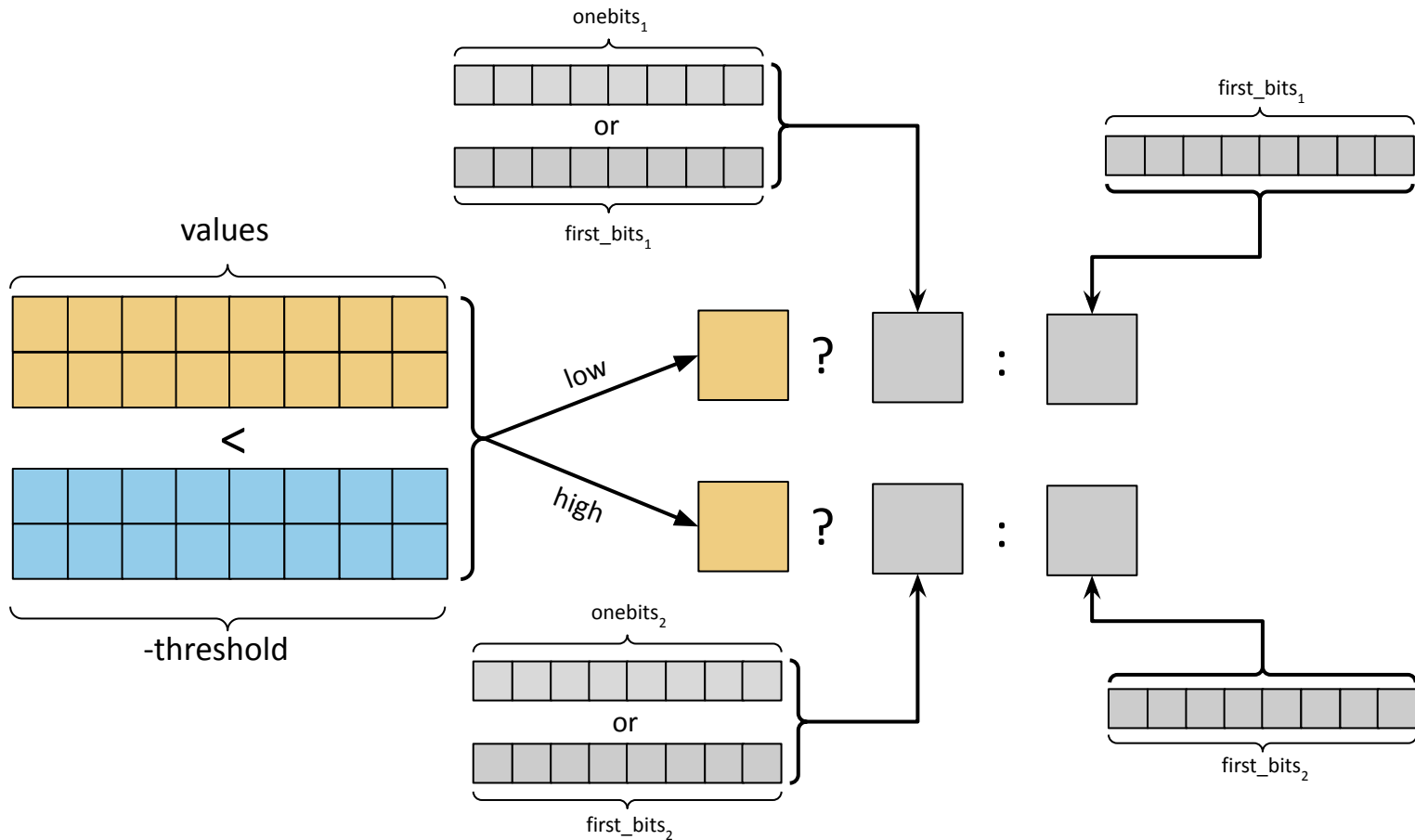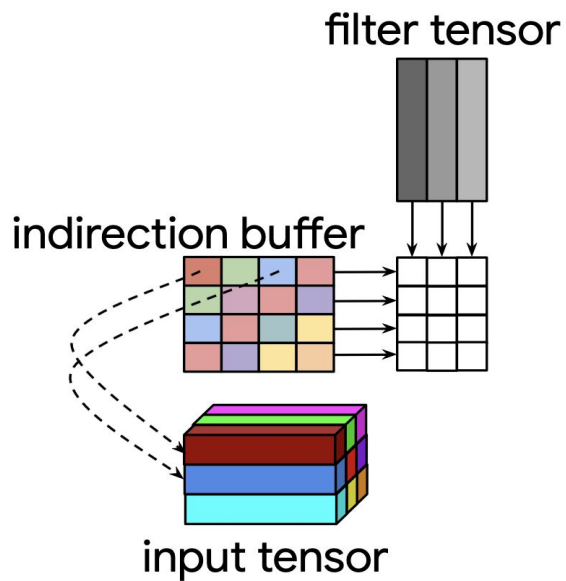| Loop Unrolling | AVX2 | AVX512 |
|:---:|:---:|:---:|

# Vectorizing Ternarize: Mixed AVX2-AVX512
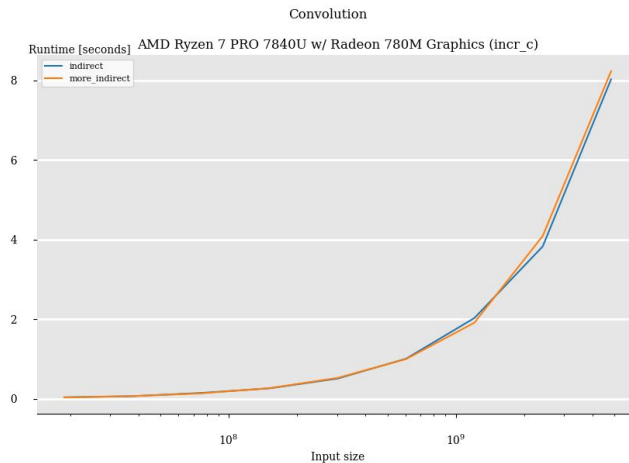


14

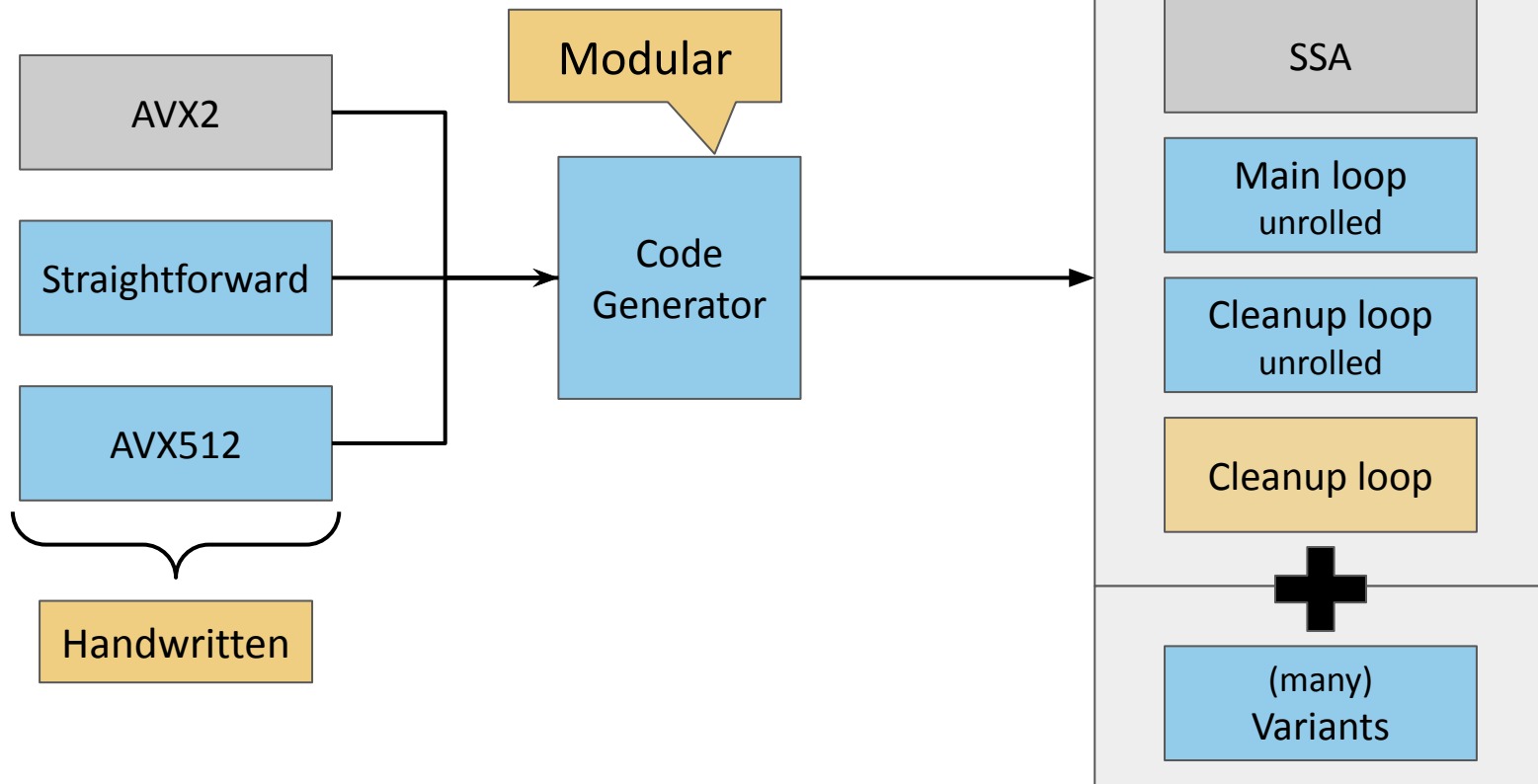# Vectorizing Ternarize: AVX512 only

# im2row + GEMM



Indirect Convolution

filter tensor

indirection buffer

input tensor

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



Convolution

Runtime [seconds] AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)

indirect
more_indirect

Input size

Convolution

Runtime [seconds] AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)

indirect
more_indirect

Input size

# Code Generation

# Code Generation - Variant: Autotuning + libpopcnt

**TNN GEMM**

**Autotuning on M and N**

```
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] = ..
```
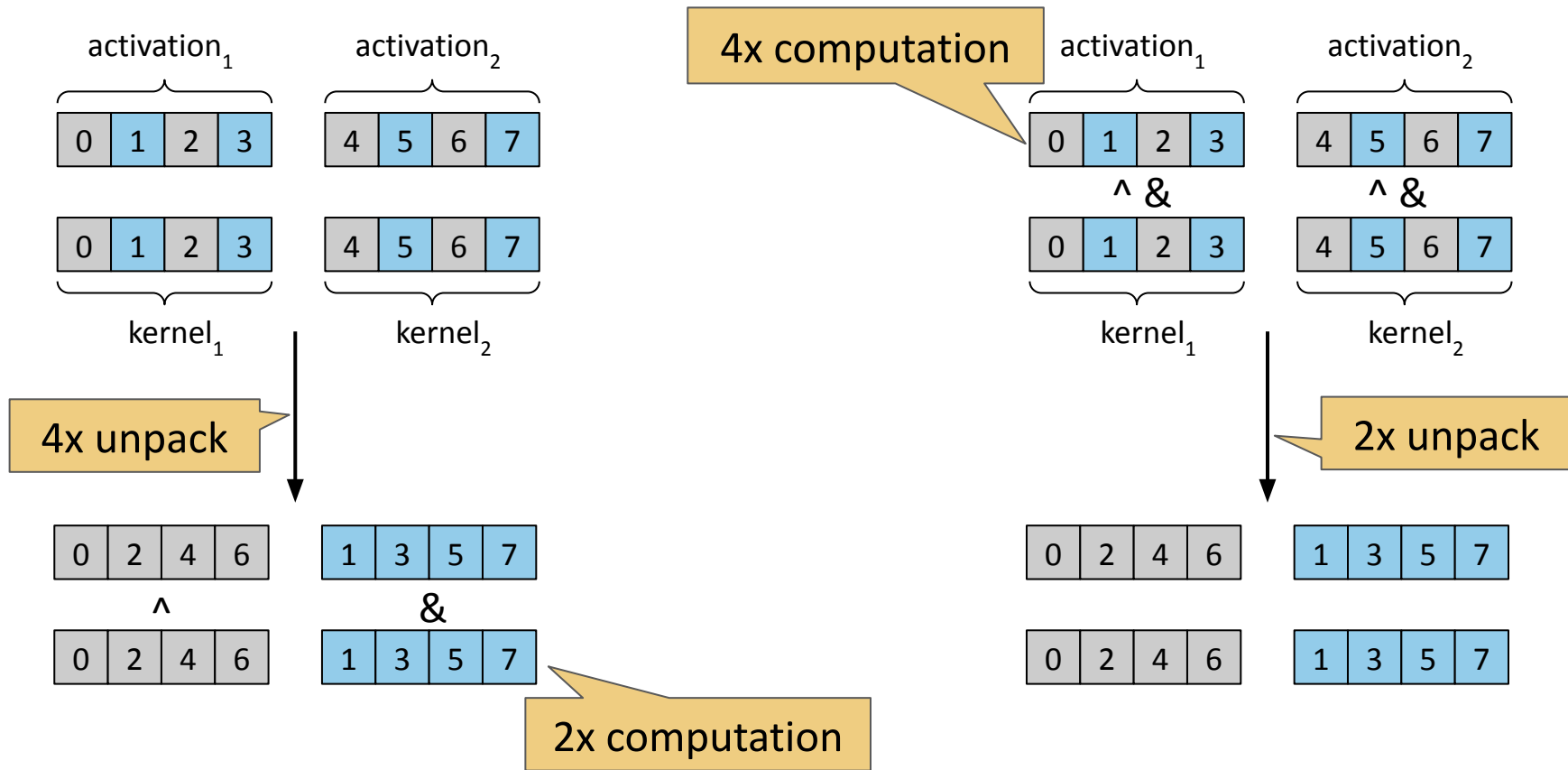
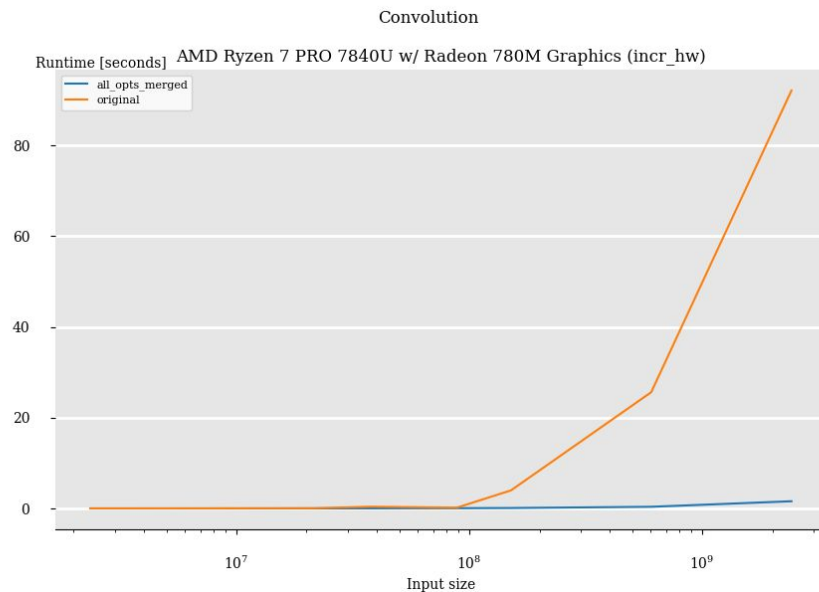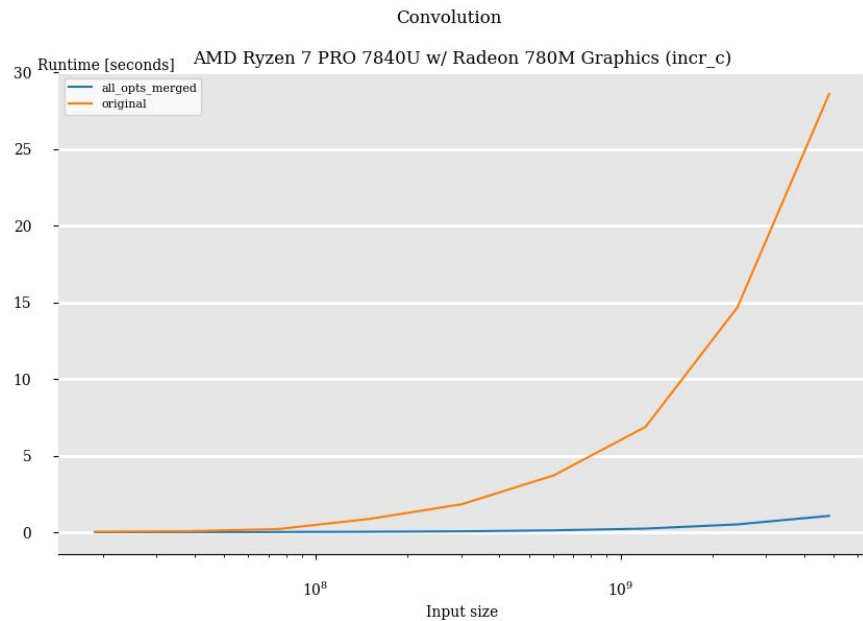**TNN GEMM**

```
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)
        ..
```
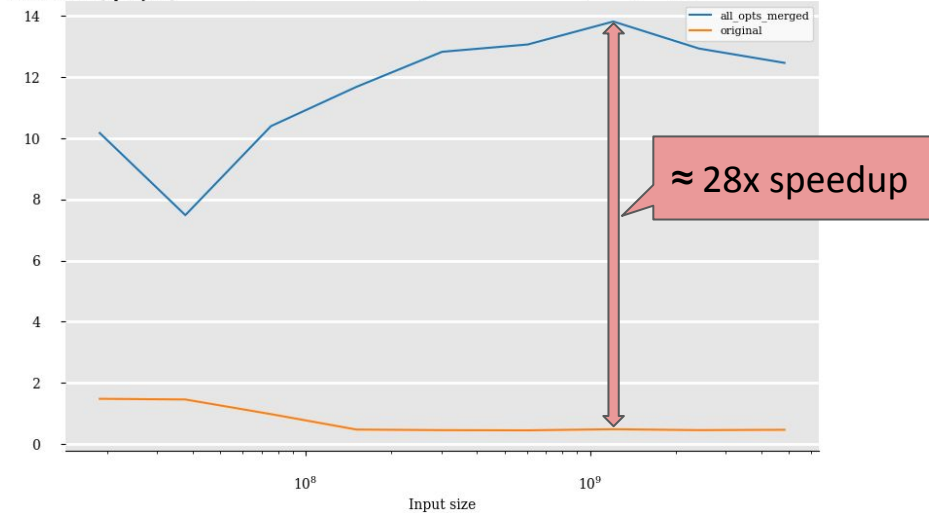
**popcount on a big vector**

# Code Generation - Variant: Less Unpacking

# Conclusion



Convolution
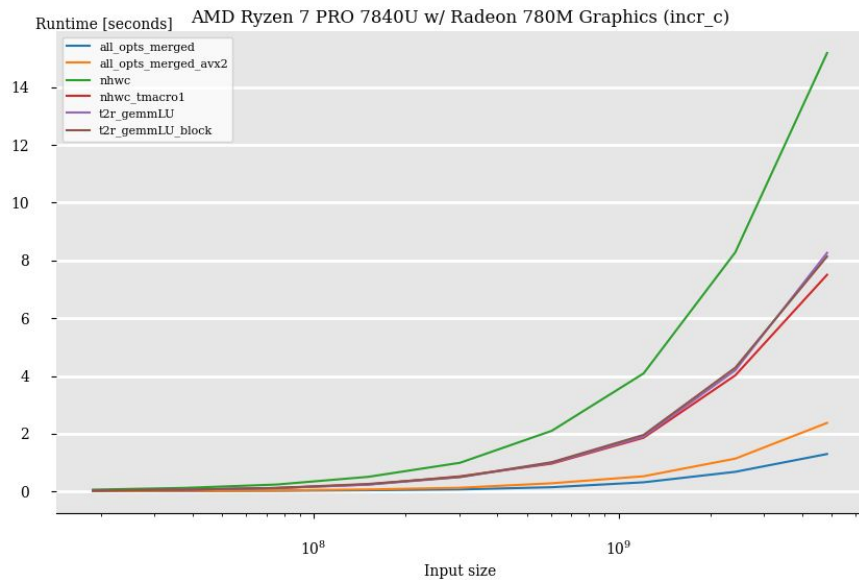AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)

Convolution
AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)

# Conclusion



≈ 28x speedup

# Conclusion



Convolution
AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_c)

Convolution
AMD Ryzen 7 PRO 7840U w/ Radeon 780M Graphics (incr_hw)