# ENABLE EDGE AI PRODUCTS YOU DREAM OF

Understanding linalg.pack and linalg.unpack

28 OCTOBER 2025
MAXIMILIAN BARTEL

roofline

# WHY IS A BLOCKED MATMUL IMPLEMENTATION NOT GETTING THE BEST PERFORMANCE?

## Logical layout



- Fast algorithm, but it generates cache misses:

- Reads from A are good

- Reads from B and C are not

## In-memory layout

# WE CAN IMPROVE MEMORY ACCESS PATTERNS BY CHANGING THE MEMORY LAYOUT

## Logical layout



- Changing the memory layout improves cache hit rate

- Optimal layout aligns perfectly with the reads
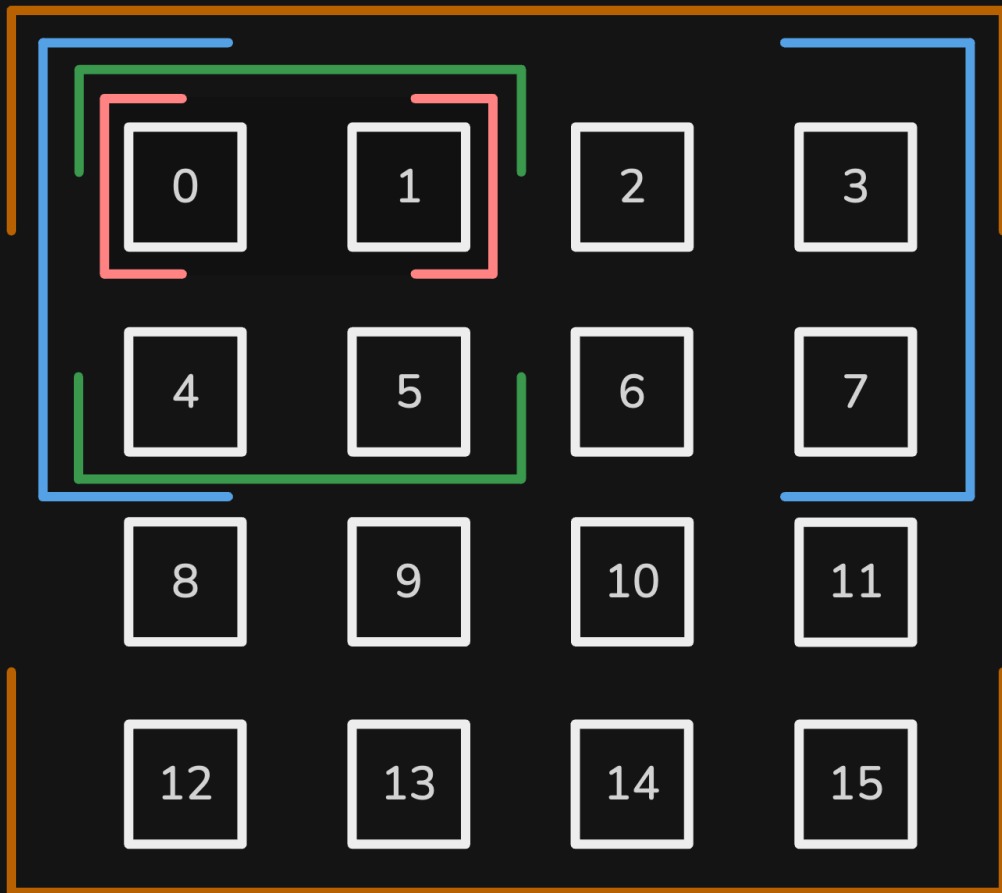
- This is called data tiling or packed layouts

## In-memory layout

# PACKING INCREASES THE DIMENSIONALITY OF THE TENSOR

## Exemplary 2D matrix
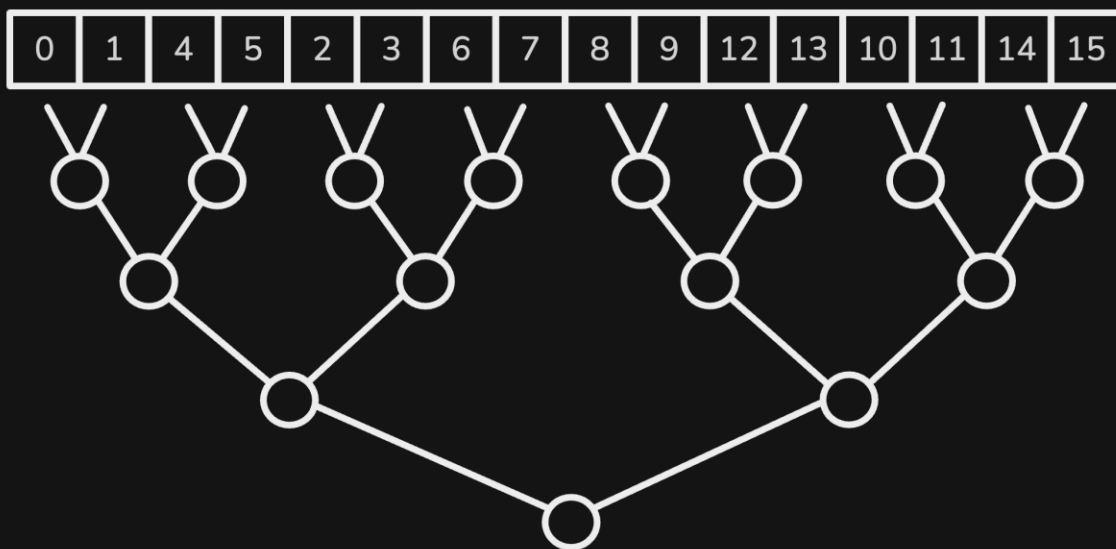


## Insights

- Starting form the original logical layout, the new dimensions are marked in different colors

- We go from 4x4 to 2x2x2x2

- This is so far just a view

# PACKING MEANS ACTIVELY FLATTEN THE LAYOUT IN MEMORY

## Flattened Tensor



## Insights

- Packing takes the view and flattens it

- The reduction tree visualizes the dimensions

- These data movements come at a cost!

EXAMPLE 1

# A PACK OP IS ALLOWED TO PAD TO A SPECIFIC TILE SIZE

Exemplary code

Insights

```
1 func.func @simple_pad_and_pack_static_tiles(
2     %input: tensor<3x1xf32>,
3     %output: tensor<1x1x5x2xf32>,
4     %pad: f32)
5         -> tensor<1x1x5x2xf32> {
6   %0 = linalg.pack %input
7         padding_value(%pad : f32)
8         inner_dims_pos = [0, 1]
9         inner_tiles = [5, 2]
10    into %output : tensor<5x1xf32> -> tensor<1x1x5x2xf32>
11  return %0 : tensor<1x1x5x2xf32>
12 }
```

EXAMPLE 2

# WE ALSO NEED TO BE ABLE TO REVERT THE PADDING

Exemplary code

Insights

```
1 func.func @unpack_as_pad(
2     %arg0: tensor<1x1x2x3xf32>, %arg1: tensor<1x2xf32>)
3         -> tensor<1x2xf32> {
4   %pack = linalg.unpack %arg0
5         inner_dims_pos = [0, 1]
6         inner_tiles = [2, 3]
7     into %arg1 : tensor<1x1x2x3xf32> -> tensor<1x2xf32>
8   return %pack : tensor<1x2xf32>
9 }
```

EXAMPLE 3

# DYNAMIC SHAPES ARE FULLY SUPPORTED BY THE UNPACK OP

Exemplary code

Insights

```
 1 func.func @unpack_fully_dynamic(
 2     %source: tensor<?x?x?x?xf32>, %dest: tensor<?x?xf32>,
 3     %tile_n : index, %tile_m : index)
 4         -> tensor<?x?xf32> {
 5   %0 = linalg.unpack %source
 6         inner_dims_pos = [0, 1]
 7         inner_tiles = [%tile_n, %tile_m]
 8     into %dest : tensor<?x?x?x?xf32> -> tensor<?x?xf32>
 9   return %0 : tensor<?x?xf32>
10 }
```

- Dynamic dimensions are fully supported by this op

- The inner_tiles can take SSA values

- A custom parser treats it as part of an attribute

roofline

EXAMPLE 3

# UNPACK WILL DECOMPOSE AND INSERT OPERATIONS TO GET THE SHAPES OF THE TENSORS

## Exemplary code

```
 1 func.func @unpack_fully_dynamic(
 2     %arg0: tensor<?x?x?x?xf32>, %arg1: tensor<?x?xf32>,
 3     %arg2: index, %arg3: index)
 4         -> tensor<?x?xf32> {
 5     ... // Get the dimensions of the packed tensor
 6     %0 = tensor.empty(%dim, %dim_1, %dim_0, %dim_2) : tensor<?x?x?x?xf32>
 7     %transposed = linalg.transpose ... permutation = [0, 2, 1, 3]
 8     %collapsed = tensor.collapse_shape %transposed [[0, 1], [2, 3]]
 9     %dim_3 = tensor.dim %arg1, %c0 : tensor<?x?xf32>
10     %dim_4 = tensor.dim %arg1, %c1 : tensor<?x?xf32>
11     %extracted_slice = tensor.extract_slice
12         %collapsed[0, 0] [%dim_3, %dim_4] [1, 1] : tensor<?x?xf32> to tensor<?x?xf32>
13     %1 = linalg.copy
14     return %1 : tensor<?x?xf32>
15 }
```

## Insights

- Decomposing dynamic shapes leads to many dim operations

- Even a padded layout is supported by dynamic shapes

- The output is parsed for its dimensions

roofline

EXAMPLE 4

# PACK ALSO SUPPORTS DYNAMIC SHAPES – BUT THE LOWERING CANNOT HANDLE IT YET

## Exemplary code

```
1 func.func @pack_fully_dynamic(
2     %source: tensor<?x?xf32>, %dest: tensor<?x?x?x?xf32>,
3     %tile_n : index, %tile_m : index, %cst_0 : f32)
4         -> tensor<?x?x?x?xf32> {
5   %0 = linalg.pack %source
6         padding_value(%cst_0 : f32)
7         inner_dims_pos = [0, 1]
8         inner_tiles = [%tile_n, %tile_m]
9     into %dest : tensor<?x?xf32> -> tensor<?x?x?x?xf32>
10   return %0 : tensor<?x?x?x?xf32>
11 }
```

## Insights

- This is a valid pack operation

- Decomposition should be like unpack

- However, it doesn't lower (yet)

- A lowering can be enabled by inser_slice supporting dynamic shapes

EXAMPLE 5

# UNIT DIMENSIONS ARE THE MOST COMMON PITFALLS AND PRODUCERS OF BUGS

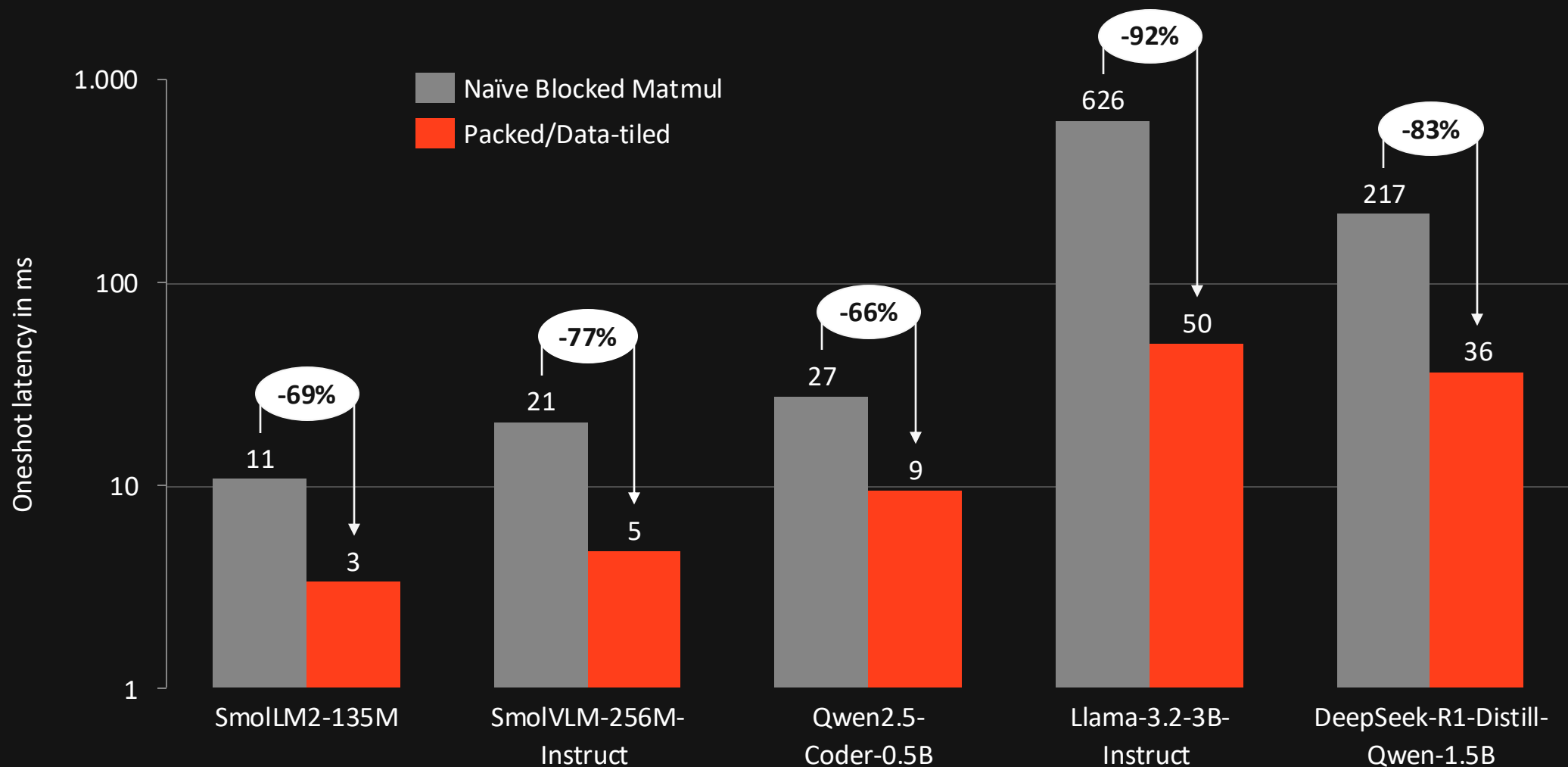## Exemplary code

```
 1 func.func @unit_dims(
 2     %arg0: tensor<1x1x1x4x1xf32>, %arg1: tensor<1x1x4xf32>)
 3         -> tensor<1x1x1x4x1xf32> {
 4   %pack = linalg.pack %arg1
 5         outer_dims_perm = [1, 2, 0]
 6         inner_dims_pos = [2, 0]
 7         inner_tiles = [4, 1]
 8     into %arg0 : tensor<1x1x4xf32> -> tensor<1x1x1x4x1xf32>
 9   return %pack : tensor<1x1x1x4x1xf32>
10 }
```
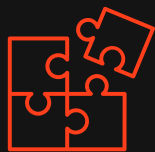
## Insights

- Most of the tricky behavior arises with unit dimensions

- This behavior is common after tiling

- Special handling is required to produce "better" IR in these cases

- This op eliminates data movements

- Untiled non-unit dim dimension in between unit dims are possible

# PACK AND UNPACK MAKE MEMORY ACCESSES VERY EFFICIENT, BUT THEY COME AT A COST

## How to use pack effectively

- Packing can degrade performance for single kernels

- Only full model compilation can fuse packing operators into producers to hide the movement cost

## What are unexplored paths?

- Explore data tiling for other kernels like convolutions

- Explore matrices as representations remove complex attributes (Triton)

roofline