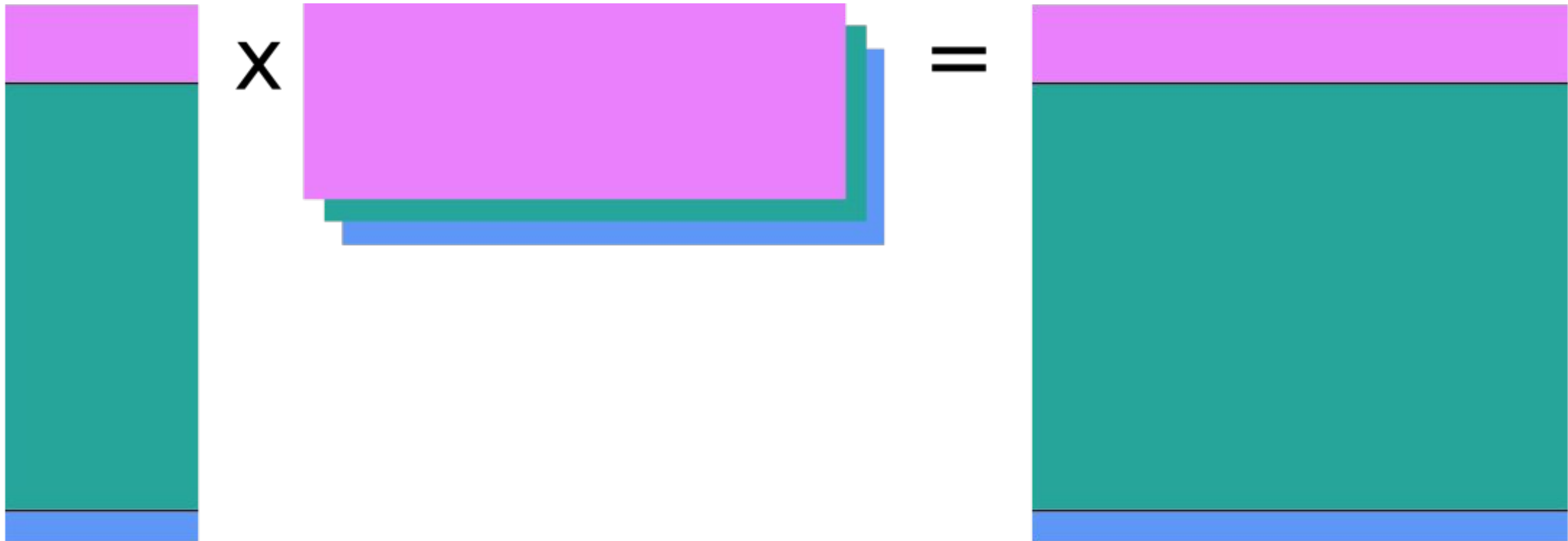


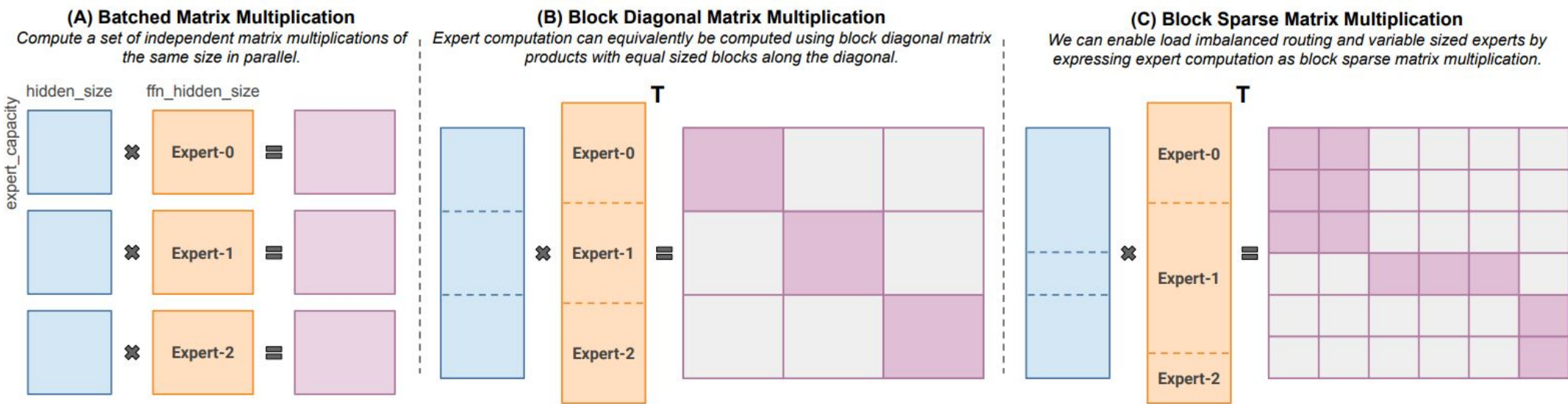
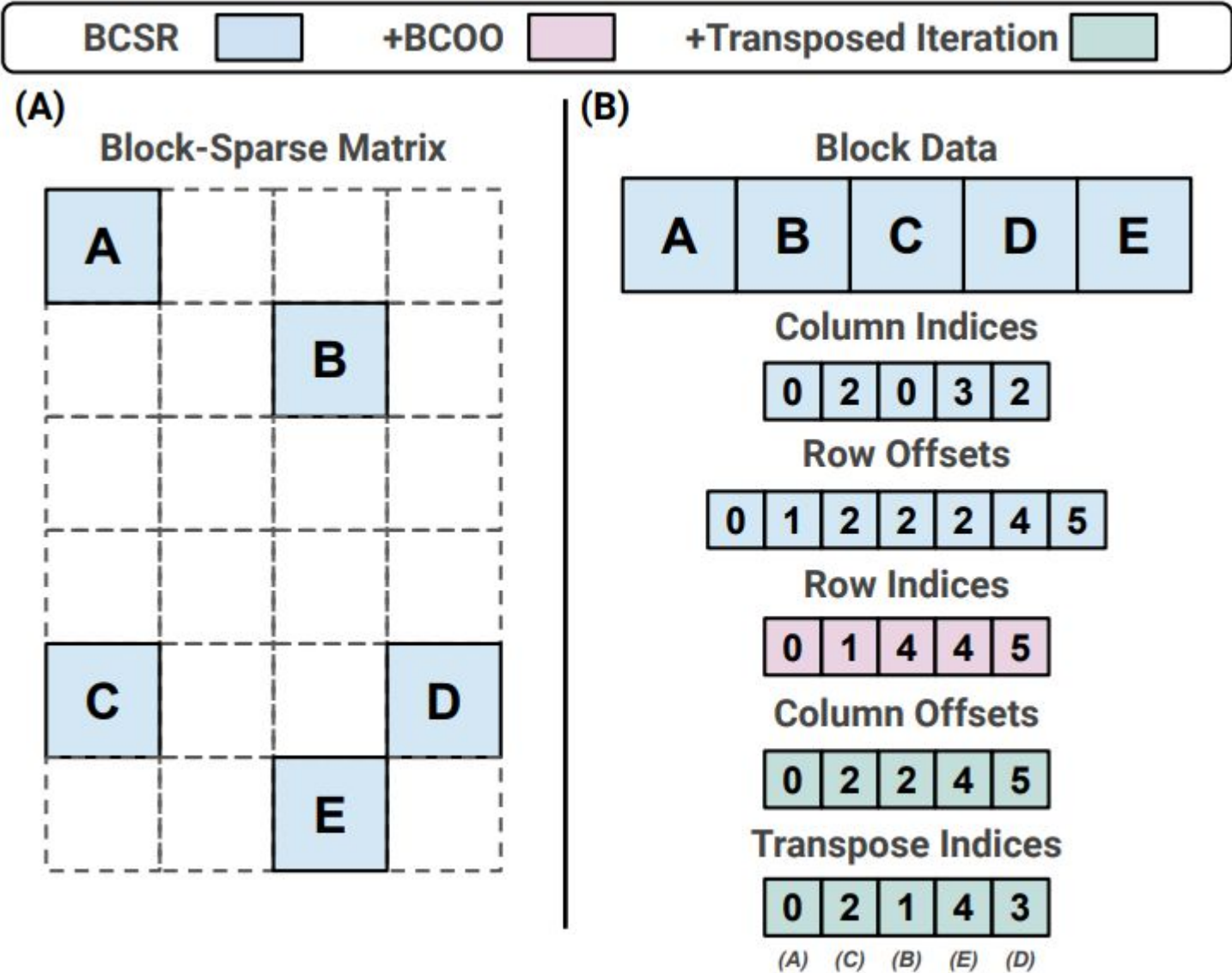
ragged dot



Intro

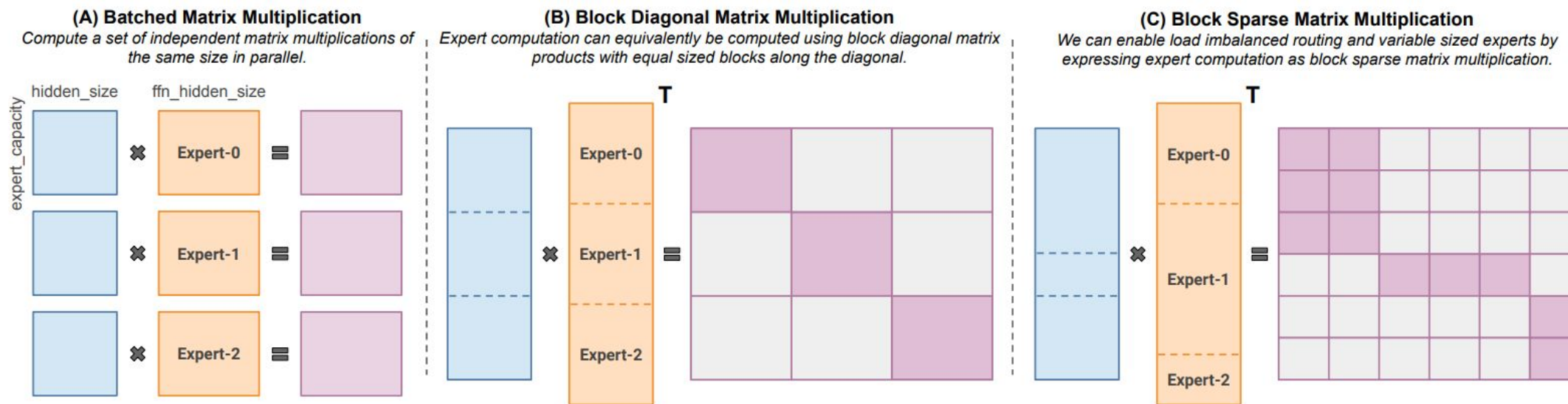
History

- *MegaBlocks: Efficient Sparse Training with Mixture-of-Experts*, Trevor Gale, Deepak Narayanan, Cliff Young, Matei Zaharia
- traditional sparsity represented using CSC or CSR (on GPU)
- modern accelerators deal badly with sparsity
- **Block-CSR - BCSR**



History

Motivation



- some possible approaches:
 - pad blocks to the maximum size (memory inefficient)
 - loop over rows individually (accelerator unfriendly)
- decent approaches on GPU
 - use NVIDIA libraries, dispatch **each expert on CPU**
 - not bad with a good compiler and/or CUDA-graphs

Base Routine

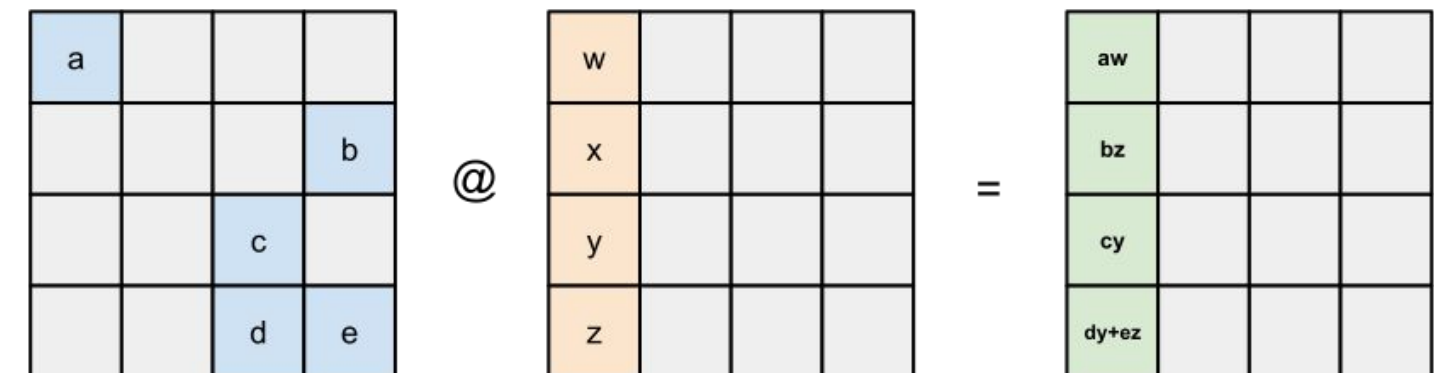
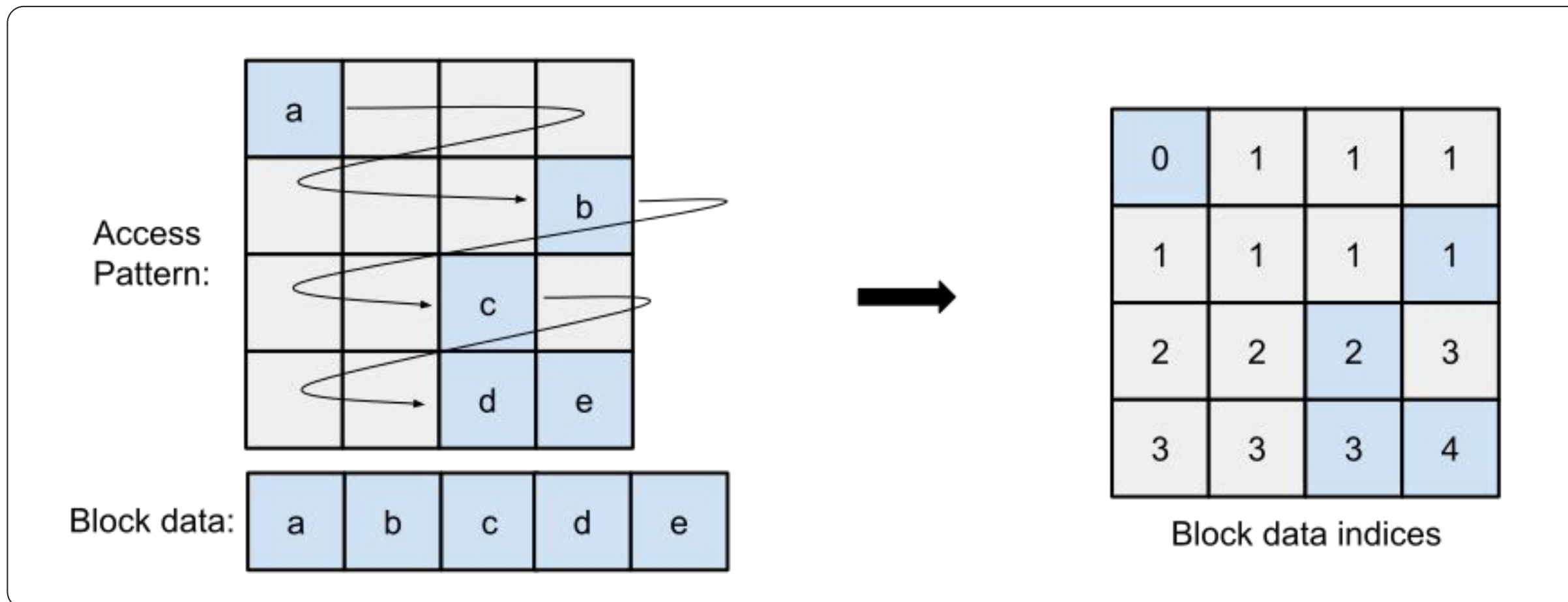


Fig: Block-sparse matrix multiplication
docs.jax.dev/en/latest/pallas/tpu/sparse.html#example-sparse-dense-matrix-multiplication

Sparse computations in pallas

pltpu.PrefetchScalarGridSpec



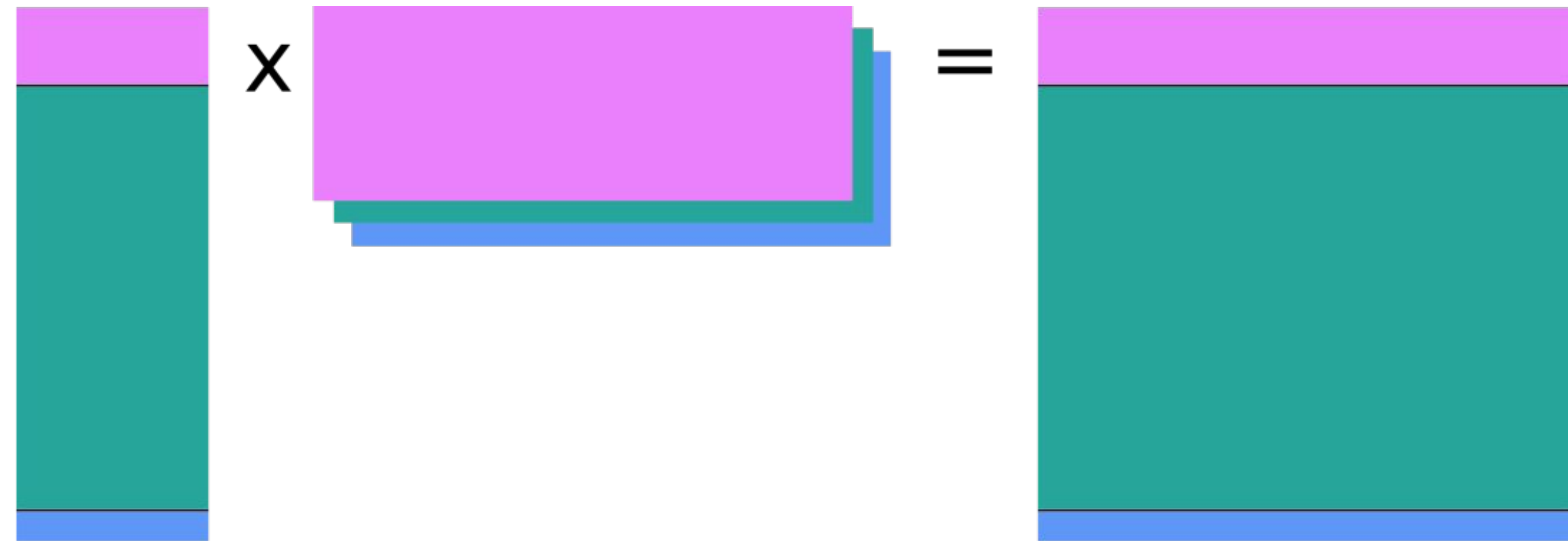
- “2, then 2” uses no BW
- read in the next block ASAP
- small penalty for visiting block at all
- can skip computation inside kernel

- pltpu.PrefetchScalarGridSpec is super useful
- generally pre-compute a metadata lookup table, then just iterate over entries
- example: `p1.BlockSpec(..., lambda i, j, scalar1_ref, scalar2_ref: (scalar1_ref[i], scalar2_ref[j]))`
- can re-implement the prefetch scalar grid using existing pallas APIs:

https://github.com/openxla/tokamax/blob/main/tokamax/_src/mosaic_tpu.py#L126

Modern “Megablox”

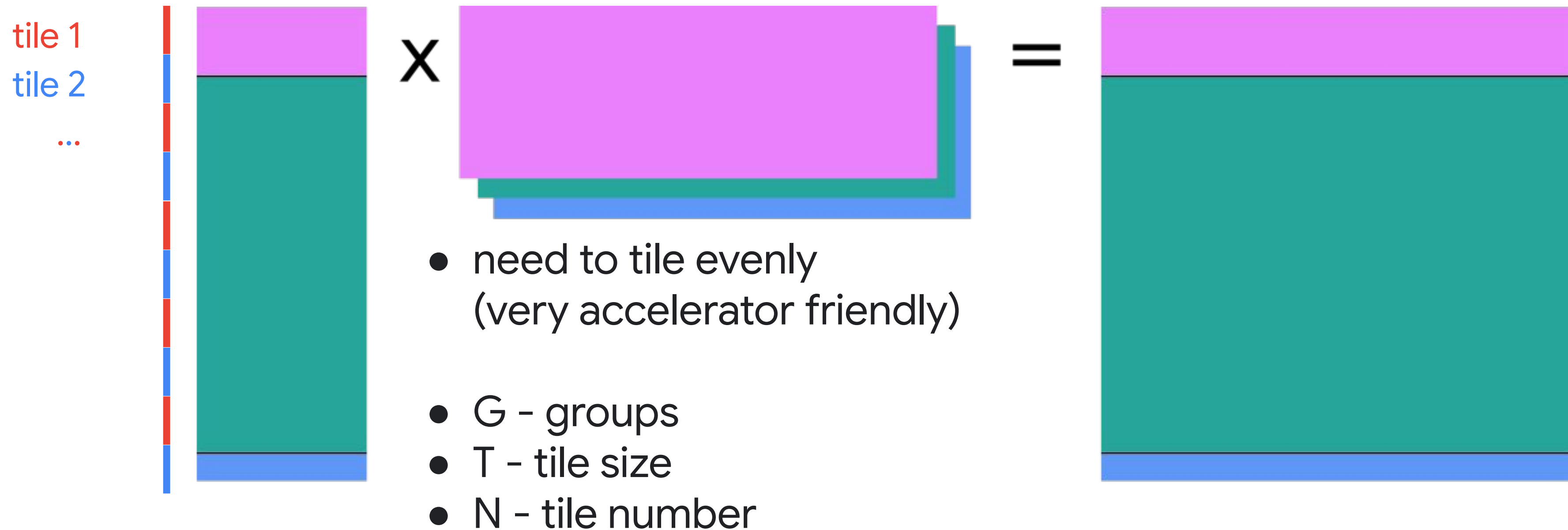
Megablox



- “ragged” representation usually means “pack stuff along one axis”
- the most dominant representation these days
- contiguous groups of rows belong to the (ordered) matrix in the stack
- LHS is 2D: $[m, k]$
- RHS is 3D: $[\text{group}, k, n]$
- OUT is 2D: $[m, n]$
- (theoretical) FLOPS are the same as for a simple matmul: $2 \cdot m \cdot k \cdot n$

Megablox

Accelerator Tiling & Tile Revisting



- probability of having to “revisit” a tile: $P(\text{revisit}) = 1 - 1 / T$
- number of tiles revisited: $G \cdot P(\text{revisit}) = G \cdot (1 - 1 / T)$
- efficiency is tile number / actual tiles visited: $N / (N + G \cdot (1 - 1 / T))$

Megablox Accelerator Efficiency

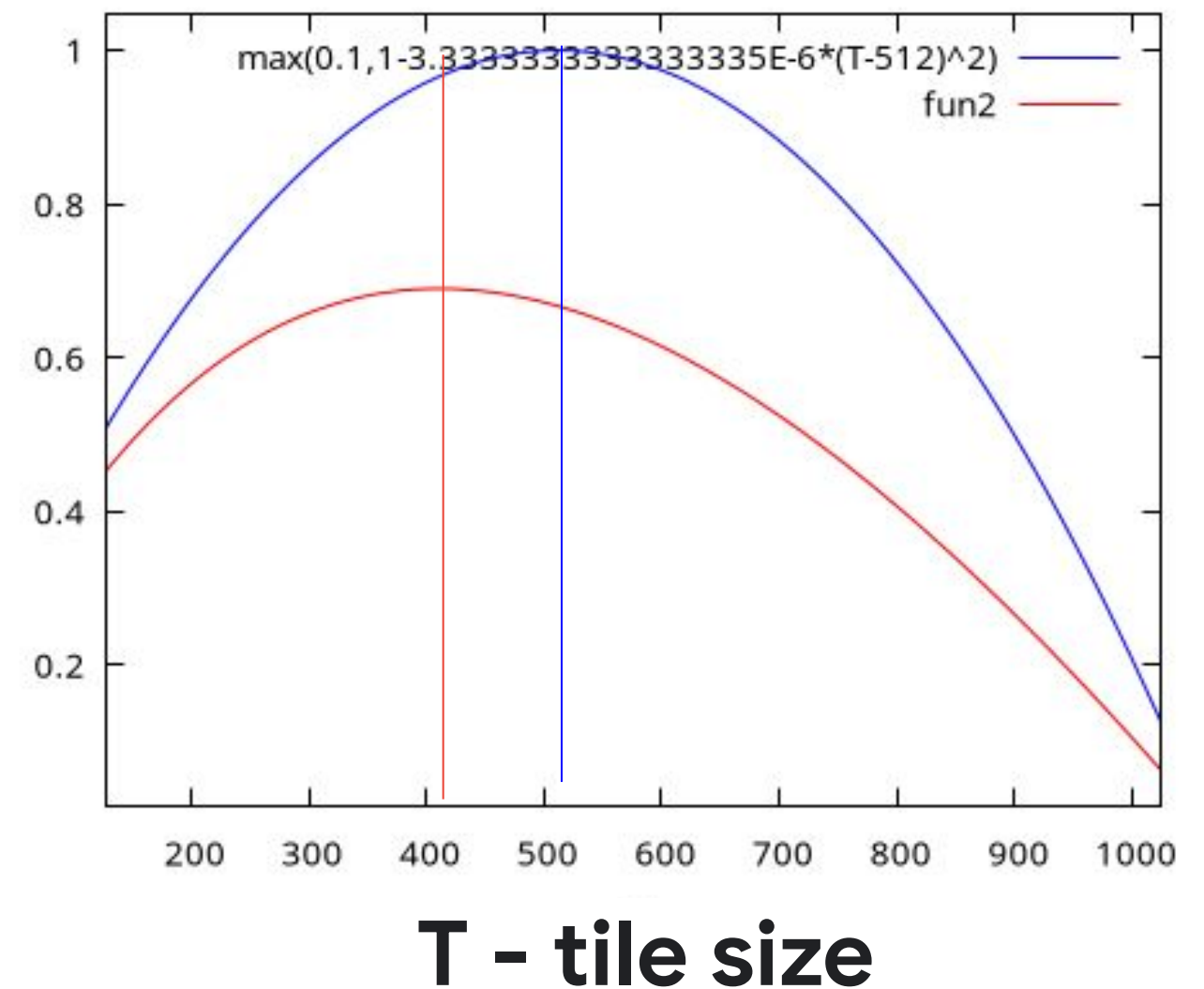
- G - groups
- T - tile size
- N - tile number

- probability of having to “revisit” a tile: $P(\text{revisit}) = 1 - 1 / T$
- number of tiles revisited: $G \cdot P(\text{revisit}) = G \cdot (1 - 1 / T)$
- efficiency is tile number / actual tiles visited: $N / (N + G \cdot (1 - 1 / T))$

megablox efficiency model

- assume optimal matmul tile size model
- smaller tile sizes: more block revisiting
- overall lower MFU
- optimal megablox tile size skews smaller

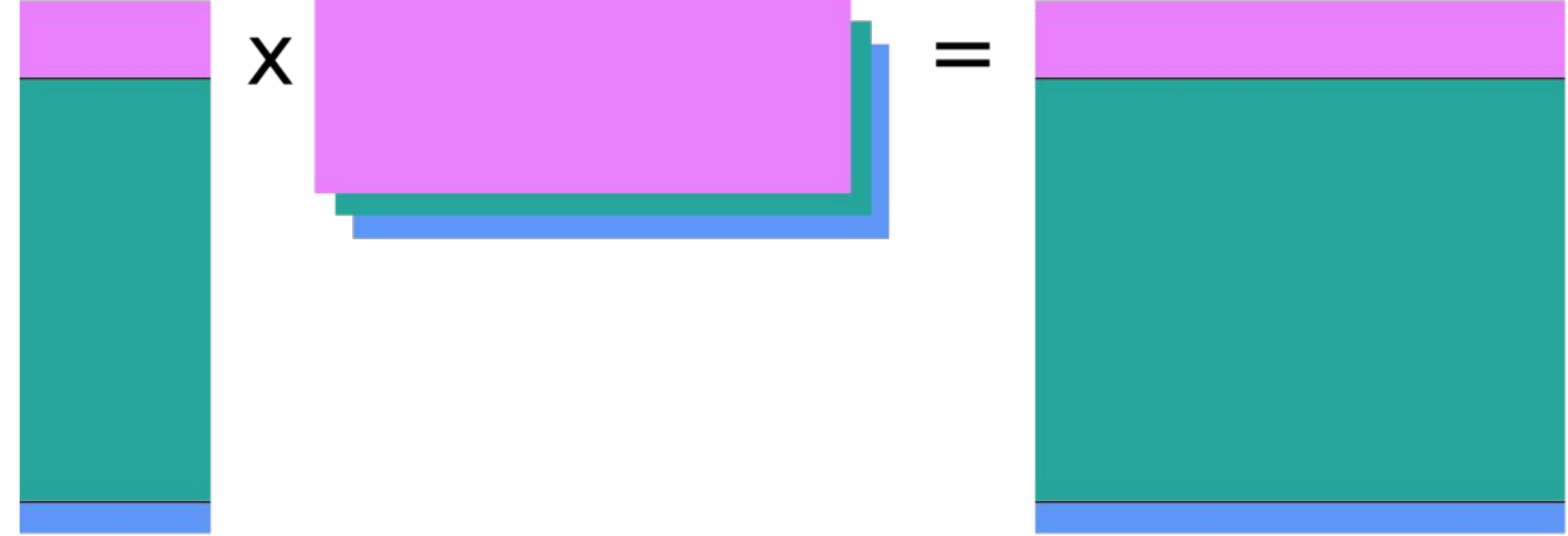
MFU efficiency (%)



Megablox

Derivatives

- LHS: $[m, k]$
- RHS: $[g, k, n]$
- OUT: $[m, n]$



d(LHS)

- needs to be $[m, k]$
- `ragged_dot(
 d(OUT),
 RHST # transpose RHS
)`
- FLOPS: $2 \cdot m \cdot k \cdot n$ (still)
- reduction along the `n-axis`

d(RHS)

- needs to be $[g, k, n]$
- `“transposed”_ragged_dot(
 LHST,
 d(OUT)
)`
- FLOPS: $2 \cdot m \cdot k \cdot n$ (still)
- lots of outer products
- reduction along: **masked** `m-axis`

Megablox

Actual metadata computation

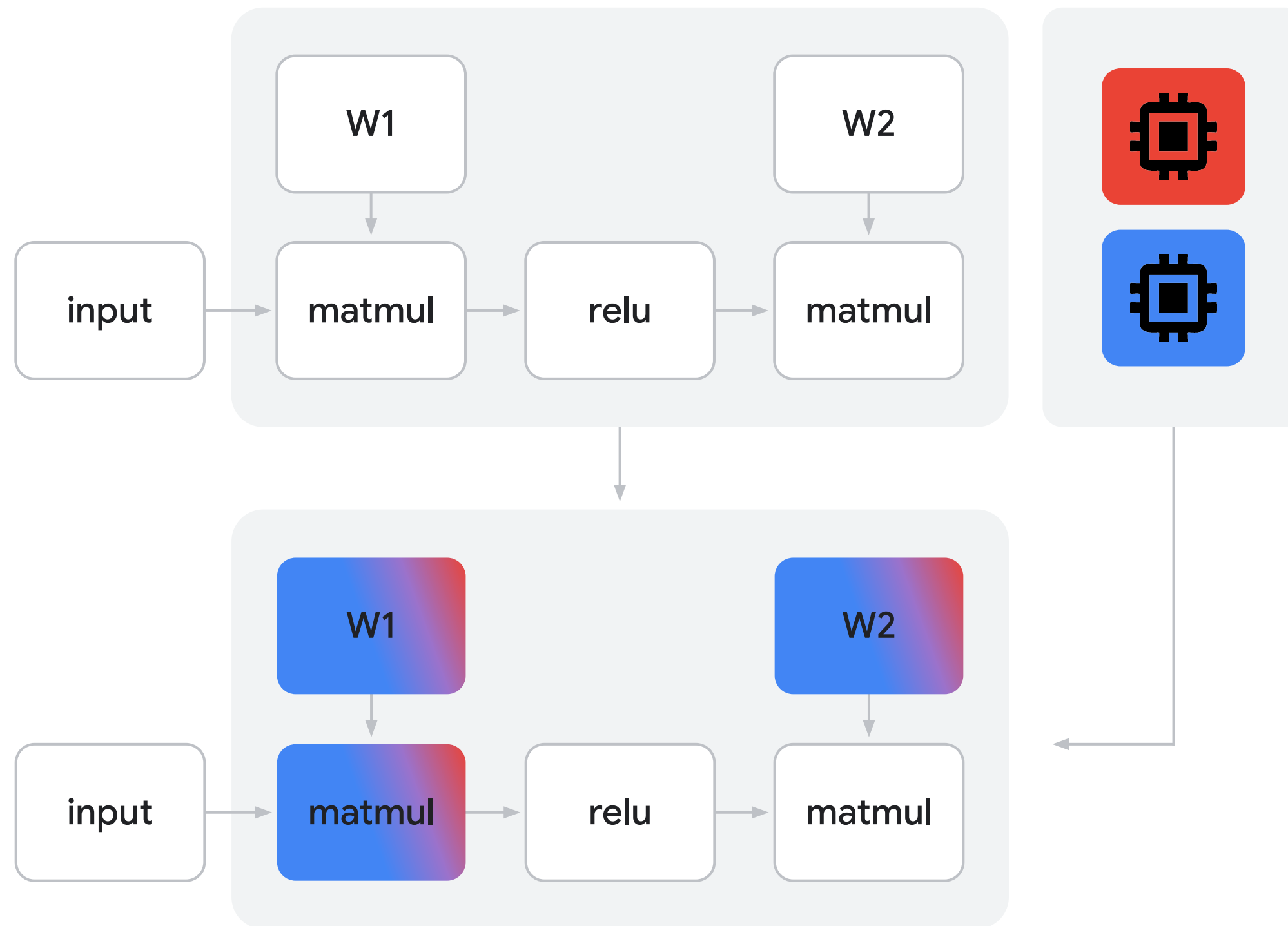
- don't visit & skip empty blocks
 - this costs kernel dispatch latency
- compute tile and expert map flat
 - sometimes tile increase
 - sometimes expert
 - just a flat lookup table
- how to do this efficiently on TPU?
 - no loops, must use e.g., cumsum
 - computing metadata blocks TC

```
1 def make_group_metadata(  
2     *,  
3     group_sizes: jnp.ndarray,  
4     m: int,  
5     tm: int,  
6     start_group: jnp.ndarray,  
7     num_nonzero_groups: int,  
8     visit_empty_groups: bool = True,  
9 ) -> GroupMetadata:  
10     """Create the metadata needed for grouped matmul computation.  
11  
12     Args:  
13     group_sizes: A 1d, jnp.ndarray with shape [num_groups] and jnp.int32 dtype.  
14     m: The number of rows in lhs.  
15     tm: The m-dimension tile size being used.  
16     start_group: The group in group sizes to start computing from. This is  
17     particularly useful for when rhs num_groups is sharded.  
18     num_nonzero_groups: Number of groups in group sizes to compute on. Useful in  
19     combination with group_offset.  
20     visit_empty_groups: If True, do not squeeze tiles for empty groups out of  
21     the metadata. This is necessary for tgmm, where we at least need to zero  
22     the output for each group.  
23  
24     Returns:  
25     tuple of:  
26     group_offsets: A 1d, jnp.ndarray with shape [num_groups+1] and jnp.int32  
27     dtype. group_offsets[i] indicates the row at which group [i] starts in  
28     the lhs matrix and group_offsets[i-1] = m.  
29     group_ids: A 1d, jnp.ndarray with shape [m_tiles + num_groups] and  
30     jnp.int32 dtype. group_ids[i] indicates which group grid index 'i' will  
31     work on.  
32     m_tile_ids: A 1d, jnp.ndarray with shape [m_tiles + num_groups] and  
33     jnp.int32. m_tile_ids[i] indicates which m-dimension tile grid index 'i'  
34     will work on.  
35     num_tiles: The number of m-dimension tiles to execute.  
36     """  
37     num_groups = group_sizes.shape[0]  
38     end_group = start_group + num_nonzero_groups - 1  
39  
40     # Calculate the offset of each group, starting at zero. This metadata is  
41     # similar to row offsets in a CSR matrix. The following properties hold:  
42     #  
43     # group_offsets.shape = [num_groups + 1]  
44     # group_offsets[0] = 0  
45     # group_offsets[num_groups] = m  
46     #  
47     # The row at which group 'i' starts is group_offsets[i].  
48     group_ends = jnp.cumsum(group_sizes)  
49     group_offsets = jnp.concatenate([jnp.zeros(1, dtype=jnp.int32), group_ends])  
50  
51     # Assign a group id to each grid index.  
52     #  
53     # If a group starts somewhere other than the start of a tile or ends somewhere  
54     # other than the end of a tile we need to compute that full tile. Calculate  
55     # the number of tiles for each group by rounding their end up to the nearest  
56     # 'tm' and their start down to the nearest 'tm'.  
57  
58     # (1) Round the group_ends up to the nearest multiple of 'tm'.  
59     #  
60     # NOTE: This does not change group_offsets[num_groups], which is m  
61     # (because we enforce m is divisible by tm).  
62     rounded_group_ends = ((group_ends + tm - 1) // tm * tm).astype(jnp.int32)  
63  
64     # (2) Round the group_starts down to the nearest multiple of 'tm'.  
65     group_starts = jnp.concatenate(  
66         [jnp.zeros(1, dtype=jnp.int32), group_ends[:-1]]  
67     )  
68     rounded_group_starts = group_starts // tm * tm  
69  
70     # (3) Calculate the number of rows in each group.  
71     #  
72     # NOTE: Handle zero-sized groups as a special case. If the start for a  
73     # zero-sized group is not divisible by 'tm' its start will be rounded down and  
74     # its end will be rounded up such that its size will become 1 tile here.  
75     rounded_group_sizes = rounded_group_ends - rounded_group_starts  
76     rounded_group_sizes = jnp.where(group_sizes == 0, 0, rounded_group_sizes)  
77  
78     # (4) Convert the group sizes from units of rows to unit of 'tm' sized tiles.  
79     #  
80     # An m-dimension tile is 'owned' by group 'i' if the first row of the tile  
81     # belongs to group 'i'. In addition to owned tiles, each group can have 0 or 1  
82     # initial partial tiles if it's first row does not occur in the first row of a  
83     # tile. The '0-th' group never has a partial tile because it always starts at  
84     # the 0-th row.  
85     #  
86     # If no group has a partial tile, the total number of tiles is equal to  
87     # 'm // tm'. If every group has a partial except the 0-th group, the total  
88     # number of tiles is equal to 'm // tm + num_groups - 1'. Thus we know that  
89     #  
90     # tiles_m <= group_tiles.sum() <= tiles_m + num_groups - 1
```

```
91     #  
92     # Where tiles_m = m // tm.  
93     #  
94     # NOTE: All group sizes are divisible by 'tm' because of the rounding in steps  
95     # (1) and (2) so this division is exact.  
96     group_tiles = rounded_group_sizes // tm  
97  
98     if visit_empty_groups:  
99         # Insert one tile for empty groups.  
100         group_tiles = jnp.where(group_sizes == 0, 1, group_tiles)  
101  
102     # Create the group ids for each grid index based on the tile counts for each  
103     # group.  
104     #  
105     # NOTE: This repeat(...) will pad group_ids with the final group id if  
106     # group_tiles.sum() < tiles_m + num_groups - 1. The kernel grid will be sized  
107     # such that we only execute the necessary number of tiles.  
108     tiles_m = _calculate_num_tiles(m, tm)  
109     group_ids = jnp.repeat(  
110         jnp.arange(num_groups, dtype=jnp.int32),  
111         group_tiles,  
112         total_repeat_length=tiles_m + num_groups - 1,  
113     )  
114  
115     # Assign an m-dimension tile id to each grid index.  
116     #  
117     # NOTE: Output tiles can only be re-visited consecutively. The following  
118     # procedure guarantees that m-dimension tile indices respect this.  
119  
120     # (1) Calculate how many times each m-dimension tile will be visited.  
121     #  
122     # Each tile is guaranteed to be visited once by the group that owns the tile.  
123     # The remaining possible visits occur when a group starts inside of a tile at  
124     # a position other than the first row. We can calculate which m-dimension tile  
125     # each group starts in by floor-dividing its offset with 'tm' and then count  
126     # tile visits with a histogram.  
127     #  
128     # To avoid double counting tile visits from the group that owns the tile,  
129     # filter these out by assigning their tile id to 'tile_m' (one beyond the max)  
130     # such that they're ignored by the subsequent histogram. Also filter out any  
131     # group which is empty.  
132     #  
133     # TODO(tgale): Invert the 'partial_tile_mask' predicates to be more clear.  
134     partial_tile_mask = jnp.logical_or(  
135         (group_offsets[:-1] % tm) == 0, group_sizes == 0  
136     )  
137  
138     # Explicitly enable tiles for zero sized groups, if specified. This covers  
139     # zero sized groups that start on a tile-aligned row and those that do not.  
140     if visit_empty_groups:  
141         partial_tile_mask = jnp.where(group_sizes == 0, 0, partial_tile_mask)  
142  
143     partial_tile_ids = jnp.where(  
144         partial_tile_mask, tiles_m, group_offsets[:-1] // tm  
145     )  
146  
147     tile_visits = (  
148         jnp.histogram(partial_tile_ids, bins=tiles_m, range=(0, tiles_m - 1))[0]  
149         + 1  
150     )  
151  
152     # Create the m-dimension tile ids for each grid index based on the visit  
153     # counts for each tile.  
154     m_tile_ids = jnp.repeat(  
155         jnp.arange(tiles_m, dtype=jnp.int32),  
156         tile_visits.astype(jnp.int32),  
157         total_repeat_length=tiles_m + num_groups - 1,  
158     )  
159  
160     # Account for sharding.  
161     #  
162     # Find the start of the groups owned by our shard and shift the group_ids and  
163     # m_tile_ids s.t. the metadata for our tiles are at the front of the arrays.  
164     #  
165     # TODO(tgale): Move this offset into the kernel to avoid these rolls.  
166     first_tile_in_shard = (group_ids < start_group).sum()  
167     group_ids = jnp.roll(group_ids, shift=-first_tile_in_shard, axis=0)  
168     m_tile_ids = jnp.roll(m_tile_ids, shift=-first_tile_in_shard, axis=0)  
169  
170     # Calculate the number of tiles we need to compute for our shard.  
171     #  
172     # Remove tile visits that belong to a group not in our shard.  
173     iota = jnp.arange(num_groups, dtype=jnp.int32)  
174     active_group_mask = jnp.logical_and(iota <= end_group, iota >= start_group)  
175     group_tiles = jnp.where(active_group_mask, group_tiles, 0)  
176     num_tiles = group_tiles.sum()  
177     return (group_offsets, group_ids, m_tile_ids), num_tiles  
178
```

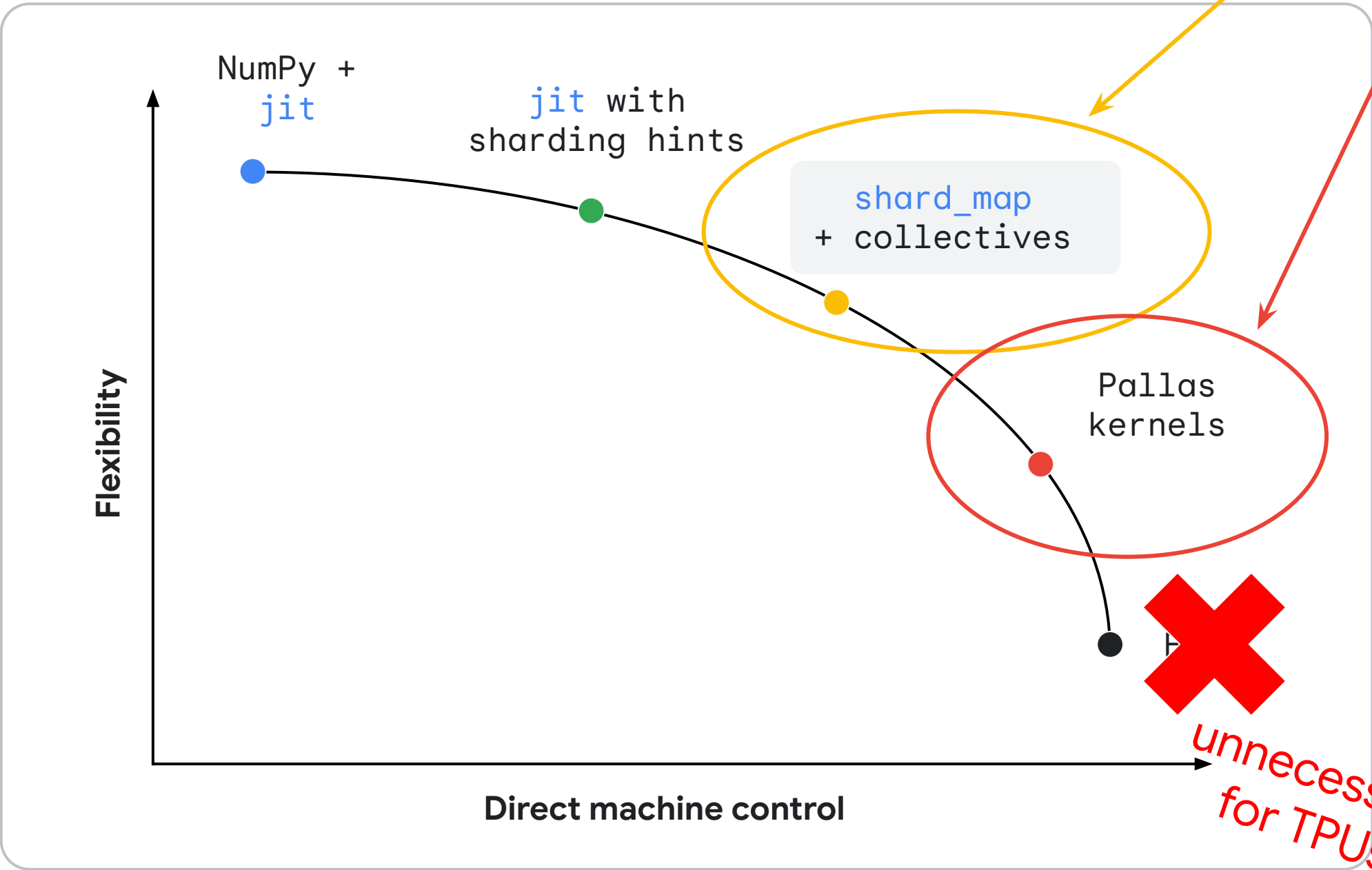
The MoE Layer in JAX

Sharded MLP (or MoE)



JAX: The escape hatch hierarchy

Flexibility vs. control



for MoE
for Megablox

Math

- Compiler, take the wheel!
- Here's a hint
- I'll handle comms
- Kernel languages
- FFI

Physics

MoE Prep

helper lambdas

token routing
(just computing the index map)

in specs and out specs for shard_map

constants for the MoE closure

```
1 def moe_block(x: jax.Array, layer: MoELayer, cfg: Config):
2     assert x.ndim == 3
3     l2p = lambda *axes: logical_to_physical(axes, cfg.rules)
4     _psc = lambda z, spec: reshard(z, P(*spec))
5     psc = lambda z, spec: _qpsc(z, spec) if is_type(z, QuantArray) else _psc(z, spec)
6
7     # we're decoding or device count does not divide total token count
8     replicated_routing = x.shape[-2] == 1 or (x.shape[-2] * x.shape[-3]) % jax.device_count() != 0
9     topk_weights, topk_idx = _route_tokens_to_moe_experts(x, layer.w_router, replicated_routing, cfg)
10    tensor_axname, expert_axname = l2p("moe_e_tp")[0], l2p("moe_e_ep")[0]
11
12    x_spec = l2p("batch", "sequence", None)
13    topk_weights_spec, topk_idx_spec = l2p("batch", "sequence", None), l2p("batch", "sequence", None)
14    out_spec = l2p("batch", "sequence", None)
15
16    we_gate_spec = l2p("moe_e_ep", None, "moe_e_tp")
17    we_up_spec = l2p("moe_e_ep", None, "moe_e_tp")
18    we_down_spec = l2p("moe_e_ep", "moe_e_tp", None)
19    we_gate = psc(layer.we_gate, we_gate_spec)
20    we_up = psc(layer.we_up, we_up_spec)
21    we_down = psc(layer.we_down, we_down_spec)
22
23    in_specs = (x_spec, we_gate_spec, we_up_spec, we_down_spec, topk_weights_spec, topk_idx_spec)
24
25    is_embedding_sharded = l2p("act_embed")[0] is not None
26    if is_embedding_sharded: # activations are sharded
27        out_spec = P(*(out_spec[:-1] + (tensor_axname,))) # override last axis name
28
29    expert_count = cfg.mesh.axis_sizes[cfg.mesh.axis_names.index(expert_axname)] if expert_axname is not None else 1
30    tensor_count = cfg.mesh.axis_sizes[cfg.mesh.axis_names.index(tensor_axname)] if tensor_axname is not None else 1
31    assert cfg.moe_num_experts % expert_count == 0
32    expert_size = cfg.moe_num_experts // expert_count
33
```

MoE Compute

sort indices

gather tokens
(potentially increase their count)

group sizes via bincount

up and gate projection
(first tensor-parallel stage)

down projection
(second tensor-parallel stage)

weight the tokens

all-gather tokens

and reduce across experts

```
38 @partial(shard_map, mesh=cfg.mesh, in_specs=in_specs, out_specs=out_spec, check_rep=False)
39 def _expert_fn(x, we_gate, we_up, we_down, topk_weights, topk_idx):
40     (b, s, d), e = x.shape, cfg.moe_experts_per_tok
41     expert_idx = jax.lax.axis_index(expert_axname) if expert_axname is not None else 0
42     tensor_idx = jax.lax.axis_index(tensor_axname) if tensor_axname is not None else 0
43     del tensor_idx
44     topk_idx = topk_idx.reshape(-1)
45     valid_group_mask = (topk_idx >= expert_size * expert_idx) & (topk_idx < expert_size * (expert_idx + 1))
46     expert_mapped_topk_idx = jnp.where(valid_group_mask, topk_idx - expert_idx * expert_size, 2**30)
47
48     sort_idx = jnp.argsort(expert_mapped_topk_idx, axis=-1) # [b * s * e]
49     isort_idx = jnp.argsort(sort_idx)
50
51     topk_idx_sort = topk_idx[sort_idx] # [b * s * e]
52     expert_mapped_topk_idx_sort = expert_mapped_topk_idx[sort_idx]
53     valid_group_mask_sort = expert_mapped_topk_idx_sort < 2**30
54     expert_mapped_topk_idx_sort = jnp.where(expert_mapped_topk_idx_sort < 2**30, expert_mapped_topk_idx_sort, 0)
55
56     # equivalent to:
57     # x_repeat = jnp.repeat(x.reshape((-1, x.shape[-1])), e, axis=0)
58     # x_repeat_sort = jnp.take_along_axis(x_repeat, sort_idx[:, None], axis=-2) # [b * s, d]
59     # x_repeat_sort = jnp.take_along_axis(
60     #     x.reshape((-1, x.shape[-1])),
61     #     sort_idx[:, None] // e,
62     #     axis=-2, # index trick to avoid jnp.repeat
63     # ) # [b * s * e, d]
64
65     group_sizes = jnp.bincount(topk_idx_sort, length=cfg.moe_num_experts)
66     group_sizes_shard = jax.lax.dynamic_slice_in_dim(group_sizes, expert_idx * expert_size, expert_size, 0)
67
68     with jax.named_scope("we_gate"):
69         ff_gate = _moe_gmm(x_repeat_sort, we_gate, group_sizes_shard, expert_mapped_topk_idx_sort, cfg)
70         ff_gate = jax.nn.silu(ff_gate)
71         ff_gate = jnp.where(valid_group_mask_sort[..., None], ff_gate, 0)
72     with jax.named_scope("we_up"):
73         ff_up = _moe_gmm(x_repeat_sort, we_up, group_sizes_shard, expert_mapped_topk_idx_sort, cfg)
74         ff_gate_up = jnp.where(valid_group_mask_sort[..., None], ff_gate * ff_up, 0)
75     with jax.named_scope("we_down"):
76         ff_out = _moe_gmm(ff_gate_up, we_down, group_sizes_shard, expert_mapped_topk_idx_sort, cfg)
77         ff_out = jnp.where(valid_group_mask_sort[..., None], ff_out, 0) # expensive
78
79     ff_out = ff_out * topk_weights.reshape(-1)[sort_idx][:, None]
80
81     with jax.named_scope("unpermute"):
82         ff_out = jnp.take_along_axis(ff_out, isort_idx[..., None], axis=-2)
83     with jax.named_scope("expert_summing"):
84         ff_out_expert = jnp.sum(ff_out.reshape((b * s, e, d)), -2)
85         ff_out_expert = ff_out_expert.astype(cfg.dtype)
86
87     with jax.named_scope("experts_collective"):
88         # collectives
89         psum_axes = tensor_axname if expert_axname is None else (expert_axname, tensor_axname)
90         ff_out_expert = jax.lax.psum(ff_out_expert, psum_axes)
91         ff_out_expert = ff_out_expert.reshape((b, s, ff_out_expert.shape[-1]))
92         return ff_out_expert
93
94     with jax.named_scope("moe_routed_expert"):
95         x = psc(x, x_spec)
96         ff_out_expert = _expert_fn(x, we_gate, we_up, we_down, topk_weights, topk_idx)[..., : x.shape[-1]]
97     return psc(ff_out_expert, l2p("batch", "sequence", "act_embed"))
98
99
100
```


Quantization

Quantizing Matmuls

most hardware computes matmuls 2x as fast in lower precision (int8 / fp8)

full-channel quantization:

- compute-bound: $([m, k], [m, 1]) @ ([k, n], [1, n]) = ([m, k] @ [k, n]) * [m, 1] * [1, n]$
- HBM BW-bound:
 - option 1 (scale out): $[m, k] @ ([k, n], [1, n]) = ([m, k] @ [k, n]) * [1, n]$
 - option 1 (scale in): $[m, k] @ ([k, n], [k, 1]) = ([m, k] * [k, 1]) @ [k, n]$

what about subchannel quantization?

- similar, but each tile size along reduction dimension receives separate scale
- $([tile_m, tile_k] @ [tile_k, tile_n]) * [tile_m, 1] * [1, tile_n]$

Quantizing Ragged Dot

LHS: $[m, k]$
RHS: $[g, k, n]$
OUT: $[m, n]$

OUT: `ragged_dot(LHS, RHS)`

`d(LHS): ragged_dot(d(OUT), RHST)`

`d(RHS): transposed_ragged_dot(LHST, d(OUT))`

-
- quantizing very similar to matmul quantization strategies
 - full-channel quantization:
 - $([m, k], [m, 1]) @ ([k, n], [1, n]) = ([m, k] @ [k, n]) * [m, 1] * [1, n]$

-
- usually want single scale (or just a few scales, subchannel) along **reduction axis**
 - reduction axis **changes** (matmuls have the same problem)

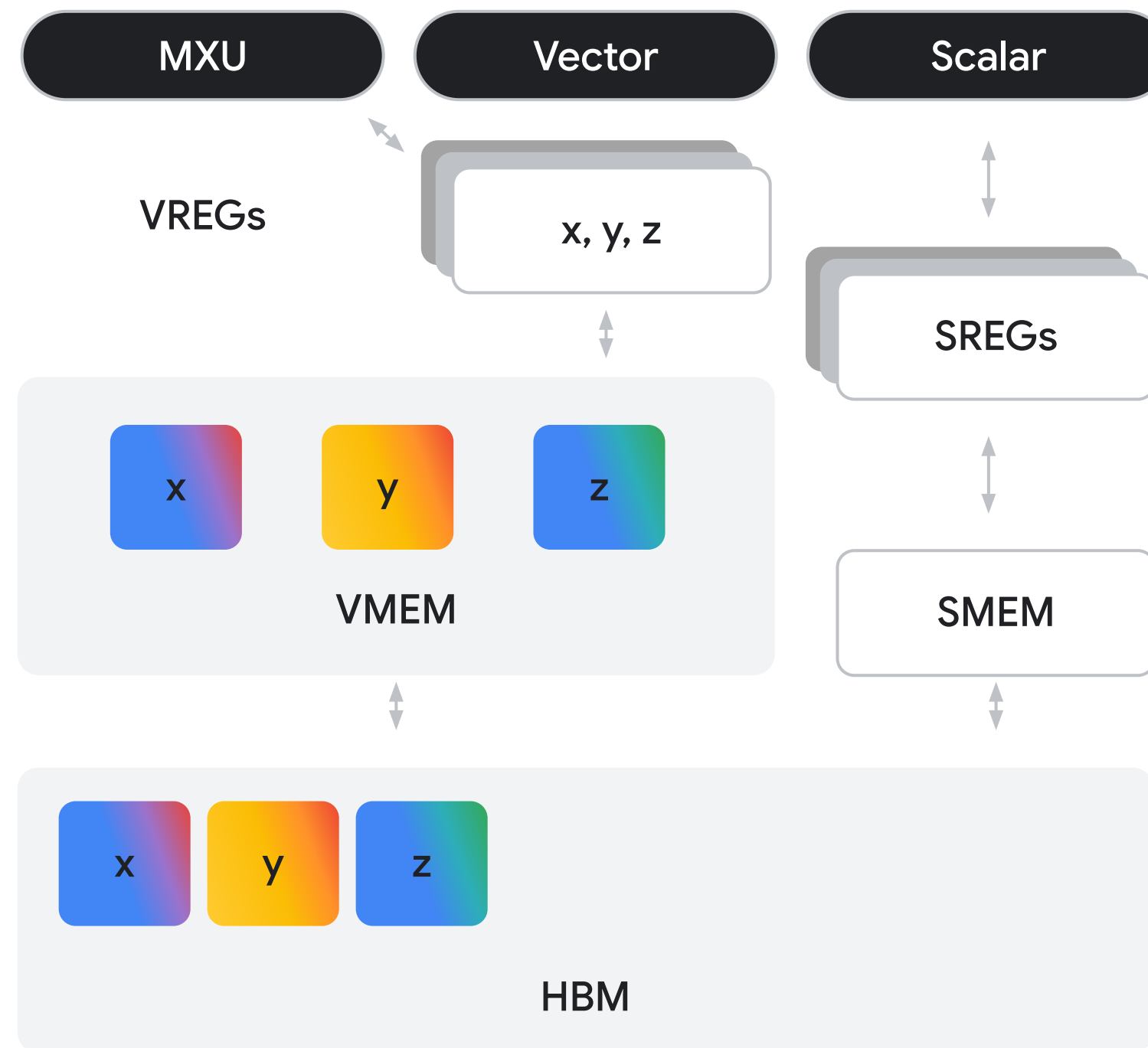
OUT: k-axis

d(LHS): n-axis

d(RHS): m-axis

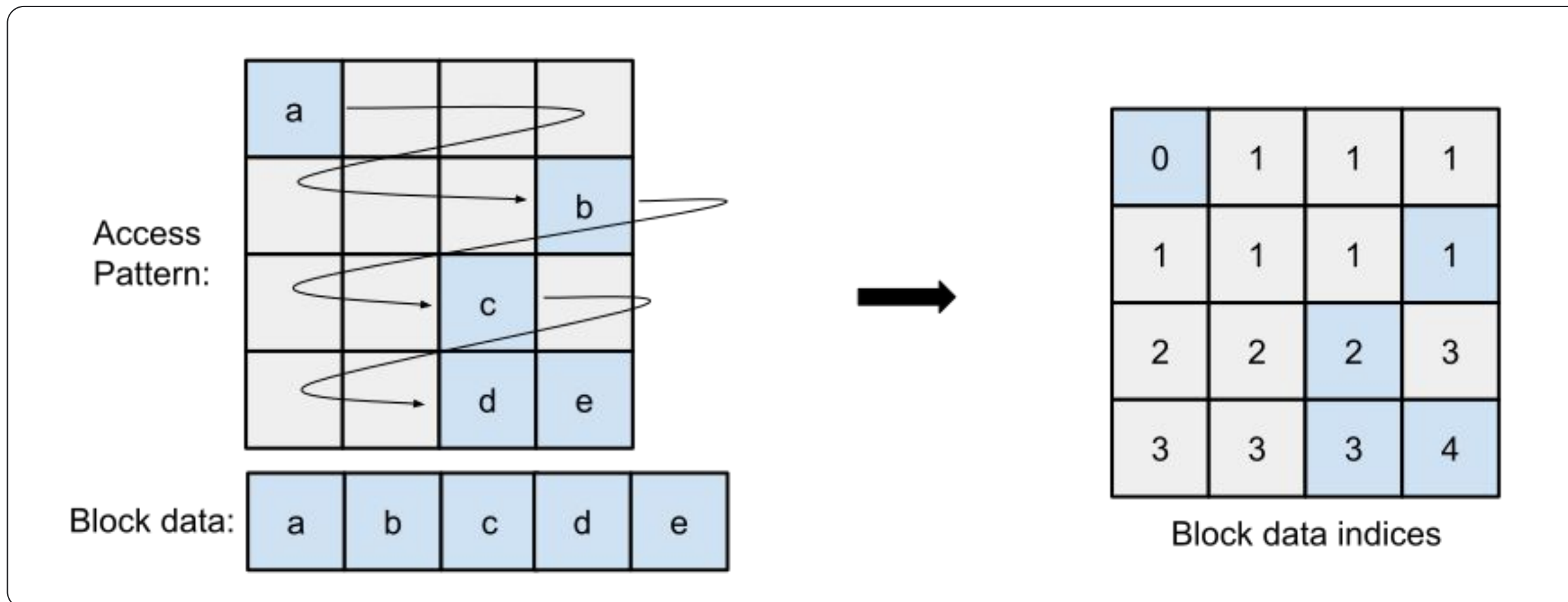
Matmuls & Megablox in Pallas

TPU Memory pipeline



Sparse computations in pallas

pltpu.PrefetchScalarGridSpec



- “2, then 2” uses no BW
- read in the next block ASAP
- small penalty for visiting block at all
- can skip computation inside kernel

- pltpu.PrefetchScalarGridSpec is super useful
- generally pre-compute a metadata lookup table, then just iterate over entries
- example: `p1.BlockSpec(..., lambda i, j, scalar1_ref, scalar2_ref: (scalar1_ref[i], scalar2_ref[j]))`
- can re-implement the prefetch scalar grid using existing pallas APIs:

https://github.com/openxla/tokamax/blob/main/tokamax/_src/mosaic_tpu.py#L126

Writing a Matmul on TPUs

- a super simple kernel
 - `p1.dot` to get correct “output” dtype
 - we’re issuing a hardware matmul
- 3D grid, no need for anything special
- **megablox** kernel is essentially this
 - apart from metadata of course
 - some row masking for tile overlap

```
def matmul_kernel(x_ref, y_ref, out_ref):
    out_ref[...] = p1.dot(x_ref[...], y_ref[...]).astype(out_ref.dtype)

def matmul(A, B, block_m, block_n, block_k):
    in_specs = [
        p1.BlockSpec((block_m, block_k), lambda i, j, k: (i, k)),
        p1.BlockSpec((block_k, block_n), lambda i, j, k: (k, j))
    ]
    out_specs = p1.BlockSpec((block_m, block_n), lambda i, j, k: (i, j))
    return p1.pallas_call(
        matmul_kernel,
        out_shape=jax.ShapeDtypeStruct((m, n), "bfloat16"),
        grid=(p1.cdiv(m, block_m), p1.cdiv(n, block_n), p1.cdiv(k, block_k)),
        in_specs=in_specs, out_specs=out_specs,
    )(A, B)
C = jax.jit(partial(matmul, block_m=2048, block_n=1024, block_k=1024))(A, B)
```

matmul: gist.github.com/rdyro/dd149ef5650f185bd96ec0666a23de9e
megablox: [github.com/openxla/tokamax ... pallas_mosaic_tpu_kernel.py](https://github.com/openxla/tokamax/blob/main/pallas_mosaic_tpu_kernel.py)

Tuning is (essentially) necessary

tune-jax

```
import tune_jax
tune_jax.logger.setLevel("INFO") # print some info for sanity
tiles = [512, 1024, 2048, 4096] # any multiple of 128 will do
hyperparams = dict(
    block_m=tiles,
    block_k=tiles,
    block_n=tiles,
)
fn = tune_jax.tune(matmul, hyperparams=hyperparams)
_ = fn(A, B) # run to tune
print(tune_jax.tabulate(fn)) # print results nicely
```

```
Compiling...: 100%|██████████| 64/64 [00:39<00:00, 1.62it/s]
Compiling...: 100%|██████████| 39/39 [00:02<00:00, 16.27it/s]
Profiling tpu: 0%|██████████| 0/5 [00:00<?, ?it/s]
Saving optimization profile to `/tmp/tuning_profile_2025-10-19_02:15:04_87s164fh`
```

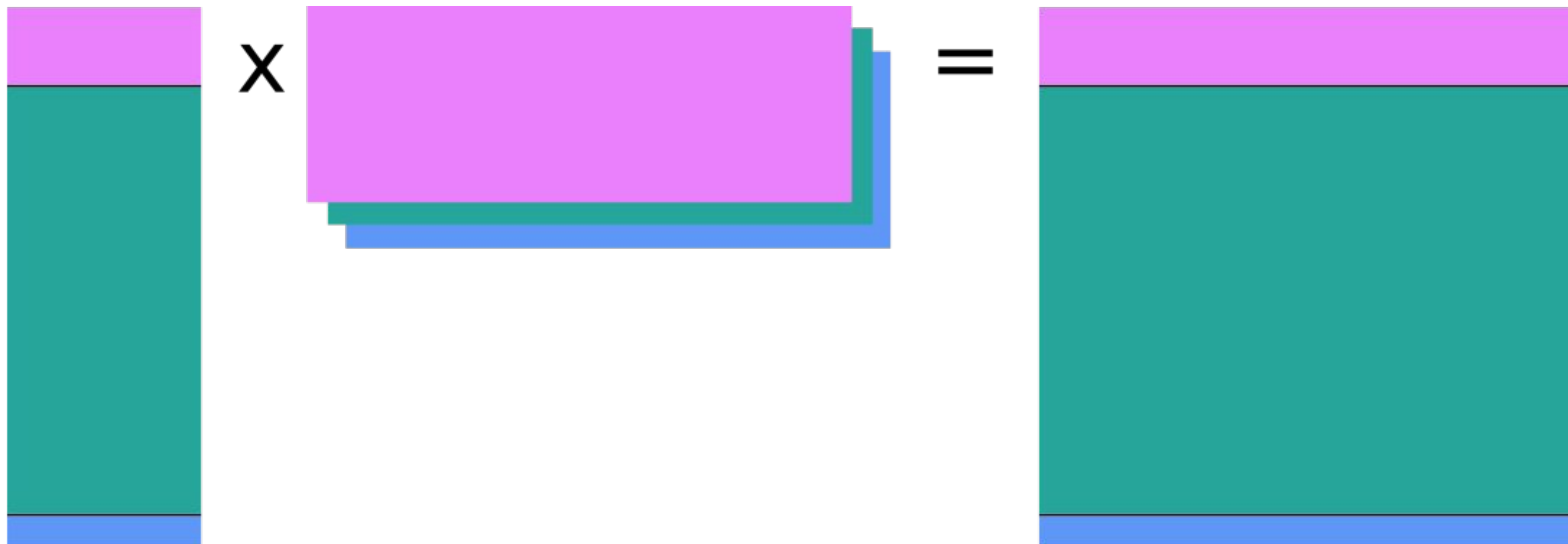
```
Profiling tpu: 100%|██████████| 5/5 [00:03<00:00, 1.66it/s]
```

id	block_m	block_k	block_n	t_mean (s)	t_std (s)
34	2048	512	2048	1.3220e-03	4.6820e-07
49	4096	512	1024	1.3353e-03	2.1772e-07
...					
4	512	1024	512	3.4113e-03	8.1699e-05
0	512	512	512	4.3446e-03	4.1993e-06

- worst result is 3.3x slower
- gap can get much wider for more complicated kernels
- tune_jax - available on PyPI
 - compiles in parallel (multi-core CPU speedup)
 - timing via automatic xprof parsing

future work & ragged dot on GPU

Future work



- collective fusions into the ragged dot kernel
- efficient quantization support
 - reusing quantized matrices in the backwards pass
 - in-kernel dynamic quantization


ragged dot on GPU


- originally in CUTLASS/CUBLAS
- kernel languages these days mostly
 - (many) triton implementations
 - e.g., (for JAX) https://github.com/rdyro/gpu_ragged_dot
 - not particularly efficient (no wgmma)
- Mosaic GPU implementation
 - matmul itself more complicated (166 lines)
 - more dynamic control over tiles
 - [github.com/jax-ml/jax ... blackwell_ragged_dot_mgpu.py](https://github.com/jax-ml/jax...blackwell_ragged_dot_mgpu.py)


Extras


jax-llm-examples
















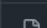


end-to-end inference


 **jax-llm-examples** Public


 main


 5 Branches


 0 Tags

 rdyro adding splash attention to qwen3 ✓	ec8398a · 2 weeks ago	 86 Commits
 .github/workflows	First draft of the GPT OSS model	2 months ago
 deepseek_r1_jax	serving cleanup	last month
 gpt_oss	enable prefill splash attention in gpt_oss	last month
 kimi_k2	Implement ring-buffer attention for all models	3 months ago
 llama3	serving cleanup	last month
 llama4	serving cleanup	last month
 misc	Add TPU creation/deletion utilities; by gpolovets1	5 months ago
 qwen3	adding splash attention to qwen3	2 weeks ago
 serving	enable prefill splash attention in gpt_oss	last month
 .gitignore	serving cleanup	last month
 .pre-commit-config.yaml	Add pre-commit, ruff linter and github CI	4 months ago
 CONTRIBUTING.md	Initial commit of repo basics	8 months ago
 LICENSE	Initial commit of repo basics	8 months ago
 README.md	First draft of the GPT OSS model	2 months ago
 multi_host_README.md	Add pre-commit, ruff linter and github CI	4 months ago
 pyproject.toml	Allow random quantized model initialization	3 months ago

 **README**

 Contributing

 Apache-2.0 license



JAX LLM examples

A collection (in progress) of example high-performance large language model implementations, written with JAX.

Current contents include:

- [DeepSeek R1](#)
- [Llama 4](#)
- [Llama 3](#)
- [Qwen 3](#)
- [Kimi K2](#)
- [OpenAI GPT OSS](#)

For multi-host cluster setup and distributed training, see [multi_host_README.md](#) and the [tpu_toolkit.sh script](#).

```
920 @partial(jax.shard_map, mesh=cfg.mesh, in_specs=in_specs, out_specs=out_spec, check_vma=False)
921 def _expert_fn(x, we_gate_up, we_gate_up_bias, we_down, we_down_bias, topk_weights, topk_idx):
922     (b, s, d), e = x.shape, cfg.moe_experts_per_tok
923     expert_idx = jax.lax.axis_index(expert_axname) if expert_axname is not None else 0
924     tensor_idx = jax.lax.axis_index(tensor_axname) if tensor_axname is not None else 0
925     topk_idx_ = topk_idx.reshape(-1)
926     valid_group_mask_ = (topk_idx_ >= expert_size * expert_idx) & (topk_idx_ < expert_size * (expert_idx + 1))
927     expert_mapped_topk_idx_ = jnp.where(valid_group_mask_, topk_idx_ - expert_idx * expert_size, 2**30)
928
929     sort_idx_ = jnp.argsort(expert_mapped_topk_idx_, axis=-1) # [b * s * e]
930     isort_idx_ = jnp.argsort(sort_idx_)
931
932     if cfg.ep_strategy == "prefill":
933         truncate_size = round(2 * sort_idx_.size / expert_count)
934         sort_idx_, isort_idx_ = sort_idx_[:truncate_size], isort_idx_[:truncate_size]
935
936     topk_idx_sort_ = topk_idx[sort_idx_] # [b * s * e]
937     expert_mapped_topk_idx_sort_ = expert_mapped_topk_idx[sort_idx_]
938     valid_group_mask_sort_ = expert_mapped_topk_idx_sort_ < 2**30
939     expert_mapped_topk_idx_sort_ = jnp.where(expert_mapped_topk_idx_sort_ < 2**30, expert_mapped_topk_idx_sort_, 0)
940
941     # equivalent to:
942     # ...
943     # x_repeat_ = jnp.repeat(x.reshape((-1, x.shape[-1])), e, axis=0)
944     # x_repeat_sort_ = jnp.take_along_axis(x_repeat_, sort_idx[:, None], axis=-2) # [b * s, d]
945     # ...
946     x_repeat_sort_ = jnp.take_along_axis(x.reshape((-1, x.shape[-1])), sort_idx[:, None] // e, axis=-2)
947     # [b * s * e, d] # "// e" is an index trick to avoid jnp.repeat
948
949     group_sizes = jnp.bincount(topk_idx_sort_, length=cfg.moe_num_experts)
950     group_sizes_shard = jax.lax.dynamic_slice_in_dim(group_sizes, expert_idx * expert_size, expert_size, 0)
951
952     with jax.named_scope("we_gate"):
953         ff_gate_up = _moe_gmm(x_repeat_sort_, we_gate_up, group_sizes_shard, expert_mapped_topk_idx_sort_, cfg)
954         ff_gate_up = ff_gate_up + we_gate_up_bias[expert_mapped_topk_idx_sort_, :]
955         ff_gate = jnp.clip(ff_gate_up[..., ::2], max=cfg.moe_gate_up_limit)
956         ff_up = jnp.clip(ff_gate_up[..., 1::2], min=-cfg.moe_gate_up_limit, max=cfg.moe_gate_up_limit)
957         ff_gate_up = (ff_up + 1) * (ff_gate * jax.nn.sigmoid(ff_gate * cfg.moe_gate_up_alpha))
958         ff_gate_up = jnp.where(valid_group_mask_sort_[..., None], ff_gate_up, 0)
959     with jax.named_scope("we_down"):
960         ff_out = _moe_gmm(ff_gate_up, we_down, group_sizes_shard, expert_mapped_topk_idx_sort_, cfg)
961         ff_out = ff_out + (tensor_idx == 0) * we_down_bias[expert_mapped_topk_idx_sort_, :]
962         ff_out = jnp.where(valid_group_mask_sort_[..., None], ff_out, 0) # expensive
963
964     if cfg.ep_strategy == "prefill":
965         rs_shape = math.ceil((ff_out.shape[-1] // tensor_count) / 256) * 256 * tensor_count
966         pad_size = rs_shape - ff_out.shape[-1]
967         ff_out = jnp.pad(ff_out, ((0, 0), (0, pad_size)))
968         ff_out = jax.lax.psum_scatter(ff_out, axis_name=tensor_axname, scatter_dimension=1, tiled=True)
969         ff_out = ff_out * topk_weights.reshape((-1, sort_idx.size // e, None))
```


Extras

nvtop: TPU support

<https://github.com/syllo/nvtop>

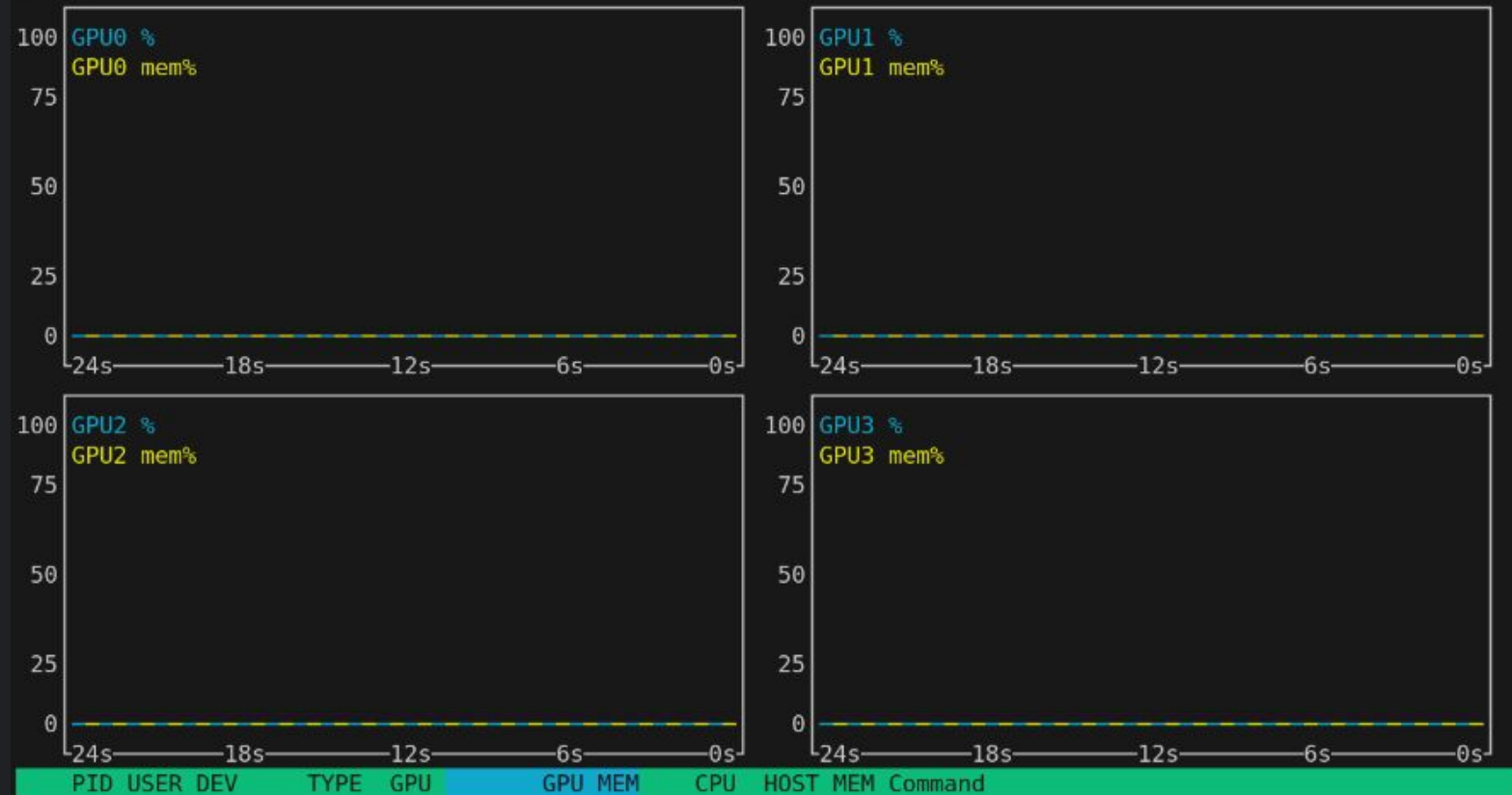
<https://github.com/rdyro/libtpuinfo>

```
Device 0 [TPU0] PCIe GEN N/A RX: N/A TX: N/A
GPU N/A MHz MEM N/A MHz TEMP N/A°C FAN N/A POW N/A W
GPU[ 0%] MEM[ 0.000Gi/31.246Gi]
```

```
Device 1 [TPU1] PCIe GEN N/A RX: N/A TX: N/A
GPU N/A MHz MEM N/A MHz TEMP N/A°C FAN N/A POW N/A W
GPU[ 0%] MEM[ 0.000Gi/31.246Gi]
```

```
Device 2 [TPU2] PCIe GEN N/A RX: N/A TX: N/A
GPU N/A MHz MEM N/A MHz TEMP N/A°C FAN N/A POW N/A W
GPU[ 0%] MEM[ 0.000Gi/31.246Gi]
```

```
Device 3 [TPU3] PCIe GEN N/A RX: N/A TX: N/A
GPU N/A MHz MEM N/A MHz TEMP N/A°C FAN N/A POW N/A W
GPU[ 0%] MEM[ 0.000Gi/31.246Gi]
```



```
wget https://github.com/rdyro/libtpuinfo/releases/download/v0.0.1/libtpuinfo-linux-x86_64.so
sudo mv libtpuinfo-linux-x86_64.so /lib/libtpuinfo.so
git clone https://github.com/Syllo/nvtop.git
cd nvtop && mkdir build && cd build && cmake -DTPU_SUPPORT=ON .. && make
sudo cp ./src/nvtop /usr/local/bin
```