# Uncovering Nested Data Parallelism and Data Reuse with FractalTensor

Microsoft Research

The Imperative implementation of
stacked RNN

```
List<List<Vector>> xss  // input sequences
List<Matrix> ws // learnable weights
Vector I  // initial state, constant
List<List<List<Vector>>> ysss //output
```

$for\ 0 \leq j < D$         // stacked $depth$
  $for\ 0 \leq k < L$    // sequence length
    $for\ 0 \leq i < N$  // batch
      $if\ j == 0\ \&\&\ k == 0$
        $s = I$
        $x = xss[i][k]$

      $elif\ j > 0\ \&\&\ k == 0$

        $s = I$
        $x = ysss[i][j-1][k]$

      $elif\ j == 0\ \&\&\ k > 0$

        $s = ysss[i][j][k-1]$
        $x = xs[i]$

      $else$
        $s = ysss[i][j-1][k]$
        $x = ysss[i][j][k-1]$
      // user-defined cell function
      $yss[i][j][k] = x@ws[j] + s$

- ## Strong expressiveness
  - Most intuitive way of implementation

- ## Less effectiveness
  - Three-level loops, many branches
  - Hard to analyze data dependency
  - Poor performance

# Tensor Operator: The Dilemma of Expressiveness and Effectiveness

Three ways to define tensor operators

1

$$\textit{for } 0 \le j < D$$
$$\quad \textit{for } 0 \le k < L$$
$$\quad\quad \textit{for } 0 \le i < N$$

**RNN Cell as Op**

The bottom level loop as an operator
    Most expressive
    Less effectiveness

2

$$\textit{for } 0 \le j < D$$
$$\quad \textit{for } 0 \le k < L$$
$$\quad\quad \textit{for } 0 \le i < N$$

**RNN Layer as Op**

Wrap the two-level loops as an operator
    Less expressive
    More effectiveness

3

$$\textit{for } 0 \le j < D$$
$$\quad \textit{for } 0 \le k < L$$
$$\quad\quad \textit{for } 0 \le i < N$$

**Stacked RNN Layers as Op**

The whole three-level loops as an operator
    Worst expressive
    Most effectiveness

# The Impact of Expressiveness on Effectiveness: Stacked RNN



The performance varies for stacked RNN layer.

As the implementations become less expressive, the performance increases, i.e., more effective



O User-defined RNN cell

---------  The hyperplane for largest coarse-grained data parallelism

# Another Example: Flash Attention

```
List<List<List<Matrix>>> qsss, ksss, vsss
List<List<List<Matrix>>> osss
for 0 ≤ i < B          // batch size
  for 0 ≤ j < H        // heads
    for 0 ≤ m < L₁     // sequence length
```

$M_t = -\overrightarrow{\inf}_{[d_2,1]}$

$S_t = \vec{0}_{[d_2,1]}$

$O_t = \vec{0}_{[d_2,d_4]}$

```
      for 0 ≤ n < L₂   // sequence length
```

$Q = qsss[i][j][m], K = ksss[i][j][n], V = vsss[i][j][n]$ `// load from DRAM`

$T_1 = Q@K^T$

$T_2 = max(T_1)$

$T_3 = exp(T_1 - T_2)$

$T_4 = T_3@V$

$T_5 = sum(T_3, dim = -1)$

$M'_t = \text{maximum}(T_2, M_t)$

$T_6 = exp(M'_t - M_t)$

$T_7 = exp(T_2 - M'_t)$

$S_t = T_6 * S_t + T_7 * T_5$

$O_t = (O_t * S_t * T_6 + T_7 * T_4) / S_t$

$osss[i][j][n] = O_t$ `// store to DRAM`

**1. online normalization instead of compute the full batch at once**

**2. reusable data cached and computed in SRAM**

The imperative implementation of FlashAttention

- The entire algorithm is implemented into one monolithic, opaque operator

- No room for analysis and optimization

- Minor adjustment requires re-implementation
  - FlashAttention 1/2/3
  - Sacrifice expressiveness for effectiveness

The dilemma between tensor operator's expressiveness and effectiveness: a common painpoint

5

# Problem

The DAG of tensor operators, to achieve **effectiveness**, often exposes only single-level parallelism, lacking **expressiveness** and resulting in insufficient global analysis
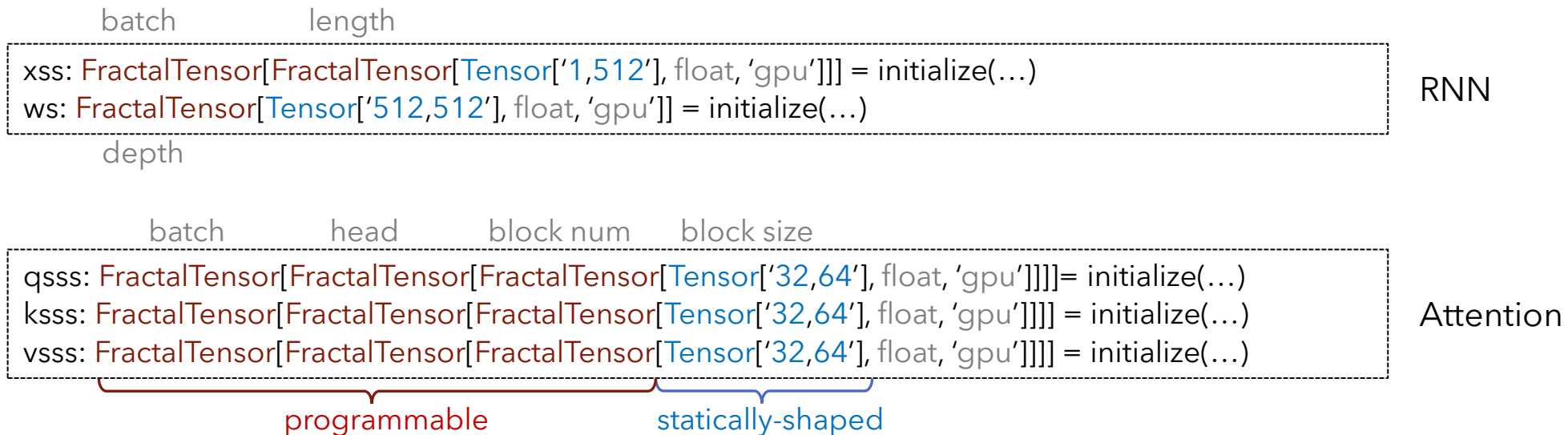
# Observation

- Diverse DNN computation patterns can be expressed by a combination of second-order array **compute operators** like map, reduce, scan, fold

- Data access patterns in DNN computation are highly stylized and can be expressed by a few first-order array **access operators**

- DNN algorithms can be expressed along tensor dimensions with **compute** and **access** operator nesting

It is possible to provide an expressive programming model for DNN, and generate effective, high-performance code.

# FractalTensor: Decompose Tensor into Nested Lists of Tensors

- FractalTensor: a list-based ADT, an element is a static-shape tensor, or a FractalTensor

- FractalTensor: decompose dimensions of a tensor into:
  - The innermost statically-shaped dimensions
  - The **programmable dimensions**

Examples: organize data as FractalTensors

batch          length

xss: FractalTensor[FractalTensor[Tensor['1,512'], float, 'gpu']]] = initialize(…)    RNN
ws: FractalTensor[Tensor['512,512'], float, 'gpu']] = initialize(…)

depth

batch        head        block num        block size

qsss: FractalTensor[FractalTensor[FractalTensor[Tensor['32,64'], float, 'gpu']]]]= initialize(…)
ksss: FractalTensor[FractalTensor[FractalTensor[Tensor['32,64'], float, 'gpu']]]] = initialize(…)    Attention
vsss: FractalTensor[FractalTensor[FractalTensor[Tensor['32,64'], float, 'gpu']]]] = initialize(…)

programmable        statically-shaped

# The FractalTensor Program

RNN: FractalTensor code

```
// N, L, D stands for batch, length, depth
xss: [N, L]float32[1,512] = … // load from storage
ws: [D]float32[512,512] = … // load from storage

// output transformed from existing FractalTensors
ysss: [N, D, L]float32[1,512] = ⋯

// map over the batch dimension of xss
ysss = map xss xs ⟹
    // scan the depth dimension of ws
    yss = ws scanl xs, (x̃s, w) ⟹
        // scan the length dimension of xss
        ys = dilate(x̃s) scanl 0, (s, x) ⟹
            // user-defined small math function
            // [1,512]=[1,512]@[512,512]+[1,512]
            y = x @ w + s
```

**Compute operators**

**Access operators**

Functional **array compute** and **access operators** are tied to programmable dimensions.

1. array compute operators
   - map, reduce, fold, scan

2. array access operators
   - contiguously linear
   - constantly strided (dilation)
   - window (convolution)
   - indirect access (gather)

1. *No explicit tensor operators*
2. *Yet analyzable: loop nest with compute and access patterns **understandable by the compiler***

# FractalTensor Code is Fully Permutable

The nested loops in FractalTensor code can be **reordered arbitrarily**[1]
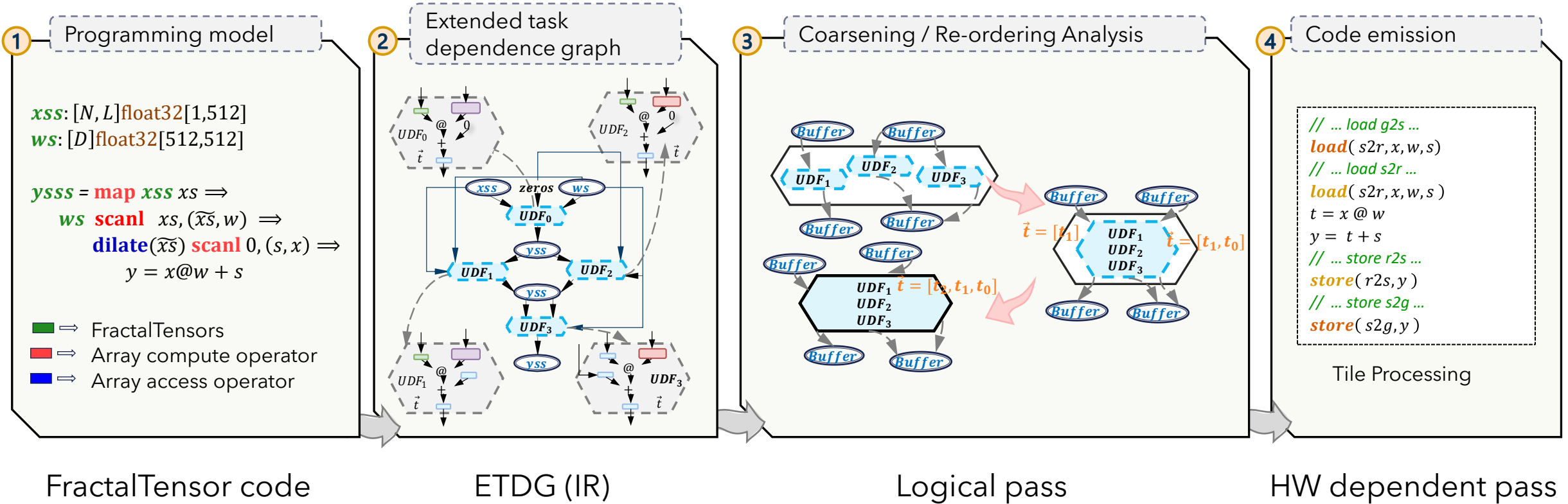
Because:

- FractalTensor code must follow SSA (single static assignment)

- Data dependence distance, regulated by array compute operators, is constant

Iteration-level data dependence can be permuted and moved to the outermost loop, allowing all inner loops to be parallel

Inner loops can then focus on data locality

1. Wolf, Michael E., and Monica S. Lam. "A loop transformation theory and an algorithm to maximize parallelism." *IEEE Transactions on Parallel & Distributed Systems* 2.04 (1991): 452-471.
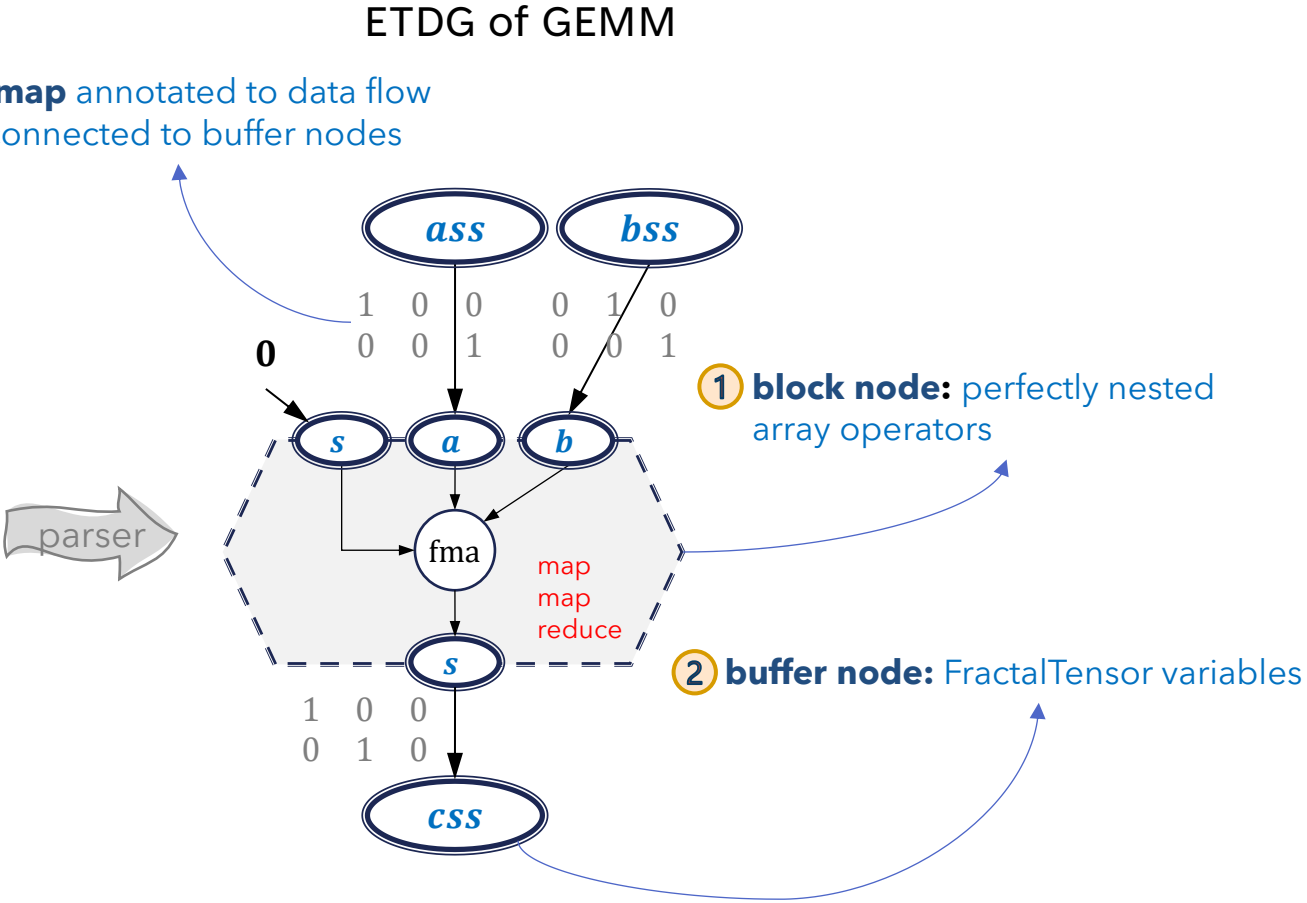
# Workflow Overview



① Programming model — FractalTensor code

② Extended task dependence graph — ETDG (IR)

③ Coarsening / Re-ordering Analysis — Logical pass

④ Code emission — HW dependent pass

# Extended Task Dependence Graph

ETDG, parsed from the FractalTensor program, reflects its nested structure.

ETDG of GEMM

③ **access map** annotated to data flow edges connected to buffer nodes

GEMM: FractalTensor code

```
1  ass:[16,8]float32[32,32] = ...
2  bss:[8,16]float32[32,32] = ...
3  css:[16,16]float32[32,32]  // output
4
5  css = ass.map as ⇒
6      css = bss.map bs ⇒
7          c = zip(as,bs).reduce 0, (s,(a,b)) ⇒
8              c = a @ b + s
```

parser

$$\begin{matrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{matrix} \qquad \begin{matrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{matrix}$$

**0**

① **block node:** perfectly nested array operators

map
map
reduce

$$\begin{matrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{matrix}$$
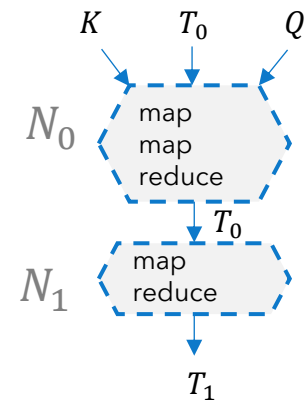
② **buffer node:** FractalTensor variables

# ETDG Transformation: Coarsening

- Multi-level loops introduce control overhead
- Coarsening reduces the overhead

$N_0$   *for $i_3$ in [0, $d_3$) // map*
     *for $i_2$ in [0, $d_2$) // map*

       *for $i_1$ in [0, $d_1$) // reduce*, 0, +

         $T_0[i_3, i_2] = T_0[i_3, i_2] + Q[i_3, i_1] * K[i_2, i_1]$

$N_1$   *for $j_2$ in [0, $d_3$) // map*
     *for $j_1$ in [0, $d_2$) // reduce*, -inf, max

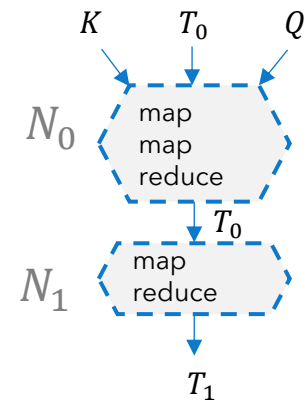       $T_1[j_2, 1] = \max(T_0[j_2, j_1 - 1], T_0[j_2, j_1])$

# ETDG Transformation: Coarsening

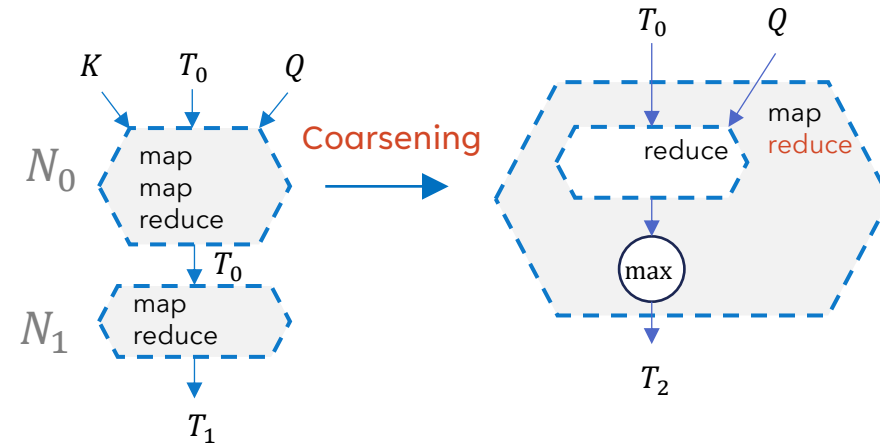- Multi-level loops introduce control overhead
- Coarsening reduces the overhead

$N_0$ | for $i_3$ in [0, $d_3$) // map
     |   for $i_2$ in [0, $d_2$) // map

Same shape

$K$    $T_0$    $Q$

        for $i_1$ in [0, $d_1$) // reduce, 0, +

$$T_0[i_3, i_2] = T_0[i_3, i_2] + Q[i_3, i_1] * K[i_2, i_1]$$

$N_0$    map / map / reduce

$T_0$

$N_1$ | for $j_2$ in [0, $d_3$) // map
     |   for $j_1$ in [0, $d_2$) // reduce, -inf, max

Same shape

$N_1$    map / reduce

$$T_1[j_2, 1] = \max(T_0[j_2, j_1 - 1], T_0[j_2, j_1])$$

$T_1$

# ETDG Transformation: Coarsening

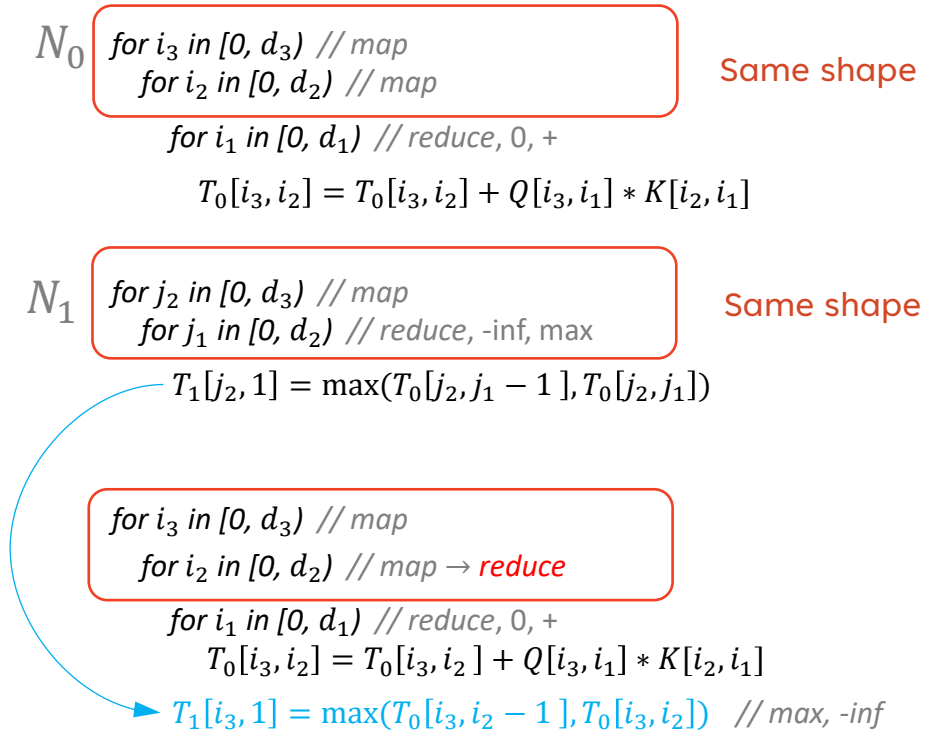- Multi-level loop nest introduce control overheads
- Coarsening reduces the overhead

$N_0$
$$\begin{aligned}
&\text{for } i_3 \text{ in } [0, d_3) \; // \; map \\
&\quad \text{for } i_2 \text{ in } [0, d_2) \; // \; map \\
&\qquad \text{for } i_1 \text{ in } [0, d_1) \; // \; reduce, 0, + \\
&\qquad\quad T_0[i_3, i_2] = T_0[i_3, i_2] + Q[i_3, i_1] * K[i_2, i_1]
\end{aligned}$$

Same shape

$N_1$
$$\begin{aligned}
&\text{for } j_2 \text{ in } [0, d_3) \; // \; map \\
&\quad \text{for } j_1 \text{ in } [0, d_2) \; // \; reduce, \text{-inf, max} \\
&\qquad T_1[j_2, 1] = \max(T_0[j_2, j_1 - 1], T_0[j_2, j_1])
\end{aligned}$$

Same shape

$$\begin{aligned}
&\text{for } i_3 \text{ in } [0, d_3) \; // \; map \\
&\quad \text{for } i_2 \text{ in } [0, d_2) \; // \; map \rightarrow reduce \\
&\qquad \text{for } i_1 \text{ in } [0, d_1) \; // \; reduce, 0, + \\
&\qquad\quad T_0[i_3, i_2] = T_0[i_3, i_2] + Q[i_3, i_1] * K[i_2, i_1] \\
&\qquad T_1[i_3, 1] = \max(T_0[i_3, i_2 - 1], T_0[i_3, i_2]) \; // \; max, \text{-inf}
\end{aligned}$$

# ETDG Transformation: Access Reordering

Enhance exploitable data parallelism and locality

- Permuate FractalTensor and move all data dependencies to the outermost dimension (Hyperplane method[1])
- Dimensions with data reuse moved to innermost to enhance locality (null space of access matrix to detect data reuse dimension[2])

RNN example

scan

original access order

*Access reordering*

new access order

scan

Nodes with the same color (a hyperplane) can be executed in parallel

1. Lamport, Leslie. "The parallel execution of DO loops." Communications of the ACM 17.2 (1974): 83-93.
2. Wolf, Michael E., and Monica S. Lam. "A data locality optimizing algorithm." Proceedings of the ACM SIGPLAN 1991 conference on Programming language design and implementation. 1991.
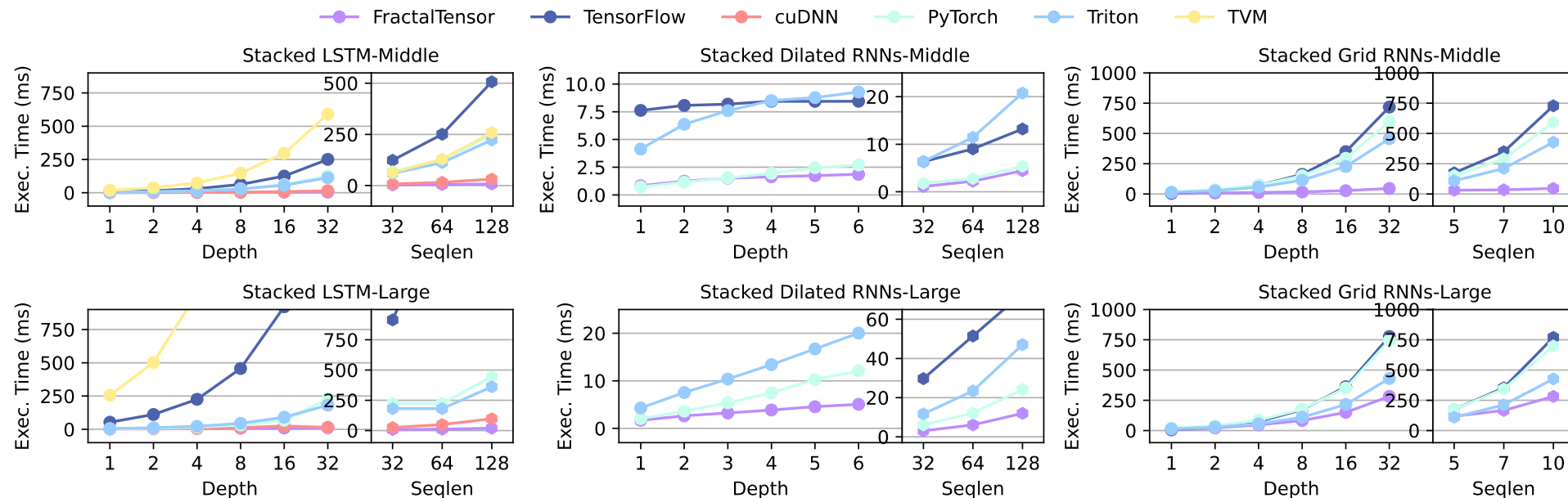
1. Hardware bottom-up tile processing library

   - BaseTile optimizes compute and memory usage aligned with TensorCore

2. Decompose buffer nodes into BaseTiles

3. Materialize access maps into load/store tiles

```
for (int k1 = 0; k1 < GIteratorA::sc1; ++k1) {
    g2s_a(gAs(k1), sA);        load tiles from global
    g2s_b(gBs(k1), sB);        to shared memory
    __copy_async();
    __syncthreads();


    for (int k2 = 0; k2 < SIteratorA::sc1; ++k2) {
        s2r_a(sAs(k2), rA);    load tiles from shared
        s2r_b(sBs(k2), rB);    memory to register


        compute::gemm(rA, rB, acc);    compute
    }
}
r2s_c(acc, sC);    store tiles from register to shared memory
__syncthreads();
s2g_c(sC, gC);     store tiles from shared to global memory
```

GEMM with our tile library

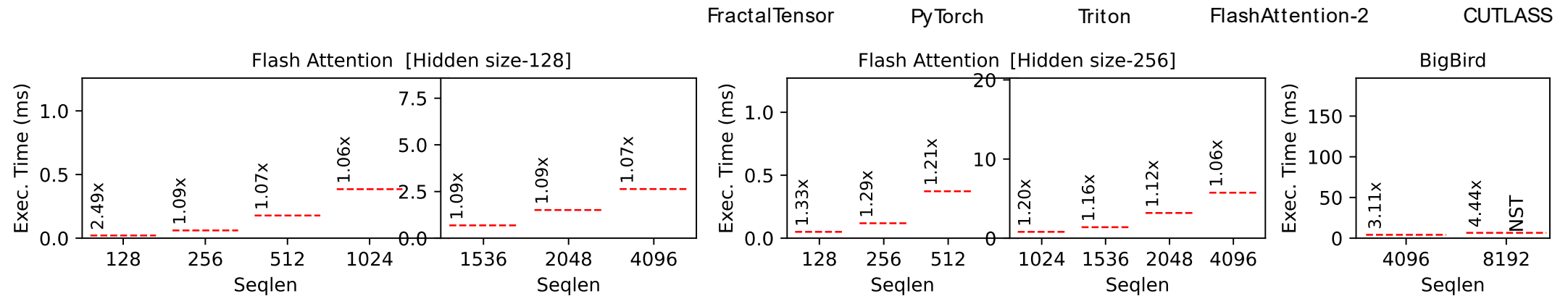# Overall Performance on Stacked RNN and the Variants



Existing practice
- Only standard stacked LSTM is optimized through the vendor library
- Slight algorithmic changes negate the optimization

In FractalTensor
- Optimizations apply to patterns commonly found in new RNN variants
- Reordering analysis identifies exploitable data parallelism, ensuring stacked RNN performance regardless with variants (e.g., changes in depth)

# Overall Performance on Flash Attention and the Variants

FractalTensor   PyTorch   Triton   FlashAttention-2   CUTLASS

**Flash Attention [Hidden size-128]**

Exec. Time (ms): 2.49x, 1.09x, 1.07x, 1.06x / 1.09x, 1.09x, 1.07x
Seqlen: 128, 256, 512, 1024 / 1536, 2048, 4096

**Flash Attention [Hidden size-256]**

Exec. Time (ms): 1.33x, 1.29x, 1.21x / 1.20x, 1.16x, 1.12x, 1.06x
Seqlen: 128, 256, 512 / 1024, 1536, 2048, 4096

**BigBird**

Exec. Time (ms): 3.11x, 4.44x, NST
Seqlen: 4096, 8192

Existing practice
- FlashAttention's online normalization algorithm is hard to express as a DAG
- Manual GPU memory optimization is complex due to TensorCore details

In FractalTensor
- Online normalization algorithm fits naturally into map and reduce operator nesting
- Tile library abstracts the hardware programming model and maximizes hardware usage
- Near-direct translation achieves *better* performance in multiple configurations

18

# Conclusion

The dilemma of tensor operator: expressiveness and effectiveness

FractalTensor can solve the dilemma by

- An ADT to capture the key characteristics of tensors in DNN

- A set of array compute and access operator to compose arbitrary nested DNN structure based on FractalTensor

- Evaluation demonstrates that FractalTensor codes can achieve both expressiveness and effectiveness