# Thesis contributions

- APPy: Annotated Parallelism for Python on GPUs
  - [CC24] Parallelize Python loops and tensor expressions on GPUs

- **ReACT: Redundancy-Aware Code Generation for Tensor Expressions**
  - **[PACT22] Redundancy elimination when fusing sparse/dense tensor operators**

- Intrepydd: Performance, Productivity, and Portability for Data Science Application Kernels
  - [Onward!20] Compile Python/NumPy to C++ with high-level optimizations

Georgia Tech

# Problem statement: desired input and output

- Desired input: operator program in Python (can be sparse)

- Desired output: fused CPU kernel with reduced redundant memory accesses and computations

```python
1  def sddmm(sp_A, B, C):
2      return sp_A * (B @ C)
3
4  def spmm_mm(sp_A, B, C):
5      return sp_A @ (B @ C)
6
7  def norm_row(sp_A):
8      return sp_A / sum(sp_A, axis=1)
```

```c
130      #pragma omp parallel
131  {
132
133      auto T = new double [D2_dimension]();
134      int jT = 0;
135      #pragma omp for schedule(static)
136      for (int32_t i = 0; i < C1_dimension; i++) {
137          for (int32_t k = 0; k < D1_dimension; k++) {
138              int32_t kC = i * C2_dimension + k;
139              jT = 0;
140              for (int32_t jB = B2_pos[i]; jB < B2_pos[(i + 1)]; jB++) {
141                  int32_t j = B2_crd[jB];
142                  int32_t jD = k * D2_dimension + j;
143                  T[jT] += C_vals[kC] * D_vals[k * D2_dimension + j];
144                  jT++;
145              }
146          }
147
148          jT = 0;
149          for (int32_t jB = B2_pos[i]; jB < B2_pos[(i + 1)]; jB++) {
150              int32_t j = B2_crd[jB];
151              A_vals[jB] += B_vals[jB] * T[jT];
152              T[jT] = 0;
153              jT++;
154          }
155      }
156
157      delete T;
158  }
```

Georgia Tech

# Limitations with State-of-the-art

- TACO
  - A code generator for arbitrary sparse/dense tensor algebra expressions
  - **maximal fusion** is implicit during code generation

- Limitations
  - Maximal fusion may introduce some types of redundant memory accesses and computations
  - Maximal fusion cannot properly fuse certain reduction expressions

```python
1   def sddmm(sp_A, B, C):
2       return sp_A * (B @ C)
3
4   def spmm_mm(sp_A, B, C):
5       return sp_A @ (B @ C)
6
7   def norm_row(sp_A):
8       return sp_A / sum(sp_A, axis=1)
```

Maximal fusion does not work because
it requires the "/" operator to be distributive over a summation

Georgia Tech.

# Redundancy types identified

- **Type 1** (Reduction Redundancy): When multiple multiply-add operations are performed instead of multiple adds followed by a single multiply (distributive law).

- **Type 2** (Loop-Invariant Redundancy): When a loop invariant expression is introduced (could be invariant in a non-innermost loop) due to maximum fusion.

- **Type 3** (Load-Store Redundancy): When some values are stored and loaded in separate loops, and the loads/stores can be eliminated after fusion --- a classical benefit of loop fusion.

- **Type 4** (Dead-Value Redundancy): When some values are computed but not used later on (e.g., when multiplying with 0s in a sparse tensor) --- another classical benefit of loop fusion.

Georgia Tech.

# (Type 1) Reduction redundancy

Input: c = b * sum(A, axis=1)

**With redundancy (due to maximal fusion)**

```
1.  for (int i = 0; i < NI; i++) {
2.    double s = 0;
3.    double bi = b[i];
4.    for (int j = 0; j < NJ; j++) {
5.      s += A[i,j] * bi;
6.      ...
7.    }
8.    ...
9.  }
```

**Without redundancy**

```
1.  for (int i = 0; i < NI; i++) {
2.    double s = 0;
3.    for (int j = 0; j < NJ; j++) {
4.      s += A[i,j];
5.      ...
6.    }
7.    s = s * B[i];
8.    ...
9.  }
```

Reduced number of multiplications in the innermost loop!

Georgia Tech.

# (Type 2) Loop-Invariant redundancy

Input: A = (B + E) * (C @ D)

### With redundancy (due to maximal fusion)

```
1.   for (int i = 0; i < NI; i++)
2.     for (int k = 0; k < NK; k++)
3.       for (int j = 0; j < NJ; j++)
4.         A[i,j] += (B[i,j] + E[i,j]) * \ (C[i,k] * D[k,j]);
```

### Without redundancy

```
1.    double* T = new double[NJ];
2.    for (int i = 0; i < NI; i++) {
3.      for (int j = 0; j < NJ; j++) {
4.        T[j] = B[i,j] + E[i,j];
5.      }
6.      for (int k = 0; k < NK; k++) {
7.        for (int j = 0; j < NJ; j++) {
8.          A[i,j] += T[j] * (C[i,k] * D[k,j]);
9.        }
10.   }
11. }
```

$B[i,j] + E[i,j]$ is no longer repeatedly calculated for different *k* iterations!

Georgia Tech

# (Type 3) Load-Store redundancy

Input: s = sum(A, axis=1); B = A / s[:, None]

## With redundancy (due to no fusion)

```
1.   double* s = new double[NI];
2.   // Operator 1
3.   for (int i = 0; i < NI; i++) {
4.      s[i] = 0;
5.      for (int j = 0; j < NJ; j++) {
6.         s[i] += A[i,j];
7.      }
8.   }
9.   // Operator 2
10.  for (int i = 0; i < NI; i++) {
11.     for (int j = 0; j < NJ; j++) {
12.        B[i,j] = A[i,j] / s[i];
13.     }
14.  }
```

## Without redundancy

```
1.   // Operator 1 and 2 fused
2.   for (int i = 0; i < NI; i++) {
3.      double s = 0;
4.      for (int j = 0; j < NJ; j++) {
5.         s += A[i,j];
6.      }
7.
8.      for (int j = 0; j < NJ; j++) {
9.         B[i,j] = A[i,j] / s;
10.     }
11.  }
```

A[i,j] and s[i] now have reduced reuse distance, which leads to better locality!

Georgia Tech.

# (Type 4) Dead-Value redundancy

Input: B = where(A < 0, alpha * A, A)

### With redundancy (due to no fusion)

```
1.   // Operator 1
2.   double* tmp = new double[NI];
3.   for (int i = 0; i < NI; i++) {
4.      tmp[i] = alpha * A[i];
5.   }
6.   // Operator 2
7.   for (int i = 0; i < NI; i++) {
8.      if (A[i] < 0) {
9.        B[i] = tmp[i];
10.     }
11.    else {
12.       B[i] = A[i];
13.    }
14. }
```

Not all values in array tmp are useful!

### Without redundancy

```
1.   // Operator 1 and 2 fused
2.   for (int i = 0; i < NI; i++) {
3.      if (A[i] < 0) {
4.        B[i] = alpha * A[i];
5.      }
6.      else {
7.        B[i] = A[i];
8.      }
9.   }
```
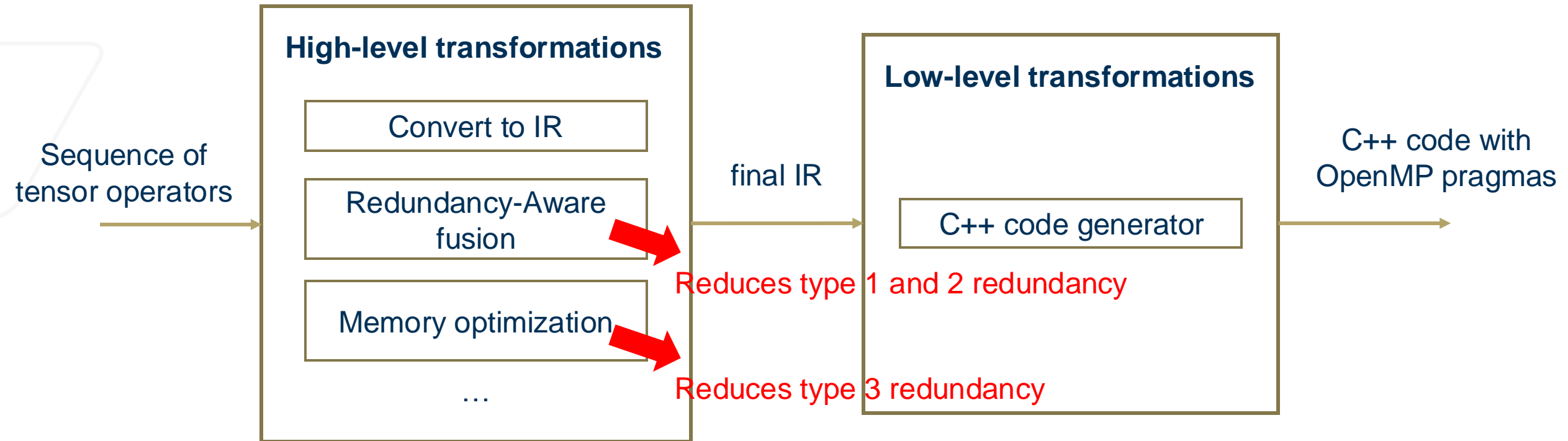
The use of tmp is now eliminated, which reduces redundant computations and memory accesses!

Georgia Tech.

# Redundancies eliminated by each approach

| Redundancy type | ReACT (this work) | TACO | SciPy |
|---|---|---|---|
| Reduction (type 1) | Yes | No | Yes |
| Loop invariant (type 2) | Yes | No | Yes |
| Load store (type 3) | Yes | Partially | No |
| Dead value (type 4) | Yes | Yes | No |

Georgia Tech

# How is ReACT able to reduce these redundancies?

Transformation passes are redundancy-aware

**High-level transformations**

Convert to IR

Sequence of tensor operators

Redundancy-Aware fusion

Reduces type 1 and 2 redundancy

Memory optimization

Reduces type 3 redundancy

…

final IR

**Low-level transformations**

C++ code generator

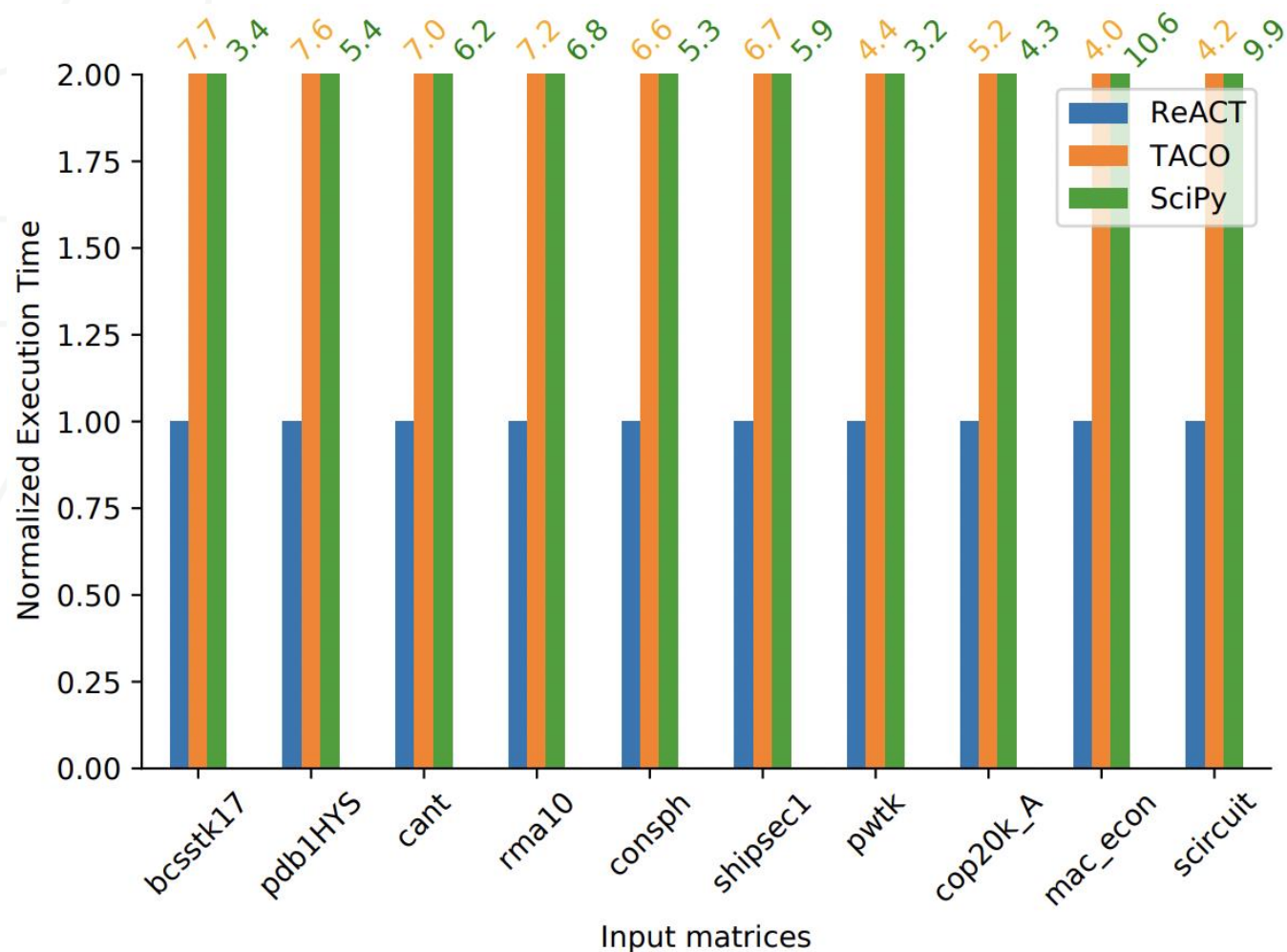C++ code with OpenMP pragmas

Georgia Tech.

# Performance evaluation

- Test machine
  - 16-core Intel(R) Xeon(R) 2.20GHz CPU
  - `OMP_NUM_THREADS` is set to 16
- Kernels (all kernels have at least 2 operators)
  - SpMM-MM (sparse-dense matmul followed by dense matmul)
  - SDDMM/Masked MM (a dense matmul followed by a dense-sparse element-wise mul)
  - Sparse-softmax (row-wise softmax on a sparse matrix)
    - Expressed using basic operators such as exp, sum, divide etc
- Sparse matrices
  - A collection of real-world matrices from SuiteSparse
  - All sparse matrices are in CSR format
- Comparisons
  - ReACT (our approach)
  - TACO (SOTA compiler)
  - SciPy.sparse (SOTA library)

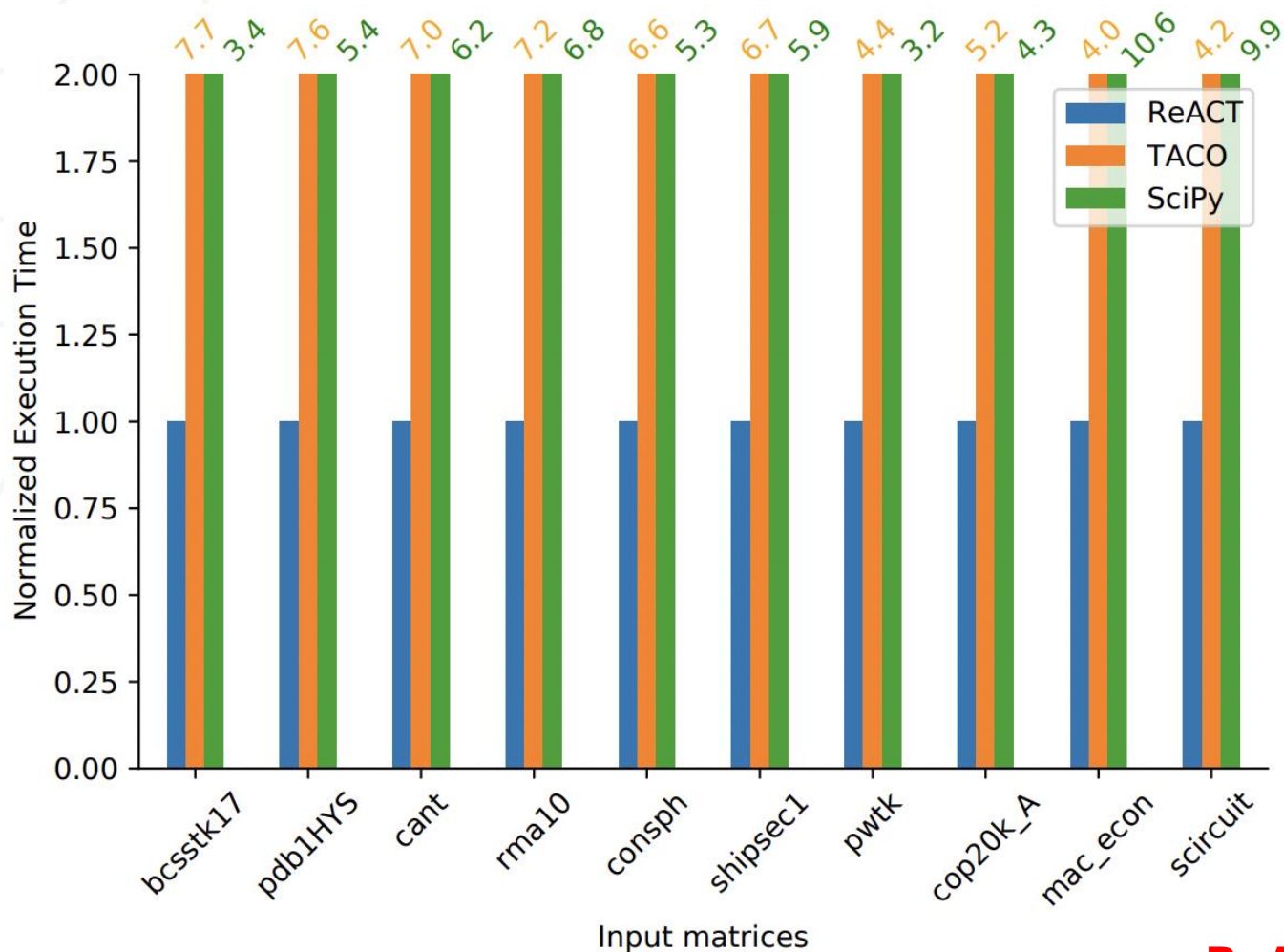Georgia Tech®

# SpMM-MM results – 5.9x faster than TACO

"No" is good here!



(b) GNN-kernel1 (NH=256, NJ=16)

| Redundancy types present | TACO output | ReACT output |
|---|---|---|
| Type 1 | Yes | No |
| Type 2 | Yes | No |
| Type 3 | No | No |
| Type 4 | No | No |

Code time complexity is reduced from $O(NNZ * NH * NJ)$ (TACO) to $O(NI * NH * NJ)$ (ReACT)

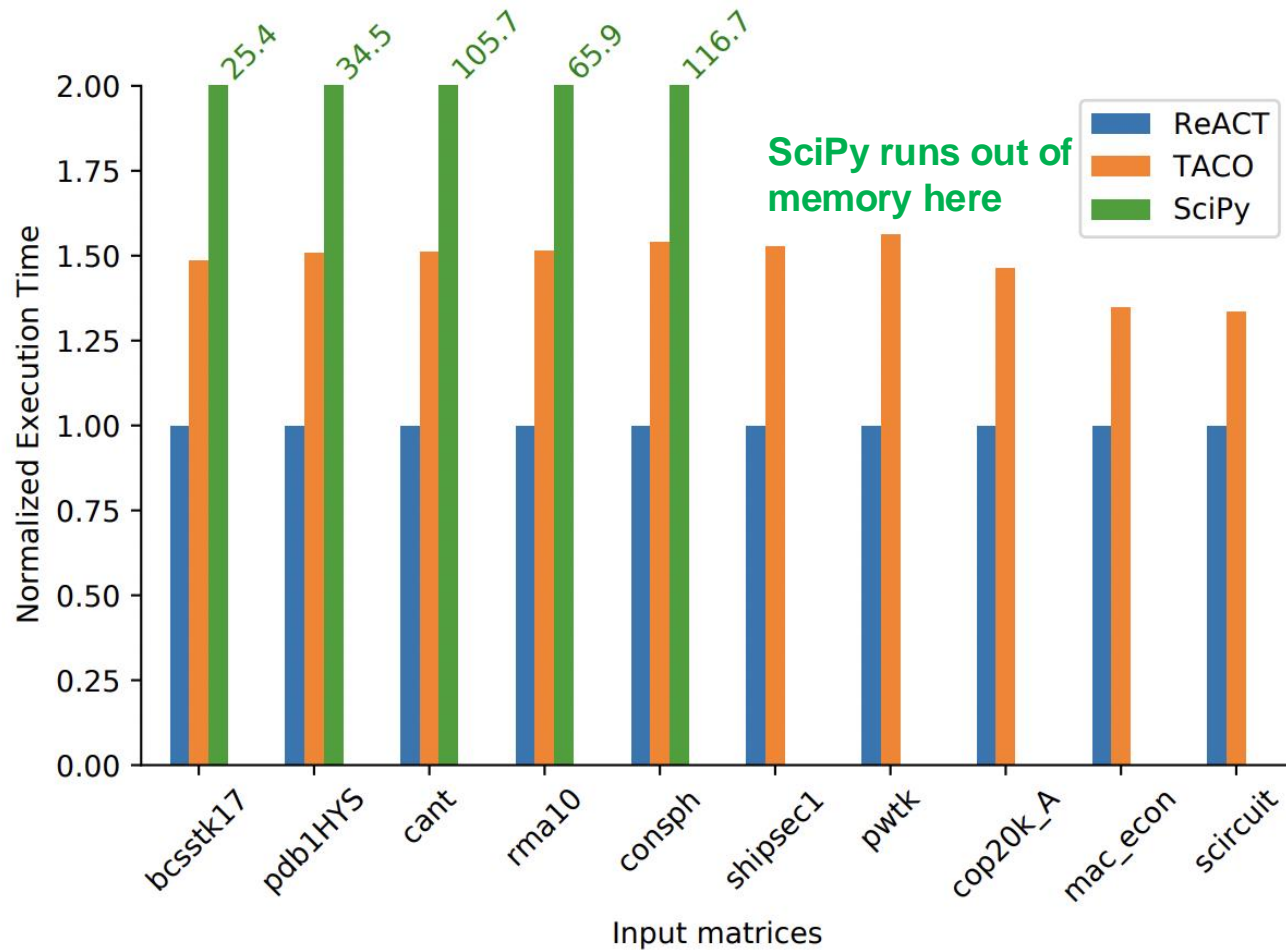Georgia Tech

# SpMM-MM results – 5.7x faster than SciPy



(b) GNN-kernel1 (NH=256, NJ=16)

| Redundancy types present | SciPy | ReACT output |
|---|---|---|
| Type 1 | No | No |
| Type 2 | No | No |
| Type 3 | No | No |
| Type 4 | Yes | No |

**ReACT has better locality + more parallelism**
**Note: SciPy uses only a single thread for its SpMM implementation**
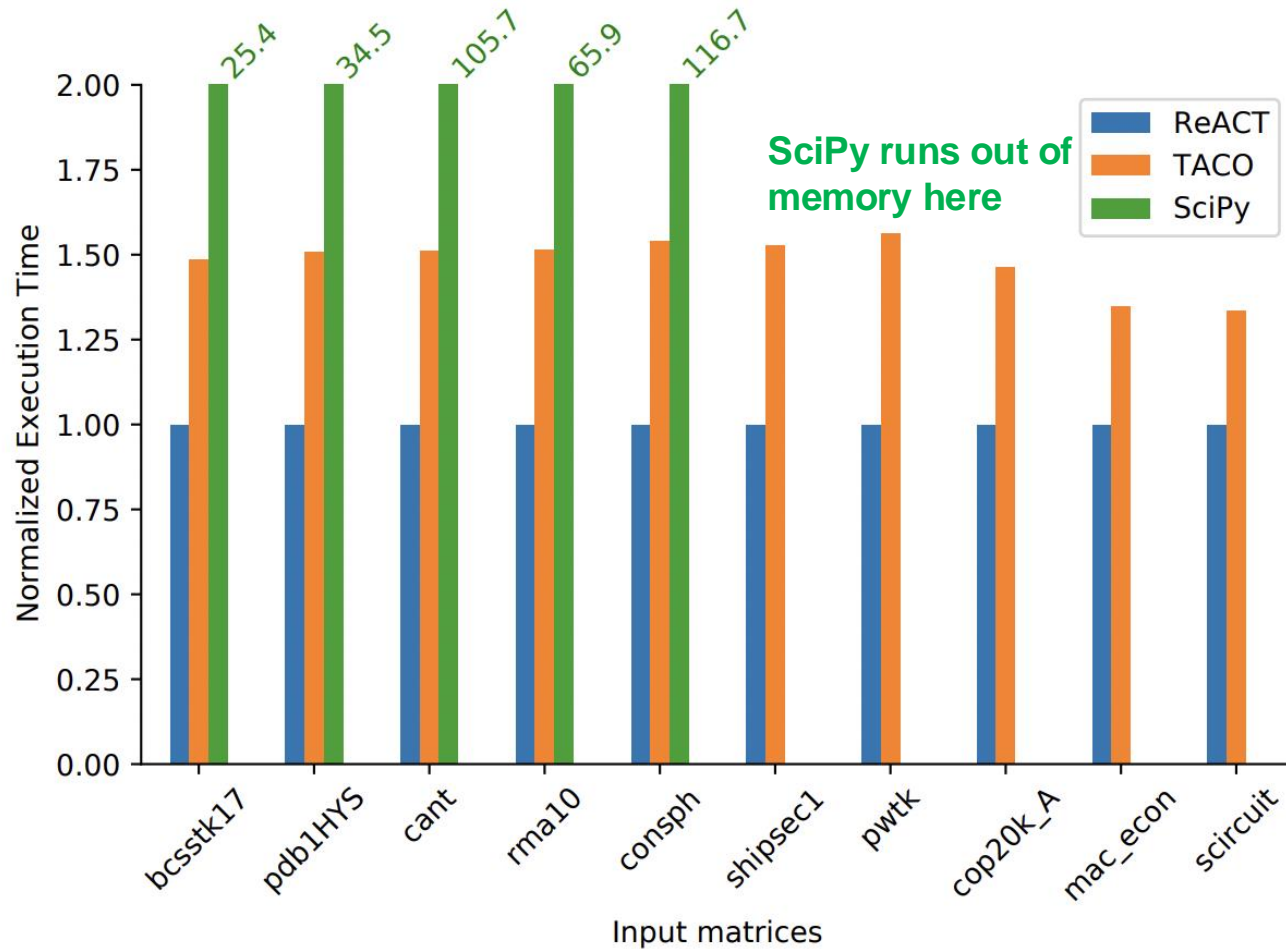
# SDDMM results – 1.5x faster than TACO



SciPy runs out of memory here

Legend: ReACT, TACO, SciPy

(a) SDDMM (NK=64)

| Redundancy types present | TACO output | ReACT output |
|---|---|---|
| Type 1 | Yes | No |
| Type 2 | No | No |
| Type 3 | No | No |
| Type 4 | No | No |

**Both the amount of memory accesses and computations are reduced by eliminating type 1 redundancy.**

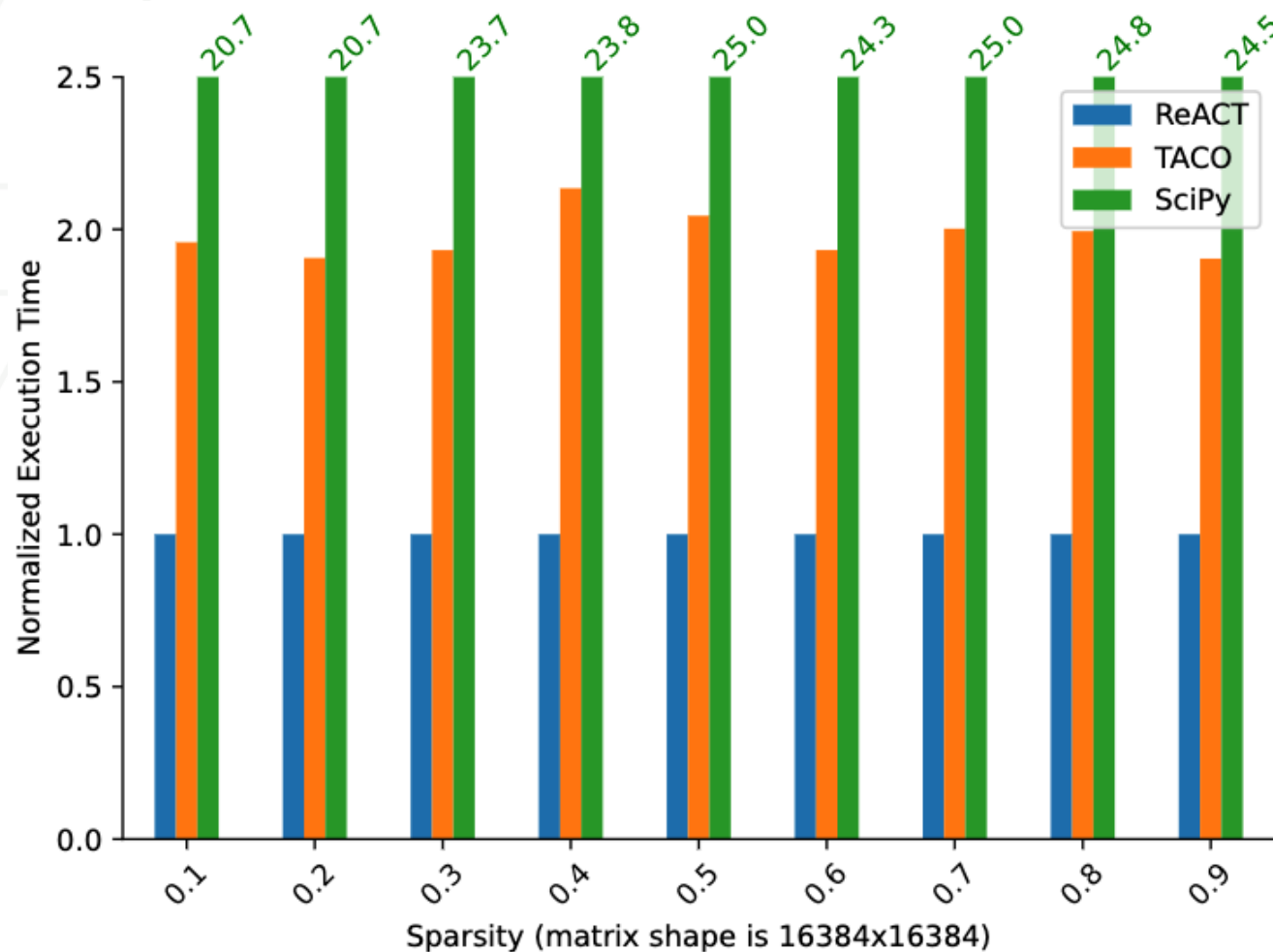Georgia Tech

# SDDMM results – 57.3x faster than SciPy



(a) SDDMM (NK=64)

| Redundancy types present | SciPy | ReACT output |
|---|---|---|
| Type 1 | No | No |
| Type 2 | No | No |
| Type 3 | Yes | No |
| Type 4 | Yes | No |

**Many redundant computations are saved by eliminating type 4 (dead value) redundancies**

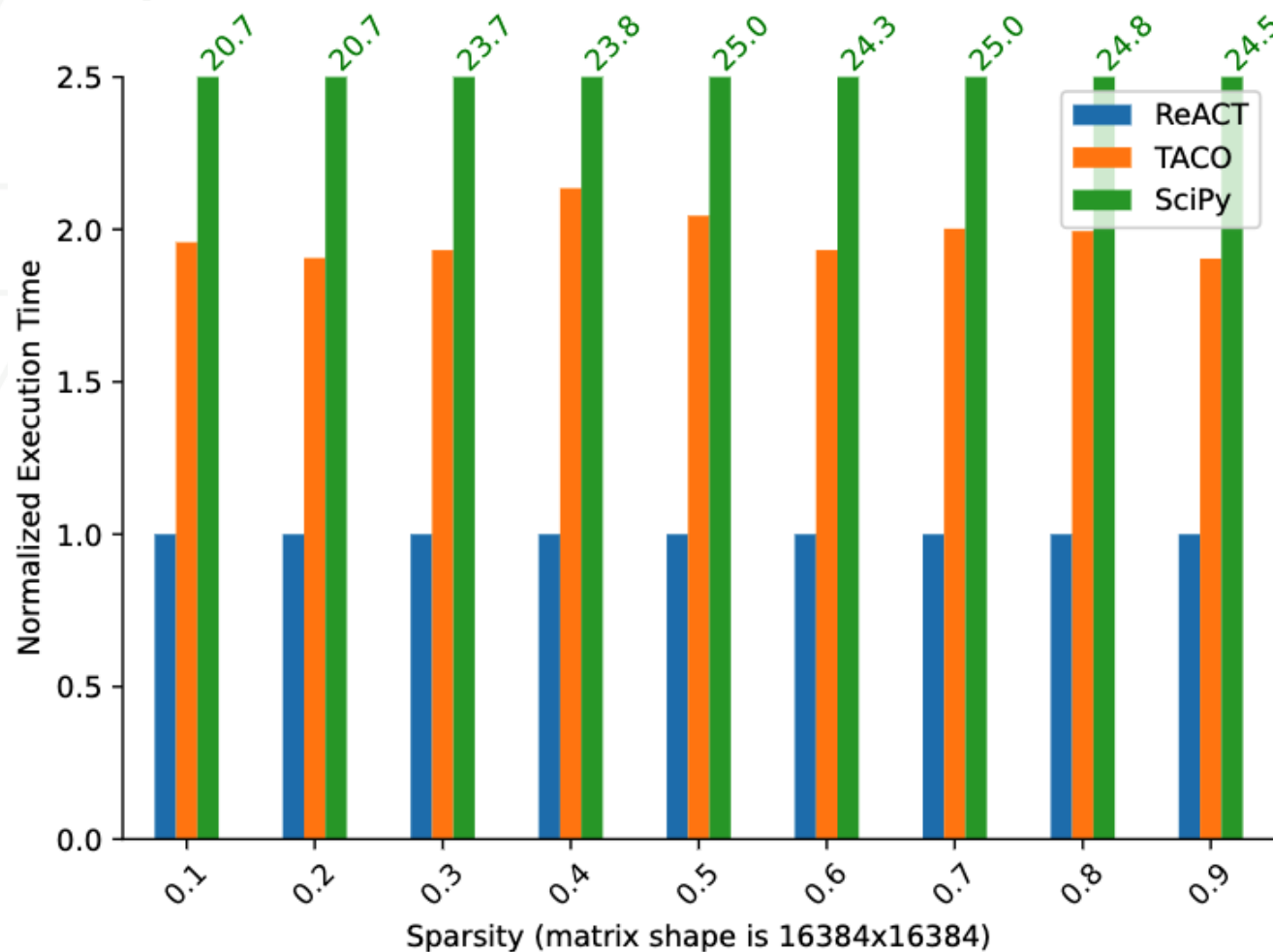Georgia Tech

# Sparse-softmax results – 2.0x faster than TACO



Normalized Execution Time vs Sparsity (matrix shape is 16384x16384). Bars for ReACT, TACO, SciPy. SciPy values shown above bars: 20.7, 20.7, 23.7, 23.8, 25.0, 24.3, 25.0, 24.8, 24.5

**Note: an AMD Ryzen 9 3900X was used for this experiment**

| Redundancy types present | TACO output | ReACT output |
|---|---|---|
| Type 1 | No | No |
| Type 2 | No | No |
| Type 3 | Yes | No |
| Type 4 | No | No |

**TACO cannot fuse it into one single kernel while ReACT does, so ReACT has better locality**

Georgia Tech

# Sparse-softmax results – 23.5x faster than SciPy



| Redundancy types present | TACO output | ReACT output |
|---|---|---|
| Type 1 | No | No |
| Type 2 | No | No |
| Type 3 | Yes | No |
| Type 4 | No | No |

**SciPy's sparse kernels are not parallelized
The operations are also not fused**

**Note: an AMD Ryzen 9 3900X was used for this experiment**

Georgia Tech

# Example: SpMM-MM

- Sparse-dense matmul followed by dense-dense matmul
  - Commonly used in graph neural networks

- Original input expression (sparse matrices are in <span style="color:red">red</span>, assuming CSR format)
  - Python: $A = \textcolor{red}{B} \ @ \ C \ @ \ D$

- Transformations
  - Step 1: convert into *index notation* statements (each statement contains one operator)
    - $S_0$: $T_{ih} = \textcolor{red}{B_{ik}} \ @ \ C_{kh}$ (sparse-dense MM)
    - $S_1$: $A_{ij} = T_{ih} \ @ \ D_{hj}$ (dense-dense MM)
    - $T_{ih}$ is compiler-generated temporary variable
  - Step 2: create an *index tree* from the index notation statements
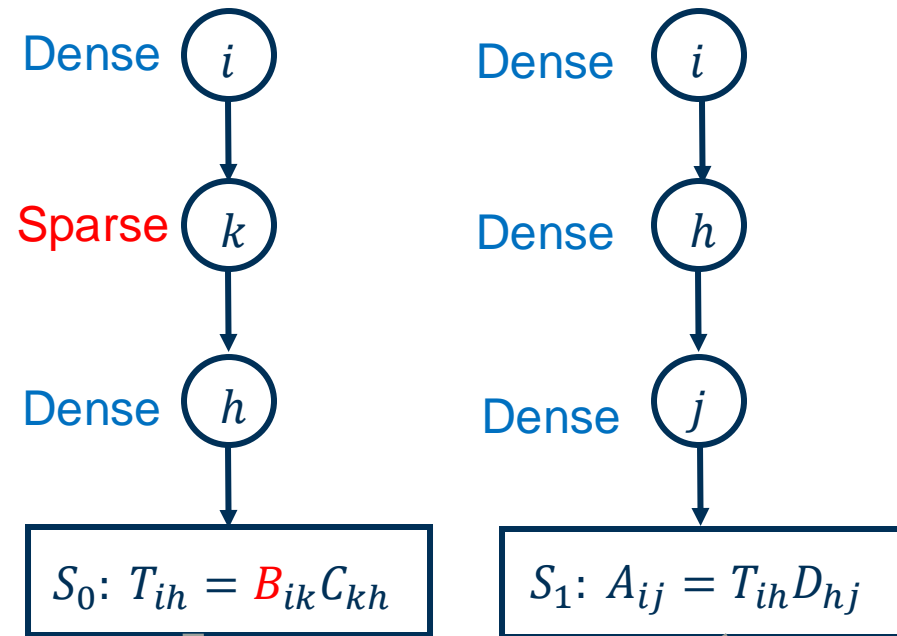    - Next slide

# Index tree of SpMM-MM

- Two operations => create two subtrees
  - $S_0$: $T_{ih} = B_{ik} @ C_{kh}$
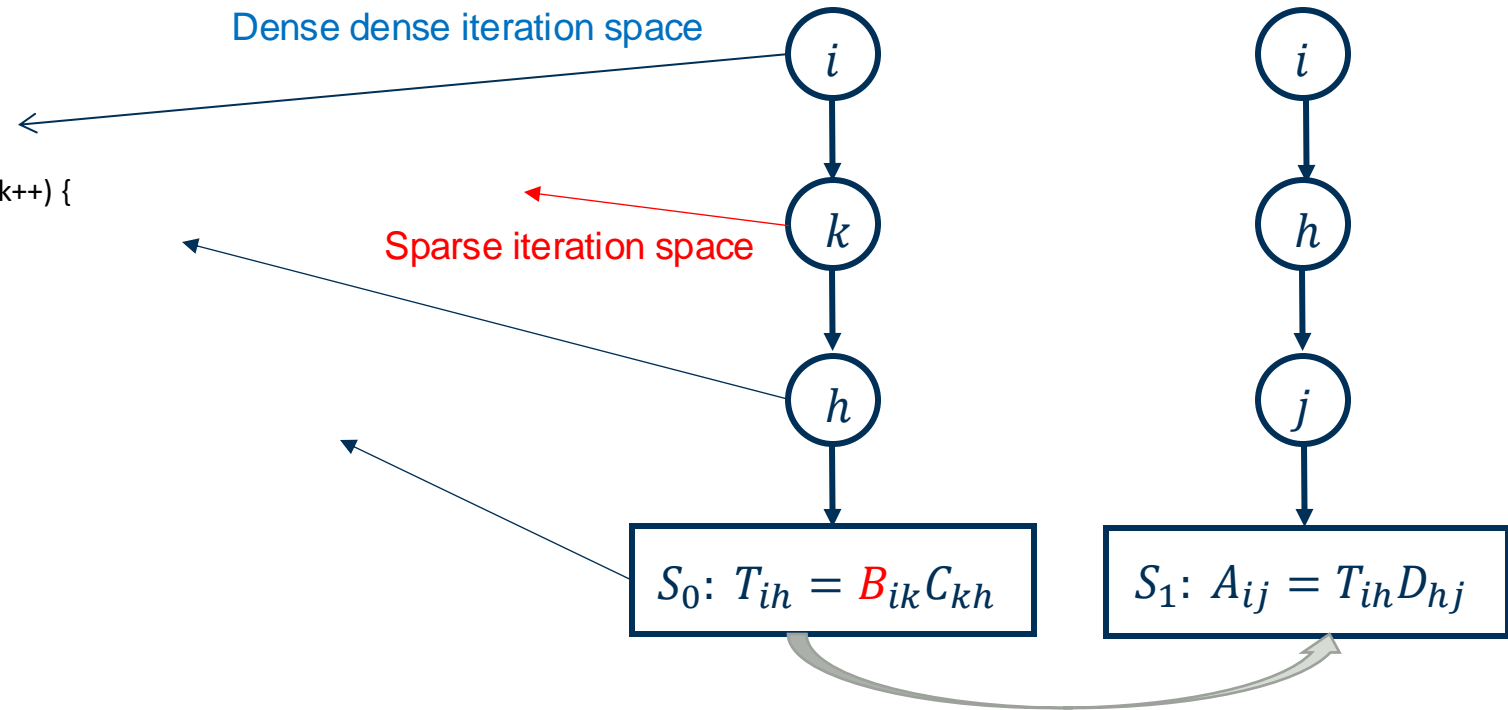  - $S_1$: $A_{ij} = T_{ih} @ D_{hj}$



*Index node*

*Compute node*     $S_0$: $T_{ih} = B_{ik}C_{kh}$     $S_1$: $A_{ij} = T_{ih}D_{hj}$

*Dependence edge*

Georgia Tech.

# SpMM-MM index trees

- Annotate each index node as "Dense" or "Sparse"



Dense $i$

Sparse $k$

Dense $h$

$S_0: T_{ih} = B_{ik}C_{kh}$

Dense $i$

Dense $h$

Dense $j$

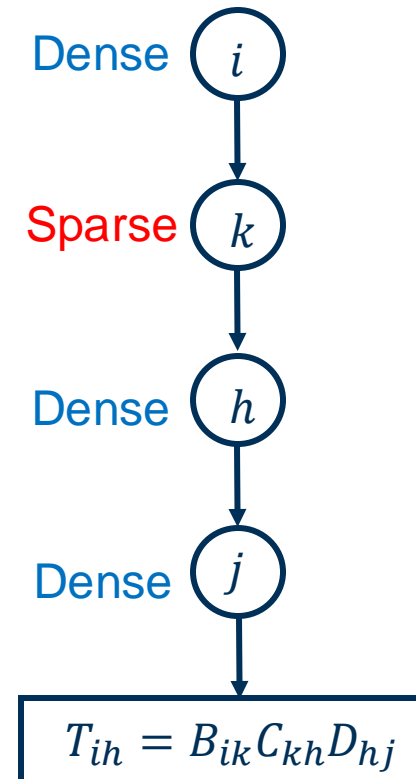$S_1: A_{ij} = T_{ih}D_{hj}$

Georgia Tech.

# Index tree corresponding loop structure



```
1.  for (int i = 0; i < NI; i++) {
2.      for (int k = B.rowptrs[i]; k < B.rowptrs[i+1]; k++) {
3.          for (int h = 0; h < NH; h++) {
4.              ...
5.              // T[i, h] += B[i, k] * C[k, h]
6.              T[i, h] += B.vals[k] * C[B.cols[k], h];
7.              ...
8.          }
9.      }
10. }
```

Dense dense iteration space

Sparse iteration space

$i$

$k$

$h$

$S_0: T_{ih} = B_{ik}C_{kh}$

$i$

$h$

$j$

$S_1: A_{ij} = T_{ih}D_{hj}$

# SpMM-MM index trees: TACO (maximal fusion )

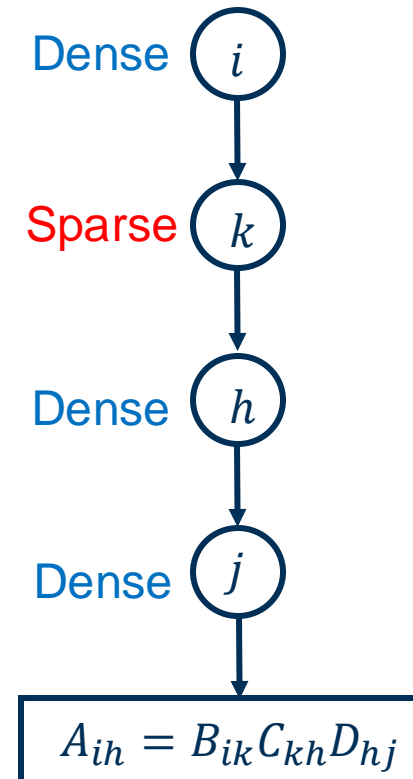- Time: <span style="color:red">Bad</span>, $O(NNZ_B * NH * NJ)$
  - Due to type 1 and 2 redundancies
- Intermediate space: <span style="color:green">Great</span>, $O(1)$
- Locality: <span style="color:green">Great</span>

Dense $\rightarrow i$

Sparse $\rightarrow k$

Dense $\rightarrow h$

Dense $\rightarrow j$

$$T_{ih} = B_{ik}C_{kh}D_{hj}$$

Georgia Tech.

# SpMM-MM index trees: TACO (maximal fusion )

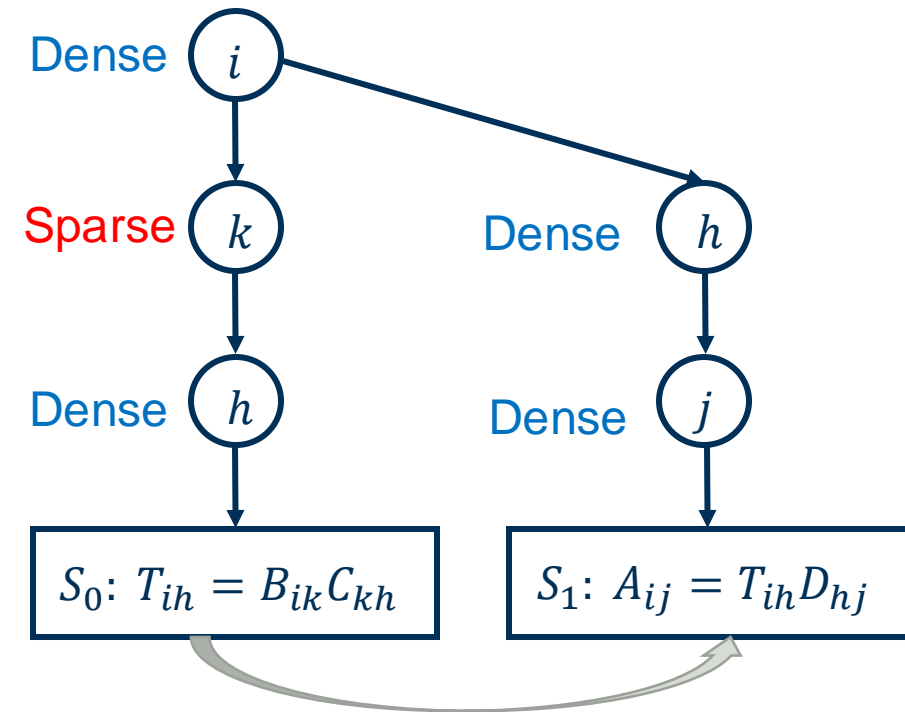## Generated code

```
1.    for (int i = 0; i < NI; i++) {
2.        for (int k = B.rowptrs[i]; k < B.rowptrs[i+1]; k++) {
3.            for (int h = 0; h < NH; h++) {
4.                for (int j = 0; j < NJ; j++) {
5.                    ...
6.                    // A[i, h] += B[i, k] * C[k, h] * D[h, j]
7.                    A[i, h] += B.vals[k] * C[B.cols[k], h] * D[h, j];
8.                    ...
9.                }
10.            }
11.        }
12.    }
```

Dense $i$

Sparse $k$

Dense $h$

Dense $j$

$$A_{ih} = B_{ik}C_{kh}D_{hj}$$

Georgia Tech.

# SpMM-MM index trees: ReACT (partial fusion)
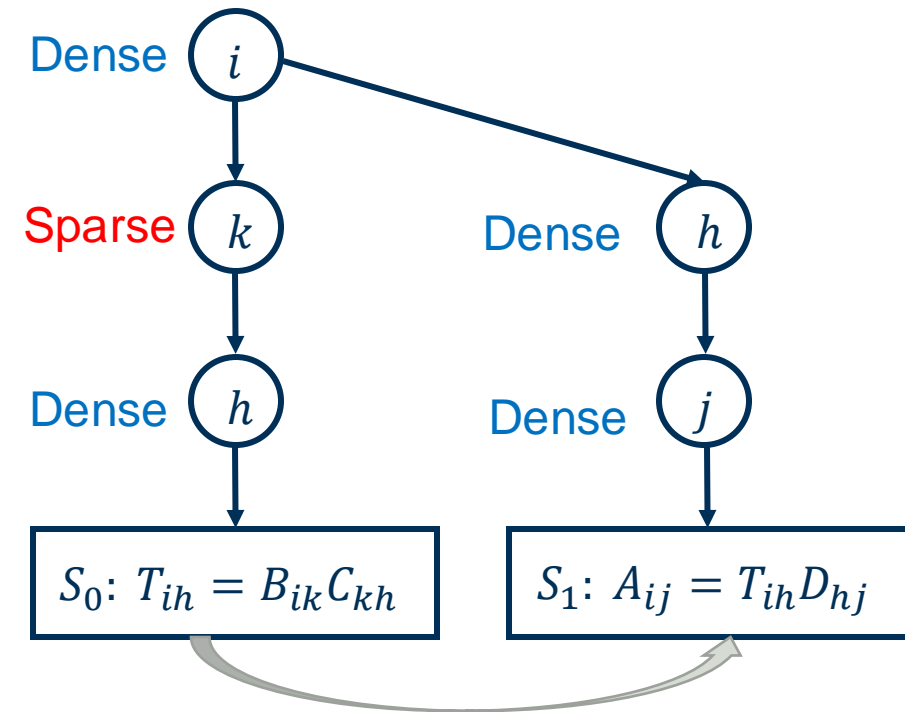
- Time: Good, $O(NNZ_B * NH + NI * NH * NJ)$
  - Typically much smaller than $O(NNZ_B * NH * NJ)$

- Intermediate space: Good, $O(NH)$
  - After memory optimization

- Locality: Good

Dense $i$

Sparse $k$    Dense $h$

Dense $h$    Dense $j$

$$S_0: T_{ih} = B_{ik}C_{kh}$$

$$S_1: A_{ij} = T_{ih}D_{hj}$$

# SpMM-MM index trees: ReACT (partial fusion)

## Generated code

```
1.   for (int i = 0; i < NI; i++) {
2.     for (int k = B.rowptrs[i]; k < B.rowptrs[i+1]; k++) {
3.       for (int h = 0; h < NH; h++) {
4.         ...
5.         // T[i, h] += B[i, k] * C[k, h]
6.         T[h] += B.vals[k] * C[B.cols[k], h];
7.         ...
8.       }
9.     }
10.    for (int h = 0; h < NH; h++) {
11.      for (int j = 0; j < NJ; j++) {
12.        ...
13.        // A[i, h] += T[i, h] * D[h, j]
14.        A[i, h] += T[h] * D[h, j];
15.        ...
16.      }
17.      T[h] = 0;
18.    }
19.  }
```



Dense $i$

Sparse $k$    Dense $h$

Dense $h$    Dense $j$

$S_0: T_{ih} = B_{ik}C_{kh}$    $S_1: A_{ij} = T_{ih}D_{hj}$

# ReACT summary

- We identify four common types of redundancies that can occur when generating code for a sequence of dense/sparse tensor operations

- We introduce ReACT, which consists of a set of redundancy-aware code generation techniques and can generate code with reduced redundancies

- Empirical evaluation on real-world applications such as SDDMM, GNN, Sparse-Softmax, and MTTKRP showed that our generated code with redundancy elimination resulted in 1.1× to orders-of-magnitude performance improvements relative to a state-of-the-art tensor algebra compiler (TACO) and library approaches such as scipy.sparse

Georgia Tech.