# Compilation of Shape Operators on Sparse Arrays

- 3 ALEXANDER J. ROOT, Stanford University, USA
- 4 BOBBY YAN, Stanford University, USA
- <sup>5</sup> PEIMING LIU, Google, USA

1 2

20

21 22

23

24

25

26

27

28

34

35

36

37

38

39

40

41

42

43

44

- <sup>6</sup> AART J.C. BIK, Google, USA
- FREDRIK KJOLSTAD, Stanford University, USA

We show how to build a compiler for a sparse array language that supports shape operators (e.g. reshaping or 9 concatenating arrays) in addition to compute operators. Existing sparse array programming systems implement 10 generic shape operators for only some sparse data structures, reduce shape operators on other data structures to 11 those, and do not support fusion. Our system compiles sparse array expressions to code that efficiently iterates 12 over reshaped views of irregular sparse data structures, without needing to materialize temporary storage for 13 intermediates. Our evaluation shows that our approach generates sparse array code competitive with popular 14 sparse array libraries: shape operators perform on average  $1.49 \times (0.408 \times -4.31 \times)$  better than hand-written 15 kernels in scipy. sparse and 28× (0.476×-274×) better than generic implementations in pydata/sparse. For 16 operators that require data structure conversions in these libraries, our generated code performs on average 17  $4.19 \times (1.99 \times -8.36 \times)$  better than scipy.sparse and  $54.6 \times (10.8 \times -259 \times)$  better than pydata/sparse. Finally, our evaluation also shows that fusing shape and compute operators improves the performance of several 18 expressions  $(1.16 \times -2.79 \times)$ . 19

# 1 INTRODUCTION

Shape operators are the type casts of array programming. The shape of an array is a key component of its type, and determines which operations on it are valid: vectors of different lengths cannot be added and matrices of incompatible dimensions cannot be multiplied. However, programmers often need to manipulate the shape of arrays, and do so via shape operators. Many important computations require explicitly manipulating array shapes, such as stacking constraint matrices in physical simulation [Solomon 2015], reshaping tensors in neural networks [Vaswani et al. 2017], and slicing images in biomedical computing [Boyse and Seidl 1996]:

$$y = \begin{bmatrix} A \\ B \\ C \end{bmatrix} x \qquad \qquad y = A \begin{bmatrix} X_{0,0} \\ \vdots \\ X_{n-1,m-1} \end{bmatrix} \qquad \qquad Y = A_{:,i}$$

Stacked matrix-vector multiply

Linear layer on a flattened matrix

2D slice of 3D MRI data

As a result, dense array programming systems [Harris et al. 2020; Iverson 1962; Paszke et al. 2019] have ample support for such operators, and they often reduce to constant-time metadata edits.

Sparse array programming systems lag behind their dense counterparts, with incomplete support for shape operators. While sparse array libraries [Abbasi 2018; Virtanen et al. 2020] support some shape operators, the implementations reduce to a small set of hand-written kernels. The conversions to intermediate data structures required by these reductions incur a significant performance cost. Sparse compilers [Ahrens et al. 2023; Kjolstad et al. 2017; Ye et al. 2023; Zhao et al. 2022] lack a complete set of shape operators.

Compiler support for shape operators is necessary for performant sparse array systems. It is not enough to simply hand-engineer a library of sparse shape operators, nor is it feasible to

Authors' addresses: Alexander J. Root, Computer Science, Stanford University, USA, ajroot@stanford.edu; Bobby Yan,
 Computer Science, Stanford University, USA, bobbyy@stanford.edu; Peiming Liu, Google, USA, peiming@google.com;
 Aart J.C. Bik, Google, USA, ajcbik@google.com; Fredrik Kjolstad, Computer Science, Stanford University, USA, kjolstad@
 stanford.edu.

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92 93

94

95 96

50 Table 1. Comparison of sparse shape operation support across several programming systems. Yellow circles 51 indicate partial support: scipy.sparse and pydata/sparse have limited sparse data format support. TACO 52 and Looplets support fusing compute operators, but not shape operators.

Programming Model	Compilation	Data Representation			Shape Operators			
		Dense	Sparse	Any dims	Slicing	Concatenation	Reshape	Fusion
numpy	×	<ul> <li>✓</li> </ul>	×	<ul> <li>✓</li> </ul>	<ul> <li></li> </ul>	<ul> <li></li> </ul>	<ul> <li></li> </ul>	×
scipy.sparse	×	<ul> <li>V</li> </ul>	•	×	<ul> <li>V</li> </ul>	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>	×
pydata/sparse	×	<ul> <li>V</li> </ul>	•	<ul> <li>✓</li> </ul>	<ul> <li>Image: A second s</li></ul>	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>	×
TACO	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>	<ul> <li>Image: A second s</li></ul>	<ul> <li>✓</li> </ul>	<ul> <li>V</li> </ul>	×	×	•
Looplets	×	× .	<ul> <li></li> </ul>	~	×	×	×	•
BURRITO (This Work)	v -	× .	<ul> <li>Image: A second s</li></ul>	~	<ul> <li>Image: A second s</li></ul>	V	~	×

implement the full Cartesian product of operations and sparse array data structures. Prior work on sparse compute operators (e.g. matrix multiplication) shows that compiler support is necessary for performance and programmer productivity, due to the scaling issues introduced with the large number of sparse data structures [Ahrens et al. 2023; Chou et al. 2018]. Compiler support for sparse shape operators enables generating fused code that is specialized to the data structures of the operands and the result. The key challenge in compiling sparse shape operators lies in generating code that iterates over, and computes on, combinations of reshaped, concatenated, and sliced views of irregular data structures representing sparse arrays.

While recent work on compilers for sparse tensor algebra [Ahrens et al. 2023; Bik et al. 2022; Kjolstad et al. 2017] addresses the sparse compute problem, they offer limited support for shape operators: TACO supports slicing unfused array operands [Henry et al. 2021] and Looplets [Ahrens et al. 2023] supports concatenation without fusion. These systems offer no unified representation of shape operators, and *no* compiler that we are aware of supports sparse reshape operations. Libraries such as scipy.sparse [Virtanen et al. 2020] and pydata/sparse [Abbasi 2018] are not feature complete, and are limited by the time required for the maintainers to implement different combinations of shape operators and sparse array data structures. Libraries also naturally do not support operator fusion.

We introduce a compiler-based approach for representing shape operators, building on ideas from TACO to support iterating over reshaped arrays. We implement these ideas in a prototype compiler named BURRITO<sup>1</sup>. We extend an existing array programming language to support shape operators (Section 3), describe compilation to a new intermediate representation (IR) for iterating over sparse data structures (Section 4), and outline the compilation process from our IR to efficient CPU code (Section 5 and Section 6). Our abstractions enable generation of specialized code for sparse array kernels containing shape operators, can generate code for a variety of data structures, and enable fusion across shape and compute operators. Table 1 summarizes how BURRITO compares to other array programming systems. Our technical contributions are:

- (1) A lowering approach from an array language with shape operators (Section 3) to a high-level loop language that expresses fusion across operators (Section 4).
- (2) A state-driven algorithm for generating optimized while loops that coiterate through reshaped iteration spaces (Section 5).
- (3) A compilation approach for generating iteration code that abstracts both data structures and reshaped iteration spaces, based on a simple iterator model (Section 6).

BURRITO generates code competitive with hand-written libraries on unfused shape operators, and generally out-performs these libraries when given sparse kernels that offer fusion opportunities.

<sup>&</sup>lt;sup>1</sup>A BURRITO is what you get when you *reshape* a TACO.



Fig. 2. The BURRITO compilation approach compiles an array language to C code in three steps through two loop-based intermediate representations over logical arrays. Section references mark our contributions.

Across five shape operators hand-written in scipy.sparse and pydata/sparse, BURRITO performs  $1.11 \times -4.31 \times$  and  $10.4 \times -274 \times$  (respectively) faster. A sixth shape operator performs  $0.408 \times$ and  $0.476 \times$  slower respectively due to vectorization, which BURRITO does not currently support. For shape operators whose library implementations convert between data structures, BURRITOgenerated code out-performs scipy.sparse by  $1.9 \times -8.36 \times$  and pydata/sparse by  $6.29 \times -287 \times$ . To evaluate the benefits of fusion, we perform a self-comparison and show that code that fuses shape and compute operators out-performs unfused code by  $1.16 \times -2.79 \times$  on some benchmarks.

# 2 OVERVIEW

107

108

117

The compilation approach we describe in this paper compiles an array language to C code. The array language is tensor index notation extended with shape operators, and operates on *logical* arrays that abstract away the details of the *physical* data structures. BURRITO separately accepts a format language that describes the sparse or dense physical data structures that implement the logical arrays, following prior work on separating algorithm from data structure [Chou et al. 2018; Kjolstad et al. 2017]. Finally, a separate scheduling language [Kjolstad et al. 2019] lets users or automated systems describe loop reordering and introduce temporaries.

125 Figure 2 shows an overview of our compilation approach, which we implement in the BURRITO 126 compiler. Compilation consists of three lowering steps through two intermediate representations 127 (IRs). The first step lowers the array language to nested forall loops over logical arrays. Scheduling 128 directives (e.g. loop reordering) apply to this IR. Each loop iterates over an expression that describes 129 the loop's iteration domain as an ordered set expression (sequence expression) over the coordinates 130 of array dimensions. Prior work expresses compute operators as intersections and unions [Kjolstad 131 et al. 2017] and further work adds complements and slicing [Henry et al. 2021]. In order to support 132 the complete set of shape operators, our work adds three new set expression operators: set product, 133 set projection, and disjoint union. 134

```
(intersection + union) + (complement + slicing) + (product + projection + disjoint union)
[Kjolstad et al. 2017] [Henry et al. 2021] [this work]
```

The forall loops are then lowered to an IR that expresses more complex control-flow. Each forall is turned into a sequence of while loops over sequence expressions that iterate over progressively simplified loop bodies as sparse arrays run out of non-zero coordinates. This lowering step requires knowledge of format properties to determine which arrays to iterate over and which arrays to randomly access into. The final step lowers the while loops to C code where logical arrays are replaced by the physical data structures described via the format language.

# **3 SHAPE OPERATORS**

BURRITO compiles a language that extends tensor index notation for tensor algebra to support shape operators. Tensor index notation is a declarative language for describing tensor algebra.

146 147

143

144

145

135

148	kernel	::=	array(I) = expr
149	expr	::=	<pre>array(I)   expr + expr   expr * expr   sum(idx, expr)   broadcast(idx, expr)  </pre>
150			$\label{eq:collapse} \mbox{collapse}(\mbox{expr, (idx, idx)} \rightarrow \mbox{idx}) \ \Big  \ \mbox{concat}((\mbox{expr, expr), (idx, idx)} \rightarrow \mbox{idx}) \ \Big $
151			<b>split</b> (expr, idx $\rightarrow$ (idx, idx)) slice(expr, idx $\rightarrow$ idx, (int, int, int))

Fig. 3. Our front-end language for expressing array algebra, which supports tensor index notation (element wise addition and multiplication, summation along an axis, and broadcasting), as well as our additional shape
 operators. Indices (idx) label array dimensions.

An example is matrix multiplication,  $A_{ij} = \sum_k B_{ik}C_{kj}$ , where each component of  $A_{ij}$  is the inner product of a row of B and a column of C. Shape operators are common in libraries like numpy and PyTorch and let users change the shape of a tensor in various ways. For example, a user can reshape (flatten) a matrix into a vector with the same number of elements or concatenate two matrices with the same number of rows into a matrix that contains columns from each of the two operands.

Figure 3 provides the syntax of the algorithm language supported by BURRITO. We first provide a high-level semantics for the shape operators by defining how each operator maps a coordinate in the logical input array(s) to a coordinate in the constructed logical array, and then provide a shape inference algorithm for detecting the validity of an array program via typing rules.

# 166 3.1 Operator Definitions

The shape of a logical *n*-dimensional logical array represents an *n*-dimensional hyper-rectangle,
where each dimension has a size. Tensor index notation labels each dimension via *index variables*.
Shape operators are used to explicitly manipulate the shapes of arrays by combining and constructing index variables, and thereby the corresponding dimensions. We define each shape operator
with respect to the coordinate mapping it induces from one or more input arrays to the constructed
array. Note that each of our operators are given as binary operators (i.e. binary concatenation), as
BURRITO supports *fusion by default*.

174 *Collapse*. The collapse operator flattens two dimensions into one, constructing a new index with a 175 size equal to the product of the size of the flattened dimensions. The shapes of the input and the 176 output contain the same number of discrete points, but the output shape has one fewer dimensions 177 than the input shape. The ordering of the flattened index variables (the second argument to collapse) 178 controls the coordinate space mapping: the ordering of coordinates in the constructed dimension is 179 equivalent to the the lexicographic ordering of the flattened dimensions. For example, a matrix B 180 with rows labeled i and columns labeled j can be row-major collapsed via a(k) = collapse(B(i, j), k)181  $(i, j) \rightarrow k$ , or column-major collapsed via  $a(k) = collapse(B(i, j), (j, i) \rightarrow k)$ . The coordinate 182 remapping follows the standard strided offset formula. 183

Concatenate. The concat operator is used to stack two arrays along a dimension. The shape of the out-184 put has the same number of discrete points as the sum of the points of the shapes of the inputs, and 185 the same number of dimensions. The ordering of arrays control the coordinate space mapping: each 186 coordinate in the first input array exists in the same location in the output array, and each coordinate 187 in the second is offset in the dimension being concatenated. For example, a program that horizon-188 tally concatenates two arrays can be expressed as:  $A(i, 1) = concat((B(i, j), C(i, k)), (j, k) \rightarrow 1)$ . 189 A contains every coordinate B contains, as well as every coordinate that C contains, with the second 190 coordinate of each element offset by the size of the dimension labeled by j. 191

Split. The split operator is the inverse of collapse, as it divides one source dimension into two constructed dimensions<sup>2</sup>. The source dimension has a size equal to the product of the sizes of

195 196 1:4

155

<sup>&</sup>lt;sup>2</sup>**split** and **collapse** are building blocks that can be used to support the *reshape* operation.

the constructed dimensions. The input and output shapes contain the same number of discrete 197 points, but the output shape has one more dimension than the input shape. The ordering of the 198 pair of constructed indices controls the coordinate space mapping: the lexicographic ordering of 199 the constructed dimensions is equivalent to the ordering of the source dimension. For example, 200 consider splitting a 1D array into 2D:  $A(i, j) = split(b(k), k \rightarrow (i, j))$ . For this row-wise splitting, 201 the coordinate 1 in b corresponds to (0, 1) in A. For a column-wise splitting: A(i, j) = split(b(k), k 202  $\rightarrow$  (j, i)), the coordinate 1 in b corresponds to (1,0) in A. This coordinate remapping is defined by 203 204 the inverse of the strided offset formula for collapse.

205 *Slice*. The slice operator extracts a uniform (possibly strided) subsequence from a dimension via 206 values for the start, end, and stride. The shape of the input contains more discrete points<sup>3</sup> but 207 has the same number of dimensions as the output shape. A coordinate in the sliced dimension is 208 discarded if: 1) it is less than the start value; 2) greater than or equal to the end value; or 3) not a 209 whole-number multiple of the stride value from the start index. If a coordinate passes these filters, 210 i.e. a coordinate c in the input dimension is sliced with a start index of s and a stride of r meets 211 the constraint:  $c = s + x \cdot r$  for  $x \in \mathbb{N}$  and c < e, then it maps to the coordinate x in the output 212 dimension. For example, if a user wished to select the even-indexed values of a vector, they could 213 write the kernel  $a(i) = \text{slice}(b(j), j \rightarrow i, (0, J, 2))$ , or if they wanted the first half of an array, 214  $a(i) = slice(b(j), j \rightarrow i, (0, J/2, 1))$ . Note that Henry et al. [2021] introduced slicing on sparse 215 array operands only, while BURRITO supports slicing intermediate computations, i.e. slicing the 216 result of a multiplication, or a reshaped array. 217

## 3.2 Shape Inference

The shape of an expression in array algebra is an essential component of the  $type^4$ , and an array compiler requires a shape inference algorithm for checking the validity of programs. For tensor index notation, shapes are easy to infer; element-wise operations produce an output with the same shape as the inputs and reductions remove a dimension from the shape of the operand. Broadcasting, though implicit in tensor index notation, is technically a shape operator, as it inserts a new dimension into the logical shape of an array. For the new shape operators defined in the prior section, new dimensions are constructed that a compiler must reason about in order to check program correctness.

As described previously, the shape of an array describes the number of dimensions and the size of each dimension. We denote shape as *S*, an unordered set of indices, each of which labels a dimension with a name and size. The set cardinality operator |i| is used to represent the size of the dimension indexed by *i*.

The full inference algorithm is provided in Figure 4, but as an illustrative example, the inference rule for **collapse** is:

a: 
$$S \quad i \in S \quad j \in S \quad k \notin S \quad |i| \cdot |j| = |k|$$
  
collapse(a, (i, j)  $\rightarrow$  k):  $(S - \{i, j\}) \cup \{k\}$ 

The first two constraints,  $i \in S$  and  $j \in S$ , require that the source indices that are being collapsed, i and j, are contained within the shape, *S*, of the array operand, a. Next,  $k \notin S$  requires that the constructed index, k, is not already in *S*, which is necessary for the uniqueness of indices in the output shape. The last constraint,  $|i| \cdot |j| = |k|$ , requires that the size of the dimension of the constructed index, |k|, is the product of the size of dimensions represented by the source indices, |i|and |j|. This means that the output array of the collapse has the same number of discrete elements

245

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

 $<sup>\</sup>overline{^{3}\text{Assuming the slice is non-trivial: a trivial slice removes no elements, i.e. slicing$ *i*with a start of 0, an end of the size of*i*, and a stride of 1.

<sup>&</sup>lt;sup>4</sup>The type of an array consists of the shape and the type of the scalar elements.

Alexander J. Root, Bobby Yan, Peiming Liu, Aart J.C. Bik, and Fredrik Kjolstad

$$a: S \quad i \in S \quad j \in S \quad k \notin S \quad |i| \cdot |j| = |k|$$

$$a(I): I \qquad \qquad \text{collapse(a, (i, j) \rightarrow k): (S - \{i, j\}) \cup \{k\}}$$

$$\begin{array}{c} \begin{array}{c} \begin{array}{c} 248\\ \\ 249\\ \\ 250\\ \end{array} \end{array} \qquad \begin{array}{c} a:S \quad i \notin S\\ \hline broadcast(i, a):S \cup \{i\} \end{array} \qquad \begin{array}{c} a:S_a \quad b:S_b \quad S_a - \{i\} = S_b - \{j\} \quad k \notin S_a \quad |i| + |j| = |k| \end{array} \\ \begin{array}{c} \begin{array}{c} a:S_a \quad b:S_b \quad S_a = S_b\\ \hline split(a, b), (i, j) \rightarrow k): (S_a - \{i\}) \cup \{k\} \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad k \notin S \quad |i| = |j| \cdot |k| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad \frac{e-s+(r-1)}{r} = |j| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad \frac{e-s+(r-1)}{r} = |j| \end{array} \\ \begin{array}{c} a:S \quad i \in S \quad j \notin S \quad \frac{e-s+(r-1)}{r} = |j| \end{array} \end{array} \end{array}$$

Fig. 4. Shape inference rules and validity constraints for shape and compute operators. An array's shape 257 is the set of indices that index it (note that the indices used to read an array are ordered, but the shape is 258 unordered). Broadcasting inserts a new index into the shape. Element-wise operations produce an expression 259 with the same shape as the inputs. Summation removes the reduced index from the shape. Collapsing 260 flattens two indices into a single constructed index, while splitting divides a single source index into two 261 constructed indices. Concatenation accepts a pair of expressions and pair of indices, concatenating the 262 two expressions along the two indices to produce a single constructed index. Concatenation requires that 263 the non-concatenated indices of each operand must match. Slicing re-labels a source index into a smaller constructed index. 264

in space as the array being collapsed. The output array has a shape where dimensions other than i and j are unchanged, and i and j are replaced by k.

Consider the program  $c(k) = a(k) * collapse(B(i, j), (i, j) \rightarrow k)$ , the element-wise multiplica-268 tion of a vector with a flattened matrix. This program is well-typed: the shape of the argument 269 to collapse is  $S_{B(i, j)} = \{i, j\}$ , and  $i \in S_{B(i, j)}$  and  $j \in S_{B(i, j)}$  are trivially true. Likewise, the 270 constructed index is not in the source shape,  $k \notin S_{B(i, j)}$ . Lastly, the dimensionality constraint is 271 satisfied if  $|\mathbf{i}| \cdot |\mathbf{j}| = |\mathbf{k}|$ . Stepping up a level, both operands to the multiplication have a shape of  $\{\mathbf{k}\}$ , 272 so the multiplication is well-typed. 273

## **SEQUENCE EXPRESSIONS**

The first compilation pass generates forall loops over logical arrays from the operators described in 276 the previous section. The IR for representing forall loop extends TACO's concrete index notation 277 (CIN) [Kjolstad et al. 2019] with a more advanced language for describing the iteration spaces of 278 loops, termed sequence expressions. CIN is a language of loops over logical arrays that supports loop 279 transformations such as loop reordering and the allocation of temporaries. These transformations 280 are essential for avoiding inefficient discordant traversals, which iterate over a data structure in an 281 order for which it was not designed [Bik et al. 1994; Gustavson 1972]. Collapsing a CSR matrix in 282 column-row order or concatenating a CSR matrix with a CSC matrix are both examples of discordant 283 traversals. It is generally more efficient to transpose an operand into a temporary matrix, and iterate 284 over the temporary in concordant order in these cases, as discordant traversals require a dense loop 285 with sparse look-ups. This section focuses on our novel contribution, sequence expressions, and 286 we refer the reader to prior work for further discussion of CIN, loop transformations, and avoiding 287 discordant traversals [Kjolstad et al. 2019]. 288

Sequence expressions are closely related to the coordinate mappings described in Section 3.1, and 289 can express both compute and shape operators. For example, addition corresponds to the union of 290 sequences and collapse corresponds to the Cartesian product of sequences, where the coordinates 291 in the resulted tuples are combined using a strided offset formula. Each shape operator corresponds 292 to an operation on the sequences of non-zeros of its operands, and these operations fundamentally 293

294

1:6

246 247

265

266

267

274



Fig. 5. Various space filling curves.

change the iteration spaces of the loops. We provide semantics for sequence combinators and a construction algorithm for generating sequence expressions from the expression language described in Section 3.1.

## 4.1 Sequence Combinator Semantics

 We provide the full grammar for CIN below, but this section focuses on our novel contribution: the grammar for sequence expressions, seq.

stmt ::= forall idx  $\in$  seq stmt | stmt; stmt | stmt where stmt | a(I) = expr | a(I) += expr

seq := idx a | seq  $\cup$  seq | seq  $\cap$  seq | seq  $\times$  seq | seq  $\sqcup$  seq |  $\pi_k$  (seq, (int, int) ) | seq[int:int:int]

Each sequence operator operates on a sequence of coordinates in a dimension of one or more arrays. Sequence operators corresponding to shape operators perform coordinate space transformations on their operand sequences. The semantics of these operators are important both for understanding how shape operators combine iteration spaces, and for compiling down to loops over irregular data structures. These operators are used to generate a sequence expression that represents the set of non-zero coordinates of the output array of a computation, for each dimension of the output array. Intuitively, each operator corresponds to an operation that changes the space filling curve of a hyperrectangle, and we provide examples in Figure 5. The semantics of these operators are as follows:

**product**:  $A \times B$  is isomorphic to the Cartesian product, but applies a strided offset formula to the produced tuples to generate a whole number coordinate value. Intuitively, product corresponds to: for each element in *A*, step through each element in *B*. Figure 5b depicts the product space-filling curve of a collapsed matrix. Formally, product is expressed as

$$a \in A \text{ and } b \in B \iff (a \cdot |B| + b) \in A \times B$$

**concatenation**:  $A \sqcup B$  is isomorphic to a disjunctive union, but produces a sequence with a range that is the sum of the input sequence ranges. Intuitively, this operator corresponds to first stepping through each element in A and then through each element in B. Figure 5c depicts the space filling curves of horizontally concatenated matrices. Formally

$$a \in A \iff a \in A \sqcup B \quad b \in B \iff (b + |A|) \in A \sqcup B$$

**projection**:  $\pi_k(A, (I, J))$  is isomorphic to the set projection operation, but first applies the inverse of the strided offset formula to elements of A to produce a tuple from a single coordinate, with the shape (I, J). The projection index, k, selects the kth element from the produced tuple. Because BURRITO's front-end allows users to express binary splitting, sequence expressions only need to support two-way projections, which can be chained if a dimension needs to be split more than once. Intuitively, a projection of a sequence produces two sequences that each iterate over sub-spaces of the original sequence. The first sequence  $\pi_0(A, (I, J))$  is a sequence of size I, and the second,  $\pi_1(A, (I, J))$  is a sequence of size J. Projection is the inverse to product, just as collapsing and 

1:7

Alexander J. Root, Bobby Yan, Peiming Liu, Aart J.C. Bik, and Fredrik Kjolstad

splitting are inverses. Figure 5e depicts the space filling curves produced by splitting a vector into two dimensions. Formally

$$a \in A \iff \frac{a}{J} \in \pi_0(A, (I, J)) \text{ and } (a\%J) \in \pi_1(A, (I, J)).$$

**slice**: A[s:e:r] is a simple filtering operation. Intuitively, slicing removes all elements from a sequence that are not a stride of r away from the start s, and any elements greater than e. Figure 5f depicts the space filling curve produced by slicing the 3rd through the 5th elements of a length 6 vector. Formally

 $a \in A$  and  $s \le a \le e$  and  $\exists x \in \mathbb{N}, a = s + r \cdot x \iff x \in A[s:e:r].$ 

#### 4.2 Sequence Combinator Construction

BURRITO's construction of sequence expressions from an array expression is defined by the function  $\Upsilon(e, i)$ .  $\Upsilon$  accepts an expression, e, from the grammar defined in Figure 3, and an index variable, i, and generates a sequence expression that represents the output coordinates in the ith dimension of e. The full definition of  $\Upsilon$  is given in Figure 7, and we provide an example for intuition.

Consider the compilation of the rightmost kernel in Figure 6, which generates a single loop over the output index, k. In order to derive the sequence expression that corresponds to k's iteration space, the recursive construction algorithm applies the following rules from Figure 7:

$$\Upsilon(a + b, i) = \Upsilon(a, i) \cup \Upsilon(b, i) \qquad \Upsilon(collapse(a, (j, k) \to 1), i) = \begin{cases} \Upsilon(a, j) \times \Upsilon(a, k) & i = k \\ \Upsilon(a, i) & i \neq k \end{cases}$$

These rules construct the sequence  $(i_A \times j_A) \cup k_b$ , the loop bounds in the rightmost loop of Figure 6. This sequence expression means that the output sequence of non-zeros is equivalent to the flattened sequence of a's non-zeros unioned with the sequence of b's non-zeros. This exactly represents the set of non-zeros in the output array. 

$$\Upsilon(\mathsf{t}(\mathsf{I}), \mathsf{i}) = \begin{cases} i_t & i \in I \\ \emptyset & i \notin I \end{cases} \qquad \qquad \Upsilon(\mathsf{collapse}(\mathsf{a}, (\mathsf{j}, \mathsf{k}) \to 1), \mathsf{i}) = \begin{cases} \Upsilon(\mathsf{a}, \mathsf{j}) \times \Upsilon(\mathsf{a}, \mathsf{k}) & i = l \\ \Upsilon(\mathsf{a}, i) & i \neq j \end{cases}$$

$$\Upsilon(\mathsf{a} + \mathsf{b}, i) = \Upsilon(\mathsf{a}, i) \cup \Upsilon(\mathsf{b}, i) \qquad \Upsilon(\mathsf{concat}((\mathsf{a}, \mathsf{b}), (j, k) \to 1), i) = \begin{cases} \Upsilon(\mathsf{a}, j) \sqcup \Upsilon(\mathsf{b}, k) & i = l \\ \Upsilon(\mathsf{a}, i) \cup \Upsilon(\mathsf{b}, i) & i \neq l \end{cases}$$

$$\Upsilon(\mathbf{a} \star \mathbf{b}, \mathbf{i}) = \Upsilon(\mathbf{a}, i) \cap \Upsilon(\mathbf{b}, \mathbf{i}) \qquad \qquad \Upsilon(\mathbf{split}(\mathbf{a}, \mathbf{j} \to (\mathbf{k}, \mathbf{l})), \mathbf{i}) = \begin{cases} \pi_0(\Upsilon(\mathbf{a}, j), (|k|, |l|)) & i = k \\ \pi_1(\Upsilon(\mathbf{a}, j), (|k|, |l|)) & i = l \\ \Upsilon(\mathbf{a}, i) \end{cases}$$

Fig. 7. Y maps an array algebra expression and an index to a sequence expression, which represents the non-zero elements being iterated over. For slice, R corresponds to the slice parameters (start, end, and stride). Reductions (sum) are not handled here because lowering rewrites reductions into loops with reduction temporaries when constructing forall loops.

# 393 5 COMPILING TO COITERATING LOOPS

We introduce a second intermediate representation that contains more complex control-flow constructs than CIN's forall loops: CFIR (Control Flow Intermediate Representation). To generate code that efficiently iterates over a sequence expression, BURRITO must reason about when a sparse sequence runs out of elements (i.e. when an array runs out of values), or does not contain a particular coordinate. CFIR allows BURRITO to represent these optimizations while still abstracting away the concrete details of the underlying physical data structures. We first define CFIR before showing how to construct it via a generalized form of iteration lattices [Henry et al. 2021; Kjolstad et al. 2017]. Iteration lattice construction reasons about format properties to determine which sequences should be iterated over, and which can be randomly-accessed into, but lattices and CFIR still abstract away details of the physical data structures underlying the logical arrays. 

# 5.1 Control Flow Intermediate Representation

 We provide the full grammar for CFIR below, whose key components are loops and conditional execution (switch statements).

```
cfir ∷= while idx ← seq (with (seq = seq)*)?) cfir | switch idx (case seq: cfir)+ |
    cfir; cfir | a(I) = expr | a(I) += expr
```

5.1.1 Loops. The while loops of CFIR iterate over a sequence expression until an operand runs out of elements. This construct is useful because a sequence expression and a loop body can be simplified if any sub-sequence is empty, and that simplification produces simpler and more efficient code. Therefore, if a sub-sequence becomes empty, control flow should transition to a loop over a simpler sequence expression with a simpler loop body. For example, consider the example program in Figure 8, which adds a flattened 2D array to a compressed 1D array. If the flattened matrix runs out of non-zero values, then the program only needs to iterate over the compressed vector. Likewise, if the compressed vector runs out of elements, the program only needs to iterate over the flattened matrix. The single CIN loop in Figure 8d is compiled to three CFIR loops in Figure 8f, where the latter two have simplified loop bodies.

CFIR loops also support *locating* into sequences that support random-access (via the with operation) in order to access array values. In certain circumstances, a sequence does not need to be



Fig. 8. Compilation of the fused collapse-and-addition in (a) with the formats in (b) produces the CIN in (d), visually represented by the space-filling curves in (c). Compilation of (d) generates the iteration lattice (e), which constructs the coiterating loops in (f). For the final compiled C output, refer to Figure 19.



(c) Concrete index notation

(e) Control-flow IR.

Fig. 9. Multiplication of a split compressed vector and a CSR matrix exposes an opportunity for the iterate/locate optimization: the *i* loop can iterate over the projection of the sparse vector's sequence, and use *i* to locate into a row of the CSR matrix, possibly skipping some rows. The *j* loop still coiterates the split sequence and the compressed columns of the CSR matrix, as both sequences are sparse. Refer to Figure 20b for the final generated C code from this example.

461 fully iterated. For example, if iterating over  $a \cap b$ , and  $a \subseteq b$ , it is sufficient to iterate over only 462 a, and randomly access into b to get the value labeled by its coordinate. This is the iterate/locate 463 optimization described by Kjolstad et al. [2017], because iterating over a sparse sequence and 464 locating into a dense sequence is asymptotically optimal. BURRITO supports the same optimization, 465 using format properties to detect which sub-sequences support fast (O(1)) location. In general, 466 if a sequence is dense, or constructed from only dense sequences (i.e. the product of two dense 467 sequences), it can and should be located into. For an example, in Figure 9, the sparse sequence  $\pi_0(k_a, [I, J])$  can be used to locate into the dense rows of B. The final C code generated from 468 469 Figure 9 is in Figure 20b, with line 4 performing the locate.

5.1.2 Conditionals. When iterating over a sequence, some sub-sequences may not contain a 471 coordinate that other sub-sequences contain. Conditional execution is required to represent when 472 a loop body should execute a simplified computation based on which sub-sequences contain a 473 particular coordinate. In CFIR, this is handled by the switch construct, which redirects computation 474 based on the value of a coordinate. For example, consider the loops in Figure 12d. When the split 475 vector still has coordinates left, but does not contain a particular *i* coordinate that the CSR matrix 476 does contain, the second loop only needs to iterate over the *j* values of the CSR matrix. This is 477 reflected in lines 16-18 of Figure 12g. 478

# 5.2 Generalized Iteration Lattices

Iteration lattices [Henry et al. 2021; Kjolstad et al. 2017] are ordered state machines that enable reasoning about multi-way merging of coordinate sets. We extend iteration lattices to represent iteration over the sequence expressions defined in Section 4, provide a construction algorithm for iteration lattices, and show how to use iteration lattices to compile a loop in CIN to CFIR loops and conditionals.

An iteration lattice consists of ordered points. A point is labeled by a sequence expression, and represents iterating over that sequence. Each point has child points that represent simplified versions of the parent's sequence expression. Edges to children are labeled by a sub-sequence whose removal from the parent's sequence expression produces the simplified sequence expression in the

1:10

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

470

479

child point. Concretely, a point representing the sequence s has an edge labeled e to a child point representing the sequence t if removing e from s (replacing it with the empty set) and performing simplification produces the sequence t. This simplification is used to construct a lattice: we provide our algorithm for lattice construction in Figure 10, which follows a recursive top-down approach that produces progressively simpler sequence expressions to iterate over.

Consider the iteration lattice in Figure 8e, which repre-496 sents the iteration pattern produced from adding a collapsed 497 498 CSR matrix, A, to a compressed vector, b. If b runs out of values (the sequence  $k_b$  becomes empty), then iteration only 499 needs to be performed over the collapsed matrix. This is 500 represented by the child  $i_A \times j_A$  of the top node in the lat-501 tice, which is transitioned to when  $k_b$  runs out of elements. 502 503 Likewise, if the matrix runs out of elements (the product becomes empty), iteration only needs to proceed over the 504 vector,  $k_b$ . This results in the leftmost and rightmost nodes 505 in the lattice, respectively. Note that the transition to the 506 node represented by  $k_b$  is labeled only  $i_A$ , not  $j_A$ . This is due 507

```
func ConstructLattice(seq):
    point = LatticePoint(seq)
    // Transition sub-sequences
    // (Section 5.2.1 and Figure 11)
    edges = χ(seq)
    for sub in edges:
        // Simplification
        // (Section 5.2.2)
        r = remove(sub, seq)
        s = simplify(r)
        l = ConstructLattice(s)
        point.add_child(sub, l)
    return point
```

Fig. 10. Top-down lattice construction.

to the semantics of product: if  $j_A$  runs out of elements,  $i_A$  needs to be stepped forward, which means that a product *only runs out* when the first operand runs out of coordinates.

5.2.1 Edge Sequences. Edges in an iteration lattice correspond to sequences that can introduce a
 state transition. These are often sequences produced by array levels, but can also correspond to
 slices or projections. This is because a slice (or projection) can run out before the sequence it is
 slicing (or projecting) runs out, and should induce a state transition when that happens.

<sup>515</sup> We give an algorithm for collecting edges from a sequence expression in Figure 11, the function  $\chi$ . <sup>516</sup> It is a recursive algorithm that generates a set of sub-sequences that correspond to state transitions. <sup>517</sup> Child lattice points are generated by replacing an edge sequence with an empty set and simplifying <sup>518</sup> the sequence expression, which we discuss further below.

Consider the expression in Figure 12, which reshapes a compressed vector into a matrix and adds it to a CSR matrix. The *j* loop over the rows of the matrices changes state when: 1) the compressed vector runs out of elements; 2) the logical slice of the compressed vector that corresponds to a single row in the logical matrix runs out; and 3) the CSR matrix row runs out. These three cases are edges in Figure 12f, as each triggers a state transition.

525 5.2.2 Sequence Simplification. Sequence simplification follows simple set rules: e.g. if one side of an 526 intersection is empty, the entire intersection is empty. Likewise, if one side of a union is empty, then 527 return the other side of the union. Our algorithm treats sparse-and-empty and dense-and-empty 528 differently. For example, if the set expression is  $a \cup b$  and a is a dense sequence, when a becomes 529 empty (dense-and-empty), the entire union is empty, because the sequence must have been fully 530 iterated if a ran out of elements. This optimization can be seen in Figure 12e, where the dense

 $\chi(i_a) = \{i_a\} \qquad \chi(x \cup y) = \chi(x) \cup \chi(y)$   $\chi(x \cap y) = \chi(x) \cup \chi(y) \qquad \chi(x \times y) = \chi(x)$   $\chi(x \sqcup y) = \chi(x) \qquad \chi(\pi_k(x, S)) = \{\pi_k(x, S)\} \cup \chi(x)$   $\chi(x[\mathbb{R}]) = \{x[\mathbb{R}]\} \cup \chi(x) \qquad \chi(s + x) = \chi(x)$ 

- Fig. 11.  $\chi$  generates a set of sub-sequences whose removal induces state transitions.
- 538 539

537

510

524



Fig. 12. Addition of a split compressed vector and a CSR matrix has the same logical iteration pattern as the multiplication in Figure 9, but requires very different loops. The iteration lattice in (e) highlights the asymmetrical sequence simplification from treating sparse-and-empty and dense-and-empty as two different simplification results: when the dense  $i_B$  iterator runs out of coordinates, the sparse  $\pi_0(k_a, [I, J])$  iterator must have run out of elements, so no lattice point is needed to represent iteration over solely  $\pi_0(k_a, [I, J])$ .

column of a CSR matrix is conterated with a sparse sequence – when the dense sequence runs out, the entire union must run out.

Sub-points. Given a lattice point p that iterates over a sequence s, the sub-points of p are 5.2.3 567 descendant points that represent sequences that are strict subsets of s. For example, in Figure 14e, 568 the sequence  $i_A \cup i_B$  has two sub-points, labeled  $i_A$  and  $i_B$ . Sub-points of p correspond to a sequence 569 that represents a strictly *simpler* state than the state represented by p. Not all descendant points 570 are sub-points due to the complexity of the concatenation operator. Consider the iteration lattice 571 in Figure 14f; the concatenation operator produces a state transition when  $k_A$  runs out, indicating 572 a transition to start iterating over the slice of  $l_B$  instead. However, this point is not considered a 573 sub-point, because it represents a completely disjoint set (it requires  $k_A$  runs out, not just that  $k_A$  is 574 missing an element). To compute the sub-points of a lattice point p, we gather all descendants that 575 represent a strict subset of *p*'s sequence. 576

#### **Compiling Lattices to Loops** 5.3

579 Iteration lattices are used to generate loops over progressively simpler loop bodies. A lattice always 580 maintains a partial ordering, because each point has a progressively simpler sequence expression. The lattice can be compiled to loops by topologically sorting the points, laying them out in order, 582 and generating a CFIR loop for each point that runs until an edge sequence runs out. A CFIR loop 583 exits when an edge sequence runs out, as this indicates that the iteration should be passed to a loop 584 over a simpler sequence expression. The body of a CFIR loop is generated as a conditional over 585 the sub-points in the lattice, where the bodies of the conditional statements contain recursively 586 simplified code. We provide the full compilation from CIN forall loops to CFIR loops in Figure 13.

559

560

561

562 563

564

565 566

577 578

581

589	Consider the CIN in Figure 12d, where the $i$	
590	loop iterates over a projection of a sparse itera-	
591	tor unioned with a dense iterator. This loop gen-	
592	erates the iteration lattice in Figure 12e, which	
593	contains two states, the initial state (top), and a	
594	state that iterates over only the dense iterator	
595	(middle left), for when the sparse iterator runs	
596	out of elements. This iteration lattice compiles	
597	to the two outer loops in Figure 12g, one of	
598	which iterates over the union, and the other	
599	handles the case where the sparse iterator runs	
600	out of elements. Within the first outer loop,	
601	there is control flow to handle the case that	
602	the sparse iterator is at the same coordinate	Fig
603	as the dense iterator, and the case where the	CF

. 13. Compilation of CIN forall loops to a chain of IR while loops over conditional statements.

dense iterator contains an element that the projection of the sparse iterator does not have. These cases within this first loop correspond to sub-states that the computation could be in, avoiding unnecessary work when the projection of the sparse iterator does not contain a coordinate.



Fig. 14. Concatenation produces a union over the shared dimensions, and iteration lattice construction for i generates loops specialized to when a coordinate in the sequence is both sequences or only one, producing simpler loop bodies in the latter case. Iteration lattice construction for *j* essentially fissions the loop over concatenated sequences into a pair of loops (lines 4-5 and 6-7 of (g)) that each coiterate a single sequence.

#### **CODE GENERATION**

In this section, we show how a simple set of primitives compose to produce iteration over any sequence expression. This final code generation pass uses a simple iterator model that builds on the indexed stream model of Kovach et al. [2023], a representation used for compiling fused sparse

```
Dense(String name, String idx,
638
                                                                              Compressed(String name, String idx,
String lb, String ub):
                      String size):
639
              Init():
    emit "int {idx}_{name} = 0;"
                                                                               Init():
    emit "int {idx}p_{name} = {lb};"
         3
                                                                         3
                                                                          4
640
              Valid():
emit "{idx}_{name} < {size}"</pre>
                                                                               Valid():
emit "{idx}p_{name} < {ub}"</pre>
         6
                                                                          6
641
642
                                                                         8
              Eval():
    emit "{idx}_{name}"
         9
                                                                               Eval():
emit "{name}_crd[{idx}p_{name}]"
                                                                         9
643
                                                                         11
644
              Equals(i):
                                                                               Equals(i):
        13
               emit "{idx}_{name} == {i}"
                                                                                emit "Eval() == {i}"
645
        14
                                                                         14
              Next():
emit "{idx}_{name}++;"
646
                                                                               Next():
    emit "{idx}p_{name}++;"
        16
                                                                         16
647
              Locate(i):
    emit "{idx}_{name} = {i};"
        18
                                                                         18
                                                                               Locate(i):
emit "{idx}p_{name} = binary_search({i},
648
        19
                                                                         19
                                                                                 {name}_crd, {lb}, {ub});'
649
                                                                         20
650
```

Fig. 15. Iterator model for a dense dimension.

Fig. 16. Iterator model for a compressed dimension.

tensor algebra that shows how to iterate over unions and intersections efficiently. We first describe the iterator model for sequence combinators, then for formats, and lastly provide the complete code generation algorithm from CFIR to C code.

#### 6.1 Iterator Model

1:14

651 652

653

654

655 656

657

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

BURRITO relies on a small set of composable primitives to implement iteration over combinations 658 of sequences: initialize, valid, evaluate, equals, next, and locate. We describe each below, giving 659 examples of the iterator model for sequence products in Figure 17 and slicing in Figure 18. Note 660 that Figure 18 shows a generalization of prior work [Henry et al. 2021], which supported slicing 661 only array dimensions, while our algorithm can generate code that iterates over the slice of *any* 662 sequence. 663

**Initialize** Initialization handles declaring any necessary iteration variables and locating the first non-zero element of a sequence.

- Valid Valid is a check that the sequence has not run out of coordinates, meaning no subsequence has run out. This is a check that no iterators have gone out of bounds, and slices and projections are still within a valid range.
- **Evaluate** Evaluation computes the current coordinate value of the sequence. For unions and intersections, this takes the minimum of the two sequences, but for sequence combinators, a sequence value (or values) must be re-mapped to the new space. For example, a product applies the strided offset formula to map two sequence coordinates to a single sequence coordinate. Likewise, a slice must re-map a sequence value to the sliced coordinate space.
- **Equals** An equality test checks whether a sequence is currently at a coordinate. For unions and intersections, equality simply means that both sequence operands are equal to the given value. For combinators, we re-map the provided value into the spaces spanned by the operand sequences. This is used to compile CFIR's conditional statements (see Figure 20a).
- **Next** The next method steps forward any sequences that are currently behind. For array levels, unions, and intersections, this performs the same conditional updates as TACO (i.e. line 25 of Figure 19). For sequence combinators, we define modular code generators for stepping forward a sequence.
- **Locate** Location moves physical iterators to point to a specific logical coordinate. Section 5.1.1 683 described cases where random-access into a sequence is asymptotically preferable, and the 684 locate interface is used to support this optimization. Many sequence operators are invertible 685

682

```
687
          Product(Seq a, Seq b):
                                                               Slice(Seq a, Expr s, Expr e, Expr r):
688
           Init():
                                                                Init():
             emit Init(a); Init(b)
                                                                 emit Init(a)
emit Locate(a,
       3
                                                           3
689
             emit while(Valid(a) && !Valid(b))
                                                           4
                                                                                 S)
       4
                                                                 emit while(Valid() && ((Eval(a)-s)%r))
             emit { Next(a); Init(b); }
                                                           5
690
                                                                 emit { Next(a); }
                                                           6
           Valid():
691
            emit Valid(a) && Valid(b)
                                                               Valid():
                                                           8
692
                                                           0
                                                                 emit Valid(a) && Eval(a) < e
           Eval():
                                                           10
      10
693
            emit (Eval(a) * |b|) + Eval(b)
                                                                Eval():
694
                                                                 emit (Eval(a) - s) / r
           Equals(i):
            emit Equals(i/|b|,a) && Equals(i%|b|,b)
                                                                Equals(i):
695
                                                           14
      14
                                                                 emit Equals((i * r) + s, a)
696
           Next():
                                                           16
                                                                Next():
            emit Next(b)
697
             emit while(Valid(a) && !Valid(b))
                                                                 emit do { Next(a) }
      18
                                                           18
             emit { Next(a); Init(b); }
                                                                 emit while(Valid() && ((Eval(a)-s)%r))
698
      19
                                                           19
      20
                                                           20
699
                                                                Locate(i):
           Locate(i):
            emit Locate(i/|b|,a); Locate(i%|b|,b)
```

emit Locate((i \* r) + s, a)

Fig. 17. Iterator interface for sequence products.

Fig. 18. Iterator interface for sequence slices.

```
1
           int i A = 0:
                                                                                       // init iA
703
           int jp_A = A_pos[i_A];
                                                                                       // init j_A
704
                                                                                       // valid i_A && !valid j_A
       3
           while ((i_A < I) && !(jp_A < A_pos[i_A+1])) {</pre>
                                                                                       // next i_A
       4
             i_A++;
705
             jp_A = A_pos[i_A];
                                                                                       // init j_A
       5
706
          }
       6
707
           int kp_b = b_pos[0];
       7
                                                                                       // init k_b
           while ((i_A < N) && (jp_A < A_pos[i_A+1]) && (kp_b < b_pos[1])) { // valid i_A \times j_A \cup k_b
       8
708
       9
             int k = min((i_A * M) + A_crd[jp_A], b_crd[kp_b]);
                                                                                       // eval (i_A \times j_A) \cup k_b
709
             if (((i_A == k / M) && (A_crd[jp_A] == k % M)) &&
       10
                                                                                       // equals k, (i_A \times j_A) \cup k_b
710
                  (k == b_crd[kp_b])) {
               c[k] = A.values[jp_A] + b.values[kp_b];
711
             } else if ((i_A == k / M) && (A_crd[jp_A] == k % M)) {
                                                                                       // equals k, i_A \times j_A
712
       14
               c[k] = A.values[jp_A];
713
             } else if (k == b_crd[kp_b]) {
                                                                                       // equals k, k_h
               c[k] = b.values[kp_b];
714
       16
             }
715
             if (k == ((i_A * M) + A_crd[jp_A])) {
                                                                                       // next k, (i_A \times j_A)
       18
716
                                                                                       // next j_A
       19
                jp_A++;
                while ((i_A < N) \& !(jp_A < A_pos[i_A+1])) \{
       20
717
                                                                                       // next i_A
       21
                 i A++:
718
                  jp_A = A_pos[i_A];
                                                                                       // init j_A
       22
719
               }
720
       24
             3
             kp_b += (k == b_crd[kp_b]);
                                                                                       // next k, k_b
721
           }
       26
```

Fig. 19. Compilation of the first loop of Figure 8f, which iterates over a collapsed CSR matrix, A, and adds it to a compressed vector, b.

and this property can be used to locate through them, and base data structure formats often support fast random access (i.e. dense or hashed formats).

727 728 729

722

723

724 725

726

700 701

702

6.1.1 *Formats.* The iterator model for array formats is equivalent to the base case of the recursive 730 algorithm for generating iteration code, as tensor dimensions are the main primitive in sequence 731 expressions. As examples, we provide the implementations of dense and compressed formats in 732 Figures 15 and 16, respectively. Extending BURRITO to support an additional sparse format only 733 requires implementing the iterator interface for that format type. Our prototype compiler supports 734

TACO's dense, compressed, and singleton level formats [Chou et al. 2018], which compose to express 736 many sparse data structures, such as CSR, CSC, COO, and variants. 737

738 6.1.2 *Combinators.* The iterator model fully composes to enable code generation for all sequence 739 combinators. As examples, we provide the implementation of the iterator model for sequence prod-740 ucts and slices in Figures 17 and 18, respectively. These implementations show the composability of this interface. We also illustrate a labeled example of the full C code generated from Figure 8 in 742 Figure 19, with line labels corresponding to the interface that generated each line. 743

#### 744 6.2 Compiling CFIR 745

We give the algorithm for compiling CFIR based on the iterator model in Figure 20a. Intuitively, 746 each loop iterates until some sub-sequence runs out (the state is permanently changed), and after 747 execution of a loop's body, steps the sequence forward. Conditional statements are used to evaluate 748 which sub-state the iterator may currently be in, and generate if-else chains. The iterator model 749 allows for a very simple code generation algorithm to compile the abstract while loops and switch 750 statements of CFIR to the complex irregular loops over physical data structures required to support 751 compute and shape operators. 752

# 7 EVALUATION

Our evaluation proves the following claims:

- (1) Generated shape operators can match the performance of hand-written shape operators.
- (2) Portability across data structures offers performance benefits over fixed-format kernels.
  - (3) Fusion of shape and compute operators can improve performance.

```
759
          func CompileCFIR(stmt):
                                                           4|int kp_a = a_pos[0];
760
          match stmt with
       2
           | While (idx, seq, locs, body) ->
                                                           5|while (kp_a < a_pos[1]) {
761
       3
            emit Init(seq)
                                                           6| int i = a_crd[kp_a] / J;
       4
762
           emit while (Valid(seq)) {
                                                           8| int i_B = i;
      5
763
            emit int idx = Eval(seq);
                                                           4| int jp_B = B_pos[i];
       6
            for a, b \in locs
764
            emit Locate(Eval(a), b)
                                                               while ((kp_a < a_pos[1]) &&</pre>
      8
765
      9
           emit CompileCFIR(body)
                                                           5
                                                                      (i == a_crd[kp_a] / J) &&
766
      10
           emit Next(seq) }
                                                                       (jp_B < B_pos[i_B+1])) \{
                                                                 int j0 = a_crd[kp_a] % J;
767
           | Switch (idx, seqs, bodies) ->
      11
           emit if (Equals(idx, seqs[0])) {
                                                           6
                                                                 int j1 = B_crd[jp_B];
      12
768
                                                                  int j = min(j0, j1);
            emit CompileCFIR(bodies[0]) }
      13
769
      14
            for s, b \in zip(seqs, bodies)[1:]
                                                          121
                                                                  if ((j == j0) && (j == j1)) {
770
             emit else if (Equals(idx, s)) {
                                                       21-22
                                                                   C[i*J + j] = a[kp_a] * B[jp_B];
                                                          13 |
             emit CompileCFIR(b) }
      16
                                                                 }
771
           | Pair (cfir0, cfir1) ->
                                                                 kp_a += (j == j0);
772
      18
          emit CompileCFIR(cfir0)
                                                          10
                                                                 jp_B += (j == j1);
773
      19
           emit CompileCFIR(cfir1)
           | Assign (array, idxs, expr) ->
      20
774
           emit CompileWrite(array, idxs)
      21
                                                               while ((kp_a < a_pos[1]) &&</pre>
775
      22
           emit CompileExpr(expr)
                                                                      (i == a_crd[kp_a] / J)) {
                                                           10
776
      23
           | Reduce (array, idxs, expr) ->
                                                                 kp_a++;
      24
           emit CompileReduction(array, idxs)
                                                               }
777
                                                             |}
      25
            emit CompileExpr(expr)
778
```

(a) Recursive codegen via the iterator model.

(b) Generated code from the CFIR in Figure 9e.

Fig. 20. (a) C code generation from CFIR. (b) shows the generated code for  $C(i, j) = split(a(k), k \rightarrow j)$ (i, j) \* B(i, j), where a is a compressed vector and B is a  $I \times J$  CSR matrix. The labels to the left in (b) correspond to the line numbers in (a) that generated that particular line of code.

1:16

741

753

754

755

756

757

758

782 783 784

779

780

We first describe the existing state-of-the-art libraries that we compare to in Section 7.1, provide 785 our benchmarking methodology in Section 7.2, and then provide evidence for the above claims in 786 the following sections. 787

#### **Baselines** 789 7.1

788

802

805

806

807

811

813

814

815

816

817 818

819

790 We compare BURRITO-generated code to scipy.sparse and pydata/sparse, the two state-of-the-791 art libraries that have the most complete set of shape and compute operators. For shape operators, 792 both of these libraries follow a simple reduction approach. Each shape operator is implemented 793 for one or a few sparse data structures, and calling a shape operator on an unsupported sparse 794 data structure incurs a conversion cost to convert to a supported format. For example, to reshape a 795 CSR matrix with scipy, the library first converts to the CSR matrix to COO (the COOrdinate list 796 data structure) and then calls COO. reshape. Likewise, when a user requests a non-standard output 797 format (i.e. reshape a COO matrix into a CSR matrix), the library will perform the operation with 798 a supported implementation and then convert the output to the requested output format. Note 799 that for sparse tensor algebra alone, BURRITO's compilation technique is equivalent to TACO's, and 800 achieves the same performance. 801

#### **Methodology and Benchmark Notation** 7.2

803 We evaluate on an Apple M1 Pro (3.2 GHz, 8 cores) with 16 GB of RAM. All benchmarks are 804 single-threaded. We use the 12 largest SuiteSparse matrices that have indices fitting in 32 (unsigned) bits. These matrices span multiple domains, including computer vision, structural engineering, economics, graph analytics, and computational fluid dynamics. Each benchmark record is the minimum time of 5 iterations. For benchmarks that require multiple matrices (i.e. concatenation), 808 we first split the SuiteSparse matrix in half, and use the halves as operands to concatenation. For 809 fusion benchmarks, we use the same matrix shifted by a small number (10), due to the need for 810 shape compatibility, as in prior evaluations of sparse tensor algebra [Kjolstad et al. 2017]. Note that BURRITO's compile times are interactive, so compilation does not introduce noticeable overhead. 812

We label benchmarks with short descriptions of their compute kernels. vstack means vertical stacking (concatenation along the first axis), hstack means horizontal stacking (concatenation along the second axis), and reshape is a collapse followed by a split. Benchmarks label operands with their array types (i.e. CSR or COO), and single letter labels ("C" or "D") correspond to randomlygenerated vectors ("Compressed" and "Dense", respectively).

#### **Comparison to Hand-written Kernels** 7.3

We compare BURRITO's generated shape operators to shape operators that scipy.sparse and 820 pydata/sparse whose implementations do not perform data structure conversions. This com-821 parison shows that a compilation-based approach can match the performance of hand-written 822 code. 823

The comparison to scipy.sparse is shown in Figure 21. On stacking COO matrices and slicing 824 CSR matrices, performance is matched. On reshaping COO and horizontally stacking CSR, BURRITO 825 out-performs scipy. Both of these performance benefits come from the fact that scipy's implemen-826 tations allocate a small number of temporary numpy arrays to perform conversions, while BURRITO 827 allocates only the output array. On vertical stacking CSR matrices, scipy out-performs BURRITO: 828 that particular operation is little more than a few memcpys, and scipy has a highly optimized (vec-829 torized) implementation of it, while the underlying C compiler used to compile BURRITO's generated 830 code did not generate the same optimized memcpys. However, in general, these benchmarks show 831 that generated code can match or exceed the performance of hand-written reshape kernels. 832



(a) Runtime speed-up over scipy.sparse.

SHARE BE STONE OF STATISTICS OF STATISTICS



(b) Runtime speed-up over pydata/sparse.

Fig. 23. Evaluation of shape operators where existing SOTA techniques perform data structure conversions on at least one input array in order to perform the shaping operation.

We also compare to pydata/sparse is shown in Figure 22. BURRITO's generated code generally severely out-performs pydata because pydata's implementions are dimensionality-agnostic (i.e. not specialized to a 2D COO, but works for an *n*D COO matrix), which is beneficial for reducing overall code size but not beneficial for code performance. This highlights the cost of generality in these systems. Like scipy, pydata out-performs BURRITO on vertically stacking CSR matrices, due to the same memcpy optimization.





Fig. 24. Runtime speed-up over scipy.sparse of shape operators applied to SuiteSparse matrices, for kernels that scipy.sparse requires data structure conversions on the output of the shape operator.



Fig. 25. Runtime speed-up over pydata/sparse of shape operators applied to SuiteSparse matrices, for kernels that pydata/sparse requires data structure conversions on the output of the shape operator.



Fig. 26. Runtime speed-up over BURRITO-generated *unfused* code. The labeled geomeans show the speed-up of BURRITO-generated *fused* code.

# 7.4 Comparison to Reduced Kernels

To show the importance of allowing code to specialize to particular input and output data structures, we evaluate shape operators with data structures that the state-of-the-art libraries perform data structure conversions to compute. In Figure 23, we evaluate BURRITO generated code against implementations that first perform a data structure conversion on an input array before applying the hand-written shape operator. Figure 24 and Figure 25 show the performance difference of applying a shape operator with an output format that is different from the standard library implementation, requiring a data structure conversion following the shape operator. In each of these benchmarks, we see that generating code specialized to particular input and output data structures (avoiding data structure conversions) achieves performance benefits. BURRITO generated code also uses less memory, as it only allocates the output array, but not any temporary arrays.

# 7.5 Shape and Compute Operator Fusion

To evaluate the performance benefits of fusing shape with compute operators, we perform a series of self-comparisons on kernels with both shape and compute operators. For relevance to the

state-of-the-art, we also compare to scipy.sparse, the faster of the two libraries that we evaluate
against in the prior sections. Figure 26 provides a performance evaluation of the speed-up gained
by fusing shape and compute kernels. The benefit of fusing such operators is largely a result from
reduced memory allocation, as there is no need to allocate expensive intermediate temporaries.
There are some cases where scipy.sparse out-performs the *unfused* BURRITO-generated code (e.g.
with the fused SpMV, where scipy.sparse's fast vstack, discussed in Section 7.3, out-performs
BURRITO's vstack), but *fused* BURRITO-generated code performs best across these benchmarks.

# 8 RELATED WORKS

*Sparse Shape Operator Compilation.* As discussed in the introduction, most prior work in sparse array compilation focuses on compiling compute operators [Bik et al. 2022; Kjolstad et al. 2017; Zhao et al. 2022], with the exceptions of two compilers: Henry et al. [2021] extended TACO to support iterating over slices of array operands, and allows computing over a slice, but does not support slicing intermediate computation, which limits fusion options; Looplets [Ahrens et al. 2023] can be used to express the concatenation of arrays, but does not express concatenation as an operator in the front-end language. BURRITO explicitly expresses shape operators in the front-end language, and has no restrictions on mixing compute and shape operators.

*Abstracting Sparse Iteration.* There are decades worth of work in abstracting sparse iteration. The database community relied on the iterator model [Graefe 1994] to implement many relational operators, and more recently, Kovach et al. [2023] introduced the stream model for iterating over sparse arrays, which supports a similar iterator interface based on a model equivalent to *init-validnext.* Our iterator model builds on Kovach et al. [2023], and we introduce additional primitives to support shape operators. Chou et al. [2018], Looplets [Ahrens et al. 2023], and SparseTIR [Ye et al. 2023] provide various abstractions over sparse data formats, but do not use these abstractions to compile shape operators.

Avoiding Discordant Traversals. Sparse tensor algebra has a long history of shaping computation to avoid discordant traversals [Bik 1996; Bik et al. 1994; Duff et al. 1990; George and Liu 1981; Gustavson 1972; Pissanetsky 1984; Tewarson 1973; Zlatev 1991], and a compiler for efficient sparse shape operators must do the same. Kjolstad et al. [2019] introduced a simple loop IR and scheduling rewrites that allow for avoiding discordant traversals, and BURRITO uses that IR (CIN) and supports the same transformations, for the same reasons.

*Sparse Array Libraries.* There are a number of sparse array libraries [Abadi et al. 2016; Paszke et al. 2019; Virtanen et al. 2020] that implement some number of shape operators. We compare to the most complete of these, scipy.sparse and pydata/sparse, in Section 7. These array libraries are feature-incomplete, generally only supporting a small number of array formats and a small number of operators. BURRITO can be used to generate custom shape operator implementations for each of these libraries, or replace them entirely.

*Array Languages.* There are decades of work in dense array programming languages [Bader and Kolda 2008; Harris et al. 2020; Henriksen et al. 2017; Iverson 1962; Liu et al. 2022; Ragan-Kelley et al. 2013; Steuwer et al. 2017], many of which support shape operators, but only for dense arrays. For many of these languages, shape operators correspond to zero-cost array metadata edits, and do not require iteration over the data like sparse shape operators do.

*Staged Compilation.* While the database community typically uses the iterator model as a runtime technique, recent work [Tahboub et al. 2018; Tahboub and Rompf 2020] applies ideas from partial evaluation to enable using the iterator model for code generation. BURRITO's code generation can be seen as an application of the same idea to compiling iteration over sequence expressions.

### 981 9 CONCLUSION

We describe the first compiler for a sparse array programming language with both shape and
 compute operators. We show how a simple declarative array language can be compiled to imperative
 loops over sequences, how to generate optimized loops that coiterate these sequences, and lastly,
 how to generate data-structure-specific code via a simple iterator model. With these ideas, sparse
 array programming moves one step closer to the completeness that dense array programming
 systems have long since achieved.

# ACKNOWLEDGMENTS

### 991 REFERENCES

988

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. TensorFlow: A System for Large-Scale Machine Learning. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation* (Savannah, GA, USA) (OSDI'16). USENIX Association, USA, 265–283.
- Hameer Abbasi. 2018. Sparse: a more modern sparse array library. In *Proceedings of the 17th python in science conference*.
   27–30.
- Willow Ahrens, Daniel Donenfeld, Fredrik Kjolstad, and Saman Amarasinghe. 2023. Looplets: A Language for Structured Coiteration. In *Proceedings of the 21st ACM/IEEE International Symposium on Code Generation and Optimization* (Montréal, QC, Canada) (*CGO 2023*). Association for Computing Machinery, New York, NY, USA, 41–54. https://doi.org/10.1145/ 3579990.3580020
- Brett W. Bader and Tamara G. Kolda. 2008. Efficient MATLAB Computations with Sparse and Factored Tensors. SIAM
   *Journal on Scientific Computing* 30, 1 (2008), 205–231.
- Aart Bik, Penporn Koanantakool, Tatiana Shpeisman, Nicolas Vasilache, Bixia Zheng, and Fredrik Kjolstad. 2022. Compiler
   Support for Sparse Tensor Computations in MLIR. ACM Trans. Archit. Code Optim. 19, 4, Article 50 (sep 2022), 25 pages. https://doi.org/10.1145/3544559
- Aart J.C. Bik. 1996. Compiler Support for Sparse Matrix Computations. Ph. D. Dissertation. Department of Computer Science,
   Leiden University. ISBN 90-9009442-3.
- Aart J.C. Bik, Peter M.W. Knijenburg, and Harry A.G. Wijshoff. 1994. Reshaping Access Patterns for Generating Sparse
   Codes. In *Proceedings of the 7th International Workshop on Languages and Compilers for Parallel Computing (LCPC '94)*.
   Springer-Verlag, Berlin, Heidelberg, 406–420.
- William E. Boyse and Andrew A. Seidl. 1996. A Block QMR Method for Computing Multiple Simultaneous Solutions to
   Complex Symmetric Systems. *SIAM Journal on Scientific Computing* 17, 1 (1996), 263–274.
- 1011Stephen Chou, Fredrik Kjolstad, and Saman Amarasinghe. 2018. Format Abstraction for Sparse Tensor Algebra Compilers.1012Proc. ACM Program. Lang. 2, OOPSLA, Article 123 (oct 2018), 30 pages. https://doi.org/10.1145/3276493
- 1013 Iain S. Duff, A.M. Erisman, and J.K. Reid. 1990. Direct Methods for Sparse Matrices. Oxford Science Publications, Oxford.
- Alan George and Joseph W.H. Liu. 1981. Computer Solution of Large Sparse Positive Definite Systems. Prentice Hall, Englewood Cliffs, New York.
- 1015G. Graefe. 1994. Volcano- An Extensible and Parallel Query Evaluation System. IEEE Trans. on Knowl. and Data Eng. 6, 11016(feb 1994), 120-135. https://doi.org/10.1109/69.273032
- Fred G. Gustavson. 1972. Some Basic Techniques for Solving Sparse Systems of Linear Equations. In Sparse Matrices and Their Applications, Donald J. Rose and Ralph A. Willoughby (Eds.). Plenum Press, New York, NY, 41–52.
- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy. *Nature* 585, 7825 (Sept. 2020), 357–362. https://doi.org/10.1038/s41586-020-2649-2
- Troels Henriksen, Niels G. W. Serup, Martin Elsman, Fritz Henglein, and Cosmin E. Oancea. 2017. Futhark: Purely
   Functional GPU-Programming with Nested Parallelism and in-Place Array Updates. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Barcelona, Spain) (*PLDI 2017*). Association
   for Computing Machinery, New York, NY, USA, 556–571. https://doi.org/10.1145/3062341.3062354
- Rawn Henry, Olivia Hsu, Rohan Yadav, Stephen Chou, Kunle Olukotun, Saman Amarasinghe, and Fredrik Kjolstad. 2021.
   Compilation of Sparse Array Programming Models. *Proc. ACM Program. Lang.* 5, OOPSLA, Article 128 (oct 2021), 29 pages.
   https://doi.org/10.1145/3485505

- Kenneth E. Iverson. 1962. A Programming Language. In *Proceedings of the May 1-3, 1962, Spring Joint Computer Conference* (San Francisco, California) (*AIEE-IRE '62 (Spring)*). Association for Computing Machinery, New York, NY, USA, 345–351.
   https://doi.org/10.1145/1460833.1460872
- Fredrik Kjolstad, Willow Ahrens, Shoaib Kamil, and Saman Amarasinghe. 2019. Tensor Algebra Compilation with Workspaces. In *Proceedings of the 2019 IEEE/ACM International Symposium on Code Generation and Optimization* (Washington, DC, USA) (*CGO 2019*). IEEE Press, New York, NY, USA, 180–192.
- Fredrik Kjolstad, Shoaib Kamil, Stephen Chou, David Lugato, and Saman Amarasinghe. 2017. The Tensor Algebra Compiler.
   *Proc. ACM Program. Lang.* 1, OOPSLA, Article 77 (oct 2017), 29 pages. https://doi.org/10.1145/3133901
- Scott Kovach, Praneeth Kolichala, Tiancheng Gu, and Fredrik Kjolstad. 2023. Indexed Streams: A Formal Intermediate Representation for Fused Contraction Programs. *Proc. ACM Program. Lang.* 7, PLDI, Article 154 (jun 2023), 25 pages. https://doi.org/10.1145/3591268
- Amanda Liu, Gilbert Louis Bernstein, Adam Chlipala, and Jonathan Ragan-Kelley. 2022. Verified Tensor-Program Optimization via High-Level Scheduling Rewrites. *Proc. ACM Program. Lang.* 6, POPL, Article 55 (jan 2022), 28 pages. https://doi.org/10.1145/3498717
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin,
   Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan
   Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. *PyTorch: An Imperative* Style, High-Performance Deep Learning Library. Curran Associates Inc., Red Hook, NY, USA.
- 1045 Sergio Pissanetsky. 1984. Sparse Matrix Technology. Academic Press, London.
- Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. 2013.
   Halide: A Language and Compiler for Optimizing Parallelism, Locality, and Recomputation in Image Processing Pipelines. In Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (Seattle, Washington, USA) (PLDI '13). Association for Computing Machinery, New York, NY, USA, 519–530. https://doi.org/10.
   1145/2491956.2462176
- Justin Solomon. 2015. Numerical Algorithms: Methods for Computer Vision, Machine Learning, and Graphics. A. K. Peters,
   Ltd., USA.
- Michel Steuwer, Toomas Remmelg, and Christophe Dubach. 2017. Lift: a functional data-parallel IR for high-performance
   GPU code generation. In *Proceedings of the 2017 International Symposium on Code Generation and Optimization, CGO* 2017, Austin, TX, USA, February 4-8, 2017, Vijay Janapa Reddi, Aaron Smith, and Lingjia Tang (Eds.). ACM, New York, NY,
   USA, 74–85. http://dl.acm.org/citation.cfm?id=3049841
- Ruby Y. Tahboub, Grégory M. Essertel, and Tiark Rompf. 2018. How to Architect a Query Compiler, Revisited. In *Proceedings* of the 2018 International Conference on Management of Data (Houston, TX, USA) (SIGMOD '18). Association for Computing
   Machinery, New York, NY, USA, 307–322. https://doi.org/10.1145/3183713.3196893
- Ruby Y. Tahboub and Tiark Rompf. 2020. Architecting a Query Compiler for Spatial Workloads. In *Proceedings of the 2020* ACM SIGMOD International Conference on Management of Data (Portland, OR, USA) (SIGMOD '20). Association for Computing Machinery, New York, NY, USA, 2103–2118. https://doi.org/10.1145/3318464.3389701
- 1060 Reginal P. Tewarson. 1973. Sparse Matrices. Academic Press, New York, NY.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia
   Polosukhin. 2017. Attention is All You Need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (Long Beach, California, USA) (*NIPS'17*). Curran Associates Inc., Red Hook, NY, USA, 6000–6010.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski,
  Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod
  Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, Ilhan Polat, Yu Feng,
  Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R.
  Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors.
  2020. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods* 17 (2020), 261–272.
- 1068 https://doi.org/10.1038/s41592-019-0686-2
- Zihao Ye, Ruihang Lai, Junru Shao, Tianqi Chen, and Luis Ceze. 2023. SparseTIR: Composable Abstractions for Sparse
   Compilation in Deep Learning. In Proceedings of the 28th ACM International Conference on Architectural Support for
   Programming Languages and Operating Systems, Volume 3 (Vancouver, BC, Canada) (ASPLOS 2023). Association for
   Computing Machinery, New York, NY, USA, 660–678. https://doi.org/10.1145/3582016.3582047
- Tuowen Zhao, Tobi Popoola, Mary Hall, Catherine Olschanowsky, and Michelle Strout. 2022. Polyhedral Specification and Code Generation of Sparse Tensor Contraction with Co-Iteration. ACM Trans. Archit. Code Optim. 20, 1, Article 16 (dec 2022), 26 pages. https://doi.org/10.1145/3566054
- 1075 Zahari Zlatev. 1991. Computational Methods for General Sparse Matrices. Kluwer Academic Publishers, Dordrecht.
- 1076 1077
- 1078

1:22