
Optimized Sparse Matrix Operations for Reverse Mode Automatic Differentiation

Nicolas Nytko¹ Ali Taghibakhshi² Tareq Uz Zaman³ Scott MacLachlan⁴ Luke N. Olson¹ Matt West²

Abstract

Sparse matrix representations are ubiquitous in computational science and machine learning, leading to significant reductions in compute time, in comparison to dense representation, for problems that have local connectivity. The adoption of sparse representation in leading ML frameworks such as PyTorch is incomplete, however, with support for both automatic differentiation and GPU acceleration missing. In this work, we present an implementation of a CSR-based sparse matrix wrapper for PyTorch with CUDA acceleration for basic matrix operations, as well as automatic differentiability. We also present several applications of the resulting sparse kernels to optimization problems, demonstrating ease of implementation and performance measurements versus their dense counterparts.

1. Introduction

The popularity of automatic differentiation (AD) support in libraries such as PyTorch (Paszke et al., 2019), Tensorflow (Abadi et al., 2016), or Jax (Bradbury et al., 2018) has exploded in recent years, allowing rapid development of models and optimization problems without requiring analytical gradients before training. These libraries allow taking the gradient of complex chains of operations by decomposing them into small, atomic operations then linking them together by copious usage of the chain rule. As an example, differentiation through traditional matrix-matrix multiplication or matrix addition is straightforward and efficient.

While dense matrix linear algebra routines are readily im-

plemented and available in these frameworks, their sparse counterparts have received less attention. In the case of PyTorch, Tensorflow, and Jax, all three have preliminary support for sparse matrices stored in the coordinate (COO) format. Yet, the use of sparse linear algebra is extensive in computational science, for example in the solution to partial differential equations (PDEs) (Briggs et al., 2000; Saad, 2003), circuit analysis (Bonfatti et al., 1973), graph neighborhood analysis (Saad, 2003; George et al., 1993), among many others, motivating the need for a more complete treatment of sparse operations. Moreover, there is growing demand for sparse kernels in the broader area of scientific machine learning (SciML), where machine learning techniques augment traditional scientific computing methods. The authors in (Greenfeld et al., 2019; Luz et al., 2020) learn methods to accelerate matrix linear solve routines, though develop various simplifications in the problems themselves due to the lack of sparse linear algebra routines. Similarly, the method presented in (Taghibakhshi et al., 2022) is constrained to training on small, dense matrices, then evaluated on arbitrarily sized inputs without any AD considerations. We thus seek an implementation of such sparse routines to allow machine learning researchers and computational scientists to develop scalable, learned methods.

The main contributions of this paper are structured as follows. We present preliminary background on AD internals and the sparse CSR matrix format in Section 2, then state the backwards-mode gradient derivations for several sparse operations and their implementation on both CPU and GPU (CUDA) processors in Section 3. To help motivate further possible uses of these sparse routines, we introduce several applications to optimization problems in Section 4, including graph neural networks (Wu et al., 2021; Kipf & Welling, 2017). Finally, we present performance results for our implementation in Section 5, underscoring the efficiency of the methods in the forward pass by comparing to existing sparse-matrix packages.

2. Background

Unless noted, $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{b}, \mathbf{x} \in \mathbb{R}^n$ represent the variables in the (sparse) system of linear equations

$$\mathbf{Ax} = \mathbf{b}. \quad (1)$$

¹Department of Computer Science, University of Illinois at Urbana-Champaign ²Department of Mechanical Science and Engineering, University of Illinois at Urbana-Champaign ³Scientific Computing Program, Memorial University of Newfoundland ⁴Department of Mathematics and Statistics, Memorial University of Newfoundland. Correspondence to: Nicolas Nytko <nnytko2@illinois.edu>.

For iterative solvers, \mathbf{x}_k denotes the approximate solution to Equation (1) at iteration k , along with residual $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$, and error $e_k = \mathbf{x} - \mathbf{x}_k$.

Given a sparse matrix \mathbf{M} , let mask (\mathbf{M}) denote the *sparsity mask* of \mathbf{M} as in

$$[\text{mask}(\mathbf{M})]_{ij} = \begin{cases} 1 & m_{ij} \neq 0 \\ 0 & \text{otherwise} \end{cases}. \quad (2)$$

Letting \odot denote the standard Hadamard product, $\mathbf{A} \odot \text{mask}(\mathbf{M})$ then represents (sparse) matrix \mathbf{A} masked to the sparsity of \mathbf{M} . The sparsity mask also leads to the identity $\mathbf{M} \odot \text{mask}(\mathbf{M}) = \mathbf{M}$, and, by convention, if \mathbf{M} is dense then $\text{mask}(\mathbf{M})$ is also dense.

2.1. Chain Rule

For background, we recall the multivariate chain rule (Johnson, 2017) in Theorem 2.1.

Theorem 2.1. *Given function $f : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ and vectors $\mathbf{t} \in \mathbb{R}^{n_t}$ and $\mathbf{x}(\mathbf{t}) \in \mathbb{R}^{n_x}$, the partial derivative $\frac{\partial z}{\partial t_i}$ for $z = f(\mathbf{x}(\mathbf{t}))$ is given by*

$$\frac{\partial z}{\partial t_i} = \sum_{j=1}^{n_x} \frac{\partial f}{\partial x_j} \frac{\partial x_j}{\partial t_i}. \quad (3)$$

For the general case, we can extend Theorem 2.1 for arbitrary tensor-valued functions, with tensors representing n -dimensional arrays (Kolda & Bader, 2009), as in Corollary 2.2.

Corollary 2.2. *Let I_1, I_2, \dots, I_N denote the indices of an N -mode tensor. Given function $f : \mathbb{R}^{I_1 \times I_2 \times \dots \times I_{N_X}} \rightarrow \mathbb{R}$, and tensors $\mathcal{T} \in \mathbb{R}^{J_1 \times J_2 \times \dots \times J_{N_T}}$ and $\mathcal{X}(\mathcal{T}) \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_{N_X}}$, the partial derivative $\frac{\partial z}{\partial t_{j_1, j_2, \dots, j_{N_T}}}$ for $z = f(\mathcal{X}(\mathcal{T}))$ is given by*

$$\frac{\partial z}{\partial t_{j_1, \dots, j_{N_T}}} = \sum_{i_1=1}^{I_1} \dots \sum_{i_{N_X}=1}^{I_{N_X}} \frac{\partial z}{\partial x_{i_1, \dots, i_{N_X}}} \frac{\partial x_{i_1, \dots, i_{N_X}}}{\partial t_{j_1, \dots, j_{N_T}}}. \quad (4)$$

From Corollary 2.2, we denote the term on the left side in the summation as the *generalized gradient*:

$$[\nabla_{\mathcal{X}}(z)]_{i_1, i_2, \dots, i_{N_X}} = \frac{\partial z}{\partial x_{i_1, i_2, \dots, i_{N_X}}}. \quad (5)$$

This represents “the gradient of z with respect to \mathcal{X} .” Moreover, the value attains the same shape (dimensionality and size in each dimension) as \mathcal{X} itself.

Likewise, from Corollary 2.2, we refer to the term on the right in the summation as the *generalized Jacobian*:

$$[J_{\mathcal{T}}(\mathcal{X})]_{(i_1, i_2, \dots, i_{N_X}), (j_1, j_2, \dots, j_{N_T})} = \frac{\partial x_{i_1, i_2, \dots, i_{N_X}}}{\partial t_{j_1, j_2, \dots, j_{N_T}}}, \quad (6)$$

where the parentheses in the indexing are used only to emphasize that the first N_X indices are used for the input, while the last N_T indices are used for the output. From this, we have that Equation (4) can be alternatively denoted as the tensor contraction $\nabla_{\mathcal{X}}(z)^T J_{\mathcal{X}}(\mathcal{T})$ over the indices i_1, i_2, \dots, i_{N_X} .

2.2. Reverse Mode Auto-Differentiation

Contemporary machine learning frameworks use *reverse-mode automatic differentiation* to compute gradient information, which is an efficient means of computing gradients of scalar-valued functions with respect to tensor-valued inputs (Baydin et al., 2015). The computation, in essence, is first done normally in a *forward pass*, with compositions of elementary functions being recorded into a *computation graph*. A second *backward pass* is then executed, tracing the computation graph backwards from the scalar output back to each input node and computing intermediate gradients at each step.

To help illustrate, we consider the function

$$f(\mathbf{x}, \mathbf{y}) = 2 \sin(\langle \mathbf{x}, \mathbf{y} \rangle), \quad (7)$$

for which we seek to evaluate the derivative with respect to the vectors \mathbf{x} and \mathbf{y} . At each node of the computation

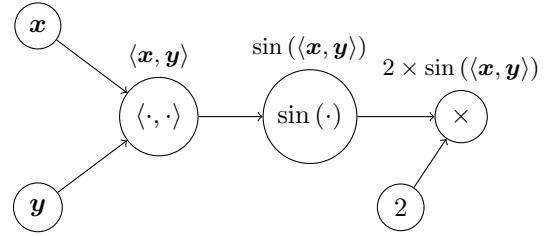


Figure 1. Example computation graph corresponding to evaluating Equation (7).

graph (see Figure 1), the gradient of the output is calculated with respect to the input of that particular node; this is then passed further back in the graph. Formally, letting z be the scalar output of our computation and $\mathcal{Y} = f_i(\mathcal{X})$ be the computation done at node i , we find the intermediate gradient with the tensor contraction

$$\nabla_{\mathcal{X}}(z) = \nabla_{\mathcal{Y}}(z)^T J_{f_i}(\mathcal{X})(\mathcal{X}). \quad (8)$$

This operation is referred to as the *vector-Jacobian product* (VJP) in automatic differentiation. In the case that functions are vector valued ($\mathbb{R}^n \rightarrow \mathbb{R}^m$) this simplifies to vector-matrix products.

2.3. Compressed Sparse Row format

The *compressed sparse row* format is a row-oriented representation for storage of sparse matrices (Saad, 2003). In

this representation, three arrays are needed to determine a single sparse matrix, $\mathbf{A} \in \mathbb{R}^{n \times m}$. Denoting $\text{nnz}(\mathbf{A})$ as the number of nonzero entries in \mathbf{A} , the three arrays become:

- `data` of size $\text{nnz}(\mathbf{A})$ storing each nonzero data entry sorted in row then column order,
- `indices` of size $\text{nnz}(\mathbf{A})$ storing the column number of each respective nonzero entry, and
- `indptr` of size $m + 1$ that stores the *offset* into `data` of the starting point of each row.

By convention, the last entry of `indptr` stores the number of nonzeros in \mathbf{A} .

$$\mathbf{C} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 3 & 2 & 0 \\ 0 & 0 & 0 & 0 \\ 7 & 0 & 0 & 5 \end{bmatrix}$$

$$\begin{aligned} \text{data}(\mathbf{C}) &= [1 \ 3 \ 2 \ 7 \ 5] \\ \text{indices}(\mathbf{C}) &= [0 \ 1 \ 2 \ 0 \ 3] \\ \text{indptr}(\mathbf{C}) &= [0 \ 1 \ 3 \ 3 \ 5] \end{aligned}$$

Figure 2. An example sparse matrix and its CSR representation.

We use the CSR representation in this work because it provides good memory access to the matrix rows, especially useful in products where the matrix is left-multiplying, for example in the sparse matrix-vector product $\mathbf{A}\mathbf{x}$ or sparse-times-dense matrix product $\mathbf{A}\mathbf{D}$, where \mathbf{A} is sparse and \mathbf{x} , \mathbf{D} are a dense vector and matrix.

3. Sparse Kernels

For this work, we implement optimized forward and backward passes (VJP implementations) for various sparse operations with an underlying CSR data structure. We propagate *sparse* gradients whereby functions that take sparse inputs will have gradients whose sparsity mask will match those of the inputs. This trade-off leads to much smaller memory and computational costs (as observed indirectly in Table 2) as, during optimization, only nonzero entries in the inputs will receive gradients. For example, optimizing a problem with a sparse-matrix input and scalar output will only lead to optimization of the nonzero entries of the input.

We next detail the forward implementations and backwards (VJP) products of several sparse operations. For derivations of the matrix reverse-mode updates, we point the reader to (Giles, 2008) and remark that with minor consideration of the sparsity patterns, the results extend to sparse matrices. In the following, not all routines have highly tuned implementations for both CPU and CUDA processors.

Table 1. Definitions of vector-Jacobian products for the different sparse operations. The vector in the VJP (denoted by \mathbf{v} or \mathbf{V}) is the intermediate gradient with respect to the output of the operation when running backpropagation.

METHOD	OPERATION	WRT	VJP
SPMV	$\mathbf{A}\mathbf{x}$	\mathbf{A}	$\mathbf{v}\mathbf{x}^T \odot \text{mask}(\mathbf{A})$
SPMV	$\mathbf{A}\mathbf{x}$	\mathbf{x}	$\mathbf{A}^T\mathbf{v}$
SPSPMM	$\mathbf{A}\mathbf{B}$	\mathbf{A}	$(\mathbf{V}\mathbf{B}^T) \odot \text{mask}(\mathbf{A})$
SPSPMM	$\mathbf{A}\mathbf{B}$	\mathbf{B}	$(\mathbf{A}^T\mathbf{V}) \odot \text{mask}(\mathbf{B})$
SPDMM	$\mathbf{A}\mathbf{B}$	\mathbf{A}	$(\mathbf{V}\mathbf{B}^T) \odot \text{mask}(\mathbf{A})$
SPDMM	$\mathbf{A}\mathbf{B}$	\mathbf{B}	$\mathbf{A}^T\mathbf{V}$
SP + SP	$\alpha\mathbf{A} + \beta\mathbf{B}$	\mathbf{A}	$\alpha\mathbf{V} \odot \text{mask}(\mathbf{A})$
SP + SP	$\alpha\mathbf{A} + \beta\mathbf{B}$	\mathbf{B}	$\beta\mathbf{V} \odot \text{mask}(\mathbf{B})$
SPSOLVE	$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$	\mathbf{A}	$-\mathbf{A}^{-T}\mathbf{v}\mathbf{x}^T \odot \text{mask}(\mathbf{A})$
SPSOLVE	$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$	\mathbf{b}	$\mathbf{A}^{-T}\mathbf{v}$

3.1. Sparse Matrix-Vector (SpMV)

The sparse matrix-vector product computes

$$\mathbf{y} \leftarrow \mathbf{A}\mathbf{x} \quad (9)$$

for sparse $\mathbf{A} \in \mathbb{R}^{n \times m}$, dense $\mathbf{x} \in \mathbb{R}^m$, and outputs $\mathbf{y} \in \mathbb{R}^n$. The CSR representation allows for direct evaluation of inner products of each row of \mathbf{A} and \mathbf{x} , and the inner products can be executed in parallel on a GPU.

For the backwards pass, we compute $\mathbf{A}^T\mathbf{v}$ (resp. $\mathbf{v}\mathbf{x}^T$) for the gradient with respect to (WRT) \mathbf{x} (resp. \mathbf{A}), for intermediate gradient value \mathbf{v} . Computing the transposed matrix-vector product does not directly align with the sparse structure, therefore we opt to use an atomic reduction-based algorithm to compute $\mathbf{A}^T\mathbf{v}$ as outlined in (Tao et al., 2014). In this form, we compute inner products of the rows of \mathbf{A} and \mathbf{v} and atomically reduce the output into their correct entries. Computing $\mathbf{v}\mathbf{x}^T$ is easily parallelized, since it is the outer product between two dense vectors that is masked to a sparse matrix, requiring computation of only nonzero entries of \mathbf{A} .

3.2. Sparse-Sparse Matrix Multiply (SpSpMM)

For the sparse-sparse matrix multiply primitive, we compute

$$\mathbf{C} \leftarrow \mathbf{A}\mathbf{B}, \quad (10)$$

where $\mathbf{A} \in \mathbb{R}^{i \times j}$, $\mathbf{B} \in \mathbb{R}^{j \times k}$, and $\mathbf{C} \in \mathbb{R}^{i \times k}$, with all three matrices are stored in CSR format. In the forward pass, we follow the parallel SpGEMM algorithm from (Dalton et al., 2015), where we compute intermediate products of each row of \mathbf{A} and the entirety of \mathbf{B} , then reduce and collect redundant entries to form \mathbf{C} .

For the backward pass, we compute $(\mathbf{V}\mathbf{B}^T) \odot \text{mask}(\mathbf{A})$

and $(A^T V) \odot \text{mask}(B)$. Noting that the product VB^T is the inner product between rows of both V and B , our data accesses aligns directly with the CSR structure. Computing $A^T V$, however, accesses columns of both A and V ; we thus follow a modified version of the SpGEMM algorithm from above that operates on columns of A instead.

3.3. Sparse-Dense Matrix Multiply (SpDMM)

For the SpDMM routine, we compute

$$C \leftarrow AB \quad (11)$$

as in the SpSpMM case except for B and C now being dense. In the forward pass, we parallelize with each fine-grained operation focusing on an entry of C . That is, we compute the respective inner product between a row of A and a column of B , which because of its dense structure does not incur any significant penalties for column accesses.

On the backward pass, we compute $(VB^T) \odot \text{mask}(A)$ (as before) and $A^T V$. For the latter, we take the sparse transposition of A and re-execute the forward routine to compute the gradient.

3.4. Sparse + Sparse

In a sparse add, we compute the linear combination of two sparse matrices as in

$$C \leftarrow \alpha A + \beta B, \quad (12)$$

which we write in a general form so that both $A + B$ and $A - B$ can be computed with the same method. A key observation is that (with slight abuse of notation),

$$\text{mask}(C) = \text{mask}(A) \cup \text{mask}(B), \quad (13)$$

meaning the computation of C can be viewed as a union over the rows of A and B . Moreover, this form can be implemented in parallel over each row.

To compute the backwards pass, we consider the gradient with respect to A , B as $\alpha V \odot \text{mask}(A)$, $\beta V \odot \text{mask}(B)$, respectively. Each can be found as the row-wise reduction from V to the sparsity mask of A or B . Because $\text{mask}(A) \subseteq \text{mask}(V)$, $\text{mask}(B) \subseteq \text{mask}(V)$, we need only to compute matching nonzero entries in both matrices. Again, this can be implemented in parallel over the rows.

3.5. Sparse Triangular Solve

A sparse triangular solve is an operation to compute the value of x in the matrix equation

$$Lx = b, \quad (14)$$

where L has a lower triangular form, i.e. $l_{ij} \neq 0$ if $i \geq j$. Without loss of generality, we can also consider upper-triangular systems U using matrix flip operations to convert $Ux = b$ into an equivalent lower-triangular system. Such a system has a (relatively) simple routine for computing the linear solve: each row depends on the intermediate values of previous rows only, and no intermediate preprocessing is needed.

For the forward pass, we elect to use the synchronization-free GPU triangular solve detailed in (Su et al., 2020) to exploit the limited parallelism that may exist in computing x .

On the backward pass, we seek $L^{-T}v$ and $-L^{-T}vx^T \odot \text{mask}(L)$ for gradients with respect to b and L , respectively. We first find $L^{-T}v$ with our existing forward triangular solve routine (except flipped to run on upper triangular systems). Then we observe that the gradient with respect to L contains the b gradient term as a masked outer product, so we can re-use the result. The masked outer-product can be executed in parallel over the nonzero entries of L .

4. Applications of Sparse Kernels

To demonstrate the use of the sparse kernels and differentiability above, we present several optimization problems that exercise sparse-matrix computations. As an example sparse matrix, we define A as

$$A = \begin{bmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & & -1 & 2 \end{bmatrix}, \quad (15)$$

which is a standard test problem in computational science and comes from the finite-difference discretization of the 1D Poisson problem (Saad, 2003; Quarteroni et al., 2006),

$$-\nabla^2 u = f \quad \text{in } \Omega, \quad (16)$$

$$u = 0 \quad \text{on } \partial\Omega, \quad (17)$$

on the domain $\Omega = (0, (N+1)^2) \subseteq \mathbb{R}$, with N being the number of interior gridpoints in the discretization — the domain is selected to cancel any constant scaling of the matrix. This results in a constant number of nonzero entries per row (3) away from the boundary.

4.1. Entry-wise Jacobi Relaxation

The Jacobi method is an iterative solver for linear systems whereby the inverse matrix is approximated by inverting the diagonal entries only (Saad, 2003). Consider solving

$$Ax = b, \quad (18)$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a sparse matrix with diagonal D . The Jacobi method is expressed as

$$\mathbf{x}^{(k+1)} = D^{-1} \left(\mathbf{b} - (\mathbf{A} - D) \mathbf{x}^{(k)} \right), \quad (19)$$

though this can often see slow convergence as the iteration tends to largely dampen *high-frequency error modes* only when \mathbf{A} comes from discretization of an elliptic differential equation (Briggs et al., 2000).

Convergence can be accelerated by separately weighting the previous value of \mathbf{x} and the approximate inverse at each iteration; here, we consider the case where this weighting is determined independently for each entry of \mathbf{x} . Denoting $\mathbf{\Omega}$ as the diagonal scaling matrix of the entries of $\mathbf{x}^{(k)}$, we then rewrite the Jacobi iteration as

$$\mathbf{x}^{(k+1)} = \mathbf{\Omega} D^{-1} \mathbf{b} + (\mathbf{I} - \mathbf{\Omega} D^{-1} \mathbf{A}) \mathbf{x}^{(k)}. \quad (20)$$

To simplify the notation, we write $j(\mathbf{x}_0, \mathbf{b}; \mathbf{\Omega})$ as the function that applies one iteration of the weighted Jacobi scheme to \mathbf{x}_0 in order to solve $\mathbf{A}\mathbf{x} = \mathbf{b}$.

We generate a finite-difference matrix \mathbf{A} from Equation (15) of size $n = 16$, then optimize the entries of $\mathbf{\Omega}$ by minimizing

$$\ell = \sum_{k=1}^K (j(\mathbf{x}_k, \mathbf{0}; \mathbf{\Omega}))^T \mathbf{A} j(\mathbf{x}_k, \mathbf{0}; \mathbf{\Omega}), \quad (21)$$

where $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K\}$ are a set of test vectors in \mathbb{R}^n with random unit-normally distributed entries. This optimizes the one-step Jacobi error when solving for a $\mathbf{0}$ right-hand side.

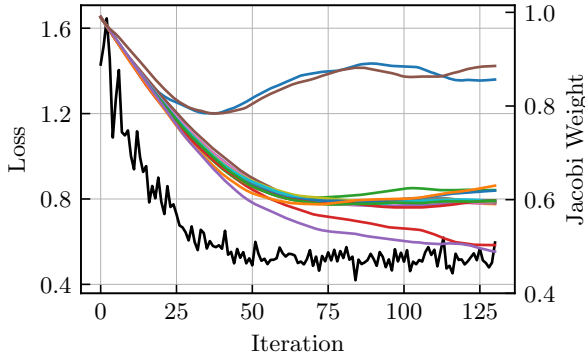


Figure 3. Loss history (in black) obtained from optimizing over the Jacobi relaxation weights. The colored lines are the history of the node-wise weights for each training iteration.

Training loss history and final values of the weights are shown in Figures 3 and 4. In Figure 4, we observe that the nodal weights at the two ends of the domain are maximized; for our problem setup this mimics the behavior of values being propagated from the boundaries inwards.

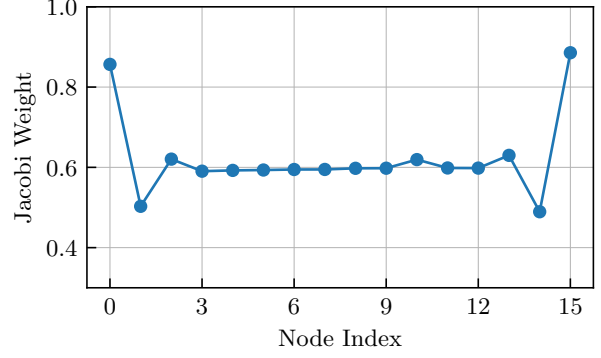


Figure 4. Final node-wise weight values for the Jacobi relaxation.

4.2. Heavyball Iteration

Polyak’s heavyball iteration (Polyak, 1964) is a two-step iterative method that uses the information from the last two iterates to generate the next approximation. It can be viewed as a gradient descent with a momentum term.

To minimize a differentiable function f , we consider the update step

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \nabla_{\mathbf{x}} f(\mathbf{x}^{(k)}) + \beta (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}), \quad (22)$$

for scalars $\alpha, \beta \in \mathbb{R}$. We then apply Equation (22) to solve systems of the form Equation (18) by defining \mathbf{A} to be

$$f(\mathbf{x}; \mathbf{A}, \mathbf{b}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}. \quad (23)$$

Next, taking the gradient of f with respect to \mathbf{x} yields

$$\nabla_{\mathbf{x}} f(\mathbf{x}; \mathbf{A}, \mathbf{b}) = \mathbf{A} \mathbf{x} - \mathbf{b}, \quad (24)$$

where we have that $\nabla_{\mathbf{x}} f = \mathbf{0}$ implies a solution to the matrix system if \mathbf{A} is SPD.

Similarly, we consider the example problem in Section 4.1 by generating a 16×16 matrix \mathbf{A} and minimizing

$$\ell = \sum_{k=1}^K (\mathbf{h}_{12}(\mathbf{x}_k, \mathbf{0}; \alpha, \beta))^T \mathbf{A} \mathbf{h}_{12}(\mathbf{x}_k, \mathbf{0}; \alpha, \beta), \quad (25)$$

where $\mathbf{h}_L(\mathbf{x}_k, \mathbf{b}; \alpha, \beta)$ is the application of L rounds of heavyball iteration to \mathbf{x}_k with right-hand-side \mathbf{b} and parameters α, β . Again, the $\{\mathbf{x}_k\}$ is a set of K vectors in \mathbb{R}^N with unit-normally distributed entries.

We optimize the error after $\frac{3}{4}N$ iterations of the heavyball method. The training and testing loss history is shown in Figure 5, and the parameter history as a function of training epoch is shown in Figure 6.

4.3. Conjugate Gradient

The conjugate gradient (CG) method is another iterative solver for matrix equations of the form in Equation (18).

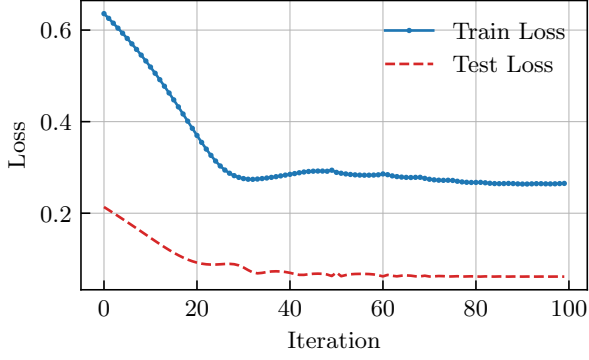


Figure 5. Loss history obtained from optimizing over the heavyball parameters.

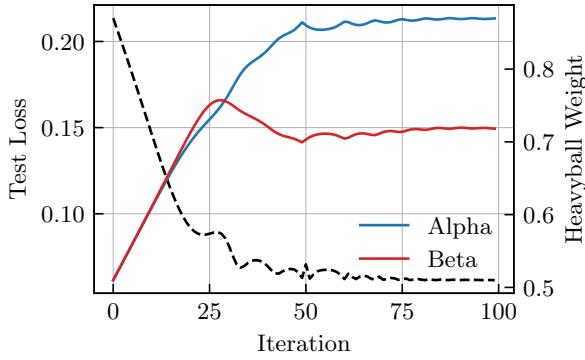


Figure 6. Parameter history for each heavyball training iteration, plotted along with the training loss.

While it typically exhibits faster convergence than both the Jacobi and heavyball methods, convergence can be accelerated further by passing an approximate inverse to \mathbf{A} , denoted by $\mathbf{M} \approx \mathbf{A}^{-1}$, with the assumption that \mathbf{M} is itself SPD, using the preconditioned conjugate gradient (PCG) method (Saad, 2003).

We extend the one-dimensional problem in Equation (15) to two dimensions via tensor product over the x and y dimensions, as in

$$\mathbf{A} = \mathbf{A}_{N_x \times N_y} = (\mathbf{A}_{N_x} \otimes \mathbf{I}_{N_y}) + (\mathbf{I}_{N_x} \otimes \mathbf{A}_{N_y}), \quad (26)$$

where \otimes denotes the standard matrix Kronecker product, and \mathbf{A}_N indicates a one-dimensional discretization with N interior grid points. For our test problem we use $N_x = N_y = 8$, resulting in a matrix shape $\mathbf{A} \in \mathbb{R}^{64 \times 64}$.

The goal is the construction of \mathbf{M} that suitably approximates the inverse to \mathbf{A} . Since, \mathbf{M} is an SPD matrix, it has a decomposition,

$$\mathbf{M} = \mathbf{L}\mathbf{L}^T. \quad (27)$$

From this, we optimize over the entries of \mathbf{L} to force the preconditioner \mathbf{M} to be symmetric, while \mathbf{M} is not directly constrained to be positive-definite.

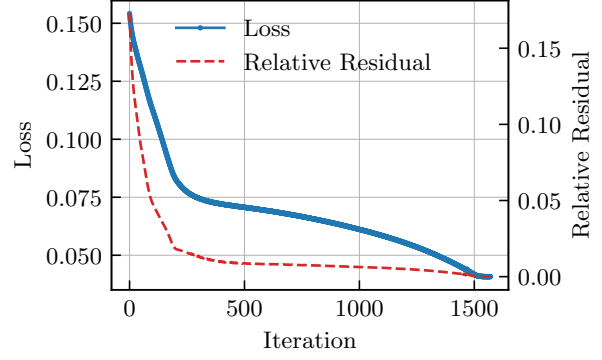


Figure 7. Loss history in constructing \mathbf{M} , the CG preconditioner. The blue line (left axis) is the loss and the red dashed line (right axis) is the relative residual.

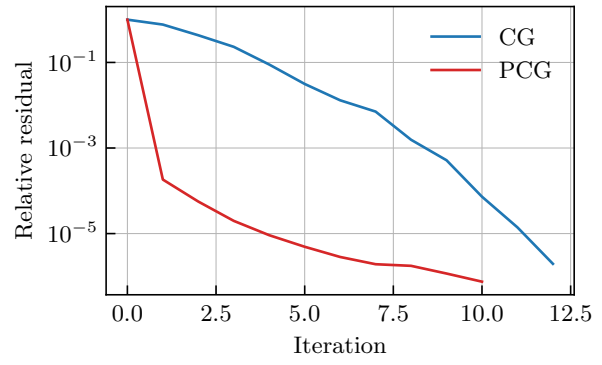


Figure 8. Relative residual history of a regular (non-preconditioned) CG and the preconditioner found by optimization.

To find the entries in \mathbf{M} , we introduce the weighted loss

$$\ell = \sum_{i=1}^{N_{it}} \left(\frac{\gamma^{N_{it}-i}}{\sum_{j=1}^{N_{it}} \gamma^{N_{it}-j}} \right) \frac{\|\mathbf{r}^{(i)}\|_2}{\|\mathbf{b}\|_2}, \quad (28)$$

where N_{it} is the number of PCG iterations run (we use $N_{it} = 4$), $\mathbf{r}^{(i)}$ is the i th residual of the iteration, \mathbf{b} is the right-hand-side, and $\gamma \in (0, 1]$ is a scaling constant (we use $\gamma = 0.6$). With this form, we minimize the overall weighted sum of the residual history with later iterates being weighted more than earlier ones. Experimentally, we observe improved convergence of the optimization problem in comparison to minimizing only over the last iterate. This also avoids the problem of vanishing gradients as the number of iterations is increased, as gradient information is used from each intermediate iterate.

We run the optimization over a lower bidiagonal \mathbf{L} , which gives a tridiagonal $\mathbf{M} = \mathbf{L}\mathbf{L}^T$. The loss and final relative residual obtained during each training iteration are displayed in Figure 7, while the residual history with and without \mathbf{M} is shown in Figure 8.

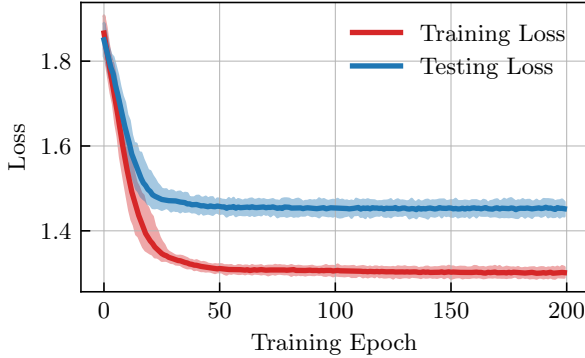


Figure 9. Loss history over training the GNN over the CiteSeer dataset. Averaged over 100 runs, lines denote mean while shaded regions are two standard deviations from mean.

4.4. Graph Neural Networks

The GCN layer (Kipf & Welling, 2017) is a graph convolutional layer that *convolves* node features on a graph, $G(\mathbf{A})$. Denoting the (weighted) adjacency matrix of the graph by $\mathbf{A} \in \mathbb{R}^{n \times n}$, and denoting $\mathbf{X}^{(i)} \in \mathbb{R}^{n \times C}$ as the C -dimensional node features at layer i , we represent the convolution of the node features as

$$\mathbf{X}^{(i+1)} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X}^{(i)} \Theta^{(i+1)}, \quad (29)$$

where $\Theta \in \mathbb{R}^{C \times F}$ is the weight matrix for the layer, $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$, and $\tilde{\mathbf{D}}$ is the diagonal matrix extracted from $\tilde{\mathbf{A}}$.

We reimplement the GCN layer using our optimized underlying sparse matrix operations. Depending on the order of operations, computing the output in Equation (29) can be seen as a series of sparse-dense matrix multiplies (multiplying from right-to-left), or as two sparse-sparse multiplies followed by dense multiplies. In the following, we employ the former method.

Using this GCN layer, we reimplement the semi-supervised CiteSeer example from (Kipf & Welling, 2017). This dataset has 3327 nodes, 4732 edges, and 6 classes. We train a two-layer GCN with ReLU and sigmoid activations after the first and second layers and use a cross-entropy loss to optimize the predicted labels on each node, which follows the same setup as in (Kipf & Welling, 2017). Likewise, between each GCN layer is a dropout layer with $p = 0.5$, and we train with an Adam optimizer using a learning rate of 0.01, L2 regularization of 5×10^{-4} , and 16-dimensional representations for the hidden node features. This is trained for 200 epochs over 100 random initializations and training history can be seen in Figure 9.

As seen in Figure 10, we are able to match the roughly 70% classification accuracy on the CiteSeer dataset.

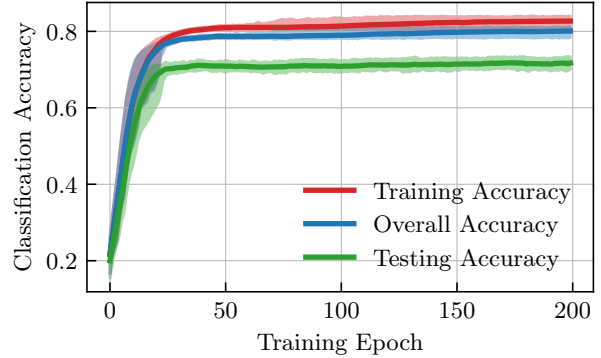


Figure 10. Classification accuracy of the networks per training epoch. Averaged over 100 runs, lines denote mean while shaded regions are two standard deviations from mean.

5. Timing Results

We present timing results for our implementation of the sparse kernels outlined in Section 3 along with a comparison of CPU and GPU (CUDA) timings in Table 2. For comparison, we have also included the respective *dense* operation in applicable cases, such as the dense matrix-vector (DMV) or dense-dense matrix-matrix product (DDMM). These dense implementations use the PyTorch built-in matrix routines. We also compare with NumPy/SciPy CPU and CuPy CUDA optimized versions for the forward passes.

Figure 11 highlights scaling results for the sparse matrix-vector product, sparse-sparse matrix product, and sparse-dense matrix product. These show the running time as a factor of the number of nonzeros in the matrix given in Equation (15), and indicate that forward and backward passes for the SpMV and SpSpMM are both linear in the nonzeros, while the forward and backward passes for the SpDMM are quadratic in the number of nonzeros. These timing results are run on a single compute node of the Narval cluster with two AMD Milan 7413 processors and four NVIDIA A100 GPUs, though for our tests we use only one CPU and one GPU. Our CPU-based implementations are single threaded only, so for fair comparison we limited the amount of threads in the dense operations to one.

6. Conclusions

In this work, we have described the implementation of a framework for automatically computing the gradient through computations with sparse matrix operations. We have stated the backwards propagation update rules and described at a high level how one may efficiently implement them on CPUs and GPUs. We demonstrated a range of applications that can be optimized using these sparse matrix primitives. Finally, we show timing and scaling results to show that our implementations are both scalable with re-

Optimized Sparse Matrix Operations for Reverse Mode Automatic Differentiation

Table 2. Timing results for CPU- and CUDA-based implementations. *Iterations* refers to the number of times the specific test was run (times listed are total wall-clock time elapsed); *N* refers to the problem size used, most usually the size of the matrix in one dimension. *GPU Speedup* is computed as CPU time / GPU time. Methods with *Sp* refer to *S*parse operations (ours), while *D* refers to *D*ense (PyTorch).

Test Name	ours CPU (s)	SciPy CPU (s)	ours GPU (s)	CuPy GPU (s)	GPU Speedup	Iter.	<i>N</i>
SpMV Forward	0.141	0.092	0.045	0.049	3.150	1000	32,768
SpMV Backward	0.595		0.198		3.003	1000	32,768
DMV Forward	143.441	143.614	2.985	2.983	48.058	1000	32,768
DMV Backward	1509.088		21.806		69.206	1000	32,768
SpSpMM Forward	0.184	0.065	0.048	0.032	3.798	100	16,384
SpSpMM Backward	2.150		0.293		7.336	100	16,384
SpDMM Forward	145.156	100.840	1.100	0.569	132.014	100	16,384
SpDMM Backward	303.903		4.923		61.736	100	16,384
DDMM Forward	190.113	202.030	0.985	0.921	193.097	2	16,384
DDMM Backward	597.621		2.846		210.020	2	16,384
Sp + Sp Forward	0.073	0.034	0.011	0.011	6.813	100	32,768
Sp + Sp Backward	0.133		0.025		5.227	100	32,768
D + D Forward	60.978	61.553	0.941	0.947	64.783	100	32,768
D + D Backward	123.519		2.596		47.580	100	32,768
SpTRSV Forward	0.030	32.746	4.710	0.767	0.006	100	32,768
SpTRSV Backward	13.462		10.997		1.224	100	32,768
DTRSV Forward	35.092	65.834	0.695	4.119	50.523	100	32,768
DTRSV Backward	863.779		3.499		246.855	100	32,768

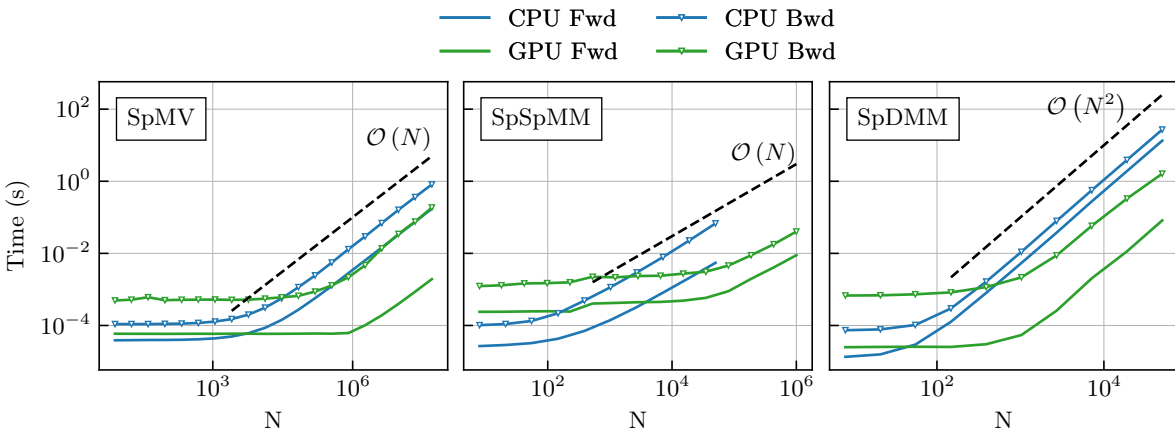


Figure 11. Scaling results for the forward and backward passes on several key sparse operations. The CPU running times are in blue, while the CUDA runtimes are in green. Triangle markers denote the backwards pass routines.

Table 3. Software versions used for testing.

Library	Backend	Version
FlexiBLAS	BLIS	0.8.1
Python	CPython	3.10.2
SciPy		1.8.0
NumPy		1.22.2
CUDA		11.7
CuPy		11.2.0
PyTorch		1.13.0

spect to dense linear algebra routines, and competitive in the forward pass with other sparse linear algebra packages.

There are a numerous directions of future work that we envision from this paper. Perhaps the most straightforward is the implementation of more complex sparse linear algebra operations such as eigensolver routines, linear least squares, etc. In an extension to the solve routine we have presented, one can consider backpropagation through the L and U factors themselves, perhaps in optimization problems where such calculations are necessary. Our CUDA LU factorization itself could also be optimized, there are recent works such as (Gaihre et al., 2022) that claim the symbolic factorization can be implemented to run in a scalable, parallel fashion on the GPU.

Acknowledgments

This research was enabled in part by support provided by ACENET (www.ace-net.ca), and the Digital Research Alliance of Canada (alliancecan.ca). The work of SM was partially supported by an NSERC Discovery Grant.

A. Software Versions

Our sparse kernels are compiled as a PyTorch extension written in C++17 using GCC, with the `-O2` flag enabled as the only non-default option. The list of library versions we use and compare against are detailed in Table 3.

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D. G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., and Zheng, X. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pp. 265–283, 2016. URL <https://www.usenix.org/system/files/conference/osdi16/osdi16-abadi.pdf>.
- Baydin, A. G., Pearlmutter, B. A., and Radul, A. A. Automatic differentiation in machine learning: a survey. *CoRR*, abs/1502.05767, 2015. URL <http://arxiv.org/abs/1502.05767>.
- Bonfatti, F., Monaco, V., and Tiberio, P. Microwave circuit analysis by sparse matrix techniques. In *1973 IEEE G-MTT International Microwave Symposium*, pp. 41–43, 1973. doi: 10.1109/GMTT.1973.1123084.
- Bradbury, J., Frostig, R., Hawkins, P., Johnson, M. J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., and Zhang, Q. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Briggs, W. L., Henson, V. E., and McCormick, S. F. *A Multigrid Tutorial, Second Edition*. Society for Industrial and Applied Mathematics, second edition, 2000. doi: 10.1137/1.9780898719505. URL <https://epubs.siam.org/doi/abs/10.1137/1.9780898719505>.
- Dalton, S., Olson, L., and Bell, N. Optimizing sparse matrix—matrix multiplication for the gpu. *ACM Trans. Math. Softw.*, 41(4), oct 2015. ISSN 0098-3500. doi: 10.1145/2699470. URL <https://doi.org/10.1145/2699470>.
- Gaihre, A., Li, X. S., and Liu, H. Gsofa: Scalable sparse symbolic lu factorization on gpus. *IEEE Transactions on Parallel and Distributed Systems*, 33:1015–1026, 4 2022. ISSN 15582183. doi: 10.1109/TPDS.2021.3090316.
- George, A., Gilbert, J. R., and Liu, J. W. H. (eds.). *Graph Theory and Sparse Matrix Computation*. Springer New York, 1993. doi: 10.1007/978-1-4613-8369-7. URL <https://doi.org/10.1007/978-1-4613-8369-7>.
- Giles, M. B. Collected matrix derivative results for forward and reverse mode algorithmic differentiation. *Lecture Notes in Computational Science and Engineering*, 64 LNCSE:35–44, 2008. ISSN 14397358. doi: 10.1007/978-3-540-68942-3_4/COVER. URL https://link.springer.com/chapter/10.1007/978-3-540-68942-3_4.
- Greenfeld, D., Galun, M., Kimmel, R., Yavneh, I., and Basri, R. Learning to optimize multigrid pde solvers, 2019.
- Johnson, J. Derivatives, backpropagation, and vectorization. <http://cs231n.stanford.edu/handouts/derivatives.pdf>, September 2017.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*,

2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.
- Kolda, T. G. and Bader, B. W. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009. doi: 10.1137/07070111X. URL <https://doi.org/10.1137/07070111X>.
- Luz, I., Galun, M., Maron, H., Basri, R., and Yavneh, I. Learning algebraic multigrid using graph neural networks, 2020.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E. Z., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. *CoRR*, abs/1912.01703, 2019. URL <http://arxiv.org/abs/1912.01703>.
- Polyak, B. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964. ISSN 0041-5553. doi: [https://doi.org/10.1016/0041-5553\(64\)90137-5](https://doi.org/10.1016/0041-5553(64)90137-5). URL <https://www.sciencedirect.com/science/article/pii/0041555364901375>.
- Quarteroni, A., Sacco, R., and Saleri, F. *Numerical Mathematics*. Texts in Applied Mathematics. Springer, Berlin, Germany, 2 edition, October 2006.
- Saad, Y. *Iterative Methods for Sparse Linear Systems*. Other Titles in Applied Mathematics. SIAM, second edition, 2003. ISBN 978-0-89871-534-7. doi: 10.1137/1.9780898718003. URL http://www-users.cs.umn.edu/~saad/IterMethBook_2ndEd.pdf.
- Su, J., Zhang, F., Liu, W., He, B., Wu, R., Du, X., and Wang, R. Capellinisptrsv: A thread-level synchronization-free sparse triangular solve on gpus. pp. 1–11. ACM, 8 2020. ISBN 9781450388160. doi: 10.1145/3404397.3404400. URL <https://dl.acm.org/doi/10.1145/3404397.3404400>.
- Taghibakhshi, A., Nytko, N., Zaman, T., MacLachlan, S., Olson, L., and West, M. Learning interface conditions in domain decomposition solvers. 5 2022. URL <http://arxiv.org/abs/2205.09833>.
- Tao, Y., Deng, Y., Mu, S., Zhu, M., Xiao, L., Ruan, L., and Huang, Z. Atomic reduction based sparse matrix-transpose vector multiplication on gpus. *Proceedings of the International Conference on Parallel and Distributed Systems - ICPADS*, 2015-April:987–992, 2014. ISSN 15219097. doi: 10.1109/PADSW.2014.7097920.
- Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., and Yu, P. S. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32:4–24, 1 2021. ISSN 2162-237X. doi: 10.1109/TNNLS.2020.2978386. URL <https://ieeexplore.ieee.org/document/9046288/>.