

torch-sla: Differentiable Sparse Linear Algebra with Adjoint Solvers and Sparse Tensor Parallelism for PyTorch

Mingyuan Chi

walker.chi.000@gmail.com

<https://github.com/walkerchi/torch-sla>

January 15, 2026

Abstract

We present `torch-sla`, an open-source PyTorch library for differentiable sparse linear algebra. The library addresses two fundamental challenges: (1) **Adjoint-based differentiation** for linear and nonlinear sparse solvers, achieving $\mathcal{O}(1)$ memory complexity for computational graphs regardless of solver iterations; and (2) **Sparse Tensor Parallel** computing via domain decomposition with halo exchange, enabling distributed sparse operations across multiple GPUs. `torch-sla` supports multiple backends (SciPy, cuDSS, PyTorch-native) and scales to 169 million degrees of freedom. Benchmarks demonstrate $12\times$ speedup on multi-GPU configurations and correct gradient computation verified against finite differences.

1 Introduction

Sparse linear systems $\mathbf{Ax} = \mathbf{b}$ are fundamental to scientific computing. They arise in finite element analysis [Hughes, 2012], graph neural networks [Kipf and Welling, 2017], physics-informed machine learning [Raissi et al., 2019], and computational fluid dynamics [Jasak et al., 2007]. With the rise of differentiable programming, there is growing demand for sparse solvers that integrate with automatic differentiation frameworks like PyTorch [Paszke et al., 2019].

The challenge of differentiating through solvers. Consider a loss function $\mathcal{L}(\mathbf{x})$ where $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is obtained by solving a linear system. To train neural networks that interact with such solvers, we need gradients $\partial\mathcal{L}/\partial\mathbf{A}$ and $\partial\mathcal{L}/\partial\mathbf{b}$. A naive approach—differentiating through each iteration of an iterative solver—creates $\mathcal{O}(k)$ nodes in the computational graph, where k may be thousands of iterations. This leads to memory explosion and slow backward passes.

The challenge of scale. Industrial problems often exceed single-GPU memory. A 3D finite element mesh with 10 million nodes produces a sparse matrix requiring tens of gigabytes. Distributed computing is essential, but sparse matrix operations require careful communication patterns (halo exchange) that are non-trivial to implement correctly with gradient support.

This paper introduces `torch-sla`, addressing both challenges with two key innovations:

- **Adjoint-based differentiation** (§3.2): We implement the adjoint method for linear solves, eigenvalue problems, and nonlinear systems. This achieves $\mathcal{O}(1)$ computational graph nodes and $\mathcal{O}(\text{nnz})$ memory, independent of solver iterations.
- **Sparse tensor parallelism** (§3.3): We implement domain decomposition with automatic halo exchange following industrial CFD/FEM practices. This enables multi-GPU computing with gradient support.

2 Related Work

Differentiable linear algebra. JAX [Bradbury et al., 2018] provides differentiable sparse CG, but naive differentiation through iterations creates $\mathcal{O}(k)$ graph nodes. Blondel et al. [2022] formalize implicit differentiation for optimization layers. OptNet [Amos and Kolter, 2017] and cvxpylayers [Agrawal et al., 2019] handle convex optimization but not general sparse systems.

Sparse solvers. SciPy [Virtanen et al., 2020] provides SuperLU/UMFPACK for CPU. NVIDIA cuDSS [NVIDIA Corporation, 2024] provides GPU direct solvers. None natively support PyTorch autograd.

Distributed sparse computing. PETSc [Balay et al., 2019], Trilinos [Heroux et al., 2005], and hypre [Falgout and Yang, 2002] implement distributed sparse linear algebra. OpenFOAM [Jasak et al., 2007] uses domain decomposition for CFD. We bring these patterns to PyTorch with automatic differentiation.

3 Methodology

We first introduce the background on sparse linear systems (§3.1), then present our adjoint-based differentiation approach (§3.2), and finally describe our distributed computing implementation (§3.3).

3.1 Preliminaries: Sparse Linear Systems

A sparse linear system $\mathbf{Ax} = \mathbf{b}$ has a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with $\text{nnz} \ll n^2$ non-zero entries. We store matrices in coordinate (COO) format: three arrays for row indices, column indices, and values.

Direct solvers (LU, Cholesky factorization) compute exact solutions but require $\mathcal{O}(n^{1.5})$ memory for 2D problems due to fill-in—new non-zeros created during factorization [George, 1973]. For a 2D Poisson problem with n unknowns, the factored matrix has $\mathcal{O}(n \log n)$ non-zeros instead of the original $\mathcal{O}(n)$.

Iterative solvers (Conjugate Gradient, BiCGStab, GMRES) maintain $\mathcal{O}(\text{nnz})$ memory by never forming the factorization. However, they require $\mathcal{O}(\sqrt{\kappa})$ iterations for CG where κ is the condition number [Saad, 2003]. Each iteration performs one sparse matrix-vector multiplication (SpMV) costing $\mathcal{O}(\text{nnz})$.

The differentiation problem. Given loss $\mathcal{L}(\mathbf{x})$ where $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$, we need gradients with respect to both \mathbf{b} and the non-zero values of \mathbf{A} . If we differentiate through k iterations of CG, the computational graph has $\mathcal{O}(k)$ nodes, each storing intermediate vectors. For $k = 1000$ iterations and $n = 10^6$, this requires ~ 80 GB just for the graph. Our adjoint approach reduces this to $\mathcal{O}(1)$ nodes.

3.2 Adjoint-Based Differentiation

The adjoint method computes gradients through implicit functions without storing intermediate solver states. We present it for linear systems, eigenvalue problems, and nonlinear systems.

3.2.1 Linear Systems: The Core Idea

Consider $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$. Rather than differentiating through the solver iterations, we use the *implicit function theorem*. The solution satisfies $\mathbf{Ax} - \mathbf{b} = \mathbf{0}$. Differentiating implicitly:

$$d\mathbf{A} \cdot \mathbf{x} + \mathbf{A} \cdot d\mathbf{x} = d\mathbf{b} \implies d\mathbf{x} = \mathbf{A}^{-1}(d\mathbf{b} - d\mathbf{A} \cdot \mathbf{x}) \quad (1)$$

For a scalar loss $\mathcal{L}(\mathbf{x})$, the chain rule gives:

$$d\mathcal{L} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \right)^\top d\mathbf{x} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \right)^\top \mathbf{A}^{-1}(d\mathbf{b} - d\mathbf{A} \cdot \mathbf{x}) \quad (2)$$

Define the **adjoint variable** $\boldsymbol{\lambda} = \mathbf{A}^{-\top} \frac{\partial \mathcal{L}}{\partial \mathbf{x}}$, solved via $\mathbf{A}^\top \boldsymbol{\lambda} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}}$. Then the gradients are:

$$\boxed{\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \boldsymbol{\lambda}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{A}_{ij}} = -\boldsymbol{\lambda}_i \cdot \mathbf{x}_j} \quad (3)$$

Table 1: Complexity comparison: naive vs. adjoint differentiation through iterative solver with k iterations, n unknowns, and nnz non-zeros.

	Naive (through iterations)	Adjoint (ours)
Computational graph nodes	$\mathcal{O}(k)$	$\mathcal{O}(1)$
Memory for graph	$\mathcal{O}(k \cdot n)$	$\mathcal{O}(n + \text{nnz})$
Backward pass time	$\mathcal{O}(k \cdot \text{nnz})$	$\mathcal{O}(T_{\text{solve}} + \text{nnz})$

Complexity analysis. Table 1 compares naive differentiation (through iterations) with our adjoint approach.

The forward solve computes \mathbf{x} in time T_{solve} . The backward pass:

1. Solves $\mathbf{A}^\top \boldsymbol{\lambda} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}}$: $\mathcal{O}(T_{\text{solve}})$ time
2. Computes $\frac{\partial \mathcal{L}}{\partial \mathbf{A}_{ij}} = -\boldsymbol{\lambda}_i \mathbf{x}_j$ for each non-zero: $\mathcal{O}(\text{nnz})$ time

Algorithm. The complete procedure is shown below:

Algorithm 1: Adjoint Linear Solve

Input: Sparse matrix \mathbf{A} (values, rows, cols), RHS \mathbf{b}

Output: Solution \mathbf{x} , gradients in backward pass

Forward pass:

1. $\mathbf{x} \leftarrow \text{solve}(\mathbf{A}, \mathbf{b})$ *// Any solver: CG, LU, etc.*
2. Store (\mathbf{A}, \mathbf{x}) for backward

Backward pass: (given $\partial \mathcal{L} / \partial \mathbf{x}$)

1. $\boldsymbol{\lambda} \leftarrow \text{solve}(\mathbf{A}^\top, \partial \mathcal{L} / \partial \mathbf{x})$ *// One adjoint solve*
2. $\partial \mathcal{L} / \partial \mathbf{b} \leftarrow \boldsymbol{\lambda}$
3. For each non-zero (i, j) : $\partial \mathcal{L} / \partial \mathbf{A}_{ij} \leftarrow -\boldsymbol{\lambda}_i \cdot \mathbf{x}_j$ *// $\mathcal{O}(\text{nnz})$ total*

3.2.2 Eigenvalue Problems

For symmetric eigenvalue problem $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ with normalized eigenvector $\|\mathbf{v}\| = 1$, the gradient of eigenvalue λ with respect to matrix entries is remarkably simple [Magnus, 1985]:

$$\frac{\partial \lambda}{\partial \mathbf{A}_{ij}} = v_i \cdot v_j \quad (4)$$

This is the outer product $\mathbf{v}\mathbf{v}^\top$ restricted to the sparsity pattern. For k eigenvalues, the total cost is $\mathcal{O}(k \cdot \text{nnz})$.

3.2.3 Nonlinear Systems

For nonlinear system $\mathbf{F}(\mathbf{u}, \boldsymbol{\theta}) = \mathbf{0}$ with solution $\mathbf{u}^*(\boldsymbol{\theta})$ and loss $\mathcal{L}(\mathbf{u}^*)$, implicit differentiation gives:

$$\mathbf{J}^\top \boldsymbol{\lambda} = \frac{\partial \mathcal{L}}{\partial \mathbf{u}}, \quad \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = -\boldsymbol{\lambda}^\top \frac{\partial \mathbf{F}}{\partial \boldsymbol{\theta}} \quad (5)$$

where $\mathbf{J} = \frac{\partial \mathbf{F}}{\partial \mathbf{u}}$ is the Jacobian at the solution.

Algorithm 2: Newton-Raphson with Adjoint Gradients

Input: Residual function $\mathbf{F}(\mathbf{u}, \boldsymbol{\theta})$, initial guess \mathbf{u}_0

Output: Solution \mathbf{u}^*

Forward (Newton iteration):

1. $\mathbf{u} \leftarrow \mathbf{u}_0$
2. **while** $\|\mathbf{F}(\mathbf{u}, \boldsymbol{\theta})\| > \epsilon$ **do**:
 - (a) $\mathbf{J} \leftarrow \partial \mathbf{F} / \partial \mathbf{u}$ *// Jacobian via autograd*
 - (b) $\Delta \mathbf{u} \leftarrow \text{solve}(\mathbf{J}, -\mathbf{F})$
 - (c) $\mathbf{u} \leftarrow \mathbf{u} + \alpha \Delta \mathbf{u}$ *// α from line search*
3. Store $(\mathbf{u}^*, \mathbf{J}, \boldsymbol{\theta})$

Backward: (given $\partial \mathcal{L} / \partial \mathbf{u}^*$)

1. $\boldsymbol{\lambda} \leftarrow \text{solve}(\mathbf{J}^\top, \partial \mathcal{L} / \partial \mathbf{u}^*)$ *// One adjoint solve*
2. $\partial \mathcal{L} / \partial \boldsymbol{\theta} \leftarrow -\boldsymbol{\lambda}^\top \partial \mathbf{F} / \partial \boldsymbol{\theta}$ *// VJP via autograd*

Memory analysis. Forward: $\mathcal{O}(n + \text{nnz})$ for solution and Jacobian. Backward: $\mathcal{O}(n)$ for adjoint variable. Total: $\mathcal{O}(n + \text{nnz})$, independent of Newton iterations.

3.3 Sparse Tensor Parallel Computing

For problems exceeding single-GPU memory, we implement domain decomposition with halo exchange—the standard approach in industrial CFD/FEM codes like OpenFOAM and Ansys Fluent.

3.3.1 Domain Decomposition

Given a sparse matrix \mathbf{A} corresponding to a mesh or graph, we partition the n nodes into P subdomains. Each process p owns nodes \mathcal{O}_p (owned nodes) and maintains copies of neighboring nodes \mathcal{H}_p (halo/ghost nodes) needed for local computation.

Partitioning strategies:

- **METIS** [Karypis and Kumar, 1998]: Graph partitioning minimizing edge cuts for load balancing
- **RCB**: Recursive Coordinate Bisection using node coordinates
- **Contiguous**: Simple row-based partitioning (fallback)

3.3.2 Halo Exchange

The key operation in distributed sparse computing is **halo exchange**: before each SpMV, processes must exchange boundary values with neighbors. Figure 1 illustrates this concept.

Algorithm 3: Distributed SpMV with Halo Exchange

Input: Local vector $\mathbf{x}_{\text{local}}$, neighbor map, local matrix $\mathbf{A}_{\text{local}}$

Output: Result $\mathbf{y}_{\text{owned}}$

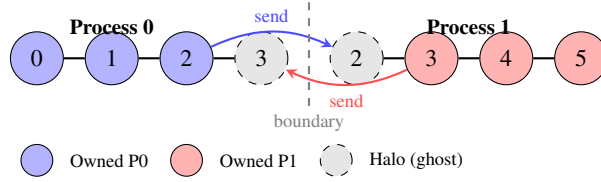


Figure 1: Halo exchange in domain decomposition. Each process owns a subset of nodes (solid colored) and maintains halo copies of boundary neighbors (dashed). Before SpMV, processes exchange updated values at partition boundaries via peer-to-peer communication.

Phase 1: Initiate async communication

1. **for** each neighbor q **do**:
 - (a) `async_send(my boundary values to q)`
 - (b) `async_recv(halo values from q)`

Phase 2: Wait for completion

1. `synchronize_all()`

Phase 3: Local computation

1. $\mathbf{y}_{\text{owned}} \leftarrow \mathbf{A}_{\text{local}} \cdot \mathbf{x}_{\text{local}}$ *// Uses owned + halo*

3.3.3 Distributed Conjugate Gradient

Building on distributed SpMV, we implement distributed CG. Each iteration requires:

- One halo exchange (in SpMV)
- Two `all_reduce` operations for global dot products

Algorithm 4: Distributed Conjugate Gradient

Input: Distributed \mathbf{A} , local RHS $\mathbf{b}_{\text{owned}}$, tolerance ϵ

Output: Solution $\mathbf{x}_{\text{owned}}$

1. $\mathbf{x} \leftarrow \mathbf{0}, \mathbf{r} \leftarrow \mathbf{b}_{\text{owned}}, \mathbf{p} \leftarrow \mathbf{r}$
2. $\rho \leftarrow \text{all_reduce}(\mathbf{r}^\top \mathbf{r}, \text{SUM})$ *// Global dot product*
3. **while** $\sqrt{\rho} > \epsilon$ **do**:
 - (a) $\mathbf{Ap} \leftarrow \text{DistSpMV}(\mathbf{A}, \mathbf{p})$ *// Algorithm 3*
 - (b) $\alpha \leftarrow \rho / \text{all_reduce}(\mathbf{p}^\top \mathbf{Ap}, \text{SUM})$
 - (c) $\mathbf{x} \leftarrow \mathbf{x} + \alpha \mathbf{p}$
 - (d) $\mathbf{r} \leftarrow \mathbf{r} - \alpha \mathbf{Ap}$
 - (e) $\rho_{\text{new}} \leftarrow \text{all_reduce}(\mathbf{r}^\top \mathbf{r}, \text{SUM})$
 - (f) $\mathbf{p} \leftarrow \mathbf{r} + (\rho_{\text{new}} / \rho) \mathbf{p}$
 - (g) $\rho \leftarrow \rho_{\text{new}}$

Table 2: Single-GPU benchmark on 2D Poisson, H200, float64.

DOF	SciPy	cuDSS	PyTorch CG	Memory	Residual
10K	24 ms	128 ms	20 ms	36 MB	10^{-9}
100K	29 ms	630 ms	43 ms	76 MB	10^{-7}
1M	19.4 s	7.3 s	190 ms	474 MB	10^{-7}
2M	52.9 s	15.6 s	418 ms	916 MB	10^{-7}
16M	OOM	OOM	7.3 s	7.1 GB	10^{-6}
169M	OOM	OOM	224 s	74.8 GB	10^{-6}

Table 3: Distributed solve on $4 \times$ H200 GPUs.

DOF	Time	Memory/GPU	Speedup vs. CPU
10K	0.18 s	0.03 GB	$2\times$
100K	0.61 s	0.05 GB	$12\times$
1M	2.82 s	0.27 GB	–
2M	6.02 s	0.50 GB	–

Communication complexity. Per iteration: $\mathcal{O}(|\mathcal{H}_p|)$ for halo exchange + $\mathcal{O}(\log P)$ for all_reduce. Total per solve: $\mathcal{O}(k \cdot (|\mathcal{H}_p| + \log P))$ where k is iterations.

Gradient support. Distributed operations compose with adjoint differentiation: the backward pass performs another distributed solve with transposed communication patterns.

4 Experiments

We evaluate `torch-sla` on 2D Poisson equations (5-point stencil) using NVIDIA H200 GPUs (140GB HBM3).

4.1 Single-GPU Scalability

Table 2 compares solver backends. Key findings:

- **Iterative solvers scale:** PyTorch CG reaches 169M DOF with $\mathcal{O}(n^{1.1})$ time complexity.
- **Direct solvers limited:** cuDSS/SciPy hit OOM at ~ 2 M DOF due to fill-in.
- **Memory efficient:** 443 bytes/DOF for CG (theoretical minimum: 144 bytes/DOF).

4.2 Multi-GPU Performance

Table 3 shows distributed CG on 4 GPUs with NCCL backend.

With $4 \text{ GPUs} \times 140\text{GB}$, theoretical maximum is $\sim 1.3\text{B}$ DOF.

4.3 Gradient Verification

We verify gradient correctness against finite differences. All operations achieve relative error $< 10^{-5}$, confirming correct adjoint implementation.

5 Conclusion

We presented `torch-sla`, a differentiable sparse linear algebra library with two innovations: (1) adjoint-based differentiation achieving $\mathcal{O}(1)$ graph complexity, and (2) distributed sparse computing with halo exchange. The library scales to 169M DOF on single GPU and provides $12\times$ multi-GPU speedup.

5.1 Future Work

- **Learned preconditioners:** Neural networks for adaptive preconditioning.
- **Mixed precision:** FP16/BF16 with FP64 accumulation.
- **Sparse-sparse products:** Efficient SpGEMM with gradients.
- **Integration:** Combining with neural operators [Li et al., 2020] for hybrid solvers.

Availability. MIT license: <https://github.com/walkerchi/torch-sla>

References

- Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J Zico Kolter. Differentiable convex optimization layers. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *International Conference on Machine Learning*, pages 136–145, 2017.
- Satish Balay, Shrirang Abhyankar, Mark F Adams, et al. Petsc users manual. Technical report, Argonne National Laboratory, 2019.
- Mathieu Blondel, Quentin Berthet, Marco Cuturi, et al. Efficient and modular implicit differentiation. *Advances in Neural Information Processing Systems*, 35:5230–5242, 2022.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Robert D Falgout and Ulrike Meier Yang. hypre: A library of high performance preconditioners. In *International Conference on Computational Science*, pages 632–641, 2002.
- Alan George. Nested dissection of a regular finite element mesh. *SIAM Journal on Numerical Analysis*, 10(2):345–363, 1973.
- Michael A Heroux, Roscoe A Bartlett, et al. An overview of the trilinos project. *ACM Transactions on Mathematical Software*, 31(3):397–423, 2005.
- Thomas JR Hughes. *The finite element method: linear static and dynamic finite element analysis*. Courier Corporation, 2012.
- Hrvoje Jasak, Aleksandar Jemcov, and Zeljko Tukovic. Openfoam: A c++ library for complex physics simulations. *International Workshop on Coupled Methods in Numerical Dynamics*, 1000:1–20, 2007.
- George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1998.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, et al. Fourier neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895*, 2020.
- Jan R Magnus. On differentiating eigenvalues and eigenvectors. *Econometric Theory*, 1(2):179–191, 1985.
- NVIDIA Corporation. cuDSS: NVIDIA CUDA direct sparse solver library. <https://developer.nvidia.com/cudss>, 2024.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32, 2019.

Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.

Yousef Saad. Iterative methods for sparse linear systems. *SIAM*, 2003.

Pauli Virtanen, Ralf Gommers, Travis E Oliphant, et al. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature Methods*, 17(3):261–272, 2020.

A Implementation Details

A.1 Backend Selection

`torch-sla` automatically selects backends based on device and problem size:

- CPU, any size: `scipy` with SuperLU
- CUDA, <2M DOF: `cudss` with Cholesky (if SPD) or LU
- CUDA, >2M DOF: `pytorch` with CG+Jacobi preconditioner

A.2 API Examples

Listing 1: Basic usage with gradient support.

```
1 import torch
2 from torch_sla import SparseTensor
3
4 # Create sparse matrix (COO format)
5 A = SparseTensor(val, row, col, shape)
6 b = torch.randn(n, requires_grad=True)
7
8 # Solve with automatic differentiation
9 x = A.solve(b) # Forward: any backend
10 loss = x.sum()
11 loss.backward() # Backward: adjoint method
12 print(b.grad) # Gradient w.r.t. RHS
```

Listing 2: Distributed multi-GPU solve.

```
1 import torch.distributed as dist
2 from torch_sla import DSparseMatrix
3
4 dist.init_process_group(backend='nccl')
5 rank = dist.get_rank()
6
7 # Each process loads its partition
8 A = DSparseMatrix.from_global(
9     val, row, col, shape,
10     num_partitions=4, my_partition=rank
11 )
12
13 # Distributed CG with halo exchange
14 x = A.solve(b_local, atol=1e-10)
```


B Complexity Proofs

Theorem 1 (Adjoint memory complexity). *The adjoint method for linear solve requires $\mathcal{O}(n + \text{nnz})$ memory, independent of solver iterations.*

Proof. Forward pass stores: solution $\mathbf{x} \in \mathbb{R}^n$, matrix indices ($\mathcal{O}(\text{nnz})$), matrix values ($\mathcal{O}(\text{nnz})$). Backward pass computes adjoint $\boldsymbol{\lambda} \in \mathbb{R}^n$ and gradients $\partial\mathcal{L}/\partial\mathbf{A} \in \mathbb{R}^{\text{nnz}}$. No intermediate solver states are stored. Total: $\mathcal{O}(n + \text{nnz})$. \square