



Full paper  
arXiv:2509.21137



Codebase  
disys-lab.github.io/meliso

## ABSTRACT

**Large-scale linear optimization** underpins decision-making across science, engineering, and industry, yet its growing computational demand increasingly exceeds the efficiency of **conventional von Neumann hardware**.

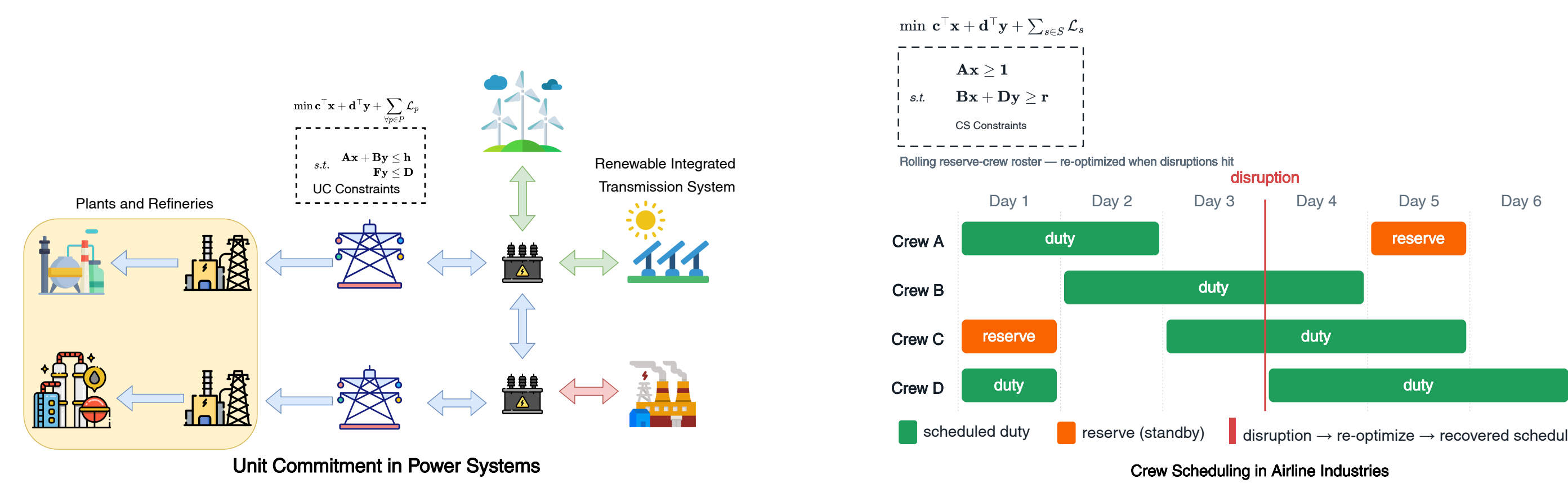
We present the **first primal-dual hybrid gradient (PDHG) linear-program solver co-designed for in-memory computing on resistive-RAM (RRAM) crossbar arrays**. Evaluated in a physics-based device simulation, the solver attains **accuracy comparable to a GPU** while reducing energy and latency by up to **three orders of magnitude**.

## MOTIVATION – Where linear programs arise

A **linear program (LP)** optimizes a linear objective function (e.g. cost, energy, or profit) subject to **linear constraints** that encode real-world limits (budgets, capacities, supply and demand):

$$\min_{x \in \mathbb{R}^n} c^T x, \quad \text{s.t.} \quad Ax \leq b; \quad l \leq x_i \leq u, \quad \forall x_i \in x$$

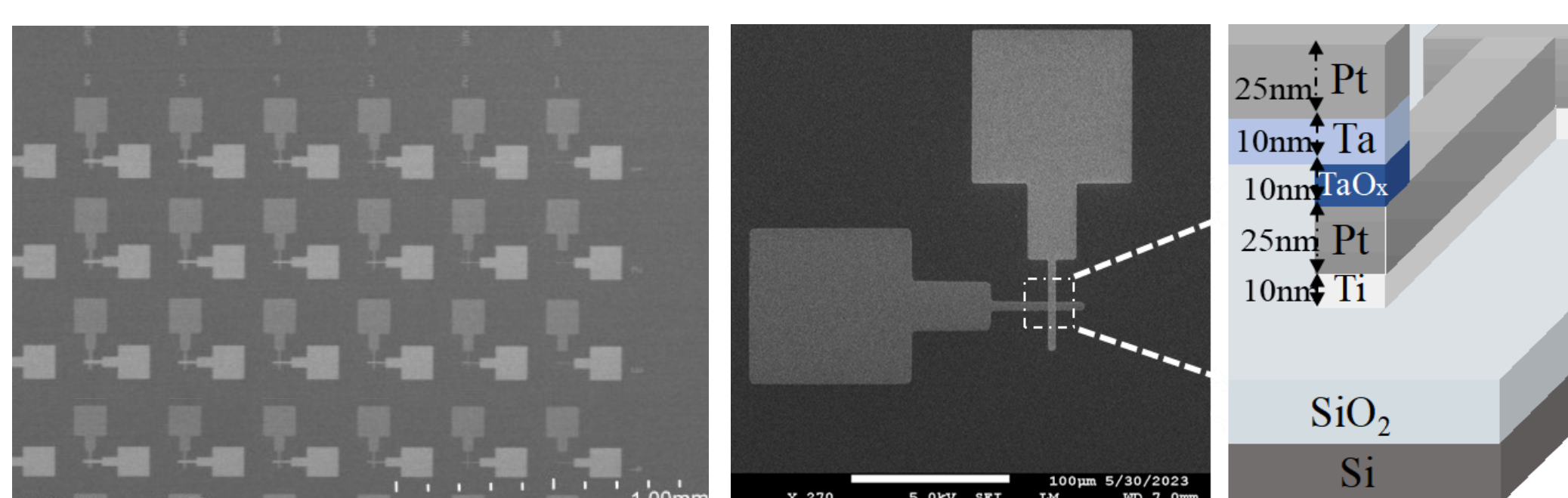
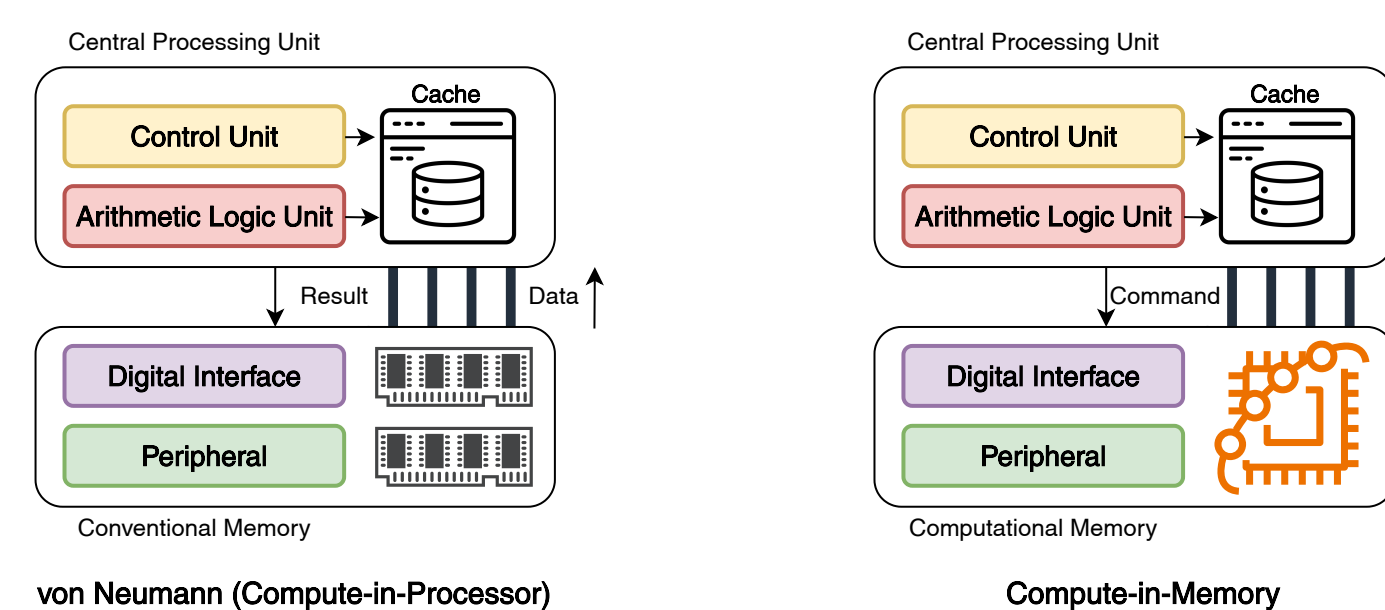
Such problems underpin optimal decision-making across modern infrastructure:



## MOTIVATION – The memory wall and in-memory computing

Conventional von Neumann architectures separate the **processor** from the **memory**. For the large matrices underlying these problems, most energy and time is consumed **transferring data** between the two – the **memory wall**.

**In-memory computing via RRAM technology** eliminates this overhead: RRAM cells encode numerical values as electrical conductances arranged in a grid (a **crossbar**); thus, applying input voltages yields a complete matrix-vector product as output currents in a **single analog step**.

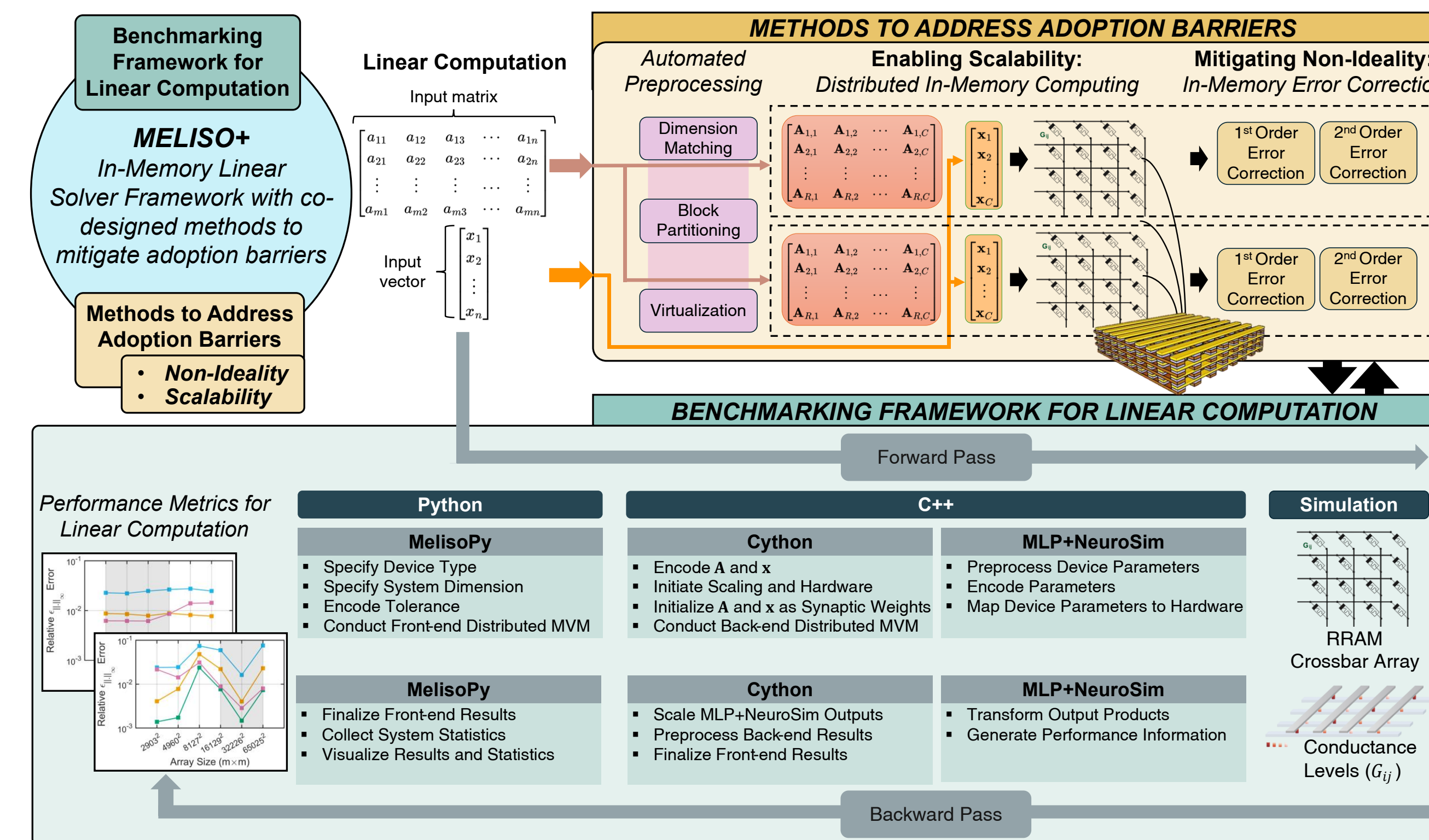


## CHALLENGES – Each generation of solver trades one limitation for another

- Exact digital solvers:** Classical methods (e.g. Gurobi) rely on repeated **matrix factorizations**. They return accurate, certified solutions but incur prohibitive memory overhead at scale.
- GPU first-order methods:** Newer methods such as PDHG exploit the abundance of GPUs to avoid factorization altogether, yet remain bound by the **von Neumann bottleneck**.
- Analog in-memory computing:** Devices such as RRAM are **noisy and non-ideal**. Additionally, the **writing operations are expensive** in RRAMs, so the problem instances (a.k.a. the matrices) **must not be re-encoded** each iteration.

## METHODOLOGY – The MELISO+ simulation framework

We first build MELISO+, an in-memory linear-solver framework that simulates RRAMs under realistic device physics to measure **accuracy, latency, and energy**. It establishes how to perform the core primitive (**matrix-vector multiplication**, abbreviated as MVM) reliably on RRAM crossbar arrays by addressing two adoption barriers: **scalability** and **device non-ideality**.



- Message Passing Interface (MPI)** framework is employed to distribute the problem instances to multiple crossbar arrays (MCAs).
- Device properties** (e.g., device-to-device and cycle-to-cycle variations) are simulated to observe their effects on MVM results.

## METHODOLOGY – In-memory optimization with PDHG

We then configure MELISO+ to solve linear programs entirely in memory using the PDHG method. Recast as a saddle-point problem,

$$\min_{l \leq x \leq u} \max_y c^T x + \langle Ax - b, y \rangle,$$

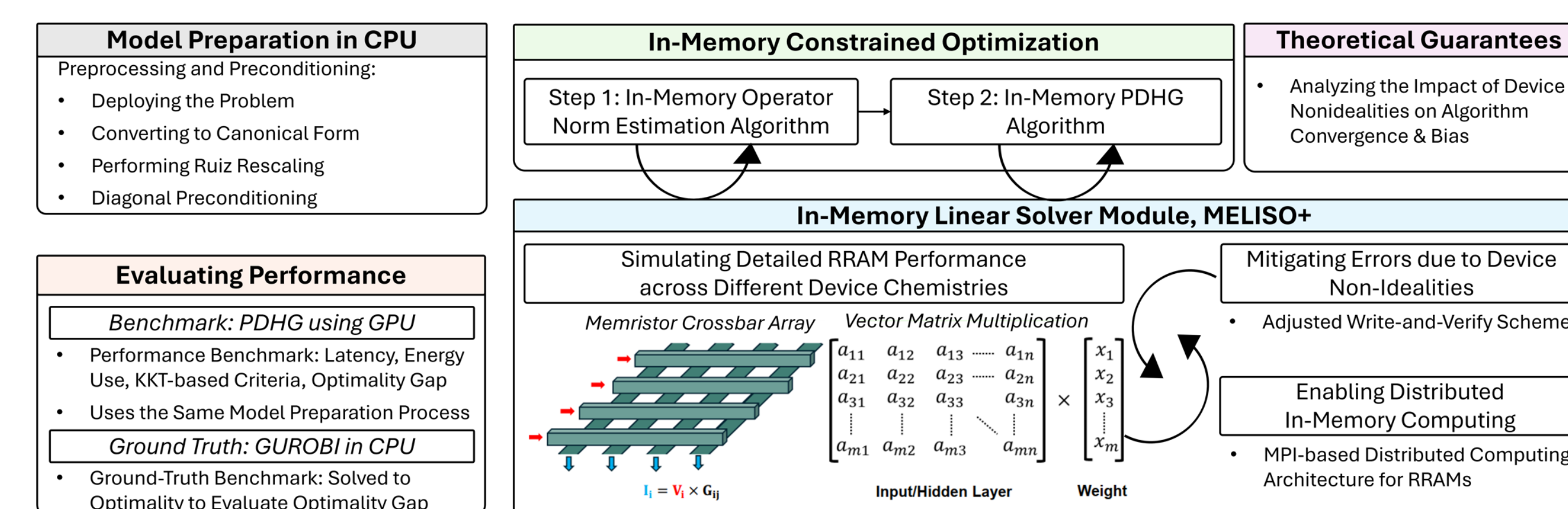
each iteration touches the matrix *only* through matrix-vector products:

$$\begin{aligned} y_{k+1} &= \Pi_Y [y_k + \sigma (Ax_k - b)], \\ x_{k+1} &= \Pi_{[l, u]} [x_k - \tau (c + A^T y_{k+1})], \\ \bar{x}_{k+1} &= x_{k+1} + \theta (x_{k+1} - x_k). \end{aligned}$$

Here, the products  $Ax_k$  and  $A^T y_{k+1}$  are the in-memory MVMs;  $\Pi$  are cheap projections,  $\{\sigma, \tau\}$  is the dual-primal step size, and  $\theta$  is the extrapolation term.

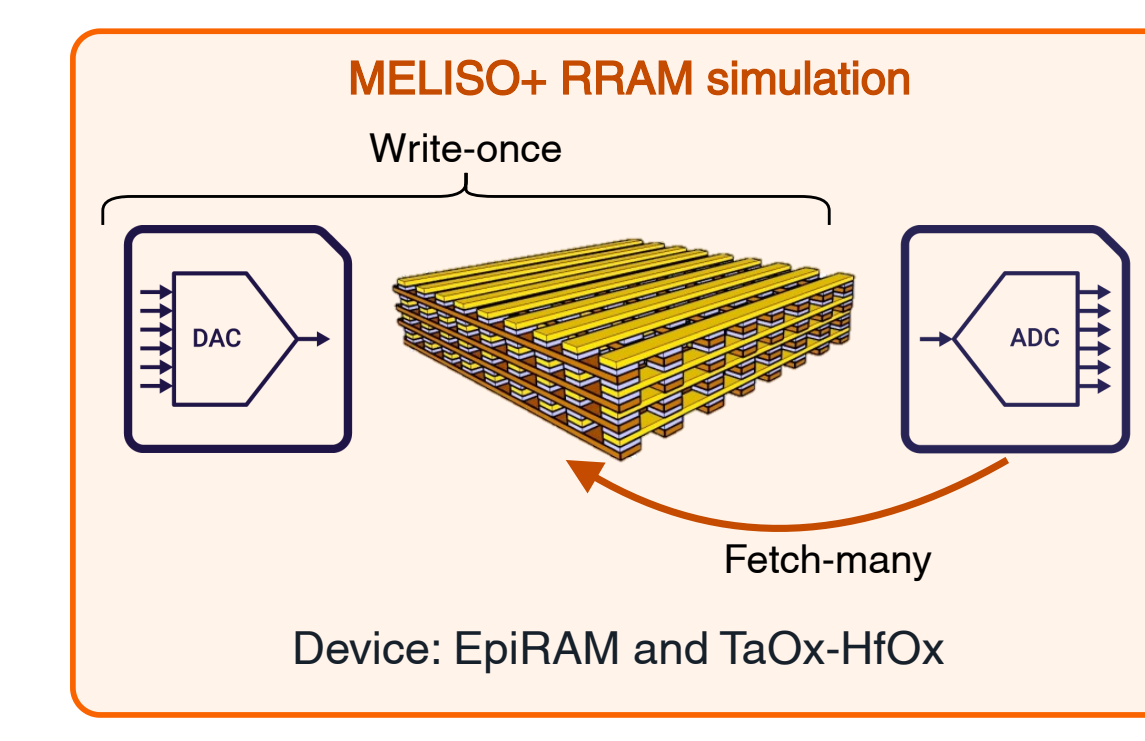
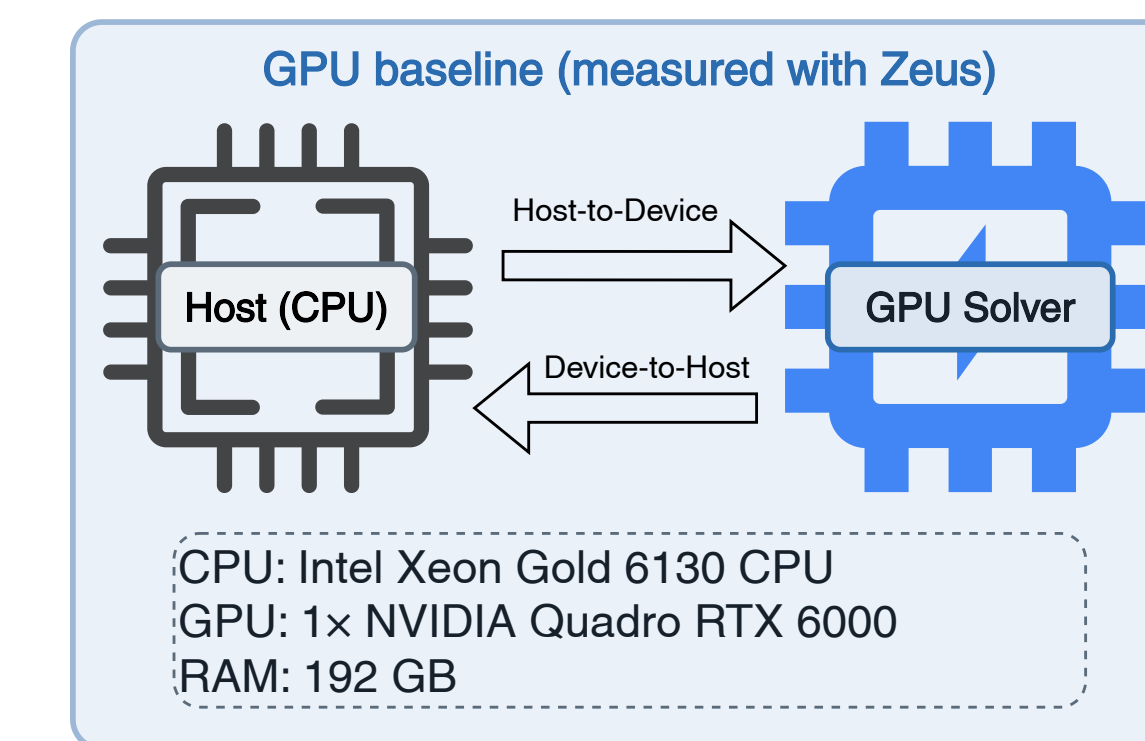
We introduce the following modifications to ensure the feasible implementation on RRAM devices:

- Write-once encoding.** The problem is cast as a single **symmetric block matrix**  $M = \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix}$  and programmed onto the crossbars **exactly once**; every iteration reuses this encoding, with only the input vectors changing.
- Distributed single accelerator.** A  $4 \times 4$  tile of  $64 \times 64$  crossbars (a  $256 \times 256$  logical array) operates as one accelerator, eliminating inter-tile data movement.
- Accelerated convergence.** The PDHG iteration is augmented with **Nesterov momentum** for adaptive step-size rule, reducing the number of iterations required to reach a target tolerance.
- Robustness by design.** Stable operator-norm estimation (**Lanczos**) and **write-and-verify** programming mitigate device non-idealities, and we **prove convergence** on noisy hardware – no worse than a digital baseline up to a small, fixed error term.

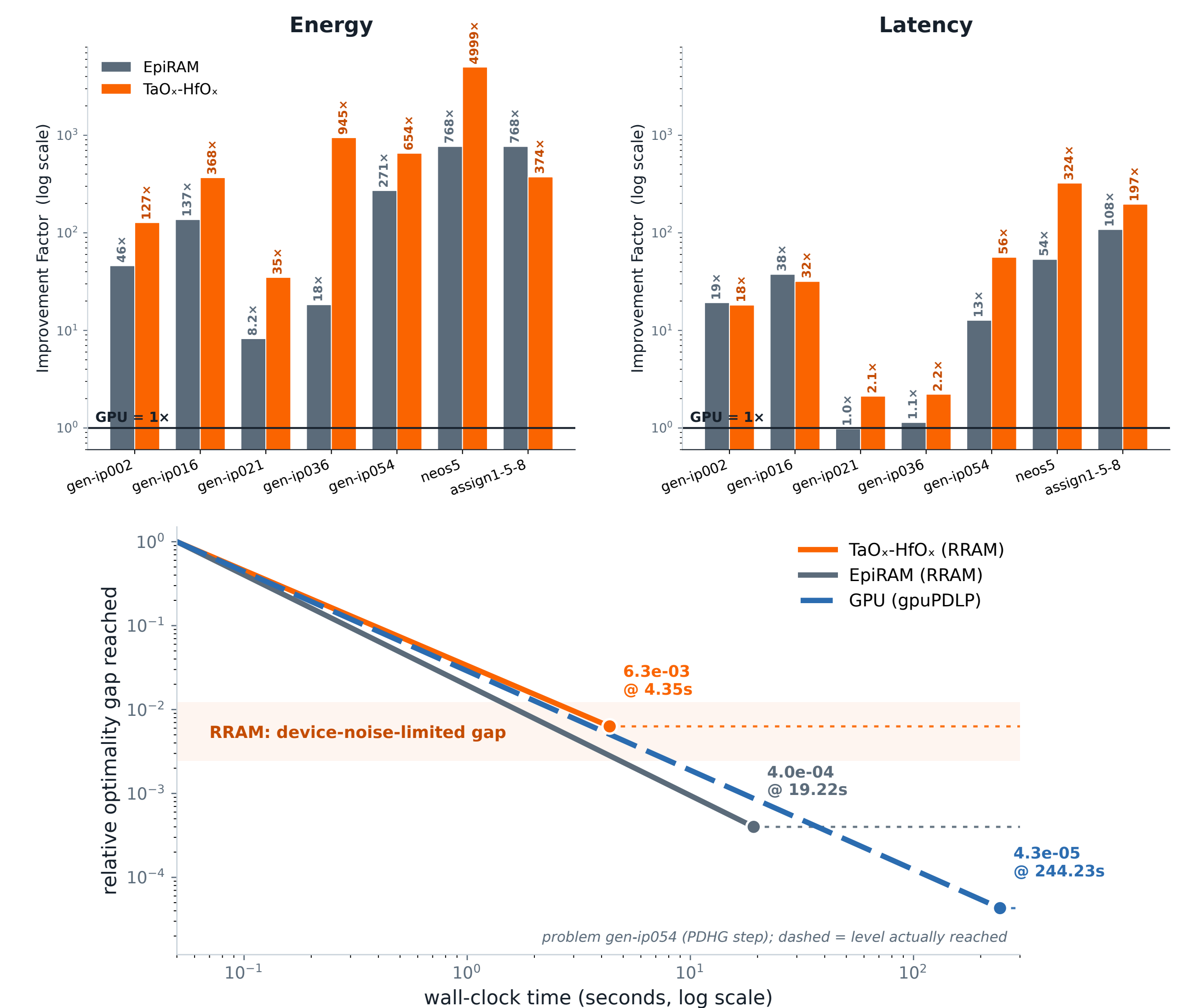


## EXPERIMENTS AND RESULTS

Benchmarked against **gpuPDL** (a GPU version of the same method), with **Gurobi** as ground truth, on real LPs from MIPLIB-2017. Both RRAM devices (**TaO<sub>x</sub>-HfO<sub>x</sub>** and **EpiRAM**) match GPU/Gurobi accuracy (relative gap  $\Delta_{rel}$ ) at a fraction of the energy and latency.



Total energy and latency measurement scheme



## CONCLUSION

- The **first PDHG linear-program solver simulated on RRAM** – a true **algorithm-hardware co-design** (write-once encoding + noise-robust iteration).
- Provable convergence** on noisy analog hardware; **orders-of-magnitude** energy and latency gains over a GPU.

## REFERENCES & ACKNOWLEDGMENTS

- Vo, Chowdhury, Ramanan et al. From GPUs to RRAMs: Distributed In-Memory Primal-Dual Hybrid Gradient Method for Solving Large-Scale Linear Optimization Problems. *SIAM PP* (2026). arXiv:2509.21137.
  - Vo, Chowdhury, Ramanan et al. Harnessing the full potential of RRAMs through scalable and distributed in-memory computing with integrated error correction. *Nature Communications Engineering* (2026). arXiv:2508.13298
  - Chowdhury et al. The Lynchpin of In-Memory Computing: A Benchmarking Framework for Vector-Matrix Multiplication in RRAMs. *ICONS* (2024). arXiv:2409.06140
  - Lu et al. A GPU Implementation of Restarted Primal-Dual Hybrid Gradient for Linear Programming in Julia. arXiv:2311.12180
- Conducted at Oklahoma State University, Argonne National Laboratory, and Wayne State University. Supported in part by the National Science Foundation Secure and Trustworthy Cyberspace (SaTC) program under Award No. 2348411.