cuPDLP.jl: A GPU Implementation of Restarted Primal-Dual Hybrid Gradient for Linear Programming in Julia

Haihao Lu*

Jinwen Yang[†]

June 2024

Abstract

In this paper, we provide an affirmative answer to the long-standing question: Are GPUs useful in solving linear programming? We present cuPDLP.jl, a GPU implementation of restarted primal-dual hybrid gradient (PDHG) for solving linear programming (LP). We show that this prototype implementation in Julia has comparable numerical performance on standard LP benchmark sets to Gurobi, a highly optimized implementation of the simplex and interior-point methods. This demonstrates the power of using GPUs in linear programming, which, for the first time, showcases that GPUs and first-order methods can lead to performance comparable to state-of-the-art commercial optimization LP solvers on standard benchmark sets.

1 Introduction

Linear programming (LP) is a fundamental optimization problem class with a long history and a vast range of applications in operation research and computer science, such as agriculture, transportation, telecommunications, economics, production and operations scheduling, strategic decision-making, etc [22, 14, 12, 29, 10, 49].

Since the 1940s, speeding up and scaling up LP has been a central topic in the optimization community, with extensive studies from both academia and industry. The current general-purpose LP solvers, such as Gurobi [41], COPT [18], CPLEX [34] and HiGHS [26], are quite mature. These classic LP solvers are based on either the simplex method or interior-point methods (IPMs), which can generally provide high-quality solutions to LP. However, further scaling up or speeding up these methods is highly challenging. The fundamental difficulty is the necessity of solving linear systems in simplex and IPMs, which requires either LU factorization (for simplex method) or Cholesky factorization (for IPMs). Such factorization-based approaches have two major drawbacks when solving large instances: (i) Storing the factorization can be quite memory-demanding. It is often the case that a sparse matrix has a much denser factorization, which is one of the reasons for classic solvers to raise an "out-of-memory" error even if the instance can be stored in memory; (ii) Both methods are highly nontrivial to exploit massive parallelization, and thus they are challenging to take advantage of modern computing architectures such as graphic processing units (GPUs) and

^{*}The University of Chicago, Booth School of Business (haihao.lu@chicagobooth.edu).

[†]The University of Chicago, Department of Statistics (jinweny@uchicago.edu).

distributed computing, due to the sequential nature of factorizations. Accordingly, all classic LP solvers are CPU-based and implemented on a single shared memory machine.

In contrast, the modern deep learning models used in practice, such as GPT-4, have trillions of variables, which are arguably significantly larger than the scale of LP commercial solvers are capable to solve¹. One commonly believed reason for the success of deep learning models is the extensive use of GPUs and distributed computing, which massively speeds up the neural network training process. Indeed, commercial solver companies, such as Gurobi, always try to use GPUs to speed up their solvers, but the early efforts were unsuccessful [20, 21]. The fundamental reason is that GPUs do not work well for solving sparse linear systems, which is the computational bottleneck of simplex or barrier method solving linear programming [48, 20, 21]².

This paper revisits this natural question:

We provide an affirmative answer to this question by presenting cuPDLP.jl, a GPU implementation of restarted primal-dual hybrid gradient (PDHG) implemented in Julia. We present an extensive numerical study comparing cuPDLP.jl with its CPU counterpart and three methods implemented in the highly-optimized commercial LP solver Gurobi on standard LP benchmark sets, which showcase that

- cuPDLP.jl, a GPU implementation of PDLP in Julia programming language, has a comparable behavior with the highly-optimized commercial LP solver Gurobi on standard LP benchmark sets.
- Compared with its CPU counterpart, cuPDLP.jl clearly has superior performance on standard benchmark sets. Furthermore, one can observe the strong correlation between the GPU speedup and the size of the instances.

cuPDLP.jl can be viewed as a CUDA implementation of its CPU counterpart PDLP [2]. PDLP has two open-sourced implementations, a prototype implementation in Julia (FirstOrderLp.jl), and a production-level C++ implementation (open-sourced through Google OR-Tools).

Different from classic LP solvers, PDLP is a first-order method (FOM) LP solver. The fundamental difference is that the computational bottleneck of FOM is (sparse) matrix-vector multiplication, in contrast to (sparse) matrix factorization in simplex or IPMs. Table 1 presents a comparison summary of simplex, barrier, and FOM LP solvers in five different dimensions: cost per iteration, number of iterations needed to solve an instance, code complexity, whether they can take advantage of massive parallelization, and whether they need extensive memory usage (beyond just storing the instance itself). As we can see, FOMs have multiple advantages in solving large instances, compared with simplex and IPMs.

In particular, thanks to the recent development of deep learning, (sparse) matrix-vector multiplication suits very well on modern GPU infrastructure. Additionally, it is also shown that PDLP has strong theoretical guarantees [4, 3, 30, 32] and superior numerical performance compared to other

¹We comment that "solving" means a different order of accuracy tolerance in LP and deep learning models, and nevertheless, there is a significant gap in scale.

²NVIDIA recently released a library, cuDSS, for direct solving of sparse linear systems on GPUs in March 2024.

	Simplex	Barrier	FOMs
Cost per iteration			•••
No. of iterations	••		••
Code complexity		••	
Massive parallelization			
Memory usage	••	••	

Table 1: Comparison of simplex, barrier, and first-order methods in five dimensions, where green, gray, and red faces represent a decreasing order of favorableness of the corresponding method in the corresponding dimension.

FOM-based solvers [2]. These are the reasons we choose PDLP as the base algorithm in our GPU implementation.

To fully unleash the potential offered by modern GPU hardware, there are two major differences and modifications to the original CPU-based PDLP: (1) Due to the slow communication between CPU and GPU, the computational framework of PDLP has to be fully implemented on GPUs. In other words, both the instances and intermediate iterates must be resident within GPU memory (see Section 3.2 for details). (2) In terms of the algorithm, the only significant difference is a major restart scheme – the CPU implementation utilizes normalized duality gap as the metric for restart [4], while we utilize KKT error as the metric for restart (see Section 3.4 for details). This avoids the trust-region algorithm to compute the normalized duality gap proposed in [4], which is sequential in nature and is not friendly for GPU implementation. Furthermore, we show in theory that restarting with KKT error enjoys the same order of linear convergence rate as restarting with normalized duality gap in Appendix A.

This methodology is not limited to solving LP, and can potentially be extended to other optimization problems, such as quadratic programming [33] and primal heuristics in mixed-integer programming [35].

1.1 Related literature

Linear programming and classic solvers. The state-of-the-art methods to solve LP are simplex methods [13] and interior-point methods [27]. Based on these methods, commercial solvers such as Gurobi [41] and COPT [18] and open-sourced solvers like HiGHS [26] can provide reliable solutions with high accuracy.

FOM-based LP solvers. There is a recent surge of research on using first-order methods (FOMs) to solve linear programming. FOMs are appealing due to their low per-iteration cost and ability of parallelization.

- PDHG-based solvers: PDLP [2]. PDLP is a general-purpose large-scale LP solver. PDLP is built upon restarted PDHG algorithm [4], with many practical algorithmic enhancements, such as preconditioning, adaptive restart, adaptive step-size, infeasibility detection, etc. Currently PDLP has three implementations: a prototype implemented in Julia (FirstOrderLp.jl), a production-level C++ implementation (open-sourced through Google OR-Tools), and an internal distributed version at Google [36]. There is also an extensive study on the theoretical results related to PDLP, such as, the linear convergence rate [30, 4], infeasibility detection [3], refined complexity analysis [32, 24], extensions to quadratic programming [33], etc.
- Matrix-free IPM solvers: ABIP [28, 15]. The core algorithm of ABIP is solving the homogeneous self-dual embedded cone programs via an interior-point method, and as a special case, it can solve LP. ABIP utilizes multiple ADMM iterations instead of one Newton step to approximately minimize the log-barrier penalty function. A recently enhanced version of ABIP (named ABIP+ [15]) includes many new enhancements, such as preconditioning, restart, and hybrid parameter tuning, on top of ABIP. It was shown that ABIP+ (developed in C) has a comparable numerical performance on the LP benchmark sets to the Julia implementation of PDLP.
- Dual-based solvers: ECLIPSE [5]. ECLIPSE is a distributed LP solver. It leverages accelerated gradient descent to solve a smoothed dual form of LP. ECLIPSE is designed specifically to solve large-scale LPs with certain decomposition structures arising from web applications. For example, ECLIPSE is used to solve real-world web applications with 10¹² decision variables at LinkedIn [44, 1].
- ADMM-based solvers: SCS [40, 39] and OSQP [47]. SCS is designed to solve convex cone programs, and OSQP is designed to solve convex quadratic programs. Both solvers are based on ADMM, and the major algorithmic difference is that SCS solves the homogeneous self-dual embedding of general conic programming, while OSQP solves the primal-dual form of quadratic program directly. LP can be solved as a special case of cone programming in SCS as well as a special case of quadratic programming in OSQP. The computational bottleneck of ADMM-based methods is solving a linear system with similar forms every iteration. Both SCS and OSQP support using either direct solve of the linear system via factorization or indirect solve of the linear system via the conjugate gradient method. The first approach still requires doing one factorization and thus may suffer from the same scaling difficulties as simplex and IPMs. The second approach usually requires several conjugate gradient steps (i.e., matrix-vector multiplications) for every iteration.

Primal-dual hybrid gradient method (PDHG). PDHG is an operator splitting method initially designed for applications in image processing [8, 11, 16, 23, 50]. PDHG exhibits a sub-linear rate on general convex-concave problems [9, 31] and achieves linear rate on a wide class of problems [30, 17]. PDLP [2] utilizes PDHG as its base algorithm. Built on the ergodic convergence of PDHG, a restart variant of PDHG exhibits an optimal linear convergence rate for solving LP [4]. Furthermore, Applegate et al. [3] shows how to extract infeasibility information of LP from PDHG iterates.

Optimization with GPUs. Previous efforts of commercial solver companies to solve LP on GPUs have not yet been successful due to the inefficiency of solving sparse linear systems on GPUs [20, 21, 48]. Recently, there have been open-sourced efforts using GPUs to solve generic

optimization problems. For example, the ADMM-based solvers SCS [39] and OSQP [47] have an implementation of indirect solves via conjugate gradient on GPUs, which avoids solving linear equations. MadNLP.jl [46] is an interior-point method based solver for nonlinear programming that can be run on GPUs. However, these solvers are designed for solving different and more general classes of optimization problems, and they may not be suitable for solving LP instances with the size we consider herein (see Appendix B for a comparison with SCS).

2 PDLP

In this section, we introduce the basics of LP and the algorithmic components of PDLP, as stated in [2]. After a brief introduction of vanilla PDHG for solving LP (Section 2.1), we summarize the enhancements of PDLP on top of PDHG that boost the practical performance (Section 2.2).

Consider LP of the general form:

$$\min_{x \in \mathbb{R}^n} \quad c^{\top} x$$
s.t. $Gx \ge h$

$$Ax = b$$

$$1 \le x \le u$$
, (1)

where $G \in \mathbb{R}^{m_1 \times n}$, $A \in \mathbb{R}^{m_2 \times n}$, $c \in \mathbb{R}^n$, $h \in \mathbb{R}^{m_1}$, $b \in \mathbb{R}^{m_2}$, $l \in (\mathbb{R} \cup \{-\infty\})^n$, $u \in (\mathbb{R} \cup \{\infty\})^n$. Notice that the constraint set is a complicated polytope that is hard to be projected onto, which makes it intractable to use standard FOMs, such as, projected gradient descent for solving (1). PDLP solves the primal-dual form of the problem by dualizing the linear constraints:

$$\min_{x \in X} \max_{y \in Y} L(x, y) := c^{\top} x - y^{\top} K x + q^{\top} y , \qquad (2)$$

where $K^{\top} = (G^{\top}, A^{\top})$ and $q^{\top} := (h^{\top}, b^{\top})$, $X := \{x \in \mathbb{R}^n : l \leq x \leq u\}$, and $Y := \{y \in \mathbb{R}^{m_1 + m_2} : y_{1:m_1} \geq 0\}$. By duality theory, it is straightforward to show that a saddle point of (2) recovers an optimal primal-dual solution to (1).

2.1 PDHG for solving LP

The base algorithm of PDLP is the primal-dual hybrid gradient (PDHG, a.k.a. Chambolle-Pock method) [8]. The update of PDHG for solving (2) is

$$\begin{cases} x^{t+1} \leftarrow \operatorname{proj}_X(x^t - \tau(c - K^\top y^t)) \\ y^{t+1} \leftarrow \operatorname{proj}_Y(y^t + \sigma(q - K(2x^{t+1} - x^t))) \end{cases},$$
(3)

where τ, σ are primal and dual step-sizes respectively. The computational bottleneck of PDHG is matrix-vector multiplication $K^{\top}y^t$ and Kx^t ; thus, PDHG is fully matrix-free to solve LP, i.e., it does not need to solve any linear equations. The primal and the dual step-size are reparameterized in PDLP as

$$\tau = \eta/\omega, \ \sigma = \eta\omega \text{ with } \eta, \omega > 0,$$

where η (called step-size) controls the scale of the step-sizes, and ω (called primal weight) balances the primal and the dual progress.

2.2 Algorithmic enhancements in PDLP

We here summarize the algorithmic enhancements in PDLP that is presented in [2]. It turns out that the numerical performance of vanilla PDHG (3) on LP is not strong enough to support a modern solver [2]. To boost the practical performance, PDLP has essentially five major enhancements on top of PDHG: preconditioning, adaptive restart, adaptive step-size, primal weight update, and infeasibility detection. The core algorithm in PDLP, i.e., restarted PDHG, is presented in Algorithm 1, followed by discussions on each algorithmic enhancement adopted.

Algorithm 1: Restarted PDHG (after preconditioning)

```
Input: Initial point z^{0,0}:
 1 Initialize outer loop counter n \leftarrow 0, total iterations k \leftarrow 0, step-size \hat{\eta}^{0,0} \leftarrow 1/\|K\|_{\infty},
       primal weight \omega^0 \leftarrow \text{InitializePrimalWeight}(c, q);
 2 repeat
 3
            t \leftarrow 0;
            repeat
  4
                  z^{n,t+1}, \eta^{n,t+1}, \hat{\eta}^{n,t+1} \leftarrow \text{AdaptiveStepPDHG}(z^{n,t}, \omega^n, \hat{\eta}^{n,t}, k);
\bar{z}^{n,t+1} \leftarrow \frac{1}{\sum_{i=1}^{t+1} \eta^{n,i}} \sum_{i=1}^{k+1} \eta^{n,i} z^{n,i};
z_c^{n,t+1} \leftarrow \text{GetRestartCandidate}(z^{n,t+1}, \bar{z}^{n,t+1});
  5
  6
                  t \leftarrow t + 1, \ k \leftarrow k + 1:
  8
            until restart or termination criteria holds;
 9
            restart the outer loop z^{n+1,0} \leftarrow z_c^{n,t}, n \leftarrow n+1;
10
            \omega^n \leftarrow \text{PrimalWeightUpdate}(z^{n,0}, z^{n-1,0}, \omega^{n-1});
11
12 until termination criteria holds;
      Output: z^{n,0}.
```

- **Preconditioning:** The efficacy of first-order methods is closely tied to the conditioning of the underlying problem. In order to mitigate the ill-posedness, PDLP employs a diagonal preconditioner to ameliorate the condition number of the original problem. Specifically, it involves the rescaling of the constraint matrix K = (G, A) to $\tilde{K} = (\tilde{G}, \tilde{A}) = D_1 K D_2$ where D_1 and D_2 are positive diagonal matrices. This rescaling ensures that the resulting matrix \tilde{K} is "well balanced". Consequently, this preconditioning step gives rise to a modified LP instance, wherein A, G, c, b, h, u and l in (1) are replaced with $\tilde{G}, \tilde{A}, \hat{x} = D_2^{-1}x, \tilde{c} = D_2c, (\tilde{b}, \tilde{h}) = D_1(b, h), \tilde{u} = D_2^{-1}u$ and $\tilde{l} = D_2^{-1}l$. In the default PDLP configuration, a combination of Ruiz rescaling [45] and the preconditioning technique proposed by Pock and Chambolle [43] is employed.
- Adaptive restarts. PDLP utilizes an adaptive restarting strategy to enhance convergence. PDLP initially selects a restart candidate at each iteration, choosing between the current iterate and the average iterate based on a greedy principle. Subsequently, various restart criteria are assessed to determine if there is a constant factor decay in the progress metric. If such decay is observed, a restart is triggered. Further details can be found in [4].

In the CPU-based PDLP, the progress metric for restarting is the normalized duality gap proposed in [4], and a trust-region algorithm is devised to compute this metric efficiently.

While the trust-region algorithm exhibits linear time complexity, it is less compatible with GPUs' massively parallel computing paradigm. Therefore, in cuPDLP.jl, we introduce a novel restart scheme based on the KKT error. Detailed discussions are deferred to Section 3.4.

• Adaptive step-size. The step-size suggested by theoretical considerations, namely $1/\|A\|_2$, turns out to be conservative in practical applications. To address this, PDLP employs a heuristic line search to determine a suitable step-size satisfying the condition:

$$\eta \le \frac{\|z^{t+1} - z^t\|_{\omega}^2}{2(y^{t+1} - y^t)^\top K(x^{t+1} - x^t)} , \tag{4}$$

where $||z||_{\omega} := \sqrt{\omega ||x||_2^2 + \frac{||y||_2^2}{\omega}}$ and ω is the current primal weight. Additional details of the adaptive step-size rule are elaborated in [2]. The inequality (4) was inspired from the $\mathcal{O}(1/k)$ convergence rate proof of PDHG [9, 31]. The empirical evidence from numerical experiments conducted in [2] attests to its consistent efficacy.

Algorithm 2: One step of PDHG using adaptive step-size heuristic

```
Function: AdaptiveStepPDHG(z^{n,t}, \omega^n, \hat{\eta}^{n,t}, k)

1 (x,y) \leftarrow z^{n,t}, \eta \leftarrow \hat{\eta}^{n,t}

2 for i = 0, 1, ... do

3 | x' \leftarrow \operatorname{proj}_X(x - \frac{\eta}{\omega^n}(c - K^\top y))

4 | y' \leftarrow \operatorname{proj}_Y(y + \eta \omega^n(q - K(2x' - x)))

5 | \bar{\eta} \leftarrow \frac{\|(x' - x, y' - y)\|_{\omega^n}^2}{2(y' - y)^\top K(x' - x)}

6 | \eta' \leftarrow \min((1 - (k + 1)^{-0.3})\bar{\eta}, (1 + (k + 1)^{-0.6})\eta)

7 | \text{if } \eta \leq \bar{\eta} \text{ then}

8 | \text{return } (x', y'), \eta, \eta'

9 | \text{end}

10 | \eta \leftarrow \eta'

11 | \text{end}
```

• Primal Weight Update: Adjusting the primal weight ω is designed to harmonize the primal and dual spaces through a heuristic approach. The update of primal weight is specific during restart occurrences, thus infrequently. More precisely, the initialization of ω involves the expression:

$$\label{eq:initializePrimalWeight} \begin{split} \text{InitializePrimalWeight}(c,q) := \begin{cases} \frac{\|c\|_2}{\|q\|_2}, & \text{if } \|c\|_2, \|q\|_2 > \epsilon_{\text{zero}} \\ 1, & \text{otherwise} \end{cases} \end{split}$$

where ϵ_{zero} denotes a small nonzero tolerance. Let $\Delta_x^n = ||x^{n,0} - x^{n-1,0}||_2$ and $\Delta_y^n = ||y^{n,0} - y^{n-1,0}||_2$. PDLP initiates the primal weight update at the beginning of each new epoch.

$$\text{PrimalWeightUpdate}(z^{n,0},z^{n-1,0},\omega^{n-1}) := \begin{cases} \exp\left(\theta\log\left(\frac{\Delta_y^n}{\Delta_x^n}\right) + (1-\theta)\omega^{n-1}\right), & \Delta_x^n, \Delta_y^n > \epsilon_{\text{zero}} \\ \omega^{n-1}, & \text{otherwise} \end{cases}$$

The intuition is to determine the primal weight ω^n in a manner that equalizes the distance to optimality in both the primal and dual domains, i.e., $\|(x^{n,t}-x^*,0)\|_{\omega^n} \approx \|(0,y^{n,t}-y^*)\|_{\omega^n}$.

Additionally, PDLP employs exponential smoothing with a parameter $\theta \in [0, 1]$ to mitigate oscillations.

• Infeasibility detection. PDLP periodically checks whether the difference of iterates $z^{n,t+1} - z^{n,t}$ or the normalized iterates $\frac{1}{t}(z^{n,t} - z^{n,0})$ provide an infeasibility certificate, and the performance of these two sequences is instance-dependent. A detailed investigation of infeasibility detection using PDHG iterates is available in [3].

3 GPU implementation of PDLP

This section presents the design of cuPDLP.jl, the GPU implementation of PDLP. Section 3.1 briefly introduces the hardware architecture and logical structure. In Section 3.2, we discuss in detail the design of cuPDLP.jl. Section 3.3 presents the implementation of basic matrix and vector operations in cuPDLP.jl, and Section 3.4 shows the restart scheme based on KKT error to adapt on GPU. The solver is available at https://github.com/jinwen-yang/cuPDLP.jl.

3.1 GPU architecture and thread hierarchy

GPUs exhibit distinct underlying architecture with CPUs. GPUs have significantly more computation cores than CPUs. For example, the GPU we used in the experiments, NVIDIA H100, has 7296 double-precision cores. However, unlike CPU cores, many GPU cores share the same control unit and must execute the same instruction simultaneously. As a result, the hardware design of GPUs encourages high bandwidth instead of a deep pipeline.

GPUs follow the single instruction multiple data (SIMD) computational paradigm; namely, the threads execute the same instruction but fetch their own data. At the heart of GPU programming is the kernel function, i.e., a program designed to execute instructions on GPUs. Upon launching the kernel function, the execution environment configures a grid of thread blocks, each block comprising an identical number of threads. A block of threads is assigned to an available streaming multiprocessor (SM) during runtime, directing them into warps, with each warp typically encompassing a set of 32 threads. Warp is the basic unit of execution in GPUs. One caveat goes that this does not mean that all thread blocks can run concurrently on SMs, and there are no guarantees on the order of block and warp execution.

This hierarchical structure is the backbone of GPUs to foster parallel execution across threads, which improves computational efficiency by reducing the latency and increasing the throughput (the data processed per time unit). The relationship between the physical architecture and logical structure is depicted in Figure 1. For more detailed discussion in the characteristic of GPUs, refer to [38].

Each GPU device has its own memory hardware, which is separate from the memory of the CPU host. A typical computational paradigm is using CPUs for IO process and exploiting GPUs for burdensome computation. Thus, data moving between CPU and GPU are needed. However, the significant cost of CPU-GPU communication, especially for large-scale problems, can be a crucial issue. Frequent data transfers between CPU and GPU can introduce dominating communication cost over any computational speed-up gained on GPU.

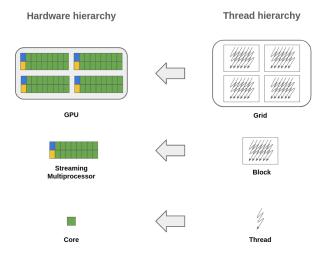


Figure 1: Illustration of the relationship between physical architecture and logical structure.

3.2 Design of cuPDLP.jl

The design of cuPDLP.jl is illustrated in Figure 2. To avoid expensive data transfers between CPU and GPU, we have designed our implementation of cuPDLP.jl to run as much as possible on the GPU.

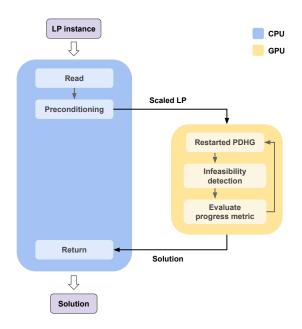


Figure 2: Illustration of the computational architecture of cuPDLP.jl.

As depicted in Figure 2, only two communications between CPU and GPU are required for cuPDLP.jl. One transfers the scaled LP instance after preconditioning from CPU to GPU while the other moves the final solution from GPU to CPU as output. Major iterations are executed

completely on GPU and there is no need to transfer any vectors before termination to alleviate expensive CPU-GPU communication.

3.3 Vector and matrix operations

Matrix-vector multiplications and vector-vector operations, the core of PDLP, can be parallelized naturally and fit well with the SIMD paradigm of GPUs. GPUs have massive parallelization capability to manipulate each coordinate of vectors on each of its core in a parallel fashion. Specifically, the constraint matrix is stored in Compressed Sparse Row (CSR) format and cuPDLP.jl uses the 1-dimensional thread configuration since most of our operations are matrix-vector and vector-vector operations. We write our custom kernels for main PDHG updates (3) in cuPDLP.jl and utilize cusparseSpMV() implemented in cuSPARSE library to do matrix-vector multiplications. In our customized kernel, each thread updates a single coordinate of the iterate vector to maximize the throughput and cusparseSpMV() uses algorithm CUSPARSE_SPMV_CSR_ALG2 to provide a deterministic result for each run.

3.4 Adaptive restart based on KKT error

PDLP utilizes normalized duality gap as the progress metric for restart, which was introduced in [4]. The normalized duality gap for LP is computed by a trust-region algorithm. Although only linear time is needed, it is nontrivial to implement the trust-region method in an efficient parallel version on GPU due to the sequential nature of the trust-region method. We use the KKT error defined in (5) as a proxy to the normalized duality gap for restarting since KKT error can be evaluated very efficiently on GPU.

$$KKT_{\omega}(z) = \sqrt{\omega^{2} \left\| \begin{pmatrix} Ax - b \\ [h - Gx]^{+} \end{pmatrix} \right\|_{2}^{2} + \frac{1}{\omega^{2}} \|c - K^{\top}y - \lambda\|_{2}^{2} + (q^{\top}y + l^{\top}\lambda^{+} - u^{\top}\lambda^{-} - c^{\top}x)^{2}}$$
 (5)

More specifically, the restart scheme is described as follows

Choosing the restart candidate.

$$z_c^{n,t+1} := \text{GetRestartCandidate}(z^{n,t+1},\bar{z}^{n,t+1}) = \begin{cases} z^{n,t+1}, & \text{KKT}_{\omega^n}(z^{n,t+1}) < \text{KKT}_{\omega^n}(\bar{z}^{n,t+1}) \\ \bar{z}^{n,t+1}, & \text{otherwise} \ . \end{cases}$$

Restart criteria. Define parameters $\beta_{\text{sufficient}} = 0.2$, $\beta_{\text{necessary}} = 0.8$ and $\beta_{\text{artificial}} = 0.36$. Denote k the total iteration counter. The algorithm restarts if one of three conditions holds:

(i) (Sufficient decay in KKT error)

$$\mathrm{KKT}_{\omega^n}(z_c^{n,t+1}) \leq \beta_{\mathrm{sufficient}} \mathrm{KKT}_{\omega^n}(z^{n,0}) \ ,$$

(ii) (Necessary decay + no local progress in KKT error)

$$\mathrm{KKT}_{\omega^n}(z_c^{n,t+1}) \leq \beta_{\mathrm{necessary}} \mathrm{KKT}_{\omega^n}(z^{n,0}) \quad \text{and} \quad \mathrm{KKT}_{\omega^n}(z_c^{n,t+1}) > \mathrm{KKT}_{\omega^n}(z^{n,t}_c) \ ,$$

(iii) (Long inner loop)

$$t \geq \beta_{\text{artificial}} k$$
.

Finally, we highlight that the only difference between the restart scheme here and that in [2, 4] is a different choice of metric: we utilize KKT residual, while [2, 4] utilized normalized duality gap.

4 Numerical experiments

In this section, we study the numerical performance of cuPDLP.jl. In particular, we compare cuPDLP.jl with the CPU implementations of PDLP [2] and three methods implemented in the commercial solver Gurobi [41], i.e., primal simplex method, dual simplex method, and interior-point method. Section 4.1 describes the setup of the experiments. Section 4.2 presents the numerical results on LP relaxations of instances from MIPLIB 2017 collection [19]. More specifically, Section 4.2.1 investigates the performance of cuPDLP.jl and Gurobi, and Section 4.2.2 compares cuPDLP.jl with different versions of PDLP, namely Julia implemented FirstOrderLp.jl and C++ implemented PDLP with single thread and multiple threads. Section 4.3 discusses the performances of different solvers on Mittelmann's LP benchmark set [37]. We present additional figures and discussions on the results in Appendix C.

4.1 Experimental setup

Benchmark datasets. We use two LP benchmark datasets in the numerical experiments, MIP Relaxations, which contain 383 instances curated from root-node LP relaxation of mixed-integer programming problems from MIPLIB 2017 collection [19] (see Section 4.2), and 49 LP instances from the Mittelmann's LP benchmark dataset [37] (see Section 4.3).

In particular, MIPLIB 2017 [19] is a collection of mixed-integer linear programming problems. We utilize the root-node LP relaxation of instances in MIPLIB as the LP benchmark set. 383 instances are selected from MIPLIB 2017 to construct MIP Relaxations based on the following criteria (similar selection criteria are used in the experiments of CPU-based PDLP [2]):

- Not tagged as numerically unstable
- Not tagged as infeasible
- Not tagged as having indicator constraints
- Finite optimal objective (if known)
- The constraint matrix has a number of nonzeros greater than 100,000
- Zero is not an optimal solution to the LP relaxation.

MIP Relaxations is further split into three classes based on the number of nonzeros (nnz) in the constraint matrix, as shown in Table 2.

	\mathbf{Small}	Medium	Large
Number of nonzeros	100K - 1M	1M - 10M	>10M
Number of instances	269	94	20

Table 2: Scales of instances in MIP Relaxations.

Mittelmann's LP is a classic dataset to benchmark LP solvers. We utilize the 49 public instances from the dataset. The number of nonzeros of the instances spans from 100 thousand to 90 million.

Software. cuPDLP.jl is implemented in an open-source Julia [7] module. cuPDLP.jl utilizes CUDA.jl [6] as the interface for working with NVIDIA CUDA GPUs using Julia. We compare cuPDLP.jl with five LP solvers: the Julia implementation of PDLP (FirstOrderLp.jl), the C++ implementation of PDLP with single thread and multiple threads (wrapped in Google OR-Tools³), the primal simplex, dual simplex and barrier method implemented in Gurobi. Crossover is disabled, and 16 threads are used in experiments of Gurobi. The running time of cuPDLP.jl and FirstOrderLp.jl is measured after pre-compilation in Julia.

Computing environment. We use NVIDIA H100-PCIe-80GB GPU, with CUDA 12.3, for running cuPDLP.jl, and we use Intel Xeon Gold 6248R CPU 3.00GHz with 160GB RAM and 16 threads for running CPU-based solvers. The experiments are performed in Julia 1.9.2 and Gurobi 11.0 (released in November 2023). Table 3 compares the theoretical peak of double-precision FLOPS (floating-point operations per second) of the CPU and GPU used in the experiments.

	CPU (16 threads)	GPU
Processor	Intel Xeon Gold 6248R CPU 3.00GHz ⁴	NVIDIA H100-PCIe ⁵
Theoretical peak (FP64)	256 GFLOPS	26 TFLOPS
Maximum memory bandwidth	137.48 GB/sec	2 TB/sec

Table 3: Comparison of CPU and GPU specifications.

Initialization. Both PDLP and cuPDLP.jl uses all-zero vectors as the initial starting points.

Optimality termination criteria. PDLP and cuPDLP.jl terminate when the relative KKT error is no greater than the termination tolerance $\epsilon \in (0, \infty)$:

$$\begin{aligned} |q^{\top}y + l^{\top}\lambda^{+} - u^{\top}\lambda^{-} - c^{\top}x| &\leq \epsilon (1 + |q^{\top}y + l^{\top}\lambda^{+} - u^{\top}\lambda^{-}| + |c^{\top}x|) \\ \left\| \begin{pmatrix} Ax - b \\ [h - Gx]^{+} \end{pmatrix} \right\|_{2} &\leq \epsilon (1 + \|q\|_{2}) \\ \|c - K^{\top}y - \lambda\|_{2} &\leq \epsilon (1 + \|c\|_{2}) \ . \end{aligned}$$

The termination criteria are checked for the original LP instance, not the preconditioned ones, so that the termination is not impacted by the preconditioning. We use $\epsilon = 10^{-4}$ for moderately accurate solutions and $\epsilon = 10^{-8}$ for high-quality solutions. We also set 10^{-4} and 10^{-8} tolerances for parameters FeasibilityTol, OptimalityTol and BarConvTol (for barrier methods) of Gurobi⁶.

Time limit. In Section 4.2, we impose a time limit of 3600 seconds on instances with small-sized and medium-sized instances and a time limit of 18000 seconds for large instances. For Section 4.3, we impose 15000 seconds as the time limit in Mittelmann's benchmark [37].

³Specifically, PDLP in master branch of Google OR-Tools is used in our experiments.

⁴See processor specifications.

 $^{^5}$ H100 has 114 SMs and 7296 FP64 cores. See H100 datasheet and H100 whitepaper for more a detailed description of H100 GPU.

⁶We comment that Gurobi barrier, Gurobi simplex and PDLP all use different termination criteria, so it can never be a fully "fair" comparison among different types of algorithms.

Shifted geometric mean. We report the shifted geometric mean of solve time to measure the performance of solvers on a certain collection of problems. More precisely, shifted geometric mean is defined as $(\prod_{i=1}^{n} (t_i + \Delta))^{1/n} - \Delta$ where t_i is the solve time for the *i*-th instance. We shift by $\Delta = 10$ and denote it SGM10. If the instance is unsolved, the solve time is always set to the corresponding time limit.

Presolves. Presolve is a step that can simplify the LP problem before solving it. It involves detecting inconsistent bounds, removing empty rows and columns from the constraint matrix, eliminating variables with equal lower and upper bounds, detecting duplicate rows, tightening bounds, etc. This step is independent of the algorithm solving the instances. In the next section, we present the comparison with Gurobi on two sets of instances, namely, the original instances without any presolve step, as well as the instances after Gurobi presolve step.

4.2 MIP Relaxations

In this section, we compare cuPDLP.jl with the commercial LP solver Gurobi and CPU implementations of PDLP on MIP Relaxations respectively. The main message is that GPU-implemented cuPDLP.jl exhibits significant speedup over its CPU-implemented counterparts, and its numerical performance is on par with Gurobi.

4.2.1 cuPDLP.jl versus Gurobi

	Small (269) (1-hour limit)		Medium (94) (1-hour limit)		Large (20) (5-hour limit)		Total (383)	
	Count	Time	Count	Time	Count	\mathbf{Time}	Count	\mathbf{Time}
cuPDLP.jl	266	8.61	92	14.80	19	111.19	377	12.02
Primal simplex (Gurobi)	268	12.56	69	188.81	11	3145.49	348	39.81
Dual simplex (Gurobi)	268	8.75	84	66.67	15	591.63	367	21.75
Barrier (Gurobi)	268	5.30	88	45.01	18	415.78	374	14.92

Table 4: Solve time in seconds and SGM10 of different solvers on instances of MIP Relaxations with tolerance 10^{-4} : cuPDLP.jl versus Gurobi without presolve.

	Small (269)			Medium (94)		Large (20)		Total (383)	
	(1-hour	limit)	(1-hour	(1-hour limit)		· limit)	· /		
	Count	Time	Count	\mathbf{Time}	Count	\mathbf{Time}	Count	\mathbf{Time}	
cuPDLP.jl	269	5.35	93	10.31	19	33.93	381	7.37	
Primal simplex (Gurobi)	269	5.67	71	121.23	19	297.59	359	20.84	
Dual simplex (Gurobi)	268	4.17	86	37.56	19	179.49	373	11.84	
Barrier (Gurobi)	269	1.21	94	15.32	20	30.70	383	4.65	

Table 5: Solve time in seconds and SGM10 of different solvers on instances of MIP Relaxations with tolerance 10^{-4} : cuPDLP.jl versus Gurobi with presolve.

Table 4-7 present a comparison between cuPDLP.jl and the commercial LP solver Gurobi. Particularly, Table 4 and Table 5 present moderate accuracy results (i.e., 10^{-4} relative KKT error),

⁷Gurobi primal and dual simplex methods may perform better for obtaining high-accuracy solution than medium-accuracy solution. This is a known effect due to the different trajectory paths when setting different tolerance levels.

	Small (269) (1-hour limit)		Medium (94) (1-hour limit)		Large (20) (5-hour limit)		Total (383)	
	Count	Time	Count	Time	Count	Time	Count	\mathbf{Time}
cuPDLP.jl	261	23.47	86	40.69	16	421.40	363	32.35
Primal simplex (Gurobi)	268	12.43	74	157.59	13	2180.23	355	36.68
Dual simplex (Gurobi)	268	8.00	83	59.93	15	687.17	366	20.40
Barrier (Gurobi)	267	6.24	88	48.62	18	438.69	373	16.46

Table 6: Solve time in seconds and SGM10 of different solvers on instances of MIP Relaxations with tolerance 10^{-8} : cuPDLP.il versus Gurobi without presolve.⁷

	Small (269) (1-hour limit)		Medium (94) (1-hour limit)		Large (20) (5-hour limit)		Total (383)	
	Count	Time	Count	Time	Count	Time	Count	\mathbf{Time}
cuPDLP.jl	264	17.53	90	30.05	19	81.07	373	22.13
Primal simplex (Gurobi)	269	5.19	75	100.03	18	171.72	362	18.11
Dual simplex (Gurobi)	268	3.53	89	27.17	19	121.94	376	9.53
Barrier (Gurobi)	269	1.34	94	16.85	20	33.48	383	5.03

Table 7: Solve time in seconds and SGM10 of different solvers on instances of MIP Relaxations with tolerance 10^{-8} : cuPDLP.jl versus Gurobi with presolve.

while Table 6 and Table 7 present the high accuracy results (i.e., 10^{-8} relative KKT error); Table 4 and Table 6 present the results on solving the original problems, while Table 5 and Table 7 present the results on LP instances after Gurobi presolve. The tables yield several noteworthy observations:

- With Gurobi presolve, cuPDLP.jl can solve 99.5% instances to medium accuracy and 97.4% instances to high accuracy within the time limit, demonstrating its reliability for solving real-world LP.
- In the case of moderate accuracy ($\epsilon = 10^{-4}$), cuPDLP.jl exhibits comparable performance to Gurobi in terms of solved count and solve time, regardless of whether to use presolve. For medium-sized and large-sized instances without presolve, cuPDLP.jl establishes an advantage over Gurobi, achieving a 3x speed-up on medium problems with 4 more instances solved and a 3.7x speed-up on large instances with one additional solved instance, respectively. With Gurobi presolve, cuPDLP.jl is able to solve 381 out of 383 instances and the solve time is comparable to the best of the three Gurobi methods (i.e., barrier methods).
- In the case of high accuracy ($\epsilon = 10^{-8}$), cuPDLP.jl has comparable performance to Gurobi primal and dual simplex method, though it is inferior to the Gurobi barrier method.
- Gurobi presolve can improve the performance of all Gurobi methods as well as the performance of cuPDLP.jl. The effect of presolve is more significant for Gurobi methods. This is expected because Gurobi presolve is designed to speed up Gurobi methods.

To summarize, these observations affirm that cuPDLP.jl attains comparable performance to Gurobi in MIP Relaxations benchmark dataset. This demonstrates that a first-order-method-based LP

solver on GPU can be on par with a strong implementation of simplex and barrier methods, even in obtaining high-accuracy solutions.

Figure 3 and Figure 4 show the number of solved instances of cuPDLP.jl and three methods in Gurobi on MIP Relaxations in a given time. The y-axes display the fraction of solved instances, and the x-axes display the wall-clock time in seconds. As shown in the left panel, when seeking solutions with moderate accuracy ($\epsilon = 10^{-4}$), cuPDLP.jl has comparable performances with Gurobi barrier after 10 seconds. It eventually has better performances on MIP Relaxation than all three methods of Gurobi without presolve. In addition, we can see the performance of cuPDLP.jl for high-quality solution ($\epsilon = 10^{-8}$), as shown in the right panel, is still comparable to Gurobi. An observation is that the number of instances Gurobi can solve for a given running time does not differ much for moderate and high accuracy; conversely, such difference is more apparent for cuPDLP.jl. This is a feature of a first-order-method-based solver. Another interesting fact is that Gurobi solves about 35% of instances within one second, which is exactly the power of Gurobi. On the other hand, cuPDLP.jl has a computational overhead of around one second due to the GPU kernel launch time.

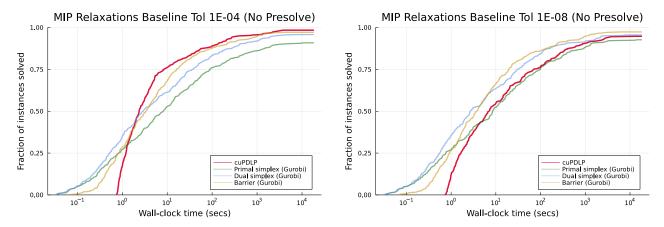


Figure 3: Number of instances solved for MIP Relaxations under moderate accuracy (left) and high accuracy (right): cuPDLP.jl versus Gurobi without presolve.

4.2.2 cuPDLP.jl versus PDLP

	Small (269) (1-hour limit)			Medium (94) (1-hour limit)		Large (20) (5-hour limit)		Total (383)	
	Count	\mathbf{Time}	Count	\mathbf{Time}	Count	\mathbf{Time}	Count	\mathbf{Time}	
cuPDLP.jl	266	8.61	92	14.80	19	111.19	377	12.02	
${f FirstOrderLp.jl}$	253	35.94	82	155.67	12	2002.21	347	66.67	
PDLP (1 thread)	256	22.69	85	98.38	15	1622.91	356	43.81	
PDLP (4 threads)	260	24.03	91	42.94	15	736.20	366	34.57	
PDLP (16 threads)	238	104.72	84	142.79	15	946.24	337	127.49	

Table 8: Solve time in seconds and SGM10 of different solvers on instances of MIP Relaxations with tolerance 10^{-4} : cuPDLP.jl versus PDLP.

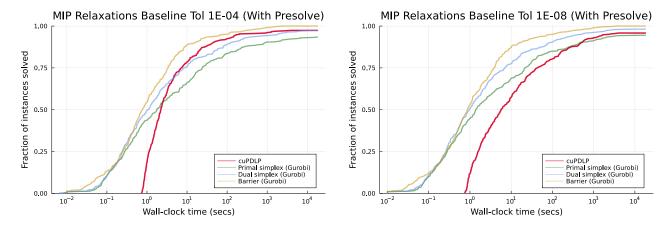


Figure 4: Number of instances solved for MIP Relaxations under moderate accuracy (left) and high accuracy (right): cuPDLP.jl versus Gurobi with presolve.

	Small (269) (1-hour limit)			Medium (94) (1-hour limit)		Large (20) (5-hour limit)		Total (383)	
	Count	\mathbf{Time}	Count	\mathbf{Time}	Count	\mathbf{Time}	Count	\mathbf{Time}	
cuPDLP.jl	261	23.47	86	40.69	16	421.40	363	32.35	
${f FirstOrderLp.jl}$	235	91.14	68	389.34	9	3552.50	312	160.63	
PDLP (1 thread)	250	49.31	73	259.04	12	3818.42	335	96.86	
PDLP (4 threads)	245	54.19	81	136.16	14	1789.54	340	83.49	
PDLP (16 threads)	214	248.34	69	403.17	14	2475.57	297	316.27	

Table 9: Solve time in seconds and SGM10 of different solvers on instances of MIP Relaxations with tolerance 10^{-8} : cuPDLP.jl versus PDLP.

Tables 8 and 9 present the performance comparison of cuPDLP.jl and its CPU implementations on MIP Relaxations with tolerances of 10^{-4} and 10^{-8} , respectively. The two tables demonstrate that GPU can significantly speed up PDLP, in particular for large instances:

- For moderate accuracy (Table 8), cuPDLP.jl demonstrates a 4x speed-up for small instances, a 10x speed-up for medium instances, and a 20x speed-up for large instances, compared to FirstOrderLp.jl, a CPU implementation of PDLP in Julia. When comparing cuPDLP.jl with the more delicate C++ implementation PDLP with multithreading support, it still exhibits a significant speedup for all sizes of problems. The significant speed-up can also be observed for high accuracy (Table 9).
- In terms of solved count, cuPDLP.jl solves significantly more instances regardless of the scales. In particular, comparing with FirstOrderLp.jl under tolerance $\epsilon = 10^{-4}$, cuPDLP.jl solves 13 more small-sized instances, 10 more medium-sized instances and 7 more large problems, with in total 30 more instances solved. Compared to PDLP with the best of 1 thread, 4 threads or 16 threads, cuPDLP.jl can solve 6 more small-sized instances and 4 more large-sized instances. The improvement is also remarkable when looking at results for high accuracy $\epsilon = 10^{-8}$.

In summary, the GPU-implemented cuPDLP.jl consistently outperforms the CPU-implemented

PDLP in numerical performance on MIP Relaxations, with cuPDLP.jl demonstrating even more pronounced advantages on instances of medium to large scale.

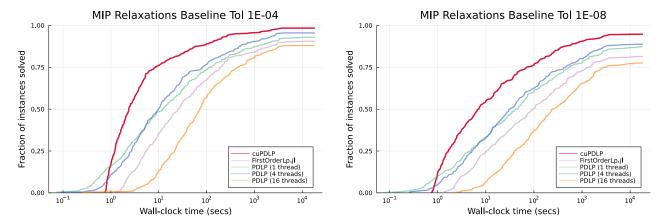


Figure 5: Number of instances solved for MIP Relaxations under moderate accuracy (left) and high accuracy (right): cuPDLP.jl versus PDLP.

Similar to Figure 3 and 4, Figure 5 demonstrates the number of solved instances of cuPDLP.jl and three CPU implementations of PDLP on MIP Relaxations in a given time. As shown in both panels, when seeking solutions with moderate accuracy ($\epsilon = 10^{-4}$) and high accuracy ($\epsilon = 10^{-8}$), cuPDLP.jl has a clear superior performance to all CPU versions of PDLP. Moreover, it is notable that cuPDLP.jl has a computational overhead of around one second due to the GPU kernel launch time.

4.3 Mittelmann's LP benchmark set

We further compare cuPDLP.jl with PDLP and Gurobi on Mittelmann's LP benchmark set. Mittelmann's LP benchmark set [37] is a standard LP benchmark to test the numerical performances of different LP solvers and includes 49 LP instances in our experiments.

Table 10, Table 11 and Table 12 summarize the results on Mittelmann's LP benchmark set. The observations are consistent with the findings on MIP Relaxations as discussed in Section 4.2. Consider the moderate accuracy. cuPDLP.jl can solve more instances against all versions of PDLP and have comparable performance to Gurobi. When solving for high-accuracy solutions, the performance of cuPDLP.jl still performs inferior to Gurobi barrier but comparable to primal and dual simplex. This again demonstrates the robust performance of GPU-implemented cuPDLP.jl.

5 Conclusion

In this paper, we present cuPDLP.jl, a GPU implementation of restarted PDHG for solving LP in Julia. The numerical experiments demonstrate that the prototype GPU implementation cuPDLP.jl can have comparable performance with commercial solvers like Gurobi and superior performance on large instances. This sheds light on using GPU to develop high-performance optimization solvers.

	Tol 1	LE-04	Tol 1E-08		
	Count	\mathbf{Time}	Count	\mathbf{Time}	
cuPDLP.jl	44	71.26	40	231.91	
Primal simplex (Gurobi)	35	1937.78	34	1715.62	
Dual simplex (Gurobi)	36	1201.48	37	1116.68	
Barrier (Gurobi)	45	108.29	44	127.58	

Table 10: Solve time in seconds and SGM10 of different solvers on instances of Mittelmann's LP benchmark set with tolerance 10^{-4} and 10^{-8} without presolve.

	Tol 1	E-04	Tol 1E-08		
	Count	\mathbf{Time}	Count	\mathbf{Time}	
cuPDLP.jl	47	52.44	44	167.15	
Primal simplex (Gurobi)	42	606.03	39	499.09	
Dual simplex (Gurobi)	40	290.69	42	212.00	
Barrier (Gurobi)	48	24.16	47	34.18	

Table 11: Solve time in seconds and SGM10 of different solvers on instances of Mittelmann's LP benchmark set with tolerance 10^{-4} and 10^{-8} with presolve.

	Tol 1	E-04	Tol 1	LE-08
	Count	\mathbf{Time}	Count	\mathbf{Time}
cuPDLP.jl	44	71.26	40	231.91
${f FirstOrderLp.jl}$	34	917.49	25	2504.79
PDLP (1 thread)	39	586.50	31	1661.20
PDLP (4 threads)	40	302.54	34	930.41
PDLP (16 threads)	39	776.22	29	2171.95

Table 12: Solve time in seconds and SGM10 of different solvers on instances of Mittelmann's LP benchmark set with tolerance 10^{-4} and 10^{-8} .

Acknowledgement

The authors would like to thank Azam Asl and Miles Lubin for the early discussions on the GPU implementation of PDLP, and David Applegate for his encouragement and support.

References

- [1] Ayan Acharya, Siyuan Gao, Borja Ocejo, Kinjal Basu, Ankan Saha, Keerthi Selvaraj, Rahul Mazumdar, Parag Agrawal, and Aman Gupta, *Promoting inactive members in edge-building marketplace*, Companion Proceedings of the ACM Web Conference 2023, 2023, pp. 945–949.
- [2] David Applegate, Mateo Díaz, Oliver Hinder, Haihao Lu, Miles Lubin, Brendan O'Donoghue, and Warren Schudy, *Practical large-scale linear programming using primal-dual hybrid gradient*, Advances in Neural Information Processing Systems **34** (2021), 20243–20257.

- [3] David Applegate, Mateo Díaz, Haihao Lu, and Miles Lubin, Infeasibility detection with primaldual hybrid gradient for large-scale linear programming, arXiv preprint arXiv:2102.04592 (2021).
- [4] David Applegate, Oliver Hinder, Haihao Lu, and Miles Lubin, Faster first-order primal-dual methods for linear programming using restarts and sharpness, Mathematical Programming 201 (2023), no. 1-2, 133–184.
- [5] Kinjal Basu, Amol Ghoting, Rahul Mazumder, and Yao Pan, *Eclipse: An extreme-scale linear program solver for web-applications*, International Conference on Machine Learning, PMLR, 2020, pp. 704–714.
- [6] Tim Besard, Christophe Foket, and Bjorn De Sutter, Effective extensible programming: unleashing julia on gpus, IEEE Transactions on Parallel and Distributed Systems 30 (2018), no. 4, 827–841.
- [7] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B Shah, *Julia: A fresh approach to numerical computing*, SIAM review **59** (2017), no. 1, 65–98.
- [8] Antonin Chambolle and Thomas Pock, A first-order primal-dual algorithm for convex problems with applications to imaging, Journal of mathematical imaging and vision 40 (2011), 120–145.
- [9] _____, On the ergodic convergence rates of a first-order primal-dual algorithm, Mathematical Programming 159 (2016), no. 1-2, 253–287.
- [10] Abraham Charnes, William W Cooper, and Merton H Miller, Application of linear programming to financial budgeting and the costing of funds, The Journal of Business 32 (1959), no. 1, 20–46.
- [11] Laurent Condat, A primal-dual splitting method for convex optimization involving lipschitzian, proximable and linear composite terms, Journal of optimization theory and applications 158 (2013), no. 2, 460–479.
- [12] Munther Dahleh and Ignacio Diaz-Bobillo, Control of uncertain systems: a linear programming approach, Prentice-Hall, Inc., 1994.
- [13] George Bernard Dantzig, *Linear programming and extensions*, vol. 48, Princeton university press, 1998.
- [14] Jerome Delson and Mohammad Shahidehpour, Linear programming applications to power system economics, planning and operations, IEEE Transactions on Power Systems 7 (1992), no. 3, 1155–1163.
- [15] Qi Deng, Qing Feng, Wenzhi Gao, Dongdong Ge, Bo Jiang, Yuntian Jiang, Jingsong Liu, Tianhao Liu, Chenyu Xue, Yinyu Ye, et al., New developments of admm-based interior point methods for linear programming and conic programming, arXiv preprint arXiv:2209.01793 (2022).
- [16] Ernie Esser, Xiaoqun Zhang, and Tony F Chan, A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science, SIAM Journal on Imaging Sciences 3 (2010), no. 4, 1015–1046.

- [17] Olivier Fercoq, Quadratic error bound of the smoothed gap and the restarted averaged primaldual hybrid gradient, (2021).
- [18] Dongdong Ge, Qi Huangfu, Zizhuo Wang, Jian Wu, and Yinyu Ye, Cardinal optimizer (copt) user guide, arXiv preprint arXiv:2208.14314 (2022).
- [19] Ambros Gleixner, Gregor Hendel, Gerald Gamrath, Tobias Achterberg, Michael Bastubbe, Timo Berthold, Philipp Christophel, Kati Jarck, Thorsten Koch, Jeff Linderoth, et al., Miplib 2017: data-driven compilation of the 6th mixed-integer programming library, Mathematical Programming Computation 13 (2021), no. 3, 443–490.
- [20] Greg Glockner, Parallel and distributed optimization with gurobi optimizer, https://www.gurobi.com/events/parallel-and-distributed-optimization-with-gurobi/, 2015.
- [21] _____, Does gurobi support gpus?, https://support.gurobi.com/hc/en-us/articles/ 360012237852-Does-Gurobi-support-GPUs-, 2023.
- [22] Peter Hazell and Pasquale Scandizzo, Competitive demand structures under risk in agricultural linear programming models, American Journal of Agricultural Economics **56** (1974), no. 2, 235–244.
- [23] Bingsheng He and Xiaoming Yuan, Convergence analysis of primal-dual algorithms for a saddle-point problem: from contraction perspective, SIAM Journal on Imaging Sciences 5 (2012), no. 1, 119–149.
- [24] Oliver Hinder, Worst-case analysis of restarted primal-dual hybrid gradient on totally unimodular linear programs, arXiv preprint arXiv:2309.03988 (2023).
- [25] Alan J Hoffman, On approximate solutions of systems of linear inequalities, Journal of Research of the National Bureau of Standards 49 (1952), 263–265.
- [26] Qi Huangfu and JA Julian Hall, Parallelizing the dual revised simplex method, Mathematical Programming Computation 10 (2018), no. 1, 119–142.
- [27] Narendra Karmarkar, A new polynomial-time algorithm for linear programming, Proceedings of the sixteenth annual ACM symposium on Theory of computing, 1984, pp. 302–311.
- [28] Tianyi Lin, Shiqian Ma, Yinyu Ye, and Shuzhong Zhang, An admm-based interior-point method for large-scale linear programming, Optimization Methods and Software **36** (2021), no. 2-3, 389–424.
- [29] Qian Liu and Garrett Van Ryzin, On the choice-based linear programming model for network revenue management, Manufacturing & Service Operations Management 10 (2008), no. 2, 288–310.
- [30] Haihao Lu and Jinwen Yang, On the infimal sub-differential size of primal-dual hybrid gradient method, arXiv preprint arXiv:2206.12061 (2022).
- [31] _____, On a unified and simplified proof for the ergodic convergence rates of ppm, pdhg and admm, arXiv preprint arXiv:2305.02165 (2023).
- [32] _____, On the geometry and refined rate of primal-dual hybrid gradient for linear programming, arXiv preprint arXiv:2307.03664 (2023).

- [33] ______, A practical and optimal first-order method for large-scale convex quadratic programming, arXiv preprint arXiv:2311.07710 (2023).
- [34] CPLEX User's Manual, *Ibm ilog cplex optimization studio*, Version **12** (1987), no. 1987-2018, 1.
- [35] Gioni Mexi, Mathieu Besançon, Suresh Bolusani, Antonia Chmiela, Ambros Gleixner, and Alexander Hoen, Scylla: a matrix-free fix-propagate-and-project heuristic for mixed-integer optimization, arXiv preprint arXiv:2307.03466 (2023).
- [36] Vahab Mirrokni, Google research, 2022 & beyond: Algorithmic advances, https://ai.googleblog.com/2023/02/google-research-2022-beyond-algorithmic.html, 2023-02-10.
- [37] Hans D Mittelmann, Decision tree for optimization software., https://plato.asu.edu/bench.html, 2023.
- [38] NVIDIA, Cuda c++ best practices guide, (2023).
- [39] Brendan O'Donoghue, Operator splitting for a homogeneous embedding of the linear complementarity problem, SIAM Journal on Optimization 31 (2021), no. 3, 1999–2023.
- [40] Brendan O'Donoghue, Eric Chu, Neal Parikh, and Stephen Boyd, *Conic optimization via operator splitting and homogeneous self-dual embedding*, Journal of Optimization Theory and Applications **169** (2016), no. 3, 1042–1068.
- [41] Gurobi Optimization et al., Gurobi optimizer reference manual, 2023.
- [42] Javier Pena, Juan C Vera, and Luis F Zuluaga, New characterizations of hoffman constants for systems of linear constraints, Mathematical Programming 187 (2021), 79–109.
- [43] Thomas Pock and Antonin Chambolle, Diagonal preconditioning for first order primal-dual algorithms in convex optimization, 2011 International Conference on Computer Vision, IEEE, 2011, pp. 1762–1769.
- [44] Rohan Ramanath, S Sathiya Keerthi, Yao Pan, Konstantin Salomatin, and Kinjal Basu, Efficient vertex-oriented polytopic projection for web-scale applications, Proceedings of the AAAI Conference on Artificial Intelligence, vol. 36, 2022, pp. 3821–3829.
- [45] Daniel Ruiz, A scaling algorithm to equilibrate both rows and columns norms in matrices, Tech. report, CM-P00040415, 2001.
- [46] Sungho Shin, François Pacaud, and Mihai Anitescu, Accelerating optimal power flow with gpus: Simd abstraction of nonlinear programs and condensed-space interior-point methods, arXiv preprint arXiv:2307.16830 (2023).
- [47] Bartolomeo Stellato, Goran Banjac, Paul Goulart, Alberto Bemporad, and Stephen Boyd, Osqp: An operator splitting solver for quadratic programs, Mathematical Programming Computation 12 (2020), no. 4, 637–672.
- [48] Kasia Świrydowicz, Eric Darve, Wesley Jones, Jonathan Maack, Shaked Regev, Michael A Saunders, Stephen J Thomas, and Slaven Peleš, *Linear solvers for power grid optimization problems: a review of gpu-accelerated linear solvers*, Parallel Computing **111** (2022), 102870.

- [49] Peng Zhou and Beng Wah Ang, Linear programming models for measuring economy-wide energy efficiency performance, Energy Policy 36 (2008), no. 8, 2911–2916.
- [50] Mingqiang Zhu and Tony Chan, An efficient primal-dual hybrid gradient algorithm for total variation image restoration, UCLA Cam Report 34 (2008), 8–34.

A Theoretical guarantees of restarted PDHG with KKT error for LP

Consider LP of the standard form:

$$\min_{x \in \mathbb{R}^n} \ c^{\top} x \quad \text{s.t. } Ax = b, \ x \ge 0 \ . \tag{6}$$

The KKT error at solution z = (x, y) is defined as the norm of the violation of KKT system of (6)

$$KKT(z) = KKT(x, y) = \left\| \begin{pmatrix} Ax - b \\ [-x]^+ \\ [A^\top y - c]^+ \\ [c^\top x - b^\top y]^+ \end{pmatrix} \right\|_2,$$

where $[x]^+ = [(x_1, ..., x_n)]^+ = (\max\{x_1, 0\}, ..., \max\{x_n, 0\})$ is the positive part of the vector x.

Algorithm 3: Restarted PDHG for solving (6)

```
Input: Initial point (x^0, y^0), outer loop counter n \leftarrow 0. Step-sizes \tau, \sigma. Restart decay \beta \in (0, 1).
```

```
1 repeat
                initialize the inner loop counter k \leftarrow 0;
   \mathbf{2}
  3
                       x^{n,k+1} \leftarrow \operatorname{proj}_{\mathbb{R}^n_+}(x^{n,k} - \tau(c - A^{\mathsf{T}}y^{n,k}));
                y^{n,k+1} \leftarrow y^{n,k} + \sigma(b - A(2x^{n,k+1} - x^{n,k}));
\bar{x}^{n,k+1} \leftarrow k = n \cdot k
   4
   \mathbf{5}
               \bar{x}^{n,k+1} \leftarrow \frac{k}{k+1} \bar{x}^{n,k} + \frac{1}{k+1} x^{n,k+1}; \\ \bar{y}^{n,k+1} \leftarrow \frac{k}{k+1} \bar{y}^{n,k} + \frac{1}{k+1} y^{n,k+1}; \\ \mathbf{until} \ \mathrm{KKT}(\bar{x}^{n,k+1}, \bar{y}^{n,k+1}) \leq \beta \mathrm{KKT}(x^{n,0}, y^{n,0}) \ \mathrm{holds};
   6
   7
   8
                initialize the initial solution (x^{n+1,0}, y^{n+1,0}) \leftarrow (\bar{x}^{n,k+1}, \bar{y}^{n,k+1});
  9
                n \leftarrow n + 1;
10
11 until (x^{n+1,0}, y^{n+1,0}) convergence;
```

Consider restarted PDHG with KKT error as restarting criteria (Algorithm 3). Denote $z^{n,k} = (x^{n,k}, y^{n,k})$ (and $\bar{z}^{n,k} = (\bar{x}^{n,k}, \bar{y}^{n,k})$, correspondingly) as the primal-dual solution pair at the *n*-th other iteration and the *k*-th inner iteration. Without loss of generality, assume the primal and dual stepsizes are equal, i.e., $\tau = \sigma =: s$ (otherwise we can rescale the primal and the dual problem so that they share the same step-size; a similar proof strategy is used in [4]). Denote $P_s = \begin{pmatrix} I & sA^\top \\ sA & I \end{pmatrix}$ and define $\|v\|_{P_s}^2 = \langle v, P_s v \rangle$. Denote $\|\cdot\|_2$ the Euclidean norm and $\operatorname{dist}_2(w, U) = \min_{u \in U} \|w - u\|_2$ the distance between point w and set U.

First, it is straight-forward to see that KKT is a sharp function, i.e.,

Proposition 1 ([25, 42]). There exists a constant $\alpha > 0$ such that for any z = (x, y), it holds that

$$\alpha \operatorname{dist}_2(z, \mathcal{Z}^*) \leq \operatorname{KKT}(z)$$
.

The next theorem proves the linear convergence rate of Algorithm 3. Notice that this is essentially the same linear convergence rate as that in [4] up to a constant.

Theorem 1. Consider $\{z^{n,k}\}$ the iterates of Algorithm 3 for solving (2). Suppose stepsize $s \leq \frac{1}{2||A||_2}$. Then it holds that

- (i) There exists an R > 0 such that for any n and k, $||z^{n,k}||_2 \le R$.
- (ii) The restart length τ^n is upper bounded by

$$\tau^n \le \frac{8\sqrt{2}\sqrt{1+R^2}}{\alpha s\beta} \ .$$

(iii) The distance to optimal solution set decays linearly

$$\operatorname{dist}(z^{n,0}, \mathcal{Z}^*) \leq \beta^{n+1} \frac{1}{\alpha} \operatorname{KKT}(z^{0,0})$$
.

Remark 1. Let stepsize $s = \frac{C}{\|A\|_2}$ where constant $C \leq \frac{1}{2}$. The iteration complexity of Algorithm 3 to find an ϵ -solution equals $\widetilde{O}\left(\frac{\|A\|_2}{\alpha}\log\frac{1}{\epsilon}\right)$, which is optimal (upto log terms) in the sense of matching the lower complexity bounds [4, Corollary 1].

Proof. (i) Let $R = \sqrt{2}||z^{0,0} - z^*||_{P_s} + ||z^*||_2$. By non-expansiveness of PDHG [9], we have

$$\|z^{n,k}-z^*\|_{P_s} \leq \|z^{n,0}-z^*\|_{P_s} = \left\|\frac{1}{\tau^{n-1}}\sum_{k=1}^{\tau^{n-1}}z^{n-1,k}-z^*\right\|_{P_s} \leq \frac{1}{\tau^{n-1}}\sum_{k=1}^{\tau^{n-1}}\left\|z^{n-1,k}-z^*\right\|_{P_s} \leq \|z^{n-1,0}-z^*\|_{P_s} \leq \ldots \leq \|z^{0,0}-z^*\|_{P_s}.$$

Thus we achieve

$$||z^{n,k}||_2 \le ||z^{n,k} - z^*||_2 + ||z^*|| \le \sqrt{2}||z^{0,0} - z^*||_{P_s} + ||z^*||_2 = R$$
.

(ii) Suppose $\tau^n > \frac{8\sqrt{2}\sqrt{1+R^2}}{\alpha\beta} =: k^*$. Note that

$$KKT(\bar{z}^{n-1,k^*}) \leq \sqrt{1+R^2} \rho_{\|\bar{z}^{n-1,k^*}-z^{n-1,0}\|_{P_s}}(\bar{z}^{n-1,k^*}) \leq \frac{4\sqrt{1+R^2}}{sk^*} \|\bar{z}^{n-1,k^*}-z^{n-1,0}\|_{P_s}$$

$$\leq \frac{8\sqrt{1+R^2}}{sk^*} \text{dist}_{P_s}(z^{n-1,0}, \mathcal{Z}^*) \leq \frac{8\sqrt{2}\sqrt{1+R^2}}{sk^*} \text{dist}_2(z^{n-1,0}, \mathcal{Z}^*) \leq \frac{8\sqrt{2}\sqrt{1+R^2}}{\alpha sk^*} KKT(z^{n-1,0})$$

$$\leq \beta KKT(z^{n-1,0}),$$

where the first and fifth inequalities follow from [4, Lemma 4], while the second and third inequalities utilize [4, Property 3]. The fourth inequality uses $P_s \leq 2I$ which is a direct consequence of stepsize $s \leq \frac{1}{2\|A\|_2}$. The last inequality follows from the definition of $k^* = \frac{8\sqrt{2}\sqrt{1+R^2}}{\alpha\beta}$.

This implies the restart condition holds at k^* . We finish the proof by noticing that this contradicts $\tau^n > k^*$.

(iii) Note that the KKT error has linearly decay due to the adaptive restart scheme

$$KKT(z^{n,0}) = KKT(\bar{z}^{n-1,\tau^{n-1}}) \le \beta KKT(z^{n-1,0}) \le \dots \le \beta^n KKT(z^{0,0})$$
,

and thus it follows from sharpness of KKT error that

$$\operatorname{dist}_2(z^{n,0},\mathcal{Z}^*) \leq \frac{1}{\alpha} \operatorname{KKT}(z^{n,0}) \leq \frac{1}{\alpha} \beta^n \operatorname{KKT}(z^{0,0}) .$$

B Comparison with GPU-based SCS

	Count	Time
SCS-GPU (Tol=1E-04)	28	1905.00
cuPDLP.jl (Tol=1E-04)	92	14.80
cuPDLP.jl (Tol=1E-08)	86	40.69

Table 13: Solve time in seconds and SGM10 of different solvers on medium-sized instances of MIP Relaxations: cuPDLP.jl versus SCS with GPU linear system solver.

In this section, we compare the results of cuPDLP.jl with GPU-based SCS which solves the linear system in each iteration via conjugate gradient method implemented on GPU. We select medium-sized instances (94 in total) in MIP Relaxations as the test set. Table 13 presents the results of GPU-based SCS under moderate tolerance 10^{-4} and cuPDLP.jl under both 10^{-4} and 10^{-8} tolerances. The advantage of cuPDLP over SCS is quite significant: GPU-based SCS only solves 28 out of 94 instances under moderate accuracy. In comparison, cuPDLP solves 92 out of 94 instances under same tolerance and 86 out of 94 instances under even higher tolerance. This showcases cuPDLP, the LP-specialized FOM-based solver, performs better at solving LP instances than general FOM-based conic solver SCS on GPU. This result is not surprising because SCS targets a more general class of optimization problems, and the CPU implementation of PDLP has showcased superior performance than CPU-based SCS [2].

C Figures

Figure 6 and Figure 7 visualize the running time comparison between Gurobi and cuPDLP.jl over the size of the instances with scatter plots. Specifically, each dot in a scatter plot is an instance both methods can solve within the time limit. The x-axes are the number of nonzeros of a MIP Relaxations instance, and the y-axes are the running time ratio of Gurobi (primal simplex, dual simplex, and barrier, respectively) over cuPDLP.jl. Though with higher variance, we can observe a strong positive correlation between the speed-up of cuPDLP.jl against Gurobi versus the number of nonzeros of the problems. We can also observe quite a few instances where cuPDLP is more than 100 times faster than Gurobi different methods.

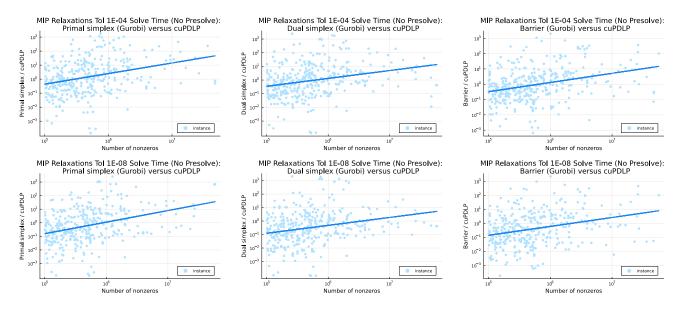


Figure 6: Ratio of Gurobi solve time over cuPDLP.jl solve time for moderate accuracy (top) and high accuracy (bottom). Presolve is not used.

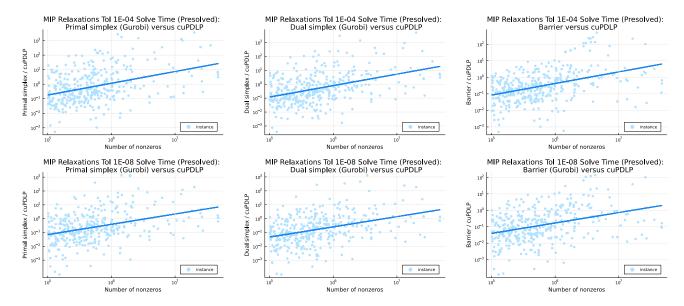


Figure 7: Ratio of Gurobi solve time over cuPDLP.jl solve time for moderate accuracy (top) and high accuracy (bottom). Presolve is used.

Furthermore, we visualize the comparison of solve time of PDLP with cuPDLP.jl in Figure 8. The y-axes are the ratio of FirstOrderLp.jl/PDLP with single thread/PDLP with 4 threads/PDLP with 16 threads solve time and cuPDLP.jl solve time, which represents the speed-up gained by cuPDLP.jl, while the x-axes is the number of nonzeros in each constraint matrix. The increasing trends in the figures imply that GPU-based cuPDLP.jl can gain more significant speed-up over CPU-based counterpart PDLP as we solve larger instances.

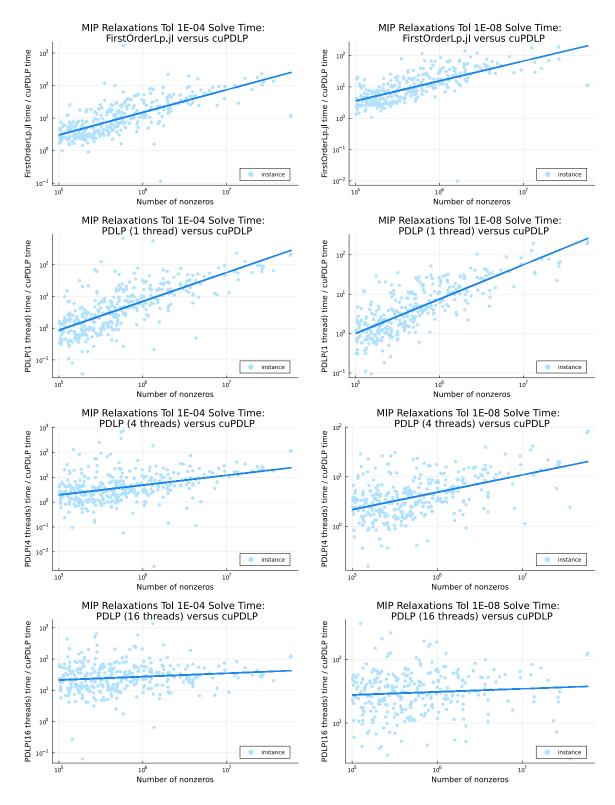


Figure 8: Ratio of PDLP solve time over cuPDLP.jl solve time for moderate accuracy (top) and high accuracy (bottom).