

DEEP DISTRIBUTED OPTIMIZATION FOR LARGE-SCALE QUADRATIC PROGRAMMING

Anonymous authors

Paper under double-blind review

ABSTRACT

Quadratic programming (QP) forms a crucial foundation in optimization, encompassing a broad spectrum of domains and serving as the basis for more advanced algorithms. Consequently, as the scale and complexity of modern applications continue to grow, the development of efficient and reliable QP algorithms becomes increasingly vital. In this context, this paper introduces a novel deep learning-aided distributed optimization architecture designed for tackling large-scale QP problems. First, we combine the state-of-the-art Operator Splitting QP (OSQP) method with a consensus approach to derive **DistributedQP**, a new method tailored for network-structured problems, with convergence guarantees to optimality. Subsequently, we unfold this optimizer into a deep learning framework, leading to **DeepDistributedQP**, which leverages learned policies to accelerate reaching to desired accuracy within a **restricted** amount of iterations. Our approach is also theoretically grounded through Probably Approximately Correct (PAC)-Bayes theory, providing generalization bounds on the expected optimality gap for unseen problems. The proposed framework, as well as its **centralized** version **DeepQP**, significantly outperform their standard optimization counterparts on a variety of tasks such as randomly generated problems, optimal control, linear regression, transportation networks and others. **Notably**, **DeepDistributedQP** demonstrates strong generalization by training on small problems and scaling to solve much larger ones (up to 50K variables and 150K constraints) using the same policy. Moreover, it achieves orders-of-magnitude improvements in wall-clock time compared to OSQP. The certifiable performance guarantees of our approach are also demonstrated, ensuring higher-quality solutions over traditional optimizers.

1 INTRODUCTION

Quadratic programming (QP) serves as a fundamental cornerstone in optimization with a wide variety of applications in machine learning (Cortes & Vapnik, 1995; Tibshirani, 1996), control and robotics (Garcia et al., 1989; Rawlings et al., 2017), signal processing (Mattingley & Boyd, 2010), finance (Cornuejols et al., 2018), and transportation networks (Mota et al., 2014) among other fields. Beyond its standalone applications, QP also acts as the core component of many advanced non-convex optimization algorithms such as sequential quadratic programming (Nocedal & Wright, 1999), trust-region methods (Conn et al., 2000), augmented Lagrangian approaches (Houska et al., 2016), mixed-integer optimization (Belotti et al., 2013), etc. For these reasons, the pursuit of more efficient QP algorithms remains an ever-evolving area of research from active set (Wolfe, 1959) and interior point methods (Nesterov & Nemirovskii, 1994) during the previous century to first-order methods such as the state-of-the-art Operator Splitting QP (OSQP) algorithm (Stellato et al., 2020).

As the scale of modern decision-making applications rapidly increases, there is an emerging interest in developing effective optimization architectures for addressing high-dimensional problems. Given the fundamental role of QP in optimization, there is a clear demand for algorithms capable of solving large-scale QPs with thousands, and potentially much more, variables and constraints. Such problems arise in diverse applications including sparse linear regression (Mateos et al., 2010) and support vector machines (Navia-Vazquez et al., 2006) with decentralized data, multi-agent control (Van Parys & Pipeleers, 2017), resource allocation (Huang et al., 2014), network flow (Mota et al., 2014), power grids (Lin et al., 2012) and image processing (Soheili & Eftekhari-Moghadam, 2020). Traditional centralized optimization algorithms are inadequate for solving such problems at

scale (see for example Fig. 1), prompting the development of distributed methods that leverage the underlying network/decentralized structure to parallelize computations. In this context, the Alternating Direction Method of Multipliers (ADMM) has gained widespread popularity as an effective approach for deriving distributed algorithms (Boyd et al., 2011; Mota et al., 2013). Nevertheless, as scale increases, such algorithms continue to face significant challenges such as their need for *meticulous tuning*, the absence of *generalization guarantees* and restrictions on the *allowed number of iterations* imposed by computational or communication limitations.

Learning-to-optimize has recently emerged as a methodology for enhancing existing optimizers or developing entirely new ones through training on sample problems (Chen et al., 2022; Amos et al., 2023). A notable approach within this paradigm is *deep unfolding*, which involves unrolling the optimizer iterations for a fixed number of steps and tuning their parameters to refine performance (Monga et al., 2021; Shlezinger et al., 2022). Our key insight is that deep unfolding is particularly well-suited for overcoming the limitations of *distributed constrained optimization*, as it can eliminate the need for extensive tuning, manage iteration restrictions and enhance generalization. However, to our best knowledge, its combination with distributed optimization has only recently been explored in Noah & Shlezinger (2024). While this framework shows promising initial results, it relies on a relatively simple setup that studies unconstrained problems, assumes local updates consisting of gradient steps, focuses solely on parameter tuning, and is not accompanied by any formal performance guarantees.

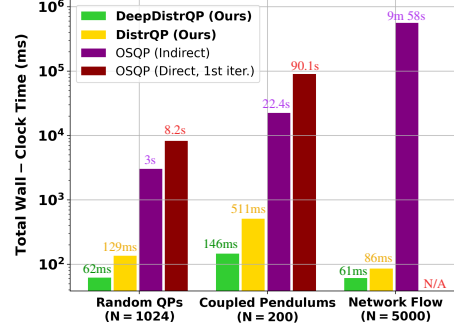


Figure 1: **Wall-clock time comparison:** DeepDistributedQP, DistributedQP (ours) and OSQP on large-scale QPs.

This paper introduces a novel deep learning-aided distributed optimization architecture for solving large-scale constrained QP problems. Our proposed approach relies on unfolding a newly introduced distributed QP algorithm as a supervised learning framework for a prescribed number of iterations. To our best knowledge, this is the first work to propose a learning-based architecture for distributed constrained optimization using ADMM, despite its widespread popularity. Our framework demonstrates remarkable performance and scalability when trained on small problems and can be effectively applied to much larger ones. Furthermore, its performance is theoretically supported by establishing guarantees based on generalization bounds from statistical learning theory. We believe that this work lays the foundation for developing learned distributed optimizers capable of handling large-scale constrained optimization problems without requiring training at such scales.

Our specific contributions can be summarized as follows:

- First, we introduce **DistributedQP**, a new decentralized method that combines the well-established OSQP solver with a consensus approach. We further prove that the algorithm is guaranteed to converge to optimality, even under varying local algorithm parameters.
- Then, we propose **DeepDistributedQP**, a deep learning-aided distributed architecture that unrolls the iterations of DistributedQP in a supervised manner, learning feedback policies for the underlying algorithm parameters. As a byproduct, we also present **DeepQP**, its centralized counterpart which corresponds to unfolding the standard OSQP solver.
- To certify the performance of the learned solver, we establish generalization guarantees on the optimality gap of the final solution of DeepDistributedQP for unseen problems using Probably Approximately Correct (PAC)-Bayes theory.
- Finally, we present an extensive experimental evaluation that validates the following:
 - For centralized QPs, DeepQP consistently outperforms OSQP requiring 1.5-3 times fewer iterations for achieving the desired accuracy.
 - DeepDistributedQP successfully scales for high-dimensional problems (up to 50K variables and 150K constraints) while being trained exclusively on much smaller ones. Furthermore, both DeepDistributedQP and DistributedQP outperform OSQP in wall-clock time by orders of magnitude as problem dimensionality increases.
 - The resulting performance bounds offer valuable guarantees on the quality of solutions produced by DeepDistributedQP for unseen problems from the same class.

2 RELATED WORK

This section provides an overview of the existing related literature from both the angles of distributed optimization and learning-to-optimize approaches.

Distributed optimization with ADMM. Distributed ADMM algorithms have emerged as a scalable approach for addressing large-scale optimization problems (Boyd et al., 2011; Mota et al., 2013). Despite their significant applicability to machine learning (Mateos et al., 2010), robotics (Shorinwa et al., 2024) and many other fields, their successful performance has been shown to be highly sensitive to the proper tuning of its underlying parameters (Xu et al., 2017; Saravanos et al., 2023). Moreover, tuning parameters for large-scale problems is often tedious and time-consuming, making it desirable to develop effective *learned* optimizers that can be trained on smaller problems instead. Furthermore, even if an distributed optimizer performs well for a specific problem instance, its generalization to new problems remains challenging to verify. These challenges constitute our main motivation for studying learning-aided distributed ADMM architectures. We also note that an ADMM-based distributed QP solver resembling a simpler version of DistributedQP was presented in Pereira et al. (2022), but it focused on multi-robot control and lacked theoretical analysis.

Learning-to-optimize. The area of learning-to-optimize methods has emerged as an effective approach for enhancing existing optimizers or even deriving new algorithmic updates through training on sample problems (Chen et al., 2022; Shlezinger et al., 2022; Amos et al., 2023). A prominent technique in this paradigm is deep unfolding, which under the realistic assumption of computational budget restrictions, unrolls a fixed number of iterations as layers of a deep learning framework and learns the optimal parameters for improving performance on a specific problem class (Monga et al., 2021; Zhang et al., 2020). Nevertheless, combining deep unfolding with distributed ADMM has only been investigated recently in Noah & Shlezinger (2024). Although this framework demonstrates promising results, it is limited to an unconstrained problem formulation, assumes gradient-based local updates, focuses exclusively on parameter tuning and lacks formal performance guarantees. A reinforcement learning algorithm for accelerating OSQP was presented in Ichnowski et al. (2021). While this approach also explores learning policies for algorithm parameters, it is limited to centralized quadratic programming, lacks guarantees and its training comes at a significant computational cost. In the context of establishing generalization bounds for learned optimizers, Sambharya & Stelato (2024) recently explored the idea of incorporating PAC-Bayes bounds in learned optimizers, yet our approach differs fundamentally, as their method employs a binary error function, whereas ours directly establishes bounds based on the optimality gap of the final solution. The works in Sucker & Ochs (2023) and Sucker et al. (2024) are also investigating generalization bounds for learned optimizers, considering the update function as a gradient step or a multi-layer perceptron, respectively.

3 DISTRIBUTED QUADRATIC PROGRAMMING

3.1 PROBLEM FORMULATION

A convex (centralized) QP problem is expressed in general as

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{q}^\top \mathbf{x} \quad \text{s.t.} \quad \mathbf{A} \mathbf{x} \leq \mathbf{b}, \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the decision vector and $\zeta = \{\mathbf{Q} \in \mathbb{S}_{++}^n, \mathbf{q} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m\}$ are the problem data.¹ As the scale of such problems increases to higher dimensions, there is often an underlying networked/decentralized structure that could be leveraged for achieving distributed computations. This work specifically aims to address problems characterized by such structures. Let $\mathbf{w} \in \mathbb{R}^n$ be the main global variable and $\mathbf{x}_i \in \mathbb{R}^{n_i}$ be local variables $i \in \mathcal{V} = \{1, \dots, N\}$. Then, assume a mapping $(i, j) \mapsto \mathcal{G}(i, j)$ from all index pairs (i, j) of local variable components $[\mathbf{x}_i]_j$ to indices $l = \mathcal{G}(i, j)$ of global components w_l ² - for an example see Fig. 2. We consider QP problems of the following *distributed consensus* form:

$$\min_{\mathbf{x}, \mathbf{w}} \sum_{i \in \mathcal{V}} \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i \quad \text{s.t.} \quad \mathbf{A}_i \mathbf{x}_i \leq \mathbf{b}_i, \quad \mathbf{x}_i = \tilde{\mathbf{w}}_i, \quad i \in \mathcal{V}, \quad (2)$$

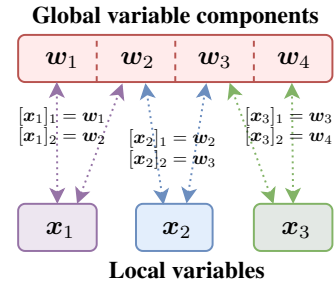


Figure 2: Example of consensus mapping \mathcal{G} in problem (2).

¹Note that equality constraints can also be captured as pairs of inequalities.

²This formulation is adopted from the standard consensus ADMM framework (Boyd et al., 2011), wherein local variables are typically associated with their respective computational nodes.

where the problem data are now given by $\zeta = \{\zeta_i\}_{i=1}^N$ with $\zeta_i = (\mathbf{Q}_i \in \mathbb{S}_{++}^{n_i}, \mathbf{q}_i \in \mathbb{R}^{n_i}, \mathbf{A}_i \in \mathbb{R}^{m_i \times n_i}, \mathbf{b}_i \in \mathbb{R}^{m_i})$. The vector $\mathbf{x} = [\{\mathbf{x}_i\}_{i \in \mathcal{V}}]$ is the concatenation of all local variables, while $\tilde{\mathbf{w}}_i \in \mathbb{R}^{n_i}$, defined as $\tilde{\mathbf{w}}_i = [\{\mathbf{w}_l\}_{l \in \mathcal{G}(q,j):q=i}]$, is the selection of global variable components that correspond to the components of \mathbf{x}_i . This form captures a wide variety of large-scale QPs found in machine learning (Mateos et al., 2010; Navia-Vazquez et al., 2006), optimal control (Van Parys & Pipeleers, 2017), transportation networks, (Mota et al., 2014), power grids (Lin et al., 2012), resource allocation (Huang et al., 2014) and many other fields.

3.2 DISTRIBUTEDQP: THE UNDERLYING OPTIMIZATION ALGORITHM

This section introduces a new distributed algorithm named **DistributedQP** for solving problems of the form (2). The proposed method can be viewed as a combination of consensus ADMM (Boyd et al., 2011) and OSQP (Stellato et al., 2020) using local iteration-varying penalty parameters.

Let us introduce the auxiliary variables $\mathbf{z}_i, \mathbf{s}_i \in \mathbb{R}^{m_i}$, such that problem (2) can be reformulated as

$$\min_{\mathbf{x}} \sum_{i \in \mathcal{V}} \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i \quad \text{s.t.} \quad \mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i, \quad \mathbf{s}_i \leq \mathbf{b}_i, \quad \mathbf{z}_i = \mathbf{s}_i, \quad \mathbf{x}_i = \tilde{\mathbf{w}}_i, \quad i \in \mathcal{V}.$$

The proposed DistributedQP algorithm is then summarized as follows, where k denotes iterations:

1. **Local updates for $\mathbf{x}_i, \mathbf{z}_i$.** For each node $i \in \mathcal{V}$, solve in parallel:

$$\begin{bmatrix} \mathbf{Q}_i + \mu_i^k \mathbf{I} & \mathbf{A}_i^\top \\ \mathbf{A}_i & -1/\rho_i^k \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x}_i^{k+1} \\ \boldsymbol{\nu}_i^{k+1} \end{bmatrix} = \begin{bmatrix} -\mathbf{q}_i + \mu_i^k \tilde{\mathbf{w}}_i - \mathbf{y}_i \\ \mathbf{z}_i - 1/\rho_i^k \boldsymbol{\lambda}_i \end{bmatrix}, \quad (3)$$

and then update in parallel

$$\mathbf{z}_i^{k+1} = \mathbf{s}_i^k + 1/\rho_i^k (\boldsymbol{\nu}_i^{k+1} - \boldsymbol{\lambda}_i^k). \quad (4)$$

2. **Local updates for \mathbf{s}_i and global update for \mathbf{w} .** For each node $i \in \mathcal{V}$, update in parallel:

$$\mathbf{s}_i^{k+1} = \Pi_{\mathbf{s}_i \leq \mathbf{b}_i} (\alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k). \quad (5)$$

In addition, each global variable component \mathbf{w}_l is updated through:

$$\mathbf{w}_l^{k+1} = \alpha^k \frac{\sum_{\mathcal{G}(i,j)=l} \mu_i^k [\mathbf{x}_i]_j}{\sum_{\mathcal{G}(i,j)=l} \mu_i^k} + (1 - \alpha^k) \mathbf{w}_l^k. \quad (6)$$

3. **Local updates for Lagrange multipliers $\boldsymbol{\lambda}_i, \mathbf{y}_i$.** For each node $i \in \mathcal{V}$, update in parallel:

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^k + \rho_i^k (\alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k - \mathbf{s}_i^{k+1}), \quad (7)$$

$$\mathbf{y}_i^{k+1} = \mathbf{y}_i^k + \mu_i^k (\alpha^k \mathbf{x}_i^{k+1} + (1 - \alpha^k) \tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_i^{k+1}). \quad (8)$$

The Lagrange multipliers $\boldsymbol{\nu}_i, \boldsymbol{\lambda}_i$ and \mathbf{y}_i correspond to the equality constraints $\mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i, \mathbf{z}_i = \mathbf{s}_i$ and $\mathbf{x}_i = \tilde{\mathbf{w}}_i$, respectively. The penalty parameters $\rho_i, \mu_i > 0$ correspond to $\mathbf{z}_i = \mathbf{s}_i$ and $\mathbf{x}_i = \tilde{\mathbf{w}}_i$, while $\alpha^k \in [1, 2)$ are over-relaxation parameters. A complete derivation is provided in Appendix A.

3.3 CONVERGENCE GUARANTEES

Prior to unrolling DistributedQP into a deep learning framework, it is particularly important to establish that the underlying optimization algorithm is well-behaved even for varying parameters, i.e., it is expected to asymptotically converge to the optimal solution. This property is especially important in deep unfolding where parameters are expected to be distinct between different iterations.

In the simpler case of $\alpha^k = 1, \rho_i^k = \rho, \mu_i^k = \mu$, the standard convergence guarantees of two-block ADMM would apply directly (Deng & Yin, 2016); for a detailed discussion, see Appendix B. Nevertheless, the introduction of local iteration-varying penalty parameters ρ_i^k, μ_i^k , as well as the over-relaxation with varying parameters α^k makes proving the convergence of this algorithm non-trivial. In the following, we provide convergence guarantees to optimality for DistributedQP.

We consider the following assumption for the penalty parameters.

Assumption 1. As $k \rightarrow \infty$, the parameters $\rho_i^k = \rho_i^{k-1}, \mu_i^k = \mu_i^{k-1}$, for all $i \in \mathcal{V}$.

Theorem 1 (Convergence guarantees for DistributedQP). *If Assumption 1 holds and $\alpha^k \in [1, 2)$, then the iterates \mathbf{w}^k converge to the optimal solution \mathbf{w}^* of problem (2), as $k \rightarrow \infty$.*

The proof of Theorem 1, as well as necessary intermediate results, are provided in Appendix C.

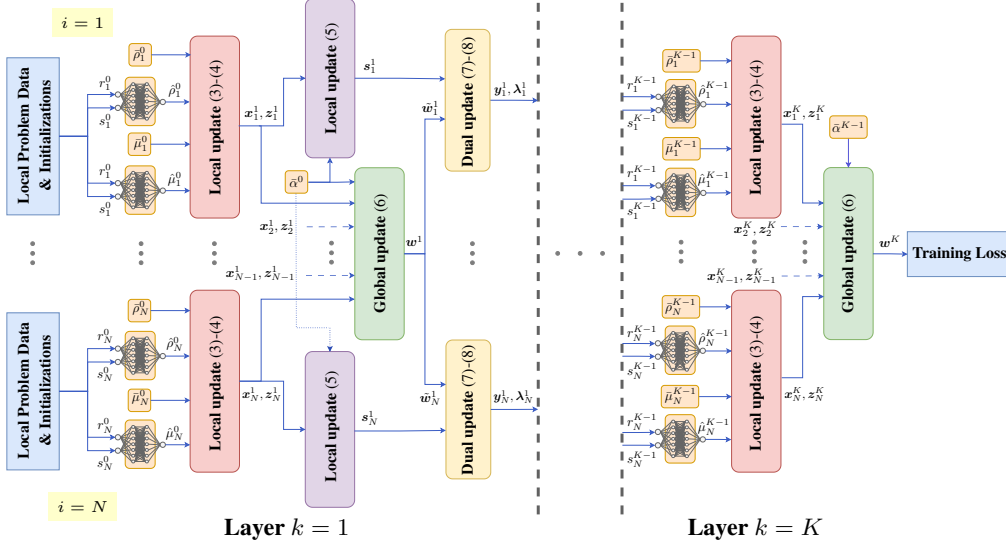


Figure 3: **The DeepDistributedQP architecture.** The proposed framework relies on unrolling the DistributedQP optimizer as a supervised deep learning framework. In particular, we interpret its iterations (3)-(8) as sequential network layers and introduce learnable components (orange blocks) to facilitate reaching the desired accuracy after a predefined number of allowed iterations.

4 THE DEEPDISTRIBUTEDQP ARCHITECTURE

The proposed DeepDistributedQP architecture emerges from unfolding the iterations of the DistributedQP optimizer into a deep learning framework. Section 4.1 illustrates the main architecture, key aspects of our methodology, as well as the centralized version DeepQP. Section 4.2 leverages implicit differentiation during backpropagation to facilitate the training of our framework.

4.1 MAIN ARCHITECTURE

Architecture overview. The **DeepDistributedQP** architecture arises from unrolling the DistributedQP optimizer within the supervised learning paradigm. (Fig. 3). This is accomplished through treating the updates (3)-(7) as blocks in sequential layers of a deep learning network. The number of layers is equal to the predefined number of allowed iterations K , with each layer corresponding to an iteration $k = 1, \dots, K$. The inputs of the network are the local problem data ζ_i and initializations $x_i^0, z_i^0, w_i^0, s_i^0, \lambda_i^0$ and y_i^0 . These are initially passed to N parallel local blocks corresponding to (3)-(4), which output the new variables x_i^1 and z_i^1 . Then, all z_i^1 are fed into N new parallel local blocks (5), yielding the new iterates s_i^1 . In the meantime, all x_i^1 are communicated to a central node that computes the new iterate w^1 through the weighted averaging step (6). Subsequently, the global variable components \tilde{w}_i are communicated back to each local node i , to perform the updates (7)-(8) which output the updated dual variables λ_i, y_i . This group of blocks is then repeated K times, yielding the output of the network which is the final global variable iterate w^K .

Learning feedback policies. Standard deep unfolding typically leverages data to learn algorithm parameters tailored for a specific problem (Shlezinger et al., 2022). From a control theoretic point of view, this process can be interpreted as seeking *open-loop* policies without the incorporating any feedback. In our setup, this would be equivalent with learning the optimal parameters $\bar{\rho}_i^k, \bar{\mu}_i^k, \bar{\alpha}^k$

$$\rho_i^k = \text{SoftPlus}(\bar{\rho}_i^k), \quad \mu_i^k = \text{SoftPlus}(\bar{\mu}_i^k), \quad \alpha^k = \text{Sigmoid}_{1,2}(\bar{\alpha}^k), \quad (9)$$

for all $i = 1, \dots, N$ and $k = 1, \dots, K$, where the $\text{SoftPlus}(\cdot)$ function is used to guarantee the positivity of ρ_i^k, μ_i^k , and the sigmoid function $\text{Sigmoid}_{1,2}(\cdot)$ restricts each α^k to lie between (1, 2).

In the meantime, the predominant practice for online adaptation of the ADMM penalty parameters relies on observing the primal and dual residuals every few iterations (Boyd et al., 2011). The widely-used rule suggests that if the ratio of primal-to-dual residuals is high, the penalty parameter ρ should be increased; conversely, if the ratio is low, ρ should be decreased. Despite its heuristic nature, this approach includes a notion of “feedback” since the current state of the optimizer is used to adapt the parameters, and as a result, it can be interpreted as a closed-loop policy. Based on this

point of view, our goal is to learn the optimal *closed-loop* policies for the local penalty parameters

$$\rho_i^k = \text{SoftPlus}\left(\underbrace{\bar{\rho}_i^k + \pi_{i,\rho}^k(r_{i,\rho}^k, s_{i,\rho}^k; \theta_{i,\rho}^k)}_{\hat{\rho}_i^k}\right), \quad \mu_i^k = \text{SoftPlus}\left(\underbrace{\bar{\mu}_i^k + \pi_{i,\mu}^k(r_{i,\mu}^k, s_{i,\mu}^k; \theta_{i,\mu}^k)}_{\hat{\mu}_i^k}\right), \quad (10)$$

where $\hat{\rho}_i^k, \hat{\mu}_i^k$ are feedback components obtained from policies $\pi_{i,\cdot}^k(r_{i,\cdot}^k, s_{i,\cdot}^k; \theta_{i,\cdot}^k)$, parameterized by fully-connected neural network layers with inputs $r_{i,\cdot}^k, s_{i,\cdot}^k$ and weights $\theta_{i,\cdot}^k$. The terms $r_{i,\cdot}^k$ and $s_{i,\cdot}^k$ represent the local primal and dual residuals of node i at layer k and are detailed in Appendix D.

Solving the local updates. The most computationally demanding block in DeepDistributedQP is solving the local updates (3), as this requires solving a linear system of size $n_i + m_i$. Similar to OSQP (Stellato et al., 2020), this can be accomplished using either a direct or an indirect method. The direct method factors the KKT matrix, solving the system via forward and backward substitution. This approach is particularly efficient when penalty parameters remain fixed, as the same factorization can then be reused across iterations. Nevertheless, at larger scales, this factorization might become impractical. In contrast, with the indirect method, we eliminate ν_i^{k+1} to solve the linear system:

$$\underbrace{(Q_i + \mu_i^k I + A_i^\top \rho_i^k A_i)}_{\tilde{Q}_i^k} x_i^{k+1} = \underbrace{-q_i + \mu_i^k \tilde{w}_i - y_i + A_i^\top \rho_i^k z_i - A_i^\top \lambda_i}_{\tilde{b}_i^k}. \quad (11)$$

This new linear system is solved for x_i^{k+1} using an iterative scheme such as the conjugate gradient (CG) method. We then substitute $\nu_i^{k+1} = \rho_i^k (A_i x_i^{k+1} - z_i) + \lambda_i$. The indirect method has three important properties that make it particularly attractive in our setup. First, its computational complexity scales better w.r.t. the dimension of the local problem, while no additional overhead is introduced by changing the penalty parameters. Second, it can be warmstarted using the solution from the previous iteration, greatly reducing the number of iterations required to converge to a solution. The final important property, which is critical for the scalability of the DeepDistributedQP, is that training with the indirect method can be much more memory efficient as shown in Section 4.2.

Training loss. Let $\mathcal{S} = \{\zeta^j\}_{j=1}^H$ be a dataset consisting of H problem instances $\zeta^j = \{(Q_i, q_i, A_i, b_i)_{i=1}^N, w^*\}_j$ subject to the known mapping \mathcal{G} of problem (2). The loss we are using for training is the average of the γ_k -scaled distances of the global iterates w_1, \dots, w_N from the known optimal solution w^* of each problem instance ζ^j , provided as

$$\ell(\mathcal{S}; \theta) = \frac{1}{H} \sum_{j=1}^H \sum_{k=1}^K \gamma_k \|w^k(\zeta^j; \theta) - w^*(\zeta^j)\|_2, \quad (12)$$

where θ corresponds to the concatenation of all learnable parameters/weights of our framework.

Centralized version. While this work primarily focuses on distributed optimization, for completeness, we also introduce **DeepQP**, the centralized version of our framework, for addressing general QPs of the form (1). In the centralized case, our framework simplifies to $N = 1$, eliminating the need for distinguishing between local and global variables. Under this simplification, the DistributedQP optimizer coincides with OSQP. Hence, DeepQP consists of unfolding the OSQP updates (see Appendix E) and learning policies for adapting its penalty and over-relaxation parameters. The resulting framework is illustrated in Fig. 4. Additional details on DeepQP are provided in Appendix E.

4.2 IMPLICIT DIFFERENTIATION

When solving for the local updates in (11) using the indirect method, it is computationally intractable to backpropagate through all CG iterations. This is especially important in the context of unfolding, as it would become necessary to unroll multiple inner CG optimization loops. To address this, we leverage the implicit function theorem (IFT) to express the solution of (11) as an implicit function

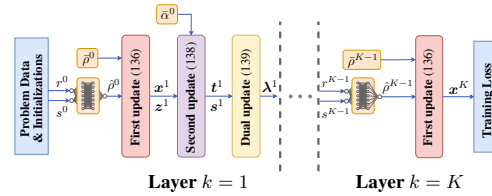


Figure 4: **The DeepQP architecture:** The centralized version of DeepDistributedQP.

of the local problem data. This allows us to compute gradients in a manner that avoids unrolling the CG iterations and requires solving a linear system with the same coefficient matrix, but with a new RHS, achieved by rerunning the CG method. This result is formalized in the following theorem.

Theorem 2 (Implicit Differentiation of Indirect Method). *Let \mathbf{x}_i^{k+1} be the unique solution to the linear system $\bar{\mathbf{Q}}_i^k \mathbf{x}_i^{k+1} = \bar{\mathbf{b}}_i^k$ in (11). Let $\nabla_{\mathbf{x}} L(\mathbf{x}_i^{k+1})$ be a backward pass vector computed through reverse-mode automatic differentiation of some loss function L . Then, the gradient of L with respect to $\bar{\mathbf{Q}}_i^k$ and $\bar{\mathbf{b}}_i^k$ is given by*

$$\begin{aligned}\nabla_{\bar{\mathbf{Q}}_i^k} L &= \frac{1}{2}(\mathbf{x}_i^{k+1} \otimes d\mathbf{x}_i^{k+1} + d\mathbf{x}_i^{k+1} \otimes \mathbf{x}_i^{k+1}), \\ \nabla_{\bar{\mathbf{b}}_i^k} L &= -d\mathbf{x}_i^{k+1},\end{aligned}$$

where $d\mathbf{x}_i^{k+1}$ is the unique solution to the linear system $\bar{\mathbf{Q}}_i^k d\mathbf{x}_i^{k+1} = -\nabla_{\mathbf{x}} L(\mathbf{x}_i^{k+1})$.

The proof is provided in Appendix F and is a straightforward application of the IFT, similar to the results established by Amos & Kolter (2017) and Agrawal et al. (2019).

5 GENERALIZATION BOUNDS

In this section, we establish guarantees on the expected performance of DeepDistributedQP. To achieve this, we leverage the PAC-Bayes framework (Alquier, 2024), a well-known statistical learning methodology for providing bounds on expected loss metrics that hold with high probability. In our case, we provide bounds on the *expected progress* of the final iterate \mathbf{w}^K towards reaching the optimal solution \mathbf{w}^* for unseen problems drawn from the same distribution as the training dataset.

Learning stochastic policies. PAC-Bayes theory is applicable to frameworks that learn weight distributions rather than fixed weights. For this reason, in order to establish such guarantees, we switch to learning a Gaussian distribution of weights $\mathcal{P} = \mathcal{N}(\mu_{\Theta}, \Sigma_{\Theta})$ based on a prior $\mathcal{P}_0 = \mathcal{N}(\mu_{\Theta}^0, \Sigma_{\Theta}^0)$. This choice is motivated by the fact that PAC-Bayes bounds include Kullback–Leibler (KL) divergence terms which can be easily evaluated and optimized for Gaussian distributions.

Generalization bound for DeepDistributedQP. To facilitate the exhibition of our performance guarantees, we provide necessary preliminaries on PAC-Bayes theory in Appendix G. To establish a generalization guarantee for DeepDistributedQP, a meaningful loss function must first be selected. This quantity will be denoted $q(\zeta; \theta)$ to differentiate from the loss used for training. To capture the progress the optimizer makes towards optimality, we propose the following *progress metric*:

$$q(\zeta; \theta) = \min \left\{ \frac{\|\mathbf{w}^K(\zeta; \theta) - \mathbf{w}^*(\zeta)\|_2}{\|\mathbf{w}^0(\zeta) - \mathbf{w}^*(\zeta)\|_2}, 1 \right\}. \quad (13)$$

This loss function measures progress by comparing the distance between the final iterate $\mathbf{w}^K(\zeta; \theta)$ and problem solution $\mathbf{w}^*(\zeta)$ with the distance between the initialization $\mathbf{w}^0(\zeta; \theta)$ and the solution. This choice satisfies the requirement of being bounded between 0 and 1 while being more informative than the indicator losses used in prior work that simply determine whether the final iterate is within a specified neighborhood of the optimal solution (Sambharya & Stellato, 2024). Moreover, this loss is invariant to the scale of the problem data since it is a relative measurement.

As in Appendix G, let $q_{\mathcal{D}}(\mathcal{P})$ be the true expected loss and $q_{\mathcal{S}}(\mathcal{P})$ the empirical expected loss. To evaluate the PAC-Bayes bounds in (148), the expectation $\mathbb{E}_{\theta \sim \mathcal{P}}[q(\zeta; \theta)]$ must be computed as part of the definition of $q_{\mathcal{S}}(\mathcal{P})$. Since no closed-form solution is available, an empirical estimate using M sampled weights $(\theta_i)_{i=1}^M$ is required to upper bound $q_{\mathcal{S}}(\mathcal{P})$ with high probability. We adopt a standard approach involving a sample convergence bound (Majumdar et al. (2021), Dziugaite & Roy (2017), Langford & Caruana (2001)). Specifically, define the empirical estimate of $q_{\mathcal{S}}(\mathcal{P})$ as:

$$\hat{q}_{\mathcal{S}}(\mathcal{P}; M) = \frac{1}{MH} \sum_{i=1}^H \sum_{j=1}^M q(\zeta_i; \theta_j). \quad (14)$$

Then, the following sample convergence bound provides an upper bound on $q_{\mathcal{S}}(\mathcal{P})$,

$$q_{\mathcal{S}}(\mathcal{P}) \leq \bar{q}_{\mathcal{S}}(\mathcal{P}; M, \epsilon) := \mathbb{D}_{\text{KL}}(\hat{q}_{\mathcal{S}}(\mathcal{P}; M) \parallel M^{-1} \log(2/\epsilon)). \quad (15)$$

with probability $1 - \epsilon$. The following theorem summarizes the PAC-Bayes bound we use to evaluate the generalization capabilities of our framework.

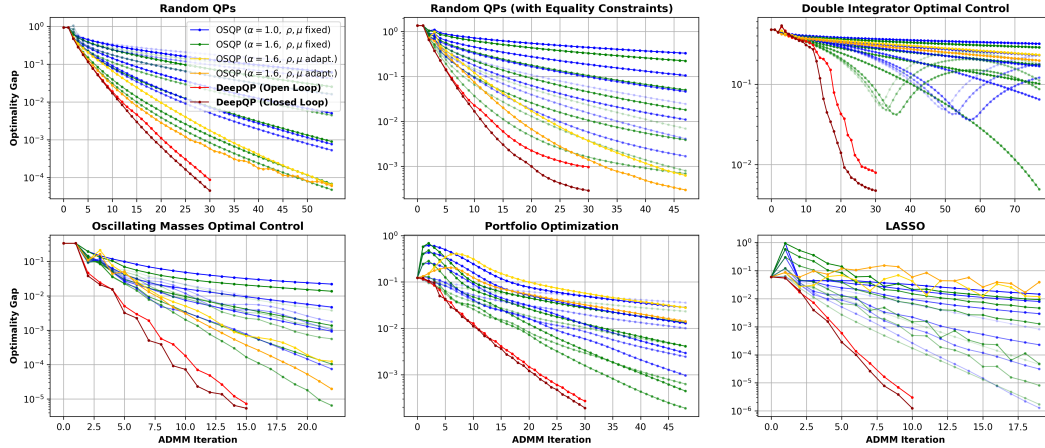


Figure 5: **Small-scale centralized comparison of DeepQP and OSQP.** Across all tested problems, DeepQP consistently outperforms OSQP (same per-iteration complexity using the indirect method).

Theorem 3 (Generalization bound for DeepDistributedQP). *For problems $\zeta \in \mathcal{Z}$ drawn from distribution \mathcal{D} , the true expected progress metric of DeepDistributedQP with policy \mathcal{P} , i.e.,*

$$q_{\mathcal{D}}(\mathcal{P}) = \mathbb{E}_{\zeta \sim \mathcal{D}} \mathbb{E}_{\theta \sim \mathcal{P}} \left[\min \left\{ \frac{\|\mathbf{w}^K(\zeta; \theta) - \mathbf{w}^*(\zeta)\|_2}{\|\mathbf{w}^0(\zeta) - \mathbf{w}^*(\zeta)\|_2}, 1 \right\} \right], \quad (16)$$

is bounded with probability at least $1 - \delta - \epsilon$ by:

$$q_{\mathcal{D}}(\mathcal{P}) \leq \mathbb{D}_{\text{KL}}^{-1} \left(\bar{q}_{\mathcal{S}}(\mathcal{P}; M, \epsilon) \left\| \left(\mathbb{D}_{\text{KL}}(\mathcal{P} \parallel \mathcal{P}_0) + \log(2\sqrt{H}/\delta) \right) / H \right\| \right), \quad (17)$$

where $\bar{q}_{\mathcal{S}}(\mathcal{P}; M, \epsilon)$ is the estimate of $q_{\mathcal{S}}(\mathcal{P}; M, \epsilon)$ described in Eq. (15).

We explain in detail how we train for optimizing the generalization bounds in Appendix H.

6 EXPERIMENTS

We conduct extensive experiments to highlight the effectiveness, scalability and generalizability of the proposed methods. Section 6.1 shows the advantageous performance of DeepQP against OSQP on a variety of centralized QPs. In Section 6.2, we address large-scale problems, showcasing the scalability of DeepDistributedQP despite being trained exclusively on much lower-dimensional instances. Additionally, we discuss the advantages of learning local policies over shared ones and evaluate the proposed generalization bounds, which provide guarantees for the performance of our framework on unseen problems. An overall discussion and potential limitations are provided in Section 6.3. All experiments were performed on a system with an RTX 4090 GPU 24GB, a 13th Gen Intel(R) Core(TM) i9-13900K and 64GB of RAM.

6.1 SMALL-SCALE CENTRALIZED EXPERIMENTS: DEEPQP VS OSQP

Setup. We begin with comparing DeepQP against OSQP for solving centralized QPs (1). The following problems are considered: i,ii) random QPs without/with equality constraints, iii, iv) optimal control for double integrator and oscillating masses, v) portfolio optimization, and vi) LASSO regression. For all problems, we set a maximum allowed amount of iterations K for DeepQP within $[10, 30]$ and examine how many iterations OSQP requires to reach the same accuracy. We train DeepQP using both open-loop and closed-loop policies and with a dataset of size $H \in [500, 2000]$. For OSQP, we consider both constant and adaptive penalty parameters ρ and we set α to be either 1.0 or 1.6. Additional details on DeepQP, OSQP and the problems can be found in Appendix I.

Performance comparison. The comparison between DeepQP and OSQP is illustrated in Fig. 5. Note that both methods share the same per-iteration complexity from solving (139). We evaluate their performance by comparing the (normalized) optimality gap $\|\mathbf{x}^k - \mathbf{x}^*\|_2 / \sqrt{n}$. For all tested problems, DeepQP provides a consistent improvement over OSQP, requiring 1.5 – 3 times fewer iterations to reach the desired accuracy. Furthermore, the advantage of incorporating feedback in the policies is shown, as closed-loop policies outperform open-loop ones in all cases.

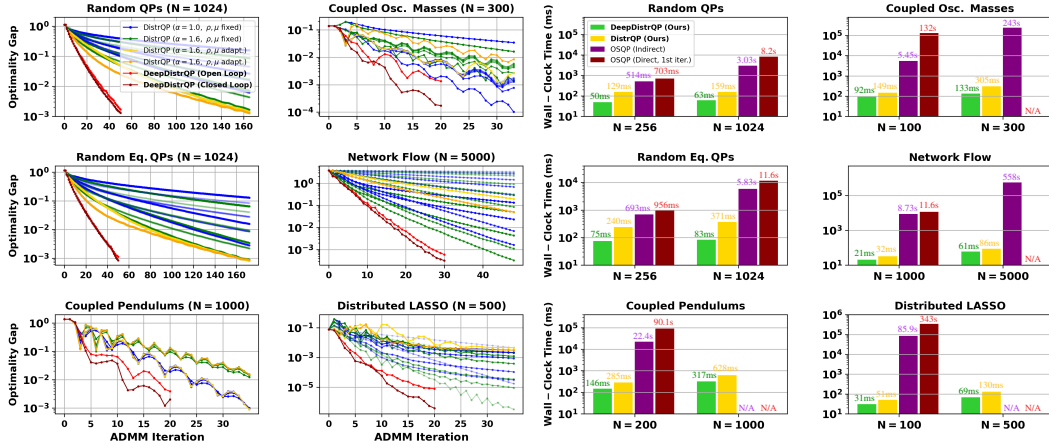


Figure 6: **Scaling DeepDistributedQP to high-dimensional problems.** Left: Comparison between DeepDistributedQP and its traditional optimization counterpart DistributedQP (same per-iteration complexity). Right: Total wall-clock time required by DeepDistributedQP, DistributedQP and OSQP (using indirect or direct method) to achieve the same accuracy.

6.2 LARGE-SCALE DISTRIBUTED EXPERIMENTS: SCALING DEEPPRODISTRIBUTEDQP

Setup. The purpose of the following analysis is to compare the performance and scalability of DeepDistributedQP (ours), DistributedQP (ours) and OSQP for large-scale QPs of the form (2). We consider the following six problems: i,ii) random networked QPs without/with equality constraints, iii, iv) multi-agent optimal control for coupled pendulums and oscillating masses, v) network flow, and vi) distributed LASSO. We select a maximum allowed number of iterations K for DeepDistributedQP within [20, 50] and examine what is the computational effort required by DistributedQP and OSQP to achieve the same accuracy measured by the optimality gap $\|w^k - w^*\|_2/\sqrt{n}$. More details about our experimental setup are provided in Appendix I.

Training on low-dimensional problems. One of the key advantages of DeepDistributedQP is that it only requires using small-scale problems for training. The training dimensions for each problem are detailed in Table 1. Both open-loop and closed-loop versions are trained using shared policies on datasets of size $H \in [500, 1000]$. We employ the shared policies version of DeepDistributedQP to enable the same policies to be applied to larger problems during testing.

Scaling to high-dimensional problems. Subsequently, we evaluate DeepDistributedQP on problems with significantly larger scale than those used during training. The maximum problem dimensions tested are shown in Table 1. On the left side of Fig. 6, we highlight the superior performance of DeepDistributedQP over its standard optimization counterpart DistributedQP (same per-iteration complexity). In all cases, the learned algorithm achieves the same level of accuracy while requiring 1.5-3.5 times fewer iterations. Additionally, the right side of Fig. 6 compares the total wall-clock time between DeepDistributedQP, DistributedQP and OSQP (using indirect or direct method). For a complete illustration, we refer the reader to Table 6 in Appendix I.5. The provided results emphasize the superior scalability of the two proposed distributed methods against OSQP for large-scale QPs, as well as the advantage of our deep learning-aided approach over traditional optimization.

Local vs shared policies. When applying a policy to a problem with the same dimensions as used during training, leveraging local policies instead of shared ones can be advantageous for better exploiting the structure of the problem. On the left side of Fig. 7, we compare the performance of local and shared policies on random QPs and coupled pendulums. For the coupled pendulums problem, which exhibits significant underlying structure, local policies demonstrate clear superiority. For the random QPs problem, where structural patterns are less pronounced, the advantage of local policies is smaller but still significant.

Performance guarantees. Next, we verify the guarantees of our framework for generalizing on unseen random QP ($N = 16$) and coupled pendulums ($N = 10$) problems. We switch from learning deterministic weights to learning stochastic ones and follow the procedure described in Appendix H with $H = 15000$ training samples, $M = 30000$ sampled weights for the bounds evaluation, $\delta = 0.009$ and $\epsilon = 0.001$. The resulting generalization bounds, illustrated in Fig. 7 (right), are expressed in terms of the expected final relative optimality gap - the progress metric

Table 1: **Training dimensions for DeepDistributedQP and maximum testing dimensions.** The metric $\text{nnz}(\mathbf{Q}, \mathbf{A})$ denotes the total number of non-zero elements in \mathbf{Q} and \mathbf{A} .

Problem Class	Training				Max Testing			
	N	n	m	$\text{nnz}(\mathbf{Q}, \mathbf{A})$	N	n	m	$\text{nnz}(\mathbf{Q}, \mathbf{A})$
Random QPs	16	160	120	4,000	1,024	10,240	9,920	300,800
Random QPs w/ Eq. Constr.	16	160	168	4,960	1,024	10,240	9,920	300,800
Coupled Pendulums	10	470	640	3,690	1,000	47,000	64,000	380,880
Coupled Osc. Masses	10	470	1,580	4,590	300	28,200	47,400	141,180
Network Flow	20	100	140	600	5,000	25,000	35,000	150,000
Distributed LASSO	10	1,100	3,000	29,000	500	50,100	150,000	1,450,000

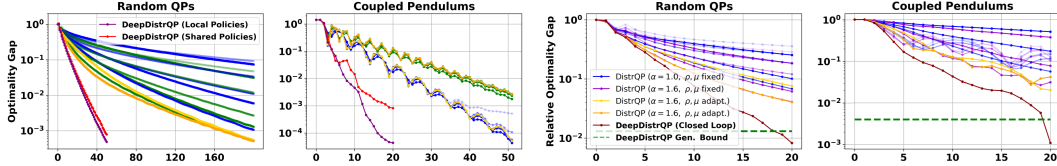


Figure 7: **Left: Local vs shared policies.** We showcase the advantage of learning local policies over shared ones. **Right: Performance guarantees.** The obtained generalization bounds guarantee the performance of DeepDistributedQP and its improvements over DistributedQP.

used for deriving bounds in Section 5, implying that with 99% probability the average performance of our framework will be bounded by this threshold. The bounds are observed to be tight compared to actual performance, underscoring their significance. Moreover, they outperform the standard optimizers, providing a strong guarantee of improved performance for DeepDistributedQP.

6.3 DISCUSSION

In which cases can we use the direct method? As illustrated in Fig. 6 and Table 6, and further discussed in Stellato et al. (2020), the indirect method is generally preferred for solving systems of the form (3) - or (137) for DeepQP/OSQP - once the problem reaches a certain scale. In this work, we adopt this approach both for training, due the memory and computational advantages outlined in Section 4.2, and evaluating DeepDistributedQP/DeepQP. However, it is worth considering whether the direct method might be advantageous during evaluation, a choice that depends on the problem scale and capabilities of the available hardware. Overall, the results of this work show that learning policies for the algorithm parameters is significantly beneficial in the context of both distributed and centralized QP assuming the indirect method is used. In future work, we wish to also explore schemes that adapt the parameters less frequently using the direct method and/or designing mechanisms to dynamically switch between the two approaches.

Limitations. One limitation of the proposed framework is its reliance on a supervised training loss, requiring a dataset of pre-solved problems. In future work, we aim to explore training through directly minimizing the problem residuals rather than the optimality gaps. Furthermore, while PAC-Bayes theory provides an important probabilistic bound on average performance, stronger guarantees may be necessary for safety-critical applications to ensure reliability and robustness.

7 CONCLUSION AND FUTURE WORK

In this work, we introduced DeepDistributedQP, a new deep learning-aided distributed optimization architecture for solving large-scale QP problems. The proposed method relies on unfolding the iterations of a novel optimizer named DistributedQP as layers of a supervised deep learning framework. The expected performance of our learned optimizer on unseen problems is also theoretically established through PAC-Bayes theory. DeepDistributedQP exhibits impressive scalability in effectively tackling large-scale optimization problems while being trained exclusively on much smaller ones. In addition, both DeepDistributedQP and Distributed significantly outperform OSQP in terms of required wall-clock time to reach the same accuracy as dimension increases. Furthermore, we showcase that the proposed PAC-Bayes bounds provide meaningful practical guarantees for the performance of DeepDistributedQP on new problems. In future work, we wish to extend the proposed framework to a semi-supervised version that relies less on pre-solved problems for training. In addition, we wish to explore incorporating more complex learnable components such as LSTMs for feedback within our framework. Finally, we wish to consider other classes of distributed constrained optimization methods outside of quadratic programming.

REFERENCES

- Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J. Zico Kolter. Differentiable convex optimization layers. *Advances in neural information processing systems*, 32, 2019.
- Pierre Alquier. User-friendly introduction to PAC-Bayes bounds. *Foundations and Trends in Machine Learning*, 17(2):174–303, 2024.
- Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *International conference on machine learning*, pp. 136–145. PMLR, 2017.
- Brandon Amos et al. Tutorial on amortized optimization. *Foundations and Trends® in Machine Learning*, 16(5):592–732, 2023.
- Pietro Belotti, Christian Kirches, Sven Leyffer, Jeff Linderoth, James Luedtke, and Ashutosh Mahajan. Mixed-integer nonlinear optimization. *Acta Numerica*, 22:1–131, 2013.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine learning*, 3(1):1–122, 2011.
- Steven W Chen, Tianyu Wang, Nikolay Atanasov, Vijay Kumar, and Manfred Morari. Large scale model predictive control with neural networks and primal active sets. *Automatica*, 135:109947, 2022.
- Andrew R Conn, Nicholas IM Gould, and Philippe L Toint. *Trust region methods*. SIAM, 2000.
- Christian Conte, Tyler Summers, Melanie N Zeilinger, Manfred Morari, and Colin N Jones. Computational aspects of distributed optimization in model predictive control. In *2012 IEEE 51st IEEE conference on decision and control (CDC)*, pp. 6819–6824. IEEE, 2012a.
- Christian Conte, Niklaus R Voellmy, Melanie N Zeilinger, Manfred Morari, and Colin N Jones. Distributed synthesis and control of constrained linear systems. In *2012 American control conference (ACC)*, pp. 6017–6022. IEEE, 2012b.
- Gerard Cornuejols, Javier Peña, and Reha Tütüncü. *Optimization methods in finance*. Cambridge University Press, 2018.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 1995.
- Wei Deng and Wotao Yin. On the global and linear convergence of the generalized alternating direction method of multipliers. *Journal of Scientific Computing*, 66:889–916, 2016.
- Gintare Karolina Dziugaite and Daniel M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint arXiv:1703.11008*, 2017.
- Gintare Karolina Dziugaite, Kyle Hsu, Waseem Gharbieh, Gabriel Arpino, and Daniel Roy. On the role of data in PAC-Bayes bounds. In *International Conference on Artificial Intelligence and Statistics*, pp. 604–612. PMLR, 2021.
- Carlos E Garcia, David M Prett, and Manfred Morari. Model predictive control: Theory and practice—a survey. *Automatica*, 25(3):335–348, 1989.
- Boris Houska, Janick Frasch, and Moritz Diehl. An augmented lagrangian based algorithm for distributed nonconvex optimization. *SIAM Journal on Optimization*, 26(2):1101–1127, 2016.
- Shaojun Huang, Qiuwei Wu, Shmuel S Oren, Ruoyang Li, and Zhaoxi Liu. Distribution locational marginal pricing through quadratic programming for congestion management in distribution networks. *IEEE Transactions on Power Systems*, 30(4):2170–2178, 2014.
- Jeffrey Ichnowski, Paras Jain, Bartolomeo Stellato, Goran Banjac, Michael Luo, Francesco Borrelli, Joseph E Gonzalez, Ion Stoica, and Ken Goldberg. Accelerating quadratic optimization with reinforcement learning. *Advances in Neural Information Processing Systems*, 34:21043–21055, 2021.

- Steven George Krantz and Harold R Parks. *The implicit function theorem: history, theory, and applications*. Springer Science & Business Media, 2002.
- John Langford and Rich Caruana. (not) bounding the true error. *Advances in Neural Information Processing Systems*, 14, 2001.
- Shieh-Shing Lin, Shih-Cheng Horng, et al. Distributed quadratic programming problems of power systems with continuous and discrete variables. *IEEE Transactions on Power Systems*, 28(1): 472–481, 2012.
- Anirudha Majumdar, Alec Farid, and Anoopkumar Sonar. PAC-Bayes control: learning policies that provably generalize to novel environments. *The International Journal of Robotics Research*, 40 (2-3):574–593, 2021.
- Gonzalo Mateos, Juan Andrés Bazerque, and Georgios B Giannakis. Distributed sparse linear regression. *IEEE Transactions on Signal Processing*, 58(10):5262–5276, 2010.
- John Mattingley and Stephen Boyd. Real-time convex optimization in signal processing. *IEEE Signal processing magazine*, 27(3):50–61, 2010.
- Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):18–44, 2021.
- Joao FC Mota. *Communication-efficient algorithms for distributed optimization*. PhD thesis, Carnegie Mellon University, 2013.
- Joao FC Mota, Joao MF Xavier, Pedro MQ Aguiar, and Markus Püschel. D-admm: A communication-efficient distributed algorithm for separable optimization. *IEEE Transactions on Signal processing*, 61(10):2718–2723, 2013.
- João FC Mota, João MF Xavier, Pedro MQ Aguiar, and Markus Püschel. Distributed optimization with local domains: Applications in mpc and network flows. *IEEE Transactions on Automatic Control*, 60(7):2004–2009, 2014.
- Angel Navia-Vazquez, D Gutierrez-Gonzalez, Emilio Parrado-Hernández, and JJ Navarro-Abellan. Distributed support vector machines. *IEEE Transactions on Neural Networks*, 17(4):1091–1097, 2006.
- Yurii Nesterov and Arkadii Nemirovskii. *Interior-point polynomial algorithms in convex programming*. SIAM, 1994.
- Yoav Noah and Nir Shlezinger. Distributed learn-to-optimize: Limited communications optimization over networks via deep unfolded distributed admm. *IEEE Transactions on Mobile Computing*, 2024.
- Jorge Nocedal and Stephen J Wright. *Numerical optimization*. Springer, 1999.
- Marcus A Pereira, Augustinos D Saravanos, Oswin So, and Evangelos A. Theodorou. Decentralized Safe Multi-agent Stochastic Optimal Control using Deep FBSDEs and ADMM. In *Proceedings of Robotics: Science and Systems*, New York City, NY, USA, June 2022. doi: 10.15607/RSS.2022.XVIII.055.
- James Blake Rawlings, David Q Mayne, Moritz Diehl, et al. *Model predictive control: theory, computation, and design*, volume 2. Nob Hill Publishing Madison, WI, 2017.
- Rajiv Sambharya and Bartolomeo Stellato. Data-driven performance guarantees for classical and learned optimizers. *arXiv preprint arXiv:2404.13831*, 2024.
- Augustinos D Saravanos, Yuichiro Aoyama, Hongchang Zhu, and Evangelos A Theodorou. Distributed differential dynamic programming architectures for large-scale multiagent control. *IEEE Transactions on Robotics*, 2023.
- Nir Shlezinger, Yonina C Eldar, and Stephen P Boyd. Model-based deep learning: On the intersection of deep learning and optimization. *IEEE Access*, 10:115384–115398, 2022.

- Ola Shorinwa, Trevor Halsted, Javier Yu, and Mac Schwager. Distributed optimization methods for multi-robot systems: Part 1—a tutorial. *IEEE Robotics & Automation Magazine*, 2024.
- Majid Soheili and Amir Masoud Eftekhari-Moghadam. Dqpfs: Distributed quadratic programming based feature selection for big data. *Journal of Parallel and Distributed Computing*, 138:1–14, 2020.
- Valeriu Soltan. Moreau-type characterizations of polar cones. *Linear Algebra and its Applications*, 567:45–62, 2019. ISSN 0024-3795. doi: <https://doi.org/10.1016/j.laa.2019.01.006>. URL <https://www.sciencedirect.com/science/article/pii/S0024379519300199>.
- Bartolomeo Stellato, Goran Banjac, Paul Goulart, Alberto Bemporad, and Stephen Boyd. OSQP: An operator splitting solver for quadratic programs. *Mathematical Programming Computation*, 12(4):637–672, 2020.
- Michael Sucker and Peter Ochs. Pac-bayesian learning of optimization algorithms. In *International Conference on Artificial Intelligence and Statistics*, pp. 8145–8164. PMLR, 2023.
- Michael Sucker, Jalal Fadili, and Peter Ochs. Learning-to-optimize with pac-bayesian guarantees: Theoretical considerations and practical implementation. *arXiv preprint arXiv:2404.03290*, 2024.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- Ruben Van Parys and Goele Pipeleers. Distributed mpc for multi-vehicle systems moving in formation. *Robotics and Autonomous Systems*, 97:144–152, 2017.
- Philip Wolfe. The simplex method for quadratic programming. *Econometrica: Journal of the Econometric Society*, pp. 382–398, 1959.
- Zheng Xu, Mario Figueiredo, and Tom Goldstein. Adaptive admm with spectral penalty parameter selection. In *Artificial Intelligence and Statistics*, pp. 718–727. PMLR, 2017.
- Kai Zhang, Luc Van Gool, and Radu Timofte. Deep unfolding network for image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3217–3226, 2020.

A COMPLETE DERIVATION OF DISTRIBUTEDQP ALGORITHM

Here, we present the detailed derivation of the DistributedQP algorithm presented in Section 3.2. We consider the over-relaxed version of ADMM (Boyd et al., 2011) with $\alpha \in [1, 2)$.

First, let us rewrite problem (2) as

$$\min_{\mathbf{x}} \sum_{i \in \mathcal{V}} \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i \quad \text{s.t.} \quad \mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i, \mathbf{z}_i \leq \mathbf{b}_i, \mathbf{x}_i = \tilde{\mathbf{w}}_i, \quad i \in \mathcal{V}. \quad (18)$$

where we have introduced the auxiliary variables \mathbf{z}_i for each $i = 1, \dots, N$. In addition, we let us define the new variables \mathbf{s}_i , $i = 1, \dots, N$, and rewrite the above problem as

$$\begin{aligned} \min_{\mathbf{x}} \quad & \sum_{i \in \mathcal{V}} \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i + \mathcal{I}_{\mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i}(\mathbf{x}_i, \mathbf{z}_i) \\ \text{s.t.} \quad & \mathbf{z}_i = \mathbf{s}_i, \mathbf{s}_i \leq \mathbf{b}_i, \mathbf{x}_i = \tilde{\mathbf{w}}_i, \quad i = 1, \dots, N. \end{aligned} \quad (19)$$

The above splitting constitutes the problem suitable for being addressed with a two-block ADMM scheme where the first block of variables consists of $\{\mathbf{x}_i, \mathbf{z}_i\}_{i=1, \dots, N}$, while the second one consists of $\{\mathbf{s}_i\}_{i=1, \dots, N}$ and \mathbf{w} . The (scaled) augmented Lagrangian (AL) for this problem is given by

$$\begin{aligned} \mathcal{L} = \sum_{i \in \mathcal{V}} \quad & \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i + \mathcal{I}_{\mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i}(\mathbf{x}_i, \mathbf{z}_i) + \mathcal{I}_{\mathbf{s}_i \leq \mathbf{b}_i}(\mathbf{s}_i) \\ & + \frac{\rho_i}{2} \left\| \mathbf{z}_i - \mathbf{s}_i + \frac{\boldsymbol{\lambda}_i}{\rho_i} \right\|_2^2 + \frac{\mu_i}{2} \left\| \mathbf{x}_i - \tilde{\mathbf{w}}_i + \frac{\mathbf{y}_i}{\mu_i} \right\|_2^2. \end{aligned} \quad (20)$$

The first block of variables is updated through

$$\{\mathbf{x}_i, \mathbf{z}_i\}_{i \in \mathcal{V}} = \arg \min \mathcal{L}(\mathbf{x}, \mathbf{z}, \mathbf{s}^k, \mathbf{w}^k, \boldsymbol{\lambda}^k, \mathbf{y}^k) \quad (21)$$

which can be decoupled to the following subproblems for each $i \in \mathcal{V}$:

$$\begin{aligned} \{\mathbf{x}_i, \mathbf{z}_i\} = \arg \min \quad & \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i + \frac{\rho_i}{2} \left\| \mathbf{z}_i - \mathbf{s}_i + \frac{\boldsymbol{\lambda}_i}{\rho_i} \right\|_2^2 + \frac{\mu_i}{2} \left\| \mathbf{x}_i - \tilde{\mathbf{w}}_i + \frac{\mathbf{y}_i}{\mu_i} \right\|_2^2 \\ \text{s.t.} \quad & \mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i. \end{aligned} \quad (22)$$

Since these problems are equality constrained QPs, we can find a closed-form solution. The optimality conditions for each subproblem are given by

$$\mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i + \mu_i(\mathbf{x}_i - \tilde{\mathbf{w}}_i) + \mathbf{y}_i + \mathbf{A}_i^\top \boldsymbol{\nu}_i = \mathbf{0} \quad (23a)$$

$$\rho_i(\mathbf{z}_i - \mathbf{s}_i) + \boldsymbol{\lambda}_i - \boldsymbol{\nu}_i = \mathbf{0} \quad (23b)$$

$$\mathbf{A}_i \mathbf{x}_i - \mathbf{z}_i = \mathbf{0} \quad (23c)$$

where $\boldsymbol{\nu}_i$ is the Lagrange multiplier corresponding to the constraint $\mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i$. Eliminating \mathbf{z}_i leads to the following system of equations

$$\begin{bmatrix} \mathbf{Q}_i + \mu_i \mathbf{I} & \mathbf{A}_i^\top \\ \mathbf{A}_i & -1/\rho_i \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x}_i^{k+1} \\ \boldsymbol{\nu}_i^{k+1} \end{bmatrix} = \begin{bmatrix} -\mathbf{q}_i + \mu_i \tilde{\mathbf{w}}_i^k - \mathbf{y}_i^k \\ \mathbf{z}_i^k - 1/\rho_i \boldsymbol{\lambda}_i^k \end{bmatrix} \quad (24)$$

with \mathbf{z}_i^{k+1} given by

$$\mathbf{z}_i^{k+1} = \mathbf{s}_i^k + \rho_i^{-1}(\boldsymbol{\nu}_i^{k+1} - \boldsymbol{\lambda}_i^k). \quad (25)$$

The second block of updates is given by

$$\{\mathbf{s}_i\}_{i \in \mathcal{V}}, \mathbf{w} = \arg \min \mathcal{L}(\mathbf{x}^{k+1}, \mathbf{z}^{k+1}, \mathbf{s}, \mathbf{w}, \boldsymbol{\lambda}^k, \mathbf{y}^k) \quad (26)$$

or more analytically

$$\{\mathbf{s}_i\}_{i \in \mathcal{V}}, \mathbf{w} = \arg \min \sum_{i \in \mathcal{V}} \frac{\rho_i}{2} \left\| \alpha \mathbf{z}_i^{k+1} + (1 - \alpha) \mathbf{s}_i^k - \mathbf{s}_i + \frac{\boldsymbol{\lambda}_i^k}{\rho_i} \right\|_2^2 \quad (27)$$

$$+ \frac{\mu_i}{2} \left\| \alpha \mathbf{x}_i^{k+1} + (1 - \alpha) \tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_i + \frac{\mathbf{y}_i^k}{\mu_i} \right\|_2^2 \quad \text{s.t.} \quad \mathbf{s}_i \leq \mathbf{b}_i \quad (28)$$

Note that this minimization can be decoupled w.r.t. all \mathbf{s}_i , $i \in \mathcal{V}$ and \mathbf{w} . In particular, each \mathbf{s}_i can be updated in parallel through

$$\mathbf{s}_i^{k+1} = \Pi_{\mathbf{s}_i \leq \mathbf{b}_i} (\alpha \mathbf{z}_i^{k+1} + (1 - \alpha) \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i). \quad (29)$$

The global variable \mathbf{w} minimization can be decoupled among the components $\ell = 1, \dots, n$, as follows

$$\mathbf{w}_\ell = \arg \min \sum_{\mathcal{G}(i,j)=\ell} \frac{\mu_i}{2} \left\| \alpha [\mathbf{x}_i^{k+1}]_j + (1 - \alpha) [\tilde{\mathbf{w}}_i^k]_j - [\tilde{\mathbf{w}}_i]_j + \frac{[\mathbf{y}_i^k]_j}{\mu_i} \right\|_2^2 \quad (30)$$

Setting the gradient to be equal to zero gives

$$\sum_{\mathcal{G}(i,j)=\ell} \mu_i \left[\alpha [\mathbf{x}_i^{k+1}]_j + (1 - \alpha) \mathbf{w}_\ell^k - \mathbf{w}_\ell^{k+1} + \frac{[\mathbf{y}_i^k]_j}{\mu_i} \right] = \mathbf{0} \quad (31)$$

which leads to

$$\sum_{(\mathcal{G}(i,j)=\ell)} \mu_i \mathbf{w}_\ell^{k+1} = \sum_{\mathcal{G}(i,j)=\ell} \mu_i \left[\alpha [\mathbf{x}_i^{k+1}]_j + (1 - \alpha) \mathbf{w}_\ell^k + \frac{[\mathbf{y}_i^k]_j}{\mu_i} \right] \quad (32)$$

which eventually gives the update rule

$$\mathbf{w}_\ell^{k+1} = \frac{\sum_{\mathcal{G}(i,j)=\ell} \alpha \mu_i [\mathbf{x}_i^{k+1}]_j + [\mathbf{y}_i^k]_j}{\sum_{\mathcal{G}(i,j)=\ell} \mu_i} + (1 - \alpha) \mathbf{w}_\ell^k. \quad (33)$$

Finally, the dual variables are updated through dual ascent steps as follows

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^k + \rho_i (\alpha \mathbf{z}_i^{k+1} + (1 - \alpha) \mathbf{s}_i^k - \mathbf{s}_i^{k+1}) \quad (34)$$

$$\mathbf{y}_i^{k+1} = \mathbf{y}_i^k + \mu_i (\alpha \mathbf{x}_i^{k+1} + (1 - \alpha) \tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_i^{k+1}). \quad (35)$$

It is important to observe that after the first iteration, the global update can be simplified to

$$\mathbf{w}_\ell^{k+1} = \alpha \frac{\sum_{\mathcal{G}(i,j)=\ell} \mu_i [\mathbf{x}_i^{k+1}]_j}{\sum_{\mathcal{G}(i,j)=\ell} \mu_i} + (1 - \alpha) \mathbf{w}_\ell^k, \quad (36)$$

since the summation

$$\begin{aligned} \sum_{\mathcal{G}(i,j)=\ell} [\mathbf{y}_i^{k+1}]_j &= \sum_{\mathcal{G}(i,j)=\ell} [\mathbf{y}_i^k]_j + \mu_i (\alpha [\mathbf{x}_i^{k+1}]_j + (1 - \alpha) [\tilde{\mathbf{w}}_i^k]_j - [\tilde{\mathbf{w}}_i^{k+1}]_j) \\ &= \sum_{\mathcal{G}(i,j)=\ell} [\mathbf{y}_i^k]_j + \mu_i (\alpha [\mathbf{x}_i^{k+1}]_j + (1 - \alpha) \mathbf{w}_\ell^k - \mathbf{w}_\ell^{k+1}) \\ &= \sum_{\mathcal{G}(i,j)=\ell} [\mathbf{y}_i^k]_j + \mu_i \left[\alpha [\mathbf{x}_i^{k+1}]_j + \cancel{(1 - \alpha) \mathbf{w}_\ell^k} \right. \\ &\quad \left. - \frac{\sum_{\mathcal{G}(u,v)=\ell} \alpha \mu_u [\mathbf{x}_u^{k+1}]_v + [\mathbf{y}_u^k]_v}{\sum_{\mathcal{G}(u,v)=\ell} \mu_u} - \cancel{(1 - \alpha) \mathbf{w}_\ell^k} \right] \\ &= \sum_{\mathcal{G}(i,j)=\ell} [\mathbf{y}_i^k]_j + \mu_i \left[\alpha [\mathbf{x}_i^{k+1}]_j - \frac{\sum_{\mathcal{G}(u,v)=\ell} \alpha \mu_u [\mathbf{x}_u^{k+1}]_v + [\mathbf{y}_u^k]_v}{\sum_{\mathcal{G}(u,v)=\ell} \mu_u} \right] \\ &= \sum_{\mathcal{G}(i,j)=\ell} [\mathbf{y}_i^k]_j + \alpha \mu_i [\mathbf{x}_i^{k+1}]_j - \frac{\sum_{\mathcal{G}(i,j)=\ell} \mu_i \left[\sum_{\mathcal{G}(u,v)=\ell} \alpha \mu_u [\mathbf{x}_u^{k+1}]_v + [\mathbf{y}_u^k]_v \right]}{\sum_{\mathcal{G}(u,v)=\ell} \mu_u} \\ &= \sum_{\mathcal{G}(i,j)=\ell} [\mathbf{y}_i^k]_j + \alpha \mu_i [\mathbf{x}_i^{k+1}]_j - \sum_{\mathcal{G}(u,v)=\ell} \alpha \mu_u [\mathbf{x}_u^{k+1}]_v + [\mathbf{y}_u^k]_v = 0. \end{aligned} \quad (37)$$

B STANDARD CONVERGENCE GUARANTEES FOR SIMPLIFIED DISTRIBUTEDQP

In the simplified case where $\rho_i^k = \rho$, $\mu_i^k = \mu$ for all $i \in \mathcal{V}$ and for all k , as well as $\alpha^k = 1$, for all k , it would be straightforward to apply the classical convergence guarantees of two-block ADMM for convex optimization problems.

Let us define the variables $\bar{\mathbf{x}} = [\{\mathbf{x}_i\}_{i \in \mathcal{V}}; \{\mathbf{z}_i\}_{i \in \mathcal{V}}]$ and $\bar{\mathbf{z}} = [\{\mathbf{s}_i\}_{i \in \mathcal{V}}; \mathbf{w}]$. Then, we can rewrite problem (19) as

$$\min f(\bar{\mathbf{x}}) + g(\bar{\mathbf{z}}) \quad \text{s.t.} \quad \bar{\mathbf{A}}\bar{\mathbf{x}} + \bar{\mathbf{B}}\bar{\mathbf{z}} = \bar{\mathbf{c}}, \quad (38)$$

where

$$f(\bar{\mathbf{x}}) = \sum_{i \in \mathcal{V}} \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i + \mathcal{I}_{\mathbf{A}_i \mathbf{x}_i = \mathbf{z}_i}(\mathbf{x}_i, \mathbf{z}_i), \quad g(\bar{\mathbf{z}}) = \sum_{i \in \mathcal{V}} \mathcal{I}_{\mathbf{s}_i \leq \mathbf{b}_i}(\mathbf{s}_i), \quad (39)$$

and $\bar{\mathbf{A}} = \text{bdiag}(\mathbf{I}, \mathbf{I})$, $\bar{\mathbf{B}} = \text{bdiag}(\mathbf{I}, \mathbf{G})$ and $\bar{\mathbf{c}} = \mathbf{0}$, with $\mathbf{G} \in \mathbb{R}^{(\sum_i n_i) \times n}$ defined such that $\mathbf{x} = \mathbf{G}\mathbf{w}$. In other words, \mathbf{G} is the matrix that represents the local-to-global variable components mapping, formally defined as $\mathbf{G} = [\mathbf{G}_1; \dots; \mathbf{G}_N]$ with each submatrix $\mathbf{G}_i \in \mathbb{R}^{n_i \times n}$ given by

$$[\mathbf{G}_i]_{u,v} = \begin{cases} 1, & \text{if } v = \mathcal{G}(i, u) \\ 0, & \text{else} \end{cases}. \quad (40)$$

Given this representation, it becomes clear that our algorithm can be framed as a two-block ADMM. Now, note that \mathbf{G} is a full column rank matrix since all global variable components \mathbf{g}_ℓ are mapped to at least one local variable component $[\mathbf{x}_i]_j$. Then, since the functions f, g are convex and the matrices $\bar{\mathbf{A}}, \bar{\mathbf{B}}$ are full column rank, it follows from Deng & Yin (2016) that the algorithm is guaranteed to converge to the optimal solution.

Nevertheless, this analysis would have only been applicable to this simplified case of the proposed DISTRQP algorithm. In Appendix C, we tackle the more complex case of varying local penalty parameters and varying relaxation parameters. A similar analysis is exhibited in Xu et al. (2017), the convergence of an adaptive relaxed variant of two-block ADMM is provided, yet that is not directly applicable to our case which includes local penalty parameters.

C PROOF OF DISTRIBUTEDQP ASYMPTOTIC CONVERGENCE

C.1 SKETCH OF PROOF

In this section, we prove the convergence guarantees for DistributedQP. To begin, we outline the following conventions. The points $\mathbf{x}^*, \mathbf{z}^*, \mathbf{s}^*, \mathbf{w}^*, \mathbf{y}^*, \boldsymbol{\lambda}^*$ are the KKT points of the problem (19). We refer to the notion of a distance function at any $(k+1)^{\text{th}}$ iteration to be representing a weighted squared norm of the difference between the variables $\mathbf{s}^{k+1}, \mathbf{w}^{k+1}, \mathbf{y}^{k+1}, \boldsymbol{\lambda}^{k+1}$ and their corresponding optimal values $\mathbf{s}^*, \mathbf{w}^*, \mathbf{y}^*, \boldsymbol{\lambda}^*$, indicating the distance from the optimal point.

We prove the convergence in the following steps.

- First, we will derive a descent relation (101), which establishes a relationship between the values of the distance function for consecutive iterations. To derive the descent relation in Lemma 4, we first introduce the relations (R1-R8) in Lemma 1-3.
- Next, we use the derived descent relation to prove the convergence in the subsection C.3 based on Assumption 1.

C.2 NECESSARY LEMMAS

Here, we present the necessary lemmas before proving the convergence of the algorithm. For notational convenience, we use

$$f_i(\mathbf{x}_i) = \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i, \quad \mathcal{C}_i = \{\mathbf{s}_i | \mathbf{s}_i \leq \mathbf{b}_i\}. \quad (41)$$

for each $i \in \mathcal{V}$.

Lemma 1. For all $i \in \mathcal{V}$, the following relationships hold at every iteration k :

$$(R1): \quad \sum_{i \in \mathcal{V}} \mathbf{G}_i^\top \mathbf{y}_i^{k+1} = \mathbf{0}, \quad (42)$$

$$(R2): \quad \alpha^k \mathbf{x}_i^{k+1} = \frac{1}{\mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) - (1 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1}, \quad (43)$$

$$(R3): \quad \alpha^k \mathbf{z}_i^{k+1} = \frac{1}{\rho_i^k} (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k) - (1 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1}, \quad (44)$$

$$(R4): \quad \boldsymbol{\lambda}_i^{k\top} (\mathbf{t}_1 - \mathbf{t}_2) = 0, \quad \text{for any } \mathbf{t}_1, \mathbf{t}_2 \in \mathcal{C}_i. \quad (45)$$

Proof. Relationship (R1) is equivalent with the argument proved in (37). Indeed, if we observe that each matrix $\mathbf{G}_i^\top \in \mathbb{R}^{n \times n_i}$ indicates the mapping from local indices (i, j) to global indices ℓ for a particular i , then we can write

$$\sum_{i \in \mathcal{V}} \mathbf{G}_i^\top \mathbf{y}_i^{k+1} = \begin{bmatrix} \sum_{\mathcal{G}(i,j)=1} [\mathbf{y}_i^{k+1}]_j \\ \vdots \\ \sum_{\mathcal{G}(i,j)=n} [\mathbf{y}_i^{k+1}]_j \end{bmatrix} = \mathbf{0}. \quad (46)$$

Relationship (R2) follows by rearranging the dual update step (8) and replacing $\tilde{\mathbf{w}}_i = \mathbf{G}_i \mathbf{w}$. Similarly, relationship (R3) follows by rearranging the dual update step (7).

In the remaining, we focus on proving (R4). Let us first rewrite the \mathbf{s}_i update (5) as

$$\mathbf{s}_i^{k+1} = \Pi_{\mathcal{C}_i} (\alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k). \quad (47)$$

Now, we consider a closed convex cone $\bar{\mathcal{C}}_i$ defined as

$$\bar{\mathcal{C}}_i = \{\mathbf{p} \mid \mathbf{p} \leq \mathbf{0}\}, \quad (48)$$

such that (47) can be rewritten as follows.

$$\mathbf{s}_i^{k+1} = \Pi_{\bar{\mathcal{C}}_i} (\alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k - \mathbf{b}_i) + \mathbf{b}_i \quad (49)$$

Let us also define $\hat{\mathbf{s}}^{k+1}$ as follows,

$$\hat{\mathbf{s}}_i^{k+1} = \alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k - \mathbf{b}_i \quad (50)$$

such that (49) is given as

$$\mathbf{s}_i^{k+1} = \Pi_{\bar{\mathcal{C}}_i} (\hat{\mathbf{s}}_i^{k+1}) + \mathbf{b}_i \quad (51)$$

Now, let us rewrite the dual update for $\boldsymbol{\lambda}_i^{k+1}$ in (7) as follows.

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^k + \rho_i^k (\alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k - \mathbf{s}_i^{k+1}) \quad (52)$$

which can be rearranged to

$$\boldsymbol{\lambda}_i^{k+1} = \rho_i^k (\boldsymbol{\lambda}_i^k / \rho_i^k + \alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k - \mathbf{s}_i^{k+1}) \quad (53)$$

Using (50), the above can be rewritten as

$$\boldsymbol{\lambda}_i^{k+1} = \rho_i^k (\hat{\mathbf{s}}_i^{k+1} + \mathbf{b}_i - \mathbf{s}_i^{k+1}) \quad (54)$$

Substituting (51) in the above, we get

$$\boldsymbol{\lambda}_i^{k+1} = \rho_i^k (\hat{\mathbf{s}}_i^{k+1} - \Pi_{\bar{\mathcal{C}}_i} (\hat{\mathbf{s}}_i^{k+1})) \quad (55)$$

By Moreau's decomposition (refer to theorems 1.1 and 1.2 from Soltan (2019)), the vector $\hat{\mathbf{s}}_i^{k+1}$ can be written as follows.

$$\hat{\mathbf{s}}_i^{k+1} = \Pi_{\bar{\mathcal{C}}_i} (\hat{\mathbf{s}}_i^{k+1}) + \Pi_{\bar{\mathcal{C}}_i^\circ} (\hat{\mathbf{s}}_i^{k+1}) \quad (56)$$

where $\bar{\mathcal{C}}_i^\circ$ is a polar cone to $\bar{\mathcal{C}}_i$.

Note that any two cone sets \mathcal{D} and \mathcal{D}° are called polar cones if for any $\mathbf{d} \in \mathcal{D}$ and $\bar{\mathbf{d}} \in \mathcal{D}^\circ$, $\mathbf{d}^\top \bar{\mathbf{d}} = 0$.

Thus, using (56) and (55), we get

$$\lambda_i^{k+1} = \rho_i^k \Pi_{\bar{\mathcal{C}}_i^o}(\hat{s}_i^{k+1}) \quad (57)$$

which implies that

$$\lambda_i^{k+1} / \rho_i^k \in \bar{\mathcal{C}}_i^o \quad (58)$$

Further, since $\bar{\mathcal{C}}_i^o$ is a cone, and $\rho_i^k > 0$, we get

$$\lambda_i^{k+1} \in \bar{\mathcal{C}}_i^o \quad (59)$$

Now, any vector $t \in \mathcal{C}_i$ satisfies the following.

$$t - b_i \in \bar{\mathcal{C}}_i \quad (60)$$

Since $\bar{\mathcal{C}}_i$ and $\bar{\mathcal{C}}_i^o$ are polar cones, and using (59), the following relation holds true by the definition of polar cones.

$$\lambda_i^{k+1T}(t - b_i) = 0 \quad \text{for all } t \in \mathcal{C}_i \quad (61)$$

Thus, for any vectors $t_1, t_2 \in \mathcal{C}_i$ and for all k , we have

$$\lambda_i^{k+1T}(t_1 - t_2) = \lambda_i^{k+1T}(t_1 - b_i - (t_2 - b_i)) = 0 \quad (62)$$

which proves the relationship (R3). \square

Lemma 2. *The following relationships hold at every iteration k :*

$$(R5): \quad (\nabla f_i(x_i^*) + y_i^*)^\top (x_i^* - x_i^{k+1}) + \lambda_i^{*\top} (z_i^* - z_i^{k+1}) = 0 \quad (63)$$

$$(R6): \quad \left[\nabla f_i(x_i^{k+1}) + y_i^{k+1} + \mu_i^k \left((1 - \alpha^k)x_i^{k+1} - (2 - \alpha^k)G_i w^k + G_i w^{k+1} \right) \right]^\top (x_i^{k+1} - x_i^*) \\ + \left[\lambda_i^{k+1} + \rho_i^k \left((1 - \alpha^k)z_i^{k+1} - (2 - \alpha^k)s_i^k + s_i^{k+1} \right) \right]^\top (z_i^{k+1} - z_i^*) = 0 \quad (64)$$

Proof. We start with proving relationship (R5). The KKT conditions for problem (19) can be written as follows. The point (x^*, z^*, s^*, w^*) is the optimum of problem (19) if and only if the following conditions are true:

$$\text{Optimality for } x_i: \quad \nabla f_i(x_i^*) + A_i^\top \nu_i^* + y_i^* = 0 \quad (65a)$$

$$\text{Optimality for } z_i: \quad -\nu_i^* + \lambda_i^* = 0 \quad (65b)$$

$$\text{Optimality for } s_i: \quad \lambda_i^* \in \mathcal{N}_{\mathcal{C}_i}(s_i^*) \Leftrightarrow \lambda_i^{*\top} (s_i - s_i^*) \leq 0 \quad \forall s_i \in \mathcal{C}_i \quad (65c)$$

$$\text{Optimality for } w: \quad \sum_{i \in \mathcal{V}} G_i^\top y_i^* = 0 \quad (65d)$$

$$\text{Constraints feasibility:} \quad \tilde{z}_i^* = s_i^* \quad (65e)$$

$$x_i^* = G_i w^* \quad (65f)$$

$$A_i x_i^* = z_i^* \quad (65g)$$

$$s_i \in \mathcal{C}_i \quad (65h)$$

From (65a), we have

$$(\nabla f_i(x_i^*) + A_i^\top \nu_i^* + y_i^*)^\top (x_i^* - x_i^{k+1}) = 0 \quad (66)$$

and similarly from (65b), we get

$$(-\nu_i^* + \lambda_i^*)^\top (z_i^* - z_i^{k+1}) = 0. \quad (67)$$

Adding the above two equations, we get

$$(\nabla f_i(x_i^*) + A_i^\top \nu_i^* + y_i^*)^\top (x_i^* - x_i^{k+1}) + (-\nu_i^* + \lambda_i^*)^\top (z_i^* - z_i^{k+1}) = 0 \quad (68)$$

which yields

$$(\nabla f_i(\mathbf{x}_i^*) + \mathbf{y}_i^*)^\top (\mathbf{x}_i^* - \mathbf{x}_i^{k+1}) + \boldsymbol{\lambda}_i^{*\top} (\mathbf{z}_i^* - \mathbf{z}_i^{k+1}) + \boldsymbol{\nu}_i^{*\top} (\mathbf{A}_i(\mathbf{x}_i^* - \mathbf{x}_i^{k+1}) - (\mathbf{z}_i^* - \mathbf{z}_i^{k+1})) = 0. \quad (69)$$

Using (65g) and the fact that $\mathbf{A}_i \mathbf{x}_i^{k+1} - \mathbf{z}_i^{k+1} = \mathbf{0}$, we can rewrite the above as follows

$$(\nabla f_i(\mathbf{x}_i^*) + \mathbf{y}_i^*)^\top (\mathbf{x}_i^* - \mathbf{x}_i^{k+1}) + \boldsymbol{\lambda}_i^{*\top} (\mathbf{z}_i^* - \mathbf{z}_i^{k+1}) = 0 \quad (70)$$

which yields (R5).

Subsequently, we proceed with proving relationship (R6). The KKT conditions for the $(k+1)$ -th update of $\mathbf{x}_i, \mathbf{z}_i$ are given by

$$\text{Optimality for } \mathbf{x}_i: \quad \nabla f_i(\mathbf{x}_i^{k+1}) + \mathbf{A}_i^\top \boldsymbol{\nu}_i^{k+1} + \mu_i^k (\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k + \mathbf{y}_i^k / \mu_i^k) = \mathbf{0} \quad (71a)$$

$$\text{Optimality for } \mathbf{z}_i: \quad -\boldsymbol{\nu}_i^{k+1} + \rho_i^k (\mathbf{z}_i^{k+1} - \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k) = \mathbf{0} \quad (71b)$$

$$\text{Constraints feasibility:} \quad \mathbf{A}_i \mathbf{x}_i^{k+1} = \mathbf{z}_i^{k+1} \quad (71c)$$

From (71a), we have

$$\left[\nabla f_i(\mathbf{x}_i^{k+1}) + \mathbf{A}_i^\top \boldsymbol{\nu}_i^{k+1} + \mu_i^k (\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k + \mathbf{y}_i^k / \mu_i^k) \right]^\top (\mathbf{x}_i^{k+1} - \mathbf{x}_i^*) = 0 \quad (72)$$

We simplify the term $\mu_i^k (\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k + \mathbf{y}_i^k / \mu_i^k)$ using (8) as follows

$$\begin{aligned} \mu_i^k (\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k + \mathbf{y}_i^k / \mu_i^k) &= \\ &= \mu_i^k \left(\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k + \mathbf{y}_i^{k+1} / \mu_i^k - \left(\alpha^k \mathbf{x}_i^{k+1} + (1 - \alpha^k) \mathbf{G}_i \mathbf{w}^k - \mathbf{G}_i \mathbf{w}^{k+1} \right) \right) \\ &= \mathbf{y}_i^{k+1} + \mu_i^k \left(\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k - \alpha^k \mathbf{x}_i^{k+1} - (1 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} \right) \\ &= \mathbf{y}_i^{k+1} + \mu_i^k \left((1 - \alpha^k) \mathbf{x}_i^{k+1} - (2 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} \right) \end{aligned} \quad (73)$$

such that (72) can be rewritten as

$$\begin{aligned} &\left[\nabla f_i(\mathbf{x}_i^{k+1}) + \mathbf{A}_i^\top \boldsymbol{\nu}_i^{k+1} + \mathbf{y}_i^{k+1} \right. \\ &\quad \left. + \mu_i^k \left((1 - \alpha^k) \mathbf{x}_i^{k+1} - (2 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} \right) \right]^\top (\mathbf{x}_i^{k+1} - \mathbf{x}_i^*) = 0 \end{aligned} \quad (74)$$

From (71b), we get

$$\left[-\boldsymbol{\nu}_i^{k+1} + \rho_i^k (\mathbf{z}_i^{k+1} - \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k) \right]^\top (\mathbf{z}_i^{k+1} - \mathbf{z}_i^*) = 0. \quad (75)$$

We simplify the term $\rho_i^k (\mathbf{z}_i^{k+1} - \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k)$ using (7) as follows

$$\begin{aligned} \rho_i^k (\mathbf{z}_i^{k+1} - \mathbf{s}_i^k + \boldsymbol{\lambda}_i^k / \rho_i^k) &= \rho_i^k \left(\mathbf{z}_i^{k+1} - \mathbf{s}_i^k + \boldsymbol{\lambda}_i^{k+1} / \rho_i^k - \left(\alpha^k \mathbf{z}_i^{k+1} + (1 - \alpha^k) \mathbf{s}_i^k - \mathbf{s}_i^{k+1} \right) \right) \\ &= \boldsymbol{\lambda}_i^{k+1} + \rho_i^k (\mathbf{z}_i^{k+1} - \mathbf{s}_i^k - \alpha^k \mathbf{z}_i^{k+1} - (1 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1}) \\ &= \boldsymbol{\lambda}_i^{k+1} + \rho_i^k \left((1 - \alpha^k) \mathbf{z}_i^{k+1} - (2 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1} \right) \end{aligned} \quad (76)$$

such that equation 75 can be rewritten as follows

$$\left[-\boldsymbol{\nu}_i^{k+1} + \boldsymbol{\lambda}_i^{k+1} + \rho_i^k \left((1 - \alpha^k) \mathbf{z}_i^{k+1} - (2 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1} \right) \right]^\top (\mathbf{z}_i^{k+1} - \mathbf{z}_i^*) = 0 \quad (77)$$

Combining (74) and (77) and using (65g) and the fact that $\mathbf{A}_i \mathbf{x}_i^{k+1} - \mathbf{z}_i^{k+1} = \mathbf{0}$, we get

$$\begin{aligned} &\left[\nabla f_i(\mathbf{x}_i^{k+1}) + \mathbf{y}_i^{k+1} + \mu_i^k \left((1 - \alpha^k) \mathbf{x}_i^{k+1} - (2 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} \right) \right]^\top (\mathbf{x}_i^{k+1} - \mathbf{x}_i^*) \\ &\quad + \left[\boldsymbol{\lambda}_i^{k+1} + \rho_i^k \left((1 - \alpha^k) \mathbf{z}_i^{k+1} - (2 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1} \right) \right]^\top (\mathbf{z}_i^{k+1} - \mathbf{z}_i^*) = 0 \end{aligned} \quad (78)$$

which yields relationship (R6). \square

Lemma 3. For $\alpha^k > 0$,

$$\begin{aligned}
 \text{(R7): } & \left(\mathbf{y}_i^{k+1} - \mathbf{y}_i^* + \mu_i^k ((1 - \alpha^k) \mathbf{x}_i^{k+1} - (2 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1}) \right)^\top (\mathbf{x}_i^{k+1} - \mathbf{x}_i^*) \\
 &= \frac{1}{2\alpha^k \mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 \\
 &+ \frac{(2 - \alpha^k) \mu_i^k}{2(\alpha^k)^2} \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{\mu_i^k}{2\alpha^k} (\|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2 \\
 &- \|\mathbf{G}_i(\mathbf{w}^k - \mathbf{w}^*)\|^2) + \frac{1}{\alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top \mathbf{G}_i(\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*) \\
 &+ \frac{1}{(\alpha^k)^2} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i((2 - \alpha^k) \mathbf{w}^{k+1} - (1 + (1 - \alpha^k)^2) \mathbf{w}^k - \alpha^k (1 - \alpha^k) \mathbf{w}^*)
 \end{aligned} \tag{79}$$

$$\begin{aligned}
 \text{(R8): } & \left(\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^* + \rho_i^k ((1 - \alpha^k) \mathbf{z}_i^{k+1} - (2 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1}) \right)^\top (\mathbf{z}_i^{k+1} - \mathbf{z}_i^*) \\
 &= \frac{1}{2\alpha^k \rho_i^k} (\|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 - \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \\
 &+ \frac{\rho_i^k}{2\alpha^k} (\|\mathbf{s}_i^{k+1} - \mathbf{s}_i^*\|^2 - \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2) + \frac{(2 - \alpha^k) \rho_i^k}{2(\alpha^k)^2} \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \\
 &+ \frac{1}{\alpha^k} (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (-(1 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1} - \alpha^k \mathbf{s}_i^*)
 \end{aligned} \tag{80}$$

Proof. Let us simplify the individual terms of the LHS of the relationship (R7). For that, we start by rewriting the term $\mathbf{x}_i^{k+1} - \mathbf{x}_i^*$ as follows using the relationship (R2) (i.e., (43)). (It should be noted that we consider $\alpha^k > 0$, thus making the division by α^k possible.)

$$\mathbf{x}_i^{k+1} - \mathbf{x}_i^* = \frac{1}{\alpha^k} \left(\frac{1}{\mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) - (1 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} - \alpha^k \mathbf{x}_i^* \right) \tag{81}$$

Using (65d), we can rewrite the above as following.

$$\mathbf{x}_i^{k+1} - \mathbf{x}_i^* = \frac{1}{\alpha^k} \left(\frac{1}{\mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) - (1 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} - \alpha^k \mathbf{G}_i \mathbf{w}^* \right) \tag{82}$$

which can be written in simplified form as

$$\mathbf{x}_i^{k+1} - \mathbf{x}_i^* = \frac{1}{\alpha^k \mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{1}{\alpha^k} \mathbf{G}_i (\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*). \tag{83}$$

Let us now simplify the following term in the LHS of the relationship (R7).

$$(1 - \alpha^k) \mathbf{x}_i^{k+1} - (2 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} = (1 - \alpha^k) (\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k) + \mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k) \tag{84}$$

We further simplify the term $(\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k)$ using the relationship (R2) (i.e., (43)) as follows.

$$\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k = \frac{1}{\alpha^k} \left(\frac{1}{\mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) - (1 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} \right) - \mathbf{G}_i \mathbf{w}^k \tag{85}$$

which can be written in a simplified form as

$$\mathbf{x}_i^{k+1} - \mathbf{G}_i \mathbf{w}^k = \frac{1}{\mu_i^k \alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{1}{\alpha^k} \mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k). \tag{86}$$

Substituting (86) in (84), we get

$$(1 - \alpha^k) \mathbf{x}_i^{k+1} - (2 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1} = \frac{(1 - \alpha^k)}{\mu_i^k \alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{1}{\alpha^k} \mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k) \tag{87}$$

Using the above result, we rewrite the following term on the LHS of the relationship (R7).

$$\begin{aligned} \mathbf{y}_i^{k+1} - \mathbf{y}_i^* + \mu_i^k ((1 - \alpha^k) \mathbf{x}_i^{k+1} - (2 - \alpha^k) \mathbf{G}_i \mathbf{w}^k + \mathbf{G}_i \mathbf{w}^{k+1}) \\ = \mathbf{y}_i^{k+1} - \mathbf{y}_i^* + \frac{(1 - \alpha^k)}{\alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{\mu_i^k}{\alpha^k} \mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k) \end{aligned} \quad (88)$$

For notational simplicity, let us consider the LHS of the relationship (R7) as $LHS(R7)$. Using (88) and (83), we get

$$\begin{aligned} LHS(R7) = \left(\mathbf{y}_i^{k+1} - \mathbf{y}_i^* + \frac{(1 - \alpha^k)}{\alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{\mu_i^k}{\alpha^k} \mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k) \right)^\top \\ \left(\frac{1}{\alpha^k \mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{1}{\alpha^k} \mathbf{G}_i (\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*) \right) \end{aligned} \quad (89)$$

which can be further rewritten as

$$\begin{aligned} LHS(R7) = \frac{1}{\alpha^k \mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{1}{\alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top \mathbf{G}_i (\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k \\ - \alpha^k \mathbf{w}^*) + \frac{(1 - \alpha^k)}{(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \frac{(1 - \alpha^k)}{(\alpha^k)^2} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i (\mathbf{w}^{k+1} \\ - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*) + \frac{1}{(\alpha^k)^2} (\mathbf{w}^{k+1} - \mathbf{w}^k)^\top \mathbf{G}_i^\top (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) \\ + \frac{\mu_i^k}{(\alpha^k)^2} (\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k))^\top \mathbf{G}_i (\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*) \end{aligned} \quad (90)$$

Let us now simplify each term on the RHS of the above equation. We start with the terms including only the variables \mathbf{y}_i^{k+1} , \mathbf{y}_i^k and \mathbf{y}_i^* . Using the fact that $a^\top b = \frac{1}{2}(\|a\|^2 + \|b\|^2 - \|a - b\|^2)$, we get

$$\frac{1}{\alpha^k \mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) = \frac{1}{2\alpha^k \mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 + \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) \quad (91)$$

Using the above result, we can write

$$\begin{aligned} \frac{1}{\alpha^k \mu_i^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) + \frac{(1 - \alpha^k)}{(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 \\ = \frac{1}{2\alpha^k \mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 + \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \frac{(1 - \alpha^k)}{(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 \\ = \frac{1}{2\alpha^k \mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 \end{aligned} \quad (92)$$

Next, we consider the following terms in the RHS of (90) involving only the variables \mathbf{w}^{k+1} , \mathbf{w}^k and \mathbf{w}^* .

$$\begin{aligned} \frac{\mu_i^k}{(\alpha^k)^2} (\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k))^\top \mathbf{G}_i (\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*) \\ = \frac{(1 - \alpha^k) \mu_i^k}{(\alpha^k)^2} \|\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{\mu_i^k}{\alpha^k} (\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k))^\top (\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^*)) \end{aligned} \quad (93)$$

Using the similar approach that is used to derive (92), we derive the following.

$$\begin{aligned} \frac{(1 - \alpha^k) \mu_i^k}{(\alpha^k)^2} \|\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{\mu_i^k}{\alpha^k} (\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k))^\top (\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^*)) \\ = \frac{(2 - \alpha^k) \mu_i^k}{2(\alpha^k)^2} \|\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{\mu_i^k}{2\alpha^k} (\|\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2 - \|\mathbf{G}_i (\mathbf{w}^k - \mathbf{w}^*)\|^2) \end{aligned} \quad (94)$$

Now, we will consider the following terms from the rest of the terms on the RHS of (90) as follows.

$$\begin{aligned}
& \frac{(1 - \alpha^k)}{(\alpha^k)^2} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i (\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*) \\
& \quad + \frac{1}{(\alpha^k)^2} (\mathbf{w}^{k+1} - \mathbf{w}^k)^\top \mathbf{G}_i^\top (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) \\
& = \frac{1}{(\alpha^k)^2} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i ((1 - \alpha^k) \mathbf{w}^{k+1} - (1 - \alpha^k)^2 \mathbf{w}^k - \alpha^k (1 - \alpha^k) \mathbf{w}^* + \mathbf{w}^{k+1} - \mathbf{w}^k) \\
& = \frac{1}{(\alpha^k)^2} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i ((2 - \alpha^k) \mathbf{w}^{k+1} - (1 + (1 - \alpha^k)^2) \mathbf{w}^k - \alpha^k (1 - \alpha^k) \mathbf{w}^*) \quad (95)
\end{aligned}$$

Substituting (92), (93), (94), and (95) in (90), we get

$$\begin{aligned}
LHS(R7) & = \frac{1}{2\alpha^k \mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 \\
& \quad + \frac{(2 - \alpha^k) \mu_i^k}{2(\alpha^k)^2} \|\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{\mu_i^k}{2\alpha^k} (\|\mathbf{G}_i (\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2 \\
& \quad - \|\mathbf{G}_i (\mathbf{w}^k - \mathbf{w}^*)\|^2) + \frac{1}{\alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top \mathbf{G}_i (\mathbf{w}^{k+1} - (1 - \alpha^k) \mathbf{w}^k - \alpha^k \mathbf{w}^*) \\
& \quad + \frac{1}{(\alpha^k)^2} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i ((2 - \alpha^k) \mathbf{w}^{k+1} - (1 + (1 - \alpha^k)^2) \mathbf{w}^k - \alpha^k (1 - \alpha^k) \mathbf{w}^*) \quad (96)
\end{aligned}$$

which proves the relationship (R7).

Subsequently, we prove the relation (R8). Using similar steps as in the derivation of the relationship (R7), we can derive the following.

$$\begin{aligned}
& \left(\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^* + \rho_i^k ((1 - \alpha^k) \mathbf{z}_i^{k+1} - (2 - \alpha^k) \mathbf{s}_i^k + \mathbf{s}_i^{k+1}) \right)^\top (\mathbf{z}_i^{k+1} - \mathbf{z}_i^*) \\
& = \frac{1}{2\alpha^k \rho_i^k} (\|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 - \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \\
& \quad + \frac{\rho_i^k}{2\alpha^k} (\|\mathbf{s}_i^{k+1} - \mathbf{s}_i^*\|^2 - \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2) + \frac{(2 - \alpha^k) \rho_i^k}{2(\alpha^k)^2} \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \quad (97) \\
& \quad + \frac{1}{\alpha^k} (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (\mathbf{s}_i^{k+1} - (1 - \alpha^k) \mathbf{s}_i^k - \alpha^k \mathbf{s}_i^*) \\
& \quad + \frac{1}{(\alpha^k)^2} (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k)^\top ((2 - \alpha^k) \mathbf{s}_i^{k+1} - (1 + (1 - \alpha^k)^2) \mathbf{s}_i^k - \alpha^k (1 - \alpha^k) \mathbf{s}_i^*)
\end{aligned}$$

Let us now simplify the last term of the RHS of the above equation as follows.

$$\begin{aligned}
& (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k)^\top ((2 - \alpha^k) \mathbf{s}_i^{k+1} - (1 + (1 - \alpha^k)^2) \mathbf{s}_i^k - \alpha^k (1 - \alpha^k) \mathbf{s}_i^*) \\
& = (1 + (1 - \alpha^k)^2) (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k)^\top (\mathbf{s}_i^{k+1} - \mathbf{s}_i^k) + \alpha^k (1 - \alpha^k) (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k)^\top (\mathbf{s}_i^{k+1} - \mathbf{s}_i^*) \quad (98)
\end{aligned}$$

From (5) and (65h), we have that the vectors $\mathbf{s}_i^k, \mathbf{s}_i^{k+1}, \mathbf{s}_i^* \in \mathcal{C}_i$. Using the relationship (R4) (i.e., (45)), the above equation gives us the following.

$$\begin{aligned}
& (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k)^\top ((2 - \alpha^k) \mathbf{s}_i^{k+1} - (1 + (1 - \alpha^k)^2) \mathbf{s}_i^k - \alpha^k (1 - \alpha^k) \mathbf{s}_i^*) \\
& = (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k)^\top ((2 - \alpha^k) \mathbf{s}_i^{k+1} - (2 + (\alpha^k)^2 - 2\alpha^k) \mathbf{s}_i^k + (-\alpha^k + (\alpha^k)^2) \mathbf{s}_i^*) \\
& = 0 \quad (99)
\end{aligned}$$

Substituting the above result in (97), we get

$$\begin{aligned}
& \left(\lambda_i^{k+1} - \lambda_i^* + \rho_i^k ((1 - \alpha^k) z_i^{k+1} - (2 - \alpha^k) s_i^k + s_i^{k+1}) \right)^\top (z_i^{k+1} - z_i^*) \\
&= \frac{1}{2\alpha^k \rho_i^k} (\|\lambda_i^{k+1} - \lambda_i^*\|^2 - \|\lambda_i^k - \lambda_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \rho_i^k} \|\lambda_i^{k+1} - \lambda_i^k\|^2 \\
&+ \frac{\rho_i^k}{2\alpha^k} (\|s_i^{k+1} - s_i^*\|^2 - \|s_i^k - s_i^*\|^2) + \frac{(2 - \alpha^k) \rho_i^k}{2(\alpha^k)^2} \|s_i^{k+1} - s_i^k\|^2 \\
&+ \frac{1}{\alpha^k} (\lambda_i^{k+1} - \lambda_i^*)^\top (s_i^{k+1} - (1 - \alpha^k) s_i^k - \alpha^k s_i^*)
\end{aligned} \tag{100}$$

which proves the relationship (R8). \square

Lemma 4. For $\alpha^k \geq 1$,

$$\begin{aligned}
& \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} (\|y_i^{k+1} - y_i^*\|^2 - \|y_i^k - y_i^*\|^2) + \mu_i^k (\|G_i(w^{k+1} - w^*)\|^2 - \|G_i(w^k - w^*)\|^2) \right. \\
& \quad \left. + \frac{1}{\rho_i^k} (\|\lambda_i^{k+1} - \lambda_i^*\|^2 - \|\lambda_i^k - \lambda_i^*\|^2) + \rho_i^k (\|s_i^{k+1} - s_i^*\|^2 - \|s_i^k - s_i^*\|^2) \right) \\
& \leq -\frac{(2 - \alpha^k)}{\alpha^k} \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|y_i^{k+1} - y_i^k\|^2 + \mu_i^k \|G_i(w^{k+1} - w^k)\|^2 + \frac{1}{\rho_i^k} \|\lambda_i^{k+1} - \lambda_i^k\|^2 \right. \\
& \quad \left. + \rho_i^k \|s_i^{k+1} - s_i^k\|^2 \right)
\end{aligned} \tag{101}$$

Proof. We start by combining the relationships (R5) and (R6) (i.e., (63) and (64)) to get the following equation.

$$\begin{aligned}
& \left(y_i^{k+1} - y_i^* + \mu_i^k ((1 - \alpha^k) x_i^{k+1} - (2 - \alpha^k) G_i w^k + G_i w^{k+1}) \right)^\top (x_i^{k+1} - x_i^*) \\
& + \left(\lambda_i^{k+1} - \lambda_i^* + \rho_i^k ((1 - \alpha^k) z_i^{k+1} - (2 - \alpha^k) s_i^k + s_i^{k+1}) \right)^\top (z_i^{k+1} - z_i^*) \\
& = -(\nabla f_i(x_i^{k+1}) - \nabla f_i(x_i^*))^\top (x_i^{k+1} - x_i^*)
\end{aligned} \tag{102}$$

Since f_i is convex, the following holds true.

$$(\nabla f_i(x_i^{k+1}) - \nabla f_i(x_i^*))^\top (x_i^{k+1} - x_i^*) \geq 0 \tag{103}$$

Using the above inequality, we can rewrite (102) as follows.

$$\begin{aligned}
& \left(y_i^{k+1} - y_i^* + \mu_i^k ((1 - \alpha^k) x_i^{k+1} - (2 - \alpha^k) G_i w^k + G_i w^{k+1}) \right)^\top (x_i^{k+1} - x_i^*) \\
& + \left(\lambda_i^{k+1} - \lambda_i^* + \rho_i^k ((1 - \alpha^k) z_i^{k+1} - (2 - \alpha^k) s_i^k + s_i^{k+1}) \right)^\top (z_i^{k+1} - z_i^*) \leq 0
\end{aligned} \tag{104}$$

Adding the above result over all the agents $i \in \mathcal{V}$, we get

$$\begin{aligned}
& \sum_{i \in \mathcal{V}} \left(y_i^{k+1} - y_i^* + \mu_i^k ((1 - \alpha^k) x_i^{k+1} - (2 - \alpha^k) G_i w^k + G_i w^{k+1}) \right)^\top (x_i^{k+1} - x_i^*) \\
& + \sum_{i \in \mathcal{V}} \left(\lambda_i^{k+1} - \lambda_i^* + \rho_i^k ((1 - \alpha^k) z_i^{k+1} - (2 - \alpha^k) s_i^k + s_i^{k+1}) \right)^\top (z_i^{k+1} - z_i^*) \leq 0
\end{aligned} \tag{105}$$

Now, we use the relationships (R7) and (R8) (i.e., (79) and (80)) to rewrite the above equation as following.

$$\begin{aligned}
0 \geq \sum_{i \in \mathcal{V}} & \left(\frac{1}{2\alpha^k \mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 \right. \\
& + \frac{(2 - \alpha^k) \mu_i^k}{2(\alpha^k)^2} \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{\mu_i^k}{2\alpha^k} (\|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2 \\
& - \|\mathbf{G}_i(\mathbf{w}^k - \mathbf{w}^*)\|^2) + \frac{1}{\alpha^k} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top \mathbf{G}_i(\mathbf{w}^{k+1} - (1 - \alpha^k)\mathbf{w}^k - \alpha^k \mathbf{w}^*) \\
& + \frac{1}{(\alpha^k)^2} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i((2 - \alpha^k)\mathbf{w}^{k+1} - (1 + (1 - \alpha^k)^2)\mathbf{w}^k - \alpha^k(1 - \alpha^k)\mathbf{w}^*) \\
& + \frac{1}{2\alpha^k \rho_i^k} (\|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 - \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2) + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \\
& + \frac{\rho_i^k}{2\alpha^k} (\|\mathbf{s}_i^{k+1} - \mathbf{s}_i^*\|^2 - \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2) + \frac{(2 - \alpha^k) \rho_i^k}{2(\alpha^k)^2} \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \\
& \left. + \frac{1}{\alpha^k} (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (-(1 - \alpha^k)\mathbf{s}_i^k + \mathbf{s}_i^{k+1} - \alpha^k \mathbf{s}_i^*) \right)
\end{aligned} \tag{106}$$

Let us now further simplify the terms on the RHS of the above equation. For that, let us start with the last term on the RHS. We have

$$\begin{aligned}
(\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (-(1 - \alpha^k)\mathbf{s}_i^k + \mathbf{s}_i^{k+1} - \alpha^k \mathbf{s}_i^*) &= (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (\mathbf{s}_i^{k+1} - \mathbf{s}_i^*) \\
&\quad - (1 - \alpha^k) (\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (\mathbf{s}_i^k - \mathbf{s}_i^*)
\end{aligned} \tag{107}$$

Using the relationship (R4) (i.e., (45)), and (65c), and using the fact that $\mathbf{s}_i^k, \mathbf{s}_i^{k+1}, \mathbf{s}_i^* \in \mathcal{C}_i$, we get

$$(\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (\mathbf{s}_i^{k+1} - \mathbf{s}_i^*) \geq 0, \tag{108}$$

$$(\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (\mathbf{s}_i^k - \mathbf{s}_i^*) \geq 0. \tag{109}$$

Thus, for $\alpha^k \geq 1$, combining (107), (108), and (109), we get

$$(\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*)^\top (-(1 - \alpha^k)\mathbf{s}_i^k + \mathbf{s}_i^{k+1} - \alpha^k \mathbf{s}_i^*) \geq 0. \tag{110}$$

Now, the following results hold based on the relationship (R1) (i.e., (42)) and (65d).

$$\sum_{i \in \mathcal{V}} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^*)^\top \mathbf{G}_i = 0, \quad \sum_{i \in \mathcal{V}} (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k)^\top \mathbf{G}_i = 0. \tag{111}$$

By substituting (110) and (111) in (106), and by rearranging the terms, we get

$$\begin{aligned}
& \sum_{i \in \mathcal{V}} \left(\frac{1}{2\alpha^k \mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \frac{\mu_i^k}{2\alpha^k} (\|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2 - \|\mathbf{G}_i(\mathbf{w}^k - \mathbf{w}^*)\|^2) \right. \\
& \quad \left. + \frac{1}{2\alpha^k \rho_i^k} (\|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 - \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2) + \frac{\rho_i^k}{2\alpha^k} (\|\mathbf{s}_i^{k+1} - \mathbf{s}_i^*\|^2 - \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2) \right) \\
& \leq - \sum_{i \in \mathcal{V}} \left(\frac{(2 - \alpha^k)}{2(\alpha^k)^2 \mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \frac{(2 - \alpha^k) \mu_i^k}{2(\alpha^k)^2} \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 \right. \\
& \quad \left. + \frac{(2 - \alpha^k)}{2(\alpha^k)^2 \rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 + \frac{(2 - \alpha^k) \rho_i^k}{2(\alpha^k)^2} \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \right)
\end{aligned} \tag{112}$$

Since, we consider $\alpha^k \geq 1$, we can multiply the above equation with $2\alpha^k$ to obtain the following.

$$\begin{aligned}
& \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \mu_i^k (\|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2 - \|\mathbf{G}_i(\mathbf{w}^k - \mathbf{w}^*)\|^2) \right. \\
& \quad \left. + \frac{1}{\rho_i^k} (\|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 - \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2) + \rho_i^k (\|\mathbf{s}_i^{k+1} - \mathbf{s}_i^*\|^2 - \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2) \right) \\
& \leq -\frac{(2 - \alpha^k)}{\alpha^k} \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \mu_i^k \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \right. \\
& \quad \left. + \rho_i^k \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \right)
\end{aligned} \tag{113}$$

□

C.3 PROOF OF THEOREM 1

Let us first rewrite the relation (101) derived in Lemma 4 for $\alpha^k \in [1, 2)$, as follows.

$$\begin{aligned}
& \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} (\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2) + \mu_i^k (\|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2 - \|\mathbf{G}_i(\mathbf{w}^k - \mathbf{w}^*)\|^2) \right. \\
& \quad \left. + \frac{1}{\rho_i^k} (\|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2 - \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2) + \rho_i^k (\|\mathbf{s}_i^{k+1} - \mathbf{s}_i^*\|^2 - \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2) \right) \\
& \leq -\frac{(2 - \alpha^k)}{\alpha^k} \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \mu_i^k \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \right. \\
& \quad \left. + \rho_i^k \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \right)
\end{aligned} \tag{114}$$

which can be rearranged to give the following.

$$\begin{aligned}
& \frac{(2 - \alpha^k)}{\alpha^k} \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \mu_i^k \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \right. \\
& \quad \left. + \rho_i^k \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \right) \\
& \leq \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} (\|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2 - \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^*\|^2) + \mu_i^k (\|\mathbf{G}_i(\mathbf{w}^k - \mathbf{w}^*)\|^2 - \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^*)\|^2) \right. \\
& \quad \left. + \frac{1}{\rho_i^k} (\|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 - \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^*\|^2) + \rho_i^k (\|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2 - \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^*\|^2) \right)
\end{aligned} \tag{115}$$

For convenience, let us define for each iteration k , the terms η_i^k , $i \in \mathcal{V}$, and η^k such that

$$\eta_i^k + 1 = \max \left(\frac{\rho_i^k}{\rho_i^{k-1}}, \frac{\rho_i^{k-1}}{\rho_i^k}, \frac{\mu_i^k}{\mu_i^{k-1}}, \frac{\mu_i^{k-1}}{\mu_i^k} \right), \quad \eta_{\max}^k = \max_{i \in \mathcal{V}} \eta_i^k, \tag{116}$$

and the term V^k as

$$\begin{aligned}
V^k = \sum_{i \in \mathcal{V}} & \left(\frac{1}{\mu_i^{k-1}} \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2 + \mu_i^{k-1} \|\mathbf{G}_i(\mathbf{w}^k - \mathbf{w}^*)\|^2 + \frac{1}{\rho_i^{k-1}} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 \right. \\
& \quad \left. + \rho_i^{k-1} \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2 \right).
\end{aligned} \tag{117}$$

Based on the definition of η_i^k in (116), we can write

$$\begin{aligned} & \frac{1}{\mu_i^k} \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2 + \mu_i^k \|G_i(\mathbf{w}^k - \mathbf{w}^*)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 + \rho_i^k \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2 \\ & \leq (\eta_i^k + 1) \left(\frac{1}{\mu_i^{k-1}} \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2 + \mu_i^{k-1} \|G_i(\mathbf{w}^k - \mathbf{w}^*)\|^2 + \frac{1}{\rho_i^{k-1}} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 \right. \\ & \quad \left. + \rho_i^{k-1} \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2 \right) \end{aligned} \quad (118)$$

Further, by adding the above result over all the agents $i \in \mathcal{V}$, and using the fact that $\eta_{\max}^k \geq \eta_i^k$ for all i , we get

$$\begin{aligned} & \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2 + \mu_i^k \|G_i(\mathbf{w}^k - \mathbf{w}^*)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 + \rho_i^k \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2 \right) \\ & \leq \sum_{i \in \mathcal{V}} (\eta_i^k + 1) \left(\frac{1}{\mu_i^{k-1}} \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2 + \mu_i^{k-1} \|G_i(\mathbf{w}^k - \mathbf{w}^*)\|^2 + \frac{1}{\rho_i^{k-1}} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 \right. \\ & \quad \left. + \rho_i^{k-1} \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2 \right) \\ & \leq (\eta_{\max}^k + 1) \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^{k-1}} \|\mathbf{y}_i^k - \mathbf{y}_i^*\|^2 + \mu_i^{k-1} \|G_i(\mathbf{w}^k - \mathbf{w}^*)\|^2 + \frac{1}{\rho_i^{k-1}} \|\boldsymbol{\lambda}_i^k - \boldsymbol{\lambda}_i^*\|^2 \right. \\ & \quad \left. + \rho_i^{k-1} \|\mathbf{s}_i^k - \mathbf{s}_i^*\|^2 \right) \\ & = (\eta_{\max}^k + 1) V^k \end{aligned} \quad (119)$$

Substituting the above result in (115), we get

$$\begin{aligned} & \frac{(2 - \alpha^k)}{\alpha^k} \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \mu_i^k \|G_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \right. \\ & \quad \left. + \rho_i^k \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \right) \leq (\eta_{\max}^k + 1) V^k - V^{k+1} \end{aligned} \quad (120)$$

Now that we have derived the above relation, we need to next prove that V^k is bounded. By the definition of V^k , we have that V^k is lower bounded by zero. Thus, we now prove that V^k is upper bounded. From (120), we have

$$V^{k+1} \leq (\eta_{\max}^k + 1) V^k, \quad (121)$$

which leads to the following relation

$$V^{k+1} \leq \prod_{l=1}^k (\eta_{\max}^l + 1) V^1 \quad (122)$$

It should be noted that based on the assumption 1, we have $(\eta_{\max}^k + 1) \rightarrow 1$, as $k \rightarrow \infty$. Based on this condition, (122) implies that V^{k+1} is upper bounded for all k , and there exists V_{\max} such that

$$V^k \leq V_{\max} < \infty, \quad \text{for all } k \quad (123)$$

Let us now consider adding the result (120) over k as follows.

$$\begin{aligned} & \sum_{k=1}^{\infty} \frac{(2 - \alpha^k)}{\alpha^k} \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \mu_i^k \|G_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \right. \\ & \quad \left. + \rho_i^k \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \right) \leq \sum_{k=1}^{\infty} (\eta_{\max}^k + 1) V^k - V^{k+1} \end{aligned} \quad (124)$$

The term on the RHS of the above equation can be further simplified as follows.

$$\sum_{k=1}^{\infty} (\eta_{\max}^k + 1) V^k - V^{k+1} = \sum_{k=1}^{\infty} \eta_{\max}^k V^k + \sum_{k=1}^{\infty} (V^k - V^{k+1}) = V^1 - V^{\infty} + \sum_{k=1}^{\infty} \eta_{\max}^k V^k \quad (125)$$

Based on the assumption 1, we have $\eta_{\max}^k \rightarrow 0$ as $k \rightarrow \infty$, which implies the following.

$$\sum_{k=1}^{\infty} \eta_{\max}^k < \infty \quad (126)$$

Using the above fact and (123), we can upper bound $\sum_{k=1}^{\infty} \eta_{\max}^k V^k$ as follows.

$$\sum_{k=1}^{\infty} \eta_{\max}^k V^k \leq \left(\sum_{k=1}^{\infty} \eta_{\max}^k \right) V_{\max} < \infty \quad (127)$$

Using the facts that V^1 is upper bounded, and V^{∞} is lower bounded by zero, and using the above equation, we get

$$V^1 - V^{\infty} + \sum_{k=1}^{\infty} \eta_{\max}^k V^k \leq V^1 + \sum_{k=1}^{\infty} \eta_{\max}^k V^k < \infty \quad (128)$$

Thus, we can rewrite (124) as following.

$$\begin{aligned} \sum_{k=1}^{\infty} \frac{(2 - \alpha^k)}{\alpha^k} \sum_{i \in \mathcal{V}} \left(\frac{1}{\mu_i^k} \|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|^2 + \mu_i^k \|\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 + \frac{1}{\rho_i^k} \|\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k\|^2 \right. \\ \left. + \rho_i^k \|\mathbf{s}_i^{k+1} - \mathbf{s}_i^k\|^2 \right) < \infty \end{aligned} \quad (129)$$

Since $\alpha^k \in [1, 2)$, we have $\frac{(2 - \alpha^k)}{\alpha^k} > 0$ for all k . Further, we have $0 < \mu_i^k, \rho_i^k < \infty$ for all k . Thus, (129) implies the following for all $i \in \mathcal{V}$.

$$\text{As } k \rightarrow \infty, \quad (\mathbf{y}_i^{k+1} - \mathbf{y}_i^k) \rightarrow 0, \quad \mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k) \rightarrow 0 \quad (130)$$

$$(\boldsymbol{\lambda}_i^{k+1} - \boldsymbol{\lambda}_i^k) \rightarrow 0, \quad (\mathbf{s}_i^{k+1} - \mathbf{s}_i^k) \rightarrow 0 \quad (131)$$

which proves the convergence of the variables $\mathbf{y}_i, \boldsymbol{\lambda}_i$ and \mathbf{s}_i . Further, we have, as $k \rightarrow \infty$,

$$\mathbf{G}_i(\mathbf{w}^{k+1} - \mathbf{w}^k) \rightarrow 0 \quad \forall i \in \mathcal{V} \implies \mathbf{G}(\mathbf{w}^{k+1} - \mathbf{w}^k) \rightarrow 0 \quad (132)$$

Since \mathbf{G} has full column rank, the above equation implies

$$(\mathbf{w}^{k+1} - \mathbf{w}^k) \rightarrow 0 \text{ as } k \rightarrow \infty, \quad (133)$$

which proves the convergence of the global variable \mathbf{w} .

Further, using the relationships (R2) and (R3) (i.e., (43) and (44)), and the convergence results (130) and (131), we obtain the following.

$$\text{As } k \rightarrow \infty, \quad (\mathbf{x}_i^{k+1} - \mathbf{x}_i^k) \rightarrow 0, \quad (\mathbf{z}^{k+1} - \mathbf{z}^k) \rightarrow 0. \quad (134)$$

Hence, we prove the convergence of the algorithm.

Now that we have proved convergence, we can verify that the limit point of convergence is the optimal solution to the problem 19. For that, we need to check if the limit point satisfies the KKT condition (65) for the problem 19. The convergence of the dual variables \mathbf{y}_i and $\boldsymbol{\lambda}_i$, and the update steps verify that the limit points have constraint feasibility (65e - 65h). The constraint feasibility of the limit points and the optimality conditions of $(k+1)$ -th update of $\mathbf{x}_i, \mathbf{z}_i$ (71) imply that the limit points satisfy the optimality conditions (65a - 65b). Further, using relations (R1) and (R4) (i.e., (42) and (45)), we can prove that the limit points also satisfy (65c - 65d).

D DETAILS ON DEEPDISTRIBUTEDQP FEEDBACK POLICIES

In DeepDistributedQP, the penalty parameters are given by

$$\rho_i^k = \text{SoftPlus}\left(\bar{\rho}_i^k + \underbrace{\pi_{i,\rho}^k(r_{i,\rho}^k, s_{i,\rho}^k; \theta_{i,\rho}^k)}_{\hat{\rho}_i^k}\right), \quad \mu_i^k = \text{SoftPlus}\left(\bar{\mu}_i^k + \underbrace{\pi_{i,\mu}^k(r_{i,\mu}^k, s_{i,\mu}^k; \theta_{i,\mu}^k)}_{\hat{\mu}_i^k}\right) \quad (135)$$

where $\bar{\rho}_i^k, \bar{\mu}_i^k$ are learnable feed-forward parameters and $\hat{\rho}_i^k, \hat{\mu}_i^k$ and the feedback parts. The latter are obtained through the learnable policies $\pi_{i,\cdot}^k(r_{i,\cdot}^k, s_{i,\cdot}^k; \theta_{i,\cdot}^k)$ parameterized by fully-connected neural network layers with inputs $r_{i,\cdot}^k, s_{i,\cdot}^k$ and weights $\theta_{i,\cdot}^k$. The analytical expressions for $r_{i,\cdot}^k, s_{i,\cdot}^k$ are provided as follows:

$$r_{i,\rho}^k = \begin{bmatrix} \|z_i^k - s_i^k\|_2 \\ \|A_i x_i^k - s_i^k\|_2 \end{bmatrix}, \quad s_{i,\rho}^k = \begin{bmatrix} \|s_i^k - s_i^{k-1}\|_2 \\ \|Q_i x_i^k + q_i + A_i^\top \lambda_i^k\|_2 \end{bmatrix} \quad (136a)$$

$$r_{i,\mu}^k = \|x_i^k - \tilde{w}_i^k\|_2, \quad s_{i,\mu}^k = \|\tilde{w}_i^k - \tilde{w}_i^{k-1}\|_2, \quad (136b)$$

being motivated by the primal and dual residuals of ADMM (Boyd et al., 2011, Section 3) and the ones used in the OSQP algorithm (Stellato et al., 2020).

E THE CENTRALIZED VERSION: DEEPPQP

The **centralized** version of DeepDistributedQP boils down to simply unfolding the iterates of the standard OSQP algorithm for solving centralized QPs (1), while applying the same principles as in Section 4.1 for DeepDistributedQP.

For convenience, we repeat the OSQP updates from Stellato et al. (2020) here:

1. *Update for (x, z)* : Solve linear system

$$\begin{bmatrix} Q + \sigma I & A^\top \\ A & -1/\rho^k I \end{bmatrix} \begin{bmatrix} x^{k+1} \\ \nu^{k+1} \end{bmatrix} = \begin{bmatrix} \sigma t^k - q \\ s^k - 1/\rho^k \lambda^k \end{bmatrix} \quad (137)$$

and update

$$z^{k+1} = s^k + 1/\rho^k (\nu^{k+1} - \lambda^k). \quad (138)$$

As explained in Stellato et al. (2020), as the scale of the system equation 137 increases, it is often preferable to solve the following system instead,

$$(Q + \sigma I + \rho^k A^\top A) x^{k+1} = \sigma x^k - q + A^\top (\rho^k z^k - y^k), \quad (139)$$

using a method such as conjugate gradient.

2. *Update for (t, s)* :

$$t^{k+1} = \alpha^k x^{k+1} + (1 - \alpha^k) t^k \quad (140a)$$

$$s^{k+1} = \Pi_C (\alpha^k z^{k+1} + (1 - \alpha^k) s^k + \lambda^k / \rho^k) \quad (140b)$$

3. *Dual update for λ* :

$$\lambda^{k+1} = \lambda^k + \rho^k (\alpha^k z^{k+1} + (1 - \alpha^k) s^k - s^{k+1}) \quad (141)$$

The DeepQP framework then emerges through unfolding the OSQP updates following the same methodology as in DeepDistributedQP. In particular, its iterations are unrolled for a prescribed amount of K iterations as shown in Fig. 4.

F PROOF OF INDIRECT METHOD IMPLICIT DIFFERENTIATION

We start by stating the implicit function theorem, whose proof is given in (Krantz & Parks, 2002).

Lemma 5 (Implicit Function Theorem). *Let $r : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ be a continuously differentiable function. Let $(\mathbf{x}_0, \boldsymbol{\theta}_0)$ be a point such that $r(\mathbf{x}_0, \boldsymbol{\theta}_0) = 0$. If the Jacobian matrix $\frac{\partial r}{\partial \mathbf{x}}(\mathbf{x}_0, \boldsymbol{\theta}_0)$ is invertible, then there exists a function $\mathbf{x}^*(\cdot)$ defined in a neighborhood of $\boldsymbol{\theta}_0$ such that $\mathbf{x}^*(\boldsymbol{\theta}_0) = \mathbf{x}_0$, and*

$$\frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = - \left(\frac{\partial r}{\partial \mathbf{x}}(\mathbf{x}^*(\boldsymbol{\theta}), \boldsymbol{\theta}) \right)^{-1} \frac{\partial r}{\partial \boldsymbol{\theta}}(\mathbf{x}^*(\boldsymbol{\theta}), \boldsymbol{\theta}). \quad (142)$$

Proof of Theorem 2. Let $\boldsymbol{\theta} = (\bar{\mathbf{Q}}_i^k, \bar{\mathbf{b}}_i^k)$ be the concatenation of all the parameters in Eq. (11). $\bar{\mathbf{Q}}_i^k$ is always positive definite since \mathbf{Q}_i is positive definite and the penalty parameters are always non-negative. Therefore, Eq. (11) has a unique solution \mathbf{x}_i^{k+1} satisfying $r(\mathbf{x}_i^{k+1}, \boldsymbol{\theta}) := \bar{\mathbf{Q}}_i^k \mathbf{x}_i^{k+1} - \bar{\mathbf{b}}_i^k = 0$. Applying Lemma 5 to this residual function yields the relationship $\frac{\partial \mathbf{x}_i^{k+1}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = -(\bar{\mathbf{Q}}_i^k)^{-1} \frac{\partial r}{\partial \boldsymbol{\theta}}(\mathbf{x}_i^{k+1}(\boldsymbol{\theta}), \boldsymbol{\theta})$.

Now, for any downstream loss function $L(\mathbf{x}_i^{k+1}(\boldsymbol{\theta}))$, we have that

$$\nabla_{\boldsymbol{\theta}} L(\mathbf{x}_i^{k+1}(\boldsymbol{\theta})) = \frac{\partial \mathbf{x}_i^{k+1}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) \nabla_{\mathbf{x}} L(\mathbf{x}_i^{k+1}(\boldsymbol{\theta})) \quad (143)$$

$$= - \frac{\partial r}{\partial \boldsymbol{\theta}}(\mathbf{x}_i^{k+1}(\boldsymbol{\theta}), \boldsymbol{\theta})^\top (\bar{\mathbf{Q}}_i^k)^{-1} \nabla_{\mathbf{x}} L(\mathbf{x}_i^{k+1}(\boldsymbol{\theta})) \quad (144)$$

$$= \frac{\partial r}{\partial \boldsymbol{\theta}}(\mathbf{x}_i^{k+1}(\boldsymbol{\theta}), \boldsymbol{\theta})^\top d\mathbf{x}_i^{k+1}, \quad (145)$$

where $d\mathbf{x}_i^{k+1}$ is the unique solution to the linear system

$$\bar{\mathbf{Q}}_i^k d\mathbf{x}_i^{k+1} = -\nabla_{\mathbf{x}} L(\mathbf{x}_i^{k+1}(\boldsymbol{\theta})).$$

Expanding the matrix multiplication in Eq. (145) yields

$$\begin{aligned} \nabla_{\bar{\mathbf{Q}}_i^k} L &= \frac{1}{2}(\mathbf{x}_i^{k+1} \otimes d\mathbf{x}_i^{k+1} + d\mathbf{x}_i^{k+1} \otimes \mathbf{x}_i^{k+1}), \\ \nabla_{\bar{\mathbf{b}}_i^k} L &= -d\mathbf{x}_i^{k+1}. \end{aligned}$$

□

G BACKGROUND ON PAC-BAYES THEORY

Here, we provide a brief overview of PAC-Bayes theory (Alquier (2024)). Consider a bounded loss function $\ell(\zeta; \theta)$. Without loss of generality, we assume that this loss is uniformly bounded between 0 and 1. PAC-Bayes theory aims to providing a probabilistic bound for the true expected loss

$$\ell_{\mathcal{D}}(\mathcal{P}) = \mathbb{E}_{\zeta \sim \mathcal{D}} \mathbb{E}_{\theta \sim \mathcal{P}} [\ell(\zeta; \theta)], \quad (146)$$

where \mathcal{D} is the data distribution — in our case, this is the distribution optimization problems are drawn from. The empirical expected loss is given by,

$$\ell_{\mathcal{S}}(\mathcal{P}) = \mathbb{E}_{\theta \sim \mathcal{P}} \left[\frac{1}{H} \sum_{j=1}^H (\zeta^j; \theta) \right], \quad (147)$$

where $\mathcal{S} = \{\zeta^j\}_{j=1}^H$ is the training dataset consisting of H problem instances.

The PAC-Bayes framework operates by forming a bound that holds in high probability on the true loss $\ell_{\mathcal{D}}(\mathcal{P})$ in terms of the empirical loss and a the deviation between the learned policy \mathcal{P} and a prior policy \mathcal{P}_0 used to as an initial guess for \mathcal{P} . This deviation is measured using the KL divergence. Importantly, \mathcal{P}_0 need not be a Bayesian prior but can be any distribution independent of the data used to train \mathcal{P} and evaluate the sample loss. Moreover, $\ell(\zeta; \theta)$ need not be the loss used to train \mathcal{P} , but can be any bounded function. This observation is useful because, both in the literature and in the sequel, it is common to use a loss function modified for practicality during training before evaluating the bound using the loss function of interest.

Specifically, the following PAC-Bayes bounds hold with probability $1 - \delta$,

$$\ell_{\mathcal{D}}(\mathcal{P}) \leq \mathbb{D}_{\text{KL}}^{-1} \left(\ell_{\mathcal{S}}(\mathcal{P}) \parallel \frac{\mathbb{D}_{\text{KL}}(\mathcal{P} \parallel \mathcal{P}_0) + \log \frac{2\sqrt{H}}{\delta}}{H} \right) \leq \ell_{\mathcal{S}}(\mathcal{P}) + \sqrt{\frac{\mathbb{D}_{\text{KL}}(\mathcal{P} \parallel \mathcal{P}_0) + \log \frac{2\sqrt{H}}{\delta}}{2H}}, \quad (148)$$

where the $\mathbb{D}_{\text{KL}}^{-1}(p \parallel c)$ is the *inverse of the KL divergence* for Bernoulli random variables $\mathcal{B}(p), \mathcal{B}(q)$:

$$\mathbb{D}_{\text{KL}}^{-1}(p \parallel c) = \sup\{q \in [0, 1] \mid \mathbb{D}_{\text{KL}}(\mathcal{B}(p) \parallel \mathcal{B}(q)) \leq c\}. \quad (149)$$

The probability δ captures the failure case that the data set \mathcal{S} is not sufficiently representative of the data distribution \mathcal{D} . In the sequel, both of the above inequalities will be used. As the first bound is tighter, it is used to evaluate the generalization capabilities of the learned optimizer. The benefit of the second, looser, bound is that its form is convenient to use during training as a regularizer. Using both bounds in this manner is a common technique in the PAC-Bayes literature (Majumdar et al. (2021), Dziugaite & Roy (2017)).

H OPTIMIZING AND EVALUATING GENERALIZATION BOUND

Two important requirements for establishing a tight PAC-Bayes bound are selecting an informative prior and optimizing the PAC-Bayes bounds in Eq. (148) instead of simply minimizing the loss function. The choice of prior \mathcal{P}_0 is particularly important because the KL divergence is unbounded and can produce a vacuous result Dziugaite et al. (2021). While the distribution \mathcal{P}_0 need not be a Bayesian prior, it must be selected independently from the data used to optimize \mathcal{P} and evaluate the bound. To select \mathcal{P}_0 , we follow a common approach in the literature and split our training set \mathcal{S} into two disjoint subsets $\mathcal{S}_0, \mathcal{S}_1$. The prior \mathcal{P}_0 is first trained using the data set \mathcal{S}_0 and the loss $\ell(\mathcal{D}; \Theta)$ discussed in Section 4.

Subsequently, the posterior \mathcal{P} is trained by minimizing the looser (i.e., rightmost) PAC-Bayes bound in Eq. (148). This bound is used for training because it is straightforward to evaluate in comparison to computing the inverse of the KL divergence, and this objective is easily interpreted as minimizing an expected loss function with a regularizer. To evaluate the loss function in the PAC-Bayes bound, parameters are sampled from \mathcal{P} using the current network weights and an empirical average is used. Once training is complete, the PAC-Bayes bound is evaluated as described in Theorem 3, i.e., by using the tighter PAC-Bayes bound in (148) and the sample convergence bound in (15).

I DETAILS ON EXPERIMENTS

This section provides further details on the problems considered in the experiments, the training of the learned optimizers, as well as the evaluation of both learned and traditional methods.

I.1 PROBLEM TYPES IN CENTRALIZED EXPERIMENTS

Random QPs. We consider randomly generated problems of the following form

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{q}^\top \mathbf{x} \quad \text{s.t.} \quad \mathbf{A} \mathbf{x} \leq \mathbf{b}, \quad \mathbf{C} \mathbf{x} = \mathbf{d}. \quad (150)$$

For each generated problem, the cost Hessian is given by $\mathbf{Q} = \mathbf{F}^\top \mathbf{F} + \gamma \mathbf{I}$, where each element of $\mathbf{F} \in \mathbb{R}^{n \times n}$ is sampled through $F_{ij} \sim \mathcal{N}(0, 1)$ and $\gamma = 1.0$. The coefficients of \mathbf{q} are also sampled as $q_i \sim \mathcal{N}(0, 1)$. The elements of the inequality constraints matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ are given by $A_{ij} \sim \mathcal{N}(0, 1)$, while $\mathbf{b} = \mathbf{A} \boldsymbol{\theta}$, where each element of $\boldsymbol{\theta} \in \mathbb{R}^n$ is sampled through $\theta_i \sim \mathcal{N}(0, 1)$. Similarly, the elements of the equality constraints matrix $\mathbf{C} \in \mathbb{R}^{p \times n}$ are given by $C_{ij} \sim \mathcal{N}(0, 1)$, while $\mathbf{d} = \mathbf{C} \boldsymbol{\xi}$, where each element of $\boldsymbol{\xi} \in \mathbb{R}^n$ is $\xi_i \sim \mathcal{N}(0, 1)$.

For random QPs without equality constraints, we set $n = 50$, $m = 40$ and $p = 0$. For random QPs with equality constraints, we set $n = 50$, $m = 25$ and $p = 20$.

Optimal control. We consider linear optimal control problems of the following form

$$\min_{\mathbf{x}, \mathbf{u}} \sum_{t=0}^{T-1} \mathbf{x}_t^\top \mathbf{Q} \mathbf{x}_t + \mathbf{u}_t^\top \mathbf{R} \mathbf{u}_t + \mathbf{x}_T^\top \mathbf{Q}_T \mathbf{x}_T \quad (151a)$$

$$\text{s.t. } \mathbf{x}_{t+1} = \mathbf{A}_d \mathbf{x}_t + \mathbf{B}_d \mathbf{u}_t, \quad t = 0, \dots, T-1, \quad (151b)$$

$$\mathbf{A}_u \mathbf{u}_t \leq \mathbf{b}_u, \quad \mathbf{A}_x \mathbf{x}_t \leq \mathbf{b}_x, \quad t = 0, \dots, T, \quad (151c)$$

$$\mathbf{x}_0 = \bar{\mathbf{x}}_0. \quad (151d)$$

where $\mathbf{x} = \{\mathbf{x}_0, \dots, \mathbf{x}_T\}$ is the state trajectory, $\mathbf{u} = \{\mathbf{u}_0, \dots, \mathbf{u}_{T-1}\}$ is the control trajectory, $\bar{\mathbf{x}}_0$ is the given initial state condition, \mathbf{Q} and \mathbf{R} are the running state and control cost matrices, \mathbf{Q}_T is the terminal state cost matrix, \mathbf{A}_d and \mathbf{B}_d are the dynamics matrices, and finally \mathbf{A}_u , \mathbf{b}_u and \mathbf{A}_x , \mathbf{b}_x are the control and state constraints coefficients, respectively.

Both the double integrator and the mass-spring problem setups are drawn from Chen et al. (2022). For the double integrator system, we have $x_t \in \mathbb{R}^2$ and $u_t \in \mathbb{R}$, with time horizon $T = 20$. The dynamics matrices are given by

$$\mathbf{A}_d = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{B}_d = \begin{bmatrix} 0.5 \\ 0.1 \end{bmatrix} \quad (152)$$

The cost matrices are $\mathbf{Q} = \mathbf{Q}_T = \mathbf{I}_2$ and $\mathbf{R} = 1.0$. The state and control constraint coefficients are given by

$$\mathbf{A}_x = \begin{bmatrix} \mathbf{I}_2 \\ -\mathbf{I}_2 \end{bmatrix}, \quad \mathbf{b}_x = [5 \quad 1 \quad 5 \quad 1]^\top, \quad \mathbf{A}_u = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad \mathbf{b}_u = [0.1 \quad 0.1]^\top. \quad (153)$$

Finally, the initial state conditions are sampled from the uniform distribution $\mathcal{U}([-1; -0.3], [1; 0.3])$.

For the oscillating masses, we have $x_t \in \mathbb{R}^{12}$ and $u_t \in \mathbb{R}^3$, with time horizon $T = 10$. The discrete-time dynamics matrices are obtained from the continuous-time ones through Euler discretization,

$$\mathbf{A}_d = \mathbf{I} + \mathbf{A}_c \Delta t, \quad \mathbf{B}_d = \mathbf{A}_c \Delta t. \quad (154)$$

The continuous-time dynamics matrices are given by

$$\mathbf{A}_c = \begin{bmatrix} \mathbf{0}_{6 \times 6} & \mathbf{I}_6 \\ a\mathbf{I}_6 + c\mathbf{L}_6 + c\mathbf{L}_6^\top & b\mathbf{I}_6 + d\mathbf{L}_6 + d\mathbf{L}_6^\top \end{bmatrix}, \quad \mathbf{B}_c = \begin{bmatrix} \mathbf{0}_{6 \times 3} \\ \mathbf{F} \end{bmatrix} \quad (155)$$

with $c = 1.0$, $d = 0.1$, $a = -2c$, $b = -2.0$. \mathbf{L}_6 is the 6×6 lower shift matrix and

$$\mathbf{F} = [\mathbf{e}_1 \quad -\mathbf{e}_1 \quad \mathbf{e}_2 \quad \mathbf{e}_3 \quad -\mathbf{e}_2 \quad \mathbf{e}_3]^\top \quad (156)$$

where $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ are the standard basis vectors in \mathbb{R}^3 .

The timestep is set as $\Delta t = 0.5$. The cost matrices are $\mathbf{Q} = \mathbf{Q}_T = \mathbf{I}_{12}$ and $\mathbf{R} = \mathbf{I}_3$. The state and control constraints are defined through

$$\mathbf{A}_x = \begin{bmatrix} \mathbf{I}_{12} \\ -\mathbf{I}_{12} \end{bmatrix}, \quad \mathbf{b}_x = 4 \cdot \mathbf{1}_{24}, \quad \mathbf{A}_u = \begin{bmatrix} \mathbf{I}_3 \\ -\mathbf{I}_3 \end{bmatrix}, \quad \mathbf{b}_u = 0.5 \cdot \mathbf{1}_6. \quad (157)$$

The initial conditions $\bar{\mathbf{x}}_0$ are sampled from $\mathcal{U}([-1, 1]^{12})$.

Portfolio optimization. We consider the same portfolio optimization problem setup as in Stellato et al. (2020). For completeness, we briefly repeat it here,

$$\max_{\mathbf{x}} \boldsymbol{\mu}^\top \mathbf{x} - \gamma(\mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x}) \quad \text{s.t.} \quad x_1 + \dots + x_n = 1, \quad \mathbf{x} \geq \mathbf{0}, \quad (158)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the assets allocation vector, $\boldsymbol{\mu} \in \mathbb{R}^n$ is the expected returns vector, $\boldsymbol{\Sigma} \in \mathbb{R}_+^n$ is the risk covariance matrix and $\gamma > 0$ is the risk aversion parameter. The matrix $\boldsymbol{\Sigma}$ is of the form $\boldsymbol{\Sigma} = \mathbf{F}\mathbf{F}^\top + \mathbf{D}$ with $\mathbf{F} \in \mathbb{R}^{d \times n}$ is the factors matrix and $\mathbf{D} \in \mathbb{R}^{n \times n}$ is a diagonal matrix involving individual asset risks. Using an auxiliary variable $\mathbf{t} = \mathbf{F}^\top \mathbf{x}$, then problem equation 158 is rewritten as

$$\min_{\mathbf{x}, \mathbf{t}} \mathbf{x}^\top \mathbf{D} \mathbf{x} + \mathbf{t}^\top \mathbf{t} - \frac{1}{\gamma} \boldsymbol{\mu}^\top \mathbf{x} \quad \text{s.t.} \quad \mathbf{t} = \mathbf{F}^\top \mathbf{x}, \quad \mathbf{1}^\top \mathbf{x} = 1, \quad \mathbf{x} \geq \mathbf{0}. \quad (159)$$

For the problems we are generating, we use $n = 250$, $k = 25$ and $\gamma = 1.0$. Each element of the expected return vector $\boldsymbol{\mu}$ is sampled through $\mu_i \sim \mathcal{N}(0, 1)$. The matrix \mathbf{F} consists of 50% non-zero elements sampled through $F_{ij} \sim \mathcal{N}(0, 1)$. Finally, the diagonal elements of \mathbf{D} are sampled with $\mathcal{D}_{ii} \sim \mathcal{U}[0, \sqrt{k}]$.

LASSO. The least absolute shrinkage and selection operator (LASSO) is a linear regression technique with an added ℓ_1 -norm regularization term to promote sparsity in the parameters (Tibshirani, 1996). We again consider the same problem setup as in Stellato et al. (2020), where the initial optimization problem

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1 \quad (160)$$

is rewritten as

$$\min_{\mathbf{x}, \mathbf{t}} (\mathbf{Ax} - \mathbf{b})^\top (\mathbf{Ax} - \mathbf{b}) + \lambda \mathbf{1}^\top \mathbf{t} \quad \text{s.t.} \quad -\mathbf{t} \leq \mathbf{x} \leq \mathbf{t}, \quad (161)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the vector of parameters, $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the data matrix, λ is the weighting parameter, and $\mathbf{t} \in \mathbb{R}^n$ are newly introduced variables. The matrix \mathbf{A} consists of 15% non-zero elements sampled through $A_{ij} \sim \mathcal{N}(0, 1)$. The true sparse vector $\mathbf{v} \in \mathbb{R}^n$ to be learned consists of 50% non-zero elements sampled through $v_i \sim \mathcal{N}(0, 1/n)$. We then construct $\mathbf{b} = \mathbf{Av} + \boldsymbol{\xi}$ where $\boldsymbol{\xi}_i \sim \mathcal{N}(0, 1)$ represents noise in the data. Finally, we set $\lambda = (1/5)\|\mathbf{A}^\top \mathbf{b}\|_\infty$. For the problems we are generating, we set $n = 100$ and $m = 10^4$.

I.2 PROBLEM TYPES IN DISTRIBUTED EXPERIMENTS

Random Networked QPs. In this family of problems, we generate random QPs with an underlying network structure. Consider an undirected graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, where \mathcal{V} and \mathcal{E} are the nodes and edges sets, respectively. Each node i is associated with a decision variable $\mathbf{x}_i \in \mathbb{R}^{n_i}$. Then, we generate problems of the following form

$$\min_{\{\mathbf{x}_i\}_{i \in \mathcal{V}}} \sum_{i \in \mathcal{V}} \frac{1}{2} \mathbf{x}_i^\top \mathbf{Q}_i \mathbf{x}_i + \mathbf{q}_i^\top \mathbf{x}_i \quad (162a)$$

$$\text{s.t.} \quad \mathbf{A}_{ij} \begin{bmatrix} \mathbf{x}_i \\ \mathbf{x}_j \end{bmatrix} \leq \mathbf{b}_{ij}, \quad \mathbf{C}_{ij} \begin{bmatrix} \mathbf{x}_i \\ \mathbf{x}_j \end{bmatrix} = \mathbf{d}_{ij}, \quad (i, j) \in \mathcal{E}, \quad (162b)$$

For each generated problem, a cost Hessian is constructed as $\mathbf{Q}_i = \mathbf{F}_i^\top \mathbf{F}_i + \gamma \mathbf{I}$, where each element of $\mathbf{F}_i \in \mathbb{R}^{n_i \times n_i}$ is sampled through $F_i^{kl} \sim \mathcal{N}(0, 1)$ and $\gamma = 1.0$. The elements of the cost coefficients vectors \mathbf{q}_i are also sampled through $q_i^k \sim \mathcal{N}(0, 1)$. The elements of the inequality constraints matrix $\mathbf{A}_{ij} \in \mathbb{R}^{m_{ij} \times (n_i + n_j)}$ are given by $A_{ij}^{kl} \sim \mathcal{N}(0, 1)$. The vectors $\mathbf{b}_{ij} \in \mathbb{R}^{m_{ij}}$ are obtained through $\mathbf{b}_{ij} = \mathbf{A}_{ij} \boldsymbol{\theta}_{ij}$, where each element of $\boldsymbol{\theta}_{ij} \in \mathbb{R}^{n_i + n_j}$ is sampled through $\theta_{ij}^k \sim \mathcal{N}(0, 1)$. In a similar manner, the elements of the equality constraints matrices $\mathbf{C}_{ij} \in \mathbb{R}^{p_{ij} \times (n_i + n_j)}$ are generated through $C_{ij}^{kl} \sim \mathcal{N}(0, 1)$, while the vectors $\mathbf{d}_{ij} \in \mathbb{R}^{p_{ij}}$ are acquired through $\mathbf{d}_{ij} = \mathbf{C}_{ij} \boldsymbol{\xi}_{ij}$, where each element of $\boldsymbol{\xi}_{ij} \in \mathbb{R}^{n_i + n_j}$ is generated with $\xi_{ij}^k \sim \mathcal{N}(0, 1)$.

It is straightforward to observe that problems of the form (162) can be cast in the form (2) by introducing the augmented node variables $\mathbf{x}_i^{\text{aug}} = [\mathbf{x}_i, \{\mathbf{x}_j\}_{j \in \mathcal{N}_i}]^\top$. The problem data can then be augmented based on this new $\mathbf{x}_i^{\text{aug}}$ to yield the desired problem structure. Most notably, the constraints can be rewritten as $\mathbf{A}_i^{\text{aug}} \mathbf{x}_i^{\text{aug}} \leq \mathbf{b}_i^{\text{aug}}$ and $\mathbf{C}_i^{\text{aug}} \mathbf{x}_i^{\text{aug}} = \mathbf{d}_i^{\text{aug}}$, respectively.

In our experiments, the underlying graph structure is a square grid. For random QPs without equality constraints, we set $n_i = 10$, $m_{ij} = 5$, and $p_{ij} = 0$. For random QPs with equality constraints, we set $n_i = 10$, $m_{ij} = 3$, and $p_{ij} = 2$ for the $N = 16$ training experiment and $p_{ij} = 1$ for the rest of the testing experiments until $N = 1,024$.

Multi-agent optimal control. We adapt the distributed MPC problem from (Conte et al., 2012a;b), which generalizes to different systems based on the choice of dynamics matrices, as described below. The optimization problem is given as

$$\min_{\mathbf{x}, \mathbf{u}} \sum_{i \in \mathcal{V}} \sum_{t=0}^{T-1} (\mathbf{x}_i^t)^\top \mathbf{Q}_i \mathbf{x}_i^t + (\mathbf{u}_i^t)^\top \mathbf{R}_i \mathbf{u}_i^t + (\mathbf{x}_i^T)^\top \mathbf{P}_i \mathbf{x}_i^T, \quad (163a)$$

$$\text{s.t.} \quad \mathbf{x}_i^{t+1} = \mathbf{A}_{ii} \mathbf{x}_i^t + \mathbf{B}_{ii} \mathbf{u}_i^t + \sum_{j \in \mathcal{N}_i} \mathbf{A}_{ij} \mathbf{x}_j^t, \quad t = 0, \dots, T-1, \quad i \in \mathcal{V} \quad (163b)$$

$$\mathbf{G}_x^i \mathbf{x}_i^t \leq \mathbf{f}_x^i, \quad \mathbf{G}_u^i \mathbf{u}_i^t \leq \mathbf{f}_u^i, \quad t = 0, \dots, T, \quad i \in \mathcal{V} \quad (163c)$$

$$\mathbf{x}_i^0 = \bar{\mathbf{x}}_i^0, \quad i \in \mathcal{V}, \quad (163d)$$

where \mathbf{x}_i^t and \mathbf{u}_i^t are the state and control for agent i at time t . Eq. (163b) describes the dynamics and the coupling between the agents, Eq. (163c) describe local inequality constraints, and Eq. (163d) describes the initial condition for each of the agents.

For the coupled pendulums, the individual state $\mathbf{x}_i^t \in \mathbb{R}^2$ for each agent consists of the angle and angular velocity of the pendulum and the control $\mathbf{u}_i^t \in \mathbb{R}^1$ is the torque. The dynamics matrices are given as

$$\mathbf{A}_{ii} = \begin{bmatrix} 1 & dt \\ -(\frac{g}{\ell} + \frac{\text{nn}(i)k}{m})dt & 1 - \frac{\text{nn}(i)c}{m}dt \end{bmatrix}, \quad \mathbf{A}_{ij} = \begin{bmatrix} 0 & 0 \\ \frac{k}{m}dt & \frac{c}{m}dt \end{bmatrix}, \quad \mathbf{B}_i = \begin{bmatrix} 0 \\ \frac{1}{m\ell^2}dt \end{bmatrix},$$

where $dt = 0.1$ is the discretization step size, $g = 9.81$ is the gravitational constant, $m = 1.0$ is the mass of each pendulum, $\ell = 0.5$ is the length of each pendulum, $\text{nn}(i)$ is the number of neighbors of agent i , $k = 0.1$ is the spring constant between each pendulum, and $c = 0.1$ is the damping constant between each pendulum. We have used the small angle assumption $\sin \theta \approx \theta$ so the dynamics are linear and therefore the optimization is convex. There are no inequality constraints for the coupled pendulums. The initial states are sampled uniformly from $\mathcal{U}[-\pi, \pi]$. Finally, we considered $N = 10$ and $T = 30$.

For the coupled oscillating masses, we adapt the same benchmark system from Chen et al. (2022) used in the non-distributed experiments. The individual state $\mathbf{x}_i^t \in \mathbb{R}^2$ for each agent consists of the displacement and velocity of the mass and the control $\mathbf{u}_i^t \in \mathbb{R}^1$ is the force acting on the mass. The dynamics matrices are

$$\mathbf{A}_{ii} = \begin{bmatrix} 1 & dt \\ -\frac{2k}{m}dt & 1 - \frac{2c}{m}dt \end{bmatrix}, \quad \mathbf{A}_{ij} = \begin{bmatrix} 0 & 0 \\ \frac{k}{m}dt & \frac{c}{m}dt \end{bmatrix}, \quad \mathbf{B}_i = \begin{bmatrix} 0 \\ \frac{1}{m}dt \end{bmatrix},$$

where $dt = 0.5$ is the discretization step size, $m = 1.0$ is the mass, $k = 0.4$ is the spring constant between each mass, and $c = 0.1$ is the damping constant between each mass. The initial states are sampled uniformly from $\mathcal{U}[-2.0, 2.0]$. Inequality constraints $-4 \leq \mathbf{x}_i^t \leq 4$ and $-0.5 \leq \mathbf{u}_i^t \leq 0.5$ are represented as

$$\mathbf{G}_x^i = \begin{bmatrix} \mathbf{I}_2 \\ -\mathbf{I}_2 \end{bmatrix}, \quad \mathbf{f}_x^i = 4 \cdot \mathbf{1}_4, \quad \mathbf{G}_u^i = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad \mathbf{f}_u^i = 0.5 \cdot \mathbf{1}_2,$$

For both the distributed MPC problems described above, the cost matrices are taken to be identity matrices: $\mathbf{Q}_i = \mathbf{I}_2$, $\mathbf{R}_i = \mathbf{I}_1$, and $\mathbf{P}_i = \mathbf{I}_2$, for all $i \in \mathcal{V}$.

The optimization Eq. (163) can be expressed in the form of Eq. (2) by defining an augmented vector consisting of the individual agent's states and controls, as well as the states and controls of its neighbors. Letting $\mathbf{z}_i = [\mathbf{x}_i^0, \mathbf{u}_i^0, \dots, \mathbf{x}_i^T]^\top$, the augmented optimization vector for each agent i is given as $\mathbf{x}_i^{\text{aug}} = [\mathbf{z}_i, \{\mathbf{z}_j\}_{j \in \mathcal{N}_i}]^\top$. The cost, dynamics, and constraint matrices can be augmented straightforwardly based on this new $\mathbf{x}_i^{\text{aug}}$. For all problems, we considered $T = 15$.

Network flow. The network flow problem is adapted from Mota (2013); Mota et al. (2014). We consider a directed regular graph with 200 nodes and 1000 directed edges $x_{ij} \in \mathcal{E}$. Each edge has an associated quadratic cost function $\phi_{ij}(x_{ij}) = \frac{1}{2}(x_{ij} - a_{ij})^2$, where a_{ij} is sampled from $[1.0, 2.0, 3.0, 4.0, 5.0, 10.0]$ with probabilities $[0.2, 0.2, 0.2, 0.2, 0.1, 0.1]$. The objective is to optimize the flow through the graph subject to equality constraints on the flow into and out of each node. Namely, the flow into each node should be equal to the flow out of the node. For node i , the flow conservation constraint is $\sum_{j \in \mathcal{E}_i^-} x_{ji} = \sum_{k \in \mathcal{E}_i^+} x_{ik}$, where \mathcal{E}_i^- is the set of all incoming edges to node i , and similarly \mathcal{E}_i^+ is the set of all outgoing edges from node i . 100 nodes are randomly selected and injected with an external flow f_k sampled identically to a_{ij} . For each of these nodes, a reachable descendant is randomly selected and an equivalent amount of flow f_k is removed from those nodes.

This problem is straightforward to express in the form Eq. (2) by considering each node as an individual agent and defining the local state vector for each agent as

$$\mathbf{x}_i = \begin{bmatrix} \{x_{ji}\}_{j \in \mathcal{E}_i^-} \\ \{x_{ik}\}_{k \in \mathcal{E}_i^+} \end{bmatrix}, \quad (164)$$

Table 2: Training and testing details for DeepQP.

Problem Class	No of layers K	Training dataset size	Epochs	Training time	Test dataset size
Random QPs	30	2,000	125	21min	1,000
Random QPs with Eq. Constraints	30	2,000	125	23min	1,000
Double Integrator	30	500	300	28min	1,000
Osc. Masses	15	500	300	48min	1,000
Portfolio Optimization	30	500	300	1h 14min	1,000
LASSO	10	500	300	20min	1,000

Table 3: Training and testing details for DeepDistributedQP.

Problem Class	No of layers K	Training dataset size	Epochs	Training time	Test dataset size
Random QPs	50	1,000	300	3h 21min	500
Random QPs with Eq. Constraints	50	500	600	3h 29min	500
Coupled Pendulums	20	500	400	1h 49min	500
Coupled Osc. Masses	20	500	600	2h 29min	500
Network Flow	30	500	600	2h 8min	500
Distributed LASSO	20	500	600	56min	500

consisting of all the incoming and outgoing edges for node i . Each agent is responsible for its own flow constraint defined by

$$\mathbf{A}_i = \begin{bmatrix} \{1\}_{j \in \mathcal{E}_i^-} & \{-1\}_{k \in \mathcal{E}_i^+} \\ \{-1\}_{j \in \mathcal{E}_i^-} & \{1\}_{k \in \mathcal{E}_i^+} \end{bmatrix}, \quad \mathbf{b}_i = \mathbf{0}, \quad (165)$$

where \mathbf{b}_i might instead contain the external injected or removed flow f_i for that node i . The augmented cost matrix \mathbf{Q}_i is zero for all incoming edges and has entries $1/2$ on the diagonal of the outgoing edges. The augmented cost vector \mathbf{q}_i contains each of the quadratic cost offsets a_{ik} :

$$\mathbf{Q}_i = \begin{bmatrix} \{0\}_{j \in \mathcal{E}_i^-} & \\ & \{\frac{1}{2}\}_{k \in \mathcal{E}_i^+} \end{bmatrix}, \quad \mathbf{q}_i = \begin{bmatrix} \{0\}_{j \in \mathcal{E}_i^-} \\ \{-a_{ik}\}_{k \in \mathcal{E}_i^+} \end{bmatrix}. \quad (166)$$

Finally, we impose the constraint $-f_{\max} \cdot \mathbf{1} \leq \mathbf{x}_i \leq f_{\max} \cdot \mathbf{1}$ on the maximum allowed flow of all edges, with $f_{\max} = 5$.

Distributed LASSO. Distributed LASSO (Mateos et al., 2010) extends LASSO to situations where the training data are distributed across different agents and agents cannot share training data with each other. It can be formulated as

$$\min_{\{\mathbf{x}_i\}_{i=1}^N, \mathbf{w}} \sum_{i=1}^N \|\mathbf{A}_i \mathbf{x}_i - \mathbf{b}_i\|_2^2 + \frac{\lambda}{N} \|\mathbf{x}_i\|_1 \quad \text{s.t.} \quad \mathbf{x}_i = \mathbf{w}, \quad i = 1, \dots, N \quad (167)$$

where $\mathbf{w} \in \mathbb{R}^{n_i}$ is a global vector of regression coefficients, $\mathbf{x}_i \in \mathbb{R}^{n_i}$ is a local copy of \mathbf{w} , $\mathbf{A}_i \in \mathbb{R}^{m_i \times n_i}$ and $\mathbf{b}_i \in \mathbb{R}^{m_i}$ are the training data available to agent i , and λ is the weighting parameter. Similarly to non-distributed LASSO, this formulation is rewritten as

$$\min \sum_{i=1}^N (\mathbf{A}_i \mathbf{x}_i - \mathbf{b}_i)^\top (\mathbf{A}_i \mathbf{x}_i - \mathbf{b}_i) + \frac{\lambda}{N} \mathbf{1}^\top \mathbf{t}_i \quad (168a)$$

$$\text{s.t.} \quad \mathbf{t}_i \leq \mathbf{x}_i \leq \mathbf{t}_i, \quad \mathbf{x}_i = \mathbf{w}, \quad \mathbf{t}_i = \mathbf{g}, \quad i = 1, \dots, N \quad (168b)$$

where $\mathbf{t}_i \in \mathbb{R}^{n_i}$ are newly-introduced variables and \mathbf{g} is the global copy of \mathbf{t}_i .

The matrix \mathbf{A}_i consists of 15% non-zero elements sampled through $\mathbf{A}_i^{kl} \sim \mathcal{N}(0, 1)$. The true sparse vector $\mathbf{v} \in \mathbb{R}^n$ to be learned consists of 50% non-zero elements sampled through $\mathbf{v}_i \sim \mathcal{N}(0, 1/n)$. We then construct $\mathbf{b} = \mathbf{A}\mathbf{v} + \boldsymbol{\xi}$ where $\boldsymbol{\xi}_i \sim \mathcal{N}(0, 1)$ represents noise in the data.

Finally, we set $\lambda = (1/5) \max_i (\|\mathbf{A}_i^\top \mathbf{b}_i\|_\infty)$. For the problems, we have $n_i = 50$ and $m_i = 5 \cdot 10^3$.

I.3 DETAILS ON TRAINING AND TESTING

Here, we discuss details regarding the training and testing of DeepQP and DeepDistributedQP in the presented experiments.

Table 4: List of OSQP penalty parameters used in centralized experiments

Problem Class	List of penalty parameters ρ
Random QPs	0.1, 0.3, ..., 3, 10
Random QPs with Eq. Constraints	0.1, 0.3, ..., 3, 10
Double Integrator	3, 5, ..., 100, 300
Osc. Masses	0.1, 0.3, ..., 3, 10
Portfolio Optimization	3, 5, ..., 100, 300
LASSO	30, 50, ..., 1000, 3000

Table 5: List of DistributedQP penalty parameters used in distributed experiments

Problem Class	List of penalty parameters ρ
Random QPs	0.1, 0.3, ..., 3, 10
Random QPs with Eq. Constraints	0.1, 0.3, ..., 3, 10
Coupled Pendulums	0.1, 0.3, ..., 3, 10
Coupled Osc. Masses	0.1, 0.3, ..., 3, 10
Network Flow	0.1, 0.3, ..., 3, 10
Distributed LASSO	30, 50, ..., 1000, 3000

Centralized experiments. Table 2 shows the number of layers K , training dataset size, number of epochs, total training time and testing dataset size for DeepQP in every centralized problem. The increased dataset size and number of epochs for RandomQPs is motivated by the fact that the structure in these problems is less clear; learning policies that exploit this structure therefore requires more examples and takes longer. In all experiments, DeepQP was trained with a batch size of 50 using the Adam optimizer with learning rate 10^{-3} . The feedback layers are set as 2×16 MLPs. DeepQP and OSQP always start with zero initializations in all comparisons. The weights of the training loss were set to $\gamma_k = \exp((k - K)/5)$ in all experiments. **Both the training and testing datasets are constructed after letting OSQP running until optimality.**

Distributed experiments. Table 3 shows the number of layers K , training dataset size, number of epochs, total training time and testing dataset size for DeepDistributedQP in every distributed problem. In all experiments, DeepDistributedQP was trained with a batch size of 50 using the Adam optimizer with learning rate 10^{-3} . The feedback layers are set as 2×16 MLPs. DeepDistributedQP and DistributedQP always start with zero initializations in all comparisons. In all experiments, the weights of the training loss were set to $\gamma_k = \exp((k - K)/5)$. **For the low-dimensional testing datasets, these datasets are constructed using OSQP. For larger scales, the testing dataset is constructed with DistributedQP instead as it is much faster (see Table 6), after ensuring convergence to optimality.**

Generalization bounds experiments. These experiments were performed on a networked random QPs problem with $N = 16$, $n_i = 10$, $m_{ij} = 5$, $p_{ij} = 0$ and on a coupled pendulums problem with $N = 10$ and the same parameters as described in the previous section. The prior was obtained through training on a small separate dataset of 500 problems for 50 epochs. The posterior was then acquired through optimizing for the generalization bound with a dataset of 15,000 problems for 100 epochs.

I.4 DETAILS ON STANDARD OPTIMIZERS

Details on OSQP. When comparing with OSQP using fixed penalty parameters, we selected the best-performing subsequence of $\{\dots, 0.1, 0.3, 0.5, 1.0, 3.0, 5.0, \dots\}$ as the penalty parameters to plot against. Table 4 shows these parameters for every centralized problem in our experiments. For equality constraints, we scaled ρ by 10^3 , as in Stellato et al. (2020). For the adaptive version, we preferred the standard heuristic adaptation rule shown in Boyd et al. (2011) with $\tau = 2.0$ and $\mu = 10.0$, instead of the OSQP adaptation scheme (Stellato et al., 2020), as it performed better in our problem instances. We hypothesize that this might be due to the fact that as scale increases the infinity norm is ignoring more information than the 2-norm. The initial ρ^0 was initialized as the median of the range of fixed penalty parameters.

Details on DistributedQP. The range of fixed penalty parameters to compare with was chosen using the same methodology as with OSQP. Table 5 shows these parameters for every distributed problem in our experiments. For the adaptive version, we used the standard heuristic adaptation rule

Table 6: **Wall-clock times and iterations for DeepDistributedQP, DistributedQP, OSQP (indirect) and OSQP (direct).** This comparison shows the total wall-clock times for DistributedQP and OSQP (indirect or direct method) required to reach the same accuracy as DeepDistributedQP. For OSQP with direct method, we only report the time for the first iteration, assuming the best-case scenario in which the factorized KKT matrix can be reused for all subsequent iterations. Both DeepDistributedQP and DistributedQP demonstrate orders-of-magnitude improvements compared to OSQP as scale increases. In addition, DeepDistributedQP maintains a significant advantage over its standard optimization counterpart in all cases.

				DeepDistrQP (ours)	DistrQP (ours)	OSQP (Indirect)		OSQP (Direct)			
Networked Random QPs											
N	n	m	$\text{nnz}(Q, A)$	Time	Iters	Time	Iters	Time (1st iter.)	Iters		
16	160	120	4,000	33.05 ms	50	141.9 ms	208	46.16 ms	29	0.86 ms	29
64	640	560	17,600	39.11 ms	50	129.2 ms	192	185.1 ms	28	23.8 ms	28
256	2,560	2,400	73,600	50.21 ms	50	128.8 ms	168	514 ms	23	703.5 ms	23
1,024	10,240	9,920	300,800	62.68 ms	50	158.9 ms	165	3.03s	23	8.20 s	23
Networked Random QPs with Equality Constraints											
N	n	m	$\text{nnz}(Q, A)$	Time	Iters	Time	Iters	Time	Iters	Time (1st iter.)	Iters
16	160	168	4,960	37.21 ms	50	138.9 ms	170	36.52 ms	19	0.76 ms	19
64	640	560	17,600	57.76 ms	50	238.1 ms	172	109.0 ms	17	26.9 ms	17
256	2,560	2,400	73,600	74.54 ms	50	239.5 ms	164	692.5 ms	17	956.0 ms	17
1,024	10,240	9,920	300,800	82.55 ms	50	371.0 ms	172	5.83 s	16	11.60 s	16
Coupled Pendulums Optimal Control											
N	n	m	$\text{nnz}(Q, A)$	Time	Iters	Time	Iters	Time	Iters	Time (1st iter.)	Iters
10	470	640	3,690	50.99 ms	20	89.81 ms	35	49.46 ms	8	4.95 ms	8
20	940	1,200	7,500	66.44 ms	20	116.7 ms	35	372.0 ms	8	199.7 ms	8
50	2,350	3,200	18,930	75.9 ms	20	142.1 ms	34	948.8 ms	8	4.38 s	8
100	4,700	6,400	37,980	101.9 ms	20	201.9 ms	35	3.97 s	9	19.91 s	9
200	9,400	12,800	76,080	146.0 ms	20	284.8 ms	34	22.41 s	8	90.07 s	8
500	23,500	32,000	190,380	204.3 ms	20	379.8 ms	36	112.9 s	9	Out of memory	
1,000	47,000	64,000	380,880	317.2 ms	20	628.2 ms	34	Out of memory		Out of memory	
Coupled Oscillating Masses Optimal Control											
N	n	m	$\text{nnz}(Q, A)$	Time	Iters	Time	Iters	Time	Iters	Time (1st iter.)	Iters
10	470	1,580	4,590	48.22 ms	20	73.58 ms	33	79.1 ms	9	178.4 ms	9
20	940	3,160	9,300	67.93 ms	20	91.53 ms	33	641.9 ms	9	2.37 s	9
50	2,350	7,900	23,430	73.92 ms	20	97.34 ms	32	1.07 s	8	28.1 s	8
100	4,700	15,800	46,980	91.93 ms	20	148.8 ms	33	5.45 s	8	132 s	8
200	9,400	31,600	94,080	109.4 ms	20	194.4 ms	34	31.8 s	8	614 s	8
300	28,200	47,400	141,180	132.8 ms	20	304.8 ms	33	243 s	8	Out of memory	
Network Flow											
N	n	m	$\text{nnz}(Q, A)$	Time	Iters	Time	Iters	Time	Iters	Time (1st iter.)	Iters
20	100	140	600	6.80 ms	30	10.68 ms	50	9.51 ms	15	0.59 ms	15
50	250	350	1,500	7.81 ms	30	13.17 ms	48	14.81 ms	16	1.30 ms	16
200	1,000	1,400	6,000	12.08 ms	30	17.61 ms	42	208.19 ms	17	61.93 ms	17
500	2,500	3,500	15,000	13.63 ms	30	19.73 ms	40	425.7 ms	17	745.2 ms	17
1,000	5,000	7,000	30,000	20.51 ms	30	31.62 ms	40	8.73 s	18	11.59 s	18
2,000	10,000	14,000	60,000	29.86 ms	30	47.22 ms	40	51.6 s	18	73.9 s	18
5,000	25,000	35,000	150,000	61.23 ms	30	85.99 ms	39	558 s	18	Out of memory	
Distributed LASSO											
N	n	m	$\text{nnz}(Q, A)$	Time	Iters	Time	Iters	Time	Iters	Time (1st iter.)	Iters
10	1,100	3,000	29,000	15.06 ms	20	28.57 ms	37	2.04 s	33	148.2 ms	33
50	5,500	15,000	145,000	24.92 ms	20	44.27 ms	38	13.74 s	31	49.21 s	31
100	10,100	30,000	290,000	30.51 ms	20	51.44 ms	35	85.92 s	32	342.9 s	32
200	20,100	60,000	580,000	40.88 ms	20	76.21 ms	36	418.9 s	32	Out of memory	
500	50,100	150,000	1,450,000	69.19 ms	20	130.24 ms	35	Out of memory		Out of memory	

shown in Boyd et al. (2011) with $\tau = 2.0$ and $\mu = 10.0$. The initial value was again always chosen as the median value of the above lists.

I.5 DETAILS ON WALL-CLOCK TIMES

In Table 6, we list the observed wall-clock times for DeepDistributedQP (ours), DistributedQP (ours) and OSQP using either the indirect or the direct method. The table presents all six studied problems with an increasing dimension. As clearly observed, DeepDistributedQP and DistributedQP demonstrate a substantially more favorable scalability than OSQP. In fact, the two algorithms can efficiently solve problems that OSQP cannot even handle due to memory overflow on our system.

1944 Finally, DeepDistributedQP also maintains a clear advantage over its standard optimization counter-
1945 part DistributedQP across all experiments which signifies the importance of learning policies for the
1946 algorithm parameters.
1947

1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

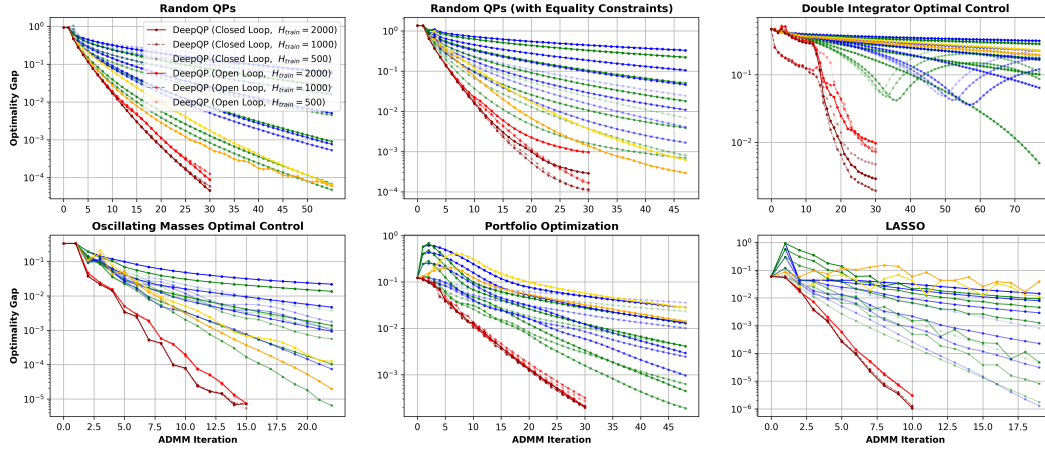


Figure 8: **Varying training dataset size for DeepQP.** The performance of DeepQP remains robust (for both open-loop and closed-loop policies) even as the training dataset size is reduced.

J ADDITIONAL EXPERIMENTS

The following experiments are dedicated into providing additional insight on exploring the performance of DeepDistributedQP and DeepQP in various testing scenarios.

J.1 VARYING TRAINING DATASET SIZE

This section provides additional insight on the amount of training data required for the proposed learned optimizers to perform well.

In Fig. 8, we compare the performance of DeepQP on the centralized problems using a training dataset size of 500, 1000 or 2000. To ensure an “equivalent total training effort”, we train these three cases for $4e$, $2e$ and e epochs, respectively, where $e = 125$ for random QPs and $e = 75$ for the other problems. This comparison highlights the robust performance of DeepQP even with a limited amount of training data. Interestingly, we also observe that training with less data but over more epochs had a beneficial effect on two out of six problems. We hypothesize that this could be attributed to the non-convex nature of training in deep learning, as well as the possibility that additional epochs might have allowed for further improvements in cases where the training of the model had not yet fully converged. Overall, we conclude that DeepQP maintains reliable performance even when training data is limited.

For the training of DeepDistributedQP, a limited training dataset of 500 sample problems was used for all problems except for the random QPs without equality constraints. For completeness, Fig. 9 presents a performance comparison of the learned optimizer when trained with 500 sample problems (600 epochs) and 1000 sample problems (300 epochs). While additional training data provides some improvement, the model trained with less sample problems still significantly outperforms the standard optimization counterparts.

J.2 CAN POLICIES TRAINED FOR SPECIFIC PROBLEMS BE APPLIED TO OTHER PROBLEMS?

The field of learning-to-optimize primarily focuses on improving the performance of an underlying optimizer for problems drawn from the same distribution as the training data (Shlezinger et al., 2022). However, this prompts an interesting question: How does a policy trained on a specific problem class perform when evaluated on a different class?

At this point, we wish to emphasize the following fact:

The proposed DeepDistributedQP framework already surpasses the expected capabilities of typical learning-to-optimize algorithms, as it is trained on small-scale problems and then successfully deployed on much higher-dimensional ones.

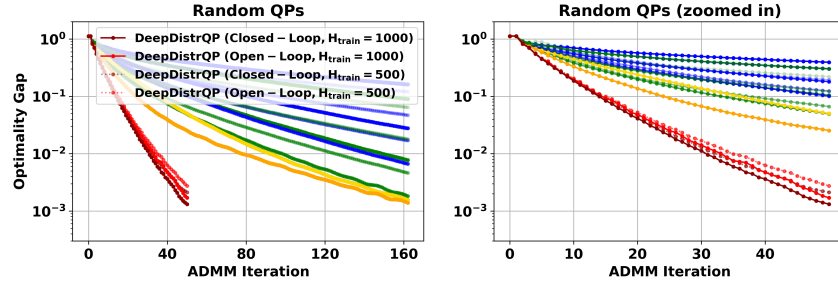


Figure 9: **Performance of DeepDistributedQP on random QPs using training dataset size 500 and 1000.** While using more training data results in a slight improvement in the performance of DeepDistributedQP, in both scenarios, the proposed learned optimizer consistently outperforms the traditional one. The right figure illustrates only the first 50 iterations.

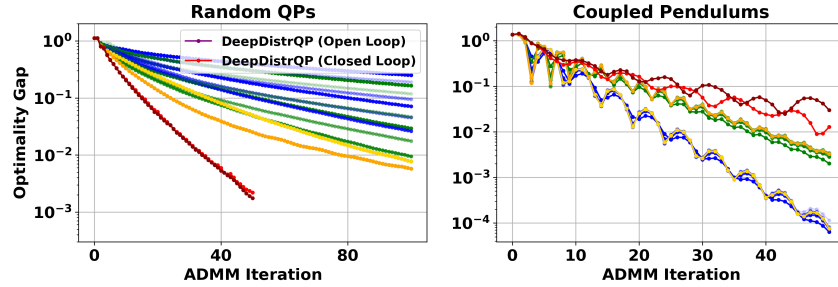


Figure 10: **Testing policies on different classes of problems.** We evaluate the policy trained on small-scale random equality-constrained QP problems ($N = 16$) on two large-scale scenarios of different problem types: random QPs without equality constraints ($N = 1,024$) and coupled pendulums ($N = 1,000$). Notably, in the first case (left), the policy demonstrates strong performance which is attributed on the fact that there is still some similarity between the training and testing setups. In the second case (right), where the testing problems differ entirely from the training setup, the performance of the learned optimizer is suboptimal but remains acceptable.

For completeness, we also conducted curiosity-driven experiments by applying the trained policies to different classes of problems than the ones used for training. In Fig. 10, we test a policy trained on small-scale random equality-constrained QPs on large-scale random QPs without equality constraints and large-scale coupled pendulums problems. In the first case, DeepDistributedQP maintains remarkable performance compared to DistributedQP due to the existing similarity between the two classes. In the second setup, where the training and testing problems are entirely different, the performance is suboptimal. Overall, these results highlight that when there is a degree of similarity between the training and testing setups, DeepDistributedQP is expected to still perform very well. In future work, we plan to explore extensions trained on a broader variety of problem classes to improve generalization on entirely different setups.

J.3 VARYING THE NUMBER OF LAYERS IN TESTING DEEPCONVEXQP

Another natural question that arises is how DeepDistributedQP can be adapted to run for more iterations than the number of layers it was originally trained for. A straightforward modification is to repeat the last layer of the framework for the extra needed iterations. In Fig. 11, we add 30 extra iterations for the random QPs and 20 for the other problems. For all cases, the closed-loop policies continue to outperform the standard optimizers. Additionally, the open-loop policies maintain strong performance in 4 out of 6 problems. In future work, we plan to incorporate the repetition of the last layer during training to further ensure robust performance when additional iterations are required.

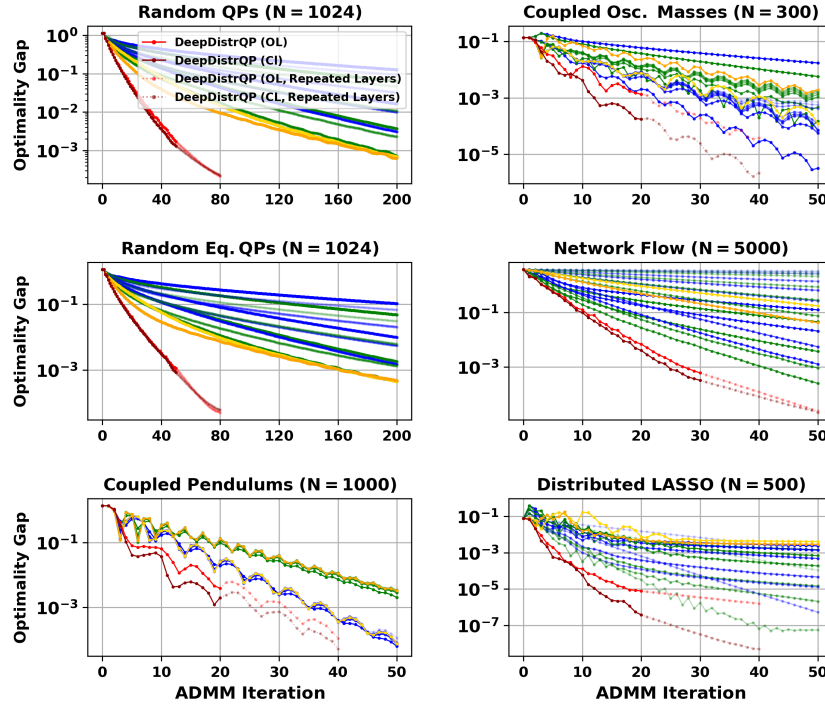


Figure 11: **Varying the number of layers while testing DeepDistributedQP.** If additional iterations are needed, DeepDistributedQP maintains strong performance by repeating its last layer for these extra iterations. Specifically, we explore adding 30 iterations for the random QPs and 20 for the rest of the problems. In all cases, the closed-loop policies continue to demonstrate superior performance, while in 4 out of 6 problems, the open-loop policies also remain advantageous.