
Efficient Algorithms for Distributed Saddle Problems

Ruichen Luo
IST Austria
Klosterneuburg, Austria
rluo@ist.ac.at

Anton Rodomanov
CISPA Helmholtz Center
Saarbrücken, Germany
anton.rodomanov@cispa.de

Sebastian U. Stich
CISPA Helmholtz Center
Saarbrücken, Germany
stich@cispa.de

Abstract

The distributed setting for Saddle Problems (SPs) has recently emerged as a framework for modern applications in machine learning and multiagent systems. Despite its relevance, the theoretical foundations of this setting have not yet been thoroughly established. In this paper, we advance this research direction by formalizing the distributed setup for SPs and providing rigorous definitions of communication and oracle costs. Further, we prove lower bounds for any distributed gradient-span algorithm, which reveals the gap from existing methods and this theoretical limit. To this end, we provide a Decoupled Method built upon a novel multi-stage reduction that reduces the SP into a sequence of decoupled minimization tasks of residual norms. Our algorithm matches the communication lower bound, thus setting the communication complexity within the gradient-span algorithms. Moreover, it yields the first strict improvement over the long-standing oracle cost of the Extragradient method for general SPs. Finally, we study the extension of distributed SP into Variational Inequality Problem (VIP), which generalizes two-player zero-sum games to multiplayer general-sum games. We show that our Decoupled Method achieves a new state-of-the-art communication complexity for this broader class.

1 Introduction

Motivation. Saddle problems (SPs) and their generalizations, variational inequality problems (VIPs), are of fundamental importance to optimization and game theory. These problems have a wide array of modern applications, including the training of GANs [Goodfellow et al., 2014], robust optimization [Ben-Tal and Nemirovski, 2002], and equilibrium computation in game theory and multiagent systems [von Neumann and Morgenstern, 1947, Rosen, 1965, Hu et al., 1998].

The growing scale of modern problems—driven by applications in machine learning, complex game dynamics, and multiagent protocols—renders reliance on a single central processor increasingly impractical. Beyond scalability, a more crucial factor is that many applications are inherently *distributed*: agents are often geographically dispersed, driven by their own individual interests, and bound by privacy constraints that prohibit the sharing of raw data or utilities. Consequently, distributed computation has become an essential regime for these problems. This perspective underlies a growing body of work in large-scale learning, game-theoretic models, and multiagent systems [McMahan et al., 2017, Zhang et al., 2024, Conitzer and Sandholm, 2004, Nisan and Segal, 2006, Hart and Mansour, 2010, Yoon et al., 2025].

In this work, we consider a natural setup where the decision variables and oracles of the SPs or VIPs are partitioned among distributed agents. For instance, in a classic saddle problem $\min_{\mathbf{x}} \max_{\mathbf{y}} F(\mathbf{x}, \mathbf{y})$, we consider one agent controls the minimizing variable \mathbf{x} while another controls the maximizing variable \mathbf{y} . This partition naturally models, for instance, the strategic autonomy of players in game theory, the interaction protocol in multiagent systems, and the physical separation of the generator and discriminator in GANs [Conitzer and Sandholm, 2004, Goodfellow et al., 2014]. Since the decision variables are coupled within their utilities, these agents must coordinate to reach a mutual equilibrium. To do so, they form a communication network that allows them to

exchange certain information, such as their current decision variables. Thus, this provides a natural distributed setup where the decision variables are separated among the different agents.

While distributed optimization is well-established for finite-sum minimization and federated learning [Schmidt et al., 2017, McMahan et al., 2017], the literature on SPs has primarily focused on extending these data-distributed paradigms [Deng and Mahdavi, 2021, Beznosikov et al., 2025]. In contrast, the study of distributed variables and oracles, which are essential to the multi-agent systems, remains a relatively new topic.

Although a few recent works have touched upon this direction [Zhang et al., 2024, Zindari et al., 2025, Yoon et al., 2025, Yoon and Loizou, 2025], they predominantly focus on algorithms tailored to specific, favorable scenarios. Consequently, a fundamental gap persists: the lack of a theoretical framework for the general distributed settings of SPs and VIPs. Existing discussions regarding performance often remain at a vague conceptual level, lacking a rigorous formalization of the distributed environment itself. Specifically, there are no standardized definitions for communication and oracle costs in this context. Without such a foundation, it can be difficult to determine the performance limits (lower bounds) or to formally compare the efficiency of different protocols.

To enable a rigorous analysis, it is essential to establish metrics that reflect the constraints of distributed multiagent systems, where network latency and bandwidth often dwarf local processing time. In this regime, the primary bottleneck is the communication cost (exchange rounds), while the computational cost (local gradient queries) is a secondary objective. Viewed through this lens, the Extragradient (EG) method [Tseng, 1995, Nemirovski, 2004] serves as the “gold standard” baseline, though the challenges associated with it differ by metric. Regarding computational complexity, consistently improving upon EG for general monotone problems has remained an elusive goal despite over two decades of research. Regarding communication complexity—a metric that has recently come into focus with the rise of distributed systems—EG similarly defines the current state-of-the-art. Surpassing this baseline in the general setting represents a new but critical open problem.

This leads to the following research questions:

- **Formalization and Limits:** How can we rigorously formalize the communication and oracle costs for distributed SPs?
- **Communication Efficiency:** Can we design an algorithm that surpasses the state-of-the-art communication bounds for distributed SPs and VIPs?
- **Oracle Efficiency:** Is it possible to consistently improve upon the long-standing oracle complexity of the EG method for general SPs?

Contributions. We answer the aforementioned questions in the affirmative, which advances the current theory of distributed SPs and VIPs.

- In [Section 2](#), we formalize the distributed saddle-point problem, gradient-span algorithms, and their communication and oracle costs. We review EG and other methods, casting them as gradient-span algorithms to analyze their costs.
- In [Section 3](#), by drawing connection to classic convex minimization, we establish the lower bounds for both communication and oracle costs for distributed gradient-span algorithms.
- In [Section 4](#), we start with a template DM-SP algorithm with a simple, one-loop communication protocol, which improves the state-of-the-art communication cost and matches the communication lower bound we constructed earlier in [Section 3](#).
- Then, continuing in [Section 4](#) and [C](#), after making the novel multi-stage reduction, we equip the template method with a concrete implementation, thereby consistently improving the (long-standing) oracle cost of EG for general SPs.
- Finally, in [Section 5](#), we extend the results to multi-agent settings. We propose DM-VIP and improve the state-of-the-art communication cost for the class of distributed VIPs.

Notations. Let $[n] \triangleq \{1, \dots, n\}$, for any positive integer n . For any finite-dimensional real vector space \mathcal{E} , we denote its Euclidean norm by $\|\cdot\|_{\mathcal{E}}$ and its dual norm by $\|\cdot\|_{\mathcal{E}^*}$. Specifically, we equip the space $\mathcal{E}_x = \mathbb{R}^{n_x}$ with the norm $\|\mathbf{x}\|_x = \langle \mathbf{P}_x \mathbf{x}, \mathbf{x} \rangle^{1/2}$, where $\mathbf{P}_x: \mathcal{E}_x \rightarrow \mathcal{E}_x^*$ is a self-adjoint positive definite operator and the dual pairing $\langle \phi_x, \mathbf{x} \rangle$ denotes $\phi_x(\mathbf{x})$. We denote its corresponding dual norm by $\|\cdot\|_{x^*}$. We assume analogous geometries for $\mathcal{E}_y = \mathbb{R}^{n_y}$, $\mathcal{E}_i = \mathbb{R}^{n_i}$ ($i \in [K]$), and $\mathcal{E}_w = \mathbb{R}^{n_w}$, associated with their respective operators \mathbf{P}_y , \mathbf{P}_i , and \mathbf{P}_w . For a function $\psi: \mathcal{E} \rightarrow \mathbb{R} \cup \{+\infty\}$, let $\text{dom } \psi$ denote its effective domain and $\partial\psi(\mathbf{z})$ its subdifferential at $\mathbf{z} \in \text{dom } \psi$. Finally, for any set of vectors S , let $\text{span } S$ denote its linear span.

Table 1: Summary of algorithms and complexity results for distributed SPs.

Method \mathcal{M}	Communication Cost	Better Oracle? ^a	Multi-Agent ^b
EG ^c	$\mathcal{O}\left(\theta \frac{L_{xy} D_x D_y}{\epsilon} + \frac{L_x D_x^2 + L_y D_y^2}{\epsilon}\right)$	No	Yes
DGDA ^c	$\mathcal{O}\left(\log \frac{1}{\epsilon}\right)$ (weakly coupled only)	No	Yes
Cat-EG ^c	$\mathcal{O}\left(\left(\frac{L_{\max} \hat{D}_x \hat{D}_y}{\epsilon} + \sqrt{\frac{L_x \hat{D}_x^2 + L_y \hat{D}_y^2}{\epsilon}}\right) \log^2\left(\frac{1}{\epsilon}\right)\right)$	Maybe	No
Cat-Cat-DAGDA ^c	$\mathcal{O}\left(\frac{L_{xy} \hat{D}_x \hat{D}_y}{\epsilon} \log^3\left(\frac{1}{\epsilon}\right)\right)$	Maybe	No
Lower Bound (Thm. 1)	$\Omega\left(\frac{L_{xy} D_x D_y}{\epsilon}\right)$	–	–
DM-SP (Thm. 2)	$\mathcal{O}\left(\theta \frac{L_{xy} D_x D_y}{\epsilon}\right)$	Yes	Yes

^a Indicates whether the method’s theoretical oracle cost outperforms the Extragradient (EG) baseline.

^b Indicates whether the method supports multi-agent extensions. ^c These methods are for the non-composite subclass $\mathcal{P}_{\text{SP}}^{\circ}$.

2 Saddle problems with distributed oracles

In the context of saddle problems, we consider two separate finite-dimensional real vector spaces, $\mathcal{E}_x = (\mathbb{R}^{n_x}, \|\cdot\|_x)$ and $\mathcal{E}_y = (\mathbb{R}^{n_y}, \|\cdot\|_y)$. We are interested in solving composite Saddle Problems (SPs) of the following form:

$$\min_{\mathbf{x} \in \text{dom } \psi_x} \max_{\mathbf{y} \in \text{dom } \psi_y} [F(\mathbf{x}, \mathbf{y}) \triangleq f(\mathbf{x}, \mathbf{y}) + \psi_x(\mathbf{x}) - \psi_y(\mathbf{y})], \quad (1)$$

where $\psi_x: \mathcal{E}_x \rightarrow \mathbb{R} \cup \{+\infty\}$ and $\psi_y: \mathcal{E}_y \rightarrow \mathbb{R} \cup \{+\infty\}$ represent relatively simple local components (such as regularizers or indicator functions for constrained sets), and $f(\cdot, \cdot)$ is a real-valued coupling function defined on an open set containing the domain $Q \triangleq \text{dom } \psi_x \times \text{dom } \psi_y$. To simplify the notation, we denote the joint variable by $\mathbf{z} \triangleq (\mathbf{x}, \mathbf{y}) \in Q$.

2.1 Problem class

Function family. We consider the function instances satisfying the following assumptions:

- (A1) For any fixed $\mathbf{y} \in \text{dom } \psi_y$, the function $f(\cdot, \mathbf{y})$ is convex; and for any fixed $\mathbf{x} \in \text{dom } \psi_x$, the function $f(\mathbf{x}, \cdot)$ is concave. The functions ψ_x and ψ_y are proper, closed, and convex.
- (A2) Let $D_x, D_y > 0$ be distance parameters, and let $\mathbf{z}^0 = (\mathbf{x}^0, \mathbf{y}^0) \in Q$ be given initial points. Relative to this initialization, Problem (1) has at least one saddle point $(\mathbf{x}^*, \mathbf{y}^*) \in Q$ that lies within the bounded initial sets $\mathbf{x}^* \in \mathcal{B}_x$ and $\mathbf{y}^* \in \mathcal{B}_y$, where

$$\mathcal{B}_x \triangleq \{\mathbf{x} \in \mathcal{E}_x \mid \|\mathbf{x}^0 - \mathbf{x}\|_x \leq D_x\} \quad \text{and} \quad \mathcal{B}_y \triangleq \{\mathbf{y} \in \mathcal{E}_y \mid \|\mathbf{y}^0 - \mathbf{y}\|_y \leq D_y\}.$$

- (A3) The function $f(\cdot, \cdot)$ is continuously differentiable over Q . Moreover, its gradients are Lipschitz continuous. That is, with Lipschitz parameters $L_x, L_{xy}, L_y > 0$, for all $\mathbf{x}, \mathbf{x}' \in \text{dom } \psi_x$ and $\mathbf{y}, \mathbf{y}' \in \text{dom } \psi_y$, we have:

$$\begin{aligned} \|\nabla_x f(\mathbf{x}', \mathbf{y}') - \nabla_x f(\mathbf{x}, \mathbf{y})\|_{x^*} &\leq L_x \|\mathbf{x}' - \mathbf{x}\|_x + L_{xy} \|\mathbf{y}' - \mathbf{y}\|_y, \\ \|\nabla_y f(\mathbf{x}', \mathbf{y}') - \nabla_y f(\mathbf{x}, \mathbf{y})\|_{y^*} &\leq L_{xy} \|\mathbf{x}' - \mathbf{x}\|_x + L_y \|\mathbf{y}' - \mathbf{y}\|_y. \end{aligned}$$

With these assumptions in place, we can formally define the function family of interest. Let \mathcal{F} denote the *function family* consisting of all instances $\mathcal{F} = (f, \psi_x, \psi_y, \mathbf{z}^0)$ that satisfy (A1) to (A3) for a fixed set of parameters $(L_x, L_{xy}, L_y, D_x, D_y)$. To facilitate our discussion, we refer to the terms $L_x D_x^2$ and $L_y D_y^2$ as the *diagonal conditioning*, and the term $L_{xy} D_x D_y$ as the *cross-coupled conditioning*.

Distributed oracle model. We consider a distributed architecture where the variables \mathbf{x} and \mathbf{y} are maintained and updated by two distinct entities: Agent x and Agent y , respectively. First-order information is acquired through (*deterministic*) *partial gradient oracles*, denoted abstractly by \mathcal{O}_x and \mathcal{O}_y . Specifically, for a given function instance \mathcal{F} from the function family \mathcal{F} with coupling function f , and for any input point $\mathbf{z} \in Q$:

- Agent x queries the oracle \mathcal{O}_x , which returns $\mathcal{O}_x^f(\mathbf{z}) = \nabla_x f(\mathbf{z})$.
- Agent y queries the oracle \mathcal{O}_y , which returns $\mathcal{O}_y^f(\mathbf{z}) = -\nabla_y f(\mathbf{z})$.

These oracles are strictly decoupled. They can be queried independently by their respective agents, differing in both query frequency and the input points evaluated. Most importantly, the agents are fully distributed, meaning an agent cannot query the counterpart’s oracle.

To provide a concrete example of this abstract setting, consider the objective $f(\mathbf{z}) = g(\mathbf{z}) + f_x(\mathbf{x}) - f_y(\mathbf{y})$, where $g(\mathbf{z})$ is a coupled global utility, while $f_x(\mathbf{x})$ and $f_y(\mathbf{y})$ are private utilities accessible only to Agents x and y , respectively. In this scenario, the partial gradient oracles take the form:

$$\mathcal{O}_x^f(\mathbf{z}) = \nabla_x g(\mathbf{z}) + \nabla f_x(\mathbf{x}) \quad \mathcal{O}_y^f(\mathbf{z}) = -\nabla_y g(\mathbf{z}) + \nabla f_y(\mathbf{y}),$$

for all $\mathbf{z} \in Q$. Due to the distributed setting, Agent x is entirely blind to the private utility f_y and can only execute \mathcal{O}_x , and vice versa.

Accuracy measure. To evaluate the quality of a candidate solution $(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in Q$, we rely on the restricted duality gap. Over the bounded domain $\mathcal{B} \triangleq \mathcal{B}_x \times \mathcal{B}_y$, the duality gap is defined as:

$$\Delta(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \triangleq \max_{\mathbf{z} \in \mathcal{B} \cap Q} [F(\bar{\mathbf{x}}, \mathbf{y}) - F(\mathbf{x}, \bar{\mathbf{y}})].$$

We say that a pair $(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in Q$ is an ϵ -saddle point of Problem (1) if $\Delta(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \leq \epsilon$. Our goal is to design an algorithm capable of finding such an ϵ -saddle point for any given $\epsilon > 0$.

We remark that for the classic problem of constrained optimization with bounded domains, one can enclose the constrained sets in the balls \mathcal{B}_x and \mathcal{B}_y with sufficiently large radius (e.g., the diameter of the constrained sets), then the restricted saddle problem in form (1) coincides with the original one.

Problem class. Finally, we formally define the overall *problem class*, denoted by $\mathcal{P}_{\text{SP}}(\mathcal{F}, \mathcal{O}_x, \mathcal{O}_y, \epsilon)$, or for short \mathcal{P}_{SP} . A specific problem instance $P \in \mathcal{P}_{\text{SP}}$ is constructed by drawing a function instance \mathcal{F} from the function family \mathcal{F} , equipping it with the specific distributed oracles $(\mathcal{O}_x^f, \mathcal{O}_y^f)$, and specifying a target accuracy $\epsilon > 0$. Solving the instance P requires an algorithm to output an ϵ -saddle point of \mathcal{F} utilizing the distributed oracles.

2.2 Distributed gradient-span algorithms

Standard gradient-based methods are commonly analyzed under the assumption that the algorithm generates iterates within the linear span of all historically evaluated gradients [Nesterov, 2004]. However, classic span conditions implicitly assume centralized, real-time access to the full problem state. This falls short for the problem class \mathcal{P}_{SP} . Because Agent x and Agent y are strictly decoupled and communication between them is typically the primary system bottleneck, they cannot continuously synchronize their variables.

Instead, practical algorithms must proceed in discrete communication rounds. Within a given round, an agent performs multiple local computational steps in isolation. For instance, Agent x generates new iterates for \mathbf{x} by taking linear combinations of its locally accumulated partial gradients $(\nabla_x f)$ and local subdifferentials $(\partial \psi_x)$. Crucially, to evaluate the partial gradient of the coupling function, Agent x also needs to know the value of \mathbf{y} . Cut off from real-time updates, it must rely on a local, delayed approximation (denoted $\hat{\mathbf{y}}$) constructed entirely from the historical messages Agent y sent before the current round started. Agent y operates symmetrically.

To rigorously analyze algorithms under these strict communication constraints, we formalize the class of *distributed gradient-span algorithms*. This mathematical framework explicitly separates communication rounds from local oracle queries, enforcing that each agent constructs its updates using only locally accessible history and delayed remote messages.

Algorithm trajectories and histories. Suppose an algorithm \mathcal{M} proceeds in T rounds. In each round $t \in \{0, \dots, T-1\}$, Agent x and Agent y generate local trajectories of lengths τ_x^t and τ_y^t , denoted respectively by:

$$\hat{Z}_x^t = \{\mathbf{z}_x^{t,l} = (\mathbf{x}^{t,l}, \hat{\mathbf{y}}^{t,l})\}_{l=0}^{\tau_x^t-1} \quad \text{and} \quad \hat{Z}_y^t = \{\mathbf{z}_y^{t,l} = (\hat{\mathbf{x}}^{t,l}, \mathbf{y}^{t,l})\}_{l=0}^{\tau_y^t-1}.$$

To track the information available to the agents, let $Z_x^{t-1} = \bigcup_{i=0}^{t-1} \hat{Z}_x^i$ and $Z_y^{t-1} = \bigcup_{i=0}^{t-1} \hat{Z}_y^i$ denote the accumulated histories from all prior rounds (where $Z_x^{-1} = Z_y^{-1} = \emptyset$). Within round t , the histories up to local step l are denoted by $Z_x^{t,l} = Z_x^{t-1} \cup \{\mathbf{z}_x^{t,i}\}_{i=0}^{l-1}$ and $Z_y^{t,l} = Z_y^{t-1} \cup \{\mathbf{z}_y^{t,i}\}_{i=0}^{l-1}$. Finally, the accumulated histories up to and including round t are given by Z_x^t and Z_y^t , and the complete trajectories across all T rounds are denoted by $\bar{Z}_x = Z_x^{T-1}$ and $\bar{Z}_y = Z_y^{T-1}$.

Definition 1. An algorithm \mathcal{M} is called a *distributed gradient-span algorithm* for problem class \mathcal{P}_{SP} if, when applied to any instance $P \in \mathcal{P}_{\text{SP}}$, its trajectories satisfy the following conditions:

1. For all rounds $t \in \{0, \dots, T-1\}$ and local steps $l \in \{0, \dots, \tau_x^t - 1\}$, the points in the x -trajectory satisfy

$$\begin{aligned} \mathbf{x}^{t,l} &\in \mathbf{x}^0 + \mathbf{P}_x^{-1} \text{span}\{\nabla_x f(\mathbf{z}) \mid \mathbf{z} \in Z_x^{t,l}\} + \mathbf{P}_x^{-1} \text{span}\{\partial\psi_x(\mathbf{x}) \mid \mathbf{z} \in Z_x^{t,l+1}\}, \\ \hat{\mathbf{y}}^{t,l} &\in \mathbf{y}^0 + \mathbf{P}_y^{-1} \text{span}\{\nabla_y f(\mathbf{z}) \mid \mathbf{z} \in Z_y^{t-1}\} + \mathbf{P}_y^{-1} \text{span}\{\partial\psi_y(\mathbf{y}) \mid \mathbf{z} \in Z_y^{t-1}\}, \end{aligned}$$

and symmetrically, for all local steps $l \in \{0, \dots, \tau_y^t - 1\}$, the points in the y -trajectory satisfy

$$\begin{aligned} \mathbf{y}^{t,l} &\in \mathbf{y}^0 + \mathbf{P}_y^{-1} \text{span}\{\nabla_y f(\mathbf{z}) \mid \mathbf{z} \in Z_y^{t,l}\} + \mathbf{P}_y^{-1} \text{span}\{\partial\psi_y(\mathbf{y}) \mid \mathbf{z} \in Z_y^{t,l+1}\}, \\ \hat{\mathbf{x}}^{t,l} &\in \mathbf{x}^0 + \mathbf{P}_x^{-1} \text{span}\{\nabla_x f(\mathbf{z}) \mid \mathbf{z} \in Z_x^{t-1}\} + \mathbf{P}_x^{-1} \text{span}\{\partial\psi_x(\mathbf{x}) \mid \mathbf{z} \in Z_x^{t-1}\}. \end{aligned}$$

2. After each round $t \in \{0, \dots, T-1\}$, \mathcal{M} returns a solution $\bar{\mathbf{z}}^{t+1} = (\bar{\mathbf{x}}^{t+1}, \bar{\mathbf{y}}^{t+1})$ such that

$$\begin{aligned} \bar{\mathbf{x}}^{t+1} &\in \mathbf{x}^0 + \mathbf{P}_x^{-1} \text{span}\{\nabla_x f(\mathbf{z}) \mid \mathbf{z} \in Z_x^t\} + \mathbf{P}_x^{-1} \text{span}\{\partial\psi_x(\mathbf{x}) \mid \mathbf{z} \in Z_x^t\}, \\ \bar{\mathbf{y}}^{t+1} &\in \mathbf{y}^0 + \mathbf{P}_y^{-1} \text{span}\{\nabla_y f(\mathbf{z}) \mid \mathbf{z} \in Z_y^t\} + \mathbf{P}_y^{-1} \text{span}\{\partial\psi_y(\mathbf{y}) \mid \mathbf{z} \in Z_y^t\}. \end{aligned}$$

We remark that the span condition imposed by [Definition 1](#) is not absolutely necessary for defining distributed algorithms or proving lower bounds. This might be avoided by a more sophisticated reasoning of information-based complexity and resisting oracles. However, we adopt the gradient-span framework in this paper since it holds naturally for the majority of practical distributed algorithms.

Communication and oracle costs. For a given instance $P \in \mathcal{P}_{\text{SP}}$, we define the *communication cost* required by a distributed gradient-span algorithm \mathcal{M} on P , denoted by $T_P^{\mathcal{M}}$, as the smallest integer $k \in \{1, \dots, T\}$ such that the generated solution $\bar{\mathbf{z}}^k$ satisfies the target accuracy. Similarly, the total number of oracle queries with respect to x and y for instance P up to this point is given by the sizes of the accumulated histories, denoted respectively by $N_{x,P}^{\mathcal{M}} = |Z_x^{T_P^{\mathcal{M}}-1}|$ and $N_{y,P}^{\mathcal{M}} = |Z_y^{T_P^{\mathcal{M}}-1}|$.

Let c_x and c_y denote the respective computational costs of evaluating a single partial gradient $\nabla_x f$ and $\nabla_y f$. The (*worst-case*) *communication cost* and *oracle cost* of algorithm \mathcal{M} over the entire problem class \mathcal{P}_{SP} are defined by taking the supremum over all instances:

$$T_{\mathcal{P}_{\text{SP}}}^{\mathcal{M}} = \sup_{P \in \mathcal{P}_{\text{SP}}} T_P^{\mathcal{M}} \quad \text{and} \quad N_{\mathcal{P}_{\text{SP}}}^{\mathcal{M}} = \sup_{P \in \mathcal{P}_{\text{SP}}} (c_x N_{x,P}^{\mathcal{M}} + c_y N_{y,P}^{\mathcal{M}}).$$

Throughout the distributed environment considered in this paper, *we treat the communication cost as the primary performance metric*, as network communication typically forms the main bottleneck. The oracle cost serves as the secondary metric to measure the local computational effort.

2.3 Existing algorithms from literature

In this section, we review existing algorithms for solving SPs and analyze their communication and oracle costs within the distributed gradient-span algorithm framework. To keep the presentation concise, we summarize the methods and their limitations below, and defer their detailed algorithmic formulations, trajectories, and complexity propositions to [Section A](#).

Extragradient (EG). The classic EG method [[Nemirovski, 2004](#), [Juditsky et al., 2011](#)] naturally fits our framework. Its distributed execution requires two communication rounds per iteration to evaluate coupled partial gradients at both the current and extrapolated points. It provides a robust and natural baseline for communication and oracle costs.

Decoupled GDA (DGDA). The DGDA method [[Zindari et al., 2025](#)] attempts to reduce communication overhead by freezing the remote variable and taking multiple local gradient steps. While it achieves a fast logarithmic communication cost, its applications are highly restrictive: it only converges for weakly coupled strongly convex-strongly concave instances. For general problem class \mathcal{P}_{SP} , the delayed remote variables cause the local updates to drift, leading the method to diverge.

Catalyst acceleration. Using a Catalyst wrapper around EG (Cat-EG) [[Lin et al., 2020](#), [Yang et al., 2020](#), [Lan and Li, 2026](#)] accelerates the algorithm's dependence on the diagonal conditioning. However, this comes with **five significant caveats**: (i) it requires a complicated, multi-loop communication protocol and careful parameter tuning; (ii) it is highly sensitive to the inexactness of the diameter estimates \hat{D}_x and \hat{D}_y ; (iii) it introduces multiplicative logarithmic factors in the complexity; (iv) under certain conditioning, its theoretical complexity can be strictly worse than the unaccelerated EG baseline; and (v) it does not support extensions to multi-agent scenarios (cf. [Section 5](#)).

Four-loop method. The Cat-Cat-DAGDA method [Wang and Li, 2020] applies double Catalyst wrappers around a decoupled accelerated GDA to further accelerate the cross-coupling term. Despite this theoretical improvement, it shares all five caveats of Cat-EG, introduces even more complicated nested loops into the communication protocol, and adds further logarithmic factors. Consequently, it serves primarily as a theoretical benchmark rather than a practical method in our setting.

Other distributed stochastic gradient methods. Some recent papers [Zhang et al., 2024, Yoon et al., 2025, Yoon and Loizou, 2025] consider distributed SPs with stochastic gradient oracles. They propose different communication-efficient approaches; however, when applied to standard deterministic oracles considered in this paper, these methods fail to outperform EG.

Consequently, as summarized in Table 1, the classic EG method remains a formidable baseline for \mathcal{P}_{SP} , and improving its communication and oracle complexity remains a significant challenge.

3 Lower complexity bounds for SPs

In this section, we present lower complexity bounds for all algorithms in the distributed gradient-span family (cf. Section 2.2). We restrict our attention to this algorithm family because it is a key framework that captures a wide range of practical algorithms, as illustrated in Section 2.3.

In the standard non-distributed setting with full-gradient oracles, the lower complexity bounds for SPs have been well-studied in the literature [Nemirovsky, 1992, Zhang et al., 2022]. These bounds are typically expressed in terms of the number of full-gradient oracle queries and does not directly translate to the distributed setting. To this end, we carefully construct the worst-case instances and analyze the trajectories of the algorithms to establish the desired lower bounds.

Theorem 1. *Let the dimensions of \mathcal{E}_x and \mathcal{E}_y be sufficiently large such that $n_x \geq \frac{2L_{xy}D_xD_y}{3\epsilon} + \sqrt{\frac{3L_xD_x^2}{8\epsilon}} + 1$ and $n_y \geq \frac{2L_{xy}D_xD_y}{3\epsilon} + \sqrt{\frac{3L_yD_y^2}{8\epsilon}} + 1$. For any distributed gradient-span algorithm \mathcal{M} for \mathcal{P}_{SP} , we have:*

$$T_{\mathcal{P}_{\text{SP}}}^{\mathcal{M}} \geq \frac{2L_{xy}D_xD_y}{3\epsilon} - 3,$$

$$N_{\mathcal{P}_{\text{SP}}}^{\mathcal{M}} \geq \frac{c_x + c_y}{9} \frac{L_{xy}D_xD_y}{\epsilon} + \frac{c_x}{3} \sqrt{\frac{3L_xD_x^2}{32\epsilon}} + \frac{c_y}{3} \sqrt{\frac{3L_yD_y^2}{32\epsilon}} - \frac{2c_x + 2c_y}{3}.$$

Main proof idea. We now provide a simplified intuition for our lower bounds by focusing on standard Euclidean geometries with $n_x = n_y = n$, though our formal analysis is more general.

To establish the communication lower bound, we construct a difficult bilinear instance $F_{xy}(\mathbf{x}, \mathbf{y}) = \langle \mathbf{A}\mathbf{x} - \mathbf{b}, \mathbf{y} \rangle$. The upper bidiagonal matrix \mathbf{A} and vector \mathbf{b} are constructed as follows:

$$\mathbf{A} \propto \begin{pmatrix} \sqrt{2/1} & -\sqrt{1/2} & 0 & \dots & 0 \\ 0 & \sqrt{3/2} & -\sqrt{2/3} & \ddots & \vdots \\ 0 & 0 & \sqrt{4/3} & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & -\sqrt{(n-1)/n} \\ 0 & 0 & \dots & 0 & \sqrt{(n+1)/n} \end{pmatrix} \quad \text{and} \quad \mathbf{A}^T \mathbf{b} \propto (1, 0, \dots, 0)^T.$$

For this instance, the Agents \mathbf{x} and \mathbf{y} have to communicate back and forth to sequentially pass the new non-zero coordinates between each other, which leads to the communication lower bound. For the oracle lower bound, we additionally leverage quadratic instances for \mathbf{x} and \mathbf{y} respectively, and take the maximum cost over these three instances. We defer the detailed proofs to Section B. \square

Remark 1. Theorem 1 presents a communication lower bound depending only on the cross-coupled conditioning. In contrast, none of the existing algorithms reviewed in Section 2.3 match this lower bound. Specifically, the communication cost of the EG baseline is suboptimal due to its dependence on the diagonal conditioning. While advanced frameworks like Cat-Cat-DAGDA successfully isolate the communication cost from the diagonal conditioning, they suffer from highly complicated nested-loop designs and introduce poly-logarithmic overheads. This gap from the theoretical limit highlights the insufficiency of existing approaches and motivates our subsequent algorithmic developments.

4 Decoupled method for SPs

When designing a communication-efficient method, the primary challenge is enabling distributed agents to compute local solutions independently despite the presence of cross-coupled functions. To address this, we propose a clean algorithmic template (or communication protocol) *that reduces an SP into a sequence of coordinate-wise computational tasks*. We highlight the key results and insights below, deferring the detailed derivation to [Section C](#).

Assembled norm. Given parameters $\alpha_x, \alpha_y > 0$ (to be specified later), we equip the joint space $\mathcal{E} = \mathcal{E}_x \times \mathcal{E}_y$ with the assembled norm:

$$\|\mathbf{z}\|_{\mathcal{E}} = \langle \mathbf{P}\mathbf{z}, \mathbf{z} \rangle^{\frac{1}{2}} = \sqrt{\alpha_x \|\mathbf{x}\|_x^2 + \alpha_y \|\mathbf{y}\|_y^2} \quad \text{for all } \mathbf{z} \in \mathcal{E}, \quad (2)$$

which corresponds to the block diagonal linear operator $\mathbf{P} = \alpha_x \mathbf{P}_x \oplus \alpha_y \mathbf{P}_y$.

Template DM-SP. [Algorithm 1](#) outlines the Decoupled Method for Saddle Problems (DM-SP), which adapts the abstract framework of the Reduced-Operator Method [[Nesterov, 2023](#)] for distributed environments. The algorithm maintains a sequence of anchor points \mathbf{v}^t and proceeds iteratively.

First, the agents decouple the joint problem by fixing the remote variable at the current anchor \mathbf{v}^t . This allows Agent x and Agent y to independently and concurrently solve their respective regularized local subproblems up to target accuracies δ_x^{t+1} and δ_y^{t+1} (Lines 4 and 5). Specifically, Agent x aims to approximately compute $\arg \min_{\mathbf{x} \in \text{dom } \psi_x} [f(\mathbf{x}, \mathbf{v}_y^t) + \frac{\alpha_x \lambda_{t+1}}{2} \|\mathbf{x} - \mathbf{v}_x^t\|_x^2 + \psi_x(\mathbf{x})]$ by finding a point \mathbf{x}^{t+1} whose regularized subgradient norm satisfies the exact mathematical bound specified in Line 4. Agent y symmetrically performs an approximate minimization for its corresponding objective $-f(\mathbf{v}_x^t, \mathbf{y}) + \frac{\alpha_y \lambda_{t+1}}{2} \|\mathbf{y} - \mathbf{v}_y^t\|_y^2 + \psi_y(\mathbf{y})$.

Following this local computation phase, the agents perform exactly two communication rounds to complete the iteration. In the first round (Line 6), the agents exchange their locally computed approximate solutions to assemble the joint intermediate point $\mathbf{z}^{t+1} = (\mathbf{x}^{t+1}, \mathbf{y}^{t+1})$. In the second round (Line 7), they use this assembled point to evaluate their local partial gradients, which they then exchange to form the full joint operator $V_{\psi}(\mathbf{z}^{t+1})$.

Finally, using this assembled operator, the agents compute a closed-form step size a_{t+1} , update the running ergodic average $\bar{\mathbf{z}}^{t+1}$, and perform an extragradient-like step to generate the next anchor \mathbf{v}^{t+1} (Lines 8 and 9). By structuring the method this way, DM-SP cleanly reduces the coupled SP into isolated coordinate-wise tasks with minimal communication overhead.

Algorithm 1 DM-SP($f, (\psi_x, \psi_y), \mathbf{z}^0, (\lambda_t)_{t \geq 1}, (\alpha_x, \alpha_y)$)

1: $\mathbf{v}^0 = (\mathbf{v}_x^0, \mathbf{v}_y^0) = \mathbf{z}^0$.

2: **for** $t = 0, 1, \dots, T - 1$ **do**

3: Let $\delta_x^{t+1} = \frac{\alpha_x \lambda_{t+1}}{2}$ and $\delta_y^{t+1} = \frac{\alpha_y \lambda_{t+1}}{2}$.

4: Agent x finds \mathbf{x}^{t+1} and $\psi'_x(\mathbf{x}^{t+1}) \in \partial \psi_x(\mathbf{x}^{t+1})$ such that

$$\|\nabla_x f(\mathbf{x}^{t+1}, \mathbf{v}_y^t) + \alpha_x \lambda_{t+1} (\mathbf{x}^{t+1} - \mathbf{v}_x^t) + \psi'_x(\mathbf{x}^{t+1})\|_{x^*} \leq \delta_x^{t+1} \|\mathbf{x}^{t+1} - \mathbf{v}_x^t\|_x.$$

5: Agent y finds \mathbf{y}^{t+1} and $\psi'_y(\mathbf{y}^{t+1}) \in \partial \psi_y(\mathbf{y}^{t+1})$ such that

$$\|-\nabla_y f(\mathbf{v}_x^t, \mathbf{y}^{t+1}) + \alpha_y \lambda_{t+1} (\mathbf{y}^{t+1} - \mathbf{v}_y^t) + \psi'_y(\mathbf{y}^{t+1})\|_{y^*} \leq \delta_y^{t+1} \|\mathbf{y}^{t+1} - \mathbf{v}_y^t\|_y.$$

6: Exchange \mathbf{x}^{t+1} and \mathbf{y}^{t+1} to assemble $\mathbf{z}^{t+1} = (\mathbf{x}^{t+1}, \mathbf{y}^{t+1})$.

7: Calculate corresponding coordinates of $V_{\psi}(\mathbf{z}^{t+1})$, then exchange to assemble:

$$V_{\psi}(\mathbf{z}^{t+1}) = (\nabla_x f(\mathbf{z}^{t+1}) + \psi'_x(\mathbf{x}^{t+1}), -\nabla_y f(\mathbf{z}^{t+1}) + \psi'_y(\mathbf{y}^{t+1})).$$

8: Let $a_{t+1} = \frac{2 \langle V_{\psi}(\mathbf{z}^{t+1}), \mathbf{v}^t - \mathbf{z}^{t+1} \rangle}{\|V_{\psi}(\mathbf{z}^{t+1})\|_{\mathcal{E}^*}^2}$ and generate solution $\bar{\mathbf{z}}^{t+1} = (\sum_{i=1}^{t+1} a_i)^{-1} \sum_{i=1}^{t+1} a_i \mathbf{z}^i$.

9: $\mathbf{v}^{t+1} = \arg \min_{\mathbf{v} \in Q} [a_{t+1} \langle V_{\psi}(\mathbf{z}^{t+1}), \mathbf{v} \rangle + \frac{1}{2} \|\mathbf{v} - \mathbf{v}^t\|_{\mathcal{E}}^2]$.

10: **end for**

We refer to [Algorithm 1](#) as a template method because we have not yet specified the implementations for the local computations in Lines 4 and 5. Provided that the trajectories of these local computations

satisfy the span conditions outlined in [Definition 1](#), the template DM-SP procedure formally qualifies as a distributed gradient-span algorithm.

Theorem 2. Consider the DM-SP template applied to \mathcal{P}_{SP} , assuming its local trajectories satisfy [Definition 1](#). With the parameter choices of $\alpha_x = \frac{L_{xy}D_y}{D_x}$, $\alpha_y = \frac{L_{xy}D_x}{D_y}$, and $\lambda_t \equiv \lambda = 2$, we have:

$$T_{\mathcal{P}_{\text{SP}}}^{\text{DM-SP}} \leq 2 + 4 \frac{L_{xy}D_xD_y}{\epsilon}.$$

Remark 2 (Communication optimality). [Theorem 2](#) shows that the communication cost of DM-SP depends only on the cross-coupled conditioning $L_{xy}D_xD_y$, independent of the diagonal conditioning. As established in [Section 3](#), this $\mathcal{O}\left(\frac{L_{xy}D_xD_y}{\epsilon}\right)$ communication cost matches the theoretical lower bound within the family of distributed gradient-span algorithms. To our knowledge, DM-SP is the first communication-optimal algorithm for distributed SPs.

Remark 3 (Robustness to inexact distance estimates). Let us consider a practical scenario where the algorithm may not have the precise values of D_x and D_y in advance, but it has access to upper estimates $\hat{D}_x \geq D_x$ and $\hat{D}_y \geq D_y$. Let

$\theta \triangleq \frac{D_x\hat{D}_y}{\hat{D}_xD_y} + \frac{D_y\hat{D}_x}{\hat{D}_yD_x}$, which quantifies the disproportionality

between the true distance parameters and their estimates. Note that $\theta \geq 2$ with equality if and only if the estimates are proportional, i.e., $\hat{D}_x/D_x = \hat{D}_y/D_y$. Now, consider the DM-SP template applied to \mathcal{P}_{SP} . With the parameter choices of $\alpha_x = \frac{L_{xy}\hat{D}_y}{\hat{D}_x}$, $\alpha_y = \frac{L_{xy}\hat{D}_x}{\hat{D}_y}$, and $\lambda_t \equiv \lambda = 2$, we have:

$$T_{\mathcal{P}_{\text{SP}}}^{\text{DM-SP}} \leq 2 + 2\theta \frac{L_{xy}D_xD_y}{\epsilon}.$$

In particular, θ provides a *scale-invariant* robustness compared to existing accelerated frameworks. As shown in [Table 1](#), the communication complexities of Cat-EG and Cat-Cat-DAGDA scale directly with the product of the estimates, $\hat{D}_x\hat{D}_y$. Consequently, if both agents conservatively overestimate their domain sizes by a uniform factor $c \gg 1$ (i.e., $\hat{D}_x = cD_x$ and $\hat{D}_y = cD_y$), the communication cost of Catalyst-based methods inflates by a massive factor of c^2 . For DM-SP, however, this uniform overestimation perfectly cancels out, yielding $\theta = 2$.

Concrete implementation. Furthermore, let us instantiate this template with the specific local solver introduced later in [Section D](#) ([Eq. \(8\)](#)) to obtain a concrete implementation of DM-SP.

Theorem 3. Consider the DM-SP algorithm equipped with the local solvers in [Eq. \(8\)](#), applied to \mathcal{P}_{SP} . With the same choices of α_x , α_y , and $(\lambda_t)_{t \geq 0}$ as in [Theorem 2](#), we have:

$$N_{\mathcal{P}_{\text{SP}}}^{\text{DM-SP}} = (c_x + c_y) \frac{2L_{xy}D_xD_y}{\epsilon} + 102 \left(\frac{L_{xy}D_xD_y}{\epsilon} \right)^{\frac{1}{2}} \left(c_x \left(\frac{L_xD_x^2}{\epsilon} \right)^{\frac{1}{2}} + c_y \left(\frac{L_yD_y^2}{\epsilon} \right)^{\frac{1}{2}} \right).$$

As detailed in [Section 2.3](#), alternative distributed algorithms either fail to outperform EG in general for \mathcal{P}_{SP} or suffer from certain theoretical and practical caveats. Consequently, the EG method remains a crucial baseline for oracle complexity, which we now compare against.

Remark 4. Let us compare the oracle cost of DM-SP against the EG baseline, which requires

$$N_{\mathcal{P}_{\text{SP}}}^{\text{EG}} = (c_x + c_y) \cdot \left(\frac{L_{xy}D_xD_y}{\epsilon} + \frac{L_xD_x^2}{\epsilon} + \frac{L_yD_y^2}{\epsilon} \right)$$

oracle costs ([Proposition 5](#)). We conclude that the computational cost of DM-SP is consistently bounded by that of EG. Furthermore, it yields a substantial improvement when

$$L_xD_x^2 + L_yD_y^2 \gg L_{xy}D_xD_y + \sqrt{L_{xy}D_xD_y} \cdot \left(\frac{c_x}{c_x + c_y} \sqrt{L_xD_x^2} + \frac{c_y}{c_x + c_y} \sqrt{L_yD_y^2} \right).$$

For instance, assuming uniform oracle costs ($c_x = c_y$), this improvement occurs when the diagonal conditioning dominates the cross-coupled conditioning, i.e., $L_xD_x^2 + L_yD_y^2 \gg \sqrt{L_{xy}D_xD_y}$. To our knowledge, DM-SP is the first method to consistently improve upon the EG oracle cost for \mathcal{P}_{SP} . We note, however, that a gap remains between this achieved oracle cost and the theoretical oracle lower bound established in [Theorem 1](#), which is a known open question even for non-distributed SPs.

Novelty. Our DM-SP is built upon a novel multi-stage reduction. We first leverage Reduced-Operator Method to reduce the problem to a Monteiro-Svaiter Subproblem (MSS). Then and crucially, we show that when this subproblem is weakly coupled, it can be solved by a Fully Decoupled Solver in one communication round. Consequently, the problem is further reduced to coordinate-wise Minimization of Residual Norms (MRNs). Finally, by exploiting the strong maximal monotonicity (cf. (A2)), the agents can apply existing accelerated methods to solve the MRNs to desired accuracy.

5 Variational inequality problems with distributed oracles

Motivation. Thus far, we have studied SPs, which naturally model two-player zero-sum games. To capture more complex multiagent interactions (such as equilibrium computation in multiplayer general-sum games, network routing, and multiagent reinforcement learning), we extend our algorithmic framework to the broader class of monotone Variational Inequality Problems (VIPs) with separable composite terms and distributed oracles. Due to lack of space, we briefly outline the problem class and our results here, while defer the detailed presentation to Section E.

Problem class. We consider a distributed multiagent setting with a star communication network over a product space $\mathcal{E} = \mathcal{E}_1 \times \dots \times \mathcal{E}_K$, where a joint decision variable is partitioned among K autonomous agents as $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_K)$. Let us consider the problem class \mathcal{P}_{VIP} as follows:

- **Operators and local components:** The problem is governed by a joint monotone operator $V(\mathbf{z}) = (V_1(\mathbf{z}), \dots, V_K(\mathbf{z}))$ and a separable local composite function $\psi(\mathbf{z}) = \sum_{i=1}^K \psi_i(\mathbf{z}_i)$ defined on $Q = \text{dom } \psi_1 \times \dots \times \text{dom } \psi_K$. We assume the solution set is bounded by local distance parameters $D_i > 0$ for each agent. Furthermore, the operator satisfies block-wise Lipschitz continuity: for any fixed \mathbf{z}_{-j} , the mapping $V_i(\mathbf{z}_j; \mathbf{z}_{-j})$ is L_{ij} -Lipschitz continuous with respect to \mathbf{z}_j .
- **Distributed oracles:** Each Agent $i \in [K]$ controls its local variable $\mathbf{z}_i \in \text{dom } \psi_i$, has access to its private function ψ_i and a partial oracle $\mathcal{O}_i(\mathbf{z}) = V_i(\mathbf{z})$.
- **Accuracy measure:** The goal is to find an ϵ -approximate solution $\bar{\mathbf{z}} \in Q$ such that the restricted gap $\Delta(\bar{\mathbf{z}}) \triangleq \sup_{\mathbf{z} \in \mathcal{B} \cap Q} [\langle V(\mathbf{z}), \bar{\mathbf{z}} - \mathbf{z} \rangle + \psi(\bar{\mathbf{z}}) - \psi(\mathbf{z})]$ satisfies $\Delta(\bar{\mathbf{z}}) \leq \epsilon$, where the bounded domain \mathcal{B} is defined by the initial point \mathbf{z}^0 and distance parameters $D_i, i \in [K]$.

Conditionings. Let $\bar{L}_{ij} \triangleq \max\{L_{ij}, L_{ji}\}$. Similar to SPs, we hereby define the *cross-coupled conditioning*, denoted by $\sum_{i \in [K]} A_i$ with $A_i \triangleq D_i \sum_{j \neq i} \bar{L}_{ij} D_j$, that quantifies the cross-dependencies between the agents over the network. The *diagonal conditioning*, denoted by $\sum_{i \in [K]} B_i$ with $B_i \triangleq L_{ii} D_i^2$, measures the self-dependency within a single agent's domain.

New state-of-the-art communication cost. In the multiagent setting, Extragradient (EG) remains the state-of-the-art method, requiring $T_{\mathcal{P}_{\text{VIP}}}^{\text{EG}} = \mathcal{O}(\sum_{i \in [K]} \frac{A_i + B_i}{\epsilon})$ communication cost, which depends on both the cross-coupled and diagonal conditionings. By extending our decoupled template to distributed VIPs, we propose DM-VIP, and establish a much better communication cost.

Theorem 4. Consider the DM-VIP algorithm for problem class \mathcal{P}_{VIP} . We have

$$T_{\mathcal{P}_{\text{VIP}}}^{\text{DM-VIP}} = \mathcal{O}\left(\sum_{i \in [K]} A_i / \epsilon\right).$$

Remark 5. The communication cost in Theorem 4 completely drops the dependence on the diagonal conditioning. Consequently, we have substantially improved the state-of-the-art EG cost when the diagonal conditioning dominates, i.e., when $\sum_i B_i \gg \sum_i A_i$.

6 Conclusion and limitations

This paper studies communication and oracle costs in distributed SPs and VIPs. For the class of SPs, we settle the communication complexity in the distributed setup within gradient-span framework, and consistently improve the long-standing oracle cost of EG method. For the class of distributed VIPs, we improve the state-of-the-art communication cost. The following directions are not addressed in this paper and are left for future work: (a) closing the gap of oracle costs; (b) showing lower bound for non-zero-sum games; and (c) showing information-theoretic lower bounds for randomized methods.

Acknowledgments and Disclosure of Funding

The authors thank Ali Zindari, Ehsan Goharshady, and Krishnendu Chatterjee for their helpful discussions and suggestions on this paper. RL acknowledges the support of ERC CoG 863818 (ForM-SMArt) and Austrian Science Fund (FWF) 10.55776/COE12. Gemini Pro 3.1 helps with part of the writing and analysis.

References

- Aharon Ben-Tal and Arkadi Nemirovski. Robust optimization—methodology and applications. *Mathematical programming*, 92(3):453–480, 2002. [1](#)
- Aleksandr Beznosikov, Valentin Samokhin, and Alexander Gasnikov. Distributed saddle point problems: lower bounds, near-optimal and robust algorithms. *Optimization Methods and Software*, pages 1–18, 2025. [1](#)
- Radu Ioan Boş and Enis Chenchene. Extra-gradient method with flexible anchoring: Strong convergence and fast residual decay. *arXiv preprint arXiv:2410.14369*, 2024. [E.4.2](#)
- Vincent Conitzer and Tuomas Sandholm. Communication complexity as a lower bound for learning in games. In *Proceedings of the twenty-first international conference on Machine learning*, page 24, 2004. [1](#)
- Yuyang Deng and Mehrdad Mahdavi. Local stochastic gradient descent ascent: Convergence analysis and communication efficiency. In *International Conference on Artificial Intelligence and Statistics*, pages 1387–1395. PMLR, 2021. [1](#)
- Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. [1](#)
- Sergiu Hart and Yishay Mansour. How long to equilibrium? the communication complexity of uncoupled equilibrium procedures. *Games and Economic Behavior*, 69(1):107–126, 2010. [1](#)
- Junling Hu, Michael P Wellman, et al. Multiagent reinforcement learning: theoretical framework and an algorithm. In *ICML*, volume 98, pages 242–250, 1998. [1](#)
- Anatoli Juditsky, Arkadi Nemirovski, et al. First order methods for nonsmooth convex large-scale optimization, ii: utilizing problems structure. *Optimization for Machine Learning*, 30(9):149–183, 2011. [2.3](#), [5](#), [E](#), [20](#), [6](#)
- Guanghui Lan and Yan Li. A novel catalyst scheme for stochastic minimax optimization. *Mathematical Programming*, pages 1–49, 2026. [2.3](#), [6](#)
- Guanghui Lan, Yuyuan Ouyang, and Zhe Zhang. Optimal and parameter-free gradient minimization methods for convex and nonconvex optimization. *arXiv preprint arXiv:2310.12139*, 2023. [C.4](#), [C.4](#), [18](#), [C.4](#)
- Tianyi Lin, Chi Jin, and Michael I Jordan. Near-optimal algorithms for minimax optimization. In *Conference on learning theory*, pages 2738–2779. PMLR, 2020. [2.3](#)
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. [1](#)
- Renato DC Monteiro and Benar Fux Svaiter. An accelerated hybrid proximal extragradient method for convex optimization and its implications to second-order methods. *SIAM Journal on Optimization*, 23(2):1092–1125, 2013. [C.2](#)
- Arkadi Nemirovski. Prox-method with rate of convergence $O(1/t)$ for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004. [1](#), [2.3](#), [E](#)
- Arkadi S Nemirovsky. Information-based complexity of linear operator equations. *Journal of Complexity*, 8(2): 153–175, 1992. [3](#)
- Yurii Nesterov. High-order reduced-gradient methods for composite variational inequalities. *arXiv preprint arXiv:2311.15154*, 2023. [4](#), [C.2](#)

- Yurii E. Nesterov. *Introductory Lectures on Convex Optimization - A Basic Course*, volume 87 of *Applied Optimization*. Springer, 2004. ISBN 978-1-4613-4691-3. doi: 10.1007/978-1-4419-8853-9. URL <https://doi.org/10.1007/978-1-4419-8853-9>. 2.2, B.1
- Noam Nisan and Ilya Segal. The communication requirements of efficient allocations and supporting prices. *Journal of Economic Theory*, 129(1):192–224, 2006. 1
- J Ben Rosen. Existence and uniqueness of equilibrium points for concave n-person games. *Econometrica: Journal of the Econometric Society*, pages 520–534, 1965. 1
- Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1):83–112, 2017. 1
- Paul Tseng. On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics*, 60(1-2):237–252, 1995. 1
- John von Neumann and Oskar Morgenstern. *Theory of games and economic behavior*. Princeton university press, 1947. 1
- Yuanhao Wang and Jian Li. Improved algorithms for convex-concave minimax optimization. *Advances in Neural Information Processing Systems*, 33:4800–4810, 2020. 2.3, A, 7
- Junchi Yang, Siqi Zhang, Negar Kiyavash, and Niao He. A catalyst framework for minimax optimization. *Advances in Neural Information Processing Systems*, 33:5667–5678, 2020. 2.3, 6
- TaeHo Yoon and Nicolas Loizou. Pearl-prox: Proximal algorithm for resolving player drift in multiplayer federated learning. In *OPT 2025: Optimization for Machine Learning*, 2025. 1, 2.3
- TaeHo Yoon, Sayantan Choudhury, and Nicolas Loizou. Multiplayer federated learning: Reaching equilibrium with less communication. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=9JX8XrTVEz>. 1, 2.3
- Junyu Zhang, Mingyi Hong, and Shuzhong Zhang. On lower iteration complexity bounds for the convex concave saddle point problems. *Mathematical Programming*, 194(1):901–935, 2022. 3
- Siqi Zhang, Sayantan Choudhury, Sebastian U Stich, and Nicolas Loizou. Communication-efficient gradient descent-ascent methods for distributed variational inequalities: Unified analysis and local updates. In *The Twelfth International Conference on Learning Representations*, 2024. 1, 2.3
- Ali Zindari, Parham Yazdkhasti, Anton Rodomanov, Tatjana Chavdarova, and Sebastian U Stich. Decoupled SGDA for games with intermittent strategy communication. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=ZYkFTSEZ6k>. 1, 2.3, A

A Detailed review of existing algorithms

In this section, we provide the detailed formulations, distributed trajectories, and formal complexity results for the algorithms summarized in [Section 2.3](#).

We recall that $\mathcal{P}_{\text{SP}}^\circ$ denotes the subclass of \mathcal{P}_{SP} where the local components ψ_x and ψ_y are zero. In this setting, the problem reduces to finding a saddle point of a smooth convex-concave function f . We note that many algorithms discussed below do not handle general composite functions and only deal with problem instances from this non-composite subclass.

As in [Remark 3](#), let us consider a practical scenario where the algorithms may not have the precise values of D_x and D_y in advance, but they have access to upper estimates $\hat{D}_x \geq D_x$ and $\hat{D}_y \geq D_y$. Let

$$\theta \triangleq \frac{D_x \hat{D}_y}{\hat{D}_x D_y} + \frac{D_y \hat{D}_x}{\hat{D}_y D_x},$$

which quantifies the disproportionality between the true distance parameters and their estimates.

Extragradient (EG). Suppose we apply EG to a non-composite problem instance $P = (f, 0, 0, \mathbf{z}^0) \in \mathcal{P}_{\text{SP}}^\circ$. The method iteratively maintains and updates two sequences of variables $\mathbf{v}^k, \mathbf{z}^k \in \mathcal{E}_x \times \mathcal{E}_y$. Initializing with $\mathbf{v}^0 = \mathbf{z}^0$, the updates for iteration $k = 0, 1, \dots, K-1$ are given by:

$$\begin{aligned} \mathbf{z}^{k+1} &= \mathbf{v}^k - \eta^k \mathbf{P}^{-1} V f(\mathbf{v}^k), \\ \mathbf{v}^{k+1} &= \mathbf{v}^k - \eta^k \mathbf{P}^{-1} V f(\mathbf{z}^{k+1}), \end{aligned}$$

where $\eta^k > 0$ is the step size at iteration k .

Although the EG method is conventionally formulated as a K -iteration loop, its distributed execution requires $2K$ communication rounds. In each iteration, the agents must synchronize twice to evaluate the coupled partial gradients at \mathbf{v}^k and at \mathbf{z}^{k+1} . Thus, the EG method corresponds to a distributed gradient-span algorithm with $T = 2K$ rounds and local step lengths $\tau_x^t = \tau_y^t = 1$, where the trajectories are given by:

$$\hat{Z}_x^{2k} = \hat{Z}_y^{2k} = \{\mathbf{v}^k\} \quad \text{and} \quad \hat{Z}_x^{2k+1} = \hat{Z}_y^{2k+1} = \{\mathbf{z}^{k+1}\},$$

for $k = 0, \dots, K-1$.

To align with [Definition 1](#), the method generates a solution $\bar{\mathbf{z}}^{t+1}$ after each round $t \in \{0, \dots, T-1\}$. Because each iteration requires two communication rounds, the algorithm effectively updates its output only every two rounds. Specifically, after completing iteration k (which corresponds to round $t = 2k+1$), it outputs the ergodic average $\bar{\mathbf{z}}^{2k+2} = \frac{1}{\sum_{i=0}^k \eta^i} (\sum_{i=0}^k \eta^i \mathbf{z}^{i+1})$. During the intermediate rounds (at $t = 2k$), it simply retains the previous solution by setting $\bar{\mathbf{z}}^{2k+1} = \bar{\mathbf{z}}^{2k}$ (where $\bar{\mathbf{z}}^0 = \mathbf{z}^0$). Since the trajectories and the generated solutions are formed by linear combinations of the evaluated gradients, they satisfy the span conditions in [Definition 1](#).

The complexity of the EG method provides a natural baseline. Translating the classic results into our distributed complexity measures yields the following upper bounds.

Proposition 5 ([Juditsky et al. 2011](#), Eq. (6.21)). *Consider the EG method applied to $\mathcal{P}_{\text{SP}}^\circ$. With the parameter choices of $\alpha_x = \frac{L_x \hat{D}_x + L_{xy} \hat{D}_y}{\hat{D}_x}$, $\alpha_y = \frac{L_y \hat{D}_y + L_{xy} \hat{D}_x}{\hat{D}_y}$, and $\eta^k \equiv 1$, we have*

$$\begin{aligned} T_{\mathcal{P}_{\text{SP}}^\circ}^{\text{EG}} &= \theta \cdot \frac{L_{xy} D_x D_y}{\epsilon} + \frac{L_x D_x^2}{\epsilon} + \frac{L_y D_y^2}{\epsilon}, \\ N_{\mathcal{P}_{\text{SP}}^\circ}^{\text{EG}} &= (c_x + c_y) \cdot \left(\theta \cdot \frac{L_{xy} D_x D_y}{\epsilon} + \frac{L_x D_x^2}{\epsilon} + \frac{L_y D_y^2}{\epsilon} \right). \end{aligned}$$

Decoupled GDA. Let us apply DGDA to an instance $P \in \mathcal{P}_{\text{SP}}^\circ$ with a fixed local trajectory length τ . The algorithm maintains local iterates $\mathbf{x}^{t,l}$ and $\mathbf{y}^{t,l}$ for round $t = 0, \dots, T-1$ and local step $l = 0, \dots, \tau$. At the beginning of round t , the agents synchronize by exchanging their latest local

iterates. Specifically, Agent x receives $\hat{\mathbf{y}}^t$ and Agent y receives $\hat{\mathbf{x}}^t$, defined as:

$$\hat{\mathbf{x}}^t \triangleq \begin{cases} \mathbf{x}^0 & \text{if } t = 0 \\ \mathbf{x}^{t-1,\tau} & \text{if } t > 0 \end{cases} \quad \text{and} \quad \hat{\mathbf{y}}^t \triangleq \begin{cases} \mathbf{y}^0 & \text{if } t = 0 \\ \mathbf{y}^{t-1,\tau} & \text{if } t > 0 \end{cases}.$$

The agents then initialize their local variables for the current round as $\mathbf{x}^{t,0} = \hat{\mathbf{x}}^t$ and $\mathbf{y}^{t,0} = \hat{\mathbf{y}}^t$. With the remote variables firmly fixed, the agents execute τ local gradient steps. For $l = 0, \dots, \tau - 1$, the local updates are given by:

$$\begin{aligned} \mathbf{x}^{t,l+1} &= \mathbf{x}^{t,l} - \eta_x \mathbf{P}_x^{-1} \nabla_x f(\mathbf{x}^{t,l}, \hat{\mathbf{y}}^t), \\ \mathbf{y}^{t,l+1} &= \mathbf{y}^{t,l} + \eta_y \mathbf{P}_y^{-1} \nabla_y f(\hat{\mathbf{x}}^t, \mathbf{y}^{t,l}), \end{aligned}$$

where $\eta_x, \eta_y > 0$ are the local step sizes.

The DGDA algorithm yields the following trajectories:

$$\hat{Z}_x^t = \{(\mathbf{x}^{t,l}, \hat{\mathbf{y}}^t)\}_{l=0}^{\tau-1} \quad \text{and} \quad \hat{Z}_y^t = \{(\hat{\mathbf{x}}^t, \mathbf{y}^{t,l})\}_{l=0}^{\tau-1}.$$

After each round $t \in \{0, \dots, T - 1\}$, the method returns the updated local iterates as the current solution $\bar{\mathbf{z}}^{t+1} = (\mathbf{x}^{t,\tau}, \mathbf{y}^{t,\tau})$. Because the trajectories and the generated solutions are formed entirely by linear combinations of the evaluated gradients, DGDA is indeed a distributed gradient-span algorithm by [Definition 1](#).

While DGDA lies perfectly in our framework, its theoretical guarantees are highly restrictive. The algorithm is only proven to converge for restricted strongly convex-strongly concave problem instances where the cross-coupling between the variables is sufficiently weak [[Zindari et al., 2025](#)]. In this narrowly defined regime, DGDA achieves a logarithmic communication complexity of $\mathcal{O}(\log \frac{1}{\epsilon})$, which is a clear improvement over the EG baseline. However, for general problem instances in $\mathcal{P}_{\text{SP}}^\circ$ with stronger coupling, the delayed remote variables cause the local updates to drift, ultimately leading the method to diverge. Consequently, DGDA fails to provide any meaningful complexity guarantee for the general problem class $\mathcal{P}_{\text{SP}}^\circ$ under consideration.

Catalyst acceleration. Catalyst methods first add small regularization terms to the objective: $f(\mathbf{x}, \mathbf{y}) + \frac{\epsilon}{4\hat{D}_x^2} \|\mathbf{x} - \mathbf{x}^0\|_x^2 - \frac{\epsilon}{4\hat{D}_y^2} \|\mathbf{y} - \mathbf{y}^0\|_y^2$, reducing the problem to a strongly convex-strongly concave one. Then, Catalyst introduces an outer loop, indexed by $k = 0, 1, \dots, K - 1$, designed to balance the conditioning between the two variables. For instance, when the conditioning of x is worse (i.e. $L_x \hat{D}_x^2 \geq L_y \hat{D}_y^2$), the method carefully maintains an extrapolation sequence $(\tilde{\mathbf{x}}^k)_{k=0}^{K-1}$ and, in each outer iteration, adds a proximal term $\frac{\lambda_x}{2} \|\mathbf{x} - \tilde{\mathbf{x}}^k\|_x^2$ to the objective. Conversely, if the conditioning of y is worse, the outer loop would instead maintain an extrapolation sequence for y and add a corresponding regularization term for y . An inner base algorithm (EG in this case) is then deployed to solve this regularized subproblem to a specified accuracy.

To simplify the notation, let $L_{\max} = \max\{L_x, L_y, L_{xy}\}$.

Proposition 6 ([Yang et al. \[2020\]](#), [Lan and Li \[2026\]](#)). *Consider the Catalyst framework equipped with the EG method as the inner solver, denoted by Cat-EG, applied to $\mathcal{P}_{\text{SP}}^\circ$. Then, we have:*

$$\begin{aligned} T_{\mathcal{P}_{\text{SP}}^\circ}^{\text{Cat-EG}} &= \mathcal{O}\left(\left(\frac{L_{\max} \hat{D}_x \hat{D}_y}{\epsilon} + \sqrt{\frac{L_x \hat{D}_x^2}{\epsilon}} + \sqrt{\frac{L_y \hat{D}_y^2}{\epsilon}}\right) \log^2\left(\frac{1}{\epsilon}\right)\right), \\ N_{\mathcal{P}_{\text{SP}}^\circ}^{\text{Cat-EG}} &= (c_x + c_y) \cdot \mathcal{O}\left(\left(\frac{L_{\max} \hat{D}_x \hat{D}_y}{\epsilon} + \sqrt{\frac{L_x \hat{D}_x^2}{\epsilon}} + \sqrt{\frac{L_y \hat{D}_y^2}{\epsilon}}\right) \log^2\left(\frac{1}{\epsilon}\right)\right). \end{aligned}$$

Now, let us explain the caveats we mentioned in [Section 2.3](#) regarding the Cat-EG method in more detail.

To explain the second caveat of Cat-EG, its sensitivity to inexact diameter estimates, we compare it with the EG baseline and our proposed method. Because Catalyst uses \hat{D}_x and \hat{D}_y to set the initial regularization, its complexity scales directly with these estimates rather than the true distances D_x and D_y . Specifically, the diagonal terms in EG or DM-SP depend strictly on the true distances, whereas in Cat-EG they scale with \hat{D}_x and \hat{D}_y . Similarly, for the cross-coupled term, EG or DM-SP depends on the true distances multiplied by the proportionality ratio θ . If the estimates are loose but

proportional (e.g., $\hat{D}_x = cD_x$ and $\hat{D}_y = cD_y$), θ remains 2, leaving the complexity unaffected by the overestimation factor c . In contrast, the Cat-EG coupled term scales with $\hat{D}_x\hat{D}_y$, meaning any overestimation of $\hat{D}_x \gg D_x$ or $\hat{D}_y \gg D_y$ directly inflates the bound. Consequently, EG and our proposed DM-SP method are much more robust to inexact distance estimates, provided the estimates are roughly proportional.

Furthermore, we note that the fourth caveat of Cat-EG (performing worse than unaccelerated EG) can be easily verified. For instance, consider a problem instance where $L_x = 10^6$, $L_{xy} = 1$, $L_y = 10^{-6}$, $D_x = 10^{-3}$, and $D_y = 10^3$. Under this conditioning, the theoretical upper bound of Cat-EG significantly exceeds that of standard EG.

Four-loop method. The Cat-Cat-DAGDA method [Wang and Li, 2020] first adds small $\mathcal{O}(\epsilon)$ regularizations to both \mathbf{x} and \mathbf{y} to ensure the objective is strongly convex-strongly concave. It then executes four nested loops, which justifies our naming convention: (i) The first loop is a Catalyst outer loop adding a proximal regularization term to the primal variable; (ii) The second loop is another Catalyst outer loop adding a proximal regularization term to the dual variable; (iii) The third loop manages communication by exchanging and freezing the remote variables, identical to DGDA; and (iv) The fourth loop performs local computations, but unlike the standard gradient steps in DGDA, it employs an accelerated gradient method (hence DAGDA) to solve the inner subproblems.

While the fourth loop of the method is a local computation loop, the first three loops all require communication rounds. Let us state the complexity results of this method without going into the tedious details of the algorithm.

Proposition 7 (Wang and Li [2020]). *Consider the Cat-Cat-DAGDA method applied to $\mathcal{P}_{\text{SP}}^\circ$. Then, we have:*

$$T_{\mathcal{P}_{\text{SP}}^\circ}^{\text{Cat-Cat-DAGDA}} = \mathcal{O}\left(\frac{L_{xy}\hat{D}_x\hat{D}_y}{\epsilon} \log^3\left(\frac{1}{\epsilon}\right)\right),$$

$$N_{\mathcal{P}_{\text{SP}}^\circ}^{\text{Cat-Cat-DAGDA}} = (c_x + c_y) \cdot \mathcal{O}\left(\left(\frac{\sqrt{L_{\max}L_{xy}}\hat{D}_x\hat{D}_y}{\epsilon} + \sqrt{\frac{L_x\hat{D}_x^2}{\epsilon}} + \sqrt{\frac{L_y\hat{D}_y^2}{\epsilon}}\right) \log^4\left(\frac{1}{\epsilon}\right)\right).$$

B Lower complexity bounds

In this section, we provide detailed proofs for the lower complexity bounds of distributed gradient-span algorithms applied to \mathcal{P}_{SP} . We focus on a subclass of unconstrained SPs with $\psi_x = \psi_y = 0$, denoted $\mathcal{P}_{\text{SP}}^\circ$. We assume the initial points are $\mathbf{x}^0 = \mathbf{0}$ and $\mathbf{y}^0 = \mathbf{0}$. Lower bounds established for this subclass hold for the general class \mathcal{P}_{SP} .

Geometry. Let us consider the spaces $\mathcal{E}_x = \mathbb{R}^{n_x}$ and $\mathcal{E}_y = \mathbb{R}^{n_y}$ equipped with the standard Euclidean inner product $\langle \cdot, \cdot \rangle$ and the corresponding Euclidean norm $\|\cdot\|$. The cases with general preconditioned norms $\|\cdot\|_x$ and $\|\cdot\|_y$ can be proven similarly.

Function subfamilies. To establish the overall lower bound, we decompose the general function family into three distinct subfamilies:

Definition 2. Let $\mathcal{F}_x = \mathcal{F}(L_x, 0, 0, D_x, 0)$ denote the set of functions

$$F_x(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{A}_x \mathbf{x} - \mathbf{b}_x\|^2,$$

where the matrix $\mathbf{A}_x \in \mathbb{R}^{n_x \times n_x}$ satisfies $\|\mathbf{A}_x\|^2 \leq L_x$, and the vector $\mathbf{b}_x \in \mathbb{R}^{n_x}$ is such that the linear system $\mathbf{A}_x \mathbf{x} = \mathbf{b}_x$ has a solution $\mathbf{x}^* \in \mathcal{E}_x$ satisfying $\|\mathbf{x}^*\| \leq D_x$.

Definition 3. Let $\mathcal{F}_y = \mathcal{F}_y(0, 0, L_y, 0, D_y)$ denote the set of functions

$$F_y(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \|\mathbf{A}_y \mathbf{y} - \mathbf{b}_y\|^2,$$

where the matrix $\mathbf{A}_y \in \mathbb{R}^{n_y \times n_y}$ satisfies $\|\mathbf{A}_y\|^2 \leq L_y$, and the vector $\mathbf{b}_y \in \mathbb{R}^{n_y}$ is such that the linear system $\mathbf{A}_y \mathbf{y} = \mathbf{b}_y$ has a solution $\mathbf{y}^* \in \mathcal{E}_y$ satisfying $\|\mathbf{y}^*\| \leq D_y$.

Definition 4. Let $\mathcal{F}_{xy} = \mathcal{F}_{xy}(0, L_{xy}, 0, D_x, D_y)$ denote the set of functions

$$F_{xy}(\mathbf{x}, \mathbf{y}) = \langle \mathbf{A}_{xy} \mathbf{x} - \mathbf{b}_{xy}, \mathbf{y} \rangle,$$

where the matrix $\mathbf{A}_{xy} \in \mathbb{R}^{n_y \times n_x}$ satisfies $\|\mathbf{A}_{xy}\| \leq L_{xy}$, and the vector $\mathbf{b}_{xy} \in \mathbb{R}^{n_y}$ is such that the linear system $\mathbf{A}_{xy}\mathbf{x} = \mathbf{b}_{xy}$ has a solution $\mathbf{x}^* \in \mathcal{E}_x$ satisfying $\|\mathbf{x}^*\| \leq D_x$.

Distributed oracles. We consider the pairs of distributed oracles $(\mathcal{O}_x, \mathcal{O}_y)$ that return the partial gradients of the functions in these subfamilies:

- For $F_x \in \mathcal{F}_x$: $\mathcal{O}_x(\mathbf{x}, \mathbf{y}) = \nabla_x F_x(\mathbf{x}, \mathbf{y}) = \mathbf{A}_x^\top (\mathbf{A}_x \mathbf{x} - \mathbf{b}_x)$ and $\mathcal{O}_y(\mathbf{x}, \mathbf{y}) = -\nabla_y F_x(\mathbf{x}, \mathbf{y}) = \mathbf{0}$.
- For $F_y \in \mathcal{F}_y$: $\mathcal{O}_x(\mathbf{x}, \mathbf{y}) = \nabla_x F_y(\mathbf{x}, \mathbf{y}) = \mathbf{0}$ and $\mathcal{O}_y(\mathbf{x}, \mathbf{y}) = -\nabla_y F_y(\mathbf{x}, \mathbf{y}) = \mathbf{A}_y^\top (\mathbf{A}_y \mathbf{y} - \mathbf{b}_y)$.
- For $F_{xy} \in \mathcal{F}_{xy}$: $\mathcal{O}_x(\mathbf{x}, \mathbf{y}) = \nabla_x F_{xy}(\mathbf{x}, \mathbf{y}) = \mathbf{A}_{xy}^\top \mathbf{y}$ and $\mathcal{O}_y(\mathbf{x}, \mathbf{y}) = -\nabla_y F_{xy}(\mathbf{x}, \mathbf{y}) = \mathbf{b}_{xy} - \mathbf{A}_{xy} \mathbf{x}$.

Accuracy measures. For any approximate solution $\bar{\mathbf{z}} = (\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in \mathcal{E}_x \times \mathcal{E}_y$, the restricted primal-dual gaps over the bounded sets evaluate to:

- For $F_x \in \mathcal{F}_x$: $\Delta_{F_x}(\bar{\mathbf{z}}) = F_x(\bar{\mathbf{x}}, \bar{\mathbf{y}}) - \min_{\mathbf{x} \in \mathcal{B}_x} F_x(\mathbf{x}, \bar{\mathbf{y}}) = \frac{1}{2} \|\mathbf{A}_x \bar{\mathbf{x}} - \mathbf{b}_x\|^2$.
- For $F_y \in \mathcal{F}_y$: $\Delta_{F_y}(\bar{\mathbf{z}}) = \max_{\mathbf{y} \in \mathcal{B}_y} F_y(\bar{\mathbf{x}}, \mathbf{y}) - F_y(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \frac{1}{2} \|\mathbf{A}_y \bar{\mathbf{y}} - \mathbf{b}_y\|^2$.
- For $F_{xy} \in \mathcal{F}_{xy}$: $\Delta_{F_{xy}}(\bar{\mathbf{z}}) = \max_{\mathbf{y} \in \mathcal{B}_y} \langle \mathbf{A}_{xy} \bar{\mathbf{x}} - \mathbf{b}_{xy}, \mathbf{y} \rangle = D_y \|\mathbf{A}_{xy} \bar{\mathbf{x}} - \mathbf{b}_{xy}\|$.

Problem subclasses. These constructions yield three distinct problem subclasses:

- $\mathcal{P}_x = (\mathcal{F}_x, \mathcal{O}_x, \mathcal{O}_y, \epsilon)$,
- $\mathcal{P}_y = (\mathcal{F}_y, \mathcal{O}_x, \mathcal{O}_y, \epsilon)$, and
- $\mathcal{P}_{xy} = (\mathcal{F}_{xy}, \mathcal{O}_x, \mathcal{O}_y, \epsilon)$.

The worst-case complexity for general SPs is bounded below by the maximum complexity among these three subclasses.

Structure of this section. We present our lower bound proofs by establishing a connection from distributed SPs to convex minimization. In [Section B.1](#), we recall and properly rescale Nesterov's construction for unconstrained convex optimization. We analyze the three problem subclasses separately in [Sections B.2.1](#) and [B.2.2](#), mapping the lower bound for each subclass to this core convex minimization problem. Finally, we combine the results to establish the overall lower bound for distributed SPs in [Section B.3](#).

B.1 The “worst function in the world” (with proper rescaling)

While the original formulation (Theorem 2.1.7 by [Nesterov \[2004\]](#)) establishes a lower bound with respect to an arbitrary initial distance $\|\mathbf{v}^0 - \mathbf{v}^*\|$, this is insufficient for the lower bound proof in this paper. We are working on a lower bound analysis with a more specific problem subclass, where the algorithms may be designed with prior knowledge of a bound D on the distance to the optimum. To this end, we present a refined version of Nesterov's proposition by properly rescaling Nesterov's original construction to ensure that the distance to the optimum is bounded by D , while maintaining the identical lower bound on the function value.

For any given matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and vector $\mathbf{b} \in \mathbb{R}^m$, we define the k -th Krylov subspace as follows:

$$\mathcal{H}^k(\mathbf{A}, \mathbf{b}) \triangleq \text{span}\{\mathbf{A}^\top \mathbf{b}, (\mathbf{A}^\top \mathbf{A}) \mathbf{A}^\top \mathbf{b}, \dots, (\mathbf{A}^\top \mathbf{A})^{k-1} \mathbf{A}^\top \mathbf{b}\}.$$

Proposition 8. *Let $L > 0$, $D > 0$, and integer $1 \leq k \leq \frac{\min\{m, n\} - 1}{2}$. Then, there exist a matrix $\mathbf{A} = \mathbf{A}(L, k) \in \mathbb{R}^{m \times n}$ with $\|\mathbf{A}\| \leq L$ and a vector $\mathbf{b} = \mathbf{b}(L, D, k) \in \mathbb{R}^m$, such that*

$$\min_{\mathbf{v} \in \mathcal{H}^k(\mathbf{A}, \mathbf{b})} \frac{1}{2} \|\mathbf{A} \mathbf{v} - \mathbf{b}\|^2 \geq \frac{3L^2 D^2}{32(k+1)^2},$$

and the linear system $\mathbf{A} \mathbf{v} = \mathbf{b}$ has a solution $\mathbf{v}^* \in \mathbb{R}^n$ satisfying $\|\mathbf{v}^*\| \leq D$.

Proof. Let $p = 2k + 1$. By the condition $k \leq \frac{\min\{m,n\}-1}{2}$, we have $p \leq \min\{m,n\}$. Let $\mathbf{M}_p \in \mathbb{R}^{p \times p}$ be the symmetric tridiagonal matrix defined as:

$$\mathbf{M}_p = \begin{bmatrix} 2 & -1 & 0 & \cdots & 0 \\ -1 & 2 & -1 & \cdots & 0 \\ 0 & -1 & 2 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & -1 \\ 0 & 0 & \cdots & -1 & 2 \end{bmatrix}.$$

Let $\mathbf{B}_p \in \mathbb{R}^{p \times p}$ be the upper bidiagonal matrix such that $\mathbf{B}_p^\top \mathbf{B}_p = \mathbf{M}_p$, defined as:

$$\mathbf{B}_p = \begin{bmatrix} \sqrt{\frac{2}{1}} & -\sqrt{\frac{1}{2}} & 0 & \cdots & 0 \\ 0 & \sqrt{\frac{3}{2}} & -\sqrt{\frac{2}{3}} & \ddots & \vdots \\ 0 & 0 & \sqrt{\frac{4}{3}} & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & -\sqrt{\frac{p-1}{p}} \\ 0 & 0 & \cdots & 0 & \sqrt{\frac{p+1}{p}} \end{bmatrix}.$$

We define the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ as the block matrix:

$$\mathbf{A} = \frac{L}{2} \begin{bmatrix} \mathbf{B}_p & \mathbf{0}_{p \times (n-p)} \\ \mathbf{0}_{(m-p) \times p} & \mathbf{0}_{(m-p) \times (n-p)} \end{bmatrix}.$$

The matrix $\mathbf{A}^\top \mathbf{A} \in \mathbb{R}^{n \times n}$ is given by the block-diagonal matrix:

$$\mathbf{A}^\top \mathbf{A} = \frac{L^2}{4} \begin{bmatrix} \mathbf{M}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

The spectral norm of \mathbf{A} is bounded as $\|\mathbf{A}\| = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} = \frac{L}{2} \sqrt{\lambda_{\max}(\mathbf{M}_p)} \leq L$.

Let $\gamma = D \sqrt{\frac{6(p+1)}{p(2p+1)}}$. We define the vector $\mathbf{u} \in \mathbb{R}^p$ by its coordinates:

$$u_i = \frac{1}{\sqrt{i(i+1)}}, \quad \text{for } i = 1, \dots, p.$$

By the structure of \mathbf{B}_p^\top , we have $\mathbf{B}_p^\top \mathbf{u} = \mathbf{e}_1^{(p)} = (1, 0, \dots, 0)^\top \in \mathbb{R}^p$. We define $\mathbf{b} \in \mathbb{R}^m$ as the block vector:

$$\mathbf{b} = \gamma \frac{L}{2} \begin{bmatrix} \mathbf{u} \\ \mathbf{0}_{m-p} \end{bmatrix}.$$

The linear system $\mathbf{A} \mathbf{v} = \mathbf{b}$ has a solution $\mathbf{v}^* \in \mathbb{R}^n$ given by the block vector:

$$\mathbf{v}^* = \gamma \begin{bmatrix} \mathbf{M}_p^{-1} \mathbf{e}_1^{(p)} \\ \mathbf{0}_{n-p} \end{bmatrix}.$$

By the structure of \mathbf{M}_p^{-1} , the coordinates of \mathbf{v}^* are $v_i^* = \gamma \frac{p+1-i}{p+1}$ for $1 \leq i \leq p$, and 0 otherwise. Its squared norm evaluates to:

$$\|\mathbf{v}^*\|^2 = \frac{\gamma^2}{(p+1)^2} \sum_{j=1}^p j^2 = \gamma^2 \frac{p(2p+1)}{6(p+1)} = D^2.$$

The Krylov subspace is defined as:

$$\mathcal{H}^k(\mathbf{A}, \mathbf{b}) = \text{span}\{\mathbf{A}^\top \mathbf{b}, (\mathbf{A}^\top \mathbf{A}) \mathbf{A}^\top \mathbf{b}, \dots, (\mathbf{A}^\top \mathbf{A})^{k-1} \mathbf{A}^\top \mathbf{b}\}.$$

Since $\mathbf{B}_p^\top \mathbf{u} = \mathbf{e}_1^{(p)}$, the initial vector evaluates to $\mathbf{A}^\top \mathbf{b} = \gamma \frac{L^2}{4} \mathbf{e}_1^{(n)}$. Successive multiplication of $\mathbf{e}_1^{(n)}$ by $\mathbf{A}^\top \mathbf{A}$ expands the non-zero support by one standard basis vector at a time. Thus, the subspace spans the first k standard basis vectors:

$$\mathcal{H}^k(\mathbf{A}, \mathbf{b}) = \text{span}\{\mathbf{e}_1^{(n)}, \mathbf{e}_2^{(n)}, \dots, \mathbf{e}_k^{(n)}\}.$$

For any $\mathbf{v} \in \mathcal{H}^k(\mathbf{A}, \mathbf{b})$, its non-zero support is confined to the first k coordinates. Let $\mathbf{v}_k \in \mathbb{R}^k$ denote these first k coordinates, such that the first p coordinates of \mathbf{v} are $\mathbf{v}_p = (\mathbf{v}_k^\top, \mathbf{0}_{p-k}^\top)^\top$. Let $\mathbf{M}_k \in \mathbb{R}^{k \times k}$ be the leading principal submatrix of \mathbf{M}_p . The squared residual norm for $\mathbf{v} \in \mathcal{H}^k(\mathbf{A}, \mathbf{b})$ evaluates to:

$$\begin{aligned} \frac{1}{2} \|\mathbf{A}\mathbf{v} - \mathbf{b}\|^2 &= \frac{L^2}{8} \|\mathbf{B}_p \mathbf{v}_p - \gamma \mathbf{u}\|^2 \\ &= \frac{L^2}{8} \left(\langle \mathbf{v}_p, \mathbf{B}_p^\top \mathbf{B}_p \mathbf{v}_p \rangle - 2\gamma \langle \mathbf{v}_p, \mathbf{B}_p^\top \mathbf{u} \rangle + \gamma^2 \langle \mathbf{u}, \mathbf{u} \rangle \right) \\ &= \frac{L^2}{8} \left(\langle \mathbf{v}_k, \mathbf{M}_k \mathbf{v}_k \rangle - 2\gamma \langle \mathbf{v}_k, \mathbf{e}_1^{(k)} \rangle + \gamma^2 \sum_{i=1}^p \frac{1}{i(i+1)} \right). \end{aligned}$$

Using $\sum_{i=1}^p \frac{1}{i(i+1)} = 1 - \frac{1}{p+1} = \frac{p}{p+1}$, we have $\gamma^2 \|\mathbf{u}\|^2 = \gamma(\frac{p}{p+1}) = \gamma v_1^*$. Denoting $v_1 = \langle \mathbf{v}_k, \mathbf{e}_1^{(k)} \rangle$, the norm simplifies to:

$$\frac{1}{2} \|\mathbf{A}\mathbf{v} - \mathbf{b}\|^2 = \frac{L^2}{8} \left(\langle \mathbf{v}_k, \mathbf{M}_k \mathbf{v}_k \rangle - 2\gamma v_1 + \gamma v_1^* \right).$$

Minimizing this residual norm over $\mathbf{v}_k \in \mathbb{R}^k$ yields the optimal solution $\mathbf{v}_k^* = \gamma \mathbf{M}_k^{-1} \mathbf{e}_1^{(k)}$. By the structure of \mathbf{M}_k^{-1} , the coordinates of \mathbf{v}_k^* are $v_{k,i}^* = \gamma \frac{k+1-i}{k+1}$ for $1 \leq i \leq k$. The minimum value of the \mathbf{v}_k -dependent terms evaluates to:

$$\langle \mathbf{v}_k^*, \mathbf{M}_k \mathbf{v}_k^* \rangle - 2\gamma v_{k,1}^* = \gamma v_{k,1}^* - 2\gamma v_{k,1}^* = -\gamma v_{k,1}^* = -\gamma^2 \frac{k}{k+1}.$$

Substituting these values gives the minimum residual norm over the Krylov subspace:

$$\begin{aligned} \min_{\mathbf{v} \in \mathcal{H}^k(\mathbf{A}, \mathbf{b})} \frac{1}{2} \|\mathbf{A}\mathbf{v} - \mathbf{b}\|^2 &= \frac{L^2 \gamma^2}{8} \left(\frac{p}{p+1} - \frac{k}{k+1} \right) \\ &= \frac{L^2 \gamma^2}{8} \left(\frac{2k+1}{2k+2} - \frac{k}{k+1} \right) \\ &= \frac{L^2 \gamma^2}{16(k+1)} \\ &= \frac{L^2 D^2}{16(k+1)} \frac{6(2k+2)}{(2k+1)(4k+3)} \\ &= \frac{3L^2 D^2}{4(8k^2 + 10k + 3)}. \end{aligned}$$

Since $8k^2 + 10k + 3 \leq 8(k+1)^2$, this residual norm is lower bounded by $\frac{3L^2 D^2}{32(k+1)^2}$. \square

B.2 Three subclasses

B.2.1 Quadratic subclasses \mathcal{P}_x and \mathcal{P}_y

In this section, we apply [Proposition 8](#) to establish lower bounds for the quadratic subclasses \mathcal{P}_x and \mathcal{P}_y .

Notice that the gradient span sequences are naturally confined to the standard Krylov subspaces. Specifically, for $\mathbf{A}_x \in \mathbb{R}^{n_x \times n_x}$ and $\mathbf{b}_x \in \mathbb{R}^{n_x}$, the subspace satisfies the algebraic progression:

$$(\mathbf{A}_x^\top \mathbf{A}_x) \mathcal{H}^m(\mathbf{A}_x, \mathbf{b}_x) \subseteq \mathcal{H}^{m+1}(\mathbf{A}_x, \mathbf{b}_x) \quad \text{and} \quad \mathbf{A}_x^\top \mathbf{b}_x \in \mathcal{H}^{m+1}(\mathbf{A}_x, \mathbf{b}_x).$$

Proposition 9. *Let \mathcal{M} be a distributed gradient-span algorithm operating on an instance $F_x \in \mathcal{P}_x$ of the form $F_x(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}_x \mathbf{x} - \mathbf{b}_x\|^2$. Let $\bar{\mathbf{x}}$ be the approximate solution generated by Agent x after evaluating N_x partial gradients. Then, $\bar{\mathbf{x}} \in \mathcal{H}^{N_x}(\mathbf{A}_x, \mathbf{b}_x)$.*

Proof. For $F_x \in \mathcal{P}_x$, we have $\psi_x = \psi_y = 0$. Let $\mathbf{z}_1, \dots, \mathbf{z}_{N_x}$ be the sequence of query points evaluated by Agent x , where $\mathbf{z}_i = (\mathbf{x}_i, \mathbf{y}_i)$. Since $\mathbf{x}^0 = \mathbf{0}$, [Definition 1](#) requires that each query

point \mathbf{x}_{m+1} and the generated solution $\bar{\mathbf{x}}$ reside in the span of historical gradients. Let $S_m \triangleq \text{span}\{\nabla_x F_x(\mathbf{z}_i) \mid 1 \leq i \leq m\} \subseteq \mathcal{E}_x$. We have $\mathbf{x}_{m+1} \in S_m$ and $\bar{\mathbf{x}} \in S_{N_x}$.

We show $S_m \subseteq \mathcal{H}^m(\mathbf{A}_x, \mathbf{b}_x)$ by induction. The base case $m = 0$ holds since $S_0 = \{\mathbf{0}\} \subseteq \mathcal{H}^0(\mathbf{A}_x, \mathbf{b}_x)$.

Assume $S_m \subseteq \mathcal{H}^m(\mathbf{A}_x, \mathbf{b}_x)$ for some $m \geq 0$. The gradient at \mathbf{x}_{m+1} evaluates to $\nabla_x F_x(\mathbf{x}_{m+1}) = \mathbf{A}_x^\top \mathbf{A}_x \mathbf{x}_{m+1} - \mathbf{A}_x^\top \mathbf{b}_x$. Since $\mathbf{x}_{m+1} \in S_m \subseteq \mathcal{H}^m(\mathbf{A}_x, \mathbf{b}_x)$, applying the algebraic progression properties yields:

$$\nabla_x F_x(\mathbf{x}_{m+1}) \in (\mathbf{A}_x^\top \mathbf{A}_x) \mathcal{H}^m(\mathbf{A}_x, \mathbf{b}_x) - \mathbf{A}_x^\top \mathbf{b}_x \subseteq \mathcal{H}^{m+1}(\mathbf{A}_x, \mathbf{b}_x).$$

Thus, $S_{m+1} = S_m + \text{span}\{\nabla_x F_x(\mathbf{x}_{m+1})\} \subseteq \mathcal{H}^{m+1}(\mathbf{A}_x, \mathbf{b}_x)$.

By induction, $S_{N_x} \subseteq \mathcal{H}^{N_x}(\mathbf{A}_x, \mathbf{b}_x)$, which implies $\bar{\mathbf{x}} \in \mathcal{H}^{N_x}(\mathbf{A}_x, \mathbf{b}_x)$. \square

Analogously, for any instance $F_y \in \mathcal{P}_y$ evaluated by Agent y , of the form $F_y(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \|\mathbf{A}_y \mathbf{y} - \mathbf{b}_y\|^2$, the approximate solution generated after N_y queries satisfies $\bar{\mathbf{y}} \in \mathcal{H}^{N_y}(\mathbf{A}_y, \mathbf{b}_y)$.

Theorem 10. *Let \mathcal{M} be a distributed gradient-span algorithm for problem class \mathcal{P}_{SP} , where $n_x \geq 2\sqrt{\frac{3L_x D_x^2}{32\epsilon}} + 1$ and $n_y \geq 2\sqrt{\frac{3L_y D_y^2}{32\epsilon}} + 1$. Then, we have*

$$N_{\mathcal{P}_x}^{\mathcal{M}} \geq c_x \left(\sqrt{\frac{3L_x D_x^2}{32\epsilon}} - 1 \right),$$

and analogously,

$$N_{\mathcal{P}_y}^{\mathcal{M}} \geq c_y \left(\sqrt{\frac{3L_y D_y^2}{32\epsilon}} - 1 \right).$$

Proof. We prove the bound for \mathcal{P}_x . Let $K = \lfloor \sqrt{\frac{3L_x D_x^2}{32\epsilon}} \rfloor$. Since $n_x \geq 2K + 1$, we apply [Proposition 8](#) to obtain $\mathbf{A}_K = \mathbf{A}(\sqrt{L_x}, K)$ and $\mathbf{b}_K = \mathbf{b}(\sqrt{L_x}, D_x, K)$ constructed in $\mathbb{R}^{n_x \times n_x}$ and \mathbb{R}^{n_x} .

We set $\mathbf{A}_x = \mathbf{A}_K$, $\mathbf{b}_x = \mathbf{b}_K$, and define $F_K(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{A}_x \mathbf{x} - \mathbf{b}_x\|^2$. By [Proposition 8](#), $\|\mathbf{A}_x\|^2 \leq L_x$, and the linear system $\mathbf{A}_x \mathbf{x} = \mathbf{b}_x$ has a solution \mathbf{x}^* satisfying $\|\mathbf{x}^*\| \leq D_x$. Hence, $F_K \in \mathcal{P}_x$.

When \mathcal{M} is applied to P_K , suppose it generates an ϵ -saddle point $\bar{\mathbf{z}}$ utilizing N_x queries to \mathcal{O}_x . If $N_x \geq K$, the bound $N_x \geq \sqrt{\frac{3L_x D_x^2}{32\epsilon}} - 1$ holds. If $N_x < K$, [Proposition 9](#) implies $\bar{\mathbf{x}} \in \mathcal{H}^{N_x}(\mathbf{A}_x, \mathbf{b}_x)$.

Bounding the restricted duality gap via [Proposition 8](#) yields:

$$\epsilon \geq \Delta_{F_K}(\bar{\mathbf{z}}) = \frac{1}{2} \|\mathbf{A}_x \bar{\mathbf{x}} - \mathbf{b}_x\|^2 \geq \min_{\mathbf{v} \in \mathcal{H}^{N_x}(\mathbf{A}_x, \mathbf{b}_x)} \frac{1}{2} \|\mathbf{A}_x \mathbf{v} - \mathbf{b}_x\|^2 \geq \frac{3L_x D_x^2}{32(N_x + 1)^2}.$$

Rearranging gives $N_x \geq \sqrt{\frac{3L_x D_x^2}{32\epsilon}} - 1$. Multiplying by c_x yields $N_{\mathcal{P}_x}^{\mathcal{M}} \geq c_x \left(\sqrt{\frac{3L_x D_x^2}{32\epsilon}} - 1 \right)$.

The proof for \mathcal{P}_y is symmetric, setting $\mathbf{A}_y = \mathbf{A}(\sqrt{L_y}, K)$ and $\mathbf{b}_y = \mathbf{b}(\sqrt{L_y}, D_y, K)$ for $K = \lfloor \sqrt{\frac{3L_y D_y^2}{32\epsilon}} \rfloor$. \square

B.2.2 Bilinear subclass \mathcal{P}_{xy}

In this section, we establish the lower bounds for the bilinear subclass \mathcal{P}_{xy} .

Consider any function from \mathcal{P}_{xy} of the form $F(\mathbf{x}, \mathbf{y}) = \langle \mathbf{A}\mathbf{x} - \mathbf{b}, \mathbf{y} \rangle$, where $\mathbf{A} \in \mathbb{R}^{n_y \times n_x}$ and $\mathbf{b} \in \mathbb{R}^{n_y}$. The partial gradients are $\nabla_x F(\mathbf{x}, \mathbf{y}) = \mathbf{A}^\top \mathbf{y}$ and $\nabla_y F(\mathbf{x}, \mathbf{y}) = \mathbf{A}\mathbf{x} - \mathbf{b}$. For any $\bar{\mathbf{z}} = (\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in \mathcal{E}_x \times \mathcal{E}_y$, the restricted primal-dual gap evaluates to the residual norm:

$$\Delta_F(\bar{\mathbf{z}}) = \max_{\mathbf{y} \in \mathcal{B}_y} \langle \mathbf{A}\bar{\mathbf{x}} - \mathbf{b}, \mathbf{y} \rangle = D_y \|\mathbf{A}\bar{\mathbf{x}} - \mathbf{b}\|.$$

We define the coupled Krylov subspaces in \mathcal{E}_y and \mathcal{E}_x as:

$$\mathcal{H}_y^k(\mathbf{A}, \mathbf{b}) \triangleq \text{span}\{\mathbf{b}, (\mathbf{A}\mathbf{A}^\top)\mathbf{b}, \dots, (\mathbf{A}\mathbf{A}^\top)^{k-1}\mathbf{b}\},$$

$$\mathcal{H}_x^k(\mathbf{A}, \mathbf{b}) \triangleq \text{span}\{\mathbf{A}^\top \mathbf{b}, (\mathbf{A}^\top \mathbf{A}) \mathbf{A}^\top \mathbf{b}, \dots, (\mathbf{A}^\top \mathbf{A})^{k-1} \mathbf{A}^\top \mathbf{b}\} \equiv \mathcal{H}^k(\mathbf{A}, \mathbf{b}).$$

By definition, these subspaces satisfy the alternating properties:

$$\mathbf{A} \mathcal{H}_x^m(\mathbf{A}, \mathbf{b}) + \text{span}\{\mathbf{b}\} = \mathcal{H}_y^{m+1}(\mathbf{A}, \mathbf{b}) \quad \text{and} \quad \mathbf{A}^\top \mathcal{H}_y^m(\mathbf{A}, \mathbf{b}) = \mathcal{H}_x^m(\mathbf{A}, \mathbf{b}) \subseteq \mathcal{H}_x^{m+1}(\mathbf{A}, \mathbf{b}).$$

Proposition 11. *Let \mathcal{M} be a distributed gradient-span algorithm operating on $F(\mathbf{x}, \mathbf{y}) = \langle \mathbf{A}\mathbf{x} - \mathbf{b}, \mathbf{y} \rangle$. Let $\bar{\mathbf{x}}$ be the approximate solution generated after T communication rounds. Then, $\bar{\mathbf{x}} \in \mathcal{H}_x^{\lceil (T-1)/2 \rceil}(\mathbf{A}, \mathbf{b})$.*

Proof. Let $S_x^t \triangleq \text{span}\{\nabla_x F(\mathbf{z}) \mid \mathbf{z} \in Z_x^t\} \subseteq \mathcal{E}_x$ and $S_y^t \triangleq \text{span}\{-\nabla_y F(\mathbf{z}) \mid \mathbf{z} \in Z_y^t\} \subseteq \mathcal{E}_y$. We show by induction that $S_x^t \subseteq \mathcal{H}_x^{\lceil t/2 \rceil}(\mathbf{A}, \mathbf{b})$ and $S_y^t \subseteq \mathcal{H}_y^{\lfloor t/2 \rfloor + 1}(\mathbf{A}, \mathbf{b})$.

The base case $t = -1$ holds since $S_x^{-1} = \{\mathbf{0}\} \subseteq \mathcal{H}_x^0(\mathbf{A}, \mathbf{b})$ and $S_y^{-1} = \{\mathbf{0}\} \subseteq \mathcal{H}_y^0(\mathbf{A}, \mathbf{b})$.

Assume the claim holds for round $t - 1$. In round t , Agent x queries points using remote variables $\hat{\mathbf{y}} \in S_y^{t-1}$. Applying the alternating properties, the evaluated gradient satisfies:

$$\nabla_x F(\mathbf{x}, \hat{\mathbf{y}}) = \mathbf{A}^\top \hat{\mathbf{y}} \in \mathbf{A}^\top \mathcal{H}_y^{\lfloor (t-1)/2 \rfloor + 1}(\mathbf{A}, \mathbf{b}) = \mathcal{H}_x^{\lfloor (t-1)/2 \rfloor + 1}(\mathbf{A}, \mathbf{b}) = \mathcal{H}_x^{\lceil t/2 \rceil}(\mathbf{A}, \mathbf{b}).$$

Thus, $S_x^t \subseteq \mathcal{H}_x^{\lceil t/2 \rceil}(\mathbf{A}, \mathbf{b})$. Similarly, Agent y queries points using remote variables $\hat{\mathbf{x}} \in S_x^{t-1}$. The evaluated gradient satisfies:

$$-\nabla_y F(\hat{\mathbf{x}}, \mathbf{y}) = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}} \in \text{span}\{\mathbf{b}\} + \mathbf{A} \mathcal{H}_x^{\lceil (t-1)/2 \rceil}(\mathbf{A}, \mathbf{b}) = \mathcal{H}_y^{\lceil (t-1)/2 \rceil + 1}(\mathbf{A}, \mathbf{b}) = \mathcal{H}_y^{\lfloor t/2 \rfloor + 1}(\mathbf{A}, \mathbf{b}).$$

Thus, $S_y^t \subseteq \mathcal{H}_y^{\lfloor t/2 \rfloor + 1}(\mathbf{A}, \mathbf{b})$. By induction, the claim holds for all t . The generated approximate solution satisfies $\bar{\mathbf{x}} \in S_x^{T-1} \subseteq \mathcal{H}_x^{\lceil (T-1)/2 \rceil}(\mathbf{A}, \mathbf{b})$. \square

Theorem 12. *Let \mathcal{M} be a distributed gradient-span algorithm for problem class \mathcal{P}_{xy} , where $\min\{n_x, n_y\} \geq \frac{2L_{xy}D_xD_y}{3\epsilon} + 1$. Then, we have*

$$T_{\mathcal{P}_{xy}}^{\mathcal{M}} \geq \frac{2L_{xy}D_xD_y}{3\epsilon} - 3.$$

Proof. Let $K = \lfloor \frac{L_{xy}D_xD_y}{3\epsilon} \rfloor$. Since $n_x, n_y \geq 2K + 1$, we apply [Proposition 8](#) to obtain $\mathbf{A}_K = \mathbf{A}(L_{xy}, K) \in \mathbb{R}^{n_y \times n_x}$ and $\mathbf{b}_K = \mathbf{b}(L_{xy}, D_x, K) \in \mathbb{R}^{n_y}$.

We set $\mathbf{A}_{xy} = \mathbf{A}_K$ and $\mathbf{b}_{xy} = \mathbf{b}_K$, and define $F_K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{A}_{xy}\mathbf{x} - \mathbf{b}_{xy}, \mathbf{y} \rangle$. By [Proposition 8](#), $\|\mathbf{A}_{xy}\| = \|\mathbf{A}_K\| \leq L_{xy}$. The linear system $\mathbf{A}_{xy}\mathbf{x} = \mathbf{b}_{xy}$ has a solution \mathbf{x}^* satisfying $\|\mathbf{x}^*\| \leq D_x$. Thus, $F_K \in \mathcal{P}_{xy}$.

When \mathcal{M} is applied to P_K , suppose it generates an ϵ -saddle point $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ after T communication rounds. If $T \geq 2K$, then $T \geq \frac{2L_{xy}D_xD_y}{3\epsilon} - 3$ holds. If $T < 2K$, [Proposition 11](#) implies $\bar{\mathbf{x}} \in \mathcal{H}_x^k(\mathbf{A}_{xy}, \mathbf{b}_{xy})$, where $k = \lceil (T-1)/2 \rceil$.

Since $\mathcal{H}_x^k(\mathbf{A}_{xy}, \mathbf{b}_{xy}) \equiv \mathcal{H}^k(\mathbf{A}_K, \mathbf{b}_K)$, the restricted duality gap on P_K evaluates to:

$$\Delta_{F_K}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = D_y \|\mathbf{A}_K \bar{\mathbf{x}} - \mathbf{b}_K\|.$$

Applying [Proposition 8](#) yields:

$$\epsilon \geq \Delta_{F_K}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \geq D_y \min_{\mathbf{v} \in \mathcal{H}^k(\mathbf{A}_K, \mathbf{b}_K)} \|\mathbf{A}_K \mathbf{v} - \mathbf{b}_K\| \geq D_y \sqrt{\frac{3L_{xy}^2 D_x^2}{2(8k^2 + 10k + 3)}} \geq \frac{L_{xy} D_x D_y}{3(k+1)}.$$

Rearranging gives $k \geq \frac{L_{xy} D_x D_y}{3\epsilon} - 1$. Because $k = \lceil (T-1)/2 \rceil$, we have $T \geq 2k - 1 \geq \frac{2L_{xy} D_x D_y}{3\epsilon} - 3$. \square

The oracle lower bounds can be proved using identical subspace confinement arguments. Rather than repeating the proof, we provide the main proposition and theorem here.

Proposition 13. *Let \mathcal{M} be a distributed gradient-span algorithm operating on $F(\mathbf{x}, \mathbf{y}) = \langle \mathbf{A}\mathbf{x} - \mathbf{b}, \mathbf{y} \rangle$. Let $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ be the solution generated after N_x and N_y queries to \mathcal{O}_x and \mathcal{O}_y , respectively. Then, $\bar{\mathbf{x}} \in \mathcal{H}_x^k(\mathbf{A}, \mathbf{b})$, where $k = \min(N_x, N_y)$.*

Theorem 14. Let \mathcal{M} be a distributed gradient-span algorithm for problem class \mathcal{P}_{xy} , where $\min\{n_x, n_y\} \geq \frac{2L_{xy}D_xD_y}{3\epsilon} + 1$. Then, we have

$$N_{\mathcal{P}_{xy}}^{\mathcal{M}} \geq (c_x + c_y) \left(\frac{L_{xy}D_xD_y}{3\epsilon} - 1 \right).$$

B.3 Proof of Theorem 1

Finally, we derive the lower bound for \mathcal{P}_{SP} by assembling the lower bounds obtained from the three subclasses.

Proof of Theorem 1. The general class of convex-concave saddle point problems \mathcal{P}_{SP} contains the three unregularized subclasses constructed in the previous sections: the x-quadratic subclass \mathcal{P}_x , the y-quadratic subclass \mathcal{P}_y , and the bilinear subclass \mathcal{P}_{xy} . The worst-case complexity for an algorithm operating over the entire class \mathcal{P}_{SP} is bounded from below by the maximum of the complexities required for these individual subclasses.

By Theorem 12, the communication complexity over the class is bounded by the communication complexity of the bilinear subclass:

$$T_{\mathcal{P}_{\text{SP}}}^{\mathcal{M}} \geq T_{\mathcal{P}_{xy}}^{\mathcal{M}} \geq \frac{2L_{xy}D_xD_y}{3\epsilon} - 3.$$

For the computational complexity, we combine the independent lower bounds established in Theorem 10 and Theorem 12. The total computational complexity is bounded by the maximum of the three individual requirements:

$$N_{\mathcal{P}_{\text{SP}}}^{\mathcal{M}} \geq \max\{N_{\mathcal{P}_{xy}}^{\mathcal{M}}, N_{\mathcal{P}_x}^{\mathcal{M}}, N_{\mathcal{P}_y}^{\mathcal{M}}\}.$$

Using the algebraic property $\max\{a, b, c\} \geq \frac{1}{3}(a + b + c)$, we obtain the lower bound:

$$\begin{aligned} N_{\mathcal{P}_{\text{SP}}}^{\mathcal{M}} &\geq \frac{1}{3} \left[N_{\mathcal{P}_{xy}}^{\mathcal{M}} + N_{\mathcal{P}_x}^{\mathcal{M}} + N_{\mathcal{P}_y}^{\mathcal{M}} \right] \\ &\geq \frac{1}{3} \left[(c_x + c_y) \left(\frac{L_{xy}D_xD_y}{3\epsilon} - 1 \right) + c_x \left(\sqrt{\frac{3L_xD_x^2}{32\epsilon}} - 1 \right) + c_y \left(\sqrt{\frac{3L_yD_y^2}{32\epsilon}} - 1 \right) \right] \\ &= \frac{c_x + c_y}{9} \frac{L_{xy}D_xD_y}{\epsilon} + \frac{c_x}{3} \sqrt{\frac{3L_xD_x^2}{32\epsilon}} + \frac{c_y}{3} \sqrt{\frac{3L_yD_y^2}{32\epsilon}} - \frac{2c_x + 2c_y}{3}. \end{aligned}$$

This establishes the stated lower bounds for the general problem class and concludes the proof. \square

C Multi-stage reduction for distributed SPs

We first introduce the general notion of composite variational inequality problem (VIP) in Section C.1, which is the backbone of our problems. Then, in Sections C.2 to C.4, we introduce the subsequent technical components of our algorithm. Unless otherwise specified, in this section, we work with (general) finite-dimensional real vector space \mathcal{E} equipped with the norm $\|\mathbf{z}\|_{\mathcal{E}} = \langle \mathbf{P}\mathbf{z}, \mathbf{z} \rangle^{\frac{1}{2}}$, where $\mathbf{P}: \mathcal{E} \rightarrow \mathcal{E}^*$ is a self-adjoint positive definite linear operator.

C.1 Preliminary: Variational inequality problems

For any operator $V(\cdot): \text{dom } \psi \rightarrow \mathcal{E}^*$ and any function $\psi(\cdot): \mathcal{E} \rightarrow \mathbb{R} \cup \{+\infty\}$, we say that $\mathbf{z}^* \in \mathcal{E}$ is a *solution* of the VIP of (V, ψ) if

$$\langle V(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle + \psi(\mathbf{z}) \geq \psi(\mathbf{z}^*), \text{ for all } \mathbf{z} \in \text{dom } \psi. \quad (3)$$

Assumption for VIPs. Let us introduce the following assumption:

(A1') The function ψ is a (simple) proper closed convex function. The operator V is continuous and monotone over $\text{dom } \psi$: that is, $\langle V(\mathbf{z}') - V(\mathbf{z}), \mathbf{z}' - \mathbf{z} \rangle \geq 0$, for all $\mathbf{z}', \mathbf{z} \in \text{dom } \psi$.

The solution of VIP defined in Eq. (3) is often called a “strong solution”. There is another notion of “weak solution”, but under (A1’), these two notions coincide. Meanwhile, under (A1’), a point $\mathbf{z}^* \in \mathcal{E}$ is the solution of (3) if and only if $\mathbf{0} \in V(\mathbf{z}^*) + \partial\psi(\mathbf{z}^*)$.

Indeed, for the saddle problem given by (f, ψ_x, ψ_y) , let us consider the VIP given by (V^f, ψ_z) , where

$$V^f(\mathbf{z}) \equiv (\nabla_x f(\mathbf{z}), -\nabla_y f(\mathbf{z})) \text{ and } \psi_z(\mathbf{z}) \equiv \psi_x(\mathbf{x}) + \psi_y(\mathbf{y}), \text{ for all } \mathbf{z} \in Q. \quad (4)$$

Moreover, for the saddle problem under (A1), the above VIP satisfies (A1’), and therefore, the solution of the VIP coincides with the saddle point.

C.2 Reduced-operator method for VIPs

Now, we introduce the Reduced-Operator Method (ROM) recently proposed in [Nesterov, 2023]. In particular, we apply this general framework in a special way so as to reduce the VIP to a sequence of Monteiro-Svaiter Subproblems (MSSs) [Monteiro and Svaiter, 2013]. The MSS asks to find a point for the regularized function such that the subgradient norm at this point is small compared to the distance from the initial point.

Monteiro-Svaiter Subproblem. Given an operator $V: \text{dom } \psi \rightarrow \mathcal{E}^*$, a function $\psi: \mathcal{E} \rightarrow \mathbb{R} \cup \{+\infty\}$, a reference point $\mathbf{v} \in \text{dom } \psi$, and a real number $\lambda > 0$, we say $(\mathbf{z}^+, \psi'(\mathbf{z}^+))$ is a *solution* of the MSS if $\mathbf{z}^+ \in \text{dom } \psi$, $\psi'(\mathbf{z}^+) \in \partial\psi(\mathbf{z}^+)$, and

$$\|V(\mathbf{z}^+) + \psi'(\mathbf{z}^+) + \lambda\mathbf{P}(\mathbf{z}^+ - \mathbf{v})\|_{\mathcal{E}^*} \leq \lambda\|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}. \quad (5)$$

We will discuss how to solve the MSSs later in Section C.3. But for now, let us assume there exists a solver

$$\mathcal{M}^{\text{MS}}(V, \psi, v, \lambda)$$

for the MSSs, which takes an MSS given by (V, ψ, v, λ) and returns a solution of it. Built upon such a solver \mathcal{M}^{MS} , we now introduce ROM in Algorithm 2. At each iteration t : the solver \mathcal{M}^{MS} returns a solution $(\mathbf{z}^{t+1}, \psi'(\mathbf{z}^{t+1}))$ for the MSS built at reference point \mathbf{v}^t ; this solution is used as a midpoint to compute subgradient $V_\psi(\mathbf{z}^{t+1})$; then the ‘extragradient-like’ step is taken with the stepsize a_{t+1} .

Algorithm 2 ROM $_{\|\cdot\|_{\mathcal{E}}}(V, \psi, \mathbf{z}^0, (\lambda_t)_{t \geq 1} \mid \mathcal{M}^{\text{MS}})$

Require: A solver \mathcal{M}^{MS} for the MSSs.

- 1: $\mathbf{v}^0 = \mathbf{z}^0$.
 - 2: **for** $t = 0, 1, \dots$ **do**
 - 3: $(\mathbf{z}^{t+1}, \psi'(\mathbf{z}^{t+1})) = \mathcal{M}^{\text{MS}}(V, \psi, \mathbf{v}^t, \lambda_{t+1})$.
 - 4: $V_\psi(\mathbf{z}^{t+1}) = V(\mathbf{z}^{t+1}) + \psi'(\mathbf{z}^{t+1})$.
 - 5: $a_{t+1} = \frac{2\langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v}^t - \mathbf{z}^{t+1} \rangle}{\|V_\psi(\mathbf{z}^{t+1})\|_{\mathcal{E}^*}^2}$.
 - 6: $\mathbf{v}^{t+1} = \arg \min_{\mathbf{v} \in \text{dom } \psi} [a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v} \rangle + \frac{1}{2} \|\mathbf{v} - \mathbf{v}^t\|_{\mathcal{E}}^2]$.
 - 7: **end for**
-

Next, let us show the convergence of ROM.

Lemma 15. ROM (Algorithm 2) ensures for all $\mathbf{v} \in \text{dom } \psi$ and for all $T \geq 1$,

$$\sum_{t=0}^{T-1} a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{z}^{t+1} - \mathbf{v} \rangle \leq \frac{1}{2} \|\mathbf{v}^0 - \mathbf{v}\|_{\mathcal{E}}^2 - \frac{1}{2} \|\mathbf{v}^T - \mathbf{v}\|_{\mathcal{E}}^2.$$

Moreover, we have $a_{t+1} \geq \frac{1}{\lambda_{t+1}}$, for all $t \geq 0$.

Proof. By the optimality of \mathbf{v}^{t+1} , we have for all $\mathbf{v} \in \text{dom } \psi$,

$$a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v} - \mathbf{v}^{t+1} \rangle + \frac{1}{2} \|\mathbf{v}^t - \mathbf{v}\|_{\mathcal{E}}^2 \geq \frac{1}{2} \|\mathbf{v}^{t+1} - \mathbf{v}\|_{\mathcal{E}}^2 + \frac{1}{2} \|\mathbf{v}^t - \mathbf{v}^{t+1}\|_{\mathcal{E}}^2,$$

and therefore,

$$\begin{aligned}
& a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v} - \mathbf{z}^{t+1} \rangle + \frac{1}{2} \|\mathbf{v}^t - \mathbf{v}\|_{\mathcal{E}}^2 \\
& \geq a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v}^{t+1} - \mathbf{z}^{t+1} \rangle + \frac{1}{2} \|\mathbf{v}^{t+1} - \mathbf{v}\|_{\mathcal{E}}^2 + \frac{1}{2} \|\mathbf{v}^t - \mathbf{v}^{t+1}\|_{\mathcal{E}}^2 \\
& = a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v}^t - \mathbf{z}^{t+1} \rangle + \frac{1}{2} \|\mathbf{v}^{t+1} - \mathbf{v}\|_{\mathcal{E}}^2 + a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v}^{t+1} - \mathbf{v}^t \rangle + \frac{1}{2} \|\mathbf{v}^t - \mathbf{v}^{t+1}\|_{\mathcal{E}}^2 \\
& \geq a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v}^t - \mathbf{z}^{t+1} \rangle + \frac{1}{2} \|\mathbf{v}^{t+1} - \mathbf{v}\|_{\mathcal{E}}^2 - \frac{a_{t+1}^2}{2} \|V_\psi(\mathbf{z}^{t+1})\|_{\mathcal{E}^*}^2 \\
& \geq \frac{1}{2} \|\mathbf{v}^{t+1} - \mathbf{v}\|_{\mathcal{E}}^2,
\end{aligned}$$

where the last inequality follows from the assignment of a_{t+1} in Line 5 of Algorithm 2. Then, the desired bound follows from summing the above inequality over t from 0 to $T - 1$.

Next, we show the lower bound for a_t . For all $t \geq 1$, we have

$$\begin{aligned}
& \langle V_\psi(\mathbf{z}^t), \mathbf{v}^{t-1} - \mathbf{z}^t \rangle - \frac{1}{2\lambda_t} \|V_\psi(\mathbf{z}^t)\|_{\mathcal{E}^*}^2 \\
& \equiv \frac{\lambda_t}{2} \|\mathbf{z}^t - \mathbf{v}^{t-1}\|_{\mathcal{E}}^2 - \frac{1}{2\lambda_t} \|V_\psi(\mathbf{z}^t) + \lambda_t \mathbf{P}(\mathbf{z}^t - \mathbf{v}^{t-1})\|_{\mathcal{E}^*}^2 \\
& \geq 0,
\end{aligned}$$

where the last inequality follows from Eq. (5). Therefore, we have

$$a_t = \frac{2 \langle V_\psi(\mathbf{z}^t), \mathbf{v}^t - \mathbf{z}^t \rangle}{\|V_\psi(\mathbf{z}^t)\|_{\mathcal{E}^*}^2} \geq \frac{1}{\lambda_t}.$$

□

C.3 Fully decoupled solver for MSSs with weak couplings

We now address the MSSs introduced in Section C.2. Specifically, we focus on the subproblems arising from applying ROM to the saddle problem of (f, ψ_x, ψ_y) . Therefore, we consider the MSS given by

$$(V^f, \psi_z, \mathbf{v}, \lambda), \quad (6)$$

where V^f and ψ_z are defined in Eq. (4), the reference point $\mathbf{v} = (\mathbf{v}_x, \mathbf{v}_y) \in \text{dom } \psi_x \times \text{dom } \psi_y$, and the assembled norm $\|\cdot\|_{\mathcal{E}}$ is associated with parameters α_x and α_y .

We say that an MSS has a *weak coupling* if $\lambda \geq \boxed{2\bar{L}_c \triangleq \frac{2L_{xy}}{\sqrt{\alpha_x \alpha_y}}}$. In this section, we will introduce

a Fully Decoupled Solver (FDS), which reduces the weakly-coupled MSSs to coordinate-wise Minimization of Residual Norms (MRNs).

Let us first define the problem of MRN. The MRN asks to find an approximate solution \mathbf{w}^+ of the VIP of (V_w, ψ_w) with the accuracy specified as follows:

Minimization of residual norm. Given an operator $V_w : \text{dom } \psi_w \rightarrow \mathcal{E}_w^*$, a function $\psi_w : \mathcal{E}_w \rightarrow \mathbb{R} \cup \{+\infty\}$, a reference point $\mathbf{v}_w \in \text{dom } \psi_w$, and an accuracy $\delta > 0$, we say $(\mathbf{w}^+, \psi_w'(\mathbf{w}^+))$ minimizes the residual norm to δ -relative distance accuracy, if $\mathbf{w}^+ \in \text{dom } \psi_w$, $\psi_w'(\mathbf{w}^+) \in \partial\psi_w(\mathbf{w}^+)$, and

$$\|V_w(\mathbf{w}^+) + \psi_w'(\mathbf{w}^+)\|_{w^*} \leq \delta \|\mathbf{w}^+ - \mathbf{v}_w\|_w.$$

Let us, again, defer the discussion of solving MRNs to Section C.4. But for now, let us assume there exist solvers

$$\mathcal{M}_x^{\text{MRN}}(V_x, \hat{\psi}_x, \mathbf{v}_x, \delta_x) \quad \text{and} \quad \mathcal{M}_y^{\text{MRN}}(V_y, \hat{\psi}_y, \mathbf{v}_y, \delta_y)$$

for the coordinate-wise MRNs in spaces \mathcal{E}_x and \mathcal{E}_y . In particular, these solvers take a coordinate-wise MRN problem and returns a point and a subgradient satisfying the desired accuracy.

Crucially, a key result in our proof is that, for an MSS with a weak coupling, we apply [Algorithm 3](#)—Fully Decoupled Solver (FDS)—that solves for each decision variable separately. We show in [Lemma 16](#) that FDS returns the correct solution *in one round*.

Algorithm 3 $\text{FDS}_{\|\cdot\|_\varepsilon}((\nabla_x f, -\nabla_y f), (\psi_x, \psi_y), \mathbf{v}, \lambda \mid (\mathcal{M}_x^{\text{MRN}}, \mathcal{M}_y^{\text{MRN}}))$

Require: Solvers $\mathcal{M}_x^{\text{MRN}}$ and $\mathcal{M}_y^{\text{MRN}}$ for the coordinate-wise MRNs.

- 1: $\hat{\psi}_x = \psi_x + \frac{\alpha_x \lambda}{2} \|\cdot - \mathbf{v}_x\|_x^2$ and $\hat{\psi}_y = \psi_y + \frac{\alpha_y \lambda}{2} \|\cdot - \mathbf{v}_y\|_y^2$.
- 2: **Agent** x and **Agent** y respectively compute

$$\begin{aligned} (\mathbf{x}^+, \hat{\psi}'_x(\mathbf{x}^+)) &= \mathcal{M}_x^{\text{MRN}}(\nabla_x f(\cdot, \mathbf{v}_y), \hat{\psi}_x, \mathbf{v}_x, \delta_x) \text{ and} \\ (\mathbf{y}^+, \hat{\psi}'_y(\mathbf{y}^+)) &= \mathcal{M}_y^{\text{MRN}}(-\nabla_y f(\mathbf{v}_x, \cdot), \hat{\psi}_y, \mathbf{v}_y, \delta_y), \end{aligned}$$

where $\delta_x = \frac{\alpha_x \lambda}{2}$ and $\delta_y = \frac{\alpha_y \lambda}{2}$.

- 3: $\psi'_x(\mathbf{x}^+) = \hat{\psi}'_x(\mathbf{x}^+) - \alpha_x \lambda \mathbf{P}_x(\mathbf{x}^+ - \mathbf{v}_x)$ and $\psi'_y(\mathbf{y}^+) = \hat{\psi}'_y(\mathbf{y}^+) - \alpha_y \lambda \mathbf{P}_y(\mathbf{y}^+ - \mathbf{v}_y)$.
 - 4: **return** $(\mathbf{z}^+, \psi'(\mathbf{z}^+))$, where $\mathbf{z}^+ = (\mathbf{x}^+, \mathbf{y}^+)$ and $\psi'(\mathbf{z}^+) = (\psi'_x(\mathbf{x}^+), \psi'_y(\mathbf{y}^+))$.
-

Lemma 16. *Consider the saddle problem given by (f, ψ_x, ψ_y) which satisfies [\(A3\)](#). For $\lambda \geq 2\bar{L}_c$, FDS ([Algorithm 3](#)) returns a solution of the MSS introduced in [Eq. \(6\)](#).*

The correctness of the FDS for SPs can be implied as a direct consequence of the correctness of FDS for the (more general) VIPs proven in [Section E](#). We will get to this result later in [Lemma 23](#).

C.4 Minimization of residual norms

We arrive at the last building component, the Minimization of Residual Norms (MRNs).

Assumptions for MRNs. Let us introduce the following assumptions:

- ($\hat{A}1$) The function ψ_w is a (simple) proper closed convex function. The operator V_w is monotone over $\text{dom } \psi_w$.
- ($\hat{A}2$) The set-valued operator $V_w + \partial\psi_w$ is μ -strongly maximally monotone over $\text{dom } \psi_w$. That is,
$$\langle V_w(\mathbf{w}') + \psi'_w(\mathbf{w}') - V_w(\mathbf{w}) - \psi'_w(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle \geq \mu \|\mathbf{w}' - \mathbf{w}\|_w^2,$$
for all $\mathbf{w}', \mathbf{w} \in \text{dom } \psi_w$, $\psi'_w(\mathbf{w}') \in \partial\psi_w(\mathbf{w}')$, and $\psi'_w(\mathbf{w}) \in \partial\psi_w(\mathbf{w})$,
- ($\hat{A}3$) The operator $V_w(\mathbf{w})$ is L -Lipschitz continuous over $\mathbf{w} \in \text{dom } \psi_w$.
- ($\hat{A}4$) The operator $V_w = \nabla f_w$, where f_w is a continuously differentiable function defined on an open set containing $\text{dom } \psi_w$.

The theoretical guarantee provided in the literature is usually based on the distance-to-solution accuracy (cf. [Definition 5](#)). We show in [Lemma 17](#) that, under strong maximal monotonicity, the relative distance accuracy required in this paper can be implied from distance-to-solution accuracy.

Definition 5 (Distance-to-solution accuracy). We say that $(\mathbf{w}^+, \psi'_w(\mathbf{w}^+))$ satisfies ξ -distance-to-solution accuracy if $\mathbf{w}^+ \in \text{dom } \psi_w$, $\psi'_w(\mathbf{w}^+) \in \partial\psi_w(\mathbf{w}^+)$, and $\|V_w(\mathbf{w}^+) + \psi'_w(\mathbf{w}^+)\|_{w^*} \leq \xi \|\mathbf{v}_w - \tilde{\mathbf{w}}\|_w$ for some $\tilde{\mathbf{w}}$ in the solution set of the VIP of (V_w, ψ_w) .

Lemma 17. *Consider the MRN problem given by $(V_w, \psi_w, \mathbf{v}_w, \delta)$ which satisfies [\(\$\hat{A}2\$ \)](#). Let $\xi \leq \frac{\mu\delta}{\mu+\delta}$. If $(\mathbf{w}^+, \psi'_w(\mathbf{w}^+))$ satisfies ξ -distance-to-solution accuracy, then $(\mathbf{w}^+, \psi'_w(\mathbf{w}^+))$ is a solution of the MRN problem.*

Proof. In view of the triangle inequality and the μ -strong maximal monotonicity of $V_w + \psi_w$, we have

$$\begin{aligned} \|\mathbf{v}_w - \tilde{\mathbf{w}}\|_w &\leq \|\mathbf{w}^+ - \mathbf{v}_w\|_w + \|\mathbf{w}^+ - \tilde{\mathbf{w}}\|_w \\ &\leq \|\mathbf{w}^+ - \mathbf{v}_w\|_w + \frac{1}{\mu} \|V_w(\mathbf{w}^+) + \psi'_w(\mathbf{w}^+)\|_{w^*} \\ &\leq \|\mathbf{w}^+ - \mathbf{v}_w\|_w + \frac{\xi}{\mu} \|\mathbf{v}_w - \tilde{\mathbf{w}}\|_w. \end{aligned}$$

Then, we have

$$\|\mathbf{v}_w - \tilde{\mathbf{w}}\|_w \leq \frac{\mu}{\mu - \xi} \|\mathbf{w}^+ - \mathbf{v}_w\|_w.$$

Therefore, we have

$$\|V_w(\mathbf{w}^+) + \psi'_w(\mathbf{w}^+)\|_{w^*} \leq \xi \|\mathbf{v}_w - \tilde{\mathbf{w}}\|_w \leq \frac{\mu\xi}{\mu - \xi} \|\mathbf{w}^+ - \mathbf{v}_w\|_w \leq \delta,$$

where the last inequality follows from the assignment $\xi \leq \frac{\mu\delta}{\mu + \delta}$. \square

We will leverage efficient existing solvers for the MRN problems. In particular, we are to deal with the specific MRNs given in Line 2 in Algorithm 3, where the corresponding coordinate-wise operators are gradients of smooth convex functions. Therefore, we can leverage the existing accelerated gradient methods from the literature.

Let us apply, for instance, the Accumulative Regularization Method (ARM) from Lan et al. [2023], whose detailed pseudocode is presented in Algorithm 4 for self-consistency. Let us denote this algorithm as

$$\text{ARM}(\nabla f_w, \psi_w, \mathbf{v}_w, \xi \mid L), \quad (7)$$

which takes an MRN instance of interest, has knowledge of the parameter L in ($\hat{\text{A3}}$), and returns a solution that satisfies ξ -distance-to-solution accuracy.

Algorithm 4 $\text{ARM}(\nabla f_w, \psi_w, \mathbf{v}_w, \xi \mid L)$

- 1: Set $\tau = 2 + \max\{0, \lceil \log_4(\frac{3L}{2\xi}) \rceil\}$ and $\sigma^{(0)} = 0$.
 - 2: Set $\sigma^{(k)} = 4^{k-3} \frac{2\xi}{3}$ and $N_k = \lceil 16 \sqrt{\frac{L}{\sigma^{(k)}}} \rceil$ for $k = 1, \dots, \tau$.
 - 3: Initialize $\bar{\mathbf{w}}^{(0)} = \mathbf{w}^{(0)} = \mathbf{v}_w$.
 - 4: **for** $k = 1, \dots, \tau$ **do**
 - 5: $\gamma^{(k)} = 1 - \frac{\sigma^{(k-1)}}{\sigma^{(k)}}$
 - 6: $\bar{\mathbf{w}}^{(k)} = (1 - \gamma^{(k)})\bar{\mathbf{w}}^{(k-1)} + \gamma^{(k)}\mathbf{w}^{(k-1)}$
 - 7: *% Begin Inner Subroutine: Nesterov's Accelerated Gradient Method*
 - 8: Initialize $\mathbf{x}_k^{(0)} = \mathbf{w}^{(k-1)}$, $\mathbf{y}_k^{(0)} = \mathbf{w}^{(k-1)}$, $t_0 = 1$, and $L_k = L + \sigma^{(k)}$.
 - 9: **for** $i = 0, \dots, N_k - 1$ **do**
 - 10: $\nabla f_w^{(k)}(\mathbf{y}_k^{(i)}) = \nabla f_w(\mathbf{y}_k^{(i)}) + \sigma^{(k)}(\mathbf{y}_k^{(i)} - \bar{\mathbf{w}}^{(k)})$
 - 11: $\mathbf{x}_k^{(i+1)} = \arg \min_{\mathbf{w} \in w} \{\langle \nabla f_w^{(k)}(\mathbf{y}_k^{(i)}), \mathbf{w} \rangle + \frac{L_k}{2} \|\mathbf{w} - \mathbf{y}_k^{(i)}\|_w^2 + \psi_w(\mathbf{w})\}$
 - 12: $t_{i+1} = \frac{1 + \sqrt{1 + 4t_i^2}}{2}$
 - 13: $\mathbf{y}_k^{(i+1)} = \mathbf{x}_k^{(i+1)} + \frac{t_i - 1}{t_{i+1}}(\mathbf{x}_k^{(i+1)} - \mathbf{x}_k^{(i)})$
 - 14: **end for**
 - 15: *% End Inner Subroutine*
 - 16: $\mathbf{w}^{(k)} = \mathbf{x}_k^{(N_k)}$
 - 17: $\psi'_w(\mathbf{w}^{(k)}) = -\nabla f_w^{(k)}(\mathbf{y}_k^{(N_k-1)}) - L_k(\mathbf{x}_k^{(N_k)} - \mathbf{y}_k^{(N_k-1)})$
 - 18: **end for**
 - 19: **return** $(\mathbf{w}^{(\tau)}, \psi'_w(\mathbf{w}^{(\tau)}))$
-

Now, we state the gradient query complexity with distance-to-solution accuracy. The original result of Lan et al. [2023] is given in projected gradient norm, which can be converted to the subgradient norm considered in this paper.

Lemma 18 (Lan et al. 2023). *Assume ($\hat{\text{A1}}$), ($\hat{\text{A3}}$), ($\hat{\text{A4}}$), and that the solution set of the VIP of (V_w, ψ_w) is non-empty. Let*

$$(\mathbf{w}^+, \psi'_w(\mathbf{w}^+)) = \text{ARM}(\nabla f_w, \psi_w, \mathbf{v}_w, \xi \mid L).$$

Then, ARM takes no more than $34 \cdot \sqrt{\frac{3L}{2\xi}}$ queries to $\nabla f_w(\cdot)$ and ensures that $(\mathbf{w}^+, \psi'_w(\mathbf{w}^+))$ satisfies ξ -distance-to-solution accuracy.

Proof. It was originally shown in [Lan et al., 2023, Theorem 3.1] that ARM takes no more than $34 \sqrt{\frac{3L}{2\xi}}$ gradient queries and obtains $\mathbf{w}^+ \in \text{dom } \psi_w$, such that there exists $\tilde{\mathbf{w}} \in \text{dom } \psi_w$ with

$\nabla f_w(\tilde{\mathbf{w}}) \in -\partial\psi_w(\tilde{\mathbf{w}})$, and

$$\|2L(\mathbf{w}^+ - \mathbf{w}')\|_w \leq \frac{2}{3}\xi\|\mathbf{v}_w - \tilde{\mathbf{w}}\|_w, \text{ where } \mathbf{w}^+ = \arg \min_{\mathbf{w} \in \text{dom } \psi_w} [\langle \nabla f_w(\mathbf{w}'), \mathbf{w} \rangle + \psi_w(\mathbf{w}) + L\|\mathbf{w} - \mathbf{w}'\|_w^2].$$

By the optimality of \mathbf{w}^+ , there exists $\psi'_w(\mathbf{w}^+) \in \partial\psi_w(\mathbf{w}^+)$ such that

$$\nabla f_w(\mathbf{w}') + \psi'_w(\mathbf{w}^+) + 2LP(\mathbf{w}^+ - \mathbf{w}') = \mathbf{0}.$$

Then, we have

$$\begin{aligned} & \|\nabla f_w(\mathbf{w}^+) + \psi'_w(\mathbf{w}^+)\|_{w^*} \\ & \leq \|\nabla f_w(\mathbf{w}') + \psi'_w(\mathbf{w}^+)\|_{w^*} + \|\nabla f_w(\mathbf{w}^+) - \nabla f_w(\mathbf{w}')\|_{w^*} \\ & = \|2LP(\mathbf{w}^+ - \mathbf{w}')\|_{w^*} + \|\nabla f_w(\mathbf{w}^+) - \nabla f_w(\mathbf{w}')\|_{w^*} \\ & \leq 3L\|\mathbf{w}^+ - \mathbf{w}'\|_w \\ & \leq \xi\|\mathbf{v}_w - \tilde{\mathbf{w}}\|_w. \end{aligned}$$

□

D Decoupled method for saddle problems: Concrete implementation

We are now back to considering the SPs in Eq. (1). Let us combine the technical components in Section C and present the final, implementable algorithm.

Implementation of DM-SP. We use the ARM in Lemma 18 for Minimization of Residual Norms:

$$\begin{aligned} \mathcal{M}_x^{\text{MRN}}(V_x, \hat{\psi}_x, \mathbf{v}_x, \delta_x) &\triangleq \text{ARM}(V_x, \hat{\psi}_x, \mathbf{v}_x, \frac{2\delta_x}{3} \mid L_x), \\ \mathcal{M}_y^{\text{MRN}}(V_y, \hat{\psi}_y, \mathbf{v}_y, \delta_y) &\triangleq \text{ARM}(V_y, \hat{\psi}_y, \mathbf{v}_y, \frac{2\delta_y}{3} \mid L_y). \end{aligned} \quad (8)$$

Consider the assembled norm $\|\cdot\|_{\mathcal{E}}$ with parameters α_x and α_y . Then, for any Monteiro-Svaiter Subproblem given by $(V^f, \psi_z, \mathbf{v}, \lambda)$, we leverage the solver

$$\text{FDS-ARM}(V^f, \psi_z, \mathbf{v}, \lambda) = \text{FDS}_{\|\cdot\|_{\mathcal{E}}}(V^f, \psi_z, \mathbf{v}, \lambda \mid (\mathcal{M}_x^{\text{MRN}}, \mathcal{M}_y^{\text{MRN}}))$$

Finally, we obtain the concrete algorithm DM-SP as follows:

$$\text{ROM}_{\|\cdot\|_{\mathcal{E}}}(V^f, \psi_z, \mathbf{z}^0, (\lambda_t)_{t \geq 1} \mid \text{FDS-ARM}). \quad (9)$$

Combining Lemmas 15 to 18, we obtain the following result on the computational cost.

Let us prove the main convergence lemma for SPs in Lemma 19.

Lemma 19. Under (A1) to (A3), for $\lambda_{t+1} \equiv \lambda \geq \frac{2L_{xy}}{\sqrt{\alpha_x \alpha_y}}$, DM-SP with the implementation in Eq. (9) takes no more than $2T$ communication rounds, no more than

$$T \cdot \left(1 + 34\sqrt{\frac{9L_x}{2\alpha_x \lambda}}\right)$$

queries to \mathcal{O}_x , and no more than

$$T \cdot \left(1 + 34\sqrt{\frac{9L_y}{2\alpha_y \lambda}}\right)$$

queries to \mathcal{O}_y , and obtains an ϵ -saddle point $\bar{\mathbf{z}}^T$, where

$$T = \left\lceil \frac{\alpha_x \lambda D_x^2 + \alpha_y \lambda D_y^2}{2\epsilon} \right\rceil.$$

Proof. By (A1), we have

$$\Delta(\bar{\mathbf{z}}^T) \leq \left(\sum_{t=0}^{T-1} a_{t+1}\right)^{-1} \max_{\mathbf{z} \in \mathcal{B} \cap \mathcal{Q}} \left[\sum_{t=0}^{T-1} a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{z}^{t+1} - \mathbf{z} \rangle \right].$$

Further, with $\lambda \geq 2\bar{L}_c$, by [Lemmas 15](#) and [16](#), we have

$$\begin{aligned} \Delta(\bar{\mathbf{z}}^T) &\leq \left(\sum_{t=0}^{T-1} a_{t+1} \right)^{-1} \max_{\mathbf{z} \in \mathcal{B} \cap Q} \left[\sum_{t=0}^{T-1} a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{z}^{t+1} - \mathbf{z} \rangle \right] \\ &\leq \left(\sum_{t=0}^{T-1} a_{t+1} \right)^{-1} \max_{\mathbf{z} \in \mathcal{B} \cap Q} \left[\frac{\alpha_x}{2} \|\mathbf{x}^0 - \mathbf{x}\|_x^2 + \frac{\alpha_y}{2} \|\mathbf{y}^0 - \mathbf{y}\|_y^2 \right] \\ &\leq \left(\sum_{t=0}^{T-1} \frac{1}{\lambda_{t+1}} \right)^{-1} \cdot \frac{1}{2} (\alpha_x D_x^2 + \alpha_y D_y^2) \leq \epsilon, \end{aligned}$$

where the last inequality follows from the assignments of $(\lambda_t)_{t \geq 1}$ and T . Therefore, the number of communication rounds is bounded by $2T$.

Now we count the number of gradient queries. By [Lemma 17](#), ARM always returns the solution with the required relative distance accuracy; and in view of [Lemma 18](#), it takes no more than $34\sqrt{\frac{3L_x}{\alpha_x \lambda}}$ gradient queries to $\nabla_x f$ and no more than $34\sqrt{\frac{3L_y}{\alpha_y \lambda}}$ gradient queries to $\nabla_y f$. Therefore, the numbers of gradient queries to $\nabla_x f$ and $\nabla_y f$ are bounded by $T \cdot (1 + 34\sqrt{\frac{3L_x}{\alpha_x \lambda}})$ and $T \cdot (1 + 34\sqrt{\frac{3L_y}{\alpha_y \lambda}})$, respectively. \square

The main convergence result for SPs in [Theorems 2](#) and [3](#) can be directly implied from [Lemma 19](#).

E Monotone composite variational inequality problems

In this section, we study variational inequality problems (VIPs) [[Nemirovski, 2004](#), [Juditsky et al., 2011](#)], a generalization of SPs that captures, for instance, multiplayer general-sum games.

E.1 Problem formulation

VIPs (with separable composite terms). Let us consider the VIP in [Eq. \(3\)](#), where $\mathcal{E} = \mathcal{E}_1 \times \cdots \times \mathcal{E}_K$ is the direct product of K finite-dimensional real vector spaces. For all $i \in [K]$: let the mapping $V_i: \text{dom } \psi \rightarrow \mathcal{E}_i^*$, and let the function $\psi_i: \mathcal{E}_i \rightarrow \mathbb{R} \cup \{+\infty\}$. We consider the decomposition of $V(\mathbf{z}) = (V_1(\mathbf{z}), \dots, V_K(\mathbf{z}))$ and $\psi(\mathbf{z}) = \psi_1(\mathbf{z}_1) + \cdots + \psi_K(\mathbf{z}_K)$, for all $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_K) \in \mathcal{E}$. Moreover, we denote $\text{dom } \psi = \text{dom } \psi_1 \times \cdots \times \text{dom } \psi_K \triangleq Q$.

Assumptions for VIPs. Let us make the following assumptions:

- (A2') Let $\mathbf{z}^0 = (\mathbf{z}_1^0, \dots, \mathbf{z}_K^0) \in Q$ be a given point. There exists $\mathbf{z}^* = (\mathbf{z}_1^*, \dots, \mathbf{z}_K^*) \in Q$ in the solution set of the VIP of (V, ψ) , such that for all $i \in [K]$: $\mathbf{z}_i^* \in \mathcal{B}_i$, where $\mathcal{B}_i \triangleq \{\mathbf{z}_i \in \mathcal{E}_i \mid \|\mathbf{z}_i^0 - \mathbf{z}_i\|_i \leq D_i\}$ and $D_i > 0$ is a given distance.
- (A3') The operator $V_i(\mathbf{z}_j; \mathbf{z}_{-j})$ is L_{ij} -Lipschitz continuous in $\mathbf{z}_j \in \text{dom } \psi_j$ for any fixed $\mathbf{z}_{-j} \in \text{dom } \psi_1 \times \cdots \times \text{dom } \psi_{j-1} \times \text{dom } \psi_{j+1} \times \cdots \times \text{dom } \psi_K$.¹

Let the operator family \mathcal{V}_{VIP} be comprised of all the operators with initialization points $((V_i)_{i \in [K]}, \mathbf{z}^0) \in \mathcal{P}_{\text{VIP}}$, such that Assumptions [\(A1\)](#) to [\(A3\)](#) are satisfied.

Notations. To simplify the notations, let us denote $\mathbf{z} \triangleq (\mathbf{z}_1, \dots, \mathbf{z}_K) \in Q$ in the context of VIPs. Let us denote

$$\bar{L}_{ij} \triangleq \max\{L_{ij}, L_{ji}\}, \quad A_i \triangleq D_i \left(\sum_{j \in [K] \setminus \{i\}} \bar{L}_{ij} D_j \right), \quad \text{and } B_i \triangleq \bar{L}_{ii} D_i^2, \quad \text{for all } i, j \in [K].$$

We refer to $\sum_i A_i$ as the diagonal conditioning, and we say that the diagonal conditioning dominates when $\sum_i A_i \gg \sum_i B_i$.

¹For all $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_K) \in Q$, we use the following notations for simplicity: $(\mathbf{z}_j; \mathbf{z}_{-j}) \triangleq \mathbf{z}$ and $\mathbf{z}_{-j} \triangleq (\mathbf{z}_1, \dots, \mathbf{z}_{j-1}, \mathbf{z}_{j+1}, \dots, \mathbf{z}_K)$.

E.2 Communication and computational costs

Distributed oracles. We consider a distributed setting with K agents very similar to the one in Sections 2.1 and 2.2. For all $i \in [K]$: Agent i controls decision variable $\mathbf{z}_i \in \text{dom } \psi_i$, has direct access to the function ψ_i , and has access to the oracle $\mathcal{O}_i(\mathbf{z}) = V_i(\mathbf{z})$ for $\mathbf{z} \in Q$. We consider per query to \mathcal{O}_i costs $c_i \geq 0$, $i \in [K]$.

Accuracy measure. We consider the following accuracy measure for VIPs:

$$\Delta(\bar{\mathbf{z}}) \triangleq \sup_{\mathbf{z} \in \mathcal{B} \cap Q} \langle V(\mathbf{z}), \bar{\mathbf{z}} - \mathbf{z} \rangle + \psi(\bar{\mathbf{z}}) - \psi(\mathbf{z}), \text{ for all } \mathbf{z} \in Q,$$

where $\mathcal{B} \triangleq \mathcal{B}_1 \times \dots \times \mathcal{B}_K$. We say that a point $\bar{\mathbf{z}} \in Q$ is an ϵ -approximate solution of the VIP if $\Delta(\bar{\mathbf{z}}) \leq \epsilon$. Our goal is to find such an ϵ -approximate solution for any $\epsilon > 0$.

Problem class. We formally define the overall *problem class*, denoted by $\mathcal{P}_{\text{VIP}}(\mathcal{V}_{\text{VIP}}, (\mathcal{O}_i)_{i \in [K]}, \epsilon)$, or for short \mathcal{P}_{VIP} . A specific problem instance $P \in \mathcal{P}_{\text{VIP}}$ is constructed by drawing an operator instance (with initial point) V from the operator family \mathcal{F} , equipping it with the specific distributed oracles $((\mathcal{O}_i)_{i \in [K]})$, and specifying a target accuracy $\epsilon > 0$. Solving the instance P requires an algorithm to output an ϵ -approximate solution of V utilizing the distributed oracles.

Distributed gradient-span algorithm, communication and oracle costs. We study the algorithm family for distributed VIPs analogous to the gradient-span framework for SPs. Suppose an algorithm \mathcal{M} proceeds in T rounds. In each round $t \in \{0, \dots, T-1\}$, each Agent $i \in [K]$ executes multiple local computational steps to generate a local trajectory of length τ_i^t , denoted by:

$$\hat{Z}_i^t = \left\{ \mathbf{z}_i^{t,l} = (\mathbf{z}_{i,1}^{t,l}, \dots, \mathbf{z}_{i,i}^{t,l}, \dots, \mathbf{z}_{i,K}^{t,l}) \right\}_{l=0}^{\tau_i^t - 1},$$

where $\mathbf{z}_{i,i}^{t,l}$ is the local variable updated by Agent i , and $\mathbf{z}_{i,j}^{t,l}$ (for $j \neq i$) is the delayed remote variable of Agent j utilized by Agent i . Let $Z_i^{t-1} = \bigcup_{r=0}^{t-1} \hat{Z}_i^r$ denote the accumulated history from all prior rounds (where $Z_i^{-1} = \emptyset$). Within round t , the history up to local step l is denoted by $Z_i^{t,l} = Z_i^{t-1} \cup \{\mathbf{z}_i^{t,r}\}_{r=0}^{l-1}$. The total history accumulated up to the end of round t is given by Z_i^t .

An algorithm \mathcal{M} is called a *distributed gradient-span algorithm* for problem class \mathcal{P}_{VIP} if, when applied to any instance $P \in \mathcal{P}_{\text{VIP}}$, its trajectories satisfy the following conditions:

1. For all rounds $t \in \{0, \dots, T-1\}$, local steps $l \in \{0, \dots, \tau_i^t - 1\}$, and agents $i \in [K]$, the updated local variable satisfies

$$\mathbf{z}_{i,i}^{t,l} \in \mathbf{z}_i^0 + \mathbf{P}_i^{-1} \text{span}\{V_i(\mathbf{z}) \mid \mathbf{z} \in Z_i^{t,l}\} + \mathbf{P}_i^{-1} \text{span}\{\partial\psi_i(\mathbf{z}_{i,i}) \mid \mathbf{z} \in Z_i^{t,l+1}\},$$

and for all $j \neq i$, the delayed remote variables satisfy

$$\mathbf{z}_{i,j}^{t,l} \in \mathbf{z}_j^0 + \mathbf{P}_j^{-1} \text{span}\{V_j(\mathbf{z}) \mid \mathbf{z} \in Z_j^{t-1}\} + \mathbf{P}_j^{-1} \text{span}\{\partial\psi_j(\mathbf{z}_{j,j}) \mid \mathbf{z} \in Z_j^{t-1}\}.$$

2. After each round $t \in \{0, \dots, T-1\}$, \mathcal{M} generates a solution $\bar{\mathbf{z}}^{t+1} = (\bar{\mathbf{z}}_1^{t+1}, \dots, \bar{\mathbf{z}}_K^{t+1})$ such that for all $i \in [K]$:

$$\bar{\mathbf{z}}_i^{t+1} \in \mathbf{z}_i^0 + \mathbf{P}_i^{-1} \text{span}\{V_i(\mathbf{z}) \mid \mathbf{z} \in Z_i^t\} + \mathbf{P}_i^{-1} \text{span}\{\partial\psi_i(\mathbf{z}_{i,i}) \mid \mathbf{z} \in Z_i^t\}.$$

For a given instance $P \in \mathcal{P}_{\text{VIP}}$, we define the *communication cost* required by a distributed gradient-span algorithm \mathcal{M} on P , denoted by $T_P^{\mathcal{M}}$, as the smallest integer $k \in \{1, \dots, T\}$ such that the generated solution $\bar{\mathbf{z}}^k$ satisfies the target accuracy. Similarly, the total number of oracle queries evaluated by Agent i for instance P up to this point is given by the size of the accumulated history, denoted by $N_{i,P}^{\mathcal{M}} = |Z_i^{T_P^{\mathcal{M}}-1}|$, for all $i \in [K]$.

Let c_i denote the computational cost of evaluating a single partial oracle V_i . The (*worst-case*) *communication cost* and *oracle cost* of algorithm \mathcal{M} over the entire problem class \mathcal{P}_{VIP} are defined by taking the supremum over all instances:

$$T_{\mathcal{P}_{\text{VIP}}}^{\mathcal{M}} = \sup_{P \in \mathcal{P}_{\text{VIP}}} T_P^{\mathcal{M}} \quad \text{and} \quad N_{\mathcal{P}_{\text{VIP}}}^{\mathcal{M}} = \sup_{P \in \mathcal{P}_{\text{VIP}}} \left(\sum_{i \in [K]} c_i N_{i,P}^{\mathcal{M}} \right).$$

We first state the classic results of the EG method in Proposition 20, which remains the state-of-the-art communication complexity bound.

Proposition 20 (Juditsky et al. 2011, Eq. (6.21)). For any target accuracy $\epsilon > 0$, the communication cost of EG is bounded by

$$\mathcal{O}\left(\sum_{i \in [K]} \frac{A_i + B_i}{\epsilon}\right),$$

and the computational cost of EG is bounded by

$$\mathcal{O}\left(\left(\sum_{i \in [K]} c_i\right)\left(\sum_{i \in [K]} \frac{A_i}{\epsilon}\right) + \left(\sum_{i \in [K]} c_i\right)\left(\sum_{i \in [K]} \frac{B_i}{\epsilon}\right)\right).$$

E.3 Decoupled method for variational inequality problems

Now, we present our DM-VIP method, which extends the DM-SP into multiplayer general-sum games.

Assembled norm. Let $\alpha_i > 0$ (to be fixed later), for all $i \in [K]$. Let the block diagonal linear operator $\mathbf{P} = \alpha_1 \mathbf{P}_1 \oplus \dots \oplus \alpha_K \mathbf{P}_K$. We consider the space \mathcal{E} to be measured by the following assembled norm: $\|\mathbf{z}\|_{\mathcal{E}} = \sqrt{\langle \mathbf{P}\mathbf{z}, \mathbf{z} \rangle}$, and its dual space $\mathcal{E}^* = \mathcal{E}_1^* \times \dots \times \mathcal{E}_K^*$ to be measured by $\|\mathbf{g}\|_{\mathcal{E}^*} = \sqrt{\langle \mathbf{g}, \mathbf{P}^{-1}\mathbf{g} \rangle}$.

DM-VIP. Let us consider the case that D_i is not known, and we use the inexact estimates \hat{D}_i instead, for all $i \in [K]$. Let us denote

$$\bar{L}_c \triangleq \sqrt{\max_{j \in [K]} \left[(\alpha_j \hat{D}_j)^{-1} \sum_{i \in [K] \setminus \{j\}} \frac{\bar{L}_{ij} (\sum_{l \in [K] \setminus \{i\}} \bar{L}_{il} \hat{D}_l)}{\alpha_i} \right]}. \quad (10)$$

Now, we introduce the (template) Decoupled Reduced-Operator Method for block composite variational inequality problems (DM-VIP), as an extended variant of DM-SP to VIPs. The pseudocode is presented in [Algorithm 5](#), provided with a solver for the minimization of residual norms.

Algorithm 5 DM-VIP($K, (V_i)_{i \in [K]}, (\psi_i)_{i \in [K]}, \mathbf{z}^0, (\lambda_t)_{t \geq 1}, (\alpha_i)_{i \in [K]} \mid (\mathcal{M}_i^{\text{MRN}})_{i \in [K]}$)

Require: A solver \mathcal{M}^{MRN} for the minimization of residual norms.

- 1: $\mathbf{v}^0 = (\mathbf{v}_1^0, \dots, \mathbf{v}_K^0) = \mathbf{z}^0$.
- 2: **for** $t = 0, 1, \dots, T - 1$ **do**
- 3: For all $i \in [K]$, Agent i computes

$$(\mathbf{z}_i^{t+1}, \psi'_i(\mathbf{z}_i^{t+1})) = \mathcal{M}_i^{\text{MRN}}(V_i(\cdot; \mathbf{v}_{-i}^t), \psi_i + \frac{\alpha_i \lambda_{t+1}}{2} \|\cdot - \mathbf{v}_i^t\|_i^2, \mathbf{v}_i^t, \frac{\alpha_i \lambda_{t+1}}{2});$$

and then, all agents exchange $\mathbf{z}^{t+1} = (\mathbf{z}_1^{t+1}, \dots, \mathbf{z}_K^{t+1})$.

- 4: The agents compute $V(\mathbf{z}^{t+1})$; and then exchange

$$V_\psi(\mathbf{z}^{t+1}) \triangleq V(\mathbf{z}^{t+1}) + \psi'(\mathbf{z}^{t+1}) \equiv V(\mathbf{z}^{t+1}) + (\psi'_1(\mathbf{z}_1^{t+1}), \dots, \psi'_K(\mathbf{z}_K^{t+1})).$$

- 5: $a_{t+1} = \frac{2\langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v}^t - \mathbf{z}^{t+1} \rangle}{\|V_\psi(\mathbf{z}^{t+1})\|_{\mathcal{E}^*}^2}$.

- 6: $\mathbf{v}^{t+1} = (\mathbf{v}_1^{t+1}, \dots, \mathbf{v}_K^{t+1}) = \arg \min_{\mathbf{v} \in Q} \left[a_{t+1} \langle V_\psi(\mathbf{z}^{t+1}), \mathbf{v} \rangle + \frac{1}{2} \|\mathbf{v} - \mathbf{v}^t\|_{\mathcal{E}}^2 \right]$.

- 7: **end for**
-

We say [Algorithm 5](#) is a template method because we have not yet specified the solver. We defer the detailed implementation to [Eq. \(17\)](#) to the end of this section.

Theorem 21. Under (A2') and (A3'), for $\lambda_{t+1} \equiv \lambda \geq 2\bar{L}_c$, DM-VIP ([Algorithm 5](#)) can be implemented with no more than

$$2T$$

communication rounds and no more than

$$T \cdot \left(1 + C_0 \cdot \frac{3L_{ii}}{\alpha_i \lambda}\right)$$

queries to V_i , for all $i \in [K]$, and obtains an ϵ -approximate solution

$$\bar{\mathbf{z}}^T = (\bar{\mathbf{z}}_1^T, \dots, \bar{\mathbf{z}}_K^T) = \left(\sum_{t=1}^T a_t \right)^{-1} \sum_{t=1}^T a_t \mathbf{z}^t,$$

where

$$T = \left\lceil \frac{\sum_i \alpha_i \lambda D_i^2}{2\epsilon} \right\rceil$$

and $C_0 > 0$ is some fixed constant.

Moreover, assume that the target accuracy $\epsilon \leq \sum_{i,j \in [K], i \neq j} \bar{L}_{ij} D_i D_j$. For the choices of $\alpha_i = \frac{\sum_{j \in [K] \setminus \{i\}} \bar{L}_{ij} \hat{D}_j}{\hat{D}_i}$ for all $i \in [K]$, and $\lambda = 2$, the communication cost is bounded by

$$2 + \sum_{i,j \in [K], i \neq j} \frac{\bar{L}_{ij} D_i D_j}{\epsilon} \left(\frac{D_i \hat{D}_j}{\hat{D}_i D_j} + \frac{\hat{D}_i D_j}{D_i \hat{D}_j} \right) \triangleq T^{\text{DM-VIP}}((\hat{D}_i)_{i \in [K]}), \quad (11)$$

and the number of queries to V_i is bounded by

$$\left(1 + 3C_0 \frac{\bar{L}_{ii} \hat{D}_i}{\sum_{j \in [K] \setminus \{i\}} \bar{L}_{ij} \hat{D}_j} \right) \left[\sum_{j,l \in [K], j \neq l} \frac{\bar{L}_{jl} D_j D_l}{\epsilon} \left(\frac{D_j \hat{D}_l}{\hat{D}_j D_l} + \frac{\hat{D}_j D_l}{D_j \hat{D}_l} \right) \right] \triangleq k_i^{\text{DM-VIP}}((\hat{D}_j)_{j \in [K]}). \quad (12)$$

Corollary 22. Equation (11) is minimized when the distance estimates $(\hat{D}_i)_{i \in [K]}$ satisfy $\frac{\hat{D}_i}{D_i} = \frac{\hat{D}_j}{D_j}$, for all $i, j \in [K]$. This results in

$$\min_{(\hat{D}_i)_{i \in [K]}} T^{\text{DM-VIP}}((\hat{D}_i)_{i \in [K]}) = 2 + \frac{2}{\epsilon} \sum_{i \in [K]} A_i, \quad (13)$$

and in the meantime, for all $i \in [K]$,

$$k_i^{\text{DM-VIP}}((\hat{D}_j)_{j \in [K]}) = \frac{2}{\epsilon} \left(1 + 3C_0 \cdot \frac{B_i}{A_i} \right) \left(\sum_{j \in [K]} A_j \right). \quad (14)$$

Remark 6 (Improved communication cost under comparable computational costs). The classic EG method takes

$$\frac{1}{\epsilon} \sum_{i \in [K]} (A_i + B_i)$$

communication rounds, and the same number of queries to V_i for all $i \in [K]$ [Juditsky et al., 2011, Eq. (6.21)]. When the gradient estimates $(\hat{D}_i)_{i \in [K]}$ satisfy $\frac{\hat{D}_i D_j}{D_i \hat{D}_j} = \Theta(1)$ for all $i, j \in [K]$, our communication complexity in Eq. (13) is consistently no worse compared to the communication complexity of EG, and is substantially faster when the ‘‘diagonal conditioning’’ dominates—that is,

$$\sum_{i \in [K]} B_i \gg \sum_{i \in [K]} A_i.$$

Moreover, our computational cost (under the same choice of parameters) is bounded by

$$\frac{2}{\epsilon} \left(\sum_{i \in [K]} c_i \right) \left(\sum_{i \in [K]} A_i \right) + \frac{3C_0}{\epsilon} \left(\sum_{i \in [K]} \frac{B_i c_i}{A_i} \right) \left(\sum_{i \in [K]} A_i \right). \quad (15)$$

Compared to the computational cost of EG, given by

$$\frac{1}{\epsilon} \left(\sum_{i \in [K]} c_i \right) \left(\sum_{i \in [K]} A_i \right) + \frac{1}{\epsilon} \left(\sum_{i \in [K]} c_i \right) \left(\sum_{i \in [K]} B_i \right),$$

our computational cost in Eq. (15) differs primarily in the second term, and consequently, may offer an advantage or disadvantage depending on the relative conditioning of A_i , B_i , and c_i for $i \in [K]$.

E.4 Detailed proofs

E.4.1 Proof for FDS

Let us provide the detailed pseudocode of FDS for VIPs in [Algorithm 6](#). Then, we prove the correctness of the solution returned by FDS.

Algorithm 6 FDS $_{\|\cdot\|_{\mathcal{E}}}$ $((V_i)_{i \in [K]}, (\psi_i)_{i \in [K]}, \mathbf{v}, \lambda \mid (\mathcal{M}_i^{\text{MRN}})_{i \in [K]})$

Require: Solver $\mathcal{M}_i^{\text{MRN}}$ for the minimization of residual norms, for all $i \in [K]$.

- 1: **for** $i \in [K]$ **do**
 - 2: $\delta_i = \frac{\alpha_i \lambda}{2}$.
 - 3: $\hat{\psi}_i = \psi_i + \frac{\alpha_i \lambda}{2} \|\cdot - \mathbf{v}_i\|_i^2$.
 - 4: $(\mathbf{z}_i^+, \psi'_i(\mathbf{z}_i^+)) = \mathcal{M}_i^{\text{MRN}}(V_i(\cdot; \mathbf{v}_{-i}), \hat{\psi}_i, \mathbf{v}_i, \delta_i)$.
 - 5: $\psi'_i(\mathbf{z}_i^+) = \hat{\psi}'_i(\mathbf{z}_i^+) - \alpha_i \lambda \mathbf{P}_i(\mathbf{z}_i^+ - \mathbf{v}_i)$.
 - 6: **end for**
 - 7: **return** $(\mathbf{z}^+, \psi'(\mathbf{z}^+))$, where $\mathbf{z}^+ = (\mathbf{z}_1^+, \dots, \mathbf{z}_K^+)$ and $\psi'(\mathbf{z}^+) = (\psi'_1(\mathbf{z}_1^+), \dots, \psi'_K(\mathbf{z}_K^+))$.
-

Lemma 23. Under (A3'), for $\lambda \geq 2\bar{L}_c$, FDS ([Algorithm 6](#)) returns the correct solution of the MS subproblem given by $(V, \psi, \mathbf{v}, \lambda)$.

Proof of Lemma 23. For all $i \in [K]$, by (A3') and then by the relative distance accuracy, we have

$$\begin{aligned}
& \|V_i(\mathbf{z}^+) + \psi'_i(\mathbf{z}_i^+) + \alpha_i \lambda \mathbf{P}_i(\mathbf{z}_i^+ - \mathbf{v}_i)\|_{i^*} \\
& \leq \|V_i(\mathbf{z}_i^+; \mathbf{v}_{-i}) + \psi'_i(\mathbf{z}_i^+) + \alpha_i \lambda \mathbf{P}_i(\mathbf{z}_i^+ - \mathbf{v}_i)\|_{i^*} + \sum_{j \in [K] \setminus \{i\}} L_{ij} \|\mathbf{z}_j^+ - \mathbf{v}_j\|_j \\
& \leq \delta_i \|\mathbf{z}_i^+ - \mathbf{v}_i\|_i + \sum_{j \in [K] \setminus \{i\}} L_{ij} \|\mathbf{z}_j^+ - \mathbf{v}_j\|_j.
\end{aligned} \tag{16}$$

Finally, we assemble the norms:

$$\begin{aligned}
& \|V(\mathbf{z}^+) + \psi'(\mathbf{z}^+) + \lambda \mathbf{P}(\mathbf{z}^+ - \mathbf{v})\|_{\mathcal{E}^*}^2 \\
& = \sum_{i \in [K]} \alpha_i^{-1} \|V_i(\mathbf{z}^+) + \psi'_i(\mathbf{z}_i^+) + \alpha_i \lambda \mathbf{P}_i(\mathbf{z}_i^+ - \mathbf{v}_i)\|_{i^*}^2 \\
& \stackrel{(16)}{\leq} \sum_{i \in [K]} \alpha_i^{-1} \left(\delta_i \|\mathbf{z}_i^+ - \mathbf{v}_i\|_i + \sum_{j \in [K] \setminus \{i\}} L_{ij} \|\mathbf{z}_j^+ - \mathbf{v}_j\|_j \right)^2 \\
& \leq \sum_{i \in [K]} 2\alpha_i^{-1} \left[\delta_i^2 \|\mathbf{z}_i^+ - \mathbf{v}_i\|_i^2 + \left(\sum_{j \in [K] \setminus \{i\}} L_{ij} \|\mathbf{z}_j^+ - \mathbf{v}_j\|_j \right)^2 \right] \\
& = \frac{\lambda^2}{2} \sum_{i \in [K]} \left(\alpha_i \|\mathbf{z}_i^+ - \mathbf{v}_i\|_i^2 \right) + 2 \sum_{i \in [K]} \left[\alpha_i^{-1} \left(\sum_{j \in [K] \setminus \{i\}} L_{ij} \|\mathbf{z}_j^+ - \mathbf{v}_j\|_j \right)^2 \right] \\
& \leq \frac{\lambda^2}{2} \|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}^2 + 2 \sum_{i \in [K]} \left[\alpha_i^{-1} \left(\sum_{l \in [K] \setminus \{i\}} L_{il} \hat{D}_l \right) \left(\sum_{j \in [K] \setminus \{i\}} \frac{L_{ij}}{\hat{D}_j} \|\mathbf{z}_j^+ - \mathbf{v}_j\|_j^2 \right) \right] \\
& = \frac{\lambda^2}{2} \|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}^2 + 2 \sum_{j \in [K]} \left[\frac{\|\mathbf{z}_j^+ - \mathbf{v}_j\|_j^2}{\hat{D}_j} \left(\sum_{i \in [K] \setminus \{j\}} \frac{L_{ij} (\sum_{l \in [K] \setminus \{i\}} L_{il} \hat{D}_l)}{\alpha_i} \right) \right] \\
& \stackrel{(10)}{\leq} \frac{\lambda^2}{2} \|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}^2 + 2 \|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}^2 \bar{L}_c^2 \\
& \leq \frac{\lambda^2}{2} \|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}^2 + \frac{\lambda^2}{2} \|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}^2 = \lambda^2 \|\mathbf{z}^+ - \mathbf{v}\|_{\mathcal{E}}^2.
\end{aligned}$$

□

E.4.2 Proof for MRN

Our algorithm is built upon [Lemma 24](#), the proof of which can be found in [[Boj and Chenchene, 2024](#), Corollary 2.4].

Lemma 24. *Assume $(\hat{A}1)$, $(\hat{A}3)$, and that the solution set of the VIP of (V_w, ψ_w) is non-empty. Then, there exists an algorithm, denoted by $(\mathbf{w}^+, \psi'_w(\mathbf{w}^+)) = \text{FEGM}(V_w, \psi_w, \mathbf{v}_w, \xi \mid L)$, which takes no more than $C_0 \cdot \frac{L}{\xi}$ operator queries and returns $(\mathbf{w}^+, \psi'_w(\mathbf{w}^+))$ that satisfies ξ -distance-to-solution accuracy, where $C_0 > 0$ is some fixed constant.*

E.4.3 Concrete implementation

We are now back to considering the VIPs. Let us use FEGM in [Lemma 24](#) for the minimization of residual norms:

$$\mathcal{M}_i^{\text{FEGM}}(\hat{V}_i, \hat{\psi}_i, \mathbf{v}_i, \delta_i) \triangleq \text{FEGM}(\hat{V}_i, \hat{\psi}_i, \mathbf{v}_i, \frac{2\delta_i}{3} \mid L_{ii}).$$

Then, for any Monteiro-Svaiter Subproblem given by $(V, \psi, \mathbf{v}, \lambda)$, we leverage the solver

$$\text{FDS-FEGM}(V, \psi, \mathbf{v}, \lambda) = \text{FDS}_{\|\cdot\|_{\mathcal{E}}}(V, \psi, \mathbf{v}, \lambda \mid (\mathcal{M}_i^{\text{FEGM}})_{i \in [K]}).$$

Finally, we obtain the concrete algorithm DM-VIP as follows:

$$\text{ROM}_{\|\cdot\|_{\mathcal{E}}}\left((V_i)_{i \in [K]}, (\psi_i)_{i \in [K]}, \mathbf{z}^0, (\lambda_t)_{t \geq 1} \mid \text{FDS-FEGM}\right). \quad (17)$$

Combining [Lemmas 15 to 17](#) and [24](#), with the implementation in [Eq. \(17\)](#), we conclude that [Theorem 21](#) holds for the constant C_0 from [Lemma 24](#). We include the complete proof below.

Proof of [Theorem 21](#). By [\(A1'\)](#), we have

$$\Delta(\bar{\mathbf{z}}^T) \leq \left(\sum_{t=0}^{T-1} a_{t+1}\right)^{-1} \max_{\mathbf{z} \in \mathcal{B} \cap \mathcal{Q}} \left[\sum_{t=0}^{T-1} a_{t+1} \langle V_{\psi}(\mathbf{z}^{t+1}), \mathbf{z}^{t+1} - \mathbf{z} \rangle \right].$$

Further, with $\lambda \geq 2\bar{L}_c$, by [Lemmas 15](#) and [16](#), we have

$$\begin{aligned} \Delta(\bar{\mathbf{z}}^T) &\leq \left(\sum_{t=0}^{T-1} a_{t+1}\right)^{-1} \max_{\mathbf{z} \in \mathcal{B} \cap \mathcal{Q}} \left[\sum_{t=0}^{T-1} a_{t+1} \langle V_{\psi}(\mathbf{z}^{t+1}), \mathbf{z}^{t+1} - \mathbf{z} \rangle \right] \\ &\leq \left(\sum_{t=0}^{T-1} a_{t+1}\right)^{-1} \left[\sum_{i \in [K]} \left(\frac{\alpha_i}{2} \max_{\mathbf{z}_i \in \mathcal{B}_i \cap \text{dom } \psi_i} \|\mathbf{z}_i^0 - \mathbf{z}_i\|_i^2\right) \right] \\ &\leq \left(\sum_{t=0}^{T-1} \frac{1}{\lambda_{t+1}}\right)^{-1} \cdot \frac{1}{2} \sum_{i \in [K]} \alpha_i D_i^2 \leq \epsilon, \end{aligned}$$

where the last inequality follows from the assignments of $(\lambda_t)_{t \geq 1}$ and T . Therefore, the number of communication rounds is bounded by $2T$.

Now we count the number of gradient queries. By [Lemma 17](#), FEGM always returns the solution with the required relative distance accuracy; and in view of [Lemma 24](#), it takes no more than $C_0 \cdot \frac{3L_{ii}}{\alpha_i \lambda}$ gradient queries to V_i , for all $i \in [K]$. Therefore, the numbers of queries to V_i are bounded by $T \cdot (1 + C_0 \cdot \frac{3L_{ii}}{\alpha_i \lambda})$, for all $i \in [K]$. \square