

# IMPLICIT MODELS: EXPRESSIVE POWER SCALES WITH TEST-TIME COMPUTE

JIALIN LIU, LISANG DING, WOTAO YIN, AND STANLEY OSHER

**ABSTRACT.** Implicit models, an emerging model class, compute outputs by iterating a single parameter block to a fixed point. This architecture realizes an infinite-depth, weight-tied network that trains with constant memory, significantly reducing memory needs for the same level of performance compared to explicit models. While it is empirically known that these compact models can often match or even exceed larger explicit networks by allocating more test-time compute, the underlying reasons are not yet well understood.

We study this gap through a non-parametric analysis of expressive power. We provide a strict mathematical characterization, showing that a simple and regular implicit operator can, through iteration, progressively express more complex mappings. We prove that for a broad class of implicit models, this process allows the model’s expressive power to grow with test-time compute, ultimately matching a much richer function class. The theory is validated across three domains: image reconstruction, scientific computing, and operations research, demonstrating that as test-time iterations increase, the complexity of the learned mapping rises, while the solution quality simultaneously improves and stabilizes.

## 1. INTRODUCTION

Many machine-learning tasks can be cast as learning a mapping  $\mathcal{F}$  from input  $\mathbf{x}$  to the desired output  $\mathbf{y}_*$ , i.e.,  $\mathbf{y}_* = \mathcal{F}(\mathbf{x})$ . An emerging alternative is the *implicit models* (or named deep equilibrium models, also fixed-point models): train an operator  $\mathcal{G}$  whose fixed point matches the target,  $\mathbf{y}_* = \mathcal{G}(\mathbf{y}_*, \mathbf{x})$  [3, 24]. At inference, we repeatedly apply the same learned operator  $\mathcal{G}$  (weight-tied across all iterations  $t$ ):

$$(1) \quad \mathbf{y}_1 = \mathcal{G}(\mathbf{y}_0, \mathbf{x}), \mathbf{y}_2 = \mathcal{G}(\mathbf{y}_1, \mathbf{x}), \mathbf{y}_3 = \mathcal{G}(\mathbf{y}_2, \mathbf{x}), \dots,$$

and expect  $\mathbf{y}_t(\mathbf{x}) \rightarrow \mathbf{y}_*(\mathbf{x}) = \mathcal{F}(\mathbf{x})$  for all  $\mathbf{x}$ . Rather than producing  $\mathbf{y}_*$  in a single feed-forward pass, implicit models reach the target through gradual equilibrium-seeking updates. Here, “*test-time compute*” refers to the computational budget spent at inference—primarily the number of iterations; increasing this budget increases runtime but not the number of learned parameters. By tailoring the structure of  $\mathcal{G}$ , implicit models have shown strong results across many domains (e.g., imaging [35], scientific computing [67], generative modeling [33, 73], LLM reasoning [32], etc.).

Behind these successes, the advantages of implicit models include: **(i)** they realize an infinite-depth, weight-tied network trainable with constant memory, which yields efficient training [29, 34]; **(ii)** they allow us to explicitly impose or “implicitly bake in” domain constraints and structure (e.g., physics, geometry, safety), see [38, 70, 97]; and, *most surprisingly*, **(iii)** *they can often match or even exceed larger explicit networks by allocating more iterations* [32, 67, 94]. Point (i) stems from the weight-tied architecture and avoiding full back-propagation. Point (ii) arises from the inherently implicit nature of many real-world, equation-based constraints. In contrast, the mechanism underlying the surprising effectiveness of (iii) remains less well understood.

We study this through the lens of expressive power—the set of input–output maps a model family can represent. We ask two questions. First, as a baseline: **(Q1)** *Do implicit models (at least) match the expressive power of explicit ones?* Concretely, for a target map  $\mathcal{F} : \mathbf{x} \mapsto \mathbf{y}_*$ , does there always exist an implicit operator  $\mathcal{G}$  such that the iterates of (1) satisfy  $\mathbf{y}_t(\mathbf{x}) \rightarrow \mathcal{F}(\mathbf{x})$  for all  $\mathbf{x}$ ? If yes, a more

insightful question follows: **(Q2)** *Do implicit models offer an expressive advantage?* In particular, can a relatively *simple* implicit operator  $\mathcal{G}$ , through iteration, represent a *complex* explicit map  $\mathcal{F}$ ? A positive answer to (Q2) would directly explain phenomenon (iii).

To our knowledge, these questions remain largely open, as prior work on the topic is limited. For instance, the universality of implicit models was discussed in specific settings as a secondary result [3, 67], and a complete theory is missing. Closer to our goal, [96] focus on expressive power, partially answering (Q2) by proving a key separation result: some mappings can be realized by an implicit model but not by much larger explicit counterparts. However, a precise characterization of the function class that implicit models can represent remains open. Our work fills this gap from a nonparametric, function-space perspective, revealing a key principle: an implicit model’s expressive power scales with test-time compute. Specifically:

- **Expressive boundary.** We identify a natural target class (locally Lipschitz mappings) and prove: simple and well-behaved (which we term “regular”) implicit models, by progressive iterations, can *express any mapping* in this class and can *only express mappings* in this class.
- **Emergent expressive power.** Our theory yields a new viewpoint on implicit models: the expressive power of a regular implicit operator is not static but *grows with iteration* (i.e., scales with test-time compute) and finally matches a much richer function class.
- **Validation across domains.** We validate our theory with case studies in a wide range of applications (e.g., image reconstruction, scientific computing, and operations research). For a representative problem in each domain, we demonstrate that as test-time iterations increase, the empirical complexity of the iterate  $\mathbf{y}_t(\mathbf{x})$  rises while solution quality improves and stabilizes.

## 2. MAIN RESULTS

We now return to (Q1): given a target map  $\mathcal{F}$ , does there exist an implicit operator  $\mathcal{G}$  whose fixed-point iteration yields  $\mathbf{y}_t(\mathbf{x}) \rightarrow \mathcal{F}(\mathbf{x})$ ? A naive construction answers “yes”: define, for  $0 < \eta < 1$ ,

$$(2) \quad \mathcal{G}(\mathbf{y}, \mathbf{x}) := (1 - \eta)\mathbf{y} + \eta\mathcal{F}(\mathbf{x}).$$

Then the fixed-point iteration reduces to  $\mathbf{y}_t = (1 - \eta)\mathbf{y}_{t-1} + \eta\mathcal{F}(\mathbf{x})$ , hence  $\mathbf{y}_t - \mathcal{F}(\mathbf{x}) = (1 - \eta)(\mathbf{y}_{t-1} - \mathcal{F}(\mathbf{x}))$ . As  $0 < \eta < 1$ , it holds that, for all  $\mathbf{x}$ ,  $\mathbf{y}_t(\mathbf{x}) - \mathcal{F}(\mathbf{x}) \rightarrow \mathbf{0}$  as  $t \rightarrow \infty$ .

However, (2) is merely a trivial averaging of  $\mathbf{y}$  and  $\mathcal{F}(\mathbf{x})$ ; learning such an implicit model is no different from learning  $\mathcal{F}$  directly. This prompts the natural follow-up: is there any *nontrivial* implicit representation that is able to indicate the expressive benefits of implicit models?

**An illustrative example.** Let  $\mathcal{F}(x) = 1/x$  on  $[-1, 1] \setminus \{0\}$ . This function is smooth (differentiable to any order) almost everywhere, but blows up near the singular point  $x = 0$ :

$$|\mathcal{F}(x)| = \left| \frac{1}{x} \right| \rightarrow \infty, \quad \left| \frac{d\mathcal{F}}{dx} \right| = \left| -\frac{1}{x^2} \right| \rightarrow \infty, \quad \text{as } x \rightarrow 0.$$

Neural networks approximating  $1/x$  on  $[-1, -\delta) \cup (\delta, 1]$  typically demands higher network complexity—i.e., increasing depth/width as  $\delta \rightarrow 0$  to capture the growing steepness near the singularity [90]. If we adopt the naive implicit form (2),  $\mathcal{G}(y, x) = (1 - \eta)y + \eta/x$ , nothing is gained: the model still inherits the singular behavior  $|\partial\mathcal{G}/\partial x| = \eta/x^2 \rightarrow \infty$ .

What would be a nontrivial implicit representation in this setting? Instead of writing  $(1/x)$  explicitly, we can regard it as the solution of the equation  $xy - 1 = 0$  (**implicit representation**). Inspired by this, we apply a fixed-point iteration to  $xy - 1 = 0$ :  $\mathcal{G}(y, x) = y - \eta(xy - 1)$ . Using the general scheme in (1), we have  $y_t = y_{t-1} - \eta(xy_{t-1} - 1)$ . Subtracting the true solution gives

$$y_t - \frac{1}{x} = y_{t-1} - \frac{1}{x} - \eta x \left( y_{t-1} - \frac{1}{x} \right) = (1 - \eta x) \left( y_{t-1} - \frac{1}{x} \right)$$

For any  $0 < \eta < 1$  and any  $x \in (0, 1]$ , we have  $0 < (1 - \eta x) < 1$  which implies  $y_t \rightarrow 1/x$ . (For  $x < 0$ , simply flip the stepsize sign,  $\eta$  to  $-\eta$ .) This implicit formulation is much simpler and more elegant: the operator  $\mathcal{G}(y, x) = y - \eta(xy - 1)$  has *no singularity* and *no blow-up*.

The example indicates: intuitively, an implicit representation can realize a complicated map with singularities via a much simpler, smoother update operator  $\mathcal{G}$ . Next, we make it precise: we formally define what we mean by “simple” versus “complex,” and characterize—beyond the  $1/x$  example—the class of target functions for which an implicit representation admits such a simple form.

**Definition 2.1** (Lipschitz continuity). Let  $(\mathbb{X}, \|\cdot\|)$  and  $(\mathbb{Y}, \|\cdot\|)$  be normed spaces, and let  $\mathcal{Q} : \mathbb{X} \rightarrow \mathbb{Y}$ . We say  $\mathcal{Q}$  is *L-Lipschitz* (globally Lipschitz) on  $\mathbb{X}$  if there exists  $L > 0$  such that

$$\|\mathcal{Q}(\mathbf{x}_1) - \mathcal{Q}(\mathbf{x}_2)\| \leq L \|\mathbf{x}_1 - \mathbf{x}_2\| \quad \text{for all } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{X},$$

and the smallest such  $L$  is the *Lipschitz constant* (or *Lipschitz modulus*), denoted as  $\text{Lip}(\mathcal{Q})$ . If the Lipschitz constant  $L < 1$ , we say  $\mathcal{Q}$  is *L-contractive* on  $\mathbb{X}$ . Given  $\mathbf{x} \in \mathbb{X}$ , we say  $\mathcal{Q}$  is *locally Lipschitz at  $\mathbf{x}$*  if there exists a neighborhood  $\mathbb{U}$  of  $\mathbf{x}$  on which  $\mathcal{Q}$  is  $L_{\mathbb{U}}$ -Lipschitz continuous for some  $L_{\mathbb{U}} > 0$ . If  $\mathcal{Q}$  is locally Lipschitz at every  $\mathbf{x} \in \mathbb{X}$ , we say  $\mathcal{Q}$  is *locally Lipschitz on  $\mathbb{X}$* .

Intuitively, Lipschitz continuity limits how quickly a function’s value can change. When a function is differentiable, its Lipschitz modulus can be characterized by the norm of its first derivative via the mean-value theorem. For example,  $\mathcal{F}(x) = 1/x$  is locally Lipschitz on  $[-1, 1] \setminus \{0\}$  but not globally Lipschitz there, since  $|\mathrm{d}\mathcal{F}/\mathrm{d}x| = 1/x^2$  is unbounded as  $x \rightarrow 0$ , causing local Lipschitz constants to blow up near the singularity. In contrast, the implicit update  $\mathcal{G}(y, x) = y - \eta(xy - 1)$  has simple partial derivatives  $|\partial\mathcal{G}/\partial x| = |\eta y|$  and  $|\partial\mathcal{G}/\partial y| = |1 - \eta x|$  without singularity.

Locally Lipschitz mappings form a much richer class than globally Lipschitz ones. Typical examples (locally Lipschitz everywhere in their domains but not globally Lipschitz on the whole set) include:  $\log x$  in  $(0, 1]$ ,  $\tan x$  in  $(-\frac{\pi}{2}, \frac{\pi}{2})$ ,  $\sqrt{x}$  in  $(0, 1]$ ,  $\Gamma(x)$  in  $\mathbb{R} \setminus \{0, -1, -2, \dots\}$ , etc.

For this reason, we refer to globally Lipschitz maps as “simple” operators and locally Lipschitz maps (which may exhibit large local slopes near certain inputs) as “complex.” Next, we formally state our main result: identifying a broad family of target functions for which implicit representations provide such simple update operators while expressing complex fixed-point mappings.

**Assumption 2.2.** Let  $\mathbb{X} \subset \mathbb{R}^d$  be bounded and  $\mathcal{F} : \mathbb{X} \rightarrow \mathbb{R}^n$  be locally Lipschitz on  $\mathbb{X}$ .

We only assume the domain  $\mathbb{X}$  is *bounded*; it need not be compact, closed, or connected. For instance,  $\mathbb{X} = [-1, 0) \cup (0, 1]$  excludes the singular point and permits  $\mathcal{F}(x) = 1/x$  to blow up at the interior gap  $x = 0$  while remaining locally Lipschitz on  $\mathbb{X}$ . Another example is  $\mathbb{X} = (-\frac{\pi}{2}, \frac{\pi}{2})$ , where  $\mathcal{F}(x) = \tan x$  is locally Lipschitz though it blows up at the boundary points  $\pm\frac{\pi}{2}$  (not in  $\mathbb{X}$ ).

We now formalize what we mean by “simple” update rules—namely, *regular implicit operators*.

**Definition 2.3** (Regular implicit operator). Let  $\mathbb{X} \subset \mathbb{R}^d$  be bounded. An operator  $\mathcal{G} : \mathbb{R}^n \times \mathbb{X} \rightarrow \mathbb{R}^n$  is *regular* if: (i) For any  $\mathbf{y} \in \mathbb{R}^n$ , the map  $\mathbf{x} \mapsto \mathcal{G}(\mathbf{y}, \mathbf{x})$  is globally Lipschitz (w.r.t.  $\mathbf{x}$ ) on  $\mathbb{X}$ , and the Lipschitz constant grows linearly w.r.t.  $\|\mathbf{y}\|$ , and (ii) For each  $\mathbf{x} \in \mathbb{X}$ , there exists  $\mu(\mathbf{x}) \in (0, 1)$ , the map  $\mathbf{y} \mapsto \mathcal{G}(\mathbf{y}, \mathbf{x})$  is  $\mu(\mathbf{x})$ -contractive on  $\mathbb{R}^n$ , and  $\mu(\mathbf{x})$  is continuous w.r.t.  $\mathbf{x}$ .

**Explanation.** A regular  $\mathcal{G}$  satisfies: (i) Fixing  $\mathbf{y}$ ,  $\mathcal{G}(\mathbf{y}, \cdot)$  is *globally Lipschitz in  $\mathbf{x}$* , this makes it a “simple” operator, and (ii) Fixing  $\mathbf{x}$ ,  $\mathcal{G}(\cdot, \mathbf{x})$  is *contractive in  $\mathbf{y}$* ; by Banach’s theorem, this yields a unique fixed point  $\mathbf{y}_*(\mathbf{x})$  for each  $\mathbf{x}$  and guarantees that iterates of (1) converge to it:  $\mathbf{y}_t(\mathbf{x}) \rightarrow \mathbf{y}_*(\mathbf{x})$ . **An example** of such a regular  $\mathcal{G}$  is the aforementioned  $\mathcal{G}(y, x) = y - \eta(xy - 1)$  on  $x \in (0, 1]$  with  $0 < \eta < 1$ . With this definition, we present our main results.

**Theorem 2.4** (Sufficiency). *Under Assumption 2.2, for any  $\mathcal{F}$  there exists a regular implicit operator  $\mathcal{G} : \mathbb{R}^n \times \mathbb{X} \rightarrow \mathbb{R}^n$  whose fixed-point map reproduces  $\mathcal{F}$ :  $\text{Fix}(\mathcal{G}(\cdot, \mathbf{x})) = \mathcal{F}(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{X}$ .*

**Theorem 2.5** (Necessity). *Let  $\mathbb{X} \subset \mathbb{R}^d$  be bounded and let  $\mathcal{G} : \mathbb{R}^n \times \mathbb{X} \rightarrow \mathbb{R}^n$  be regular. Then, for every  $\mathbf{x} \in \mathbb{X}$ , the map  $\mathbf{y} \mapsto \mathcal{G}(\mathbf{y}, \mathbf{x})$  has a unique fixed point  $\mathbf{y}_*(\mathbf{x})$ , and the resulting fixed-point map  $\mathbf{y}_*(\mathbf{x})$  must be locally Lipschitz on  $\mathbb{X}$ .*

Proofs are deferred to Appendix A. Theorem 2.4 provides an affirmative answer to (Q1) and (Q2) posed in the introduction. It proves that for any locally Lipschitz target  $\mathcal{F}$  on a bounded domain, there exists a *regular* implicit operator  $\mathcal{G}$ , whose iterations converge to the target  $\mathbf{y}_t(\mathbf{x}) \rightarrow \mathcal{F}(\mathbf{x})$  for all  $\mathbf{x}$ .

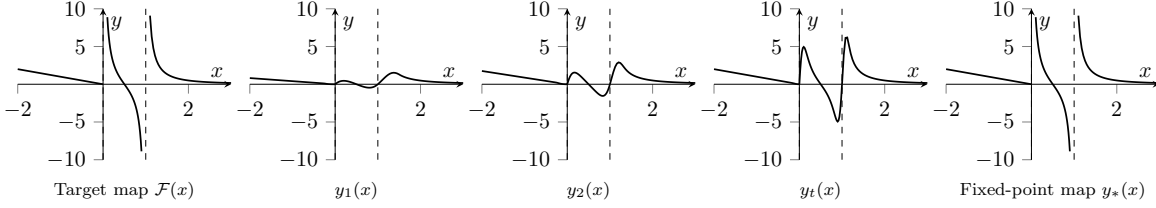


FIGURE 1. (Conceptual diagram) A simple implicit update expresses a complex map via iteration.

This demonstrates that the expressive power of implicit models **not only matches** that of explicit models **but also provides a distinct expressive benefit**: *a relatively simple (regular) implicit representation can yield a complex fixed-point mapping*. Complementarily, Theorem 2.5 shows the boundary is tight: fixed points induced by any regular  $\mathcal{G}$  are *necessarily* locally Lipschitz. Together, the two results give an exact expressivity characterization for regular implicit models.

**What does our theory imply?** Take any locally Lipschitz target  $\mathcal{F}$  (e.g., the curve in Fig. 1). Our results guarantee the existence of a *regular* implicit operator  $\mathcal{G}$  such that the iteration  $\mathbf{y}_t = \mathcal{G}(\mathbf{y}_{t-1}, \mathbf{x})$  with  $\mathbf{y}_0 = \mathbf{0}$  converges:  $\mathbf{y}_t(\mathbf{x}) \rightarrow \mathcal{F}(\mathbf{x})$ . Consider the first iterate:

$$\mathbf{y}_1(\mathbf{x}) = \mathcal{G}(\mathbf{0}, \mathbf{x}) \quad \implies \quad \text{Lip}(\mathbf{y}_1) = \sup_{\mathbf{x}, \mathbf{x}'} \frac{\|\mathcal{G}(\mathbf{0}, \mathbf{x}) - \mathcal{G}(\mathbf{0}, \mathbf{x}')\|}{\|\mathbf{x} - \mathbf{x}'\|} = \text{Lip}(\mathcal{G}(\mathbf{0}, \cdot)).$$

Thus, *after a single iteration,  $\mathbf{y}_1(\cdot)$  is exactly as smooth as  $\mathcal{G}(\mathbf{0}, \cdot)$* . By Definition 2.3,  $\mathcal{G}(\mathbf{0}, \cdot)$  is globally Lipschitz in  $\mathbf{x}$  with a *uniform* constant, hence the map  $\mathbf{y}_1(\cdot)$  inherits the same property. With more iterations,  $\mathbf{y}_t$  approaches  $\mathcal{F}$ , so the Lipschitz constant of  $\mathbf{y}_t(\cdot)$  approaches that of  $\mathcal{F}(\cdot)$ :

$$\lim_{t \rightarrow \infty} \frac{\|\mathbf{y}_t(\mathbf{x}) - \mathbf{y}_t(\mathbf{x}')\|}{\|\mathbf{x} - \mathbf{x}'\|} = \frac{\|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{x}')\|}{\|\mathbf{x} - \mathbf{x}'\|}.$$

If  $\mathcal{F}$  has singularities (where local slopes are large or even unbounded), then the effective Lipschitz constants of  $\mathbf{y}_t(\cdot)$  necessarily grow with  $t$  to match that complexity—illustrated in Fig. 1. In short: *a regular operator  $\mathcal{G}$  gains expressive power through iteration (“test-time compute”)*, ultimately matching a richer function class that can include singularities and unbounded Lipschitz behavior.

**Generalization.** Someone may ask: does a large Lipschitz constant of the fixed-point map  $\mathbf{y}_*(\mathbf{x})$  imply sensitivity or poor generalization (cf. [71])? Our view is that this sensitivity is *inherent to the target  $\mathcal{F}$* , not to the implicit representation—any faithful model, explicit or implicit, must track  $\mathcal{F}$ ’s sharp variations. Our case studies in Section 3 confirm this: the target  $\mathcal{F}$  in many tasks is indeed steep somewhere and the effective Lipschitz growth as accuracy improves. Crucially, the implicit formulation can realize such targets with a *simple* operator  $\mathcal{G}$ , which regularizes training and supports good generalization in practice.

**Insights for practitioners.** A substantial line of work (e.g., [24, 39, 45, 80, 95]) enforces a global Lipschitz bound on the fixed-point map  $\mathbf{y}_*(\mathbf{x})$ . Typically, the model is parameterized as  $\mathcal{G}(\mathbf{y}, \mathbf{x}) = \sigma(\mathbf{A}\mathbf{y} + \mathbf{B}\mathbf{x} + \mathbf{b})$ , and by imposing specific algebraic structure on  $\mathbf{A}$  and  $\mathbf{B}$ , one ensures that  $\mathbf{y}_*(\mathbf{x})$  is globally Lipschitz in  $\mathbf{x}$ . While this indeed improves robustness, our theory shows it **constrains expressivity and undercuts the unique advantage of implicit models**. Our recommendation is different: rather than imposing uniform Lipschitz constraints, incorporate case-by-case *domain-specific knowledge, priors, or constraints* (as illustrated in our case studies Sec. 3). This method provides effective regularization, leading to robustness and strong test performance while unlocking the full power of implicit models—representing complex maps with relatively simple operators.

### 3. CASE STUDIES

In this section, we present three case studies. For each task, we (i) show that the target map satisfies Assumption 2.2, bringing it under our theory; (ii) specify a  $\mathcal{G}$  architecture infused with domain priors

or constraints; (iii) verify—without any explicit Lipschitz control, using only vanilla training—that the learned  $\mathcal{G}$  is regular (i.e.,  $\mathcal{G}$  is Lipschitz in  $\mathbf{x}$  with a modest constant and the iterates converge  $\mathbf{y}_t(\mathbf{x}) \rightarrow \mathcal{F}(\mathbf{x})$ ; and (iv) test whether increasing test-time iterations enables  $\mathcal{G}$  to realize progressively more complex mappings.

**3.1. Case Study 1: Image Reconstruction (Inverse problems).** Inverse problems in imaging seek to recover an image  $\mathbf{y}_* \in \mathbb{R}^n$  from partial, noisy measurements  $\mathbf{x} = \mathbf{A}\mathbf{y}_* + \mathbf{n} \in \mathbb{R}^d$  ( $d < n$ ), where  $\mathbf{A}$  is a known linear operator and  $\mathbf{n}$  is noise. A common prior is that  $\mathbf{y}_*$  lies near a smooth data manifold  $\mathbb{M} \subset \mathbb{R}^n$ . To recover  $\mathbf{y}_*$ , a standard estimator solves

$$(3) \quad \min_{\mathbf{y} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{y}\|^2 + \frac{\alpha}{2} \text{dist}^2(\mathbf{y}, \mathbb{M}),$$

or, equivalently, a variable-splitting surrogate

$$(4) \quad \min_{\mathbf{y}, \mathbf{z} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{y}\|^2 + \frac{\alpha}{2} \text{dist}^2(\mathbf{z}, \mathbb{M}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{z}\|^2.$$

Next we will show that, under mild assumptions, both (3) and (4) admit a unique minimizer for each  $\mathbf{x}$  in a bounded set, and the solution map  $\mathbf{x} \mapsto \hat{\mathbf{y}}(\mathbf{x})$  is *locally Lipschitz*. Hence the reconstruction target falls within Assumption 2.2 and is covered by our expressivity results in Section 2.

**Assumption 3.1.** Let  $\mathbb{M} \subset \mathbb{R}^n$  be a compact,  $\mathcal{C}^2$ , embedded (possibly nonconvex) submanifold with positive reach  $\tau > 0$ . Assume the forward operator  $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^d$  is  $(\mu, L)$ -bi-Lipschitz when restricted to  $\mathbb{M}$  and let  $\sigma_{\max}$  denote the maximal singular value of  $\mathbf{A}$ .

These assumptions are modest: they are standard in prior work and supported by existing theory. Formal definitions (reach and bi-Lipschitz continuity) and relevant literature appear in Appendix C.

**Definition 3.2.** Define the admissible set of observations  $\mathbf{x}$  for (3) and (4):

$$\mathbb{X} := \left\{ \mathbf{x} : \mathbf{x} = \mathbf{A}\mathbf{y}_* + \mathbf{n}, \quad \text{for some } \mathbf{y}_* \in \mathbb{M}, \quad \|\mathbf{n}\| < \frac{1}{80} \frac{\mu^5}{\sigma_{\max}^2 L^2} \tau \right\}$$

**Theorem 3.3.** Under Assumption 3.1, there exists  $\alpha > 0$  for all  $\mathbf{x} \in \mathbb{X}$  such that the minimization problem (3) yields a unique minimizer  $\hat{\mathbf{y}}$ . Let  $\mathcal{F}_{1a} : \mathbf{x} \mapsto \hat{\mathbf{y}}$  denote the associated solution map from input  $\mathbf{x}$  to the recovery  $\hat{\mathbf{y}}$ . Then  $\mathcal{F}_{1a}$  is locally Lipschitz continuous on  $\mathbb{X}$ .

**Theorem 3.4.** Under Assumption 3.1, there exist  $\alpha, \beta > 0$  for all  $\mathbf{x} \in \mathbb{X}$  such that the minimization problem (4) yields a unique minimizer  $(\hat{\mathbf{y}}, \hat{\mathbf{z}})$ . Let  $\mathcal{F}_{1b} : \mathbf{x} \mapsto \hat{\mathbf{y}}$  denote the associated solution map from input  $\mathbf{x}$  to the recovery  $\hat{\mathbf{y}}$ . Then  $\mathcal{F}_{1b}$  is locally Lipschitz continuous on  $\mathbb{X}$ .

**Corollary 3.5.** There must be a regular implicit operator  $\mathcal{G}(\mathbf{y}, \mathbf{x})$  such that  $\text{Fix}(\mathcal{G}(\cdot, \mathbf{x})) = \mathcal{F}_{1a}(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{X}$ . The same conclusion holds for  $\mathcal{F}_{1b}(\mathbf{x})$ .

Proofs of the theorems are deferred to Appendix C, and Corollary 3.5 follows immediately from Theorems 2.4, 3.3, and 3.4. This corollary guarantees the existence of regular implicit models  $\mathcal{G}$  for image reconstruction. Next, we present how to implement  $\mathcal{G}$  in this context.

**Problem-specific  $\mathcal{G}$ .** We adopt *algorithm-inspired* designs that mirror classical solvers for (3) and (4). Parameterizing these iterative solvers gives problem-tailored implicit models. In particular,

- Option I (PGD-style). To solve (3), if  $\mathbb{M}$  were known, one would use proximal gradient descent (PGD):  $\mathbf{y}_{t+1} = \text{prox}_{\sigma}(\mathbf{y}_t - \gamma \mathbf{A}^\top(\mathbf{A}\mathbf{y}_t - \mathbf{x}))$ , with parameters  $\sigma, \gamma > 0$ , where  $\text{prox}_{\sigma}$  is the proximal map of  $(\sigma/2)\text{dist}^2(\mathbf{y}, \mathbb{M})$  (see Appendix C.1). In practice, we replace  $\text{prox}_{\sigma}$  by a learnable neural network denoiser  $\mathcal{H}_{\theta, \sigma}$  (parameters  $\theta$  and noise level input  $\sigma$ ) and obtain

$$(5) \quad \mathcal{G}_{\Theta}(\mathbf{y}, \mathbf{x}) = \mathcal{H}_{\theta, \sigma}(\mathbf{y} - \gamma \mathbf{A}^\top(\mathbf{A}\mathbf{y} - \mathbf{x})), \quad \Theta = \{\theta, \sigma, \gamma\}.$$

- Option II (HQS-style). For (4), a standard solver is half-quadratic splitting (HQS, see Appendix C.2). Similar to Option I, we replace the proximal map by a learned module and obtain

$$(6) \quad \mathcal{G}_{\Theta}(\mathbf{y}, \mathbf{x}) = \mathcal{H}_{\theta, \sigma} \left( (\mathbf{A}^\top \mathbf{A} + \beta \mathbf{I})^{-1} (\mathbf{A}^\top \mathbf{x} + \beta \mathbf{y}) \right), \quad \Theta = \{\theta, \sigma, \beta\}.$$

Here we follow the long-standing “plug-in denoiser” idea from Plug-and-Play (PnP) methods [92], which replaces a proximal operator with an off-the-shelf denoiser inside an iterative solver (see brief bibliography in Appendix C.2). Unlike PnP, one can also train the *entire*  $\mathcal{G}_\Theta$  as an implicit model, in both PGD-style [19, 35, 86, 95, 101, 104] and HQS-style [36] formulations. We adopt the latter.

**Questions.** Given the parameterizations in (5) and (6), we examine: (i) are these  $\mathcal{G}_\Theta$  operators Lipschitz with respect to  $\mathbf{x}$ ; and (ii) do they, as our theory predicts, realize progressively more complex input–output mappings over iterations despite having simple per-iteration operators?

**Experiment settings.** We study image deblurring,  $\mathbf{x} = \mathbf{A}(\mathbf{y}_*) + \mathbf{n}$ , where  $\mathbf{A}$  is a motion-blur operator and  $\mathbf{n}$  is additive Gaussian noise. Using BSDS500 [66], we construct 200 training, 100 validation, and 200 test pairs  $(\mathbf{x}, \mathbf{y}_*)$ , yielding datasets  $\mathbb{D}_{\text{inv,train}}$ ,  $\mathbb{D}_{\text{inv,val}}$ , and  $\mathbb{D}_{\text{inv,test}}$ . Implementation details (data preprocessing, model choices, and training) are in Appendix F.

For evaluation, we analyze 100 iterations of the learned dynamics,  $\mathbf{y}_{t+1}(\mathbf{x}) = \mathcal{G}_\Theta(\mathbf{y}_t(\mathbf{x}), \mathbf{x})$ ,  $0 \leq t \leq 99$  and  $\mathbf{y}_0 = \mathbf{0}$ , on the test set  $\mathbb{D}_{\text{inv,test}} = \{(\mathbf{x}_i, \mathbf{y}_i^*)\}_{i=1}^{200}$ . For each  $i$ , we create 5 perturbed ground truths  $\mathbf{y}_{i,j}^*$ ,  $1 \leq j \leq 5$ , and for each  $\mathbf{y}_{i,j}^*$ , we apply  $\mathbf{A}$ , add noise, and then obtain  $\mathbf{x}_{i,j}$ . The perturbed pairs  $\{(\mathbf{x}_{i,j}, \mathbf{y}_{i,j}^*)\}_{i,j}$  form the perturbed dataset  $\mathbb{D}'_{\text{inv,test}}$ . Details appear in Appendix F. We track two metrics, including an empirical Lipschitz estimate and reconstruction quality in PSNR (i.e., Peak Signal-to-Noise Ratio, higher PSNR means more accurate reconstruction, see appendix):

$$L_t := \max_{1 \leq i \leq 200} \max_{1 \leq j \leq 5} \frac{\|\mathbf{y}_t(\mathbf{x}_i) - \mathbf{y}_t(\mathbf{x}_{i,j})\|}{\|\mathbf{x}_i - \mathbf{x}_{i,j}\|}, \quad \text{and} \quad P_t(i, j) := \text{PSNR}(\mathbf{y}_t(\mathbf{x}_{i,j}), \mathbf{y}_{i,j}^*),$$

for  $1 \leq i \leq 200, 0 \leq j \leq 5$ , where  $j = 0$  means the original (unperturbed) sample,  $\mathbf{x}_{i,0} := \mathbf{x}_i, \mathbf{y}_{i,0}^* := \mathbf{y}_i^*$ . Here,  $L_t$  estimates how complex the  $t$ -th iterate map  $\mathbf{y}_t(\cdot)$  is, while  $P_t$  measures the reconstruction quality on *both* the original dataset  $\mathbb{D}_{\text{inv,test}}$  and the perturbed set  $\mathbb{D}'_{\text{inv,test}}$ .

**Experiment results.** *(i)* Results in Figure 2 support our theory. Figure 2a plots  $L_t$  versus  $t$ , while Figure 2b reports the mean  $\pm$  std of  $\{P_t(i, j)\}_{i,j}$  versus  $t$ . At  $t = 1$ , the mapping  $\mathbf{y}_1(\mathbf{x}) = \mathcal{G}_\Theta(\mathbf{0}, \mathbf{x})$  reflects a *single application of  $\mathcal{G}_\Theta$  and exhibits low Lipschitz constant*:  $L_1 = 0.140$  for PGD and  $L_1 = 0.436$  for HQS. As  $t$  increases,  $\mathbf{y}_t$  approaches the fixed point and  $L_t$  grows substantially, saturating around  $\approx 5.0$  for both models (Figure 2a). Meanwhile, the PSNR rises and stabilizes, indicating that  $\mathbf{y}_t(\mathbf{x})$  converges toward the ground truth (Figure 2b). Thus, the increase in  $L_t$  does not reflect divergence or instability; rather, *it captures the greater complexity of the underlying target mapping  $\mathbf{x} \mapsto \mathbf{y}_*$ , which is progressively expressed through iteration*. *(ii)* We also provide a comparison (both visually and quantitatively) to an explicit model in Figure 3. This baseline uses the identical DRUnet and is trained on the deblurring dataset with an end-to-end MSE loss. A visual inspection reveals that implicit models, particularly implicit HQS (6), produce sharper images with better-recovered textures and fewer artifacts than the explicit baseline. This perceptual advantage is corroborated by the quantitative metrics, where the DEQ-HQS model achieves a significant PSNR gain of over 2dB on average across the entire test set. *(iii)* Additional experiments showing a small implicit model outperforming larger explicit ones appear in Appendix F.

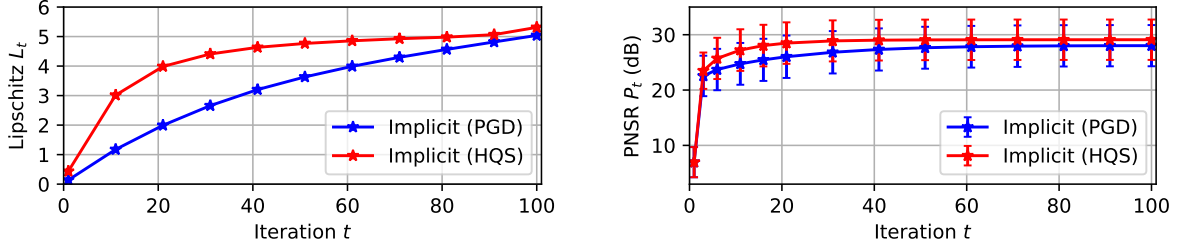
**3.2. Case Study 2: Scientific Computing.** The Navier-Stokes (NS) equations are foundational to computational fluid dynamics. We focus on the 2D steady-state incompressible case on a periodic domain  $\Omega := [0, 2\pi]^2$ :

$$(7) \quad (\mathbf{u} \cdot \nabla) \mathbf{u} + \nabla p = \nu \Delta \mathbf{u} + \mathbf{f}, \quad \nabla \cdot \mathbf{u} = 0 \quad \text{on } \Omega$$

where  $\mathbf{u} : \Omega \rightarrow \mathbb{R}^2$  is the velocity field,  $p : \Omega \rightarrow \mathbb{R}$  is the pressure,  $\nu > 0$  is the viscosity,  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^2$  is the external force. Solving NS equations refers to determining  $\mathbf{u}$  solves (7) given  $\mathbf{f}$ . Although global existence/smoothness of the solution given general forcings is famously open, classical results guarantee well-posedness under suitable conditions on  $\mathbf{f}$ .

**Theorem 3.6** ([91]). *There exists a constant  $c > 0$  depending only on  $\Omega$  such that, if  $\|\mathbf{f}\|_{L^2(\Omega)} \leq c\nu^2$ , then (7) admits a unique solution  $\mathbf{u}_*(\mathbf{f})$ . Let  $\mathbb{H}$  denote the space of admissible forcings<sup>1</sup>, and set*

<sup>1</sup>Details regarding the function spaces are provided in Appendix D.



(A) Empirical Lipschitz  $L_t$  of  $\mathbf{y}_t(\cdot)$  vs. iteration.  $L_t$  starts small at  $t=1$  and grows to a plateau ( $\sim 5$ ), indicating increasing expressivity of  $\mathbf{y}_t(\cdot)$ .

(B) Reconstruction quality  $P_t$  (mean  $\pm$  std over the original and perturbed test samples) increases and stabilizes:  $\mathbf{y}_t(\mathbf{x})$  converges toward the truth.

FIGURE 2. Validation on image deblurring. Iterating a simple operator  $\mathcal{G}_\Theta$  produces a complex fixed-point mapping: Lipschitz (a) grows, while accuracy (b) improves and stabilizes.

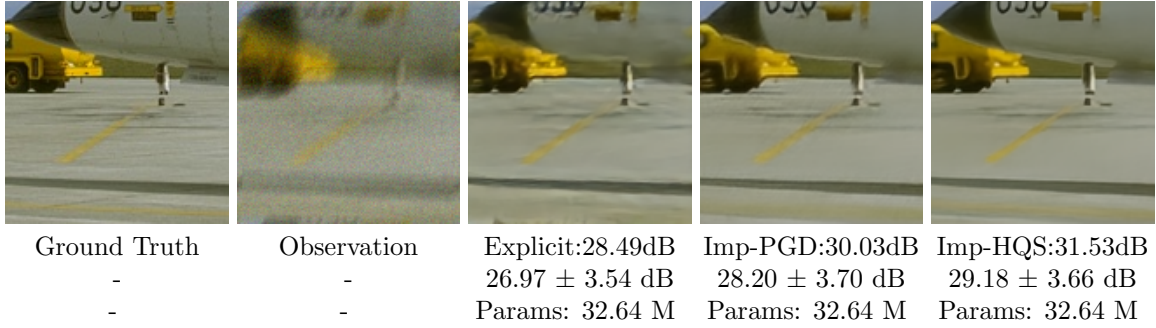


FIGURE 3. Visual results for deblurring. The top PSNR values (28.49, 30.03, or 31.53 dB) correspond to the single visualized image; the second line shows the average ( $\pm$  std) over all test samples.

$\mathbb{B}_\nu := \{f \in \mathbb{H} : \|f\|_{L^2(\Omega)} \leq c\nu^2\}$ . Then there exists a subset  $\mathbb{H}_\nu \subset \mathbb{B}_\nu$  that is dense in  $\mathbb{B}_\nu$ , on which the solution map  $f \mapsto u_*(f)$  is locally Lipschitz.

**Vorticity form.** Let  $\omega := \nabla \times u$  (and hence  $\omega_* := \nabla \times u_*$ ). Under periodic boundary and zero-mean conditions, one can recover the velocity  $u$  from vorticity  $\omega$  by solving a Poisson equation [64]. We hence focus on the solution map in vorticity:  $f \mapsto \omega_*$ .

While Theorem 3.6 gives a local Lipschitz result in function spaces, our expressivity results (Section 2) are stated for finite-dimensional spaces. To bridge this gap, we discretize the NS equations.

**Discretization.** Partition  $\Omega$  into  $N_h$  cells  $\Omega_h := \{C_i\}_{i=1}^{N_h}$  and define the cell-average restriction  $\mathcal{R}_h(f)|_C := \frac{1}{|C|} \int_C f(\xi) d\xi$  (similarly for  $\omega$ ). We work with the discrete forcings and vorticities:

$$\mathbf{x} := \mathcal{R}_h(f) \in \mathbb{R}^{N_h \times 2}, \quad \mathbf{y} := \mathcal{R}_h(\omega) \in \mathbb{R}^{N_h}$$

and aim to learn  $\mathbf{x} \mapsto \mathbf{y}_*$  where  $\mathbf{y}_* := \mathcal{R}_h(\omega_*)$  is the discrete solution in vorticity form. Back to the continuum setting, let the lifting operator  $\mathcal{E}_h$  be the piecewise-constant reconstruction  $\mathcal{E}_h(\mathbf{x}) := \sum_{C \in \Omega_h} x_C \mathbf{1}_C$ , and let  $\mathcal{P}$  be the orthogonal projection onto divergence-free, zero-mean fields.

**Corollary 3.7.** *The mapping  $\mathcal{F}_2 : \mathbf{x} \mapsto \mathbf{y}_*$  is locally Lipschitz continuous on  $\mathbb{X}_{\nu,h} := \{\mathbf{x} \in \mathbb{R}^{N_h \times 2} : \mathcal{P}(\mathcal{E}_h(\mathbf{x})) \in \mathbb{H}_\nu\}$ . Consequently, there exists a regular implicit operator  $\mathcal{G}(\mathbf{y}, \mathbf{x})$  satisfying  $\text{Fix}(\mathcal{G}(\cdot, \mathbf{x})) = \mathcal{F}_2(\mathbf{x})$  on any bounded subset of  $\mathbb{X}_{\nu,h}$ .*

The corollary instantiates our expressivity theory for steady-state NS, guaranteeing the existence of a *regular* implicit model  $\mathcal{G}$ . As in the image-reconstruction case, we now (i) choose a problem-specific parameterization of  $\mathcal{G}$  and (ii) verify our theory numerically on this architecture.

**Problem-tailored parameterization.** We use [67] as our code base. In particular,

$$(8) \quad \mathbf{z}_* = \mathcal{G}_\Theta(\mathbf{z}_*, \mathcal{Q}_\Phi(\mathbf{x})), \quad \mathbf{y}_* = \mathcal{Q}_\Psi(\mathbf{z}_*).$$

The core  $\mathcal{G}_\Theta$  is implemented as a Fourier Neural Operator (FNO) [57], and both the encoder  $\mathcal{Q}_\Phi$  and decoder  $\mathcal{Q}_\Psi$  use pointwise MLPs<sup>2</sup>. Details appear in Appendix G.

There is a growing literature on AI for scientific computing (e.g., [49, 56, 57, 63, 77, 78]; see [48, 50] for surveys). We adopt (8) because it learns solution *operators* rather than per-instance solutions, yielding discretization/mesh-resolution invariance, fast global mixing via Fourier layers, and strong parameter efficiency.

**Experiments.** We use the dataset of [67] with viscosity  $\nu = 0.01$ , which provides 4500 training pairs and 500 test pairs  $(\mathbf{x}, \mathbf{y}^*)$ , where  $\mathbf{x}$  is the discretized force and  $\mathbf{y}^*$  is the corresponding vorticity; we denote these sets by  $\mathbb{D}_{\text{pde,train}}$  and  $\mathbb{D}_{\text{pde,test}}$ . Details are given in Appendix G.

We test iteration-wise behavior for 50 steps starting from  $\mathbf{z}_0 = \mathbf{0}$ :  $\mathbf{z}_{t+1} = \mathcal{G}_\Theta(\mathbf{z}_t, \mathcal{Q}_\Phi(\mathbf{x}))$  for  $0 \leq t \leq 49$ , and  $\mathbf{y}_t(\mathbf{x}) = \mathcal{Q}_\Psi(\mathbf{z}_t)$ . Analogous to the inverse-problem study, we augment the test set with perturbations. For each  $(\mathbf{x}_i, \mathbf{y}_i^*) \in \mathbb{D}_{\text{pde,test}}$ , we construct 15 perturbed vorticities  $\{\mathbf{y}_{i,j}^*\}_{j=1}^{15}$ ; we then compute compatible forces  $\{\mathbf{x}_{i,j}\}_{j=1}^{15}$  by evaluating the NS operator (see Appendix G for details). The perturbed test set is  $\mathbb{D}'_{\text{pde,test}} = \{(\mathbf{x}_{i,j}, \mathbf{y}_{i,j}^*) : 1 \leq i \leq 500, 1 \leq j \leq 15\}$ . Across iterations we report an empirical Lipschitz estimate  $L_t$  and relative reconstruction error  $E_t$ :

$$L_t := \max_{1 \leq i \leq 500} \max_{1 \leq j \leq 15} \frac{\|\mathbf{y}_t(\mathbf{x}_i) - \mathbf{y}_t(\mathbf{x}_{i,j})\|}{\|\mathbf{x}_i - \mathbf{x}_{i,j}\|}, \quad \text{and} \quad E_t(i, j) := \frac{\|\mathbf{y}_t(\mathbf{x}_{i,j}) - \mathbf{y}_{i,j}^*\|}{\|\mathbf{y}_{i,j}^*\| + \epsilon},$$

for  $1 \leq i \leq 500, 0 \leq j \leq 15$ , where  $j = 0$  means the original (unperturbed) sample,  $\mathbf{x}_{i,0} := \mathbf{x}_i, \mathbf{y}_{i,0}^* := \mathbf{y}_i^*$ . Therefore,  $E_t$  evaluates accuracy on both  $\mathbb{D}_{\text{pde,test}}$  and  $\mathbb{D}'_{\text{pde,test}}$ .

The results in Figure 4 align with our theory. At  $t = 1$ , the mapping  $\mathbf{y}_1(\mathbf{x})$  reflects a single application of  $\mathcal{G}_\Theta$  and exhibits low Lipschitz constant:  $L_1 = 23.1$ . As iterations proceed toward the fixed point, the complexity grows markedly:  $L_t$  increases to  $\approx 367$  by  $t = 50$  (Figure 4a). Meanwhile, the relative error  $E_t$  decreases monotonically and stabilizes at  $0.078 \pm 0.028$  (Figure 4b), indicating convergence to a good approximation of  $\mathbf{y}_*$ . Thus, *the learned operator  $\mathcal{G}_\Theta$  is simple (Lipschitz in  $\mathbf{x}$ ), while additional test-time iterations let  $\mathbf{y}_t$  realize progressively more complex mappings*. In addition, a comparison with an explicit baseline (vanilla FNO) in Figure 5 shows the implicit model produces more accurate solutions, both visually and quantitatively. Additional experiments showing a small implicit model outperforming larger explicit ones appear in Appendix G.

**3.3. Case Study 3: Operations Research.** Linear programming (LP) is fundamental to operations research, of which a general form is given by

$$(9) \quad \min_{\mathbf{y} \in \mathbb{R}^n} \mathbf{c}^\top \mathbf{y}, \quad \text{s.t.} \quad \mathbf{A}\mathbf{y} \circ \mathbf{b}, \quad \mathbf{l} \leq \mathbf{y} \leq \mathbf{u}.$$

Here,  $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m, \mathbf{c} \in \mathbb{R}^n, \mathbf{l} \in \mathbb{R}^n, \mathbf{u} \in \mathbb{R}^n$ , and  $\circ \in \{=, \leq\}^m$  denotes componentwise relations, i.e., each  $\circ_i \in \{=, \leq\}$  specifies whether  $(\mathbf{A}\mathbf{y})_i$  equals or is bounded above by  $b_i$ . Let  $\mathbf{x} := (\mathbf{A}, \mathbf{b}, \mathbf{c}, \circ, \mathbf{l}, \mathbf{u})$  as the input that describes the LP in (9). To define the solution mapping  $\mathcal{F}_3$  that maps  $\mathbf{x}$  to the solution of LP, we require feasibility and boundedness (which ensure an optimal solution [6]). Accordingly, let

$$\mathbb{X} := \{(\mathbf{A}, \mathbf{b}, \mathbf{c}, \circ, \mathbf{l}, \mathbf{u}) : \text{The resulting LP is feasible and bounded}\}$$

Within  $\mathbb{X}$ , there are some LPs where the solution mapping is not single-valued or not continuous. By excluding these LPs, it forms a subset  $\mathbb{X}_{\text{sub}} \subset \mathbb{X}$  on which  $\mathcal{F}_3$  is single-valued and locally Lipschitz. The strict definition of  $\mathbb{X}_{\text{sub}}$  and the proof of Theorem 3.8 are provided in Appendix E.

<sup>2</sup>Introducing additional encoder and decoder is common in practice. Compared to the vanilla formulation  $\mathbf{y}_* = \mathcal{G}(\mathbf{y}_*, \mathbf{x})$ , it does not change our expressivity results in Section 2. Details appear in Appendix B.

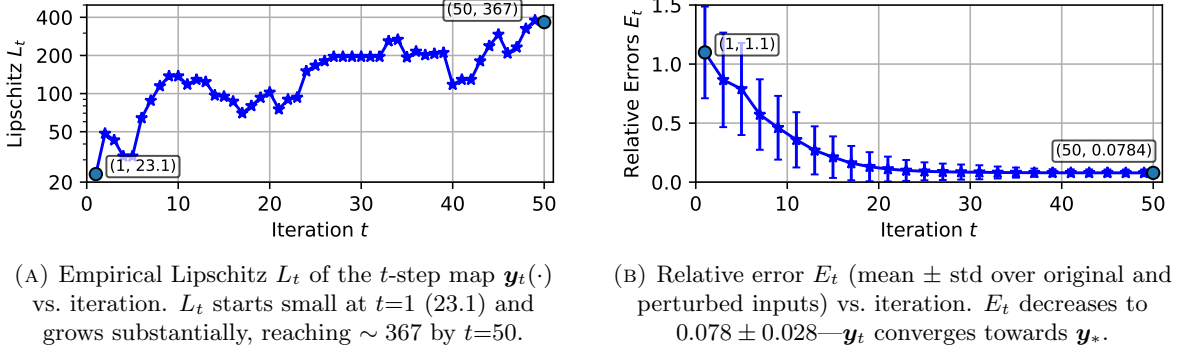


FIGURE 4. Validation on the steady Navier–Stokes task. Iterating a simple operator  $\mathcal{G}_\Theta$  yields a complex fixed-point mapping: Lipschitz constant (a) increases, while error (b) decreases.

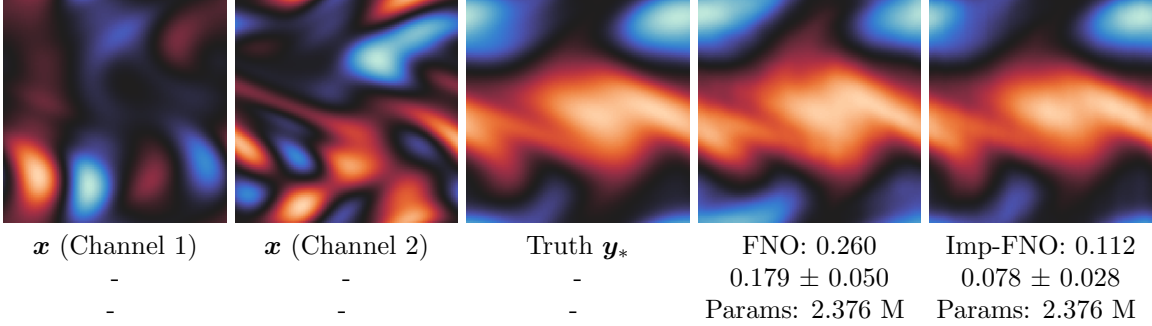


FIGURE 5. Visual results for NS equations. The top value (0.260 and 0.112) corresponds to the single visualized sample; the second line shows the average relative error ( $\pm$  std) over all test samples.

**Theorem 3.8.** *There is a subset  $\mathbb{X}_{\text{sub}} \subset \mathbb{X}$  that is dense in  $\mathbb{X}$ , on which each LP admits a unique solution  $\mathbf{y}_*$ , and the solution map  $\mathcal{F}_3 : \mathbf{x} \mapsto \mathbf{y}_*$  is locally Lipschitz continuous on  $\mathbb{X}_{\text{sub}}$ .*

**Corollary 3.9.** *There must be a regular implicit model  $\mathcal{G}(\mathbf{y}, \mathbf{x})$  that satisfies  $\text{Fix}(\mathcal{G}(\cdot, \mathbf{x})) = \mathcal{F}_3(\mathbf{x})$  on any bounded subsets of  $\mathbb{X}_{\text{sub}}$ .*

Corollary 3.9 follows immediately from Theorems 2.4 and 3.8. It indicates the existence of implicit models with desired properties that solves LP. As in the previous case studies, we now (i) choose a problem-specific parameterization of  $\mathcal{G}$  and (ii) verify the theory numerically on this architecture.

**Implicit GNN parameterization.** We model the implicit operator  $\mathcal{G}$  for LP with a *graph neural network* (GNN). First, express an LP instance  $\mathbf{x} = (\mathbf{A}, \mathbf{b}, \mathbf{c}, \mathbf{o}, \mathbf{l}, \mathbf{u})$  as a bipartite graph (Figure 6). We create  $n$  variable nodes  $\{V_j\}_{j=1}^n$  and  $m$  constraint nodes  $\{W_i\}_{i=1}^m$ . Node features collect the data of the LP: each  $V_j$  stores  $(c_j, l_j, u_j)$ ; each  $W_i$  stores  $(b_i, o_i)$ . We connect  $W_i$  to  $V_j$  if  $A_{ij} \neq 0$ , and place  $A_{ij}$  on that edge as its feature. Given this representation, an (explicit) GNN can map the LP to a solution, i.e.,  $\mathbf{y}_* = \text{GNN}(\mathbf{x})$  where  $\mathbf{x}$  denotes the graph-encoded LP. This approach was proposed in [30] in the context of mixed-integer linear programs, and [17] subsequently showed that (explicit) GNNs offer a universal framework for representing LPs. Built on this, we propose an *implicit* GNN:

$$(10) \quad \mathbf{z}_* = \mathcal{G}_\Theta(\mathbf{z}_*, \mathcal{Q}_\Phi(\mathbf{x})), \quad \mathbf{y}_* = \mathcal{Q}_\Psi(\mathbf{z}_*)$$

where  $\mathcal{G}_\Theta$  is the core GNN,  $\mathcal{Q}_\Phi$  encodes instance-specific (static) features from  $\mathbf{x}$ , and  $\mathcal{Q}_\Psi$  decodes per-variable states to the solution. Both  $\mathcal{Q}_\Phi$  and  $\mathcal{Q}_\Psi$  are small MLPs shared across all nodes. At inference, we repeatedly call  $\mathcal{G}_\Theta$  with initialization  $\mathbf{z}_0 = \mathbf{0}$ :  $\mathbf{z}_t = \mathcal{G}_\Theta(\mathbf{z}_{t-1}, \mathcal{Q}_\Phi(\mathbf{x}))$  for  $t = 1, 2, \dots, T$ ,

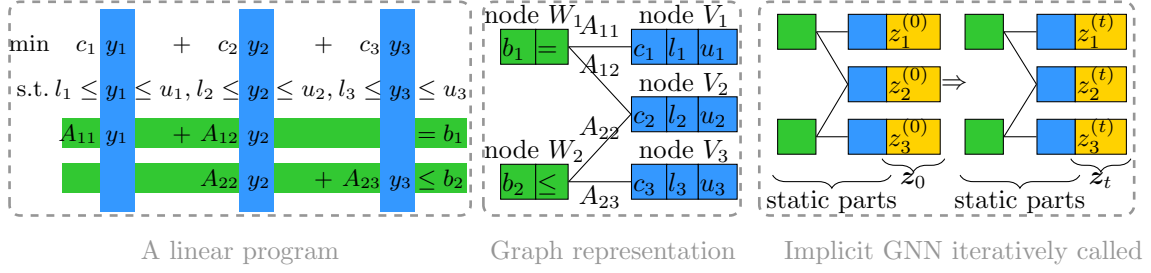


FIGURE 6. The graph representation of LP and implicit GNN applied on this graph

and finally output  $\mathbf{y}_t = \mathcal{Q}_\Psi(\mathbf{z}_t)$ . Relative to prior work, our only modification is to attach to each variable node an additional *dynamic* state  $z_j^{(t)} \in \mathbb{R}$ . Details appear in Appendix H.

There is a rapidly growing literature on implicit GNNs with diverse applications and theories [4, 11, 12, 37, 58, 69, 72, 100, 103]. Our LP case study is complementary to this line of work: rather than adopting a particular implicit-GNN architecture, we start from a standard explicit GNN for LP and convert it into a fixed-point formulation tailored to linear programs.

In addition, recent work on ML for LP has progressed quickly (e.g., [17, 26, 51, 54, 55, 62, 76, 84]); for broader context, see overviews on Learning to Optimize [16] and AI for Operations Research [25]. Our goal here is not to outperform all state-of-the-arts; rather, we aim to establish a foundational point: converting a standard LP-GNN into an implicit form yields the benefits predicted by our theory—implicit models can be simple at the update level yet expressive at the fixed-point map, and their performance improves predictably with test-time iteration.

**Experiments.** We sample LP instances  $\mathbf{x} = (\mathbf{A}, \mathbf{b}, \mathbf{c}, \circ, \mathbf{l}, \mathbf{u})$ , solve it to obtain an optimal solution  $\mathbf{y}_*$ , and form 2,500 training pairs and 1,000 test pairs like  $(\mathbf{x}, \mathbf{y}_*)$ , denoted  $\mathbb{D}_{\text{LP,train}}$  and  $\mathbb{D}_{\text{LP,test}}$ . We also create five perturbed test sets  $\{\mathbb{D}_{\text{LP,test}}^{(j)}\}_{j=1}^5$  by altering exactly one block among  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, \mathbf{l}, \text{ or } \mathbf{u})$ . For each  $(\mathbf{x}_i, \mathbf{y}_i^*) \in \mathbb{D}_{\text{LP,test}}$  and each perturbation type  $j$ , we form a perturbed instance  $\mathbf{x}_{i,j}$ , solve it to obtain  $\mathbf{y}_{i,j}^*$ , and collect  $\mathbb{D}_{\text{LP,test}}^{(j)} = \{(\mathbf{x}_{i,j}, \mathbf{y}_{i,j}^*)\}_{i=1}^{1000}$ . Details in Appendix H. We report:

$$L_t(j) := \max_{1 \leq i \leq 1000} \frac{\|\mathbf{y}_t(\mathbf{x}_i) - \mathbf{y}_t(\mathbf{x}_{i,j})\|}{\|\mathbf{x}_i - \mathbf{x}_{i,j}\|}, \quad \text{and} \quad E_t(i, j) := \frac{\|\mathbf{y}_t(\mathbf{x}_{i,j}) - \mathbf{y}_{i,j}^*\|}{\|\mathbf{y}_{i,j}^*\| + \epsilon},$$

for  $1 \leq i \leq 1000, 0 \leq j \leq 5$ , where  $j = 0$  denotes the unperturbed pair  $(\mathbf{x}_{i,0}, \mathbf{y}_{i,0}^*) := (\mathbf{x}_i, \mathbf{y}_i^*)$ .

Results support our theory. **(i)** Figure 7a plots the five curves  $L_t(j)$  (one for each perturbation type). At  $t = 1$ , a single application of (10) yields relatively small empirical Lipschitz constants for *all* perturbation modes. As iterations proceed toward the fixed point, Lipschitz constants grow markedly. **(ii)** Figure 7b reports the mean $\pm$ std of  $E_t(i, j)$ :  $E_t$  decreases and stabilizes at 0.146, indicating that the growth of  $L_t$  reflects the higher intrinsic complexity of the solution mapping  $\mathbf{y}_*(\mathbf{x})$  rather than divergence or instability. **(iii)** Table 1 contrasts implicit and explicit GNNs. At matched embedding sizes, implicit GNNs match or beat explicit ones—most clearly at small/mid sizes (4/8/16). In addition, a smaller implicit model can outperform a larger explicit model on training error. For example, implicit-4 vs. explicit-8 (0.203 vs. 0.233) and implicit-8 vs. explicit-16 (0.162 vs. 0.183). This supports our theory that iterating a simple implicit operator can yield strong expressivity.

**Discussion on generalization.** While generalization is not our main focus, a trend in Table 1 is informative: explicit GNNs improve as width increases from 4 to 8 but then *overfit* (test error significantly rises at 16/32), whereas implicit GNNs improve from 4 to 8 to 16 and only tick up slightly at 32. We attribute this to: (i) LP constraints  $\mathbf{A}\mathbf{y} \circ \mathbf{b}$  in (9) are specified implicitly rather than as an explicit set; implicit models align naturally with such a structure, and (ii) while fixed-point maps  $\mathbf{y}_*(\mathbf{x})$  can be sensitive to inputs  $\mathbf{x}$ , the implicit formulation allows us realize them via a simpler, smaller operator  $\mathcal{G}$ , which “implicitly” regularizes training and support good generalization in practice.

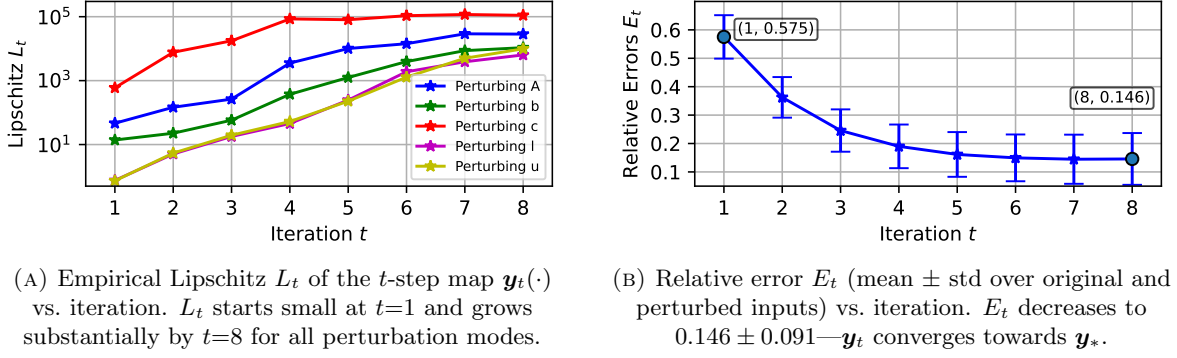


FIGURE 7. Numerical validation on the linear-program task.

TABLE 1. Comparison between explicit GNNs and implicit GNNs on the LP task.

	Emb. size	4	8	16	32
Exp-GNNs	# Params.	580	2,088	7,888	30,624
	Err (Train)	$0.387 \pm 0.103$	$0.233 \pm 0.084$	$0.183 \pm 0.070$	$0.112 \pm 0.049$
	Err (Test)	$0.397 \pm 0.107$	$0.273 \pm 0.104$	$0.283 \pm 0.111$	$0.318 \pm 0.122$
	Emb. size	4	8	16	32
Imp-GNNs	# Params.	722	2,350	8,390	31,606
	Err (Train)	$0.203 \pm 0.107$	$0.162 \pm 0.094$	$0.131 \pm 0.080$	$0.118 \pm 0.073$
	Err (Test)	$0.218 \pm 0.117$	$0.177 \pm 0.105$	$0.152 \pm 0.098$	$0.156 \pm 0.109$

#### 4. CONCLUSIONS

We establish a sharp expressivity boundary for regular implicit models: as iterations increase, their expressive power grows, and the resulting fixed points can represent exactly the class of locally Lipschitz maps. In three case studies, per-iteration Lipschitz estimates grew toward the target’s complexity while accuracy improved and stabilized. Overall, both theory and evidence show that iterating a simple operator is a principled route to powerful models, clarifying how fixed-point architectures can match or surpass large explicit networks.

#### REFERENCES

- [1] Eddie Aamari, Jisu Kim, Frédéric Chazal, Bertrand Michel, Alessandro Rinaldo, and Larry Wasserman, *Estimating the reach of a manifold*, Electronic Journal of Statistics **13** (2019), no. 1, 1359–1399.
- [2] Iskander Azangulov, George Deligiannidis, and Judith Rousseau, *Convergence of diffusion models under the manifold hypothesis in high-dimensions*, arXiv preprint arXiv:2409.18804 (2024).
- [3] Shaojie Bai, J Zico Kolter, and Vladlen Koltun, *Deep equilibrium models*, Advances in neural information processing systems **32** (2019).
- [4] Justin Baker, Qingsong Wang, Cory D Hauck, and Bao Wang, *Implicit graph neural networks: A monotone operator viewpoint*, International conference on machine learning, 2023, pp. 1521–1548.
- [5] Richard G Baraniuk and Michael B Wakin, *Random projections of smooth manifolds*, Foundations of computational mathematics **9** (2009), no. 1, 51–77.
- [6] Dimitris Bertsimas and John N Tsitsiklis, *Introduction to linear optimization*, Vol. 6, Athena scientific Belmont, MA, 1997.
- [7] Stephen P Boyd and Lieven Vandenbergh, *Convex optimization*, Cambridge university press, 2004.
- [8] Gregory T Buzzard, Stanley H Chan, Suhas Sreehari, and Charles A Bouman, *Plug-and-play unplugged: Optimization-free reconstruction using consensus equilibrium*, SIAM Journal on Imaging Sciences **11** (2018), no. 3, 2001–2020.
- [9] Emmanuel J Candes and Terence Tao, *Near-optimal signal recovery from random projections: Universal encoding strategies?*, IEEE transactions on information theory **52** (2006), no. 12, 5406–5425.
- [10] Stanley H Chan, Xiran Wang, and Omar A Elgandy, *Plug-and-play admm for image restoration: Fixed-point convergence and applications*, IEEE Transactions on Computational Imaging **3** (2016), no. 1, 84–98.

- [11] Qi Chen, Yifei Wang, Yisen Wang, Jianlong Chang, Qi Tian, Jiansheng Yang, and Zhouchen Lin, *Efficient and scalable implicit graph neural networks with virtual equilibrium*, 2022 IEEE International Conference on Big Data (Big Data), 2022, pp. 864–873.
- [12] Qi Chen, Yifei Wang, Yisen Wang, Jiansheng Yang, and Zhouchen Lin, *Optimization-induced graph implicit nonlinear diffusion*, International conference on machine learning, 2022, pp. 3648–3661.
- [13] Tianlong Chen, Xiaohan Chen, Wuyang Chen, Howard Heaton, Jialin Liu, Zhangyang Wang, and Wotao Yin, *Learning to optimize: A primer and a benchmark*, Journal of Machine Learning Research **23** (2022), no. 189, 1–59.
- [14] Tianlong Chen, Weiyi Zhang, Zhou Jingyang, Shiyu Chang, Sijia Liu, Lisa Amini, and Zhangyang Wang, *Training stronger baselines for learning to optimize*, Advances in Neural Information Processing Systems **33** (2020), 7332–7343.
- [15] Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin, *Theoretical linear convergence of unfolded ista and its practical weights and thresholds*, Advances in Neural Information Processing Systems **31** (2018).
- [16] Xiaohan Chen, Jialin Liu, and Wotao Yin, *Learning to optimize: A tutorial for continuous and mixed-integer optimization*, Science China Mathematics **67** (2024), no. 6, 1191–1262.
- [17] Ziang Chen, Jialin Liu, Xinshang Wang, and Wotao Yin, *On representing linear programs by graph neural networks*, The eleventh international conference on learning representations, 2023.
- [18] Kenneth L Clarkson, *Tighter bounds for random projections of manifolds*, Proceedings of the twenty-fourth annual symposium on computational geometry, 2008, pp. 39–48.
- [19] Christian Daniele, Silvia Villa, Samuel Vaiter, and Luca Calatroni, *Deep equilibrium models for poisson imaging inverse problems via mirror descent*, arXiv preprint arXiv:2507.11461 (2025).
- [20] Manfredo P Do Carmo, *Differential geometry of curves and surfaces: revised and updated second edition*, Courier Dover Publications, 2016.
- [21] David L Donoho and Carrie Grimes, *Image manifolds which are isometric to euclidean space*, Journal of mathematical imaging and vision **23** (2005), no. 1, 5–24.
- [22] Asen L Dontchev and R Tyrrell Rockafellar, *Characterizations of strong regularity for variational inequalities over polyhedral convex sets*, SIAM Journal on Optimization **6** (1996), no. 4, 1087–1105.
- [23] ———, *Implicit functions and solution mappings*, Vol. 543, Springer, 2009.
- [24] Laurent El Ghaoui, Fangda Gu, Bertrand Travacca, Armin Askari, and Alicia Tsai, *Implicit deep learning*, SIAM Journal on Mathematics of Data Science **3** (2021), no. 3, 930–958.
- [25] Zhenan Fan, Bissan Ghaddar, Xinglu Wang, Linzi Xing, Yong Zhang, and Zirui Zhou, *Artificial intelligence for operations research: Revolutionizing the operations research process*, arXiv preprint arXiv:2401.03244 (2024).
- [26] Zhenan Fan, Xinglu Wang, Oleksandr Yakovenko, Abdullah Ali Sivas, Owen Ren, Yong Zhang, and Zirui Zhou, *Smart initial basis selection for linear programs*, International conference on machine learning, 2023, pp. 9650–9664.
- [27] Herbert Federer, *Curvature measures*, Transactions of the American Mathematical Society **93** (1959), no. 3, 418–491.
- [28] Gerald B. Folland, *Advanced calculus*, 2nd ed., University of Washington, Seattle, WA, 2023. Freely available at <https://sites.math.washington.edu/~folland/AdvCalc/>.
- [29] Samy Wu Fung, Howard Heaton, Qiuwei Li, Daniel McKenzie, Stanley Osher, and Wotao Yin, *Jfb: Jacobian-free backpropagation for implicit networks*, Proceedings of the aaai conference on artificial intelligence, 2022.
- [30] Maxime Gasse, Didier Chételat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi, *Exact combinatorial optimization with graph convolutional neural networks*, Advances in neural information processing systems **32** (2019).
- [31] Ruturaj G Gavaskar and Kunal N Chaudhury, *Plug-and-play ista converges with kernel denoisers*, IEEE Signal Processing Letters **27** (2020), 610–614.
- [32] Jonas Geiping, Sean McLeish, Neel Jain, John Kirchenbauer, Siddharth Singh, Brian R Bartoldson, Bhavya Kaikhura, Abhinav Bhatele, and Tom Goldstein, *Scaling up test-time compute with latent reasoning: A recurrent depth approach*, arXiv preprint arXiv:2502.05171 (2025).
- [33] Zhengyang Geng, Ashwini Pople, and J Zico Kolter, *One-step diffusion distillation via deep equilibrium models*, Advances in Neural Information Processing Systems **36** (2023), 41914–41931.
- [34] Zhengyang Geng, Xin-Yu Zhang, Shaojie Bai, Yisen Wang, and Zhouchen Lin, *On training implicit models*, Advances in Neural Information Processing Systems **34** (2021), 24247–24260.
- [35] Davis Gilton, Gregory Ongie, and Rebecca Willett, *Deep equilibrium architectures for inverse problems in imaging*, IEEE Transactions on Computational Imaging **7** (2021), 1123–1133.
- [36] Alexandros Gkillas, Dimitris Ampeliotis, and Kostas Berberidis, *An optimization-based deep equilibrium model for hyperspectral image deconvolution with convergence guarantees*, arXiv preprint arXiv:2306.06378 (2023).
- [37] Fangda Gu, Heng Chang, Wenwu Zhu, Somayeh Sojoudi, and Laurent El Ghaoui, *Implicit graph neural networks*, Advances in neural information processing systems **33** (2020), 11984–11995.
- [38] Alper Güngör, Baris Askin, Damla Alptekin Soydan, Can Barış Top, Emine Ulku Saritas, and Tolga Cukur, *Deq-mpi: A deep equilibrium reconstruction with learned consistency for magnetic particle imaging*, IEEE Transactions on Medical Imaging **43** (2023), no. 1, 321–334.

- [39] Aaron Havens, Alexandre Araujo, Siddharth Garg, Farshad Khorrani, and Bin Hu, *Exploiting connections between lipschitz structures for certifiably robust deep equilibrium models*, Advances in Neural Information Processing Systems **36** (2023), 21658–21674.
- [40] Yutong He, Qiulin Shang, Xinneng Huang, Jialin Liu, and Kun Yuan, *A mathematics-inspired learning-to-optimize framework for decentralized optimization*, arXiv preprint arXiv:2410.01700 (2024).
- [41] Chinmay Hegde and Richard G Baraniuk, *Signal recovery on incoherent manifolds*, IEEE Transactions on Information Theory **58** (2012), no. 12, 7204–7214.
- [42] Samuel Hurault, Arthur Leclaire, and Nicolas Papadakis, *Gradient step denoiser for convergent plug-and-play*, International conference on learning representations, 2022.
- [43] ———, *Proximal denoiser for convergent plug-and-play optimization with nonconvex regularization*, International conference on machine learning, 2022, pp. 9483–9505.
- [44] Mark A Iwen and Mauro Maggioni, *Approximation of points on low-dimensional manifolds via random linear projections*, Information and Inference: A Journal of the IMA **2** (2013), no. 1, 1–31.
- [45] Saber Jafarpour, Alexander Davydov, Anton Proskurnikov, and Francesco Bullo, *Robust implicit networks via non-euclidean contractions*, Advances in Neural Information Processing Systems **34** (2021), 9857–9868.
- [46] Ulugbek S Kamilov, Charles A Bouman, Gregory T Buzzard, and Brendt Wohlberg, *Plug-and-play methods for integrating physical and learned models in computational imaging: Theory, algorithms, and applications*, IEEE Signal Processing Magazine **40** (2023), no. 1, 85–97.
- [47] Ulugbek S Kamilov, Hassan Mansour, and Brendt Wohlberg, *A plug-and-play priors approach for solving nonlinear imaging inverse problems*, IEEE Signal Processing Letters **24** (2017), no. 12, 1872–1876.
- [48] George Em Karniadakis, Ioannis G. Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang, *Physics-informed machine learning*, Nature Reviews Physics **3** (2021), no. 6, 422–440.
- [49] Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar, *Neural operator: Learning maps between function spaces with applications to pdes*, Journal of Machine Learning Research **24** (2023), 1–97.
- [50] Nikola B. Kovachki, Samuel Lanthaler, and Andrew M. Stuart, *Operator learning: Algorithms and analysis*, arXiv preprint arXiv:2402.15715 (2024).
- [51] Yufei Kuang, Xijun Li, Jie Wang, Fangzhou Zhu, Meng Lu, Zhihai Wang, Jia Zeng, Houqiang Li, Yongdong Zhang, and Feng Wu, *Accelerate presolve in large-scale linear programming via reinforcement learning*, IEEE Transactions on Pattern Analysis and Machine Intelligence (2025).
- [52] Gunther Leobacher and Alexander Steinicke, *Existence, uniqueness and regularity of the projection onto differentiable manifolds*, Annals of global analysis and geometry **60** (2021), no. 3, 559–587.
- [53] Anat Levin, Yair Weiss, Frédo Durand, and William T. Freeman, *Understanding and evaluating blind deconvolution algorithms*, 2009 IEEE conference on computer vision and pattern recognition (cvpr), 2009, pp. 1964–1971.
- [54] Bingheng Li, Linxin Yang, Yupeng Chen, Senmiao Wang, Haitao Mao, Qian Chen, Yao Ma, Akang Wang, Tian Ding, Jiliang Tang, and Ruoyu Sun, *PDHG-unrolled learning-to-optimize method for large-scale linear programming*, Forty-first international conference on machine learning, 2024.
- [55] Qian Li, Tian Ding, Linxin Yang, Minghui Ouyang, Qingjiang Shi, and Ruoyu Sun, *On the power of small-size graph neural networks for linear programming*, The thirty-eighth annual conference on neural information processing systems, 2024.
- [56] Zongyi Li, Daniel Zhengyu Huang, Burigede Liu, and Anima Anandkumar, *Fourier neural operator with learned deformations for pdes on general geometries*, Journal of Machine Learning Research **24** (2023), no. 388, 1–26.
- [57] Zongyi Li, Nikola Borislavov Kovachki, Kamyar Azizzadenesheli, Burigede liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar, *Fourier neural operator for parametric partial differential equations*, International conference on learning representations, 2021.
- [58] Junchao Lin, Zenan Ling, Zhanbo Feng, Jingwen Xu, Minxuan Liao, Feng Zhou, Tianqi Hou, Zhenyu Liao, and Robert C Qiu, *Ignn-solver: A graph neural solver for implicit graph neural networks*, arXiv preprint arXiv:2410.08524 (2024).
- [59] Jialin Liu, Xiaohan Chen, Zhangyang Wang, and Wotao Yin, *Alista: Analytic weights are as good as learned weights in lista*, International conference on learning representations (iclr), 2019.
- [60] Jialin Liu, Xiaohan Chen, Zhangyang Wang, Wotao Yin, and HanQin Cai, *Towards constituting mathematical structures for learning to optimize*, International conference on machine learning, 2023, pp. 21426–21449.
- [61] Jiaming Liu, Salman Asif, Brendt Wohlberg, and Ulugbek Kamilov, *Recovery analysis for plug-and-play priors using the restricted eigenvalue condition*, Advances in Neural Information Processing Systems **34** (2021), 5921–5933.
- [62] Tianhao Liu, Shanwen Pu, Dongdong Ge, and Yinyu Ye, *Learning to pivot as a smart expert*, Proceedings of the aaai conference on artificial intelligence, 2024, pp. 8073–8081.
- [63] Lu Lu, Pengzhan Jin, Guofei Pang, Zhongqiang Zhang, and George Em Karniadakis, *Learning nonlinear operators via deepnet based on the universal approximation theorem of operators*, Nature Machine Intelligence **3** (2021), no. 3, 218–229.
- [64] Andrew J Majda, Andrea L Bertozzi, and A Ogawa, *Vorticity and incompressible flow. cambridge texts in applied mathematics*, Appl. Mech. Rev. **55** (2002), no. 4, B77–B78.

- [65] Olvi L Mangasarian and T-H Shiau, *Lipschitz continuity of solutions of linear inequalities, programs and complementarity problems*, SIAM Journal on Control and Optimization **25** (1987), no. 3, 583–595.
- [66] D. Martin, C. Fowlkes, D. Tal, and J. Malik, *A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics*, Proc. 8th int’l conf. computer vision, 2001 July, pp. 416–423.
- [67] Tanya Marwah, Ashwini Pople, J Zico Kolter, Zachary Lipton, Jianfeng Lu, and Andrej Risteski, *Deep equilibrium based neural operators for steady-state pdes*, Advances in Neural Information Processing Systems **36** (2023), 15716–15737.
- [68] Florent Nacry and Lionel Thibault, *Distance function associated to a prox-regular set*, Set-Valued and Variational Analysis (2022), 1–20.
- [69] Matthieu Nastorg, Michele-Alessandro Bucci, Thibault Faney, Jean-Marc Gratien, Guillaume Charpiat, and Marc Schoenauer, *An implicit gnn solver for poisson-like problems*, Computers & Mathematics with Applications **176** (2024), 270–288.
- [70] Alex Oshin, Hassan Almubarak, and Evangelos A. Theodorou, *Differentiable robust model predictive control*, Robotics: Science and systems (rss), 2024.
- [71] Chirag Pabbaraju, Ezra Winston, and J Zico Kolter, *Estimating lipschitz constants of monotone deep equilibrium models*, International conference on learning representations, 2021.
- [72] Junyoung Park, Jinhyun Choo, and Jinkyoo Park, *Convergent graph solvers*, International conference on learning representations, 2022.
- [73] Ashwini Pople, Zhengyang Geng, and J Zico Kolter, *Deep equilibrium approaches to diffusion models*, Advances in Neural Information Processing Systems **35** (2022), 37975–37990.
- [74] René Poliquin, R Rockafellar, and Lionel Thibault, *Local differentiability of distance functions*, Transactions of the American mathematical Society **352** (2000), no. 11, 5231–5249.
- [75] Peter Potapchik, Iskander Azangulov, and George Deligiannidis, *Linear convergence of diffusion models under the manifold hypothesis*, arXiv preprint arXiv:2410.09046 (2024).
- [76] Chendi Qian, Didier Chételat, and Christopher Morris, *Exploring the power of graph neural networks in solving linear optimization problems*, International conference on artificial intelligence and statistics, 2024, pp. 1432–1440.
- [77] Md Ashiqur Rahman, Zachary E. Ross, and Kamyar Azizzadenesheli, *U-no: U-shaped neural operators*, arXiv preprint arXiv:2204.11127 (2022).
- [78] Maziar Raissi, Paris Perdikaris, and George Em Karniadakis, *Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations*, Journal of Computational Physics **378** (2019), 686–707.
- [79] Arash Rasti-Meymandi, Aboozar Ghaffari, and Emad Fatemizadeh, *Plug and play augmented hqs: Convergence analysis and its application in mri reconstruction*, Neurocomputing **518** (2023), 1–14.
- [80] Max Revay, Ruigang Wang, and Ian R Manchester, *Lipschitz bounded equilibrium networks*, arXiv preprint arXiv:2010.01732 (2020).
- [81] Stephen M Robinson, *Strongly regular generalized equations*, Mathematics of Operations Research **5** (1980), no. 1, 43–62.
- [82] Sam T Roweis and Lawrence K Saul, *Nonlinear dimensionality reduction by locally linear embedding*, science **290** (2000), no. 5500, 2323–2326.
- [83] Ernest Ryu, Jialin Liu, Sicheng Wang, Xiaohan Chen, Zhangyang Wang, and Wotao Yin, *Plug-and-play methods provably converge with properly trained denoisers*, International conference on machine learning, 2019, pp. 5546–5557.
- [84] Shinsaku Sakaue and Taihei Oki, *Generalization bound and learning methods for data-driven projections in linear programming*, Advances in Neural Information Processing Systems **37** (2024), 12825–12846.
- [85] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini, *The graph neural network model*, IEEE transactions on neural networks **20** (2008), no. 1, 61–80.
- [86] Vineet R Shenoy, Suhas Lohit, Hassan Mansour, Rama Chellappa, and Tim K Marks, *Recovering pulse waves from video using deep unrolling and deep equilibrium models*, arXiv preprint arXiv:2503.17269 (2025).
- [87] Yu Sun, Brendt Wohlberg, and Ulugbek S Kamilov, *An online plug-and-play algorithm for regularized image reconstruction*, IEEE Transactions on Computational Imaging **5** (2019), no. 3, 395–408.
- [88] Julián Tachella, Matthieu Terris, Samuel Hurault, Andrew Wang, Dongdong Chen, Minh-Hai Nguyen, Maxime Song, Thomas Davies, Leo Davy, Jonathan Dong, Paul Escande, Johannes Hertrich, Zhiyuan Hu, Tobías I. Liaudat, Nils Laurent, Brett Levac, Mathurin Massias, Thomas Moreau, Thibaut Modrzyk, Brayan Monroy, Sebastian Neumayer, Jérémy Scanvic, Florian Sarron, Victor Sechaud, Georg Schramm, Romain Vo, and Pierre Weiss, *Deep-inverse: A python package for solving imaging inverse problems with deep learning*, Technical Report 2505.20160, arXiv, 2025. <https://arxiv.org/abs/2505.20160>.
- [89] Rong Tang and Yun Yang, *Adaptivity of diffusion models to manifold structures*, International conference on artificial intelligence and statistics, 2024, pp. 1648–1656.
- [90] Matus Telgarsky, *Neural networks and rational functions*, International conference on machine learning, 2017, pp. 3387–3393.
- [91] Roger Temam, *Navier–stokes equations and nonlinear functional analysis*, SIAM, 1995.

- [92] Singanallur V Venkatakrishnan, Charles A Bouman, and Brendt Wohlberg, *Plug-and-play priors for model based reconstruction*, 2013 ieee global conference on signal and information processing, 2013, pp. 945–948.
- [93] Michael B Wakin, *Manifold-based signal recovery and parameter estimation from compressive measurements*, arXiv preprint arXiv:1002.1247 (2010).
- [94] Haixin Wang, Jianlong Chang, Yihang Zhai, Xiao Luo, Jinan Sun, Zhouchen Lin, and Qi Tian, *Lion: Implicit vision prompt tuning*, Proceedings of the aaai conference on artificial intelligence, 2024.
- [95] Ezra Winston and J Zico Kolter, *Monotone operator equilibrium networks*, Advances in neural information processing systems **33** (2020), 10718–10728.
- [96] Zhoutong Wu, Yimu Zhang, Cong Fang, and Zhouchen Lin, *Separation and bias of deep equilibrium models on expressivity and learning dynamics*, Advances in Neural Information Processing Systems **37** (2024), 32476–32511.
- [97] Xingyu Xie, Qiuhaio Wang, Zenan Ling, Xia Li, Guangcan Liu, and Zhouchen Lin, *Optimization induced equilibrium networks: An explicit optimization perspective for understanding equilibrium models*, IEEE Transactions on Pattern Analysis and Machine Intelligence **45** (2022), no. 3, 3604–3616.
- [98] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka, *How powerful are graph neural networks?*, International conference on learning representations, 2019.
- [99] Chengda Yang, *Nonlinear image recovery with half-quadratic regularization*, IEEE transactions on Image Processing **4** (1995), no. 7, 932–946.
- [100] Yongyi Yang, Tang Liu, Yangkun Wang, Zengfeng Huang, and David Wipf, *Implicit vs unfolded graph neural networks*, Journal of Machine Learning Research **26** (2025), no. 82, 1–46.
- [101] Youhao Yu and Richard M Dansereau, *Msdq-deq-net: Deep equilibrium model (deq) with multiscale dilated convolution for image compressive sensing (cs)*, IET Signal Processing **2024** (2024), no. 1, 6666549.
- [102] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte, *Plug-and-play image restoration with deep denoiser prior*, IEEE Transactions on Pattern Analysis and Machine Intelligence **44** (2021), no. 10, 6360–6376.
- [103] Yongjian Zhong, Hieu Vu, Tianbao Yang, and Bijaya Adhikari, *Efficient and effective implicit dynamic graph neural network*, Proceedings of the 30th acm sigkdd conference on knowledge discovery and data mining, 2024, pp. 4595–4606.
- [104] Zihao Zou, Jiaming Liu, Brendt Wohlberg, and Ulugbek S Kamilov, *Deep equilibrium learning of explicit regularization functionals for imaging inverse problems*, IEEE Open Journal of Signal Processing **4** (2023), 390–398.

## APPENDIX A. PROOFS OF MAIN RESULTS

*Proof of Theorem 2.4.* Given any  $\mathcal{F}$  satisfying Assumption 2.2, the existence of  $\mathcal{G}$  is proved by the following construction:

$$(11) \quad \mathcal{G}(\mathbf{y}, \mathbf{x}) = \mathcal{F}(\mathbf{x}) + (1 - \varepsilon(\mathbf{x}))(\mathbf{y} - \mathcal{F}(\mathbf{x})).$$

The proof will be done by choosing a function  $\varepsilon : \mathbb{X} \rightarrow \mathbb{R}$  such that

- Functions  $\varepsilon(\mathbf{x})$  and  $\varepsilon(\mathbf{x})\mathcal{F}(\mathbf{x})$  are both globally Lipschitz continuous on  $\mathbb{X}$ .
- $0 < \varepsilon(\mathbf{x}) < 1$  for any  $\mathbf{x} \in \mathbb{X}$ .

The existence of such a  $\varepsilon$  function is deferred to Theorem A.1. Now let’s suppose such a  $\varepsilon(\mathbf{x})$  is given and finish the whole proof. First let’s check the contractivity of  $\mathcal{G}$  in (11) as  $\mathbf{x}$  fixed. For any  $\mathbf{y}, \hat{\mathbf{y}} \in \mathbb{R}^n$ , it holds that

$$\mathcal{G}(\mathbf{y}, \mathbf{x}) - \mathcal{G}(\hat{\mathbf{y}}, \mathbf{x}) = (1 - \varepsilon(\mathbf{x}))(\mathbf{y} - \mathcal{F}(\mathbf{x})) - (1 - \varepsilon(\mathbf{x}))(\hat{\mathbf{y}} - \mathcal{F}(\mathbf{x})) = (1 - \varepsilon(\mathbf{x}))(\mathbf{y} - \hat{\mathbf{y}}).$$

Since  $0 < \varepsilon(\mathbf{x}) < 1$  for  $\mathbf{x} \in \mathbb{X}$ , we conclude that  $\mathcal{G}(\cdot, \mathbf{x})$  is a contractor for  $\mathbf{x} \in \mathbb{X}$ . In addition, the continuity of the contractive factor  $(1 - \varepsilon(\mathbf{x}))$  is directly resulted from the continuity of  $\varepsilon(\mathbf{x})$ . Finally, we check the Lipschitz continuity as  $\mathbf{y}$  fixed. For any  $\mathbf{x}, \hat{\mathbf{x}} \in \mathbb{X}$  and any  $\mathbf{y} \in \mathbb{R}^n$ , it holds that

$$\begin{aligned} & \mathcal{G}(\mathbf{y}, \mathbf{x}) - \mathcal{G}(\mathbf{y}, \hat{\mathbf{x}}) \\ &= \left( \mathcal{G}(\mathbf{y}, \mathbf{x}) - \mathbf{y} \right) - \left( \mathcal{G}(\mathbf{y}, \hat{\mathbf{x}}) - \mathbf{y} \right) \\ &= \left( \mathcal{F}(\mathbf{x}) - \mathbf{y} + (1 - \varepsilon(\mathbf{x}))(\mathbf{y} - \mathcal{F}(\mathbf{x})) \right) - \left( \mathcal{F}(\hat{\mathbf{x}}) - \mathbf{y} + (1 - \varepsilon(\hat{\mathbf{x}}))(\mathbf{y} - \mathcal{F}(\hat{\mathbf{x}})) \right) \\ &= -\varepsilon(\mathbf{x})(\mathbf{y} - \mathcal{F}(\mathbf{x})) + \varepsilon(\hat{\mathbf{x}})(\mathbf{y} - \mathcal{F}(\hat{\mathbf{x}})) \\ &= \left( -\varepsilon(\mathbf{x}) + \varepsilon(\hat{\mathbf{x}}) \right) \mathbf{y} + \left( \varepsilon(\mathbf{x})\mathcal{F}(\mathbf{x}) - \varepsilon(\hat{\mathbf{x}})\mathcal{F}(\hat{\mathbf{x}}) \right) \end{aligned}$$

With a fixed  $\mathbf{y} \in \mathbb{R}^n$ , the Lipschitz continuity of  $\mathcal{G}(\mathbf{y}, \cdot)$  follows from the Lipschitz continuity of  $\varepsilon(\mathbf{x})$  and  $\varepsilon(\mathbf{x})\mathcal{F}(\mathbf{x})$ . In particular, by denoting the Lipschitz constants of  $\varepsilon(\mathbf{x})$  and  $\varepsilon(\mathbf{x})\mathcal{F}(\mathbf{x})$  as  $L_\varepsilon$  and  $L_{\varepsilon\mathcal{F}}$  respectively, we have

$$\|\mathcal{G}(\mathbf{y}, \mathbf{x}) - \mathcal{G}(\mathbf{y}, \hat{\mathbf{x}})\| \leq L_\varepsilon \|\mathbf{x} - \hat{\mathbf{x}}\| \cdot \|\mathbf{y}\| + L_{\varepsilon\mathcal{F}} \|\mathbf{x} - \hat{\mathbf{x}}\| \leq (L_\varepsilon \|\mathbf{y}\| + L_{\varepsilon\mathcal{F}}) \|\mathbf{x} - \hat{\mathbf{x}}\|$$

where the Lipschitz constant of  $\mathcal{G}$ ,  $L := L_\varepsilon \|\mathbf{y}\| + L_{\varepsilon\mathcal{F}}$ , grows linearly w.r.t.  $\|\mathbf{y}\|$ , which finishes the whole proof.  $\square$

Below we provide the core theorem used in the proof of Theorem 2.4 and its proof.

**Theorem A.1.** *For any  $\mathcal{F}$  and  $\mathbb{X}$  satisfying Assumption 2.2, there exists a function  $\varepsilon : \mathbb{X} \rightarrow \mathbb{R}$  such that  $0 < \varepsilon(\mathbf{x}) < 1$  for  $\mathbf{x} \in \mathbb{X}$ , and  $\varepsilon(\mathbf{x})$  and  $\varepsilon(\mathbf{x})\mathcal{F}(\mathbf{x})$  are both globally Lipschitz continuous on  $\mathbb{X}$ .*

*Proof.* Let  $\bar{\mathbb{X}}$  be the closure of set  $\mathbb{X}$ . In this proof, we will first extend  $\mathcal{F}$  to  $\bar{\mathbb{X}}$ , construct the  $\varepsilon$  function on  $\bar{\mathbb{X}}$ , and finally prove the global Lipschitz continuity of  $\varepsilon(\mathbf{x})$  and  $\varepsilon(\mathbf{x})\mathcal{F}(\mathbf{x})$  on  $\bar{\mathbb{X}}$ .

**Step 1: Extension to  $\bar{\mathbb{X}}$ .** First we extend  $\mathcal{F}$  to  $\bar{\mathbf{x}} \in \bar{\mathbb{X}} \setminus \mathbb{X}$  by the limit relative to  $\mathbb{X}$ :

$$\mathcal{F}(\bar{\mathbf{x}}) = \begin{cases} \lim_{\mathbb{X} \ni \mathbf{x} \rightarrow \bar{\mathbf{x}}} \mathcal{F}(\mathbf{x}), & \text{if } \lim_{\mathbb{X} \ni \mathbf{x} \rightarrow \bar{\mathbf{x}}} \mathcal{F}(\mathbf{x}) \text{ exists,} \\ 0, & \text{otherwise.} \end{cases}$$

Note that even if  $\mathcal{F}$  is continuously extendable to  $\bar{\mathbf{x}}$ , it is still possible that  $\mathcal{F}$  is not locally Lipschitz continuous at the point  $\bar{\mathbf{x}}$ . A simple example is the function  $\sqrt{x}$ , which is continuous as  $x \geq 0$  and locally Lipschitz continuous for all points  $x > 0$  but NOT locally Lipschitz at  $x = 0$ . We collect all these points (where  $\mathcal{F}$  is not locally Lipschitz) into the set  $\mathbb{D}(\mathcal{F})$ :

$$\mathbb{D}(\mathcal{F}) := \{\mathbf{x} \in \bar{\mathbb{X}} : \mathcal{F} \text{ is not locally Lipschitz continuous at } \mathbf{x}\}$$

For brevity, we will use  $\mathbb{D}$  to denote  $\mathbb{D}(\mathcal{F})$ . It holds that  $\mathbb{D}$  is a closed set (ref. to Lemma A.2) and  $\mathbb{D} \subset \bar{\mathbb{X}} \setminus \mathbb{X}$ .

**Step 2: Constructing a function  $\varepsilon : \bar{\mathbb{X}} \rightarrow \mathbb{R}_{\geq 0}$ .** Now let's define a set including all points that are very "safe", i.e., sufficiently far from the discontinuity set  $\mathbb{D}$ . In particular, given a positive real number  $r > 0$ , the set  $\mathbb{D}_r$  is define by

$$\mathbb{D}_r := \{\mathbf{x} \in \bar{\mathbb{X}} : d(\mathbf{x}, \mathbb{D}) \geq r\},$$

where  $d(\mathbf{x}, \mathbb{D})$  means the distance of  $\mathbf{x}$  and  $\mathbb{D}$ , and the closedness of  $\mathbb{D}_r$  can be derived from the continuity of the distance function. Since  $\mathbb{D}_r \subset \bar{\mathbb{X}}$  and  $\bar{\mathbb{X}}$  is compact,  $\mathbb{D}_r$  must be compact. Note that  $\mathbb{D}_r$  and  $\mathbb{D}$  are disjoint, hence  $\mathcal{F}$  is locally Lipschitz continuous everywhere on  $\mathbb{D}_r$ . Thanks to the fact that local Lipschitz continuity on a compact set implies global Lipschitz continuity (ref to Lemma A.3), we can conclude that  $\mathcal{F}$  is bounded and globally Lipschitz continuous on  $\mathbb{D}_r$  for all  $r > 0$ . Therefore, the following two supremums exist, as long as the cardinality (number of elements) of  $\mathbb{D}_r$  is large enough:

$$h_1(r) = \begin{cases} \sup_{\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{D}_r, \mathbf{x}_1 \neq \mathbf{x}_2} \frac{\|\mathcal{F}(\mathbf{x}_1) - \mathcal{F}(\mathbf{x}_2)\|}{\|\mathbf{x}_1 - \mathbf{x}_2\|}, & \text{card}(\mathbb{D}_r) \geq 2, \\ 0, & \text{otherwise.} \end{cases}$$

$$h_2(r) = \begin{cases} \sup_{\mathbf{x} \in \mathbb{D}_r} \|\mathcal{F}(\mathbf{x})\|, & \text{card}(\mathbb{D}_r) \geq 1, \\ 0, & \text{otherwise.} \end{cases}$$

Here, both  $h_1$  and  $h_2$  are non-negative and monotone non-increasing on  $(0, +\infty)$ . Then we define:

$$\hat{h}(r) = \frac{1}{h_1(r) + h_2(r) + 1}.$$

It has the following properties:

- Bounded:  $0 < \hat{h}(r) \leq 1$  as  $r > 0$ .
- Monotone :  $\hat{h}(r_1) \leq \hat{h}(r_2)$  as  $0 < r_1 \leq r_2$ . (Due to the monotonicity of  $h_1$  and  $h_2$ )

- Naturally extended to  $r = 0$ :  $\lim_{r \rightarrow 0_+} \hat{h}(r)$  exists. (Due to the monotonicity of  $\hat{h}$ )
- $\hat{h}(r)h_i(r) < 1$  for  $r \geq 0$  and  $i = 1, 2$ .

These properties implies that  $\hat{h}$  is Riemann integrable on  $[0, +\infty)$ . Then we can define the following function:

$$\hat{\varepsilon}(r) := \int_0^r \hat{h}(s) ds$$

with the following properties:

- $\hat{\varepsilon}(0) = 0$ .
- Monotone increasing. This is a straightforward result of the fact that  $\hat{h}(s) > 0$  for  $s > 0$ .
- Strictly positive as  $r > 0$ . This is also straightforward as  $\hat{h}(s) > 0$  for  $s > 0$ .
- 1-Lipschitz continuous on  $[0, +\infty)$ . For any  $r_1, r_2$  with  $0 \leq r_1 \leq r_2 < +\infty$ , we have

$$|\hat{\varepsilon}(r_1) - \hat{\varepsilon}(r_2)| = \hat{\varepsilon}(r_2) - \hat{\varepsilon}(r_1) = \int_{r_1}^{r_2} \hat{h}(s) ds \leq \left( \sup_{r \geq 0} \hat{h}(r) \right) |r_1 - r_2| = |r_1 - r_2|.$$

With such a  $\hat{\varepsilon}(r)$ , we can define  $\varepsilon(\mathbf{x})$  by

$$\varepsilon(\mathbf{x}) = \frac{\hat{\varepsilon}(d(\mathbf{x}, \mathbb{D}))}{1 + \hat{\varepsilon}(d(\mathbf{x}, \mathbb{D}))}.$$

It holds that  $\varepsilon(\mathbf{x}) = 0$  for  $\mathbf{x} \in \mathbb{D}$  and  $0 < \varepsilon(\mathbf{x}) < 1$  for  $\mathbf{x} \in \bar{\mathbb{X}} \setminus \mathbb{D}$ . As  $\mathbb{D} \subset \bar{\mathbb{X}} \setminus \mathbb{X}$ , we have  $0 < \varepsilon(\mathbf{x}) < 1$  for  $\mathbf{x} \in \bar{\mathbb{X}}$ .

**Step 3: Establishing the Lipschitz continuity.** Since the distance function  $d(\mathbf{x}, \mathbb{D})$  is 1-Lipschitz continuous [27, Theorem 4.8 (1)], the Lipschitz continuity of  $\hat{\varepsilon}$  implies the Lipschitz continuity of  $\varepsilon$ . In particular, for all  $\mathbf{x}_1, \mathbf{x}_2 \in \bar{\mathbb{X}}$ , it holds that

$$\begin{aligned} & \left| \varepsilon(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2) \right| \\ &= \left| \frac{\hat{\varepsilon}(d(\mathbf{x}_1, \mathbb{D}))}{1 + \hat{\varepsilon}(d(\mathbf{x}_1, \mathbb{D}))} - \frac{\hat{\varepsilon}(d(\mathbf{x}_2, \mathbb{D}))}{1 + \hat{\varepsilon}(d(\mathbf{x}_2, \mathbb{D}))} \right| \quad \left( \frac{x}{1+x} \text{ is 1-Lipschitz as } \left( \frac{x}{1+x} \right)' = \frac{1}{(1+x)^2} \right) \\ &\leq \left| \hat{\varepsilon}(d(\mathbf{x}_1, \mathbb{D})) - \hat{\varepsilon}(d(\mathbf{x}_2, \mathbb{D})) \right| \quad (\text{Lipschitz continuity of } \hat{\varepsilon}) \\ &\leq \left| d(\mathbf{x}_1, \mathbb{D}) - d(\mathbf{x}_2, \mathbb{D}) \right| \leq \|\mathbf{x}_1 - \mathbf{x}_2\| \quad (\text{Lipschitz continuity of } d) \end{aligned}$$

Therefore, to complete the whole proof, it's enough to show the global Lipschitz continuity of  $\varepsilon\mathcal{F}$  on  $\bar{\mathbb{X}}$ . As  $\bar{\mathbb{X}}$  is compact, and thanks to Lemma A.3, it's enough to show  $\varepsilon\mathcal{F}$  is locally Lipschitz everywhere on  $\bar{\mathbb{X}}$ .

First, we consider the local Lipschitz continuity of  $\varepsilon\mathcal{F}$  on  $\bar{\mathbb{X}} \setminus \mathbb{D}$ . Due to Lemma A.2,  $\bar{\mathbb{X}} \setminus \mathbb{D}$  must be open relative to  $\bar{\mathbb{X}}$ . For any  $\mathbf{x} \in \bar{\mathbb{X}} \setminus \mathbb{D}$ , there must be a small enough  $r > 0$  such that  $\mathbb{U} := \mathbb{B}(\mathbf{x}, r) \cap \bar{\mathbb{X}} \subset \bar{\mathbb{X}} \setminus \mathbb{D}$ . Pick  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{U}$ . For any  $\mathbf{x}_1, \mathbf{x}_2$ , it holds that

$$\begin{aligned} & \|\varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)\mathcal{F}(\mathbf{x}_2)\| \\ (12) \quad &= \|\varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_1) - \varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_2) + \varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_2) - \varepsilon(\mathbf{x}_2)\mathcal{F}(\mathbf{x}_2)\| \\ &\leq \varepsilon(\mathbf{x}_1)\|\mathcal{F}(\mathbf{x}_1) - \mathcal{F}(\mathbf{x}_2)\| + |\varepsilon(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)| \cdot \|\mathcal{F}(\mathbf{x}_2)\|. \end{aligned}$$

Since both  $\varepsilon$  and  $\mathcal{F}$  are locally Lipschitz and locally bounded everywhere on  $\bar{\mathbb{X}} \setminus \mathbb{D}$ , they must be Lipschitz and bounded within  $\mathbb{U}$ . Then the local Lipschitz continuity of  $\varepsilon\mathcal{F}$  at  $\mathbf{x}$  immediately follows from (12). Note that  $\mathbf{x}$  is arbitrarily picked from  $\bar{\mathbb{X}} \setminus \mathbb{D}$ , hence  $\varepsilon\mathcal{F}$  is locally Lipschitz everywhere on  $\bar{\mathbb{X}} \setminus \mathbb{D}$ .

Next, we consider the local Lipschitz continuity of  $\varepsilon\mathcal{F}$  on  $\mathbb{D}$ . For any  $\mathbf{x} \in \mathbb{D}$ , we consider its neighborhood  $\mathbb{U} := \mathbb{B}(\mathbf{x}, 1) \cap \bar{\mathbb{X}}$  and pick  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{U}$ . Then we need to consider three cases. The

first case is both  $\mathbf{x}_1, \mathbf{x}_2$  belong to the discontinuity set  $\mathbb{D}$ :  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{D}$ . In this case, it holds that  $\varepsilon(\mathbf{x}_1) = \varepsilon(\mathbf{x}_2) = 0$  and hence

$$\left\| \varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)\mathcal{F}(\mathbf{x}_2) \right\| = 0 \leq \|\mathbf{x}_1 - \mathbf{x}_2\|.$$

The second case is that one of the point is in  $\mathbb{D}$  while the other is not, we suppose  $\mathbf{x}_1 \in \mathbb{D}, \mathbf{x}_2 \in \bar{\mathbb{X}} \setminus \mathbb{D}$ , then

$$\begin{aligned} & \left\| \varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)\mathcal{F}(\mathbf{x}_2) \right\| \\ &= \left\| \varepsilon(\mathbf{x}_2)\mathcal{F}(\mathbf{x}_2) \right\| \leq \hat{\varepsilon}(d(\mathbf{x}_2, \mathbb{D})) \|\mathcal{F}(\mathbf{x}_2)\| \\ &= \left( \int_0^{d(\mathbf{x}_2, \mathbb{D})} \hat{h}(s) ds \right) \|\mathcal{F}(\mathbf{x}_2)\| \\ &\leq \hat{h}(d(\mathbf{x}_2, \mathbb{D})) \cdot d(\mathbf{x}_2, \mathbb{D}) \cdot \|\mathcal{F}(\mathbf{x}_2)\| && \text{(Monotonicity of } \hat{h}) \\ &\leq \hat{h}(d(\mathbf{x}_2, \mathbb{D})) \cdot d(\mathbf{x}_2, \mathbb{D}) \cdot h_2(d(\mathbf{x}_2, \mathbb{D})) && \text{(Definition of } h_2) \\ &< d(\mathbf{x}_2, \mathbb{D}) && (\hat{h}(r) \cdot h_2(r) < 1 \text{ as } r \geq 0) \\ &= d(\mathbf{x}_2, \mathbb{D}) - d(\mathbf{x}_1, \mathbb{D}) \leq \|\mathbf{x}_1 - \mathbf{x}_2\| \end{aligned}$$

Finally, we consider the last case where  $\mathbf{x}_1, \mathbf{x}_2 \in \bar{\mathbb{X}} \setminus \mathbb{D}$ . Without loss of generality, we assume

$$0 < d(\mathbf{x}_1, \mathbb{D}) \leq d(\mathbf{x}_2, \mathbb{D}).$$

Then the definition of  $h_1$  and  $h_2$  implies that

$$\begin{aligned} & \|\mathcal{F}(\mathbf{x}_1) - \mathcal{F}(\mathbf{x}_2)\| \\ &\leq \max \left( h_1(d(\mathbf{x}_1, \mathbb{D})), h_1(d(\mathbf{x}_2, \mathbb{D})) \right) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| \\ &= h_1(d(\mathbf{x}_1, \mathbb{D})) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|, \end{aligned}$$

and

$$\|\mathcal{F}(\mathbf{x}_2)\| \leq h_2(d(\mathbf{x}_2, \mathbb{D})).$$

Consequently, applying (12) and the above inequalities, we have

$$\begin{aligned} & \left\| \varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)\mathcal{F}(\mathbf{x}_2) \right\| \\ &\leq \varepsilon(\mathbf{x}_1) \|\mathcal{F}(\mathbf{x}_1) - \mathcal{F}(\mathbf{x}_2)\| + |\varepsilon(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)| \cdot \|\mathcal{F}(\mathbf{x}_2)\| \\ &\leq \varepsilon(\mathbf{x}_1) \cdot h_1(d(\mathbf{x}_1, \mathbb{D})) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| + |\varepsilon(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)| \cdot h_2(d(\mathbf{x}_2, \mathbb{D})) \\ &\leq \hat{\varepsilon}(d(\mathbf{x}_1, \mathbb{D})) \cdot h_1(d(\mathbf{x}_1, \mathbb{D})) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| + \left| \hat{\varepsilon}(d(\mathbf{x}_1, \mathbb{D})) - \hat{\varepsilon}(d(\mathbf{x}_2, \mathbb{D})) \right| \cdot h_2(d(\mathbf{x}_2, \mathbb{D})) \\ &= \left( \int_0^{d(\mathbf{x}_1, \mathbb{D})} \hat{h}(s) ds \right) \cdot h_1(d(\mathbf{x}_1, \mathbb{D})) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| + \left( \int_{d(\mathbf{x}_1, \mathbb{D})}^{d(\mathbf{x}_2, \mathbb{D})} \hat{h}(s) ds \right) \cdot h_2(d(\mathbf{x}_2, \mathbb{D})) \\ &\leq d(\mathbf{x}_1, \mathbb{D}) \cdot \hat{h}(d(\mathbf{x}_1, \mathbb{D})) \cdot h_1(d(\mathbf{x}_1, \mathbb{D})) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| \\ &\quad + \left| d(\mathbf{x}_1, \mathbb{D}) - d(\mathbf{x}_2, \mathbb{D}) \right| \cdot \hat{h}(d(\mathbf{x}_2, \mathbb{D})) \cdot h_2(d(\mathbf{x}_2, \mathbb{D})) \\ &< d(\mathbf{x}_1, \mathbb{D}) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| + \left| d(\mathbf{x}_1, \mathbb{D}) - d(\mathbf{x}_2, \mathbb{D}) \right| \\ &\leq d(\mathbf{x}_1, \mathbb{D}) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| + \|\mathbf{x}_1 - \mathbf{x}_2\| \end{aligned}$$

The last inequality results from  $\hat{h}(r) \cdot (h_1(r) + h_2(r)) < 1$  for all  $r > 0$ . And the above inequalities imply

$$\left\| \varepsilon(\mathbf{x}_1)\mathcal{F}(\mathbf{x}_1) - \varepsilon(\mathbf{x}_2)\mathcal{F}(\mathbf{x}_2) \right\| \leq (\text{diam}(\bar{\mathbb{X}}) + 1) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|.$$

Combining all the results together, we have  $\varepsilon\mathcal{F}$  is locally  $(\text{diam}(\overline{\mathbb{X}}) + 1)$ -Lipschitz at any  $\mathbf{x} \in \overline{\mathbb{X}}$ . Then the compactness of  $\overline{\mathbb{X}}$  concludes the global Lipschitz continuous of  $\varepsilon\mathcal{F}$ , which finishes the whole proof.  $\square$

Follows are some lemmas (as well as their proofs) that we used in the proof of Theorem A.1.

**Lemma A.2.** *Let  $\mathbb{T} \subset \mathbb{R}^d$  be closed and let  $\mathcal{F} : \mathbb{T} \rightarrow \mathbb{R}^n$ . Denote by  $\mathbb{D}(\mathcal{F}) \subset \mathbb{T}$  the set of points at which  $\mathcal{F}$  is not locally Lipschitz. Then  $\mathbb{D}(\mathcal{F})$  is closed (in  $\mathbb{T}$ , hence in  $\mathbb{R}^d$ ).*

*Proof.* Recall that  $\mathcal{F}$  is locally Lipschitz (relative to  $\mathbb{T}$ ) at  $\mathbf{x} \in \mathbb{T}$  if there exist  $r > 0$  and  $L > 0$  such that

$$\|\mathcal{F}(\mathbf{u}) - \mathcal{F}(\mathbf{v})\| \leq L \|\mathbf{u} - \mathbf{v}\| \quad \text{for all } \mathbf{u}, \mathbf{v} \in \mathbb{T} \cap \mathbb{B}(\mathbf{x}, r).$$

Let  $\mathbb{G} := \mathbb{T} \setminus \mathbb{D}(\mathcal{F})$  be the set of points where  $\mathcal{F}$  is locally Lipschitz. We first show that  $\mathbb{G}$  is relatively open in  $\mathbb{T}$ . Fix  $\mathbf{x} \in \mathbb{G}$  and choose  $r, L$  as above. If  $\mathbf{x}' \in \mathbb{T} \cap \mathbb{B}(\mathbf{x}, r/2)$ , then  $\mathbb{B}(\mathbf{x}', r/2) \subset \mathbb{B}(\mathbf{x}, r)$ ; hence the same  $L$  works on  $\mathbb{T} \cap \mathbb{B}(\mathbf{x}', r/2)$ , so  $\mathcal{F}$  is locally Lipschitz at  $\mathbf{x}'$ . Therefore  $\mathbb{T} \cap \mathbb{B}(\mathbf{x}, r/2) \subset \mathbb{G}$ , proving that  $\mathbb{G}$  is open in  $\mathbb{T}$ . Consequently,  $\mathbb{D}(\mathcal{F}) = \mathbb{T} \setminus \mathbb{G}$  is closed in  $\mathbb{T}$ . Since  $\mathbb{T}$  is closed in  $\mathbb{R}^d$ , every set closed in  $\mathbb{T}$  is also closed in  $\mathbb{R}^d$ . Hence  $\mathbb{D}(\mathcal{F})$  is closed in  $\mathbb{R}^d$  as well.  $\square$

**Lemma A.3.** *Let  $\mathbb{T}$  be a compact set. If  $\mathcal{F}$  is locally Lipschitz everywhere on  $\mathbb{T}$ , then it must be globally Lipschitz on  $\mathbb{T}$ .*

*Proof.* Assume, to the contrary, that  $\mathcal{F}$  is not globally Lipschitz on  $\mathbb{T}$ . Then we can choose sequences  $\{\mathbf{x}_k\}_{k \geq 1}, \{\mathbf{y}_k\}_{k \geq 1} \subset \mathbb{T}$  such that

$$(13) \quad \frac{\|\mathcal{F}(\mathbf{x}_k) - \mathcal{F}(\mathbf{y}_k)\|}{\|\mathbf{x}_k - \mathbf{y}_k\|} \xrightarrow{k \rightarrow \infty} \infty.$$

Local Lipschitzness implies continuity of  $\mathcal{F}$  on  $\mathbb{T}$ , so by compactness  $\mathcal{F}$  is bounded: there exists  $C < \infty$  with  $\|\mathcal{F}(\mathbf{z})\| \leq C$  for all  $\mathbf{z} \in \mathbb{T}$ . Consequently,

$$\|\mathcal{F}(\mathbf{x}_k) - \mathcal{F}(\mathbf{y}_k)\| \leq 2C \quad \text{for all } k,$$

and therefore (13) forces  $\|\mathbf{x}_k - \mathbf{y}_k\| \rightarrow 0$ .

By sequential compactness of  $\mathbb{T}$ , passing to a subsequence (not relabeled) we may assume  $\mathbf{x}_k \rightarrow \mathbf{x} \in \mathbb{T}$ ; since  $\|\mathbf{x}_k - \mathbf{y}_k\| \rightarrow 0$ , we also have  $\mathbf{y}_k \rightarrow \mathbf{x}$ . Since  $\mathcal{F}$  is locally Lipschitz at  $\mathbf{x}$ , for  $k$  large enough we have

$$\frac{\|\mathcal{F}(\mathbf{x}_k) - \mathcal{F}(\mathbf{y}_k)\|}{\|\mathbf{x}_k - \mathbf{y}_k\|} \leq L,$$

for some  $L > 0$ , which contradicts (13). Therefore  $\mathcal{F}$  must be globally Lipschitz on  $\mathbb{T}$ .  $\square$

Note that Lemmas A.2 and A.3 are standard in real analysis, we provide concise proofs of them for the sake of completeness. Next we provide proof of Theorem 2.5.

*Proof of Theorem 2.5.* Let  $\overline{\mathbb{X}}$  be the closure of  $\mathbb{X}$ . In this proof, we will first extend the operator  $\mathcal{G}$  to  $\mathbb{R}^n \times \overline{\mathbb{X}}$ , and then analyze its properties on this closed domain.

**Step 1: Extension to  $\overline{\mathbb{X}}$ .** For any  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathcal{G}(\mathbf{y}, \mathbf{x})$  is globally Lipschitz continuous on  $\mathbb{X}$ , hence its extension is naturally define by

$$\mathcal{G}(\mathbf{y}, \bar{\mathbf{x}}) := \lim_{\mathbb{X} \ni \mathbf{x} \rightarrow \bar{\mathbf{x}}} \mathcal{G}(\mathbf{y}, \mathbf{x}), \quad \text{for all } \bar{\mathbf{x}} \in \overline{\mathbb{X}} \setminus \mathbb{X}.$$

Different from the proof of Theorem A.1 where  $\mathcal{F}$  might be not locally Lipschitz at  $\bar{\mathbf{x}}$  even if it is continuous at  $\bar{\mathbf{x}}$ , here the extended  $\mathcal{G}$  must be Lipschitz at  $\bar{\mathbf{x}}$  and hence Lipschitz on the overall set  $\overline{\mathbb{X}}$ . This can be verified by examining the difference quotient for  $\mathbf{x}_1 \neq \mathbf{x}_2$  and  $\mathbf{y} \in \mathbb{R}^n$ :

$$\Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \mathbf{x}_2] := \frac{\|\mathcal{G}(\mathbf{y}, \mathbf{x}_1) - \mathcal{G}(\mathbf{y}, \mathbf{x}_2)\|}{\|\mathbf{x}_1 - \mathbf{x}_2\|}$$

Let  $\mathcal{G}(\mathbf{y}, \cdot)$ 's Lipschitz constant on  $\mathbb{X}$  be  $L(\mathbf{y}) := \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \mathbb{X}} \Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \mathbf{x}_2]$ . For any  $\mathbf{x}_1 \in \mathbb{X}$  and  $\bar{\mathbf{x}}_2 \in \bar{\mathbb{X}} \setminus \mathbb{X}$ , it holds that

$$\Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \bar{\mathbf{x}}_2] = \lim_{\mathbb{X} \ni \mathbf{x}_2 \rightarrow \bar{\mathbf{x}}_2} \Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \mathbf{x}_2] \leq \sup_{\mathbf{x}_2 \in \mathbb{X}: \mathbf{x}_2 \neq \mathbf{x}_1} \Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \mathbf{x}_2] \leq L(\mathbf{y})$$

For any  $\bar{\mathbf{x}}_1 \neq \bar{\mathbf{x}}_2 \in \bar{\mathbb{X}} \setminus \mathbb{X}$ , we have

$$\Delta \mathcal{G}[\mathbf{y}; \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] = \lim_{\mathbb{X} \ni \mathbf{x}_1 \rightarrow \bar{\mathbf{x}}_1} \lim_{\mathbb{X} \ni \mathbf{x}_2 \rightarrow \bar{\mathbf{x}}_2} \Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \mathbf{x}_2] \leq \sup_{\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{X}: \mathbf{x}_2 \neq \mathbf{x}_1} \Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \mathbf{x}_2] = L(\mathbf{y})$$

Therefore, we obtain an upper bound for  $\mathcal{G}(\mathbf{y}, \cdot)$ 's Lipschitz constant on  $\bar{\mathbb{X}}$ :

$$\begin{aligned} & \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \bar{\mathbb{X}}} \Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \mathbf{x}_2] \\ &= \max \left( \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \mathbb{X}} \Delta \mathcal{G}[\mathbf{y}; \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], \sup_{\mathbf{x}_1 \in \mathbb{X}, \mathbf{x}_2 \in \bar{\mathbb{X}} \setminus \mathbb{X}} \Delta \mathcal{G}[\mathbf{y}; \mathbf{x}_1, \bar{\mathbf{x}}_2], \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \bar{\mathbb{X}} \setminus \mathbb{X}} \Delta \mathcal{G}[\mathbf{y}; \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \right) \\ &\leq \max(L(\mathbf{y}), L(\mathbf{y}), L(\mathbf{y})) = L(\mathbf{y}) \end{aligned}$$

That is, for any  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathcal{G}(\mathbf{y}, \cdot)$  is globally Lipschitz on  $\bar{\mathbb{X}}$ , and the Lipschitz constant is the same with that of  $\mathbb{X}$ .

In the other hand, let's consider the Lipschitz constant (contraction constant) w.r.t.  $\mathbf{y}$  when fixing  $\bar{\mathbf{x}} \in \bar{\mathbb{X}} \setminus \mathbb{X}$ :

$$\mu(\bar{\mathbf{x}}) = \lim_{\mathbb{X} \ni \mathbf{x} \rightarrow \bar{\mathbf{x}}} \mu(\mathbf{x})$$

Since  $0 < \mu(\mathbf{x}) < 1$  for  $\mathbf{x} \in \mathbb{X}$ , by taking limit, we have  $0 \leq \mu(\bar{\mathbf{x}}) \leq 1$ . For those  $\bar{\mathbf{x}}$  with  $\mu(\bar{\mathbf{x}}) < 1$ , the operator  $\mathcal{G}(\cdot, \bar{\mathbf{x}})$  is still contractive. But if  $\mu(\bar{\mathbf{x}}) = 1$ , the operator  $\mathcal{G}(\cdot, \bar{\mathbf{x}})$  is not contractive.

**Step 2: Defining  $\mathbb{D}$  and  $\mathbb{D}_r$ .** We collect all points  $\mathbf{x} \in \bar{\mathbb{X}}$  where the operator  $\mathcal{G}(\cdot, \mathbf{x})$  is not contractive:

$$\mathbb{D} := \{\mathbf{x} \in \bar{\mathbb{X}} : \mu(\mathbf{x}) = 1\}$$

and define a "safe" set that is sufficiently far from  $\mathbb{D}$ :

$$\mathbb{D}_r := \{\mathbf{x} \in \bar{\mathbb{X}} : d(\mathbf{x}, \mathbb{D}) \geq r\}.$$

Note that  $\bar{\mathbb{X}} \setminus \mathbb{D} = \bigcup_{r>0} \mathbb{D}_r$  and  $\mathbb{X} \subset \bar{\mathbb{X}} \setminus \mathbb{D}$ . We obtain

$$\mathbb{X} \subset \bigcup_{r>0} \mathbb{D}_r.$$

For any  $\mathbb{D}_r$  with  $r > 0$ , we can obtain a uniform contraction of the operator  $\mathcal{G}(\cdot, \mathbf{x})$ : There is a constant  $\mu_r \in (0, 1)$  such that

$$(14) \quad \|\mathcal{G}(\mathbf{y}_1, \mathbf{x}) - \mathcal{G}(\mathbf{y}_2, \mathbf{x})\| \leq \mu_r \|\mathbf{y}_1 - \mathbf{y}_2\|$$

for all  $\mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^n$  and  $\mathbf{x} \in \mathbb{D}_r$ , which follows immediately from the continuity of  $\mu(\mathbf{x})$  and the compactness of  $\mathbb{D}_r$ . By the Banach fixed-point theorem, the operator  $\mathcal{G}(\cdot, \mathbf{x})$  must have a unique fixed point  $\mathbf{y}_*$  for each  $\mathbf{x} \in \mathbb{D}_r$ .

To complete the proof of Theorem 2.5, thanks to the fact that  $\mathbb{X} \subset \bigcup_{r>0} \mathbb{D}_r$ , it's enough to show that: For any  $\mathbb{D}_r$  with  $r > 0$ , there is a constant  $C_r$  such that

$$(15) \quad \|\mathbf{y}_*(\mathbf{x}_1) - \mathbf{y}_*(\mathbf{x}_2)\| \leq C_r \|\mathbf{x}_1 - \mathbf{x}_2\|$$

holds for all  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{D}_r$ . In the following steps, we will show (15).

**Step 3: A controllable sequence.** Fix  $\mathbf{x} \in \mathbb{D}_r$ . By defining a sequence  $\{\mathbf{y}_k(\mathbf{x})\}_{k \geq 0} \subset \mathbb{R}^n$ :

$$\mathbf{y}_{k+1}(\mathbf{x}) = \mathcal{G}(\mathbf{y}_k(\mathbf{x}), \mathbf{x}), \quad \mathbf{y}_0 \text{ is constant for all } \mathbf{x},$$

we are able to estimate the upper bound of  $\|\mathbf{y}_*(\mathbf{x})\|$ . In particular, we decompose  $\mathbf{y}_0 - \mathbf{y}_*$  by a series:

$$\mathbf{y}_0 - \mathbf{y}_* = \lim_{k \rightarrow \infty} (\mathbf{y}_0 - \mathbf{y}_k) = \sum_{k=0}^{\infty} (\mathbf{y}_k - \mathbf{y}_{k+1})$$

Thanks to (14), we have

$$\|\mathbf{y}_k(\mathbf{x}) - \mathbf{y}_{k+1}(\mathbf{x})\| \leq \mu_r \|\mathbf{y}_{k-1}(\mathbf{x}) - \mathbf{y}_k(\mathbf{x})\| \cdots \leq \mu_r^k \|\mathbf{y}_0 - \mathbf{y}_1(\mathbf{x})\| = \mu_r^k \|\mathbf{y}_0 - \mathcal{G}(\mathbf{y}_0, \mathbf{x})\|$$

for all  $\mathbf{x} \in \mathbb{D}_r$ . Therefore, it holds that

$$\begin{aligned} \|\mathbf{y}_0 - \mathbf{y}_*(\mathbf{x})\| &\leq \sum_{k=0}^{\infty} \|\mathbf{y}_k(\mathbf{x}) - \mathbf{y}_{k+1}(\mathbf{x})\| \\ &\leq \left( \sum_{k=0}^{\infty} \mu_r^k \right) \|\mathbf{y}_0 - \mathcal{G}(\mathbf{y}_0, \mathbf{x})\| = \frac{1}{1 - \mu_r} \|\mathbf{y}_0 - \mathcal{G}(\mathbf{y}_0, \mathbf{x})\| \end{aligned}$$

Now we can conclude the boundedness of  $\|\mathbf{y}_*(\mathbf{x})\|$  for  $\mathbf{x} \in \mathbb{D}_r$  by the compactness of  $\mathbb{D}_r$ :

$$\|\mathbf{y}_*(\mathbf{x})\| \leq \underbrace{\|\mathbf{y}_0\| + \frac{1}{1 - \mu_r} \sup_{\mathbf{x} \in \mathbb{D}_r} \|\mathbf{y}_0 - \mathcal{G}(\mathbf{y}_0, \mathbf{x})\|}_{\text{defined as } M_r \geq 0.}$$

With the same argument, we have  $\|\mathbf{y}_k(\mathbf{x})\| \leq M_r$  for all  $k \geq 0$  and  $\mathbf{x} \in \mathbb{D}_r$ . It implies that

$$L(\mathbf{y}_k(\mathbf{x})) \leq L_1 + L_2 M_r$$

for some  $L_1, L_2 > 0$  as  $L(\mathbf{y})$  grows linearly w.r.t.  $\|\mathbf{y}\|$ . Consequently, we can estimate an upper bound for the Lipschitz constant of  $\mathbf{y}_k(\mathbf{x})$ . In particular, for  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{D}_r$ , it holds that

$$\begin{aligned} &\|\mathbf{y}_{k+1}(\mathbf{x}_1) - \mathbf{y}_{k+1}(\mathbf{x}_2)\| \\ &= \|\mathcal{G}(\mathbf{y}_k(\mathbf{x}_1), \mathbf{x}_1) - \mathcal{G}(\mathbf{y}_k(\mathbf{x}_2), \mathbf{x}_2)\| \\ &= \|\mathcal{G}(\mathbf{y}_k(\mathbf{x}_1), \mathbf{x}_1) - \mathcal{G}(\mathbf{y}_k(\mathbf{x}_2), \mathbf{x}_1) + \mathcal{G}(\mathbf{y}_k(\mathbf{x}_2), \mathbf{x}_1) - \mathcal{G}(\mathbf{y}_k(\mathbf{x}_2), \mathbf{x}_2)\| \\ &\leq \|\mathcal{G}(\mathbf{y}_k(\mathbf{x}_1), \mathbf{x}_1) - \mathcal{G}(\mathbf{y}_k(\mathbf{x}_2), \mathbf{x}_1)\| + \|\mathcal{G}(\mathbf{y}_k(\mathbf{x}_2), \mathbf{x}_1) - \mathcal{G}(\mathbf{y}_k(\mathbf{x}_2), \mathbf{x}_2)\| \\ &\leq \mu_r \|\mathbf{y}_k(\mathbf{x}_1) - \mathbf{y}_k(\mathbf{x}_2)\| + (L_1 + L_2 M_r) \|\mathbf{x}_1 - \mathbf{x}_2\| \end{aligned}$$

For simplicity, let  $L_r := L_1 + L_2 M_r$ ,  $a_k := \|\mathbf{y}_k(\mathbf{x}_1) - \mathbf{y}_k(\mathbf{x}_2)\|$ , and  $h := \|\mathbf{x}_1 - \mathbf{x}_2\|$ . By recursively applying  $a_{k+1} \leq \mu_r a_k + L_r h$  and  $a_0 = 0$ , we have

$$\|\mathbf{y}_k(\mathbf{x}_1) - \mathbf{y}_k(\mathbf{x}_2)\| = a_k \leq (\mu_r)^k a_0 + (\mu_r^{k-1} + \cdots + \mu_r + 1) L_r h \leq \frac{1}{1 - \mu_r} L_r h = \frac{L_r}{1 - \mu_r} \|\mathbf{x}_1 - \mathbf{x}_2\|.$$

**Step 4: Final proof.** As  $\mathcal{G}(\cdot, \mathbf{x})$  is a contractor w.r.t.  $\mathbf{y}$  for any  $\mathbf{x} \in \mathbb{D}_r$ , it holds that  $\mathbf{y}_k(\mathbf{x}) \rightarrow \mathbf{y}_*(\mathbf{x})$  for any  $\mathbf{x} \in \mathbb{D}_r$ . (Here, as for the ‘‘convergence,’’ we mean the pointwise convergence, which is enough here. We don’t need stronger conditions like the uniform convergence.) For the above  $\mathbf{x}_1, \mathbf{x}_2$ , there is a  $K$  such that

$$\|\mathbf{y}_k(\mathbf{x}_1) - \mathbf{y}_*(\mathbf{x}_1)\| \leq \frac{L_r}{1 - \mu_r} \|\mathbf{x}_1 - \mathbf{x}_2\|, \quad \|\mathbf{y}_k(\mathbf{x}_2) - \mathbf{y}_*(\mathbf{x}_2)\| \leq \frac{L_r}{1 - \mu_r} \|\mathbf{x}_1 - \mathbf{x}_2\|$$

for  $k \geq K$ . Combining the above results, we obtain

$$\begin{aligned} \|\mathbf{y}_*(\mathbf{x}_1) - \mathbf{y}_*(\mathbf{x}_2)\| &\leq \|\mathbf{y}_*(\mathbf{x}_1) - \mathbf{y}_k(\mathbf{x}_1)\| + \|\mathbf{y}_k(\mathbf{x}_1) - \mathbf{y}_k(\mathbf{x}_2)\| + \|\mathbf{y}_k(\mathbf{x}_2) - \mathbf{y}_*(\mathbf{x}_2)\| \\ &\leq \frac{3L_r}{1 - \mu_r} \|\mathbf{x}_1 - \mathbf{x}_2\| \end{aligned}$$

By letting  $C_r = 3L_r/(1 - \mu_r)$ , we get (15), which completes the proof.  $\square$

**Remark:** Our result relaxes two uniformity requirements in [23, Thm. 1A.4]: (i) the contraction modulus  $\mu(\mathbf{x})$  is allowed to vary with  $\mathbf{x}$  (it only needs to be continuous in  $\mathbf{x}$ ), rather than being a single global constant; and (ii) for each  $\mathbf{y}$ , the mapping  $\mathbf{x} \mapsto \mathcal{G}(\mathbf{y}, \mathbf{x})$  is Lipschitz on  $\mathbb{X}$  with a constant that may grow linearly in  $\|\mathbf{y}\|$ , instead of being uniformly bounded in  $\mathbf{y}$ . Because these bounds are not uniform, we conclude only local (as opposed to global) Lipschitz continuity of the fixed-point map  $\mathbf{x} \mapsto \mathbf{y}_*(\mathbf{x})$  on  $\mathbb{X}$ .

## APPENDIX B. A VARIANT ARCHITECTURE

In practice, many works use a variant of the vanilla model  $\mathbf{y}_* = \mathcal{G}(\mathbf{y}_*, \mathbf{x})$ :

$$(16) \quad \mathbf{z}_* = \mathcal{G}(\mathbf{z}_*, \mathcal{Q}_1(\mathbf{x})), \quad \mathbf{y}_* = \mathcal{Q}_2(\mathbf{z}_*)$$

where  $\mathcal{G}$  is the core implicit model,  $\mathcal{Q}_1$  is an encoding network and  $\mathcal{Q}_2$  is a decoding (readout).

At inference, one iterates  $\mathbf{z}_t = \mathcal{G}(\mathbf{z}_{t-1}, \mathcal{Q}_1(\mathbf{x}))$  for  $1 \leq t \leq T$  and finally  $\mathbf{y}_T = \mathcal{Q}_2(\mathbf{z}_T)$ . This often improves empirical performance but does not alter the expressivity in Theorems 2.4–2.5.

**Corollary B.1.** *Under Assumption 2.2, for any  $\mathcal{F}$  there exists a regular implicit operator  $\mathcal{G}$  and globally Lipschitz maps  $\mathcal{Q}_1, \mathcal{Q}_2$  such that  $\mathcal{Q}_2(\text{Fix}(\mathcal{G}(\cdot, \mathcal{Q}_1(\mathbf{x})))) = \mathcal{F}(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{X}$ . Conversely, for any regular implicit operator  $\mathcal{G}$  any globally Lipschitz  $\mathcal{Q}_1, \mathcal{Q}_2$ , the fixed point  $\mathbf{z}_*$  defined by (16) exists uniquely and the induced map  $\mathbf{x} \mapsto \mathbf{y}_*$  must be locally Lipschitz on  $\mathbb{X}$ .*

*Proof.* The claim follows directly from Theorems 2.4–2.5.

*Sufficiency.* Given any locally Lipschitz target  $\mathcal{F}$  on  $\mathbb{X}$ , Theorem 2.4 ensures the existence of a regular  $\mathcal{G}$  whose fixed-point map equals  $\mathcal{F}$ . Taking  $\mathcal{Q}_1, \mathcal{Q}_2$  as both identity maps recovers the sufficiency statement with globally Lipschitz  $\mathcal{Q}_1, \mathcal{Q}_2$ .

*Necessity.* Suppose  $\mathcal{G}$  is regular and  $\mathcal{Q}_1, \mathcal{Q}_2$  are globally Lipschitz. Then the composite update  $\mathcal{G}(\mathbf{z}, \mathcal{Q}_1(\mathbf{x}))$  is still regular in  $\mathbf{z}$  and  $\mathbf{x}$ . By Theorem 2.5, for every  $\mathbf{x} \in \mathbb{X}$ , there is a unique fixed point  $\mathbf{z}_*(\mathbf{x})$  and the map  $\mathbf{x} \mapsto \mathbf{z}_*(\mathbf{x})$  is locally Lipschitz on  $\mathbb{X}$ . Finally, applying the globally Lipschitz readout  $\mathcal{Q}_2$  preserves local Lipschitz continuity, so  $\mathbf{x} \mapsto \mathbf{y}_*$  is locally Lipschitz as claimed. The proof is finished.  $\square$

## APPENDIX C. PROOFS OF THEOREMS FOR INVERSE PROBLEMS

This section proves that the target solution mappings,  $\mathcal{F}_{1a}$  and  $\mathcal{F}_{1b}$ , are single-valued and locally Lipschitz on their domain, as stated in Theorems 3.3 and 3.4. Before the proofs, we first provide some definitions that used in Assumption 3.1.

Given a close subset  $\mathbb{M} \subset \mathbb{R}^n$ , its *reach*  $\tau$  is defined in [27]:

$$\tau := \sup\{r > 0 : \forall \mathbf{y} \in \mathbb{R}^n \text{ with } \text{dist}(\mathbf{y}, \mathbb{M}) < r, \\ \text{there exists a unique } \mathbf{z} \in \mathbb{M} \text{ such that } \|\mathbf{y} - \mathbf{z}\| = \text{dist}(\mathbf{y}, \mathbb{M})\}.$$

A set with positive reach is also called a “prox-regular” set in the literature [74].

The Bi-Lipschitz condition refers to: for some  $0 < \mu \leq L < +\infty$ , it holds that

$$(17) \quad \mu \|\mathbf{y}_1 - \mathbf{y}_2\| \leq \|\mathbf{A}\mathbf{y}_1 - \mathbf{A}\mathbf{y}_2\| \leq L \|\mathbf{y}_1 - \mathbf{y}_2\| \quad \forall \mathbf{y}_1, \mathbf{y}_2 \in \mathbb{M}.$$

According to the definition, it holds that  $0 < \mu \leq L \leq \sigma_{\max} < +\infty$ . This condition ensures  $\mathbf{A}$  can be viewed as an injective mapping when restricted to  $\mathbb{M}$ , which is important for the recovery guarantee.

*Remark for Assumption 3.1.* The assumption that data (particularly images) lies on a smooth manifold has a long and influential history [21, 82], and it is still widely used in recent literature. The compactness of the data manifold can be achieved by standard techniques like normalization. In addition, reach is an important concept for manifold to ensure the uniqueness of its projection [1, 27]. The overall assumptions on manifolds, smoothness, compactness and positive reach, is typically used in recent literature regarding image and signal processing [2, 75, 89]. The on-manifold bi-Lipschitz condition does *not* require  $\mathbf{A}$  to be globally invertible; it merely rules out ill-posedness *restricted to*  $\mathbb{M}$ . This is closely related to Johnson–Lindenstrauss (JL)–type embeddings in compressive sensing: e.g., [5] shows that random matrices are bi-Lipschitz on low-dimensional manifolds with high probability, and JL-style conditions are widely analyzed and used [9, 18, 41, 44, 93].

*Proof of Theorem 3.3.* For simplicity, we first denote the objective functions in (3) as  $F_{1a}(\mathbf{y})$ :

$$F_{1a}(\mathbf{y}) := \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{y}\|^2 + \frac{\alpha}{2} \text{dist}^2(\mathbf{y}, \mathbb{M})$$

Then we introduce some definitions that will be useful in our proof:

$$\mathbb{U}_r(\mathbb{M}) := \{\mathbf{y} \in \mathbb{R}^n : \text{dist}(\mathbf{y}, \mathbb{M}) < r\}, \quad \bar{\mathbb{U}}_r(\mathbb{M}) := \{\mathbf{y} \in \mathbb{R}^n : \text{dist}(\mathbf{y}, \mathbb{M}) \leq r\}$$

Here,  $\mathbb{U}_r(\mathbb{M})$  is an open tubular neighborhood of the manifold  $\mathbb{M}$  and  $\bar{\mathbb{U}}_r(\mathbb{M})$  is its closure. As  $r = \tau$ , the open set  $\mathbb{U}_r(\mathbb{M})$  is named as the *reach tube* of  $\mathbb{M}$ , denoted as  $\mathbb{U}_\tau(\mathbb{M})$ . As introduced in [27], within the reach tube, some nice properties of the distance function and projection mapping can be utilized. For any  $\mathbf{y} \in \mathbb{U}_\tau(\mathbb{M})$  or any  $\mathbf{y} \in \bar{\mathbb{U}}_r(\mathbb{M})$  with  $r < \tau$ , the projection mapping

$$\mathbf{p}(\mathbf{y}) := \arg \min_{\mathbf{z} \in \mathbb{M}} \|\mathbf{z} - \mathbf{y}\|$$

is single valued and well defined, and  $\text{dist}(\mathbf{y}, \mathbb{M}) = \|\mathbf{y} - \mathbf{p}(\mathbf{y})\|$ .

**Step 1: Existence of minimizers of  $F_{1a}$ .** As  $\mathbf{x} \in \mathbb{X}$ , there must be an underlying  $\mathbf{y}_* \in \mathbb{M}$  (hence  $\mathbf{y}_* \in \bar{\mathbb{U}}_r(\mathbb{M})$ ) and  $\mathbf{n}$  such that  $\|\mathbf{x} - \mathbf{A}\mathbf{y}_*\| = \|\mathbf{n}\|$ . Therefore, it holds that

$$F_{1a}(\mathbf{y}_*) = \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{y}_*\|^2 + \frac{\alpha}{2} \text{dist}^2(\mathbf{y}_*, \mathbb{M}) = \frac{1}{2} \|\mathbf{n}\|^2 + 0 = \frac{1}{2} \|\mathbf{n}\|^2$$

In the other hand, for any point outside the tube:  $\mathbf{y} \notin \bar{\mathbb{U}}_r(\mathbb{M})$ , the objective value is lower bounded by:

$$F_{1a}(\mathbf{y}) \geq 0 + \frac{\alpha}{2} \text{dist}^2(\mathbf{y}, \mathbb{M}) > \frac{\alpha}{2} r^2$$

As long as we have large enough  $\alpha$ :

$$(18) \quad \alpha \geq \frac{\|\mathbf{n}\|^2}{r^2},$$

we can ensure  $F_{1a}(\mathbf{y}) > F_{1a}(\mathbf{y}_*)$  for all  $\mathbf{y} \notin \bar{\mathbb{U}}_r(\mathbb{M})$ , which implies  $\inf_{\mathbf{y} \in \mathbb{R}^n} F_{1a}(\mathbf{y}) = \inf_{\mathbf{y} \in \bar{\mathbb{U}}_r(\mathbb{M})} F_{1a}(\mathbf{y})$ . As  $\mathbb{M}$  is compact,  $\bar{\mathbb{U}}_r(\mathbb{M})$  must be compact as well. Consequently, the infimum of  $F$  is attainable, which concludes the existence of the minimizer of  $F_{1a}$ , denoted by  $\hat{\mathbf{y}}$ , and  $\hat{\mathbf{y}} \in \bar{\mathbb{U}}_r(\mathbb{M})$ . Finally, we have the conclusion: It holds for all  $r > 0$  that, condition (18) ensures the existence of  $\hat{\mathbf{y}}$  and  $\hat{\mathbf{y}} \in \bar{\mathbb{U}}_r(\mathbb{M})$ .

**Step 2: Bound of minimizers of  $F_{1a}$ .** For any  $\mathbf{y} \in \mathbb{U}_\tau(\mathbb{M})$ , the projection  $\mathbf{p}(\mathbf{y})$  is uniquely defined, hence we have

$$\begin{aligned} \|\mathbf{A}\mathbf{y} - \mathbf{x}\| &= \|\mathbf{A}\mathbf{y} - \mathbf{A}\mathbf{y}_* - \mathbf{n}\| = \|\mathbf{A}\mathbf{y} - \mathbf{A}\mathbf{p}(\mathbf{y}) + \mathbf{A}\mathbf{p}(\mathbf{y}) - \mathbf{A}\mathbf{y}_* - \mathbf{n}\| \\ &\geq \|\mathbf{A}\mathbf{p}(\mathbf{y}) - \mathbf{A}\mathbf{y}_*\| - \|\mathbf{A}\mathbf{y} - \mathbf{A}\mathbf{p}(\mathbf{y})\| - \|\mathbf{n}\| \\ &\geq \mu \|\mathbf{p}(\mathbf{y}) - \mathbf{y}_*\| - \sigma_{\max} \|\mathbf{y} - \mathbf{p}(\mathbf{y})\| - \|\mathbf{n}\| \end{aligned}$$

According to the conclusion in Step 1, as long as

$$(19) \quad \alpha \geq \frac{\|\mathbf{n}\|^2}{r^2} > \frac{\|\mathbf{n}\|^2}{\tau^2},$$

it holds that the minimizer  $\hat{\mathbf{y}}$  exists and  $\hat{\mathbf{y}} \in \bar{\mathbb{U}}_r(\mathbb{M})$  for some  $r < \tau$  and hence  $\hat{\mathbf{y}} \in \mathbb{U}_\tau(\mathbb{M})$ , which allows us to use the above inequalities at the beginning of Step 2. Now we aim to establish an upper bound for  $\|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\|$  by contradiction. Suppose

$$\mu \|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\| > \sigma_{\max} \|\hat{\mathbf{y}} - \mathbf{p}(\hat{\mathbf{y}})\| + 2\|\mathbf{n}\|$$

we will obtain

$$\|\mathbf{A}\hat{\mathbf{y}} - \mathbf{x}\| \geq \mu \|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\| - \sigma_{\max} \|\hat{\mathbf{y}} - \mathbf{p}(\hat{\mathbf{y}})\| - \|\mathbf{n}\| > \|\mathbf{n}\|,$$

which implies

$$F_{1a}(\hat{\mathbf{y}}) = \frac{1}{2} \|\mathbf{A}\hat{\mathbf{y}} - \mathbf{x}\|^2 + \frac{\alpha}{2} \text{dist}^2(\hat{\mathbf{y}}, \mathbb{M}) > \frac{1}{2} \|\mathbf{n}\|^2 + 0 = F_{1a}(\mathbf{y}_*).$$

This contradicts with the definition of  $\hat{\mathbf{y}}$ : the minimizer of function  $F_{1a}$ . Therefore, we obtain:

$$\mu \|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\| \leq \sigma_{\max} \|\hat{\mathbf{y}} - \mathbf{p}(\hat{\mathbf{y}})\| + 2\|\mathbf{n}\| \leq \sigma_{\max} r + 2\|\mathbf{n}\|$$

which is equivalent to

$$\|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\| \leq \frac{\sigma_{\max}}{\mu} r + \frac{2}{\mu} \|\mathbf{n}\|$$

and implies that

$$(20) \quad \|\hat{\mathbf{y}} - \mathbf{y}_*\| \leq \|\hat{\mathbf{y}} - \mathbf{p}(\hat{\mathbf{y}})\| + \|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\| \leq \left(1 + \frac{\sigma_{\max}}{\mu}\right)r + \frac{2}{\mu}\|\mathbf{n}\|$$

holds for all  $\hat{\mathbf{y}}$  that minimizes  $F_{1a}(\mathbf{y})$ .

**Step 3: Positive definiteness of the Hessian of  $F_{1a}$ .** To prove the uniqueness of the solution, we will establish the strict convexity of the objective function  $F_{1a}(\mathbf{y})$  within a neighborhood around any point of  $\mathbb{M}$ . To achieve this, we establish the positive definiteness of the Hessian of  $F_{1a}(\mathbf{y})$  in this step.

For any  $\mathbf{y} \in \mathbb{U}_\tau(\mathbb{M})$ , the projection mapping is single valued and the objective function can be written as

$$F_{1a}(\mathbf{y}) = \frac{1}{2} \underbrace{\|\mathbf{x} - \mathbf{A}\mathbf{y}\|^2}_{f(\mathbf{y})} + \frac{\alpha}{2} \underbrace{\|\mathbf{y} - \mathbf{p}(\mathbf{y})\|^2}_{g(\mathbf{y})}$$

The smoothness of  $\mathbb{M}$  implies the smoothness of  $g$  and of the projection mapping, and hence we can take first and second orders of derivatives on  $g$  [52, Theorem 2]. Thanks to [27, Theorem 4.8], the gradient and Hessian of  $g$  are given by

$$\nabla g(\mathbf{y}) = 2(\mathbf{y} - \mathbf{p}(\mathbf{y})), \quad \nabla^2 g(\mathbf{y}) = 2(\mathbf{I} - D\mathbf{p}(\mathbf{y})),$$

where  $D\mathbf{p}$  denotes the Jacobian of the projection mapping. The overall Hessian of  $F_{1a}$  is provided by

$$(21) \quad \nabla^2 F_{1a}(\mathbf{y}) = \mathbf{A}^\top \mathbf{A} + \alpha(\mathbf{I} - D\mathbf{p}(\mathbf{y})).$$

To further present the properties of the above Hessian, we introduce a space decomposition according to  $\mathbf{p}(\mathbf{y})$ :

$$\mathbb{R}^n = \mathbb{T}_{\mathbf{p}(\mathbf{y})}(\mathbb{M}) \oplus \mathbb{N}_{\mathbf{p}(\mathbf{y})}(\mathbb{M})$$

where  $\mathbb{T}_{\mathbf{p}(\mathbf{y})}(\mathbb{M})$  denotes the tangent space of  $\mathbb{M}$  at the point  $\mathbf{p}(\mathbf{y}) \in \mathbb{M}$ , and  $\mathbb{N}_{\mathbf{p}(\mathbf{y})}(\mathbb{M})$  represents the normal space. According to [52, Theorem C and Definition 7], the matrix  $D\mathbf{p}(\mathbf{y})$  is actually restricted to the tangent space. In other words, for any decomposition  $\mathbf{h}$  with  $\mathbf{h} = \mathbf{h}_T + \mathbf{h}_N$  where  $\mathbf{h}_T \in \mathbb{T}_{\mathbf{p}(\mathbf{y})}(\mathbb{M})$  and  $\mathbf{h}_N \in \mathbb{N}_{\mathbf{p}(\mathbf{y})}(\mathbb{M})$ , it holds that

$$(22) \quad D\mathbf{p}(\mathbf{y})\mathbf{h}_N = \mathbf{0}, \quad D\mathbf{p}(\mathbf{y})\mathbf{h}_T \in \mathbb{T}_{\mathbf{p}(\mathbf{y})}(\mathbb{M}).$$

In addition, function  $g(\mathbf{y})$  is  $(\frac{s}{\tau-s})$ -weakly convex where  $\tau$  is the reach of  $\mathbb{M}$  and  $s = \text{dist}(\mathbf{y}, \mathbb{M})$  [68, Section 5], and hence the spectrum of  $\nabla^2 g$  can be lower bounded by

$$(23) \quad \langle \mathbf{h}_T, \nabla^2 g(\mathbf{y})\mathbf{h}_T \rangle \geq -\frac{2s}{\tau-s} \|\mathbf{h}_T\|^2,$$

Now, let's turn to the first term in the Hessian:  $\mathbf{A}^\top \mathbf{A}$ . It can be shown using the JL condition (17) that, the spectrum of  $\mathbf{A}^\top \mathbf{A}$  restricted to the tangent space can also be lower bounded. In particular, we pick an arbitrary tangent vector  $\mathbf{h}_T \in \mathbb{T}_{\mathbf{p}(\mathbf{y})}(\mathbb{M})$ . According to the definition of tangent space, there must be a curve  $\gamma : (-\delta, \delta) \rightarrow \mathbb{M}$  with  $\delta > 0$ ,  $\gamma(0) = \mathbf{p}(\mathbf{y})$ , and  $\gamma'(0) = \mathbf{h}_T$ . For any  $0 \leq t < \delta$ ,  $\gamma(t) \in \mathbb{M}$ . By applying condition (17) with the pair  $(\gamma(t), \gamma(0))$  and divide by  $t^2$ , we have

$$\mu^2 \frac{\|\gamma(t) - \gamma(0)\|^2}{t^2} \leq \frac{\|\mathbf{A}\gamma(t) - \mathbf{A}\gamma(0)\|^2}{t^2} \leq L^2 \frac{\|\gamma(t) - \gamma(0)\|^2}{t^2}$$

By differentiability and the continuity of the operator  $\mathbf{A}$ , it holds that

$$\lim_{t \rightarrow 0} \frac{\gamma(t) - \gamma(0)}{t} = \mathbf{h}_T, \quad \lim_{t \rightarrow 0} \frac{\mathbf{A}\gamma(t) - \mathbf{A}\gamma(0)}{t} = \mathbf{A}\mathbf{h}_T$$

which implies

$$(24) \quad \mu^2 \|\mathbf{h}_T\|^2 \leq \|\mathbf{A}\mathbf{h}_T\|^2 \leq L^2 \|\mathbf{h}_T\|^2.$$

Combining (21), (22), (23), and (24), we have

$$\langle \mathbf{h}, \nabla^2 F_{1a}(\mathbf{y})\mathbf{h} \rangle$$

$$\begin{aligned}
 &= \underbrace{\langle \mathbf{h}_T, \mathbf{A}^\top \mathbf{A} \mathbf{h}_T \rangle}_{\geq \mu^2 \|\mathbf{h}_T\|^2} + 2 \underbrace{\langle \mathbf{h}_T, \mathbf{A}^\top \mathbf{A} \mathbf{h}_N \rangle}_{\geq 0} + \underbrace{\langle \mathbf{h}_N, \mathbf{A}^\top \mathbf{A} \mathbf{h}_N \rangle}_{\geq 0} \\
 &\quad + \alpha \underbrace{\langle \mathbf{h}_T, (\mathbf{I} - D\mathbf{p}(\mathbf{y})) \mathbf{h}_T \rangle}_{\geq -\frac{s}{\tau-s} \|\mathbf{h}_T\|^2} + 2\alpha \underbrace{\langle \mathbf{h}_T, (\mathbf{I} - D\mathbf{p}(\mathbf{y})) \mathbf{h}_N \rangle}_{=\langle \mathbf{h}_T, \mathbf{h}_N \rangle=0} + \alpha \underbrace{\langle \mathbf{h}_N, (\mathbf{I} - D\mathbf{p}(\mathbf{y})) \mathbf{h}_N \rangle}_{=\|\mathbf{h}_N\|^2} \\
 &\geq \left( \mu^2 - \alpha \frac{s}{\tau-s} \right) \|\mathbf{h}_T\|^2 + \alpha \|\mathbf{h}_N\|^2 - 2 \|\mathbf{A} \mathbf{h}_T\| \cdot \|\mathbf{A} \mathbf{h}_N\| \\
 &\geq \left( \mu^2 - \alpha \frac{s}{\tau-s} \right) \|\mathbf{h}_T\|^2 + \alpha \|\mathbf{h}_N\|^2 - 2L \|\mathbf{h}_T\| \cdot \sigma_{\max} \|\mathbf{h}_N\| \\
 &= \begin{bmatrix} \|\mathbf{h}_T\| & \|\mathbf{h}_N\| \end{bmatrix} \begin{bmatrix} \mu^2 - \alpha \frac{s}{\tau-s} & -\sigma_{\max} L \\ -\sigma_{\max} L & \alpha \end{bmatrix} \begin{bmatrix} \|\mathbf{h}_T\| \\ \|\mathbf{h}_N\| \end{bmatrix}
 \end{aligned}$$

Therefore, to ensure  $\langle \mathbf{h}, \nabla^2 F_{1a}(\mathbf{y}) \mathbf{h} \rangle > 0$  for any  $\mathbf{h} \neq \mathbf{0}$ , it's enough to ensure the  $2 \times 2$  matrix to be positive definite:

$$(25) \quad \mu^2 - \alpha \frac{s}{\tau-s} > 0 \quad \text{and} \quad \alpha \left( \mu^2 - \alpha \frac{s}{\tau-s} \right) - \sigma_{\max}^2 L^2 > 0.$$

In other words, (25) will guarantee the positive definiteness of  $\nabla^2 F_{1a}(\mathbf{y})$  for all  $\mathbf{y} \in \bar{\mathbb{U}}_s(\mathbb{M})$  and any  $s < \tau$ .

**Step 4: Uniqueness of minimizers of  $F_{1a}$ .** In this step, we will combine the results from Steps 2 and 3. Then we are able to prove that the objective function  $F_{1a}(\mathbf{y})$  is strictly convex in a neighborhood of its minimizers, which implies the uniqueness of the minimizer. To achieve this, it's enough to ensure

$$(26) \quad \|\hat{\mathbf{y}} - \mathbf{y}_*\| \leq s$$

for all  $\hat{\mathbf{y}} \in \arg \min_{\mathbf{y}} F_{1a}(\mathbf{y})$ , where  $s$  satisfies (25). With this condition (26), it holds that

$$\hat{\mathbf{y}} \in \mathbb{B}(\mathbf{y}_*, s) \subset \bar{\mathbb{U}}_s(\mathbb{M}).$$

Along with the fact that  $\mathbb{B}(\mathbf{y}_*, s)$  is convex and that  $\nabla^2 F_{1a}(\mathbf{y})$  is positive definite for all  $\mathbf{y} \in \bar{\mathbb{U}}_s(\mathbb{M})$ ,  $F_{1a}$  is strictly convex within  $\mathbb{B}(\mathbf{y}_*, s)$  [7, Section 3.1.4]. As all minimizers of the strict convex function belong to this convex set,  $\mathbb{B}(\mathbf{y}_*, s)$ , the minimizer  $\hat{\mathbf{y}}$  must be unique.

Now the question is: How to guarantee (26)? According to (20), Condition (19) along with

$$(27) \quad \left( 1 + \frac{\sigma_{\max}}{\mu} \right) r + \frac{2}{\mu} \|\mathbf{n}\| \leq s$$

can guarantee (26). Finally, it's enough to choose  $\alpha$ ,  $s$ , and  $r$  such that (19), (25), and (27) are satisfied together. In particular, we choose

$$s = \frac{4}{\mu} \|\mathbf{n}\|, \quad r = \frac{1}{\sigma_{\max}} \|\mathbf{n}\|, \quad \alpha = \frac{2\sigma_{\max}^2 L^2}{\mu^2}$$

where  $\alpha$  merely depends on  $\mathbf{A}$  and  $\mathbb{M}$  but is independent of  $\mathbf{x}$ . Such a parameter choice implies (27):

$$\left( 1 + \frac{\sigma_{\max}}{\mu} \right) r + \frac{2}{\mu} \|\mathbf{n}\| \leq 2 \frac{\sigma_{\max}}{\mu} r + \frac{2}{\mu} \|\mathbf{n}\| = \frac{2}{\mu} \|\mathbf{n}\| + \frac{2}{\mu} \|\mathbf{n}\| = s.$$

As  $\|\mathbf{n}\| < \frac{1}{20} \frac{\mu^5}{\sigma_{\max}^2 L^2} \tau$ , it holds that

$$s = \frac{4}{\mu} \|\mathbf{n}\| < \frac{\mu^4}{5\sigma_{\max}^2 L^2} \tau \implies \frac{s}{\tau-s} < \frac{\frac{4}{5\sigma_{\max}^2 L^2} \tau}{\tau - \frac{4}{5\sigma_{\max}^2 L^2} \tau} \leq \frac{\frac{\mu^4}{5\sigma_{\max}^2 L^2} \tau}{\tau - \frac{1}{5} \tau} = \frac{\mu^4}{4\sigma_{\max}^2 L^2}$$

and therefore (25) is satisfied:

$$\mu^2 - \alpha \frac{s}{\tau-s} > \mu^2 - \frac{2\sigma_{\max}^2 L^2}{\mu^2} \frac{\mu^4}{4\sigma_{\max}^2 L^2} = \frac{1}{2} \mu^2 > 0$$

and

$$\alpha \left( \mu^2 - \alpha \frac{s}{\tau - s} \right) > \frac{2\sigma_{\max}^2 L^2}{\mu^2} \cdot \frac{1}{2} \mu^2 = \sigma_{\max}^2 L^2.$$

Finally, by choosing  $\alpha$  as before, condition (19) is satisfied:

$$\alpha = 2\sigma_{\max}^2 \cdot \frac{L^2}{\mu^2} \geq 2\sigma_{\max}^2 = \frac{2\|\mathbf{n}\|^2}{r^2}, \quad r = \frac{1}{\sigma_{\max}} \|\mathbf{n}\| < \frac{1}{\sigma_{\max}} \cdot \frac{1}{20} \frac{\mu^5}{\sigma_{\max}^2 L^2}^\tau < \tau,$$

which finishes the proof of the uniqueness of minimizers of  $F_{1a}$ .

**Step 5: Local Lipschitz continuity of  $F_{1a}$ .** Previous results from Steps 1-4 indicate that, for any  $\mathbf{x} \in \mathbb{X}$ , there is a unique  $\hat{\mathbf{y}}(\mathbf{x})$  that minimizes  $F_{1a}$ , but the continuity of  $\hat{\mathbf{y}}$  w.r.t.  $\mathbf{x}$  has not been established. In this step, we will show this continuity via the implicit function theorem. Firstly, as  $\hat{\mathbf{y}}$  minimizes  $F_{1a}$ , by first-order optimality conditions for smooth minimization, it holds that

$$\nabla F_{1a}(\hat{\mathbf{y}}) = \underbrace{\mathbf{A}^\top (\mathbf{A}\hat{\mathbf{y}} - \mathbf{x}) + \alpha(\hat{\mathbf{y}} - \mathbf{p}(\hat{\mathbf{y}}))}_{=: \mathcal{H}(\mathbf{x}, \hat{\mathbf{y}})} = \mathbf{0}$$

Now, let's pick a point  $\mathbf{x}_0$  from  $\mathbb{X}$ . Previous results from Steps 1-4 indicate that, operator  $\mathcal{H}(\mathbf{x}, \mathbf{y})$  is continuously differentiable within a neighborhood of  $(\mathbf{x}_0, \hat{\mathbf{y}}(\mathbf{x}_0))$ , and its Jacobian matrix w.r.t.  $\mathbf{y}$

$$D_{\mathbf{y}} \mathcal{H}(\mathbf{x}, \mathbf{y}) = \nabla^2 F_{1a}(\mathbf{y})$$

is positive definite within that neighborhood of  $(\mathbf{x}_0, \hat{\mathbf{y}}(\mathbf{x}_0))$ . Therefore, we are able to apply the implicit function theorem [28, Theorem 3.9] and conclude that  $\hat{\mathbf{y}}(\mathbf{x})$  is Lipschitz continuous within a neighborhood of  $\mathbf{x}_0$ . This argument applies for any points  $\mathbf{x}_0$  in  $\mathbb{X}$ . Therefore,  $\hat{\mathbf{y}} = \mathcal{F}_{1a}(\mathbf{x})$  is locally Lipschitz continuous on  $\mathbb{X}$ .  $\square$

The proof line of Theorem 3.4 largely follows the proof of Theorem 3.3. Here we will highlight the difference of proofs between the two theorems, so that Theorem 3.4 will be rigorously proved without too much redundancy.

*Proof of Theorem 3.4.* For simplicity, we denote the objective function in (4) as  $F_{1b}(\mathbf{y}, \mathbf{z})$ :

$$F_{1b}(\mathbf{y}, \mathbf{z}) := \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{y}\|^2 + \frac{\alpha}{2} \text{dist}^2(\mathbf{z}, \mathbb{M}) + \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|^2,$$

and we will study its properties analogously to  $F_{1a}$ .

**Step 1: Existence of minimizers of  $F_{1b}$ .** For any  $r > 0$ , as

$$\alpha \geq \frac{\|\mathbf{n}\|^2}{r^2}, \quad \beta \geq \frac{\|\mathbf{n}\|^2}{r^2},$$

it holds that

$$(28) \quad \inf_{\mathbf{y}, \mathbf{z}} F_{1b}(\mathbf{y}, \mathbf{z}) = \inf_{(\mathbf{y}, \mathbf{z}): \text{dist}(\mathbf{z}, \mathbb{M}) \leq r \text{ and } \|\mathbf{z} - \mathbf{y}\| \leq r} F_{1b}(\mathbf{y}, \mathbf{z}).$$

This can be proved by contradiction: (I) Suppose  $F_{1b}(\hat{\mathbf{y}}, \hat{\mathbf{z}})$  is lower than the right-hand-side of (28) and  $\text{dist}(\hat{\mathbf{z}}, \mathbb{M}) > r$ , we have

$$F_{1b}(\hat{\mathbf{y}}, \hat{\mathbf{z}}) \geq 0 + \frac{\|\mathbf{n}\|^2}{2r^2} \text{dist}^2(\hat{\mathbf{z}}, \mathbb{M}) + 0 > \frac{1}{2} \|\mathbf{n}\|^2 = F_{1b}(\mathbf{y}_*, \mathbf{y}_*)$$

which contradicts with the hypothesis regarding  $(\hat{\mathbf{y}}, \hat{\mathbf{z}})$ . (II) Suppose  $F_{1b}(\hat{\mathbf{y}}, \hat{\mathbf{z}})$  is lower than the right-hand-side of (28) and  $\|\hat{\mathbf{z}} - \hat{\mathbf{y}}\| > r$ , we have

$$F_{1b}(\hat{\mathbf{y}}, \hat{\mathbf{z}}) \geq 0 + 0 + \frac{\|\mathbf{n}\|^2}{2r^2} \|\hat{\mathbf{z}} - \hat{\mathbf{y}}\|^2 > \frac{1}{2} \|\mathbf{n}\|^2 = F_{1b}(\mathbf{y}_*, \mathbf{y}_*)$$

which also derives a contradiction. Arguments in (I) and (II) together prove (28). Similar to the proof of Theorem 3.3, (28) implies the existence of minimizers of  $F_{1b}$  (i.e., minimizers are attainable.)

**Step 2: Bound of minimizers of  $F_{1b}$ .** To extend the proof regarding  $F_{1a}$  to  $F_{1b}$ , we consider the following inequality that holds for all  $\mathbf{y}, \mathbf{z} \in \mathbb{U}_\tau(\mathbb{M})$

$$\|\mathbf{y} - \mathbf{p}(\mathbf{y})\| \leq \|\mathbf{y} - \mathbf{p}(\mathbf{z})\| \leq \|\mathbf{y} - \mathbf{z}\| + \|\mathbf{z} - \mathbf{p}(\mathbf{z})\| = \|\mathbf{y} - \mathbf{z}\| + \text{dist}(\mathbf{z}, \mathbb{M}) \leq 2r.$$

Therefore, we need  $2r < \tau$  and

$$(29) \quad \alpha \geq \frac{\|\mathbf{n}\|^2}{r^2} > \frac{4\|\mathbf{n}\|^2}{\tau^2}, \quad \beta \geq \frac{\|\mathbf{n}\|^2}{r^2} > \frac{4\|\mathbf{n}\|^2}{\tau^2}$$

to ensure  $\hat{\mathbf{y}}, \hat{\mathbf{z}} \in \mathbb{U}_\tau(\mathbb{M})$ . Following the same argument as the proof of Theorem 3.3, the above condition (29) implies

$$\|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\| \leq \frac{\sigma_{\max}}{\mu}(2r) + \frac{2}{\mu}\|\mathbf{n}\|$$

and hence

$$(30) \quad \|\hat{\mathbf{y}} - \mathbf{y}_*\| \leq \|\hat{\mathbf{y}} - \mathbf{p}(\hat{\mathbf{y}})\| + \|\mathbf{p}(\hat{\mathbf{y}}) - \mathbf{y}_*\| \leq 2 \left(1 + \frac{\sigma_{\max}}{\mu}\right) r + \frac{2}{\mu}\|\mathbf{n}\|$$

and

$$(31) \quad \|\hat{\mathbf{z}} - \mathbf{y}_*\| \leq \|\hat{\mathbf{z}} - \hat{\mathbf{y}}\| + \|\hat{\mathbf{y}} - \mathbf{y}_*\| \leq \left(3 + 2\frac{\sigma_{\max}}{\mu}\right) r + \frac{2}{\mu}\|\mathbf{n}\|$$

holds for all  $(\hat{\mathbf{y}}, \hat{\mathbf{z}})$  that minimizes  $F_{1b}(\mathbf{y}, \mathbf{z})$ .

**Step 3: Positive definiteness of the Hessian of  $F_{1b}$ .** Function  $F_{1b}(\mathbf{y}, \mathbf{z})$ 's Hessian matrix is of size  $2n \times 2n$  and can be written as a  $2 \times 2$  block w.r.t.  $\mathbf{y}$  and  $\mathbf{z}$ :

$$\nabla^2 F_{1b}(\mathbf{y}, \mathbf{z}) = \begin{bmatrix} \mathbf{A}^\top \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \alpha \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - D\mathbf{p}(\mathbf{z}) \end{bmatrix} + \beta \begin{bmatrix} \mathbf{I} & -\mathbf{I} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix}$$

For any  $\mathbf{h} = [\mathbf{u}^\top \ \mathbf{v}^\top]^\top \in \mathbb{R}^{2n}$ , the quadratic form  $\langle \mathbf{h}, \nabla^2 F_{1b}(\mathbf{y}, \mathbf{z}) \mathbf{h} \rangle$  can be calculated through:

$$\langle \mathbf{h}, \nabla^2 F_{1b}(\mathbf{y}, \mathbf{z}) \mathbf{h} \rangle = \mathbf{u}^\top \mathbf{A}^\top \mathbf{A} \mathbf{u} + \alpha \mathbf{v}^\top (\mathbf{I} - D\mathbf{p}(\mathbf{z})) \mathbf{v} + \beta \|\mathbf{u} - \mathbf{v}\|^2$$

Decompose  $\mathbf{u} = \mathbf{u}_T + \mathbf{u}_N$  and  $\mathbf{v} = \mathbf{v}_T + \mathbf{v}_N$  in  $\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M}) \oplus \mathbb{N}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})$ . Using the same argument as the proof of Theorem 3.3, we have

$$\begin{aligned} & \langle \mathbf{h}, \nabla^2 F_{1b}(\mathbf{y}, \mathbf{z}) \mathbf{h} \rangle \\ & \geq \left( \mu^2 \|\mathbf{u}_T\|^2 - 2\sigma_{\max} L \|\mathbf{u}_T\| \|\mathbf{u}_N\| \right) + \alpha \left( -\frac{s}{\tau - s} \|\mathbf{v}_T\|^2 + \|\mathbf{v}_N\|^2 \right) + \beta \|\mathbf{u} - \mathbf{v}\|^2 \end{aligned}$$

which implies

$$\begin{aligned} & \langle \mathbf{h}, \nabla^2 F_{1b}(\mathbf{y}, \mathbf{z}) \mathbf{h} \rangle \\ & \geq \left( \mu^2 \|\mathbf{u}_T\|^2 - 2\sigma_{\max} L \|\mathbf{u}_T\| \|\mathbf{u}_N\| \right) \\ & \quad + \alpha \left( -\frac{s}{\tau - s} \|\mathbf{v}_T\|^2 + \|\mathbf{v}_N\|^2 \right) + \beta \left( \|\mathbf{u}_T - \mathbf{v}_T\|^2 + \|\mathbf{u}_N - \mathbf{v}_N\|^2 \right) \\ & \geq \left( \mu^2 \|\mathbf{u}_T\|^2 - 2\sigma_{\max} L \|\mathbf{u}_T\| \|\mathbf{u}_N\| \right) \\ & \quad + \alpha \left( -\frac{s}{\tau - s} \|\mathbf{v}_T\|^2 + \|\mathbf{v}_N\|^2 \right) + \beta \left( (\|\mathbf{u}_T\| - \|\mathbf{v}_T\|)^2 + (\|\mathbf{u}_N\| - \|\mathbf{v}_N\|)^2 \right) \\ & = [\|\mathbf{u}_T\| \quad \|\mathbf{u}_N\| \quad \|\mathbf{v}_T\| \quad \|\mathbf{v}_N\|] \underbrace{\begin{bmatrix} \mu^2 + \beta & -\sigma_{\max} L & -\beta & -\beta \\ -\sigma_{\max} L & \beta & & -\beta \\ -\beta & & \beta - \alpha \frac{s}{\tau - s} & \\ & -\beta & & \alpha + \beta \end{bmatrix}}_{=: \mathbf{B}} \begin{bmatrix} \|\mathbf{u}_T\| \\ \|\mathbf{u}_N\| \\ \|\mathbf{v}_T\| \\ \|\mathbf{v}_N\| \end{bmatrix} \end{aligned}$$

To ensure the positive definiteness of  $\nabla^2 F_{1b}(\mathbf{y}, \mathbf{z})$ , it's enough to ensure  $\mathbf{B} \succ \mathbf{0}$ . For simplicity, we define

$$\theta := \alpha \frac{s}{\tau - s}, \quad \mathbf{B}_1 := \begin{bmatrix} \mu^2 + \beta & -\sigma_{\max} L \\ -\sigma_{\max} L & \beta \end{bmatrix} \quad \mathbf{B}_2 := \begin{bmatrix} -\beta & \\ & -\beta \end{bmatrix} \quad \mathbf{B}_3 := \begin{bmatrix} \beta - \theta & \\ & \alpha + \beta \end{bmatrix}$$

Then  $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_2 \\ \mathbf{B}_2^\top & \mathbf{B}_3 \end{bmatrix}$  is positive definite if and only if  $\mathbf{B}_3$  and its Schur complement  $\mathbf{S}$  are both positive definite:

$$\mathbf{B}_3 \succ \mathbf{0}, \quad \mathbf{S} = \mathbf{B}_1 - \mathbf{B}_2 \mathbf{B}_3^{-1} \mathbf{B}_2^\top \succ \mathbf{0}$$

As  $\mathbf{B}_2$  and  $\mathbf{B}_3$  are both diagonal, so  $\mathbf{B}_2 \mathbf{B}_3^{-1} \mathbf{B}_2^\top$  is straight forward to calculate:  $\mathbf{B}_2 \mathbf{B}_3^{-1} \mathbf{B}_2^\top = \text{diag}\left(\frac{\beta^2}{\beta-\theta}, \frac{\beta^2}{\alpha+\beta}\right)$ . Then the Schur complement can be calculated:

$$\mathbf{S} = \begin{bmatrix} \mu^2 + \beta - \frac{\beta^2}{\beta-\theta} & -\sigma_{\max} L \\ -\sigma_{\max} L & \beta - \frac{\beta^2}{\alpha+\beta} \end{bmatrix} = \begin{bmatrix} \mu^2 - \frac{\beta\theta}{\beta-\theta} & -\sigma_{\max} L \\ -\sigma_{\max} L & \frac{\alpha\beta}{\alpha+\beta} \end{bmatrix}$$

Note that  $\mathbf{B}_3 \succ \mathbf{0}$  if.f  $\beta > \theta$ . Therefore,  $\mathbf{B} \succ \mathbf{0}$  if.f.

$$(32) \quad \beta > \theta, \quad \mu^2 > \frac{\beta\theta}{\beta-\theta}, \quad \left(\mu^2 - \frac{\beta\theta}{\beta-\theta}\right) \frac{\alpha\beta}{\alpha+\beta} > \sigma_{\max}^2 L^2,$$

where  $\theta = \alpha \frac{s}{\tau-s}$ . Finally, we obtain that (32) ensures  $\nabla^2 F_{1b}(\mathbf{y}, \mathbf{z}) \succ \mathbf{0}$  for all  $\mathbf{y} \in \mathbb{R}^n$  and all  $\mathbf{z} \in \bar{\mathbb{U}}_s(\mathbb{M})$  with  $s < \tau$ .

**Step 4: Uniqueness of minimizers of  $F_{1b}$ .** Comparable to the Step 4 in Theorem 3.3, we need  $\|\hat{\mathbf{z}} - \mathbf{y}_*\| \leq s$  for all  $(\hat{\mathbf{y}}, \hat{\mathbf{z}}) \in \arg \min F_{1b}(\mathbf{y}, \mathbf{z})$ . Based on (31), it's enough to guarantee

$$(33) \quad \left(3 + 2 \frac{\sigma_{\max}}{\mu}\right) r + \frac{2}{\mu} \|\mathbf{n}\| \leq s$$

Now we choose

$$s = \frac{4}{\mu} \|\mathbf{n}\|, \quad r = \frac{2}{5\sigma_{\max}} \|\mathbf{n}\|$$

which directly satisfies (33). As  $\|\mathbf{n}\| < \frac{1}{76} \frac{\mu^5}{\sigma_{\max}^2} L^2 \tau$ , we have

$$s = \frac{4}{\mu} \|\mathbf{n}\| < \frac{1}{19} \frac{\mu^4}{\sigma_{\max}^2} L^2 \tau, \quad \frac{s}{\tau-s} < \frac{\frac{1}{19} \frac{\mu^4}{\sigma_{\max}^2} L^2 \tau}{\tau - \frac{1}{19} \frac{\mu^4}{\sigma_{\max}^2} L^2 \tau} \leq \frac{\frac{1}{19} \frac{\mu^4}{\sigma_{\max}^2} L^2 \tau}{\tau - \frac{1}{19} \tau} = \frac{1}{18} \frac{\mu^4}{\sigma_{\max}^2} L^2$$

As long as we take

$$\alpha = \frac{9\sigma_{\max}^2 L^2}{\mu^2}, \quad \beta \geq \max\left(\alpha, \frac{3}{2}\mu^2\right)$$

it holds that

$$\theta = \alpha \frac{s}{\tau-s} < \frac{9\sigma_{\max}^2 L^2}{\mu^2} \frac{1}{18} \frac{\mu^4}{\sigma_{\max}^2} L^2 = \frac{1}{2}\mu^2$$

which implies  $\beta > 3\theta$  and hence  $\beta > \theta$ . Moreover, we can verify the remaining part of (32):

$$\begin{aligned} \frac{\beta\theta}{\beta-\theta} &< \frac{\beta\theta}{\beta-\beta/3} = \frac{3}{2}\theta < \frac{3}{4}\mu^2 < \mu^2, \\ \left(\mu^2 - \frac{\beta\theta}{\beta-\theta}\right) \frac{\alpha\beta}{\alpha+\beta} &> \left(\mu^2 - \frac{3}{4}\mu^2\right) \frac{\alpha\beta}{\beta+\beta} = \frac{1}{8}\mu^2\alpha = \frac{1}{8}\mu^2 \cdot \frac{9\sigma_{\max}^2 L^2}{\mu^2} > \sigma_{\max}^2 L^2. \end{aligned}$$

which finishes the proof of (32). Finally, it's enough to verify (29):

$$2r \leq \frac{\|\mathbf{n}\|}{\sigma_{\max}} \leq \frac{1}{76} \frac{\mu^5}{\sigma_{\max}^3} L^2 \tau < \tau, \quad \frac{\|\mathbf{n}\|^2}{r^2} = \frac{25}{4} \sigma_{\max}^2 \leq \alpha \leq \beta,$$

which finishes Step 4, and concludes the uniqueness of  $(\hat{\mathbf{y}}, \hat{\mathbf{z}})$ .

**Step 5: Local Lipschitz continuity of  $\mathcal{F}_{1b}$ .** By largely following Step 5 in the proof of Theorem 3.3 and changing  $\nabla^2 F_{1a}(\mathbf{y})$  to  $\nabla^2 F_{1b}(\mathbf{y}, \mathbf{z})$ , one can directly conclude that the mapping  $\mathcal{F}_{1b}$  is locally Lipschitz continuous on  $\mathbb{X}$ .  $\square$

**C.1. Proximal operator near a manifold.** We collect here the definition and basic properties of the proximal map used in the main text and relate them to the convergence condition proposed in [83].

**Theorem C.1** (Contractivity of the proximal residual near a  $\mathcal{C}^2$  manifold). *Let  $\mathbb{M} \subset \mathbb{R}^n$  be a compact  $\mathcal{C}^2$  embedded submanifold with reach  $\tau > 0$ . For  $\sigma > 0$  define, for each  $\mathbf{z} \in \mathbb{U}_\tau(\mathbb{M})$ ,*

$$\phi_\sigma(\mathbf{y}, \mathbf{z}) := \frac{\sigma}{2} \text{dist}^2(\mathbf{y}, \mathbb{M}) + \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|^2.$$

*Then  $\phi_\sigma$  must yield a unique minimizer, and hence we are able to define*

$$\text{prox}_\sigma(\mathbf{z}) := \arg \min_{\mathbf{y}} \phi_\sigma(\mathbf{y}, \mathbf{z}), \quad \mathcal{S}_\sigma(\mathbf{z}) := \text{prox}_\sigma(\mathbf{z}) - \mathbf{z}.$$

*Then  $\mathcal{S}_\sigma$  is a contractive operator within a tubular neighborhood of  $\mathbb{M}$ . In particular, it holds that*

$$(34) \quad \|\mathcal{S}_\sigma(\mathbf{z}) - \mathcal{S}_\sigma(\mathbf{z}')\| \leq \frac{\sigma}{1 + \sigma} \|\mathbf{z} - \mathbf{z}'\|$$

*for all  $\mathbf{z}, \mathbf{z}' \in \mathbb{U}_r(\mathbb{M})$  where  $r \leq \tau/4$  and  $\|\mathbf{z} - \mathbf{z}'\| \leq \tau/4$ .*

Relation to plug-and-play (PnP): Condition (A) of [83] assumes a (nearly) contractive *denoiser residual*—precisely the kind of property (34) guarantees for the proximal residual  $\text{prox}_\sigma - \mathbf{I}$  on a neighborhood of  $\mathbb{M}$ . In practice,  $\mathbb{M}$  is unknown; one therefore learns a parameterized operator (e.g., a neural network) whose residual is constrained to be (nearly)  $\sigma$ -contractive and plugs it into PGD/HQS in place of the exact proximal map. Whereas [83] posits Condition (A) to ensure convergence, Theorem C.1 shows this condition arises naturally when the prior corresponds to the manifold-penalty  $\frac{\sigma}{2} \text{dist}^2(\cdot, \mathbb{M})$ .

*Proof of Theorem C.1.* We first note that, for any  $\mathbf{y}$ , if  $\|\mathbf{y} - \mathbf{z}\| > \|\mathbf{z} - \mathbf{p}(\mathbf{z})\|$ , then it holds that

$$\phi_\sigma(\mathbf{p}(\mathbf{z}), \mathbf{z}) = 0 + \frac{1}{2} \|\mathbf{z} - \mathbf{p}(\mathbf{z})\|^2 < \frac{\sigma}{2} \text{dist}^2(\mathbf{y}, \mathbb{M}) + \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|^2 = \phi_\sigma(\mathbf{y}, \mathbf{z})$$

which implies

$$\inf_{\mathbf{y}} \phi_\sigma(\mathbf{y}, \mathbf{z}) = \inf_{\mathbf{y}: \|\mathbf{y} - \mathbf{z}\| \leq \|\mathbf{z} - \mathbf{p}(\mathbf{z})\|} \phi_\sigma(\mathbf{y}, \mathbf{z})$$

Let  $r = \|\mathbf{z} - \mathbf{p}(\mathbf{z})\|$ . We further notice that, for any  $\mathbf{y}$  with  $\|\mathbf{y} - \mathbf{z}\| = s \leq r$ , we are able to define  $\tilde{\mathbf{y}}$

$$\tilde{\mathbf{y}} := \frac{r-s}{r} \mathbf{z} + \frac{s}{r} \mathbf{p}(\mathbf{z})$$

which satisfies  $\mathbf{p}(\tilde{\mathbf{y}}) = \mathbf{p}(\mathbf{z})$  and hence it holds that

$$\begin{aligned} \text{dist}(\tilde{\mathbf{y}}, \mathbb{M}) &= \|\tilde{\mathbf{y}} - \mathbf{p}(\mathbf{z})\| = \|\mathbf{z} - \mathbf{p}(\mathbf{z})\| - \|\tilde{\mathbf{y}} - \mathbf{z}\| \\ &< \|\mathbf{z} - \mathbf{p}(\mathbf{y})\| - \|\tilde{\mathbf{y}} - \mathbf{z}\| \\ &\leq \|\mathbf{z} - \mathbf{y}\| + \|\mathbf{y} - \mathbf{p}(\mathbf{y})\| - \|\tilde{\mathbf{y}} - \mathbf{z}\| \\ &= s + \|\mathbf{y} - \mathbf{p}(\mathbf{y})\| - s = \|\mathbf{y} - \mathbf{p}(\mathbf{y})\| = \text{dist}(\mathbf{y}, \mathbb{M}) \end{aligned}$$

which implies

$$\phi_\sigma(\tilde{\mathbf{y}}, \mathbf{z}) = \frac{\sigma}{2} \text{dist}^2(\tilde{\mathbf{y}}, \mathbb{M}) + \frac{1}{2} \|\tilde{\mathbf{y}} - \mathbf{z}\|^2 < \frac{\sigma}{2} \text{dist}^2(\mathbf{y}, \mathbb{M}) + \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|^2 = \phi_\sigma(\mathbf{y}, \mathbf{z})$$

Consequently, we conclude that minimizing  $\phi_\sigma$  is equal to minimizing it over the line segment between  $\mathbf{z}$  and its projection  $\mathbf{p}(\mathbf{z})$ :

$$\inf_{\mathbf{y}} \phi_\sigma(\mathbf{y}, \mathbf{z}) = \inf_{\xi \in [0,1]} \phi_\sigma(\xi \mathbf{z} + (1 - \xi) \mathbf{p}(\mathbf{z}), \mathbf{z}).$$

Now define  $\psi(\xi) = \phi_\sigma(\xi \mathbf{z} + (1 - \xi)\mathbf{p}(\mathbf{z}), \mathbf{z})$ . We have

$$\begin{aligned}\psi(\xi) &= \frac{\sigma}{2} \left\| \left( \xi \mathbf{z} + (1 - \xi)\mathbf{p}(\mathbf{z}) \right) - \mathbf{p}(\mathbf{z}) \right\|^2 + \frac{1}{2} \left\| \left( \xi \mathbf{z} + (1 - \xi)\mathbf{p}(\mathbf{z}) \right) - \mathbf{z} \right\|^2 \\ &= \frac{\sigma}{2} \xi^2 \|\mathbf{z} - \mathbf{p}(\mathbf{z})\|^2 + \frac{1}{2} (1 - \xi)^2 \|\mathbf{z} - \mathbf{p}(\mathbf{z})\|^2 \\ &= \left( \sigma \xi^2 + (1 - \xi)^2 \right) \cdot \frac{1}{2} \|\mathbf{z} - \mathbf{p}(\mathbf{z})\|^2\end{aligned}$$

Therefore,  $\inf_{\xi \in [0,1]} \psi(\xi)$  is attainable, and the minimizer is  $\xi_* = \frac{1}{1+\sigma}$ , which implies  $\phi_\sigma$  must yield a unique minimizer at

$$\mathbf{y}_* = \frac{\mathbf{z} + \sigma \mathbf{p}(\mathbf{z})}{1 + \sigma}.$$

Consequently, we have

$$\mathcal{S}_\sigma(\mathbf{z}) = \mathbf{y}_* - \mathbf{z} = \frac{\sigma}{1 + \sigma} (\mathbf{p}(\mathbf{z}) - \mathbf{z})$$

and hence

$$D\mathcal{S}_\sigma(\mathbf{z}) = \frac{\sigma}{1 + \sigma} (D\mathbf{p}(\mathbf{z}) - \mathbf{I}).$$

According to [52, Theorem C],  $D\mathbf{p}(\mathbf{z})$  is actually restricted to the tangent space  $\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})$ :

$$D\mathbf{p}(\mathbf{z}) = \left( \mathbf{I}_{\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})} - r \mathcal{L}_{\mathbf{p}(\mathbf{z}), \mathbf{v}} \right)^{-1} P_{\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})}$$

where  $r = \|\mathbf{p}(\mathbf{z}) - \mathbf{z}\|$ ,  $\mathbf{v} = (\mathbf{p}(\mathbf{z}) - \mathbf{z})/r$ , and  $\mathcal{L}_{\mathbf{p}(\mathbf{z}), \mathbf{v}}$  is the shape operator in direction  $\mathbf{v}$  at  $\mathbf{p}(\mathbf{z})$ . The shape operator's eigenvalues  $\kappa_1, \dots, \kappa_d$  (In this context,  $d$  means the dimension of the tangent space) are the principal curvatures of  $\mathbb{M}$  [20], which implies the eigenvalues of  $D\mathbf{p}(\mathbf{z})$ , when restricted to the tangent space, are

$$\frac{1}{1 - r\kappa_1}, \dots, \frac{1}{1 - r\kappa_d}.$$

All the curvatures are bounded by the reciprocal of the reach:  $|\kappa_i| \leq 1/\tau$  [1]. Therefore, it holds that

$$\frac{\tau}{\tau + r} \mathbf{I} \Big|_{\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})} \preceq D\mathbf{p}(\mathbf{z}) \Big|_{\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})} \preceq \frac{\tau}{\tau - r} \mathbf{I} \Big|_{\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})}.$$

Moreover, as  $D\mathbf{p}(\mathbf{z})$  is restricted to and acts only on the tangent space  $\mathbb{T}_{\mathbf{p}(\mathbf{z})}(\mathbb{M})$ , we have  $\mathbf{0} \preceq D\mathbf{p}(\mathbf{z}) \preceq \frac{\tau}{\tau - r} \mathbf{I}$ , which implies

$$-\mathbf{I} \preceq D\mathbf{p}(\mathbf{z}) - \mathbf{I} \preceq \frac{r}{\tau - r} \mathbf{I}.$$

For  $r \leq \tau/2$ , we have  $\frac{r}{\tau - r} \leq 1$  and hence  $\|D\mathcal{S}_\sigma(\mathbf{z})\| \leq \frac{\sigma}{1 + \sigma}$ . As long as  $\mathbf{z}, \mathbf{z}' \in \mathbb{U}_r(\mathbb{M})$  where  $r \leq \tau/4$  and  $\|\mathbf{z} - \mathbf{z}'\| \leq \tau/4$ , the two points  $\mathbf{z}, \mathbf{z}'$  can be included in a convex subset (actually a ball) of  $\mathbb{U}_r(\mathbb{M})$  with  $r = \tau/2$ . By the mean value theorem, we finish the proof of (34).  $\square$

## C.2. Discussions regarding PnP.

Derivation of HQS.. Consider (4):

$$\min_{\mathbf{y}, \mathbf{z} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{y}\|^2 + \frac{\alpha}{2} \text{dist}^2(\mathbf{z}, \mathbb{M}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{z}\|^2.$$

A typically method to solve it is applying block coordinate descent on it, which is also named ‘‘Half-quadratic-splitting (HQS)’’ in the literature [99]:

$$\begin{aligned}\mathbf{y}_{t+1} &= \arg \min_{\mathbf{y} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{A}\mathbf{y} - \mathbf{x}\|^2 + \frac{\beta}{2} \|\mathbf{y} - \mathbf{z}_t\|^2 = (\mathbf{A}^\top \mathbf{A} + \beta \mathbf{I})^{-1} (\mathbf{A}^\top \mathbf{x} + \beta \mathbf{z}_t) \\ \mathbf{z}_{t+1} &= \arg \min_{\mathbf{z} \in \mathbb{R}^n} \frac{\alpha}{2} \text{dist}^2(\mathbf{z}, \mathbb{M}) + \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}_{t+1}\|^2 = \text{prox}_\sigma(\mathbf{y}_{t+1}) \quad (\text{let } \sigma = \alpha/\beta)\end{aligned}$$

Similarly, we can parameterize  $\text{prox}_\sigma$  as a neural network  $\mathcal{H}_{\theta,\sigma}$ . Therefore, HQS suggests an implicit model

$$\mathcal{G}_\Theta(\mathbf{z}, \mathbf{x}) = \mathcal{H}_{\theta,\sigma} \left( (\mathbf{A}^\top \mathbf{A} + \beta \mathbf{I})^{-1} (\mathbf{A}^\top \mathbf{x} + \beta \mathbf{z}) \right)$$

where  $\Theta = \{\theta, \sigma, \beta\}$  includes all trainable parameters, which derives (6).

**Bibliographical notes.** Here we adopt the long-standing “plug-in denoiser” idea. It originated with Plug-and-Play (PnP) ADMM, which replaces a proximal operator with an off-the-shelf denoiser inside ADMM [92]. The framework has since been developed and analyzed extensively—see, e.g., [8, 10, 47, 87] and the recent survey [46]. In the PGD setting, one pretrains  $\mathcal{H}$  for Gaussian denoising and plugs it into (5) [31, 43, 61, 83]. The same plug-in idea applies to HQS via (6) [42, 79, 102]. In contrast to training a denoiser off-the-shelf and plugging it in, one can train the *entire*  $\mathcal{G}_\Theta$  via deep equilibrium methods for the target task (the approach closest to this paper) in both PGD-style [19, 35, 86, 95, 101, 104] and HQS-style [36].

#### APPENDIX D. PROOFS REGARDING NS EQUATIONS

To rigorously state and prove the theorems, we present some definitions here. First, We denote by  $H^m(\Omega)$  the Sobolev space of functions which are in  $L^2(\Omega)$  together with all their derivatives of order  $\leq m$ . Then  $H_p^m(\Omega) \subset H^m(\Omega)$  is the collection of functions in  $H^m(\Omega)$  that satisfies the periodic boundary condition on  $\Omega$  with zero mean (ref. to [91, Remark 1.1]). Then, we can define the spaces considered in this paper:

$$\mathbb{H} := \left\{ u \in \{H_p^0(\Omega)\}^2 : \nabla \cdot u = 0 \right\}, \quad \mathbb{V} := \left\{ u \in \{H_p^1(\Omega)\}^2 : \nabla \cdot u = 0 \right\}$$

For the NS equation (7), we consider  $f \in \mathbb{H}$  and  $u \in \mathbb{V}$ . Moreover, we denote  $\mathbb{V}'$  as the dual space of  $\mathbb{V}$  and have

$$\mathbb{V} \subset \mathbb{H} \subset \mathbb{V}'.$$

We then equip  $\mathbb{H}$  with the standard  $L^2$  inner product and norm for vector fields:

$$\langle u, v \rangle_{\mathbb{H}} := \int_{\Omega} \langle u(\xi), v(\xi) \rangle d\xi, \quad \|u\|_{\mathbb{H}} := \sqrt{\langle u, u \rangle_{\mathbb{H}}} = \left( \int_{\Omega} \|u(\xi)\|^2 d\xi \right)^{1/2} = \|u\|_{L^2(\Omega)}$$

The space  $\mathbb{V}$  is equipped with the  $L^2$  norm on the first-order derivatives of  $u$ . In particular,

$$\begin{aligned} \langle u, v \rangle_{\mathbb{V}} &:= \sum_{i=1}^2 \int_{\Omega} \left\langle \frac{\partial u}{\partial \xi_i}(\xi), \frac{\partial v}{\partial \xi_i}(\xi) \right\rangle d\xi \\ \|u\|_{\mathbb{V}} &:= \sqrt{\langle u, u \rangle_{\mathbb{V}}} = \left( \sum_{i=1}^2 \int_{\Omega} \left\| \frac{\partial u}{\partial \xi_i}(x) \right\|^2 d\xi \right)^{1/2} = \|\nabla u\|_{L^2(\Omega)} \end{aligned}$$

and  $\|\cdot\|_{\mathbb{V}'}$  is defined as the dual norm of  $\|\cdot\|_{\mathbb{V}}$ . By Poincare and Cauchy-Schwartz inequalities, we have

$$\|v\|_{\mathbb{H}} \leq c_1 \|v\|_{\mathbb{V}}, \quad \forall v \in \mathbb{V}$$

and

$$\|v\|_{\mathbb{V}'} \leq c_2 \|v\|_{\mathbb{H}}, \quad \forall v \in \mathbb{H}$$

where  $c_1, c_2$  are constants depending on the domain  $\Omega$ . The above definitions and results are standard in the literature and we largely follow the notation in [91, Section 2].

*Proof of Theorem 3.6.* [91, Theorem 10.1] states that, for any  $f \in \mathbb{V}'$ , if  $\|f\|_{\mathbb{V}'} \leq c_0 \nu^2$  (with  $c_0 > 0$  depending only on  $\Omega$ ), then the steady NS problem (7) has a unique solution  $u_*$ . Since  $\mathbb{H} \subset \mathbb{V}'$  and  $\|f\|_{\mathbb{V}'} \leq c_2 \|f\|_{\mathbb{H}}$ , this yields uniqueness on

$$\mathbb{H}_{\nu}^{(1)} := \left\{ f \in \mathbb{H} : \|f\|_{\mathbb{H}} \leq \frac{c_0}{c_2} \nu^2 \right\}.$$

Moreover, by [91, Theorem 10.4], there exists an open dense set  $\mathbb{H}_\nu^{(2)} \subset \mathbb{H}$  such that, on each connected component of  $\mathbb{H}_\nu^{(2)}$ , the solution  $u_*$  depends  $C^\infty$  on  $f$ ; in particular,  $f \mapsto u_*$  is locally Lipschitz there. Define  $\mathbb{H}_\nu := \mathbb{H}_\nu^{(1)} \cap \mathbb{H}_\nu^{(2)}$ . Since  $\mathbb{H}_\nu^{(2)}$  is open and dense in  $\mathbb{H}$ , the set  $\mathbb{H}_\nu$  is dense in  $\mathbb{H}_\nu^{(1)}$ . On  $\mathbb{H}_\nu$  the solution is unique and the map  $f \mapsto u_*$  is locally Lipschitz. This completes the proof.  $\square$

Before moving to Corollary 3.7, let's reclarify lifting and projection operators: Let the lifting (or extension) operator  $\mathcal{E}_h : \mathbb{R}^{N_h \times 2} \rightarrow \{L^2(\Omega)\}^2$  be the piecewise-constant reconstruction  $\mathcal{E}_h(\mathbf{x}) := \sum_{C \in \Omega_h} x_C \mathbf{1}_C$ , and let  $\mathcal{P} : \{L^2(\Omega)\}^2 \rightarrow \mathbb{H}$  be the orthogonal projection onto divergence-free, zero-mean fields. Then we move on to Corollary 3.7.

*Proof of Corollary 3.7.* The mapping  $\mathcal{F}_2 : \mathbf{x} \mapsto \mathbf{y}_*$  can be viewed as a composition of multiple mappings: We first map  $\mathbf{x} \in \mathbb{R}^{N_h \times 2}$  to a continuous version  $f \in \mathbb{H}$  by  $\mathcal{P} \circ \mathcal{E}_h$ , then  $f$  can be mapped to its corresponding solution  $u_*$  by a Locally Lipschitz operator as stated in Theorem 3.6. Here we denote this mapping by  $\mathcal{S} : f \mapsto u_*$ . Then  $u_*$  is mapped to  $\omega_*$  by vorticity:  $\nabla \times u_*$ , and finally  $\omega_*$  can be mapped to  $\mathbf{y}_*$  by a restriction operator  $\mathcal{R}_h$ :

$$\mathcal{F}_2 = \mathcal{R}_h \circ (\nabla \times) \circ \mathcal{S} \circ \mathcal{P} \circ \mathcal{E}_h.$$

Then let's analyze the norm of the above operators one by one. Firstly, the restriction operator  $\mathcal{R}_h$  has a norm no greater than 1 as:

$$\begin{aligned} \|\mathcal{R}_h(\omega)\|_{\ell_h^2}^2 &= \sum_{C \in \Omega_h} |C| \left| \frac{1}{|C|} \int_C \omega(\xi) d\xi \right|^2 \\ &\leq \sum_{C \in \Omega_h} \frac{1}{|C|} \left( \int_C |\omega(\xi)| d\xi \right)^2 \leq \sum_{C \in \Omega_h} \int_C |\omega(\xi)|^2 d\xi = \|\omega\|_{L^2(\Omega)}^2 \end{aligned}$$

Note that  $\mathcal{R}_h$  is a linear operator, hence its bounded norm immediately leads to its bounded Lipschitz constant:

$$\|\mathcal{R}_h(\omega) - \mathcal{R}_h(\omega')\|_{\ell_h^2}^2 = \|\mathcal{R}_h(\omega - \omega')\|_{\ell_h^2}^2 \leq \|\omega - \omega'\|_{L^2(\Omega)}^2.$$

Second, the curl operator  $\nabla \times$  must be a bounded linear operator because the solution  $u_* \in \mathbb{V}$ , where first-order derivatives must be  $L^2$ . Third, the solution mapping  $\mathcal{S}$  has been discussed in Theorem 3.6, it is a nonlinear operator, but it is locally Lipschitz continuous. Fourth, the projection operator  $\mathcal{P}$  must be linear and have a norm no greater than 1. Finally, the lifting operator is linear and has a bounded norm as:

$$\|\mathcal{E}_h(\mathbf{x})\|_{L^2(\Omega)}^2 = \sum_{C \in \Omega_h} |C| |x_C|^2 = \|\mathbf{x}\|_{\ell_h^2}^2$$

Therefore, except for the nonlinear operator  $\mathcal{S}$ , the other four operators are all linear and bounded and hence are globally Lipschitz continuous. As long as we can show that the input of  $\mathcal{S}$  must be taken from the unique solution regime  $\mathbb{H}_\nu$ , we will complete the proof that  $\mathcal{F}_2$  is locally Lipschitz everywhere on  $\mathbb{X}_{\nu,h}$ . This can be proved because  $\mathbf{x} \in \mathbb{X}_{\nu,h}$  implies  $\mathcal{P}(\mathcal{E}_h(\mathbf{x})) \in \mathbb{H}_\nu$ . Finally, by applying Theorem 2.4, we conclude the existence of  $\mathcal{G}$  described in Corollary 3.7, which finishes the entire proof.  $\square$

#### APPENDIX E. PROOFS REGARDING LINEAR PROGRAMMING

Although Lipschitz continuity of LP solution maps has been studied (e.g., [23, 65]), we are not aware of a reference that states Theorem 3.8 in the precise form needed here—particularly allowing perturbations of  $\mathbf{A}$  (rather than treating  $\mathbf{A}$  as fixed). For completeness, we therefore include a self-contained discussion and proof.

To work with a standard form, we rewrite the general-form problem (9) in standard form. Suppose there are  $p$  equality constraints and  $q$  inequality constraints. Without loss of generality, we assume  $\circ_i$  equals to “=” for  $1 \leq i \leq p$  and  $\circ_i$  equals to “ $\leq$ ” for  $p+1 \leq i \leq m$ . Then we denote  $\mathbf{A}_p$  as the first  $p$  rows of matrix  $\mathbf{A}$  and  $\mathbf{A}_q$  as the remaining part:

$$\mathbf{A}_p := \mathbf{A}[1 : p, :], \quad \mathbf{A}_q := \mathbf{A}[p+1 : m, :]$$

And therefore the general form LP (9) can be written as

$$\min_{\mathbf{y} \in \mathbb{R}^n} \mathbf{c}^\top \mathbf{y}, \quad \text{s.t. } \mathbf{A}_p \mathbf{y} = \mathbf{b}_p, \mathbf{A}_q \mathbf{y} \leq \mathbf{b}_q, \mathbf{l} \leq \mathbf{y} \leq \mathbf{u}.$$

Let  $\hat{\mathbf{y}} := \mathbf{y} - \mathbf{l}$ ,  $\mathbf{s} := \mathbf{b}_q - \mathbf{A}_q \mathbf{y}$ , and  $\mathbf{t} := \mathbf{u} - \mathbf{y}$ , the above problem can be transformed to

$$\min_{\mathbf{y} \in \mathbb{R}^n} \mathbf{c}^\top \hat{\mathbf{y}}, \quad \text{s.t. } \begin{bmatrix} \mathbf{A}_p & & \\ \mathbf{A}_q & \mathbf{I} & \\ \mathbf{I} & & \mathbf{I} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{y}} \\ \mathbf{s} \\ \mathbf{t} \end{bmatrix} = \begin{bmatrix} \mathbf{b}_p - \mathbf{A}_p \mathbf{l} \\ \mathbf{b}_q - \mathbf{A}_q \mathbf{l} \\ \mathbf{u} - \mathbf{l} \end{bmatrix}, \quad \hat{\mathbf{y}} \geq \mathbf{0}, \mathbf{s} \geq \mathbf{0}, \mathbf{t} \geq \mathbf{0}$$

By letting

$$\tilde{\mathbf{c}} := \begin{bmatrix} \mathbf{c} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \quad \tilde{\mathbf{A}} := \begin{bmatrix} \mathbf{b}_p - \mathbf{A}_p \mathbf{l} \\ \mathbf{b}_q - \mathbf{A}_q \mathbf{l} \\ \mathbf{u} - \mathbf{l} \end{bmatrix}, \quad \tilde{\mathbf{b}} := \begin{bmatrix} \mathbf{b}_p - \mathbf{A}_p \mathbf{l} \\ \mathbf{b}_q - \mathbf{A}_q \mathbf{l} \\ \mathbf{u} - \mathbf{l} \end{bmatrix}, \quad \tilde{\mathbf{y}} := \begin{bmatrix} \hat{\mathbf{y}} \\ \mathbf{s} \\ \mathbf{t} \end{bmatrix}$$

The problem is equivalently expressed in standard form as

$$\min_{\tilde{\mathbf{y}}} \tilde{\mathbf{c}}^\top \tilde{\mathbf{y}}, \quad \text{s.t. } \tilde{\mathbf{A}} \tilde{\mathbf{y}} = \tilde{\mathbf{b}}, \tilde{\mathbf{y}} \geq \mathbf{0}.$$

In fact, every LP can be rewritten in an equivalent standard form. While concepts such as basic feasible solutions, degeneracy, and complementary slackness are most naturally and cleanly stated in standard form, each admits a closely related analogue (with minor adjustments) for the general form. Accordingly—without loss of generality and to keep the focus on core ideas—we carry out the proof in the standard-form setting:

$$\min_{\mathbf{y}} \mathbf{c}^\top \mathbf{y}, \quad \text{s.t. } \mathbf{A} \mathbf{y} = \mathbf{b}, \mathbf{y} \geq \mathbf{0},$$

with dual

$$\min_{\mathbf{z}} \mathbf{b}^\top \mathbf{z}, \quad \text{s.t. } \mathbf{A}^\top \mathbf{z} \leq \mathbf{c}.$$

Here, we follow the standard settings in the literature:  $\mathbf{y}, \mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{z}, \mathbf{b} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(\mathbf{A}) = m$  (ensured by preprocessing with removing redundant equalities), and  $m \leq n$ . In this context, we define the domain of LP that we work on:

$$\mathbb{X} := \{(\mathbf{A}, \mathbf{b}, \mathbf{c}) : \text{The resulting standard LP is feasible and bounded}\}$$

Note that, to match the rest of the paper, we reserve  $\mathbf{x}$  for machine learning model inputs (in this context, it is  $\mathbf{x} = (\mathbf{A}, \mathbf{b}, \mathbf{c})$ ) and hence write the primal LP variable as  $\mathbf{y}$  and the dual LP variable as  $\mathbf{z}$ . This departs from the common  $(\mathbf{x}, \mathbf{y})$  convention. Note also that in the main text the symbol  $\mathbf{z}$  denotes a latent variable; here, in the appendix regarding LP's technical details, it denotes the dual variable. These meanings are unrelated and should be clear from context.

Now let's present some definitions used in this appendix. Fix a *basis* by selecting an index set  $B \subset \{1, 2, \dots, n\}$  with  $|B| = m$  such that the  $m \times m$  submatrix  $\mathbf{B} := \mathbf{A}[:, B]$  is *nonsingular*. Let  $N = \{1, 2, \dots, n\} \setminus B$  be the complement of the basis and let  $\mathbf{N} := \mathbf{A}[:, N]$ . Then the equality constraints read

$$\mathbf{B} \mathbf{y}_B + \mathbf{N} \mathbf{y}_N = \mathbf{b}$$

Setting  $\mathbf{y}_N = \mathbf{0}$  yields  $\mathbf{y}_B = \mathbf{B}^{-1} \mathbf{b}$ . Such a  $\mathbf{y} = [\mathbf{y}_B, \mathbf{0}]$  is called a *basic solution*. If additionally  $\mathbf{y}_B \geq \mathbf{0}$ , this basic solution is feasible, then it is called a *basic feasible solution* (*BFS*). On the dual side, we define the slack variable  $\mathbf{s}$  and its sub-vector restricted to  $B$  and  $N$ :

$$\mathbf{s} := \mathbf{c} - \mathbf{A}^\top \mathbf{z}, \quad \mathbf{s}_B := \mathbf{c}_B - \mathbf{B}^\top \mathbf{z}, \quad \mathbf{s}_N := \mathbf{c}_N - \mathbf{N}^\top \mathbf{z}.$$

A pair  $(\mathbf{y}, \mathbf{z})$  is primal-dual optimal (i.e., satisfies KKT for LP) iff

$$(35) \quad \mathbf{A} \mathbf{y} = \mathbf{b}, \quad \mathbf{c} = \mathbf{A}^\top \mathbf{z} + \mathbf{s}, \quad \mathbf{y} \odot \mathbf{s} = \mathbf{0}, \quad \mathbf{y} \geq \mathbf{0}, \mathbf{s} \geq \mathbf{0}$$

for some  $\mathbf{s} \in \mathbb{R}^n$ . If, in addition, there exists a basis  $B$  such that

$$(36) \quad \mathbf{y}_B \geq \mathbf{0}, \mathbf{y}_N = \mathbf{0}, \mathbf{s}_B = \mathbf{0}, \mathbf{s}_N \geq \mathbf{0},$$

then the tuple  $(\mathbf{y}, \mathbf{z}, \mathbf{s})$  is called an optimal BFS with a complementary dual. By the fundamental theorem of linear programming, any feasible instance with finite optimal value  $(\mathbf{A}, \mathbf{b}, \mathbf{c}) \in \mathbb{X}$  admits an optimal BFS with a complementary dual satisfying (35) and (36) together [6].

While conditions (35) and (36) are enough to ensure the existence of the optimal basic solutions, they are not enough to ensure that the optimal solution is unique and local Lipschitz continuous w.r.t. the inputs  $(\mathbf{A}, \mathbf{b}, \mathbf{c})$ . To ensure these points, we present two additional conditions based on (35) and (36):

$$(37) \quad \mathbf{y}_B > \mathbf{0} \quad (\text{Non-degeneracy})$$

$$(38) \quad \mathbf{s}_N > \mathbf{0} \quad (\text{Strict complementary slackness})$$

All the conditions together are enough to the uniqueness and local Lipschitz continuity. Let's introduce a set consisting of all "good" LP instances:

$$\mathbb{X}_{\text{sub}} := \{(\mathbf{A}, \mathbf{b}, \mathbf{c}) \in \mathbb{X} : \text{The LP yields a tuple } (\mathbf{y}, \mathbf{z}, \mathbf{s}) \text{ satisfying (35), (36), (37) and (38)}.\}$$

With all the preparations, we can prove Theorem 3.8 now. Actually, proving Theorem 3.8 in the context of standard-form LP is equivalent to proving the following two theorems.

**Theorem E.1.** *For any LP  $(\mathbf{A}, \mathbf{b}, \mathbf{c}) \in \mathbb{X}_{\text{sub}}$ , it must yield a unique optimal solution  $\mathbf{y}_*$ , and the solution mapping  $(\mathbf{A}, \mathbf{b}, \mathbf{c}) \mapsto \mathbf{y}_*$  is locally Lipschitz continuous everywhere on  $\mathbb{X}_{\text{sub}}$ .*

**Theorem E.2.**  *$\mathbb{X}_{\text{sub}}$  is a dense subset of  $\mathbb{X}$ .*

Theorem E.1 follows from [22], which develops Robinson's notion of strong regularity [81] for nonlinear programs. For completeness—and to keep notation consistent with linear programming—we restate the relevant lemma in an LP-adapted form and then verify its hypotheses for LP. We begin by quoting the result from [22].

**Lemma E.3** ([22]). *Consider a parameteric nonlinear program:*

$$\begin{aligned} \min_{\mathbf{y} \in \mathbb{R}^n} \quad & \mathbf{c}^\top \mathbf{y} + g_0(\mathbf{w}, \mathbf{y}) \\ \text{s.t.} \quad & g_i(\mathbf{w}, \mathbf{y}) = u_i, \quad 1 \leq i \leq r \\ & g_i(\mathbf{w}, \mathbf{y}) \leq u_i, \quad r+1 \leq i \leq d \end{aligned}$$

where  $g_i (0 \leq i \leq d)$  are all  $\mathcal{C}^2$  functions, and  $\mathbf{c}, \mathbf{w}$  and  $\mathbf{u} = [u_1, \dots, u_d]^\top$  are parameters to describe the program, and consider its Lagrangian with multipliers  $\lambda = [\lambda_1, \dots, \lambda_d] \in \mathbb{R}^d$  given by

$$L(\mathbf{w}, \mathbf{y}, \lambda) = g_0(\mathbf{w}, \mathbf{y}) + \sum_{i=1}^d \lambda_i g_i(\mathbf{w}, \mathbf{y}).$$

Let  $(\bar{\mathbf{y}}, \bar{\lambda})$  be a KKT point at  $(\bar{\mathbf{c}}, \bar{\mathbf{w}}, \bar{\mathbf{u}})$ , and define the index sets at  $(\bar{\mathbf{y}}, \bar{\lambda})$

$$\begin{aligned} I_1 &= \left\{ r+1 \leq i \leq d : g_i(\bar{\mathbf{w}}, \bar{\mathbf{y}}) = u_i, \bar{\lambda}_i > 0 \right\} \cup \left\{ 1, \dots, r \right\}, \\ I_2 &= \left\{ r+1 \leq i \leq d : g_i(\bar{\mathbf{w}}, \bar{\mathbf{y}}) = u_i, \bar{\lambda}_i = 0 \right\}, \\ I_3 &= \left\{ r+1 \leq i \leq d : g_i(\bar{\mathbf{w}}, \bar{\mathbf{y}}) < u_i, \bar{\lambda}_i = 0 \right\}. \end{aligned}$$

If the following conditions hold:

- The constraint gradients  $\nabla_{\mathbf{y}} g_i(\bar{\mathbf{w}}, \bar{\mathbf{y}})$  for  $i \in I_1 \cup I_2$  are linearly independent; and
- It holds that

$$\langle \mathbf{y}', \nabla_{\mathbf{y}}^2 L(\bar{\mathbf{w}}, \bar{\mathbf{y}}, \bar{\lambda}) \mathbf{y}' \rangle > 0$$

for all  $\mathbf{y}' \neq \mathbf{0}$  in the subspace  $\mathbb{M} = \left\{ \mathbf{y}' : \mathbf{y}' \perp \nabla_{\mathbf{y}} g_i(\bar{\mathbf{w}}, \bar{\mathbf{y}}) \text{ for all } i \in I_1 \right\}$ ,

then the KKT solution map  $(\mathbf{c}, \mathbf{w}, \mathbf{u}) \mapsto (\mathbf{y}, \lambda)$  is locally single-valued and Lipschitz around  $(\bar{\mathbf{c}}, \bar{\mathbf{w}}, \bar{\mathbf{u}}, \bar{\mathbf{y}}, \bar{\lambda})$ .

*Proof of Theorem E.1.* Taking  $r = m$  and  $d = m + n$ . Let  $\mathbf{a}_i^\top$  be the  $i$ -th row of  $\mathbf{A}$  in standard LP, and let

$$g_i(\mathbf{w}, \mathbf{y}) = \begin{cases} \mathbf{a}_i^\top \mathbf{y}, & i = 1, \dots, m, \\ -y_{i-m}, & i = m+1, \dots, m+n, \end{cases} \quad u_i = \begin{cases} b_i, & i = 1, \dots, m, \\ 0, & i = m+1, \dots, m+n, \end{cases}$$

with  $\mathbf{w}$  collecting the coefficients of  $\mathbf{A}$ . The Lagrangian in Lemma E.3 becomes

$$L(\mathbf{w}, \mathbf{y}, \lambda) = \mathbf{c}^\top \mathbf{y} + \sum_{i=1}^m \lambda_i \mathbf{a}_i^\top \mathbf{y} + \sum_{j=1}^n \lambda_{m+j} (-y_j).$$

Introduce the usual dual/primal-slack variables

$$\mathbf{z} := -\lambda_{1:m} \in \mathbb{R}^m, \quad \mathbf{s} := \lambda_{m+1:m+n} \in \mathbb{R}_{\geq 0}^n,$$

to rewrite stationarity as  $\nabla_{\mathbf{y}} L = \mathbf{c} - \mathbf{A}^\top \mathbf{z} - \mathbf{s} = \mathbf{0}$ , i.e.,  $\mathbf{s} = \mathbf{c} - \mathbf{A}^\top \mathbf{z}$ . Primal feasibility is  $\mathbf{A}\mathbf{y} = \mathbf{b}$ ,  $\mathbf{y} \geq 0$ ; dual feasibility is  $\mathbf{s} \geq \mathbf{0}$ ; and complementarity is  $\mathbf{y} \odot \mathbf{s} = \mathbf{0}$ . Thus the KKT system in Lemma E.3 coincides with the standard LP KKT conditions.

Assume  $(\mathbf{A}, \mathbf{b}, \mathbf{c}) \in \mathbb{X}_{\text{sub}}$ , i.e., the LP admits a tuple  $(\bar{\mathbf{y}}, \bar{\mathbf{z}}, \bar{\mathbf{s}})$  satisfying (35), (36), (37) and (38) (A nondegenerate and strict complementary basic point). In this context, the index sets  $I_1, I_2, I_3$  at  $(\bar{\mathbf{y}}, \bar{\mathbf{z}}, \bar{\mathbf{s}})$  become:

$$\begin{aligned} I_1 &= \{1, \dots, m\} \cup \{m+j : \bar{y}_j = 0, \bar{s}_j > 0\}, \\ I_2 &= \{m+j : \bar{y}_j = 0, \bar{s}_j = 0\}, \\ I_3 &= \{m+j : \bar{y}_j > 0, \bar{s}_j = 0\}. \end{aligned}$$

which implies:

- For each  $j$ , either  $\bar{y}_j > 0$  or  $\bar{s}_j > 0$ , which implies  $I_2 = \emptyset$ .
- $I_3$  is substantially the basis set:  $I_3 = \{m+j : j \in B\}$
- $I_1$  includes all the indices in the complement of basis:  $I_1 = \{1, \dots, m\} \cup \{m+j : j \in N\}$

To verify the hypotheses of Lemma E.3, we examine the gradients:

$$\{\nabla_{\mathbf{y}} g_i\}_{i \in I_1} = \{\mathbf{a}_i\}_{i=1}^m \cup \{-\mathbf{e}_j\}_{j \in N}$$

In the context of standard LP,  $|N| = n - m$ . Hence,  $\{\nabla_{\mathbf{y}} g_i\}_{i \in I_1}$  consists of  $n$  vectors in  $\mathbb{R}^n$ . Now we create a matrix  $\mathbf{G}$  by stacking these vectors as rows:

$$\mathbf{G} := \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_m^\top \\ \mathbf{e}_{j_1}^\top \\ \vdots \\ \mathbf{e}_{j_{n-m}}^\top \end{bmatrix}$$

By properly permuting the columns of  $\mathbf{G}$ , it becomes

$$\tilde{\mathbf{G}} = \begin{bmatrix} \mathbf{B} & \mathbf{N} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$$

where  $\mathbf{I}$  represents the identity matrix in  $\mathbb{R}^{n-m}$ . Since  $\mathbf{B}$  (the basis matrix) and  $\mathbf{I}$  are both nonsingular,  $\tilde{\mathbf{G}}$  (and hence  $\mathbf{G}$ ) must be nonsingular. Therefore, the rows of  $\mathbf{G}$  are linearly independent, i.e.,  $\{\nabla_{\mathbf{y}} g_i\}_{i \in I_1}$  is linearly independent. With  $I_2 = \emptyset$ , the first hypothesis of Lemma E.3 holds. Moreover, because these gradients  $\{\nabla_{\mathbf{y}} g_i\}_{i \in I_1}$  span  $\mathbb{R}^n$ , the  $\mathbb{M}$  subspace must be trivial:  $\mathbb{M} = \{\mathbf{0}\}$ . Therefore, the second hypothesis of Lemma E.3 is automatically satisfied.

By Lemma E.3, the KKT solution map is locally single-valued and Lipschitz around the given point, which yields the desired local uniqueness and Lipschitz dependence of  $\mathbf{y}_*$  on  $(\mathbf{A}, \mathbf{b}, \mathbf{c})$  for every  $(\mathbf{A}, \mathbf{b}, \mathbf{c}) \in \mathbb{X}_{\text{sub}}$ .  $\square$

Theorem E.2 can be proved by fundamental concepts in real analysis.

*Proof of Theorem E.2.* To prove  $\mathbb{X}_{\text{sub}}$  is dense in  $\mathbb{X}$ , it's enough to show that: For any  $(\mathbf{A}, \mathbf{b}, \mathbf{c}) \in \mathbb{X}$ , one can always create a sequence of LP  $\{(\mathbf{A}_k, \mathbf{b}_k, \mathbf{c}_k)\}_{k \geq 1} \subset \mathbb{X}_{\text{sub}}$  such that

$$\mathbf{A}_k \rightarrow \mathbf{A}, \mathbf{b}_k \rightarrow \mathbf{b}, \mathbf{c}_k \rightarrow \mathbf{c}.$$

Now let's fix  $(\mathbf{A}, \mathbf{b}, \mathbf{c}) \in \mathbb{X}$ . As we previously discussed, there must be a tuple  $(\mathbf{y}, \mathbf{z}, \mathbf{s})$  satisfying (35) and (36). Define:

$$\mathbf{y}_k := \mathbf{y} + \frac{1}{k} \mathbf{e}_B, \quad \mathbf{s}_k := \mathbf{s} + \frac{1}{k} \mathbf{e}_N, \quad \mathbf{z}_k := \mathbf{z}$$

so that  $(\mathbf{y}_k, \mathbf{z}_k, \mathbf{s}_k)$  must satisfy the nondegeneracy and strict complementary slackness: (36), (37), and (38). Accordingly, define

$$\mathbf{A}_k := \mathbf{A}, \quad \mathbf{b}_k := \mathbf{A}_k \mathbf{y}_k, \quad \mathbf{c}_k := \mathbf{A}_k^\top \mathbf{z}_k + \mathbf{s}_k$$

Then one can verify that the tuple  $(\mathbf{y}, \mathbf{z}, \mathbf{s})$  satisfies (35), (36), (37) and (38) for the LP instance  $(\mathbf{A}_k, \mathbf{b}_k, \mathbf{c}_k)$ , hence  $(\mathbf{A}_k, \mathbf{b}_k, \mathbf{c}_k) \in \mathbb{X}_{\text{sub}}$  for all  $k \geq 1$ . Finally, such a perturbed LP instance can be arbitrarily close to  $(\mathbf{A}, \mathbf{b}, \mathbf{c})$  as  $k \rightarrow \infty$ :

$$\begin{aligned} \|\mathbf{A}_k - \mathbf{A}\| &= 0 \\ \|\mathbf{b}_k - \mathbf{b}\| &= \left\| \mathbf{A} \left( \frac{1}{k} \mathbf{e}_B \right) \right\| \leq \frac{1}{k} \|\mathbf{A}\| \|\mathbf{e}_B\| = \frac{\sqrt{m}}{k} \|\mathbf{A}\| \rightarrow 0 \\ \|\mathbf{c}_k - \mathbf{c}\| &= \left\| \frac{1}{k} \mathbf{e}_N \right\| = \frac{\sqrt{n-m}}{k} \rightarrow 0 \end{aligned}$$

which finishes the proof.  $\square$

## APPENDIX F. EXPERIMENT DETAILS REGARDING IMAGE RECONSTRUCTION

This section complements the main text with additional implementation and dataset details for the inverse-problem experiments.

**Experiment settings.** We consider an image deblurring task,  $\mathbf{x} = \mathbf{A}(\mathbf{y}_*) + \mathbf{n}$ , where  $\mathbf{A}$  is the blur operator and  $\mathbf{n}$  is the Gaussian noise ( $\sigma = 0.03$ ). We use a motion-blur operator, and the blur kernel is the first of the eight kernels from [53]. Ground-truth images  $\mathbf{y}_*$  come from BSDS500 [66]. We follow the official splits (200 train / 100 validation / 200 test) and apply a random  $128 \times 128$  crop to each image. For each  $\mathbf{y}_*$ , we generate the corresponding  $\mathbf{x}$  by applying  $\mathbf{A}$  and adding noise. The resulting pairs  $(\mathbf{x}, \mathbf{y}_*)$  form three datasets  $\mathbb{D}_{\text{inv,train}}$ ,  $\mathbb{D}_{\text{inv,val}}$ , and  $\mathbb{D}_{\text{inv,test}}$  for training, validation, and testing, respectively. In both PGD and HQS style parameterizations ((5) and (6)), the operator  $\mathcal{H}$  is implemented with DRUNet [102].

**Training.** We initialize  $\mathcal{H}$  using pretrained weights from the Deepinv library [88] and then fine-tune the full implicit models on the BSDS500 training set for this deblurring task. Training follows the vanilla Jacobian-based implicit differentiation and is implemented on top of the official Deepinv framework. All models were trained with Adam (learning rate  $10^{-4}$ , batch size 3). Explicit baselines were trained for 20 epochs, and the implicit models for 10 epochs. After each epoch we evaluated on the validation set and saved the checkpoint; the final model used for testing is the one with the lowest validation loss. These epoch budgets were sufficient for validation-loss convergence.

**PSNR.** PSNR (Peak Signal-to-Noise Ratio) is defined between a reference  $\mathbf{y}^*$  and reconstructed image  $\mathbf{y}$  as

$$\text{PSNR}(\mathbf{y}, \mathbf{y}^*) := 10 \log_{10} \left( \frac{n \cdot \text{MAX}^2}{\|\mathbf{y} - \mathbf{y}^*\|^2} \right)$$

where  $n$  is the dimension of  $\mathbf{y}$  and  $\mathbf{y}^*$ , and MAX means the max possible pixel value (e.g., 255 for 8-bit, or 1 if images are in  $[0, 1]$ ). In our context, it is 1. Higher PSNR means better (more accurate) reconstruction.

**Standard test set.** Evaluation uses the 200 images from the official BSDS500 test split, randomly cropped to  $128 \times 128$ . Let  $\mathbb{D}_{\text{inv,test}} = \{(\mathbf{x}_i, \mathbf{y}_i^*)\}_{i=1}^{200}$ , where

$$\mathbf{x}_i = \mathbf{A}(\mathbf{y}_i^*) + \mathbf{n}_i, \quad \mathbf{n}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad \sigma = 0.03.$$

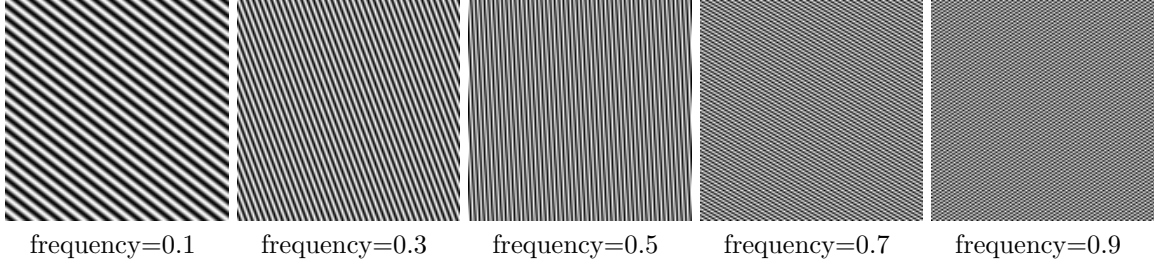


FIGURE 8. Visualized perturbations for inverse problems.

Here  $\mathbf{y}_i^*$  denotes the clean (ground-truth) image and  $\mathbf{x}_i$  its corresponding blurred–noisy observation under the forward model  $\mathbf{A}$ .

**Perturbed test set.** To empirically validate our theory, we created a perturbed version of the test set. To create a diverse and representative set of perturbations, we generate perturbations that correspond to different frequency levels. Image frequencies represent different levels of detail, where low frequencies capture smooth, large-scale areas, and high frequencies capture sharp edges and fine textures. By probing the model with perturbations across this spectrum, we can comprehensively evaluate its behavior.

Specifically, we construct each perturbation by targeting a singular vector of the forward operator  $\mathbf{A}$ . Because  $\mathbf{A}$  is (circular) convolution, its singular vectors are Fourier modes. For each image  $\mathbf{y}_i^*$  and each frequency magnitude  $f \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$ , we first identify the 2D discrete Fourier frequencies and sort them by their geometric distance from the origin. We then select the frequency coordinate  $(u, v)$  at the  $f$ -th percentile of this sorted list. A one-hot tensor is created in the Fourier domain with a value of 1.0 at the chosen  $(u, v)$  position and zeros elsewhere. This sparse frequency representation is transformed back into the image domain by applying the adjoint of the blur operator,  $\mathbf{A}^\top$ . These perturbations are visualized in Figure 8. Adding them to  $\mathbf{y}_i^*$  respectively yields perturbed clean images  $\mathbf{y}_{i,j}^*$  ( $j = 1, \dots, 5$ ); we then form the corresponding observation

$$\mathbf{x}_{i,j} = \mathbf{A}(\mathbf{y}_{i,j}^*) + \mathbf{n}_i, \quad \mathbf{n}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

The perturbed evaluation set is

$$\mathbb{D}'_{\text{inv,test}} = \{(\mathbf{x}_{i,j}, \mathbf{y}_{i,j}^*) : 1 \leq i \leq 200, 1 \leq j \leq 5\}.$$

For convenience we also define the unperturbed index  $j = 0$  by  $\mathbf{x}_{i,0} := \mathbf{x}_i$  and  $\mathbf{y}_{i,0}^* := \mathbf{y}_i^*$ .

**Platform.** All experiments were run on a workstation with eight Quadro RTX 6000 GPUs.

**Additional Experiments.** Implicit models often excel on imaging tasks, but a natural question is whether simply stacking more explicit layers (i.e., deepening the model) can close the gap. To probe this, we construct explicit counterparts to implicit models by untying the parameters across iterations:

$$\min_{\Theta} \mathbb{E}_{\mathbf{x}} \ell(\mathbf{y}_T, \mathbf{y}_*), \quad \text{s.t. } \mathbf{y}_t = \mathcal{G}_{\Theta(t)}(\mathbf{y}_{t-1}, \mathbf{x}), \quad t = 1, \dots, T$$

where each block  $\mathcal{G}_{\Theta(t)}$  has the same architecture as in the implicit case (PGD or HQS), but  $\Theta(t)$  are separate for each  $t$ . This is equivalent to stacking  $T$  blocks to form a deeper explicit model with more learnable parameters. Unlike implicit models (which can use different iteration counts at train vs. test), these explicit models must use the same  $T$  for both. We evaluate  $T \in \{1, 2, 3\}$  to compare against the corresponding implicit models.

Table 2 reports results on image deblurring. Across both PGD and HQS, deepening explicit models roughly triples parameters but yields only small PSNR gains; the implicit models, at the same parameter count as the shallow explicit baseline, outperform even the deepest explicit versions.

TABLE 2. Deeper explicit models vs implicit models for image deblurring. Here  $(\times 1), (\times 1), (\times 1)$  represents  $T = 1, 2, 3$ , respectively. Results are obtained on all 200 test samples.

		Explicit ( $\times 1$ )	Explicit ( $\times 2$ )	Explicit ( $\times 3$ )	Implicit
PGD	# Params.	32.641 M	65.282 M	97.923 M	32.641 M
	Avg. PSNR	27.14 dB	27.64 dB	27.87 dB	28.21 dB
HQS	# Params.	32.641 M	65.282 M	97.923 M	32.641 M
	Avg. PSNR	26.94 dB	28.02 dB	28.26 dB	29.18 dB

#### APPENDIX G. EXPERIMENT DETAILS REGARDING SCIENTIFIC COMPUTING

**Model structure and training.** Given cell-averaged forces  $\mathbf{x} \in \mathbb{R}^{H \times W \times 2}$  and vorticities  $\mathbf{y} \in \mathbb{R}^{H \times W \times 1}$ , where  $H$  means the height and  $W$  means the width, we learn

$$\mathbf{z}_* = \mathcal{G}_\Theta(\mathbf{z}_*, \mathcal{Q}_\Phi(\mathbf{x})), \quad \mathbf{y}_* = \mathcal{Q}_\Psi(\mathbf{z}_*),$$

where  $\mathbf{z}_* \in \mathbb{R}^{H \times W \times C}$  is a latent field with  $C$  channels. At inference, we iterate

$$\mathbf{z}_t = \mathcal{G}_\Theta(\mathbf{z}_{t-1}, \mathcal{Q}_\Phi(\mathbf{x})),$$

for  $1 \leq t \leq T$  and finally call  $\mathbf{y}_T = \mathcal{Q}_\Psi(\mathbf{z}_T)$ .

The projection  $\mathcal{Q}_\Phi$  is a *pointwise* linear encoder applied at each grid cell to lift into  $C$  channels. In particular,  $\mathbf{g} = \mathcal{Q}_\Phi(\mathbf{x})$  reads

$$\mathbf{g} = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 \in \mathbb{R}^{H \times W \times C}$$

where  $\Phi = (\mathbf{W}_1, \mathbf{b}_1)$  are learnable parameters.

The core map  $\mathcal{G}_\Theta(\mathbf{z}, \mathbf{g})$  stacks  $L$  identical FNO layers with input injection:

$$\begin{aligned} \mathbf{z}^{(0)} &= \mathbf{z} \\ \mathbf{z}^{(l)} &= \sigma \left( \mathbf{g} + \mathbf{W}_2^{(l)} \mathbf{z}^{(l-1)} + \mathbf{b}_2^{(l)} + \text{IFFT}(\mathbf{R}^{(l)} \cdot \text{FFT}(\mathbf{z}^{(l-1)})) \right), \quad l = 1, 2, \dots, L, \\ \mathcal{G}_\Theta(\mathbf{z}, \mathbf{g}) &= \mathbf{z}^{(L)} \end{aligned}$$

where  $\Theta = \{\mathbf{W}_2^{(l)}, \mathbf{b}_2^{(l)}, \mathbf{R}^{(l)}\}_{l=1}^L$  are learnable parameters. Each layer: (i) performs a global spectral convolution on  $\mathbf{z}$ : take an FFT of the  $C$ -channel tensor, keep only a small set of low Fourier modes. Suppose the number of retained Fourier modes is  $K \times K$  (2D FFT),  $\text{FFT}(\mathbf{z}) \in \mathbb{C}^{K \times K \times C}$ . For each retained mode  $(k_1, k_2)$  multiply the  $C$ -dimensional channel vector by a learnable dense matrix  $\mathbf{R}_{k_1, k_2}^{(l)} \in \mathbb{C}^{C \times C}$  (mixing channels) and hence the overall matrix is of size  $\mathbf{R}^{(l)} \in \mathbb{C}^{K \times K \times C \times C}$ , then apply an inverse FFT; (ii) adds a local pointwise transform, adds the injected encoder features  $\mathcal{Q}_\Phi(\mathbf{x})$ , and applies a nonlinearity. This realizes a resolution-invariant, globally receptive operator that naturally respects periodic boundary conditions.

Finally, we decode with the pointwise readout  $\mathcal{Q}_\Psi$  (a small per-cell two-layer MLP) to produce  $\mathbf{y} \in \mathbb{R}^{H \times W \times 1}$  where  $\Psi$  are learnable parameters.

All samples use  $H = W = 128$ . Unless stated otherwise, we set the latent width  $C = 32$ , retain  $K = 12$  Fourier modes per dimension in the FNO blocks, and use  $L = 3$  FNO layers inside  $\mathcal{G}_\Theta$ . Training differentiates implicitly through the fixed point, and the fixed-point solver uses Anderson acceleration. We optimize with Adam (learning rate  $5 \times 10^{-3}$ , batch size 16). For explicit baselines, we train for 500 epochs, which suffices for the training loss to converge.

**Perturbed data generation.** In this paragraph, we describe how we generate perturbed samples in  $\mathbb{D}'_{\text{pde, test}}$ . We take the dataset of [67] as the unperturbed set  $\mathbb{D}_{\text{pde, test}}$  and create perturbations by linearizing the steady NS equation (7). Each sample  $(f, \omega)$  comprises a forcing term  $f$  and its vorticity solution  $\omega$ . Directly prescribing  $f$  and solving for  $\omega$  is computationally costly; following [67], we instead prescribe  $\omega$  and obtain the corresponding  $f$  by evaluating the PDE operator (not by solving the PDE). In our setting, the base samples are given; thus we first construct a solution

TABLE 3. Explicit FNO vs Implicit FNO. The average is taken over all 500 test samples.

	Explicit (L=3)	Explicit (L=6)	Explicit (L=12)	Implicit (L=3)
Num of params	2.376 M	4.155 M	7.713 M	2.376 M
Avg. Relative Error	0.1787	0.1585	0.1547	0.0785

perturbation  $\delta\omega$  and then compute the induced forcing perturbation  $\delta f$  via the linearization, yielding the perturbed pair  $(f + \delta f, \omega + \delta\omega)$ .

Note that, while the dataset is discrete, we use the continuous notation  $f, \omega, u$  in this section to ease reading and to remain consistent with the PDE literature. In addition, we use  $\xi = (\xi_1, \xi_2)$  as the special domain variable to keep consistent with our main text, and use  $k = (k_1, k_2)$  as the frequency domain variable.

(Generate  $\delta\omega$ ). Fix a target wavenumber  $k_* \in \mathbb{N}$  and a desired  $L^2$ -magnitude  $\eta > 0$ . We construct  $\delta\omega$  by

$$\delta\omega(\xi_1, \xi_2) = A \sin(k_* \xi_1 + k_* \xi_2), \quad A \text{ chosen so that } \|\delta\omega\|_{L^2(\Omega)} = \eta.$$

The wavenumber is selected from a user-specified frequency percentile  $p_{\text{freq}}$  relative to the maximum resolvable frequency  $k_{\text{max}} = H/2 = W/2$ , namely

$$k_* = p_{\text{freq}} \times k_{\text{max}} \quad (\text{rounded to the nearest integer mode}).$$

In our code we set the grid size  $H = W = 128$ , the perturbation strength  $\eta = 0.01$ , and choose

$$p_{\text{freq}} \in \{0.01, 0.02, 0.03, 0.04, 0.05, 0.1, 0.3, 0.5, 0.7, 0.9, 0.95, 0.96, 0.97, 0.98, 0.99\}.$$

Accordingly, each original sample yields 15 perturbed samples.

(Generate velocity from vorticity). Given a scalar vorticity field  $\omega$  (and its perturbation  $\delta\omega$ ), we recover the corresponding velocities  $u$  via a streamfunction  $\psi$  given by

$$u = (\partial_2 \psi, -\partial_1 \psi), \quad \omega = -\Delta \psi.$$

Hence  $\psi$  is obtained by solving the Poisson equation  $\Delta \psi = -\omega$ , after which we obtain  $u$ . On a periodic grid, these operators are implemented efficiently in the Fourier domain.

(Linearization and the perturbed vorticity forcing  $\delta g$ ). Applying the (scalar) curl “ $\nabla \times$ ” to both sides of (7) yields the steady vorticity form

$$(u \cdot \nabla) \omega - \nu \Delta \omega = g, \quad g = \nabla \times f = \partial_1 f_2 - \partial_2 f_1,$$

where  $f = (f_1, f_2)$  is the body force and  $g$  is its curl. Introducing perturbations  $(\delta u, \delta\omega, \delta g)$  and expanding

$$(u + \delta u) \cdot \nabla (\omega + \delta\omega) - \nu \Delta (\omega + \delta\omega) = g + \delta g,$$

then subtracting the base equation and discarding higher-order terms gives the first-order relation

$$\delta g = (u \cdot \nabla) \delta\omega + (\delta u \cdot \nabla) \omega - \nu \Delta \delta\omega.$$

Again, for numerical implementation on a periodic grid, the differential operators are applied efficiently in the Fourier domain.

(Recover the vector force  $\delta f$  from its curl  $\delta g$ ). We recover a periodic  $\delta f = (\delta f_1, \delta f_2)$  satisfying  $\nabla \times \delta f = \delta g$  by solving a Poisson equation for an auxiliary streamfunction  $\psi$  and obtain  $\delta f$  exactly as in “Generate velocity from vorticity.”

**Additional Experiments.** A natural question is how stacking more explicit layers compares with an implicit model. We therefore increased the depth of the FNO (explicit) and contrasted it with an implicit FNO of comparable size. Table 3 shows that (i) making FNO deeper—hence adding parameters—only yield modest accuracy gains, and (ii) the implicit model achieves markedly better performance (lowest relative error) across all settings.

Note: These findings are broadly consistent with [67]. We follow their setup with two minor deviations: we use a smaller training batch size (16) due to hardware limits, and while we keep

$T = 24$  training iterations for the implicit model, at inference we run  $T = 50$ , because we observe that the trained implicit models remain stable and often benefit from additional fixed-point iterations at test time.

#### APPENDIX H. EXPERIMENT DETAILS REGARDING LP

**GNN model details.** We implement (10):

$$\mathbf{z}_* = \mathcal{G}_\Theta(\mathbf{z}_*, \mathcal{Q}_\Phi(\mathbf{x})), \quad \mathbf{y}_* = \mathcal{Q}_\Psi(\mathbf{z}_*)$$

with an  $L$ -layer message-passing GNN [85, 98] on the bipartite graph. Let  $\mathcal{N}(i)$  (resp.  $\mathcal{N}(j)$ ) be the neighbors of constraint node  $W_i$  (resp. variable node  $V_j$ ). With shared MLPs across all nodes and edges, the GNN structure is given by:

$$\begin{aligned} \text{Input-embedding:} \quad & W_i^{(0)} = \text{MLP}_{\phi_1}(b_i, \circ_i), \\ & V_j^{(0)} = \text{MLP}_{\phi_2}(c_j, l_j, u_j, z_{\text{in},j}) \\ \text{Message-passing } (1 \leq l \leq L-1): \quad & W_i^{(l)} = \text{MLP}_{\theta_1^{(l)}} \left( W_i^{(l-1)}, \sum_{j \in \mathcal{N}(i)} A_{ij} \cdot \text{MLP}_{\theta_2^{(l)}}(V_j^{(l-1)}) \right), \\ & V_j^{(l)} = \text{MLP}_{\theta_3^{(l)}} \left( V_j^{(l-1)}, \sum_{i \in \mathcal{N}(j)} A_{ij} \cdot \text{MLP}_{\theta_4^{(l)}}(W_i^{(l-1)}) \right) \\ \text{Output-embedding:} \quad & z_{\text{out},j} = \text{MLP}_{\theta_5}(V_j^{(L)}) \end{aligned}$$

We write this compactly as follows.

$$\mathbf{z}_{\text{out}} = \mathcal{G}_\Theta(\mathbf{z}_{\text{in}}, \mathcal{Q}_\Phi(\mathbf{x}))$$

where  $\Theta = \{\{\theta_1^{(l)}\}_{l=1}^{L-1}, \{\theta_2^{(l)}\}_{l=1}^{L-1}, \{\theta_3^{(l)}\}_{l=1}^{L-1}, \{\theta_4^{(l)}\}_{l=1}^{L-1}, \theta_5\}$  are trainable parameters in the GNN,  $\Phi = \{\phi_1, \phi_2\}$  includes the trainable parameters of the input embedding. The input  $\mathbf{x}$  includes all static information  $\mathbf{x} := (\mathbf{A}, \mathbf{b}, \mathbf{c}, \circ, \mathbf{l}, \mathbf{u})$ . Finally, the output embedding  $\mathbf{y} = \mathcal{Q}_\Psi(\mathbf{z})$  is given by

$$y_j = \text{MLP}_\Psi(z_j)$$

for every variable node  $j$ . All MLPs in  $\mathcal{G}_\Theta$ ,  $\mathcal{Q}_\Phi$ , and  $\mathcal{Q}_\Psi$  use two layers with ReLU activations. We sweep widths (or embedding sizes) in  $\{4, 8, 16, 32\}$  and report results in the main text.

Note that  $l$  is the layer index within the GNN structure, not the iteration number  $t$ . All parameters in  $\Theta$  are independent of the iteration number, so this GNN can be applied iteratively.  $\mathbf{x}$  is the static features and  $\mathbf{z}$  is the dynamic feature. In addition, removing the dynamic input  $z_{\text{in}}$  and decoding directly to  $\mathbf{y}$  recovers the standard (explicit) GNN baseline.

**Dataset generation.** We largely follow [17] to construct the training set  $\mathbb{D}_{\text{LP}, \text{train}}$  and test set  $\mathbb{D}_{\text{LP}, \text{test}}$ , drawing  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, \circ, \mathbf{l}, \mathbf{u})$  i.i.d. from the same distribution. Each LP has 50 variables and 10 constraints. The matrix  $\mathbf{A}$  is sparse with 100 nonzeros whose locations are chosen uniformly at random and whose values are sampled from a standard normal distribution. Entries of  $\mathbf{b}$  and  $\mathbf{c}$  are sampled i.i.d. from  $\text{Unif}[-1, 1]$ , after which  $\mathbf{c}$  is scaled by 0.01. Variable bounds  $\mathbf{l}, \mathbf{u}$  are sampled coordinatewise from  $\mathcal{N}(0, 10)$ ; whenever  $l_j > u_j$  we swap them. Constraint types are sampled independently with  $\Pr(\circ_i = “\leq”) = 0.7$  and  $\Pr(\circ_i = “=” ) = 0.3$ . Under this generator, the feasibility probability is approximately 0.53; we retain only feasible instances, yielding 2,500 LPs for training and 1,000 for testing. Solutions are computed with `scipy.optimize`.

To build the perturbed datasets  $\mathbb{D}_{\text{LP}, \text{test}}^{(j)}$ , we perturb one component at a time while holding the others fixed. For  $\mathbf{c}$ , draw  $\delta \mathbf{c}$  with i.i.d. standard normal entries, normalize, and scale to magnitude  $10^{-4}$ :

$$\mathbf{c}' = \mathbf{c} + 10^{-4} \times \frac{\delta \mathbf{c}}{\|\delta \mathbf{c}\|}.$$

We apply the same procedure to  $\mathbf{b}$ ,  $\mathbf{l}$ , and  $\mathbf{u}$ . For  $\mathbf{A}$ , we perturb only existing nonzeros to preserve the sparsity pattern: let  $\mathbb{S} = \{(i_k, j_k)\}_{k=1}^{\text{nnz}(\mathbf{A})}$  be the nonzero locations and draw  $\delta\mathbf{a} \in \mathbb{R}^{|\mathbb{S}|}$  i.i.d. standard normal; normalize and scale so  $\|\delta\mathbf{a}\| = 10^{-4}$ , then set

$$\mathbf{A}'_{i_k, j_k} = \mathbf{A}_{i_k, j_k} + (\delta\mathbf{a})_k \text{ for } (i_k, j_k) \in \mathbb{S}, \quad \mathbf{A}'_{i, j} = \mathbf{A}_{i, j} \text{ otherwise.}$$

This yields five perturbed versions (perturbing  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{l}$ , or  $\mathbf{u}$  separately). We evaluate the estimated Lipschitz constants  $L_t$  and relative errors  $E_t$  on each version and report the results in the main text.

**Training method.** To train our implicit GNNs, we employ a two-stage curriculum strategy. The model is trained by unrolling its iterative updates for a fixed number of steps,  $T$ , and minimizing the loss on the final output:

$$\begin{aligned} & \min_{\Theta, \Phi, \Psi} \sum_{(\mathbf{x}, \mathbf{y}_*) \in \mathbb{D}_{\text{LP, train}}} \ell(\mathbf{y}_T, \mathbf{y}_*) \\ & \text{s.t. } \mathbf{z}_0 = \mathbf{0} \\ & \quad \mathbf{z}_t = \mathcal{G}_{\Theta}(\mathbf{z}_{t-1}, \mathcal{Q}_{\Phi}(\mathbf{x})), \quad t = 1, 2, \dots, T \\ & \quad \mathbf{y}_T = \mathcal{Q}_{\Psi}(\mathbf{z}_T) \end{aligned}$$

We set the final unroll horizon to  $T = 6$ , as we observed no significant improvements with longer sequences. Training directly with  $T = 6$  is inefficient, so we adopt a two-stage curriculum. This approach is a standard practice in the Learning to Optimize field for training implicit or unrolled models that solve optimization problems [13]. This approach begins with a shorter unroll horizon and a larger learning rate, using the trained model to warm-start the subsequent stage with a longer horizon and a reduced learning rate. This strategy is often described as “layerwise training” [15, 59] or “curriculum learning” [14, 40, 60]. In our settings: Stage 1 uses  $T = 3$  with a learning rate 0.01; Stage 2 uses  $T = 6$  with a learning rate  $10^{-4}$ . Both stages use Adam optimizer.

For a fair comparison, the non-iterative explicit GNNs are trained using the same two-stage learning rate schedule. This regimen proved effective, as the training errors for our explicit baselines surpassed those reported in prior work [17].

At the inference time,  $T$  can be chosen as the unroll length in the training stage, or moderately longer. In our experiments, we use  $T = 8$  at the inference time, as we do not observe significant improvement with a larger number of iterations.

*Remark.* We train with unrolling rather than classic deep equilibrium/implicit-differentiation, but we still refer to the model as “implicit.” The architecture is identical in both cases—a weight-tied update  $\mathbf{z}_t = \mathcal{G}_{\Theta}(\mathbf{z}_{t-1}, \mathcal{Q}_{\Phi}(\mathbf{x}))$  with a readout  $\mathbf{y}_T = \mathcal{Q}_{\Psi}(\mathbf{z}_T)$ ; only the training procedure differs (backprop through  $T$  steps vs. implicit differentiation). Since our focus is expressivity rather than training mechanics, we treat both weight-tied unrolled models and DEQs as the same class of implicit models.

(JL) DEPARTMENT OF STATISTICS AND DATA SCIENCE, UNIVERSITY OF CENTRAL FLORIDA, ORLANDO, FL 32826.  
Email address: jialin.liu@ucf.edu

(LD) DEPARTMENT OF MATHEMATICS, UNIVERSITY OF CALIFORNIA, LOS ANGELES, LOS ANGELES, CA 90095  
Email address: lisangding@ucla.edu

(WY) DECISION INTELLIGENCE LAB, DAMO ACADEMY, ALIBABA US, BELLEVUE, WA 98004.  
Email address: wotao.yin@alibaba-inc.com

(SO) DEPARTMENT OF MATHEMATICS, UNIVERSITY OF CALIFORNIA, LOS ANGELES, LOS ANGELES, CA 90095  
Email address: sjo@math.ucla.edu