

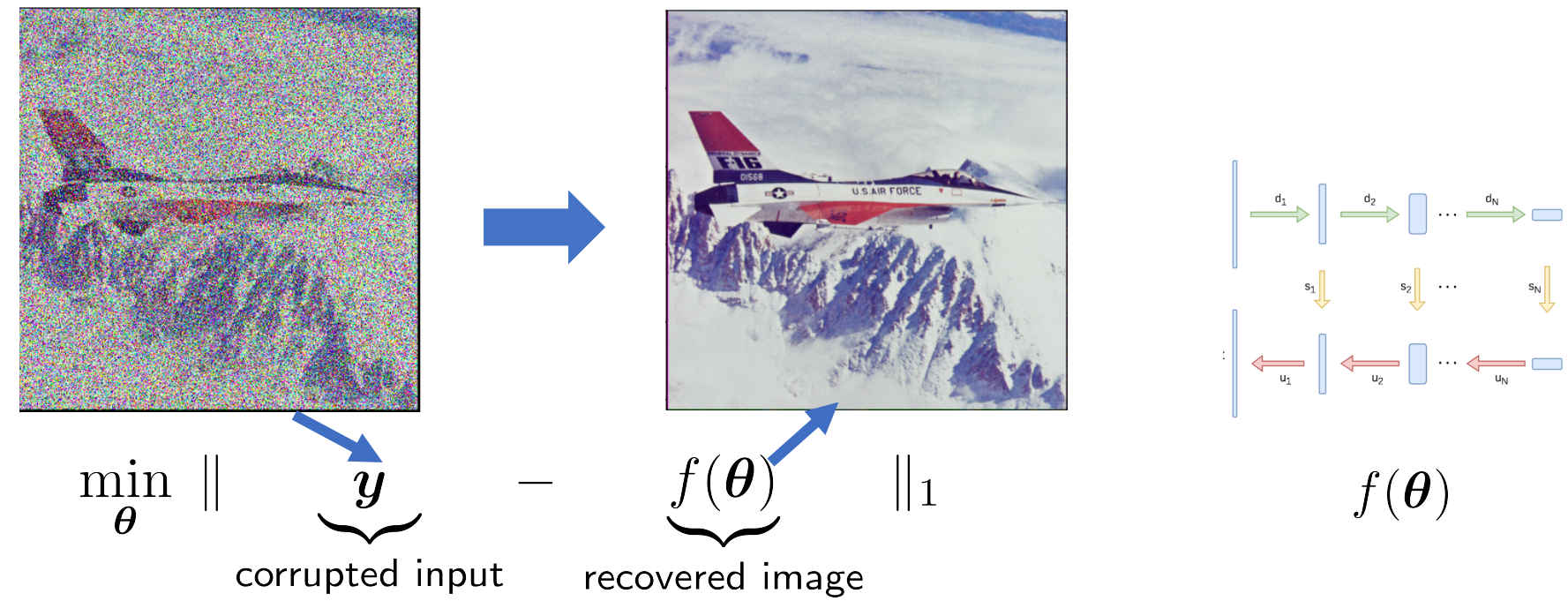
Robust Recovery via the Implicit Bias of Discrepant Learning Rates for Double Over-parameterization

Chong You[†] Zihui Zhu[‡] Qing Qu^{*} Yi Ma[†]
[†]University of California, Berkeley [‡]Denver University ^{*}New York University

Introduction

Deep neural networks are highly **over-parameterized**

- **Background:** Image recovery via *deep image prior* (DIP) [1]



Idea: network architecture (i.e., $f(\theta)$) encodes priors for **clean** images

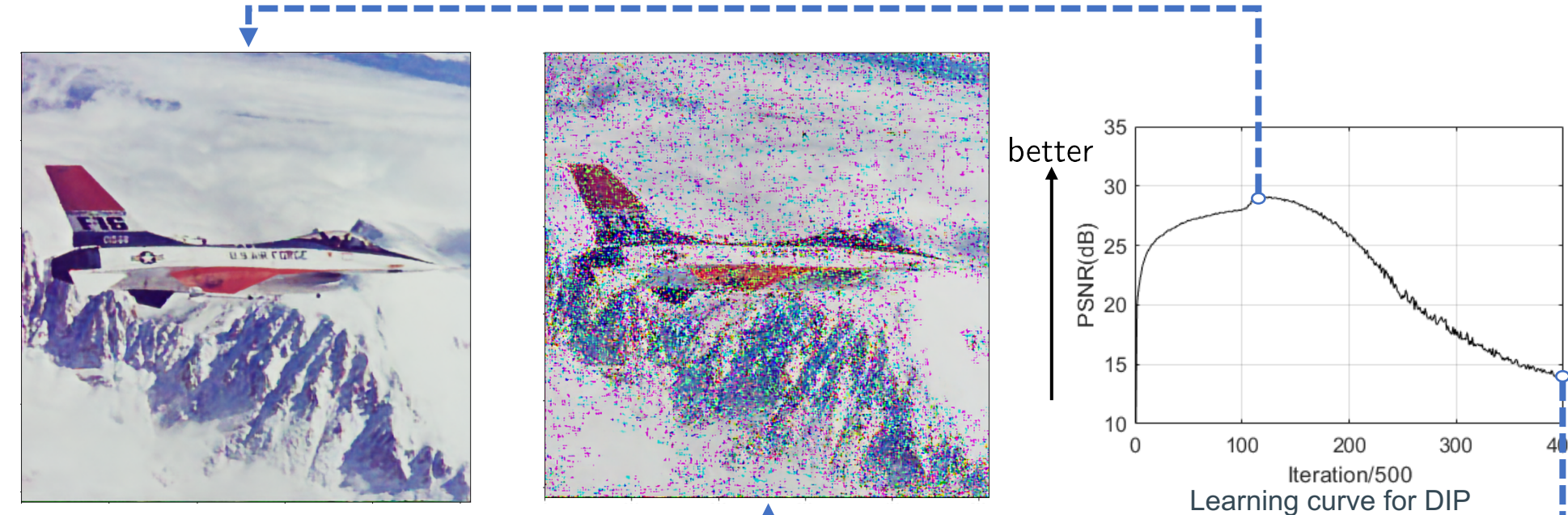
- Over-parameterization in DIP:

~ 2 million \gg ~ 0.1 million
 (#parameters in $f(\theta)$) (# pixels in an image)

In principle, $f(\theta)$ can generate **any** image!
 (i.e., not only the desired **clean** images, but also undesired **corrupted** images)

Challenge

Over-parameterization **causes** overfitting



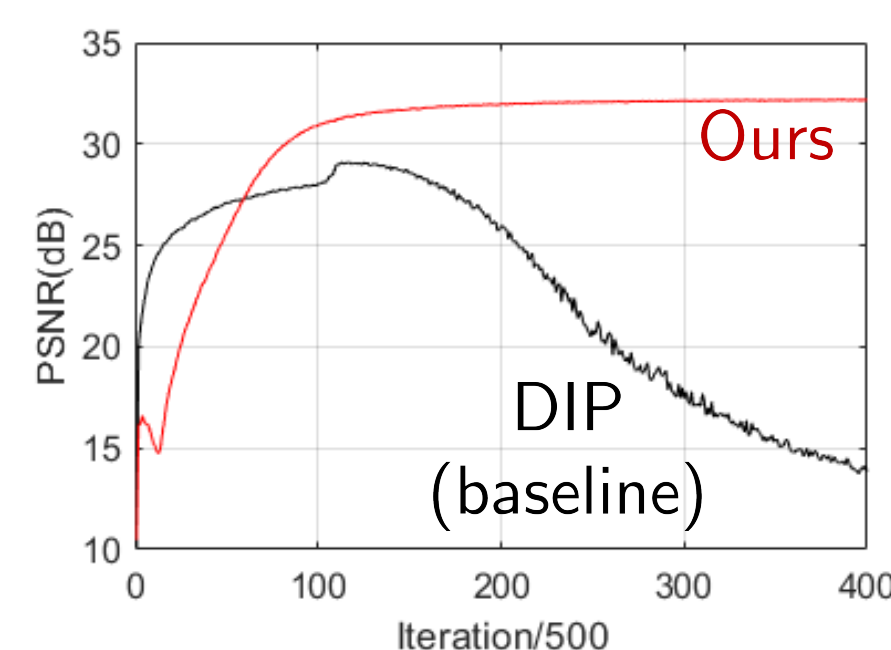
Early termination solution (impractical!) Global solution: $f(\theta) \approx y$ (overfitting!)

Contribution

Over-parameterization **WITHOUT** overfitting!

- **Model:** Double over-parameterization of both the image and the corruption
- **Algorithm:** Discrepant learning rates for different model parameters
- **Theory:** Correctness for low-rank recovery

No need to terminate early
No need to tune network width
No need to finetune parameters



Model: Double Over-Parameterization (DOP)

Goal: Recover an image $X_* \in \mathbb{R}^{H \times W \times 3}$ from $y = X_* + s_*$, where s_* is sparse noise

$$\ell(\theta, g, h) = \min_{\theta, g, h} \|y - (f(\theta) + g \odot g - h \odot h)\|_2^2 \quad (1)$$



- $f(\theta)$: U-shaped neural network (as in DIP) that *over-parameterizes* the image X
- $g \odot g - h \odot h$ (\odot - Hadamard product) *over-parameterizes* sparse corruption s

(Seemingly) Problematic:

Even more parameters than DIP,
 even more global solutions than DIP.
 Hence, even more prone to **overfitting** than DIP?

Algorithm: Discrepant Learning Rates

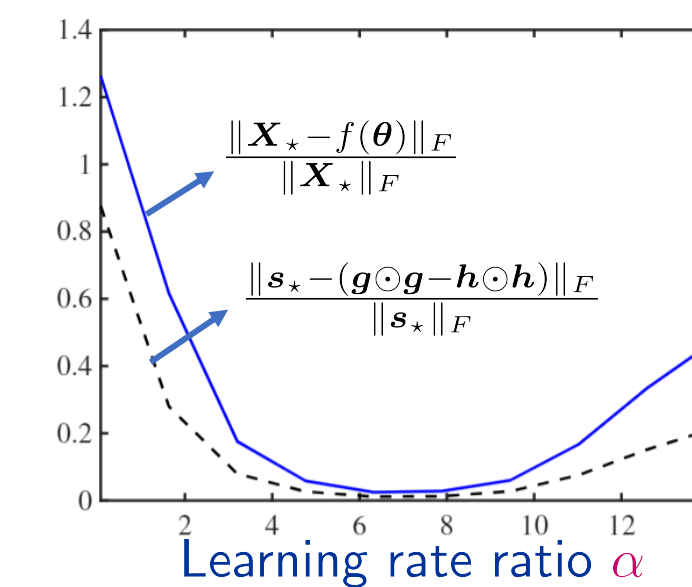
(Genuinely) a Blessing: Proper choice of algorithm leads to the desired solutions!

- **Observation:** Learning rate ratio α for θ and $\{g, h\}$ in a gradient descent

$$\begin{aligned} \theta_{k+1} &= \theta_k - \tau \cdot \nabla_{\theta} \ell(\theta_k, g_k, h_k) \\ \begin{bmatrix} g_{k+1} \\ h_{k+1} \end{bmatrix} &= \begin{bmatrix} g_k \\ h_k \end{bmatrix} - \alpha \cdot \tau \cdot \begin{bmatrix} \nabla_g \ell(\theta_k, g_k, h_k) \\ \nabla_h \ell(\theta_k, g_k, h_k) \end{bmatrix} \end{aligned} \quad (2)$$

controls the quality of the solution

- **No parameter tuning:** Empirically, best α does NOT depend on a) test image X_* , b) sparsity of s_* , and c) width of network $f(\theta)$



Theory: Insights from Low-rank Modeling

Consider the case $X_* \in \mathbb{R}^{n \times n}$ with rank $r \ll n$. Let $f(\theta) = UU^T$, where $U \in \mathbb{R}^{n \times r'}$

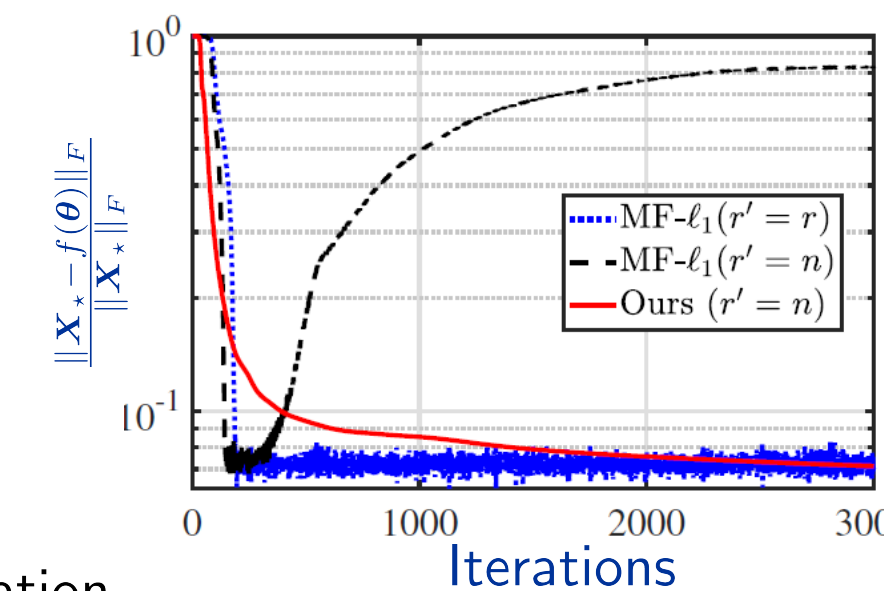
- **Classical method** based on ℓ_1 loss:

$$\min_{U \in \mathbb{R}^{n \times r'}} \|y - UU^T\|_1$$

- works for exact-param. models (i.e., $r' = r$)
- fails for over-param. models (e.g., $r' = n$)

Failure of classical robust loss methods for **modern** over-parameterized models!

- **Our method** based on double over-parameterization



Theorem: Algorithm (2) for solving problem (1) with $r' = n$ converges to a solution to

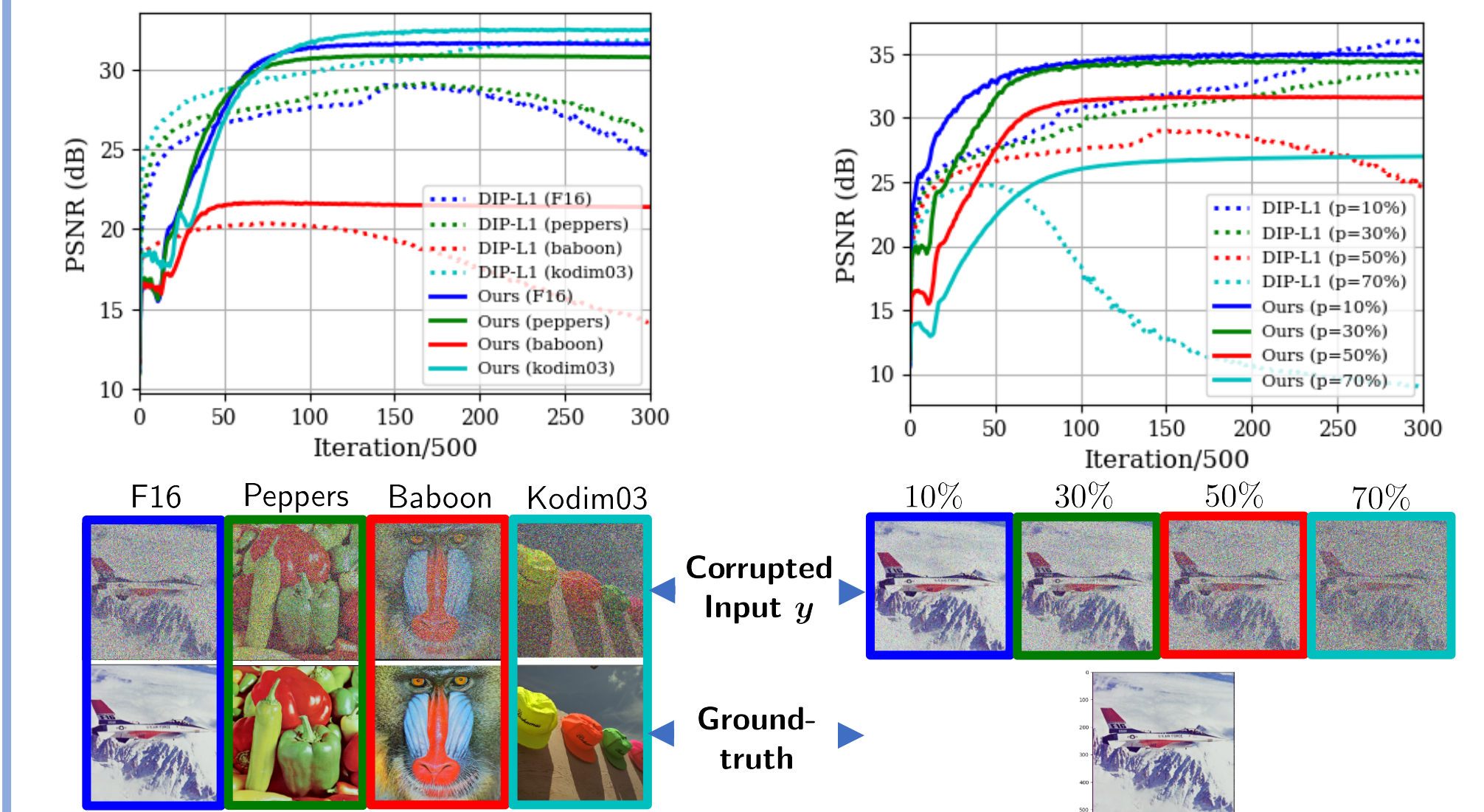
$$\min_{X, s} \|X\|_* + 1/\alpha \cdot \|s\|_1 \quad \text{s.t. } y = X + s$$

Combined with [3], setting $\alpha = \sqrt{n}$ produces desired solutions (**No tuning parameters!**)

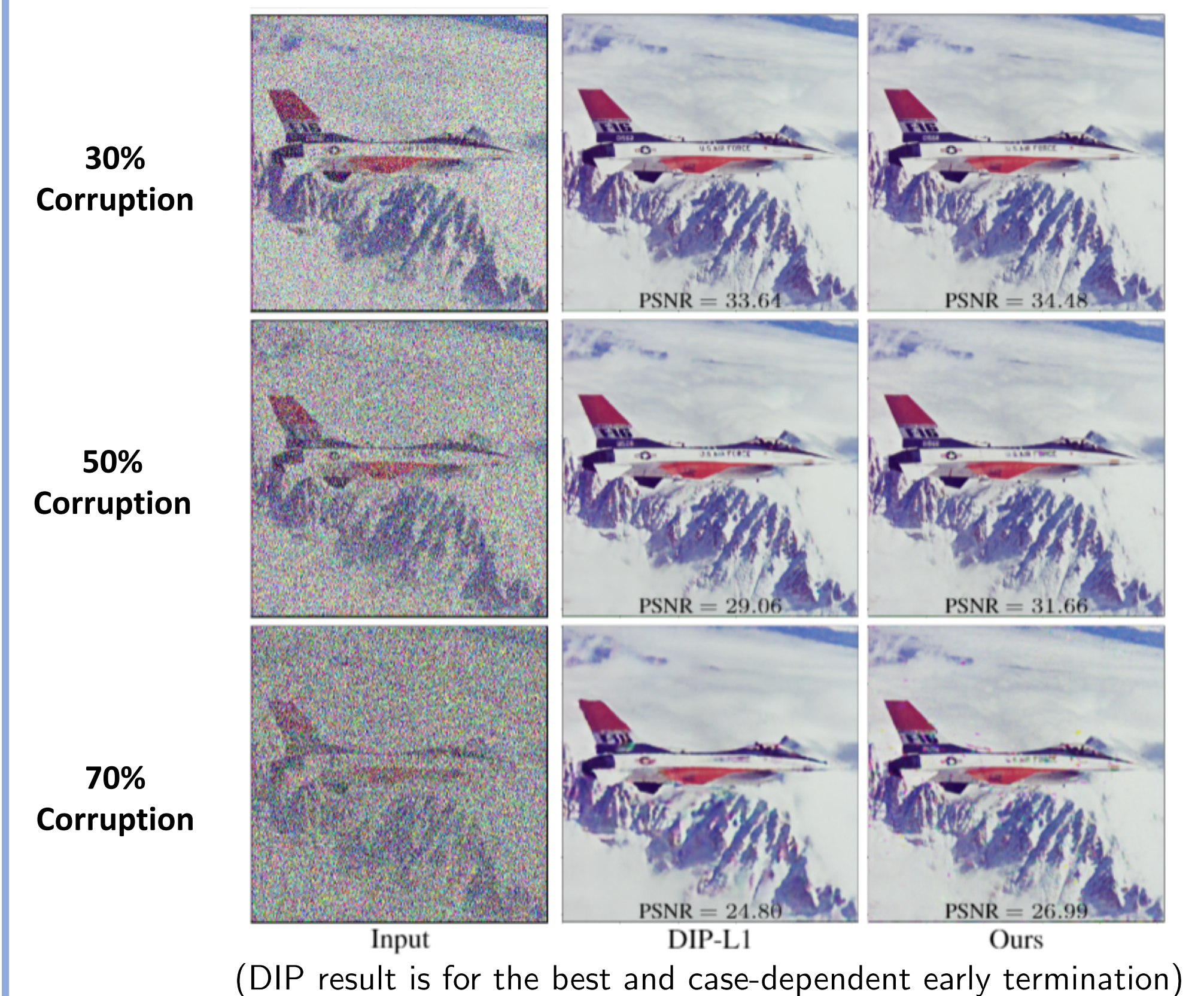
Experiments

Learning curves for varying images (left) and varying corruption levels (right)

- For DIP, best termination varies for different images and corruption levels
- For Ours, no need to terminate and no need to tune any parameters



Visualization of results with varying corruption levels



References

- [1] Ulyanov, Vedaldi, Lempitsky (2018). "Deep Image Prior." In: CVPR
- [2] Li, Zhu, So, Vidal (2019). "Robust Low-rank Matrix Recovery." In: SIAM Journal on Optimization
- [3] Candes, Li, Ma, Wright (2011). "Robust Principal Component Analysis?" In: Journal of ACM