# Instance-Optimal Compressed Sensing via Posterior Sampling

**Ajil Jalal** [1]   **Sushrut Karmalkar** [2]   **Alexandros G. Dimakis** [1]   **Eric Price** [2]

## Abstract

We characterize the measurement complexity of compressed sensing of signals drawn from a known prior distribution, even when the support of the prior is the entire space (rather than, say, sparse vectors). We show for Gaussian measurements and *any* prior distribution on the signal, that the posterior sampling estimator achieves near-optimal recovery guarantees. Moreover, this result is robust to model mismatch, as long as the distribution estimate (e.g., from an invertible generative model) is close to the true distribution in Wasserstein distance. We implement the posterior sampling estimator for deep generative priors using Langevin dynamics, and empirically find that it produces accurate estimates with more diversity than MAP.

## 1. Introduction

The goal of compressed sensing is to recover a structured signal from a relatively small number of linear measurements. The setting of such linear inverse problems has numerous and diverse applications ranging from Magnetic Resonance Imaging (Lustig et al., 2008; 2007), neuronal spike trains (Hegde et al., 2009) and efficient sensing cameras (Duarte et al., 2008). Estimating a signal in $\mathbb{R}^n$ would in general require $n$ linear measurements, but because real-world signals are structured—i.e., compressible—one is often able to estimate them with $m \ll n$ measurements.

Formally, we would like to estimate a "signal" $x^* \in \mathbb{R}^n$ from noisy linear measurements,

$$y = Ax^* + \xi$$

for a measurement matrix $A \in \mathbb{R}^{m \times n}$ and noise vector $\xi \in \mathbb{R}^m$. We will focus on the i.i.d. Gaussian setting, where $A_{ij} \sim \mathcal{N}(0, \frac{1}{m})$ and $\xi_i \sim \mathcal{N}(0, \frac{\sigma^2}{m})$, and one would like to recover $\widehat{x}$ from $(A, y)$ such that

$$\|x^* - \widehat{x}\| \le C\sigma \qquad (1)$$

with high probability for some constant $C$. When $x^*$ is $k$-sparse, this was shown by Candés, Romberg, and Tao (Candes et al., 2006) to be possible for $m$ at least $O(k \log \frac{n}{k})$.

Over the past 15 years, compressed sensing has been extended in a wide variety of remarkable ways, including by generalizing from sparsity to other signal structures, such as those given by trees (Chen & Huang, 2012), graphs (Xu et al., 2011), manifolds (Chen et al., 2010; Xu & Hassibi, 2008), or deep generative models (Bora et al., 2017; Asim et al., 2019). These are all essentially frequentist approaches to the problem: they define a small *set* of "structured" signals $x$, and ask for recovery of every such signal.

Such set-based approaches have limitations. For example, (Bora et al., 2017) uses the structure given by a deep generative model $G : \mathbb{R}^k \to \mathbb{R}^n$; with $O(kd \log n)$ measurements for $d$-layer networks, accurate recovery is guaranteed for every signal $x^*$ near the range of $G$. But this completely ignores the *distribution* over the range. Generative models like Glow (Kingma & Dhariwal, 2018) and pixelRNN (Oord et al., 2016) have seed length $k = n$ and range equal to the entire $\mathbb{R}^n$. Yet because these models are designed to approximate reality, and real images can be compressed, we know that compressed sensing is possible in principle.

This leads to the question: Given signals drawn from some *distribution* $R$, can we characterize the number of linear measurements necessary for recovery, with both upper and lower bounds? Such a Bayesian approach has previously been considered for sparsity-inducing product distributions (Aeron et al., 2010; Zhou et al., 2014) but not general distributions.

Second, suppose that we don't know the real distribution $R$, but instead have an approximation $P$ of $R$ (e.g., from a GAN or invertible generative model). In what sense should $P$ approximate $R$ for compressed sensing with good guarantees to be possible?

Figure 1: Reconstruction results on FFHQ for Gaussian measurements (here $n = 256 \times 256 \times 3 = 196,608$ pixels), using an NCSNv2 model. Each column shows the reconstruction obtained as the number of measurements $m$ varies. The top row shows reconstructions by MAP, the middle row shows reconstruction by Deep-Decoder, and the bottom row shows reconstructions by Langevin dynamics, which is the practical implementation of our proposed posterior sampling estimator.

## 1.1. Contributions.

Our main theorem is that posterior sampling is a near optimal recovery algorithm for *any* distribution. Moreover, it is sufficient to learn the distribution in Wasserstein distance.

**Theorem 1.1.** *Let $R$ be an arbitrary distribution over an $\ell_2$ ball of radius $r$. Suppose that there exists an algorithm that uses an arbitrary measurement matrix $A \in \mathbb{R}^{m \times n}$ with noise level $\sigma$ and finds a reconstruction $\widehat{x}$ such that*

$$\|x^* - \widehat{x}\| \lesssim \sigma \text{ with probability } \geq 1 - \delta.$$

*Then posterior sampling (see Definition 1.3) with respect to $R$ using $m' \geq O\left(m \log\left(1 + \frac{mr^2\|A\|_\infty^2}{\sigma^2}\right) + \log\frac{1}{\delta}\right)$ Gaussian measurements of noise level $\sigma$ will output $\widehat{x}$ satisfying*

$$\|x^* - \widehat{x}\| \lesssim \sigma \text{ with probability } \geq 1 - O(\delta).$$

*Moreover, the same holds for posterior sampling with respect to any distribution $P$ satisfying $\mathcal{W}_p(R, P) \lesssim \sigma\delta^{1/p}$ for some $p \geq 1$.*

This theorem comprises three main contributions: the introduction of posterior sampling as *a new algorithm* for recovery with a generative prior; an *upper bound* on the sample

complexity of the algorithm in terms of an approximate covering number that we introduce; and an *instance-optimal lower bound* in terms of the same approximate covering number that (unlike previous lower bounds in compressed sensing) applies to *any* distribution of input signals.

**Contribution 1: Approximate covering numbers.** The covering number of a set is the smallest number of balls that can cover the entire set. Standard compressed sensing is closely tied to the covering number $N_\eta(S)$ of the set $S$ of possible signals $x$; for example, the set of unit-norm $k$-sparse vectors has $\log N_\eta = \Theta(k \log \frac{n}{k})$, which is precisely why Candés, Romberg, and Tao use this many linear measurements to achieve (1).

For distributions, we need a different concept of covering number. As a motivating example, consider a distribution $R$ induced by a trivial *linear* generative model, $x = \Sigma z$ where $z \sim \mathcal{N}(0, I_n)$ and $\Sigma$ is a fixed $n \times n$ matrix. Further suppose the singular values $\sigma_i$ of $\Sigma$ are Zipfian, so $\sigma_i = 1/i$. In this case, $R$'s support is $\mathbb{R}^n$, so covering the entire support of $R$ is infeasible. Instead we could denote by $\text{Cov}_{\eta, 0.01}(R)$ the minimum number of $\eta$-radius balls needed to cover 99%
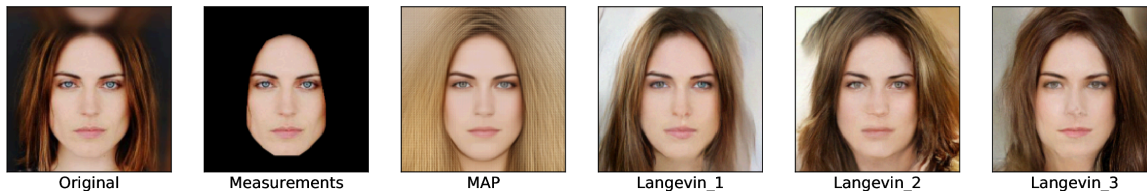
Figure 2: Reconstruction results for inpainting on CelebA-HQ using Glow. The first column shows the original image, second column shows the measurements by removing the hair and background, the third column shows reconstruction by MAP, and the last three columns show samples from posterior sampling via Langevin dynamics. MAP produces the same washed out image all the time, whereas posterior sampling produces images with diversity.

of $R$. An elementary calculation shows

$$\log \mathrm{Cov}_{\eta, 0.01}(R) = \Theta(1/\eta^2),$$

which is (up to constants) precisely the number of linear measurements you need to estimate $x$ to within $\eta$.

We show that an *approximate covering number* characterizes the measurement complexity of compressed sensing a general distribution $R$, and that recovery by *posterior sampling* achieves this bound.

**Definition 1.2.** *Let $R$ be a distribution on $\mathbb{R}^n$. For some parameters $\eta > 0, \delta \in [0, 1]$, we define the $(\eta, \delta)$-approximate covering number of $R$ as*

$$\mathrm{Cov}_{\eta, \delta}(R) := \min \left\{ k : R\left[\cup_{i=1}^k \mathcal{B}(x_i, \eta)\right] \geq 1 - \delta, x_i \in \mathbb{R}^n \right\}$$

*where $\mathcal{B}(x, \eta)$ is the $\ell_2$ ball of radius $\eta$ centered at $x$.*

When $\delta = 0$, this is $N_\eta(\mathrm{supp}\, R)$, the standard covering number of the support of $R$. Having $\delta > 0$ allows meaningful results for full-support distributions that are concentrated on smaller sets. This also generalizes our previous results in (Bora et al., 2017), which depend on the covering numbers of low-dimensional generative models.

**Contribution 2: Recovery algorithm.** The recovery algorithm we consider is posterior sampling:

**Definition 1.3.** *Given an observation $y$, the* posterior sampling *recovery algorithm with respect to $P$ outputs $\widehat{x}$ according to the posterior distribution $P(\cdot \mid y)$.*

**Contribution 3: Sample complexity upper bound.** Our main positive result is that posterior sampling achieves the guarantees of equation (1) for *general* distributions $R$, with $O(\log \mathrm{Cov}_{\sigma, \delta}(R))$ measurements. Not only this, but the algorithm is robust to model mismatch: posterior sampling with respect to $P \neq R$ still works, as long as $P$ and $R$ are close in Wasserstein distance:

**Theorem 1.4** (Upper bound). *Let $P$, $R$ be distributions with $\mathcal{W}_1(P, R) \leq \sigma$. Let $x^* \sim R$, let $y$ be Gaussian measurements with noise level $\sigma$, and let $\widehat{x} \sim P(\cdot|y)$. For any*

$\eta \geq \sigma$, *with*

$$m \geq O(\log \mathrm{Cov}_{\eta, 0.01}(R))$$

*measurements, the guarantee $\|\widehat{x} - x^*\| \leq C\eta$ is satisfied for some universal constant $C$ with $97\%$ probability over the signal $x$, measurement matrix $A$, noise $\xi$, and recovery algorithm $\widehat{x}$.*

**Contribution 4: Sample complexity lower bound.** Our second main result lower bounds the sample complexity for *any* distribjution. This is, to our knowledge, the first lower bound for compressed sensing that applies to arbitrary distributions $R$. Most lower bounds in the area are minimax, and only apply to specific "hard" distributions $R$ (Price & Woodruff, 2011; Candes & Davenport, 2013; Iwen & Tewfik, 2010); the closest result we are aware of is (Aeron et al., 2010), which characterizes product distributions.

**Theorem 1.5** (Lower bound). *Let $R$ be any distribution over an $\ell_2$ ball of radius $r$, and consider any method to achieve $\|\widehat{x} - x^*\| \leq \eta$ with $99\%$ probability, using an arbitrary measurement matrix $A \in \mathbb{R}^{m \times n}$ with noise level $\sigma$. This must have*

$$m \geq \frac{C'}{\log(1 + \frac{mr^2\|A\|_\infty^2}{\sigma^2})} \log \mathrm{Cov}_{C'\eta, 0.04}(R).$$

*for some constant $C' > 0$.*

Note that Theorem 1.4 and 1.5 directly give Theorem 1.1. For more precisely stated and general versions of these results, including dependence on the failure probability $\delta$, see Theorems 3.4 and 4.1.

## 1.2. Related Work

Generative priors have shown great promise in compressed sensing and other inverse problems, starting with (Bora et al., 2017), who generalized the theoretical framework of compressive sensing and restricted eigenvalue conditions (Tibshirani, 1996; Donoho, 2006; Bickel et al., 2009; Candes, 2008; Hegde et al., 2008; Baraniuk & Wakin, 2009; Baraniuk et al., 2010; Eldar & Mishali, 2009) for signals lying

on the range of a deep generative model (Goodfellow et al., 2014; Kingma & Welling, 2013).

Lower bounds in (Kamath et al., 2019; Liu & Scarlett, 2019; Jalali & Yuan, 2019) established that the sample complexities in (Bora et al., 2017) are order optimal. The approach in (Bora et al., 2017) has been generalized to tackle different inverse problems such as robust compressed sensing (Jalal et al., 2020), phase retrieval (Hand et al., 2018; Aubin et al., 2019; Jagatap & Hegde, 2019), blind image deconvolution (Asim et al., 2018), seismic inversion (Mosser et al., 2020), one-bit recovery (Qiu et al., 2019; Liu et al., 2020), and blind demodulation (Hand & Joshi, 2019). Alternate algorithms for reconstruction include sparse deviations from generative models (Dhar et al., 2018), task-aware compressed sensing (Kabkab et al., 2018), PnP (Pandit et al., 2019; Fletcher et al., 2018b;a), iterative projections (Mardani et al., 2018), OneNet (Rick Chang et al., 2017) and Deep Decoder (Heckel & Hand, 2018; Heckel & Soltanolkotabi, 2020). The complexity of optimization algorithms using generative models have been analyzed for ADMM (Gómez et al., 2019), PGD (Hegde, 2018), layer-wise inversion (Lei et al., 2019), and gradient descent (Hand & Voroninski, 2017). Experimental results in (Asim et al., 2019; Whang et al., 2020; Lindgren et al., 2020) show that invertible models have superior performance in comparison to low dimensional models. See (Ongie et al., 2020) for a more detailed survey on deep learning techniques for compressed sensing. A related line of work has explored learning-based approaches to tackle classical problems in algorithms and signal processing (Aamand et al., 2019; Indyk et al., 2019; Metzler et al., 2017; Hsu et al., 2018).

Lower bounds for $\ell_2/\ell_2$ recovery of sparse vectors can be found in (Scarlett & Cevher, 2016; Price & Woodruff, 2011; Aeron et al., 2010; Iwen & Tewfik, 2010; Candes & Davenport, 2013), and these are related to the lower bound in (1.5). The closest result is that of (Aeron et al., 2010), which characterizes the probability of error and $\ell_2$ error of the reconstruction via covering numbers of the probability distribution. Their approach uses the rate distortion function of a scalar random variable $\mathbf{x}$, and provides guarantees for the product measure generated via an i.i.d. sequence of $\mathbf{x}$. A Shannon theory for compressed sensing was pioneered by (Wu & Verdú, 2012; Wu, 2011). The $\delta-$Minkowski dimension of a probability measure used in (Wu & Verdú, 2012; Wu, 2011; Pesin, 2008) can be derived from our $(\varepsilon, \delta)-$covering number by taking the limit $\varepsilon \to 0$. (Reeves & Gastpar, 2012) contains a related theory of rate distortion for compressed sensing. There is also related work in the statistical physics community under different assumptions on the signal structure (Zdeborová & Krzakala, 2016; Barbier et al., 2019).

## 2. Background and Notation

In this section, we introduce a few concepts that we will use throughout the paper. $\| \cdot \|$ refers to the $\ell_2$ norm unless specified otherwise. The metric we use to quantify the similarity between distributions is the Wasserstein distance. For two probability distributions $\mu, \nu$ supported on $\Omega$, and for any $p \geq 1$, the Wasserstein-$p$ (Villani, 2008; Arjovsky et al., 2017) and Wasserstein-$\infty$ (Champion et al., 2008) distances are defined as:

$$\mathcal{W}_p(\mu, \nu) := \inf_{\gamma \in \Pi(\mu, \nu)} \left( \mathbb{E}_{(u,v) \sim \gamma} \left[ \|u - v\|^p \right] \right)^{1/p},$$

$$\mathcal{W}_\infty(\mu, \nu) := \inf_{\gamma \in \Pi(\mu, \nu)} \left( \gamma\text{-ess}\sup_{(u,v) \in \Omega^2} \|u - v\| \right),$$

where $\Pi(\mu, \nu)$ denotes the set of joint distributions whose marginals are $\mu, \nu$. The above definition says that if $\mathcal{W}_\infty(\mu, \nu) \leq \varepsilon$, and $(u, v) \sim \gamma$, then $\|u - v\| \leq \varepsilon$ almost surely.

We say that $y$ is generated from $x^*$ by a Gaussian measurement process with $m$ measurements and noise level $\sigma$, if $y = Ax^* + \xi$ where $\xi \sim \mathcal{N}(0, \frac{\sigma^2}{m} I_m)$ and $A \in \mathbb{R}^{m \times n}$ with $A_{ij} \sim \mathcal{N}(0, 1/m)$.

## 3. Upper Bound

### 3.1. Two-Ball Case

For simplicity, we will first demonstrate our proof techniques in the simple setting where $R = P$, the measurements are noiseless, and the ground truth distribution $P$ is supported on two disjoint balls (illustrated in Figure 3). In this example, two $\eta$ radius balls can cover the whole space, so the parameters in Theorem 1.4 will be $\sigma = 0$ and $\text{Cov}_{\eta,0}(P) = 2$. Applying Theorem 1.4 on $P$ tells us that a constant number of measurements is sufficient for posterior sampling to get $O(\eta)$-close to the ground truth, i.e., to return an element of the correct ball. We will now prove this claim.

Let $B_0, B_{\tilde{x}}$ denote $\eta$-radius balls centered at $0, \tilde{x} \in \mathbb{R}^n$ respectively. Suppose $P = 0.5P_0 + 0.5P_1$, where $P_0, P_1$, are uniform distributions on $B_0, B_{\tilde{x}}$. The centers of the balls are separated by a distance $d \gg \eta$.

The ground truth $x^*$ will be sampled from $P$. For a fixed matrix $A \in \mathbb{R}^{m \times n}$ with $m \ll n$, let the noiseless measurements be $y = Ax^*$ and let $H_0, H_1$, denote the distributions over $\mathbb{R}^m$ induced by the projection of $P_0, P_1$, by $A$.

Given $A, y$, we sample the reconstruction $(\widehat{x})$ according to the posterior density

$$p(\widehat{x}|y) = c_y p_0(\widehat{x}|y) + (1 - c_y) p_{\tilde{x}}(\widehat{x}|y),$$

where $c_y$ is the posterior probability that $y$ is a projection of $x^*$ drawn from the $P_0$ component of $P$. Note that $c_y$ depends on $y$.
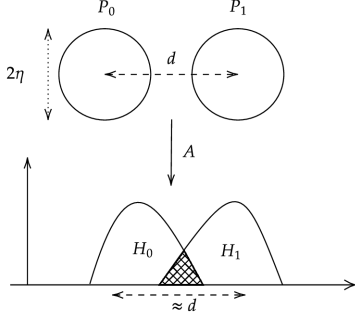
Figure 3: Illustrative example for the upper bound. The signal $x^*$ is drawn from a mixture of two well-separated balls. The observations $y = Ax^*$ are then drawn from a mixture of two distributions $H_0, H_1$ that may overlap. The probability that posterior sampling outputs something from the wrong ball is proportional to the (shaded) overlap between these distributions, which is atmost $1 - TV(H_0, H_1)$.

Since the balls $B_0 \& B_{\tilde{x}}$ are well separated, the ground truth and the reconstruction are far apart if and only if they lie in different balls, i.e., $\{x^* \in B_0, \hat{x} \in B_{\tilde{x}}\}$, or vice versa. It turns out quite generally that the probability of this event is bounded by how similar the distributions $H_0, H_1$ are:

**Lemma 3.1.** *For $c \in [0, 1]$, let $H := (1 - c)H_0 + cH_1$ be a mixture of two absolutely continuous distributions $H_0, H_1$ admitting densities $h_0, h_1$. Let $y$ be a sample from the distribution $H$, such that $y|z^* \sim H_{z^*}$ where $z^* \sim Bernoulli(c)$.*

*Define $\hat{c}_y = \frac{ch_1(y)}{(1-c)h_0(y) + ch_1(y)}$, and let $\hat{z}|y \sim Bernoulli(\hat{c}_y)$ be the posterior sampling of $z^*$ given $y$. Then we have*

$$\Pr_{z^*, y, \hat{z}}[z^* = 0, \hat{z} = 1] \leq 1 - TV(H_0, H_1).$$

The proof of this, as well as all parts of the upper bound, can be found in Appendix A.

In our current example, this gives us

$$\Pr[x^* \in B_0, \hat{x} \in B_{\tilde{x}}] \leq 1 - TV(H_0, H_1) \text{ and}$$
$$\Pr[x^* \in B_{\tilde{x}}, \hat{x} \in B_0] \leq 1 - TV(H_0, H_1).$$

Since $B_0$ and $B_{\tilde{x}}$ are balls of radius $\eta$, a union bound of the above two probabilities gives:

$$\Pr\left[\|x^* - \hat{x}\| > 2\eta\right] \leq \Pr\left[x^* \in B_0, \hat{x} \in B_{\tilde{x}}\right] +$$
$$\Pr\left[x^* \in B_{\tilde{x}}, \hat{x} \in B_0\right],$$
$$\leq 2\left(1 - TV\left(H_0, H_1\right)\right). \quad (2)$$

If $A$ is a Gaussian random matrix, the Johnson-Lindenstrauss (JL) Lemma tells us that it will preserve distances between vectors with high probability . This does

not necessarily mean that every point in the distribution $P$ will be preserved in norm. Still, we show that, since $P_0$ and $P_1$ have well-separated supports, their projected distributions $H_0$ & $H_1$ have very high TV distance. This also holds more generally, between any distribution on a ball and any distribution far from the ball and in the presence of noise.

**Lemma 3.2.** *Let $y$ be generated from $x^*$ by a Gaussian measurement process with noise level $\sigma$. For a fixed $\tilde{x} \in \mathbb{R}^n$, and parameters $\eta > 0, c \geq 4e^2$, let $P_{out}$ be a distribution supported on the set*

$$S_{\tilde{x}, out} := \{x \in \mathbb{R}^n : \|x - \tilde{x}\| \geq c(\eta + \sigma)\}.$$

*Let $P_{\tilde{x}}$ be a distribution which is supported within an $\eta-$radius ball centered at $\tilde{x}$.*

*For a fixed $A$, let $H_{\tilde{x}}$ denote the distribution of $y$ when $x^* \sim P_{\tilde{x}}$. Let $H_{out}$ denote the corresponding distribution of $y$ when $x^* \sim P_{out}$. Then we have:*

$$\mathbb{E}_A\left[TV(H_{\tilde{x}}, H_{out})\right] \geq 1 - 4e^{-\frac{m}{2}\log\left(\frac{c}{4e^2}\right)}.$$

By Markov's inequality, the expectation bound also gives a high probability bound over $A$.

For our current example, the above result implies that with probability $1 - e^{-\Omega(m)}$ over $A$, we have

$$TV(H_0, H_1) \geq 1 - e^{-\Omega(m)}. \quad (3)$$

Substituting equation (3) in equation (2), we have

$$\Pr\left[\|x^* - \hat{x}\| > 2\eta\right] \leq 2e^{-\Omega(m)}.$$

This shows that posterior sampling will produce a reconstruction which is close to the ground truth with overwhelmingly high probability for the two-ball example.

### 3.2. Going beyond two balls

The two-ball example leaves three main questions unanswered:

1. How do we handle distributions over larger collections of balls?

2. How do we handle mismatch between the distribution of reality ($R$) and the model ($P$)?

3. How do we handle having a $\delta$ probability of lying outside any ball?

**Unions of many balls.** The first question is relatively easy to answer: if $\text{Cov}_{\eta, 0}(R) \leq e^{o(m)}$, you can cover $R$ with a small number of balls, and essentially apply Lemma 3.2 with a union bound. There are a few details (e.g., Lemma 3.2

shows you will not confuse any ball with faraway balls, but you might confuse it with nearby balls) but solving them is straightforward. This shows that, if $P = R$ and $\log \mathrm{Cov}_{\eta,0}(R)$ is bounded, then posterior sampling works well with $1 - e^{-\Omega(m)}$ probability.

**Distribution mismatch in $\mathcal{W}_\infty$.** The above assumes we resample with respect to the true distribution $R$. But we only have a learned estimate $P$ of $R$. We would like to show that observing samples from $R$ and resampling according to $P$ gives good results. We first show that resampling signals drawn from $R$ with respect to $P$ is not much worse than resampling signals drawn from $P$ with respect to $P$, if $P$ and $R$ are close in $\mathcal{W}_\infty$.

**Lemma 3.3.** *Let $R, P$, denote arbitrary distributions over $\mathbb{R}^n$ such that $\mathcal{W}_\infty(R, P) \le \varepsilon$.*

*Let $x^* \sim R$ and $z^* \sim P$ and let $y$ and $u$ be generated from $x^*$ and $z^*$ via a Gaussian measurement process with $m$ measurements and noise level $\sigma$. Let $\widehat{x} \sim P(\cdot|y, A)$ and $\widehat{z} \sim P(\cdot|u, A)$. For any $d > 0$, we have*

$$\Pr_{x^*, A, \xi, \widehat{x}} [\|x^* - \widehat{x}\| \ge d + \varepsilon] \le$$
$$e^{-\Omega(m)} + e^{\left(\frac{4\varepsilon(\varepsilon + 2\sigma)m}{2\sigma^2}\right)} \Pr_{z^*, A, \xi, \widehat{z}} [\|z^* - \widehat{z}\| \ge d].$$

The idea is that with $\sigma$ Gaussian noise, measurements of a signal from $R$ aren't too different in distribution from measurements of the corresponding nearby signal from $P$.

Now, if $\mathcal{W}_\infty(R, P) \ll \sigma$, we would be nearly done: Lemma 3.3 says the situation is within $e^{o(m)}$ of the $R = P$ case, which we already know gives accurate recovery with $O(\log \mathrm{Cov}_{\eta,0}(P))$ measurements.

**Residual mass.** There are just two main issues remaining: we want to depend on $\log \mathrm{Cov}_{\eta,\delta}$ rather than $\log \mathrm{Cov}_{\eta,0}$, and we only want to require a bound on $\mathcal{W}_1(R, P)$ not $\mathcal{W}_\infty(R, P)$. By Markov's inequality, these issues are very similar: we want to allow both $R$ and $P$ to have a small constant probability of behaving badly. To address this, we note the existence of two distributions $R'$ and $P'$, which are only $\delta$-far in TV from $R$ and $P$ respectively, such that $R'$ and $P'$ do have a small cover & are close in $\mathcal{W}_\infty$. We show that, because posterior sampling would work with $R'$ and $P'$, it also works with $R$ and $P$. This leads to our full upper bound:

**Theorem 3.4.** *Let $\delta \in [0, 1/4)$, $p \ge 1$, and $\varepsilon, \eta > 0$ be parameters. Let $R, P$ be arbitrary distributions over $\mathbb{R}^n$ satisfying $\mathcal{W}_p(R, P) \le \varepsilon$.*

*Let $x^* \sim R$ and suppose $y$ is generated by a Gaussian measurement process from $x^*$ with noise level $\sigma \gtrsim \varepsilon/\delta^{1/p}$ and $m \ge O(\min(\log \mathrm{Cov}_{\eta,\delta}(R), \log \mathrm{Cov}_{\eta,\delta}(P)))$ mea-*

*surements. Given $y$ and the fixed matrix $A$, let $\widehat{x}$ output of posterior sampling with respect to $P$.*

*Then there exists a universal constant $c > 0$ such that with probability at least $1 - e^{-\Omega(m)}$ over $A, \xi$,*

$$\Pr_{x^* \sim R, \widehat{x} \sim P(\cdot|y)} [\|x^* - \widehat{x}\| \ge c\eta + c\sigma] \le 2\delta + 2e^{-\Omega(m)}.$$

Note that we can get a high-probability result by setting $p = \infty$: if $m \ge O(\log \mathrm{Cov}_{\eta,0}(R))$ and $\mathcal{W}_\infty(R, P) \le \sigma$, the error is $O(\sigma + \eta)$ with $1 - e^{-\Omega(m)}$ probability.

# 4. Lower Bound

In the previous section, we showed, for any distribution $R$ of signals, that $O(\log \mathrm{Cov}(R))$ measurements suffice for posterior sampling to recover most signals well. Now we show the converse: for any distribution of signals $R$, any algorithm for recovery must use $\Omega(\log \mathrm{Cov}(R))$ measurements.

**Theorem 4.1.** *Let $R$ be a distribution supported on a ball of radius $r$ in $\mathbb{R}^n$, and $x^* \sim R$. Let $y = Ax^* + \xi$, where $A$ is any matrix, and $\xi \sim \mathcal{N}(0, \frac{\sigma^2}{m} I_m)$. Assuming $\delta < 0.1$, if there exists a recovery scheme that uses $y$ and $A$ as inputs and guarantees*

$$\|\widehat{x} - x^*\| \le O(\eta),$$

*with probability $\ge 1 - \delta$, then we have*

$$m \ge \frac{0.15}{\log\left(1 + \frac{mr^2\|A\|_\infty^2}{\sigma^2}\right)} \left(\log \mathrm{Cov}_{3\eta, 4\delta}(R) + \log 6\delta - O(1)\right).$$

*If $A$ is an i.i.d. Gaussian matrix where each element is drawn from $\mathcal{N}(0, 1/m)$, then the above bound can be improved to:*

$$m \ge \frac{0.15}{\log\left(1 + \frac{r^2}{\sigma^2}\right)} \left(\log \mathrm{Cov}_{3\eta, 4\delta}(R) + \log 6\delta - O(1)\right).$$

This Theorem is proven using information theory, as an almost direct consequence of the following three Lemmas.

First, the measurement process reveals a limited amount of information:

**Lemma 4.2.** *Consider the setting of Theorem (4.1). If $A$ is a deterministic matrix, we have*

$$I(y; x^*) \le \frac{m}{2} \log\left(1 + \frac{mr^2\|A\|_\infty^2}{\sigma^2}\right).$$

*If $A$ is a Gaussian matrix, then $I(y; x^*|A) \le \frac{m}{2} \log\left(1 + \frac{r^2}{\sigma^2}\right)$.*

Second, since $x^* \to y \to \widehat{x}$ is a Markov chain, we can directly apply the Data Processing Inequality (Cover & Thomas, 2012).

**Lemma 4.3.** *Consider the setting of Theorem (4.1). If $A$ is a deterministic matrix, we have $I(x^*; \widehat{x}) \leq I(y; x^*)$.*

*If $A$ is a random matrix, then $I(x^*; \widehat{x}) \leq I(y; x^*|A)$.*

Finally, successful recovery must yield a large amount of information:

**Lemma 4.4** (Fano variant). *Let $(x, \widehat{x})$ be jointly distributed over $\mathbb{R}^n \times \mathbb{R}^n$, where $x \sim R$ and $\widehat{x}$ satisfies*

$$\Pr[\|x - \widehat{x}\| \leq \eta] \geq 1 - \delta.$$

*Then for any $\tau \leq 1 - 3\delta, \delta < 1/3$, we have*

$$0.99\tau(1 - 2\delta) \log \mathrm{Cov}_{3\eta, \tau+3\delta}(R) \leq I(x; \widehat{x}) + 1.98.$$

In order to complete the proof of Theorem 4.1, we need an additional counting argument to remove the extra $\tau$ term that appears in the left hand side of Lemma 4.4.

The proofs can be found in Appendix B.

# 5. Experiments

In this section we discuss our algorithm for posterior sampling, discuss why existing algorithms can fail, and show our empirical evaluation of posterior sampling versus baselines.

## 5.1. Datasets and Models

We perform our experiments on the CelebA-HQ (Liu et al., 2018; Karras et al., 2017) and FlickrFaces-HQ (Karras et al., 2019) datasets. For the CelebA dataset, we run experiments using a Glow generative model (Kingma & Dhariwal, 2018). For the FlickrFaces-HQ dataset, we use the NCSNv2 model (Song & Ermon, 2020). Both models have output size $256 \times 256 \times 3$. Details about our experiments are in Appendix C.

## 5.2. Langevin Dynamics

**Glow trained on CelebA-HQ** We first consider the Glow generative model, whose distribution $P$ is induced by the random variable $G(z)$, where $G : \mathbb{R}^n \to \mathbb{R}^n$ is a fixed deterministic generative model, and $z \sim \mathcal{N}(0, I_n)$ . Sampling from $p(z|y)$ is easier than sampling from $p(x|y)$, since it is easier to compute and we observe that sampling mixes quicker. Note that sampling $\widehat{z} \sim p(z|y)$ and setting $\widehat{x} = G(\widehat{z})$ is equivalent to sampling $\widehat{x} \sim p(x|y)$.

In order to sample from $p(z|y)$, we use *Langevin dynamics*, which samples from a given distribution by moving a random initial sample along a vector field given by the

distribution. Langevin dynamics tells us that if we sample $z_0 \sim \mathcal{N}(0, 1)$, and run the following iterative procedure:

$$z_{t+1} \leftarrow z_t + \frac{\alpha_t}{2}\nabla_z \log p(z_t|y) + \sqrt{\alpha_t}\zeta_t, \quad \zeta_t \sim \mathcal{N}(0, I),$$

then $p(z|y)$ is the stationary distribution of $z_t$ as $t \to \infty$ and $\alpha_t \to 0$. Unfortunately, this algorithm is slow to mix, as observed in (Song & Ermon, 2019). We instead use an annealed version of the algorithm, where in step $t$ we pretend that $p(z \mid y)$ has noise scale $\sigma_t \geq \sigma$ instead of $\sigma$. This gives

$$\log p_t(z|y) = \left(-\frac{\|y - AG(z)\|^2}{2\sigma_t^2/m} - \frac{\|z\|^2}{2}\right) + \log c(y), \tag{4}$$

where $c(y)$ is a constant that depends only on $y$. Since we only care about the gradient of $\log p(z|y)$, we can ignore this constant $c(y)$. By taking a decreasing sequence of $\sigma_t$ that approach the true value of $\sigma$, we can anneal Langevin dynamics and sample from $p(z|y)$. Please refer to Appendix C for more details about how $\sigma_t$ varies.

**NCSNv2 trained on FFHQ** We also consider the NCSNv2 model, which takes as input the image $x$, and outputs $\nabla_x \log p(x)$. This model is designed such that sampling from its marginal involves running Langevin dynamics. Since we have access to $\nabla_x \log p(x)$, and if we know the functional form of $p(y|x)$, we can easily compute $\nabla_x \log p(x|y)$, and run Langevin dynamics via

$$x_{t+1} \leftarrow x_t + \frac{\alpha_t}{2}\nabla_x \log p(x_t|y) + \sqrt{\alpha_t}\zeta_t, \ \zeta_t \sim \mathcal{N}(0, I).$$

Notice that we can also run MAP using this model. This can be achieved by simply following the gradient, and not adding noise: $x_{t+1} \leftarrow x_t + \frac{\alpha_t}{2}\nabla_x \log p(x_t|y)$.

This model also requires annealing, and we follow the schedule prescribed by (Song & Ermon, 2020). Please see Appendix C for more details.

## 5.3. MAP and Modified-MAP

The most relevant baseline for our algorithm is MAP, which was shown to be state-of-the-art for compressed sensing using generative priors (Asim et al., 2019).

Given access to a generative model $G$ such that the image $x = G(z)$, and $q(z)$ is the prior of $z$, the MAP estimate is

$$\widehat{z} := \arg\min_z \frac{\|y - AG(z)\|^2}{2\sigma^2/m} - \log q(z), \tag{5}$$

and set the estimate to be $\widehat{x} = G(\widehat{z})$. Typically, $q(z)$ is a standard Gaussian for many generative models. If one has
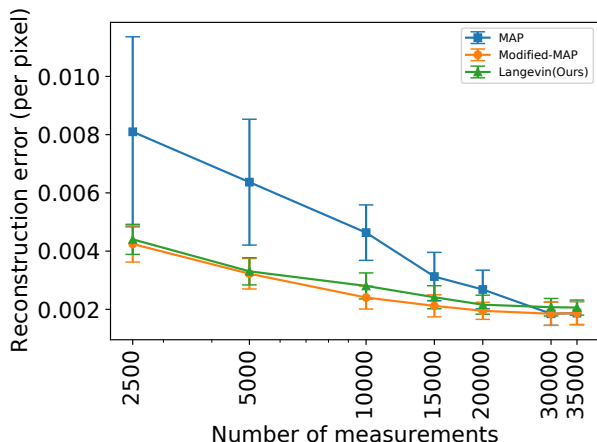
(a) $\|x^* - \widehat{x}\|^2 / n$

(b) Reconstructions for $m = 20,000$ measurements.

Figure 4: We compare our algorithm with the MAP baseline on the CelebA-HQ dataset, where the number of pixels is $n = 256 \times 256 \times 3 = 196,608$. In Figure (a) we show a plot of the per-pixel reconstruction error as we vary the number of measurements $m$. In Figure (b) we show reconstructions obtained by each algorithm for $m = 20,000$ measurements. We show original images (top row), reconstructions by MAP (second row), Modified-MAP (third row), and Langevin dynamics (ours, bottom row). Note that MAP produces several artefacts that are not seen in Modified-MAP or Langevin dynamics. In these experiments, modified-MAP picks hyperparameters based on the reconstruction error evaluated on some validation images, while MAP and Langevin dynamics pick hyperparameters that maximize the posterior likelihood. Here MAP, modified-MAP, and Langevin dynamics all use the same Glow model.

access to $p(x)$, such as in NCSNv2 (Song et al., 2019), it is possible to also do MAP in $x$-space.

One may modify this algorithm and introduce hyperparameters for better reconstructions. We call such algorithms *modified-MAP*. For example, (Asim et al., 2019) introduce a parameter $\gamma > 0$ that weights the prior, and their estimate is

$$\widehat{z}_{modified} := \arg\min_{z} \|y - AG(z)\|^2 - \gamma \log q(z), \quad (6)$$

Other examples of hyper-parameters include early stopping to avoid "over-fitting" to the measurements, and choosing optimization parameters such that the reconstruction error is minimized on a validation set of images. Then these hyper-parameters are used for evaluating reconstruction error on a different test.

### 5.4. Experimental Results

MAP estimation does not work on general distributions: as an extreme example, if $R$ is a mixture of some continuous distribution 99% of the time, and the all-zero image 1% of the time, it will always output the all-zero image, which is wrong 99% of the time. More generally, looking for high-likelihood *points* rather than *regions* means it prefers sharp but very narrow maxima to wide, but slightly shorter, maxima. Posterior sampling prefers the opposite. We now study this empirically.

**CelebA.** In Figure 4, we show the performance of our proposed algorithm for compressed sensing on CelebA-HQ with Glow. The baselines we consider are MAP, and modified-MAP. MAP directly optimizes the objective defined in Eqn (5) while Modified-MAP optimizes (6). The MAP baseline in Figure 4 tries to maximize the posterior likelihood, and hence hyperparameters are selected so that the posterior is optimized. In contrast, what we term the modified-MAP algorithm was proposed by (Asim et al., 2019), and this algorithm picks hyperparameters that minimize reconstruction error on a holdout set of images. These hyperparameters are significantly worse at optimizing the MAP objective, but lead to more accurate recovered images, presumably due to some sort of implicit regularization. This modified-MAP method has shown to be state-of-the-art for compressed sensing on CelebA (Asim et al., 2019).

We find that our algorithm is competitive with respect to modified-MAP, and beats MAP when the measurements are $< 35,000$.

**FFHQ.** In Figure 5, we show the performance of our proposed algorithm for compressed sensing on FlickrFaces-HQ with the NCSNv2 generative model. We consider MAP and Deep-Decoder (Heckel & Hand, 2018) as the baselines. Note that the NCSNv2 model was designed for Langevin dynamics, and we adapt it to MAP. Hence, we choose the Deep-Decoder as a second baseline, as it has been shown to
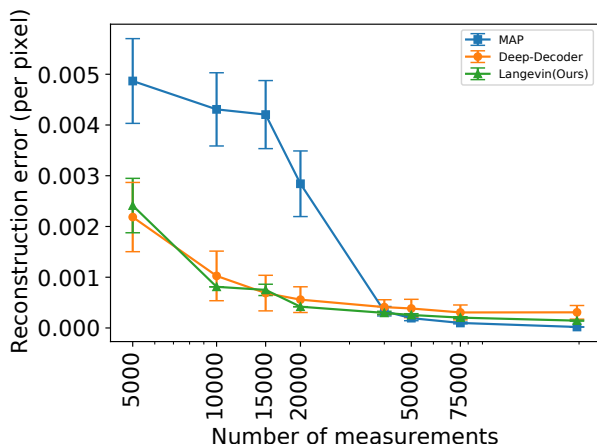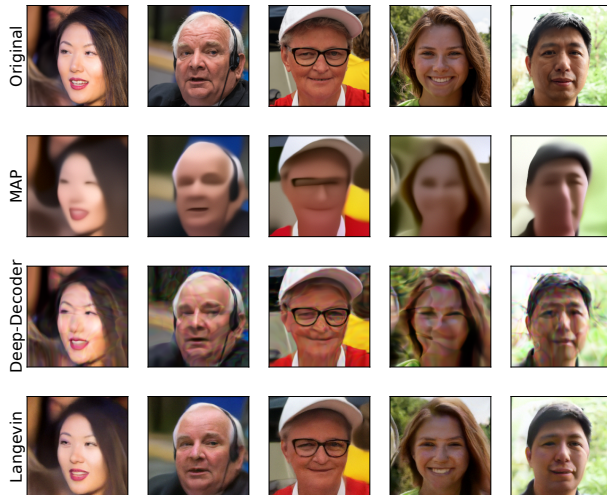
(a) $\|x^* - \widehat{x}\|^2/n$

(b) Reconstructions for $m = 5,000$ measurements.

Figure 5: We compare our algorithm with the MAP and Deep-Decoder baselines on the FFHQ dataset, where the number of pixels is $n = 256 \times 256 \times 3 = 196,608$. Figure (a) plots per-pixel reconstruction error as we vary the number of measurements $m$. Figure (b) shows original images (top row), reconstructions by MAP (second row), Deep-Decoder (third row), and Langevin dynamics (bottom row). Langevin dynamics is the practical implementation of our proposed posterior sampling estimator. Note that although Deep Decoder and Langevin achieve similar value of reconstruction errors, Langevin produces images with higher perceptual quality, as can be seen in Figure (b).

match state-of-the-art (Asim et al., 2019).

We observe that for $m < 40,000$ measurements, Langevin dynamics beats MAP, and is competitive with Deep-Decoder. In Figure 1 we visually compare the reconstruction quality as the number of measurements increases. Note that although Langevin and Deep-Decoder have similar reconstruction errors in Fig 5a, the images in Fig 1 produced by Langevin dynamics have better perceptual quality. Also see Fig 5b for more examples of reconstructions at $m = 5,000$ measurements.

**Inpainting.** In order to highlight the difference in diversity between images produced by MAP and Langevin dynamics, we evaluate them on the inverse problem of inpainting missing pixels. As shown in Figure 2, when the hair and background of a ground truth image is removed, MAP produces a single "most likely" reconstruction, while Langevin produces diverse images that satisfy the measurements. Each column for Langevin dynamics in Figure 2 corresponds to a run starting from a random initial point. We do not observe any change in MAP reconstructions as we vary the initial point.

We believe that the MAP reconstruction, while in some sense a highly likely reconstruction, is abnormally "washed out" and indistinct; analogous to how zero is the most likely sample from $N(0, I_d)$, yet is extremely atypical of the distribution. We see this quantitatively in that the corresponding $\|z\|^2/n$ for MAP is 0.007, even though samples from $R$

almost surely have $\|z\|^2/n \approx 1$, as do those of Langevin.

# 6. Conclusion

This paper studies the problem of compressed sensing a signal from a distribution $R$. We have shown that the measurement complexity is closely characterized by the log approximate covering number of $R$. Moreover, this recovery guarantee can be achieved by posterior sampling, even with respect to a distribution $P \neq R$ that is close in Wasserstein distance. Our experiments using Langevin dynamics to approximate posterior sampling match state-of-the-art recovery with a theoretically grounded algorithm.

This measurement complexity is inherent to the true distribution of images in the domain, and can't be improved. But perhaps it can be estimated: one open question is whether $\log \text{Cov}_{\eta,\delta}(P)$ can be estimated or bounded when $P$ is given by a neural network generative model.

# 7. Acknowledgements

# References

Aamand, A., Indyk, P., and Vakilian, A. (learned) frequency estimation algorithms under zipfian distribution. *arXiv preprint arXiv:1908.05198*, 2019.

Aeron, S., Saligrama, V., and Zhao, M. Information theoretic bounds for compressed sensing. *IEEE Transactions on Information Theory*, 56(10):5111–5130, 2010.

Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.

Asim, M., Shamshad, F., and Ahmed, A. Blind image deconvolution using deep generative priors. *arXiv preprint arXiv:1802.04073*, 2018.

Asim, M., Ahmed, A., and Hand, P. Invertible generative models for inverse problems: mitigating representation error and dataset bias. *arXiv preprint arXiv:1905.11672*, 2019.

Aubin, B., Loureiro, B., Baker, A., Krzakala, F., and Zdeborová, L. Exact asymptotics for phase retrieval and compressed sensing with random generative priors. *arXiv preprint arXiv:1912.02008*, 2019.

Baraniuk, R. G. and Wakin, M. B. Random projections of smooth manifolds. *Foundations of computational mathematics*, 9(1):51–77, 2009.

Baraniuk, R. G., Cevher, V., Duarte, M. F., and Hegde, C. Model-based compressive sensing. *IEEE Transactions on Information Theory*, 56(4):1982–2001, 2010.

Barbier, J., Krzakala, F., Macris, N., Miolane, L., and Zdeborová, L. Optimal errors and phase transitions in high-dimensional generalized linear models. *Proceedings of the National Academy of Sciences*, 116(12):5451–5460, 2019.

Bickel, P. J., Ritov, Y., and Tsybakov, A. B. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.

Bora, A., Jalal, A., Price, E., and Dimakis, A. G. Compressed sensing using generative models. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 537–546. JMLR. org, 2017.

Candes, E. J. The restricted isometry property and its implications for compressed sensing. *Comptes rendus mathematique*, 346(9-10):589–592, 2008.

Candes, E. J. and Davenport, M. A. How well can we estimate a sparse vector? *Applied and Computational Harmonic Analysis*, 34(2):317–323, 2013.

Candes, E. J., Romberg, J. K., and Tao, T. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 59(8):1207–1223, 2006.

Champion, T., De Pascale, L., and Juutinen, P. The $\infty$-Wasserstein distance: Local solutions and existence of optimal transport maps. *SIAM Journal on Mathematical Analysis*, 40(1):1–20, 2008.

Chen, C. and Huang, J. Compressive sensing mri with wavelet tree sparsity. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 25*, pp. 1115–1123. Curran Associates, Inc., 2012.

Chen, M., Silva, J., Paisley, J., Wang, C., Dunson, D., and Carin, L. Compressive sensing on manifolds using a nonparametric mixture of factor analyzers: Algorithm and performance bounds. *IEEE Transactions on Signal Processing*, 58(12):6140–6155, 2010.

Cover, T. M. and Thomas, J. A. *Elements of information theory*. John Wiley & Sons, 2012.

Dhar, M., Grover, A., and Ermon, S. Modeling sparse deviations for compressed sensing using generative models. *arXiv preprint arXiv:1807.01442*, 2018.

Donoho, D. L. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.

Duarte, M. F., Davenport, M. A., Takhar, D., Laska, J. N., Sun, T., Kelly, K. F., and Baraniuk, R. G. Single-pixel imaging via compressive sampling. *IEEE signal processing magazine*, 25(2):83–91, 2008.

Eldar, Y. C. and Mishali, M. Robust recovery of signals from a structured union of subspaces. *IEEE Transactions on Information Theory*, 55(11):5302–5316, 2009.

Fletcher, A. K., Pandit, P., Rangan, S., Sarkar, S., and Schniter, P. Plug-in estimation in high-dimensional linear inverse problems: A rigorous analysis. In *Advances in Neural Information Processing Systems*, pp. 7440–7449, 2018a.

Fletcher, A. K., Rangan, S., and Schniter, P. Inference in deep networks in high dimensions. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pp. 1884–1888. IEEE, 2018b.

Gómez, F. L., Eftekhari, A., and Cevher, V. Fast and provable admm for learning with generative priors. *arXiv preprint arXiv:1907.03343*, 2019.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.

Hand, P. and Joshi, B. Global guarantees for blind demodulation with generative priors. In *Advances in Neural Information Processing Systems*, pp. 11531–11541, 2019.

Hand, P. and Voroninski, V. Global guarantees for enforcing deep generative priors by empirical risk. *arXiv preprint arXiv:1705.07576*, 2017.

Hand, P., Leong, O., and Voroninski, V. Phase retrieval under a generative prior. In *Advances in Neural Information Processing Systems*, pp. 9136–9146, 2018.

Heckel, R. and Hand, P. Deep decoder: Concise image representations from untrained non-convolutional networks. *arXiv preprint arXiv:1810.03982*, 2018.

Heckel, R. and Soltanolkotabi, M. Compressive sensing with un-trained neural networks: Gradient descent finds the smoothest approximation. *arXiv preprint arXiv:2005.03991*, 2020.

Hegde, C. Algorithmic aspects of inverse problems using generative models. In *2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pp. 166–172. IEEE, 2018.

Hegde, C., Wakin, M., and Baraniuk, R. G. Random projections for manifold learning. In *Advances in neural information processing systems*, pp. 641–648, 2008.

Hegde, C., Duarte, M. F., and Cevher, V. Compressive sensing recovery of spike trains using a structured sparsity model. In *SPARS'09-Signal Processing with Adaptive Sparse Structured Representations*, 2009.

Hsu, C.-Y., Indyk, P., Katabi, D., and Vakilian, A. Learning-based frequency estimation algorithms. 2018.

Indyk, P., Vakilian, A., and Yuan, Y. Learning-based low-rank approximations. In *Advances in Neural Information Processing Systems*, pp. 7400–7410, 2019.

Iwen, M. and Tewfik, A. Adaptive group testing strategies for target detection and localization in noisy environments. 2010.

Jagatap, G. and Hegde, C. Phase retrieval using untrained neural network priors. 2019.

Jalal, A., Liu, L., Dimakis, A. G., and Caramanis, C. Robust compressed sensing using generative models. *Advances in Neural Information Processing Systems*, 33, 2020.

Jalali, S. and Yuan, X. Solving linear inverse problems using generative models. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pp. 512–516. IEEE, 2019.

Kabkab, M., Samangouei, P., and Chellappa, R. Task-aware compressed sensing with generative adversarial networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

Kamath, A., Karmalkar, S., and Price, E. Lower bounds for compressed sensing with generative models. *arXiv preprint arXiv:1912.02938*, 2019.

Karras, T., Aila, T., Laine, S., and Lehtinen, J. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.

Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4401–4410, 2019.

Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pp. 10215–10224, 2018.

Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Lei, Q., Jalal, A., Dhillon, I. S., and Dimakis, A. G. Inverting deep generative models, one layer at a time. In *Advances in Neural Information Processing Systems*, pp. 13910–13919, 2019.

Lindgren, E. M., Whang, J., and Dimakis, A. G. Conditional sampling from invertible generative models with applications to inverse problems. *arXiv preprint arXiv:2002.11743*, 2020.

Liu, Z. and Scarlett, J. Information-theoretic lower bounds for compressive sensing with generative models. *arXiv preprint arXiv:1908.10744*, 2019.

Liu, Z., Luo, P., Wang, X., and Tang, X. Large-scale celeb-faces attributes (celeba) dataset. *Retrieved August*, 15: 2018, 2018.

Liu, Z., Gomes, S., Tiwari, A., and Scarlett, J. Sample complexity bounds for 1-bit compressive sensing and binary stable embeddings with generative priors. *arXiv preprint arXiv:2002.01697*, 2020.

Lustig, M., Donoho, D., and Pauly, J. M. Sparse mri: The application of compressed sensing for rapid mr imaging. *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 58(6):1182–1195, 2007.

Lustig, M., Donoho, D. L., Santos, J. M., and Pauly, J. M. Compressed sensing mri. *IEEE signal processing magazine*, 25(2):72–82, 2008.

Mardani, M., Gong, E., Cheng, J. Y., Vasanawala, S. S., Zaharchuk, G., Xing, L., and Pauly, J. M. Deep generative adversarial neural networks for compressive sensing mri. *IEEE transactions on medical imaging*, 38(1):167–179, 2018.

Metzler, C., Mousavi, A., and Baraniuk, R. Learned d-amp: Principled neural network based compressive image recovery. In *Advances in Neural Information Processing Systems*, pp. 1772–1783, 2017.

Mosser, L., Dubrule, O., and Blunt, M. J. Stochastic seismic waveform inversion using generative adversarial networks as a geological prior. *Mathematical Geosciences*, 52(1): 53–79, 2020.

Ongie, G., Jalal, A., Metzler, C. A., Baraniuk, R. G., Dimakis, A. G., and Willett, R. Deep learning techniques for inverse problems in imaging. *arXiv preprint arXiv:2005.06001*, 2020.

Oord, A. v. d., Kalchbrenner, N., and Kavukcuoglu, K. Pixel recurrent neural networks. *arXiv preprint arXiv:1601.06759*, 2016.

Pandit, P., Sahraee-Ardakan, M., Rangan, S., Schniter, P., and Fletcher, A. K. Inference with deep generative priors in high dimensions. *arXiv preprint arXiv:1911.03409*, 2019.

Pesin, Y. B. *Dimension theory in dynamical systems: contemporary views and applications*. University of Chicago Press, 2008.

Polyanskiy, Y. and Wu, Y. Lecture notes on information theory. *Lecture Notes for ECE563 (UIUC) and*, 6(2012-2016):7, 2014.

Price, E. and Woodruff, D. P. (1+ eps)-approximate sparse recovery. In *2011 IEEE 52nd Annual Symposium on Foundations of Computer Science*, pp. 295–304. IEEE, 2011.

Qiu, S., Wei, X., and Yang, Z. Robust one-bit recovery via relu generative networks: Improved statistical rates and global landscape analysis. *arXiv preprint arXiv:1908.05368*, 2019.

Reeves, G. and Gastpar, M. The sampling rate-distortion tradeoff for sparsity pattern recovery in compressed sensing. *IEEE Transactions on Information Theory*, 58(5): 3065–3092, 2012.

Rick Chang, J., Li, C.-L., Poczos, B., Vijaya Kumar, B., and Sankaranarayanan, A. C. One network to solve them all–solving linear inverse problems using deep projection models. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5888–5897, 2017.

Scarlett, J. and Cevher, V. Limits on support recovery with probabilistic models: An information-theoretic framework. *IEEE Transactions on Information Theory*, 63(1): 593–620, 2016.

Shannon, C. E. A mathematical theory of communication. *Bell system technical journal*, 27(3):379–423, 1948.

Song, G., Fan, Z., and Lafferty, J. Surfing: Iterative optimization over incrementally trained deep networks. In *Advances in Neural Information Processing Systems*, pp. 15008–15017, 2019.

Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, pp. 11918–11930, 2019.

Song, Y. and Ermon, S. Improved techniques for training score-based generative models. *arXiv preprint arXiv:2006.09011*, 2020.

Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288, 1996.

Villani, C. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.

Whang, J., Lei, Q., and Dimakis, A. G. Compressed sensing with invertible generative models and dependent noise. *arXiv preprint arXiv:2003.08089*, 2020.

Wu, Y. *Shannon theory for compressed sensing*. Citeseer, 2011.

Wu, Y. and Verdú, S. Optimal phase transitions in compressed sensing. *IEEE Transactions on Information Theory*, 58(10):6241–6263, 2012.

Xu, W. and Hassibi, B. Compressed sensing over the grassmann manifold: A unified analytical framework. In *2008 46th Annual Allerton Conference on Communication, Control, and Computing*, pp. 562–567. IEEE, 2008.

Xu, W., Mallada, E., and Tang, A. Compressive sensing over graphs. In *2011 Proceedings IEEE INFOCOM*, pp. 2087–2095. IEEE, 2011.

Zdeborová, L. and Krzakala, F. Statistical physics of inference: Thresholds and algorithms. *Advances in Physics*, 65(5):453–552, 2016.

Zhou, Z., Liu, K., and Fang, J. Bayesian compressive sensing using normal product priors. *IEEE Signal Processing Letters*, 22(5):583–587, 2014.