

## Supplementary Materials

### A Posterior samples facilitate calibrated/fair detection

Say that  $s \in \{1, 0\}$  denotes the presence or absence of a particular pathology (e.g., brain tumor), and say that we train a soft classifier  $c(\cdot)$  on ground-truth images  $\mathbf{x}$  and calibrate it such that

$$c(\mathbf{x}) = \Pr\{s = 1|\mathbf{x}\}. \quad (\text{A.1})$$

Now say that we observe a distorted/corrupted/incomplete measurement  $\mathbf{y} = \mathcal{M}(\mathbf{x})$ . We would like to infer the probability that the pathology is present given  $\mathbf{y}$ , i.e., compute  $\Pr\{s = 1|\mathbf{y}\}$ . Note that

$$\Pr\{s = 1|\mathbf{y}\} = \int \Pr\{s = 1, \mathbf{x}|\mathbf{y}\} d\mathbf{x} = \int \Pr\{s = 1|\mathbf{x}, \mathbf{y}\} p(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (\text{A.2})$$

$$= \int \Pr\{s = 1|\mathbf{x}\} p(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \int c(\mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \mathbb{E}\{c(\mathbf{x})|\mathbf{y}\} \quad (\text{A.3})$$

$$= \lim_{P \rightarrow \infty} \frac{1}{P} \sum_{i=1}^P c(\hat{\mathbf{x}}_i) \text{ for } \hat{\mathbf{x}}_i \sim \text{i.i.d. } p(\mathbf{x}|\mathbf{y}), \quad (\text{A.4})$$

where the last equality follows from the law of large numbers. So, given access to many independent posterior samples  $\{\hat{\mathbf{x}}_i\}$ , equation (A.4) says that we can simply plug them into our calibrated classifier  $c(\cdot)$  and average the result to compute  $\Pr\{s = 1|\mathbf{y}\}$ . Conversely, if we have access to only the posterior mean  $\hat{\mathbf{x}}_{\text{mmse}} = \mathbb{E}\{\mathbf{x}|\mathbf{y}\}$ , then because

$$\Pr\{s = 1|\mathbf{y}\} = \mathbb{E}\{c(\mathbf{x})|\mathbf{y}\} \neq c(\mathbb{E}\{\mathbf{x}|\mathbf{y}\}) = c(\hat{\mathbf{x}}_{\text{mmse}}) \quad (\text{A.5})$$

for any non-linear  $c(\cdot)$ , the plug-in probability estimate will be incorrect. In fact, there exists no point estimate  $\hat{\mathbf{x}}$  that gives the correct  $\Pr\{s = 1|\mathbf{y}\}$  for general  $c(\cdot)$ .

Although above we defined  $c(\cdot)$  as a (soft) binary pathology classifier, the same results hold if we define  $c(\cdot)$  as a K-ary classifier of any protected attribute, such as race, gender, etc. This implies that, if we have a machine-learning system that has been calibrated to classify fairly on clean ground-truth data  $\mathbf{x}$ , then the use of posterior samples  $\{\hat{\mathbf{x}}_i\}$  enables it to classify fairly on distorted/corrupted/incomplete measurements  $\mathbf{y} = \mathcal{M}(\mathbf{x})$ , whereas the use of generic point-estimates  $\hat{\mathbf{x}}$  does not.

### B Proof of Proposition 3.1

Here we prove Proposition 3.1. To begin, for an  $N$ -pixel image, we rewrite (8)-(9) as

$$\mathcal{L}_{1,P}(\boldsymbol{\theta}) = \sum_{j=1}^N \mathbb{E}_{\mathbf{y}} \left\{ \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P | \mathbf{y}} \left\{ |x_j - \frac{1}{P} \sum_{i=1}^P \hat{x}_{ij}| \mid \mathbf{y} \right\} \right\} \quad (\text{B.1})$$

$$\mathcal{L}_{\text{SD},P}(\boldsymbol{\theta}) = \sum_{j=1}^N \mathbb{E}_{\mathbf{y}} \left\{ \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_P | \mathbf{y}} \left\{ \frac{\gamma_P}{P} \sum_{i=1}^P |\hat{x}_{ij} - \frac{1}{P} \sum_{k=1}^P \hat{x}_{kj}| \mid \mathbf{y} \right\} \right\}, \quad (\text{B.2})$$

where  $x_j \triangleq [\mathbf{x}]_j$ ,  $\hat{x}_{ij} \triangleq [\hat{\mathbf{x}}_i]_j$ , and

$$\gamma_P \triangleq \sqrt{\frac{\pi P}{2(P-1)}}. \quad (\text{B.3})$$

To simplify the notation in the sequel, we will consider an arbitrary fixed value of  $j$  and use the abbreviations

$$x_j \rightarrow X, \quad \hat{x}_{ij} \rightarrow \hat{X}_i. \quad (\text{B.4})$$

Recall that  $\mathbf{x}$  and  $\{\hat{\mathbf{x}}_i\}$  are mutually independent when conditioned on  $\mathbf{y}$  because the code vectors  $\{\mathbf{z}_i\}$  are generated independently of both  $\mathbf{x}$  and  $\mathbf{y}$ . In the context of Proposition 3.1, we also assume that the vector elements  $x_j$  and  $\hat{x}_{ij}$  are independent Gaussian when conditioned on  $\mathbf{y}$ . This implies that we can make the notational shift

$$p_{\mathbf{x}|\mathbf{y}}(x_j|\mathbf{y}) \rightarrow \mathcal{N}(X; \mu_0, \sigma_0^2), \quad p_{\hat{\mathbf{x}}|\mathbf{y}}(\hat{x}_{ij}|\mathbf{y}) \rightarrow \mathcal{N}(\hat{X}_i; \mu, \sigma^2), \quad (\text{B.5})$$

where  $X$  and  $\{\hat{X}_i\}$  are mutually independent. With this simplified notation, we note that  $[\hat{\mathbf{x}}_{\text{mmse}}]_j \rightarrow \mu_0$ , and that mode collapse corresponds to  $\sigma = 0$ .

Furthermore, if  $\theta$  can completely control  $(\mu, \sigma)$ , then (12) can be rewritten as

$$(\mu_*, \sigma_*) = \arg \min_{\mu, \sigma} \{ \mathcal{L}_{1,P}(\mu, \sigma) - \beta_{\text{SD}} \mathcal{L}_{\text{SD},P}(\mu, \sigma) \} \Rightarrow \begin{cases} \mu_* = \mu_0 \\ \sigma_* = \sigma_0 \end{cases} \quad (\text{B.6})$$

with

$$\mathcal{L}_{1,P}(\mu, \sigma) = \mathbb{E}_{X, \hat{X}_1, \dots, \hat{X}_P} \{ |X - \frac{1}{P} \sum_{i=1}^P \hat{X}_i| \} \quad (\text{B.7})$$

$$\mathcal{L}_{\text{SD},P}(\mu, \sigma) = \mathbb{E}_{\hat{X}_1, \dots, \hat{X}_P} \{ \frac{\gamma_P}{P} \sum_{i=1}^P |\hat{X}_i - \frac{1}{P} \sum_{k=1}^P \hat{X}_k| \}. \quad (\text{B.8})$$

Although  $\sigma_*$  must be positive, it turns out that we do not need to enforce this in the optimization (B.6) because it will arise naturally.

To further analyze (B.7) and (B.8), we define

$$\hat{\mu} \triangleq \frac{1}{P} \sum_{i=1}^P \hat{X}_i \quad (\text{B.9})$$

$$\hat{\sigma} \triangleq \frac{\gamma_P}{P} \sum_{i=1}^P |\hat{X}_i - \hat{\mu}|. \quad (\text{B.10})$$

The quantity  $\hat{\mu}$  can be recognized as the unbiased estimate of the mean  $\mu$  of  $\hat{X}_i$ , and we now show that  $\hat{\sigma}$  is an unbiased estimate of the SD  $\sigma$  of  $\hat{X}_i$  in the case that  $\hat{X}_i$  is Gaussian. To see this, first observe that the i.i.d.  $\mathcal{N}(\mu, \sigma^2)$  property of  $\{\hat{X}_i\}$  implies that  $\hat{X}_i - \hat{\mu} = (1 - \frac{1}{P})\hat{X}_i - \frac{1}{P} \sum_{k \neq i} \hat{X}_k$  is Gaussian with mean zero and variance  $(1 - \frac{1}{P})^2 \sigma^2 + \frac{P-1}{P^2} \sigma^2 = \frac{P-1}{P} \sigma^2$ . The variable  $|\hat{X}_i - \hat{\mu}|$  is thus half-normal distributed with mean  $\sqrt{\frac{2(P-1)}{\pi P}} \sigma^2$  [48]. Because  $\{\hat{X}_i\}$  are i.i.d., the variable  $\frac{1}{P} \sum_{i=1}^P |\hat{X}_i - \hat{\mu}|$  has the same mean as  $|\hat{X}_i - \hat{\mu}|$ . Finally, multiplying  $\frac{1}{P} \sum_{i=1}^P |\hat{X}_i - \hat{\mu}|$  by  $\gamma_P$  yields  $\hat{\sigma}$  from (B.10), and multiplying its mean using the expression for  $\gamma_P$  from (B.3) implies

$$\mathbb{E}\{\hat{\sigma}\} = \sigma, \quad (\text{B.11})$$

and so  $\hat{\sigma}$  is an unbiased estimator of  $\sigma$ , the SD of  $\hat{X}_i$ .

With the above definitions of  $\hat{\mu}$  and  $\hat{\sigma}$ , the optimization cost in (B.6) can be written as

$$\mathcal{L}_{1,P}(\mu, \sigma) - \beta_{\text{SD}} \mathcal{L}_{\text{SD},P}(\mu, \sigma) = \mathbb{E}_{X, \hat{X}_1, \dots, \hat{X}_P} \{ |X - \hat{\mu}| \} - \beta_{\text{SD}} \mathbb{E}_{\hat{X}_1, \dots, \hat{X}_P} \{ \hat{\sigma} \} \quad (\text{B.12})$$

$$= \mathbb{E}_{X, \hat{X}_1, \dots, \hat{X}_P} \{ |X - \hat{\mu}| \} - \beta_{\text{SD}} \sigma, \quad (\text{B.13})$$

where in the last step we exploited the unbiased property of  $\hat{\sigma}$ . To proceed further, we note that the i.i.d. Gaussian property of  $\{\hat{X}_i\}$  implies  $\hat{\mu} \sim \mathcal{N}(\mu, \sigma^2/P)$ , after which the mutual independence of  $\{\hat{X}_i\}$  and  $X$  yields

$$X - \hat{\mu} \sim \mathcal{N}(\mu_0 - \mu, \sigma_0^2 + \sigma^2/P). \quad (\text{B.14})$$

Taking the absolute value of a Gaussian random yields a folded-normal random variable [48]. Using the mean and variance in (B.14), the expressions in [48] yield

$$\begin{aligned} \mathbb{E}_{X, \hat{X}_1, \dots, \hat{X}_P} \{ |X - \hat{\mu}| \} &= \sqrt{\frac{2(\sigma_0^2 + \sigma^2/P)}{\pi}} \exp\left(-\frac{(\mu_0 - \mu)^2}{2(\sigma_0^2 + \sigma^2/P)}\right) \\ &\quad + (\mu_0 - \mu) \operatorname{erf}\left(\frac{\mu_0 - \mu}{\sqrt{2(\sigma_0^2 + \sigma^2/P)}}\right). \end{aligned} \quad (\text{B.15})$$

Thus the optimization cost (B.13) can be written as

$$\begin{aligned} J(\mu, \sigma) &= \sqrt{\frac{2(\sigma_0^2 + \sigma^2/P)}{\pi}} \exp\left(-\frac{(\mu - \mu_0)^2}{2(\sigma_0^2 + \sigma^2/P)}\right) \\ &\quad + (\mu - \mu_0) \operatorname{erf}\left(\frac{\mu - \mu_0}{\sqrt{2(\sigma_0^2 + \sigma^2/P)}}\right) - \beta_{\text{SD}} \sigma. \end{aligned} \quad (\text{B.16})$$

Since  $J(\cdot, \cdot)$  is convex, the minimizer  $(\mu_*, \sigma_*) = \arg \min_{\mu, \sigma} J(\mu, \sigma)$  satisfies  $\nabla J(\mu_*, \sigma_*) = (0, 0)$ . To streamline the derivation, we define

$$c \triangleq \sqrt{2(\sigma_0^2 + \sigma^2/P)/\pi}, \quad s \triangleq \sqrt{\sigma_0^2 + \sigma^2/P} \quad (\text{B.17})$$

so that

$$J(\mu, \sigma) = c \exp\left(-\frac{(\mu - \mu_0)^2}{2s^2}\right) + (\mu - \mu_0) \operatorname{erf}\left(\frac{\mu - \mu_0}{\sqrt{2}s^2}\right) - \beta_{\text{SD}}\sigma. \quad (\text{B.18})$$

Because  $c$  and  $s$  are invariant to  $\mu$ , we get

$$\frac{\partial J(\mu, \sigma)}{\partial \mu} = -c \exp\left(-\frac{(\mu - \mu_0)^2}{2s^2}\right) \frac{\mu - \mu_0}{s^2} + \operatorname{erf}\left(\frac{\mu - \mu_0}{\sqrt{2}s^2}\right) + (\mu - \mu_0) \frac{2}{\sqrt{\pi}} \exp\left(-\frac{(\mu - \mu_0)^2}{2s^2}\right), \quad (\text{B.19})$$

which equals zero if and only if  $\mu = \mu_0$ . Thus we have determined that  $\mu_* = \mu_0$ . Plugging  $\mu_* = \mu_0$  into (B.16), we find

$$J(\mu_*, \sigma) = \sqrt{2(\sigma_0^2 + \sigma^2/P)/\pi} - \beta_{\text{SD}}\sigma. \quad (\text{B.20})$$

Taking the derivative with respect to  $\sigma$ , we get

$$\frac{\partial J(\mu_*, \sigma)}{\partial \sigma} = \sqrt{\frac{2}{\pi P(P\sigma_0^2/\sigma^2 + 1)}} - \beta_{\text{SD}} \quad (\text{B.21})$$

$$= \sqrt{\frac{2}{\pi P(P\sigma_0^2/\sigma^2 + 1)}} - \sqrt{\frac{2}{\pi P(P + 1)}}, \quad (\text{B.22})$$

where in the last step we applied the value of  $\beta_{\text{SD}}$  from (11). It can now be seen that  $\frac{\partial J(\mu_*, \sigma)}{\partial \sigma} = 0$  if and only if  $\sigma = \sigma_0$ , which implies that  $\sigma_* = \sigma_0$ . Thus we have established (B.6), which completes the proof of Proposition 3.1.

## C Derivation of Proposition 3.2

Here we prove Proposition 3.2. To start, we establish some notation and conditional-mean properties:

$$\begin{aligned} \hat{\mathbf{x}}_{\text{mmse}} &\triangleq \mathbb{E}_{\mathbf{x}|\mathbf{y}}\{\mathbf{x}|\mathbf{y}\} \\ \mathbf{e}_{\text{mmse}} &\triangleq \mathbf{x} - \hat{\mathbf{x}}_{\text{mmse}}, & \mathbf{0} &= \mathbb{E}_{\mathbf{x}|\mathbf{y}}\{\mathbf{e}_{\text{mmse}}|\mathbf{y}\} \\ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) &\triangleq G_{\boldsymbol{\theta}}(\mathbf{z}_i, \mathbf{y}), & \bar{\mathbf{x}}(\boldsymbol{\theta}) &\triangleq \mathbb{E}_{\mathbf{z}_i|\mathbf{y}}\{\hat{\mathbf{x}}_i(\boldsymbol{\theta})|\mathbf{y}\} \\ \hat{\mathbf{x}}_{(P)}(\boldsymbol{\theta}) &\triangleq \frac{1}{P} \sum_{i=1}^P \hat{\mathbf{x}}_i(\boldsymbol{\theta}), & \bar{\mathbf{x}}(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}}\{\hat{\mathbf{x}}_{(P)}(\boldsymbol{\theta})|\mathbf{y}\} \\ \mathbf{d}_i(\boldsymbol{\theta}) &\triangleq \hat{\mathbf{x}}_i(\boldsymbol{\theta}) - \bar{\mathbf{x}}(\boldsymbol{\theta}), & \mathbf{0} &= \mathbb{E}_{\mathbf{z}_i|\mathbf{y}}\{\mathbf{d}_i(\boldsymbol{\theta})|\mathbf{y}\} \forall \boldsymbol{\theta} \\ \mathbf{d}_{(P)}(\boldsymbol{\theta}) &\triangleq \frac{1}{P} \sum_{i=1}^P \mathbf{d}_i(\boldsymbol{\theta}), & \mathbf{0} &= \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}}\{\mathbf{d}_{(P)}(\boldsymbol{\theta})|\mathbf{y}\} \forall \boldsymbol{\theta} \end{aligned} \quad (\text{C.1})$$

Our first step is to write (14) as

$$\mathcal{L}_{2,P}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{y}} \left\{ \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ \|\mathbf{x} - \hat{\mathbf{x}}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \right\} \right\}. \quad (\text{C.2})$$

Leveraging the fact that  $\hat{\mathbf{x}}_{\text{mmse}}$  and  $\bar{\mathbf{x}}(\boldsymbol{\theta})$  are deterministic given  $\mathbf{y}$ , we write the inner term in (C.2) as

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ \|\mathbf{x} - \hat{\mathbf{x}}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \right\} \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ \|\hat{\mathbf{x}}_{\text{mmse}} + \mathbf{e}_{\text{mmse}} - \bar{\mathbf{x}}(\boldsymbol{\theta}) - \mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \right\} \end{aligned} \quad (\text{C.3})$$

$$\begin{aligned} &= \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ \|\hat{\mathbf{x}}_{\text{mmse}} - \bar{\mathbf{x}}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \right\} \\ &\quad + 2 \operatorname{Re} \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ (\hat{\mathbf{x}}_{\text{mmse}} - \bar{\mathbf{x}}(\boldsymbol{\theta}))^H (\mathbf{e}_{\text{mmse}} - \mathbf{d}_{(P)}(\boldsymbol{\theta})) | \mathbf{y} \right\} \\ &\quad + \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ \|\mathbf{e}_{\text{mmse}} - \mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \right\} \end{aligned} \quad (\text{C.4})$$

$$= \|\hat{\mathbf{x}}_{\text{mmse}} - \bar{\mathbf{x}}(\boldsymbol{\theta})\|_2^2 + 2 \operatorname{Re} \left[ (\hat{\mathbf{x}}_{\text{mmse}} - \bar{\mathbf{x}}(\boldsymbol{\theta}))^H \underbrace{\mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \{ (\mathbf{e}_{\text{mmse}} - \mathbf{d}_{(P)}(\boldsymbol{\theta})) | \mathbf{y} \}}_{= \mathbf{0}} \right]$$

$$+ \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ \|\mathbf{e}_{\text{mmse}} - \mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \right\} \quad (\text{C.5})$$

$$= \|\hat{\mathbf{x}}_{\text{mmse}} - \mathbb{E}_{\mathbf{z}_i|\mathbf{y}} \{ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) | \mathbf{y} \} \|_2^2 + \mathbb{E}_{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_P|\mathbf{y}} \left\{ \|\mathbf{e}_{\text{mmse}} - \mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \right\}. \quad (\text{C.6})$$

where in (C.5) we used the fact that  $\mathbf{d}_{(P)}$  and  $\mathbf{e}_{\text{mmse}}$  are both zero-mean when conditioned on  $\mathbf{y}$ . We now leverage the fact that  $\{z_i\}$  are independent of  $\mathbf{x}$  and  $\mathbf{y}$  to write

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}, z_1, \dots, z_P | \mathbf{y}} \{ \|\mathbf{e}_{\text{mmse}} - \mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \\ &= \mathbb{E}_{\mathbf{x}, z_1, \dots, z_P | \mathbf{y}} \{ \|\mathbf{e}_{\text{mmse}}\|_2^2 | \mathbf{y} \} + 2 \operatorname{Re} \mathbb{E}_{\mathbf{x}, z_1, \dots, z_P | \mathbf{y}} \{ \mathbf{e}_{\text{mmse}}^H \mathbf{d}_{(P)}(\boldsymbol{\theta}) | \mathbf{y} \} + \mathbb{E}_{\mathbf{x}, z_1, \dots, z_P | \mathbf{y}} \{ \|\mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \\ &= \mathbb{E}_{\mathbf{x} | \mathbf{y}} \{ \|\mathbf{e}_{\text{mmse}}\|_2^2 | \mathbf{y} \} + 2 \operatorname{Re} \underbrace{\mathbb{E}_{\mathbf{x} | \mathbf{y}} \{ \mathbf{e}_{\text{mmse}} | \mathbf{y} \}}_{= \mathbf{0}} \underbrace{\mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \mathbf{d}_{(P)}(\boldsymbol{\theta}) | \mathbf{y} \}}_{= \mathbf{0}} + \mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|\mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \}. \end{aligned} \quad (\text{C.7})$$

$$(\text{C.8})$$

Finally, we can leverage the fact that  $\{z_i\}$  are i.i.d. to write

$$\mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|\mathbf{d}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} = \mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|\frac{1}{P} \sum_{i=1}^P \mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \quad (\text{C.9})$$

$$= \frac{1}{P^2} \sum_{i=1}^P \mathbb{E}_{z_i | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \quad (\text{C.10})$$

$$= \frac{1}{P} \mathbb{E}_{z_i | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \text{ for any } i \quad (\text{C.11})$$

$$= \frac{1}{P} \mathbb{E}_{z_i | \mathbf{y}} \{ \operatorname{tr}[\mathbf{d}_i(\boldsymbol{\theta}) \mathbf{d}_i(\boldsymbol{\theta})^H] | \mathbf{y} \} \quad (\text{C.12})$$

$$= \frac{1}{P} \operatorname{tr} [ \mathbb{E}_{z_i | \mathbf{y}} \{ \mathbf{d}_i(\boldsymbol{\theta}) \mathbf{d}_i(\boldsymbol{\theta})^H | \mathbf{y} \} ] \quad (\text{C.13})$$

$$= \frac{1}{P} \operatorname{tr} [ \operatorname{Cov}_{z_i | \mathbf{y}} \{ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) | \mathbf{y} \} ]. \quad (\text{C.14})$$

Combining (C.2), (C.6), (C.8), and (C.14), we get the bias-variance decomposition

$$\mathcal{L}_{2,P}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{y}} \left\{ \|\hat{\mathbf{x}}_{\text{mmse}} - \mathbb{E}_{z_i | \mathbf{y}} \{ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) | \mathbf{y} \} \|_2^2 + \frac{1}{P} \operatorname{tr} [ \operatorname{Cov}_{z_i | \mathbf{y}} \{ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) | \mathbf{y} \} ] + \mathbb{E}_{\mathbf{x} | \mathbf{y}} \{ \|\mathbf{e}_{\text{mmse}}\|_2^2 | \mathbf{y} \} \right\}. \quad (\text{C.15})$$

We now see that if  $\boldsymbol{\theta}$  has complete control over the  $\mathbf{y}$ -conditional mean and covariance of  $\hat{\mathbf{x}}_i(\boldsymbol{\theta})$ , then minimizing (C.15) over  $\boldsymbol{\theta}$  will cause

$$\mathbb{E}_{z_i | \mathbf{y}} \{ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) | \mathbf{y} \} = \hat{\mathbf{x}}_{\text{mmse}} \quad (\text{C.16})$$

$$\operatorname{Cov}_{z_i | \mathbf{y}} \{ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) | \mathbf{y} \} = \mathbf{0}, \quad (\text{C.17})$$

which proves Proposition 3.2.

## D Derivation of (19)

To show that the expression for  $\mathcal{L}_{\text{var},P}$  in (19) holds, we first rewrite (18) as

$$\mathcal{L}_{\text{var},P}(\boldsymbol{\theta}) = \frac{1}{P-1} \sum_{i=1}^P \mathbb{E}_{\mathbf{y}} \{ \mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|\hat{\mathbf{x}}_i(\boldsymbol{\theta}) - \hat{\mathbf{x}}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \} \quad (\text{D.1})$$

where the definitions from (C.1) imply

$$\begin{aligned} & \mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|\hat{\mathbf{x}}_i(\boldsymbol{\theta}) - \hat{\mathbf{x}}_{(P)}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \\ &= \mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|\bar{\mathbf{x}}(\boldsymbol{\theta}) + \mathbf{d}_i(\boldsymbol{\theta}) - \mathbf{d}_{(P)}(\boldsymbol{\theta}) - \bar{\mathbf{x}}(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \end{aligned} \quad (\text{D.2})$$

$$= \mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta}) - \frac{1}{P} \sum_{j=1}^P \mathbf{d}_j(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \quad (\text{D.3})$$

$$= \mathbb{E}_{z_1, \dots, z_P | \mathbf{y}} \{ \|(1 - \frac{1}{P})\mathbf{d}_i(\boldsymbol{\theta}) - \frac{1}{P} \sum_{j \neq i} \mathbf{d}_j(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \quad (\text{D.4})$$

$$= (1 - \frac{1}{P})^2 \mathbb{E}_{z_i | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} + \frac{P-1}{P^2} \mathbb{E}_{z_i | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \quad (\text{D.5})$$

$$= \frac{P-1}{P} \mathbb{E}_{z_i | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \text{ for any } i. \quad (\text{D.6})$$

For (D.5), we leveraged the zero-mean and i.i.d. nature of  $\{\mathbf{d}_i(\boldsymbol{\theta})\}$  conditioned on  $\mathbf{y}$ . By plugging (D.6) into (D.1), we get

$$\mathcal{L}_{\text{var},P}(\boldsymbol{\theta}) = \frac{1}{P} \sum_{i=1}^P \mathbb{E}_{\mathbf{y}} \{ \mathbb{E}_{z_i | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \} \quad (\text{D.7})$$

$$= \mathbb{E}_{\mathbf{y}} \{ \mathbb{E}_{z_i | \mathbf{y}} \{ \|\mathbf{d}_i(\boldsymbol{\theta})\|_2^2 | \mathbf{y} \} \} \text{ for any } i \quad (\text{D.8})$$

$$= \mathbb{E}_{\mathbf{y}} \{ \operatorname{tr} [ \operatorname{Cov}_{z_i | \mathbf{y}} \{ \hat{\mathbf{x}}_i(\boldsymbol{\theta}) | \mathbf{y} \} ] \}, \quad (\text{D.9})$$

where (D.8) follows because  $\{\mathbf{d}_i(\boldsymbol{\theta})\}$  are i.i.d. when conditioned on  $\mathbf{y}$  and (D.9) follows from manipulations similar to those used for (C.14).



### E Proof of Proposition 3.3

Here we prove Proposition 3.3. Recall from (C.1) that  $\hat{\mathbf{x}}_{\text{mmse}} \triangleq \mathbb{E}\{\mathbf{x}|\mathbf{y}\}$  and  $\mathbf{e}_{\text{mmse}} \triangleq \mathbf{x} - \hat{\mathbf{x}}_{\text{mmse}}$ . To reduce clutter, we will abbreviate  $\mathbf{e}_{\text{mmse}}$  by  $\mathbf{e}$  in this appendix. Also, for true-posterior samples  $\hat{\mathbf{x}}_i \sim p_{\mathbf{x}|\mathbf{y}}(\cdot|\mathbf{y})$ , we define

$$\hat{\mathbf{e}}_i \triangleq \hat{\mathbf{x}}_i - \hat{\mathbf{x}}_{\text{mmse}}. \quad (\text{E.1})$$

Then using  $\hat{\mathbf{x}}_{(P)} \triangleq \frac{1}{P} \sum_{i=1}^P \hat{\mathbf{x}}_i$  and from  $\mathcal{E}_P$  from (20), we have

$$\mathcal{E}_P = \mathbb{E}\{\|\hat{\mathbf{x}}_{(P)} - \mathbf{x}\|^2|\mathbf{y}\} \quad (\text{E.2})$$

$$= \mathbb{E}\{\|(\frac{1}{P} \sum_{i=1}^P \hat{\mathbf{x}}_i) - \mathbf{x}\|^2|\mathbf{y}\} \quad (\text{E.3})$$

$$= \mathbb{E}\{\|\frac{1}{P} \sum_{i=1}^P (\hat{\mathbf{x}}_i - \mathbf{x})\|^2|\mathbf{y}\} \quad (\text{E.4})$$

$$= \frac{1}{P^2} \mathbb{E}\{\|\sum_{i=1}^P (\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_{\text{mmse}} + \hat{\mathbf{x}}_{\text{mmse}} - \mathbf{x})\|^2|\mathbf{y}\} \quad (\text{E.5})$$

$$= \frac{1}{P^2} \mathbb{E}\{\|\sum_{i=1}^P (\hat{\mathbf{e}}_i - \mathbf{e})\|^2|\mathbf{y}\} \quad (\text{E.6})$$

$$= \frac{1}{P^2} \mathbb{E}\{\sum_{i=1}^P (\hat{\mathbf{e}}_i - \mathbf{e})^H \sum_{j=1}^P (\hat{\mathbf{e}}_j - \mathbf{e})|\mathbf{y}\} \quad (\text{E.7})$$

$$= \frac{1}{P^2} \sum_{i=1}^P \sum_{j=1}^P \mathbb{E}\{(\hat{\mathbf{e}}_i - \mathbf{e})^H (\hat{\mathbf{e}}_j - \mathbf{e})|\mathbf{y}\} \quad (\text{E.8})$$

$$= \frac{1}{P^2} \sum_{i=1}^P \mathbb{E}\{(\hat{\mathbf{e}}_i - \mathbf{e})^H (\hat{\mathbf{e}}_i - \mathbf{e})|\mathbf{y}\} + \frac{1}{P^2} \sum_{i=1}^P \sum_{j \neq i} \mathbb{E}\{(\hat{\mathbf{e}}_i - \mathbf{e})^H (\hat{\mathbf{e}}_j - \mathbf{e})|\mathbf{y}\} \quad (\text{E.9})$$

$$= \frac{1}{P^2} \sum_{i=1}^P [\mathbb{E}\{\|\hat{\mathbf{e}}_i\|^2|\mathbf{y}\} - 2 \text{Re} \mathbb{E}\{\hat{\mathbf{e}}_i^H \mathbf{e}|\mathbf{y}\} + \mathbb{E}\{\|\mathbf{e}\|^2|\mathbf{y}\}] \\ + \frac{1}{P^2} \sum_{i=1}^P \sum_{j \neq i} \text{Re} [\mathbb{E}\{\hat{\mathbf{e}}_i^H \hat{\mathbf{e}}_j|\mathbf{y}\} - \mathbb{E}\{\hat{\mathbf{e}}_i^H \mathbf{e}|\mathbf{y}\} - \mathbb{E}\{\mathbf{e}^H \hat{\mathbf{e}}_j|\mathbf{y}\} + \mathbb{E}\{\|\mathbf{e}\|^2|\mathbf{y}\}] \quad (\text{E.10})$$

$$= \frac{1}{P^2} \sum_{i=1}^P \mathbb{E}\{\|\hat{\mathbf{e}}_i\|^2|\mathbf{y}\} + \frac{1}{P} \mathbb{E}\{\|\mathbf{e}\|^2|\mathbf{y}\} + \frac{P(P-1)}{P^2} \mathbb{E}\{\|\mathbf{e}\|^2|\mathbf{y}\}, \quad (\text{E.11})$$

where certain terms vanished because the i.i.d. and zero-mean properties of  $\{\mathbf{e}, \hat{\mathbf{e}}_1, \dots, \hat{\mathbf{e}}_P\}$  imply

$$\mathbb{E}\{\hat{\mathbf{e}}_i^H \hat{\mathbf{e}}_j|\mathbf{y}\} = \mathbb{E}\{\hat{\mathbf{e}}_i|\mathbf{y}\}^H \mathbb{E}\{\hat{\mathbf{e}}_j|\mathbf{y}\} = 0 \quad (\text{E.12})$$

$$\mathbb{E}\{\hat{\mathbf{e}}_i^H \mathbf{e}|\mathbf{y}\} = \mathbb{E}\{\hat{\mathbf{e}}_i|\mathbf{y}\}^H \mathbb{E}\{\mathbf{e}|\mathbf{y}\} = 0 \quad (\text{E.13})$$

$$\mathbb{E}\{\mathbf{e}^H \hat{\mathbf{e}}_j|\mathbf{y}\} = \mathbb{E}\{\mathbf{e}|\mathbf{y}\}^H \mathbb{E}\{\hat{\mathbf{e}}_j|\mathbf{y}\} = 0. \quad (\text{E.14})$$

Finally, note that  $\mathbb{E}\{\|\mathbf{e}\|^2|\mathbf{y}\} = \mathcal{E}_{\text{mmse}}$  from (C.1). Furthermore, because  $\{\mathbf{x}, \hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_P\}$  are independent samples of  $p_{\mathbf{x}|\mathbf{y}}(\cdot|\mathbf{y})$  under the assumptions of Proposition 3.3, we have  $\mathbb{E}\{\|\mathbf{e}\|^2|\mathbf{y}\} = \mathbb{E}\{\|\hat{\mathbf{e}}_i\|^2|\mathbf{y}\}$  and so (E.11) becomes

$$\mathcal{E}_P = \frac{1}{P^2} \sum_{i=1}^P \mathcal{E}_{\text{mmse}} + \frac{1}{P} \mathcal{E}_{\text{mmse}} + \frac{P(P-1)}{P^2} \mathcal{E}_{\text{mmse}} = \frac{P+1}{P} \mathcal{E}_{\text{mmse}}. \quad (\text{E.15})$$

This result holds for any  $P \geq 1$ , which implies the ratio

$$\frac{\mathcal{E}_1}{\mathcal{E}_P} = \frac{2P}{P+1}. \quad (\text{E.16})$$

### F CFID implementation details

With the Gaussian approximation described in Section 4.1, where  $p_{\mathbf{x}|\mathbf{y}}$  and  $p_{\hat{\mathbf{x}}|\mathbf{y}}$  are approximated by  $\mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}, \boldsymbol{\Sigma}_{\mathbf{xx}|\mathbf{y}})$  and  $\mathcal{N}(\boldsymbol{\mu}_{\hat{\mathbf{x}}|\mathbf{y}}, \boldsymbol{\Sigma}_{\hat{\mathbf{x}}\hat{\mathbf{x}}|\mathbf{y}})$ , respectively, the CWD in (24) reduces to

$$\text{CFID} \triangleq \mathbb{E}_{\mathbf{y}} \left\{ \|\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}} - \boldsymbol{\mu}_{\hat{\mathbf{x}}|\mathbf{y}}\|_2^2 + \text{tr} [\boldsymbol{\Sigma}_{\mathbf{xx}|\mathbf{y}} + \boldsymbol{\Sigma}_{\hat{\mathbf{x}}\hat{\mathbf{x}}|\mathbf{y}} - 2(\boldsymbol{\Sigma}_{\mathbf{xx}|\mathbf{y}}^{1/2} \boldsymbol{\Sigma}_{\hat{\mathbf{x}}\hat{\mathbf{x}}|\mathbf{y}} \boldsymbol{\Sigma}_{\mathbf{xx}|\mathbf{y}}^{1/2})^{1/2}] \right\}. \quad (\text{F.1})$$

The values in (F.1) are computed using

$$\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}} = \boldsymbol{\mu}_{\mathbf{x}} + \boldsymbol{\Sigma}_{\mathbf{xy}} \boldsymbol{\Sigma}_{\mathbf{yy}}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{y}}) \quad (\text{F.2})$$

$$\boldsymbol{\Sigma}_{\mathbf{xx}|\mathbf{y}} = \boldsymbol{\Sigma}_{\mathbf{xx}} - \boldsymbol{\Sigma}_{\mathbf{xy}} \boldsymbol{\Sigma}_{\mathbf{yy}}^{-1} \boldsymbol{\Sigma}_{\mathbf{xy}}^{\top} \quad (\text{F.3})$$

$$\boldsymbol{\mu}_{\hat{\mathbf{x}}|\mathbf{y}} = \boldsymbol{\mu}_{\hat{\mathbf{x}}} + \boldsymbol{\Sigma}_{\hat{\mathbf{x}}\mathbf{y}} \boldsymbol{\Sigma}_{\mathbf{yy}}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{y}}) \quad (\text{F.4})$$

$$\boldsymbol{\Sigma}_{\hat{\mathbf{x}}\hat{\mathbf{x}}|\mathbf{y}} = \boldsymbol{\Sigma}_{\hat{\mathbf{x}}\hat{\mathbf{x}}} - \boldsymbol{\Sigma}_{\hat{\mathbf{x}}\mathbf{y}} \boldsymbol{\Sigma}_{\mathbf{yy}}^{-1} \boldsymbol{\Sigma}_{\hat{\mathbf{x}}\mathbf{y}}^{\top}. \quad (\text{F.5})$$

Plugging (F.2)-(F.5) into (F.1), the CFID can be written as [22, Lemma 2]:

$$\begin{aligned} \text{CFID} = & \|\mu_{\underline{x}} - \mu_{\widehat{\underline{x}}}\|_2^2 + \text{tr} \left[ (\Sigma_{\underline{xy}} - \Sigma_{\widehat{\underline{xy}}}) \Sigma_{\underline{yy}}^{-1} (\Sigma_{\underline{xy}} - \Sigma_{\widehat{\underline{xy}}})^\top \right] \\ & + \text{tr} \left[ \Sigma_{\underline{xx}|\underline{y}} + \Sigma_{\widehat{\underline{xx}}|\underline{y}} - 2(\Sigma_{\underline{xx}|\underline{y}}^{1/2} \Sigma_{\widehat{\underline{xx}}|\underline{y}} \Sigma_{\underline{xx}|\underline{y}}^{1/2})^{1/2} \right], \end{aligned} \quad (\text{F.6})$$

where  $\Sigma_{\underline{yy}}^{-1}$  is typically implemented using a pseudo-inverse.

We now detail how the means and covariances in (F.6) are computed. We start with a dataset  $\{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^n$  of truth/measurement pairs. For each  $\mathbf{y}_t$ , we generate a set of  $P$  posterior samples  $\{\widehat{\mathbf{x}}_{ti}\}_{i=1}^P$ . We merge these samples with  $P$  repetitions of  $\mathbf{x}_t$  and  $\mathbf{y}_t$  to obtain  $\{(\mathbf{x}_{ti}, \mathbf{y}_{ti}, \widehat{\mathbf{x}}_{ti})\}_{i=1}^P$  for  $t = 1 \dots n$ . These terms are processed by a feature-generating network to yield the feature embeddings  $\{(\mathbf{x}_{ti}, \mathbf{y}_{ti}, \widehat{\mathbf{x}}_{ti})\}_{i=1}^P$ , which are then packed into matrices  $\underline{\mathbf{X}}$ ,  $\underline{\mathbf{Y}}$ , and  $\widehat{\underline{\mathbf{X}}}$  with  $Pn$  rows. We used the VGG-16 feature-generating network [49] for our MRI experiments, since [41] found that it gave results that correlated much better with radiologists' perceptions, while we used the standard Inception-v3 network [50] for our inpainting experiments. The embeddings are then used to compute the sample-mean values

$$\mu_{\underline{x}} \triangleq \frac{1}{Pn} \mathbf{1}^\top \underline{\mathbf{X}}, \quad \mu_{\underline{y}} \triangleq \frac{1}{Pn} \mathbf{1}^\top \underline{\mathbf{Y}}, \quad \mu_{\widehat{\underline{x}}} \triangleq \frac{1}{Pn} \mathbf{1}^\top \widehat{\underline{\mathbf{X}}}. \quad (\text{F.7})$$

We then subtract the sample mean from each row of  $\underline{\mathbf{X}}$ ,  $\underline{\mathbf{Y}}$ , and  $\widehat{\underline{\mathbf{X}}}$  to give the zero-mean embedding matrices  $\underline{\mathbf{X}}_{\text{zm}} \triangleq \underline{\mathbf{X}} - \mathbf{1} \mu_{\underline{x}}^\top$ ,  $\underline{\mathbf{Y}}_{\text{zm}} \triangleq \underline{\mathbf{Y}} - \mathbf{1} \mu_{\underline{y}}^\top$ , and  $\widehat{\underline{\mathbf{X}}}_{\text{zm}} \triangleq \widehat{\underline{\mathbf{X}}} - \mathbf{1} \mu_{\widehat{\underline{x}}}^\top$ , which are then used to compute the sample covariance matrices

$$\Sigma_{\underline{xx}} \triangleq \frac{1}{Pn} \underline{\mathbf{X}}_{\text{zm}}^\top \underline{\mathbf{X}}_{\text{zm}}, \quad \Sigma_{\underline{yy}} \triangleq \frac{1}{Pn} \underline{\mathbf{Y}}_{\text{zm}}^\top \underline{\mathbf{Y}}_{\text{zm}}, \quad \Sigma_{\widehat{\underline{xx}}} \triangleq \frac{1}{Pn} \widehat{\underline{\mathbf{X}}}_{\text{zm}}^\top \widehat{\underline{\mathbf{X}}}_{\text{zm}} \quad (\text{F.8a})$$

$$\Sigma_{\underline{xy}} \triangleq \frac{1}{Pn} \underline{\mathbf{X}}_{\text{zm}}^\top \underline{\mathbf{Y}}_{\text{zm}}, \quad \Sigma_{\widehat{\underline{xy}}} \triangleq \frac{1}{Pn} \widehat{\underline{\mathbf{X}}}_{\text{zm}}^\top \underline{\mathbf{Y}}_{\text{zm}}. \quad (\text{F.8b})$$

We plug the sample statistics from (F.7)-(F.8) into (F.2)-(F.5), which yields the statistics needed to compute the CFID in (F.6). In [22], the authors use  $P = 1$  in all of their experiments. To be consistent with how we evaluated the other metrics, we use  $P = 32$  unless otherwise noted.

## G MR imaging details

We now give details on magnetic resonance (MR) image recovery. Suppose that the goal is to recover the  $N$ -pixel MR image  $\mathbf{i} \in \mathbb{C}^N$  from the multicoil measurements  $\{\mathbf{k}_c\}_{c=1}^C$ , where [39]

$$\mathbf{k}_c = \mathbf{M} \mathbf{F} \mathbf{S}_c \mathbf{i} + \mathbf{n}_c. \quad (\text{G.1})$$

In (G.1),  $C$  refers to the number of coils,  $\mathbf{k}_c \in \mathbb{C}^M$  are the measurements from the  $c$ th coil,  $\mathbf{M} \in \mathbb{R}^{M \times N}$  is a sub-sampling operator containing rows from  $\mathbf{I}_N$ —the  $N \times N$  identity matrix,  $\mathbf{F} \in \mathbb{C}^{N \times N}$  is the unitary 2D discrete Fourier transform,  $\mathbf{S}_c \in \mathbb{C}^{N \times N}$  is a diagonal matrix containing the sensitivity map of the  $c$ th coil, and  $\mathbf{n}_c \in \mathbb{C}^M$  is noise. From (G.1), it can be seen that the MR measurements are collected in the spatial Fourier domain, otherwise known as the “k-space.” The sensitivity maps  $\{\mathbf{S}_c\}$  are estimated from  $\{\mathbf{k}_c\}$  using ESPIRiT [40] (in our case via SigPy [51]), which yields maps with the property  $\sum_{c=1}^C \mathbf{S}_c^H \mathbf{S}_c = \mathbf{I}_N$ . The ratio  $R \triangleq \frac{N}{M}$  is known as the acceleration rate.

There are different ways that one could apply the generative posterior sampling framework to multicoil MR image recovery. One is to configure the generator to produce posterior samples  $\widehat{\mathbf{i}}$  of the complex image  $\mathbf{i}$ . Another is to configure the generator to produce posterior samples  $\widehat{\mathbf{x}}$  of the stack  $\mathbf{x} \triangleq [\mathbf{x}_1^\top, \dots, \mathbf{x}_C^\top]^\top$  of “coil images”  $\mathbf{x}_c \triangleq \mathbf{S}_c \mathbf{i}$  and later coil-combining them to yield a complex image estimate  $\widehat{\mathbf{i}} \triangleq [\mathbf{S}_1^H, \dots, \mathbf{S}_C^H] \widehat{\mathbf{x}}$ . We take the latter approach. Furthermore, rather than feeding our generator with k-space measurements  $\mathbf{k}_c$ , we choose to feed it with aliased coil images  $\mathbf{y}_c \triangleq \mathbf{F}^H \mathbf{M}^\top \mathbf{k}_c$ . Writing (G.1) in terms of the coil images, we obtain

$$\mathbf{y}_c = \mathbf{F}^H \mathbf{M}^\top \mathbf{M} \mathbf{F} \mathbf{x}_c + \mathbf{w}_c, \quad (\text{G.2})$$

where  $\mathbf{w}_c \triangleq \mathbf{F}^H \mathbf{M}^\top \mathbf{n}_c$ . Then we can stack  $\{\mathbf{y}_c\}$  and  $\{\mathbf{w}_c\}$  column-wise into vectors  $\mathbf{y}$  and  $\mathbf{w}$ , and set  $\mathbf{A} = \mathbf{I}_C \otimes \mathbf{F}^H \mathbf{M}^\top \mathbf{M} \mathbf{F} \in \mathbb{C}^{NC \times NC}$ , to obtain the formulation  $\mathbf{y} = \mathbf{A} \mathbf{x} + \mathbf{w}$  described in Section 1.

To train our generator, we assume to have access to paired training examples  $\{(\mathbf{x}_t, \mathbf{y}_t)\}$ , where  $\mathbf{x}_t$  is a stack of coil images and  $\mathbf{y}_t$  is the corresponding stack of k-space coil measurements. The fastMRI multicoil dataset [33] provides  $\{(\mathbf{x}_t, \mathbf{k}_t)\}$ , from which we can easily obtain  $\{(\mathbf{x}_t, \mathbf{y}_t)\}$ .

## H Data-consistency

In this section, we describe a data-consistency procedure that can be optionally used when our cGAN is used to solve a *linear* inverse problem, i.e., to recover  $\mathbf{x}$  from  $\mathbf{y}$  under a model of the form

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}, \quad (\text{H.1})$$

where  $\mathbf{A}$  is a known linear operator and  $\mathbf{w}$  is unknown noise. The motivation is that, in some applications, such as medical imaging or inpainting, the end user may feel comfortable knowing that the generated samples  $\{\hat{\mathbf{x}}_i\}$  are consistent with the measurements  $\mathbf{y}$  in that

$$\mathbf{y} = \mathbf{A}\hat{\mathbf{x}}_i. \quad (\text{H.2})$$

When (H.2) holds,  $\mathbf{A}^+\mathbf{y} = \mathbf{A}^+\mathbf{A}\hat{\mathbf{x}}_i$  must also hold, where  $(\cdot)^+$  denotes the pseudo-inverse. The quantity  $\mathbf{A}^+\mathbf{A}$  can be recognized as the orthogonal projection matrix associated with the row space of  $\mathbf{A}$ . So, (H.2) requires the component of  $\hat{\mathbf{x}}_i$  in the row space of  $\mathbf{A}$  to equal  $\mathbf{A}^+\mathbf{y}$ , while placing no constraints on the component of  $\hat{\mathbf{x}}_i$  in the nullspace of  $\mathbf{A}$ . This suggests the following data-consistency procedure:

$$\hat{\mathbf{x}}_i = (\mathbf{I} - \mathbf{A}^+\mathbf{A})\hat{\mathbf{x}}_i^{\text{raw}} + \mathbf{A}^+\mathbf{y}. \quad (\text{H.3})$$

where  $\hat{\mathbf{x}}_i^{\text{raw}}$  is the raw generator output. We note that a version of this idea for point estimation was proposed in [3].

The data-consistency procedure (H.3) ensures that any generative method will generate only the component of  $\mathbf{x}$  that lies in the nullspace of  $\mathbf{A}$ . Consequently, (H.3) is admissible only when  $\mathbf{A}$  has a non-trivial nullspace. Also, because no attempt is made to remove the noise  $\mathbf{w}$  in  $\mathbf{y}$ , this approach is recommended only for low-noise applications. For high-noise applications, an extension based on the dual-decomposition approach [52] would be more appropriate, but we leave this to future work.

When applying (H.3) to the MRI formulation in Appendix G, we note that  $\mathbf{A} = \mathbf{I}_C \otimes \mathbf{F}^H \mathbf{M}^T \mathbf{M} \mathbf{F}$  is an orthogonal projection matrix, and so  $\mathbf{I} - \mathbf{A}^+\mathbf{A} = \mathbf{I} - \mathbf{A} = \mathbf{I} \otimes \mathbf{F}^H (\mathbf{I} - \mathbf{M}^T \mathbf{M}) \mathbf{F}$ .

## I Implementation details

The code for our model can be found here: <https://github.com/matt-bendel/rcGAN>.

### I.1 MRI

#### I.1.1 cGAN training

At each training iteration, our cGAN’s generator takes in  $n_{\text{batch}}$  measurement samples  $\mathbf{y}_t$  and  $P_{\text{train}}$  code vectors for every  $\mathbf{y}_t$ , and performs an optimization step on the loss

$$\mathcal{L}_G(\boldsymbol{\theta}) \triangleq \beta_{\text{adv}} \mathcal{L}_{\text{adv}}(\boldsymbol{\theta}, \phi) + \mathcal{L}_{1, P_{\text{train}}}(\boldsymbol{\theta}) - \beta_{\text{SD}} \mathcal{L}_{\text{SD}, P_{\text{train}}}(\boldsymbol{\theta}), \quad (\text{I.1})$$

where by default we use  $\beta_{\text{adv}} = 1\text{e-}5$ ,  $n_{\text{batch}} = 36$ ,  $P_{\text{train}} = 2$ , and update  $\beta_{\text{SD}}$  via (23) using  $P_{\text{val}} = 8$ . Then, using the  $P_{\text{train}} n_{\text{batch}}$  generator outputs, our cGAN’s discriminator performs an optimization step on the loss

$$\mathcal{L}_D(\phi) = -\mathcal{L}_{\text{adv}}(\boldsymbol{\theta}, \phi) + \alpha_1 \mathcal{L}_{\text{gp}}(\phi) + \alpha_2 \mathcal{L}_{\text{drift}}(\phi), \quad (\text{I.2})$$

with gradient penalty  $\mathcal{L}_{\text{gp}}$  from [24]. As per [29],  $\mathcal{L}_{\text{drift}}$  is a drift penalty,  $\alpha_1 = 10$ ,  $\alpha_2 = 0.001$ , and one discriminator update was used per generator update. The models were trained for 100 epochs using the Adam optimizer [53] with a learning rate of  $1\text{e-}3$ ,  $\beta_1 = 0$ , and  $\beta_2 = 0.99$ , as in [29]. Running PyTorch on a server with 4 Tesla A100 GPUs, each with 82 GB of memory, the training of an MRI cGAN took approximately 1 day.

Adler and Öktem’s cGAN [8] uses generator loss  $\beta_{\text{adv}} \mathcal{L}_{\text{adv}}^{\text{adler}}(\boldsymbol{\theta}, \phi)$ , where  $\mathcal{L}_{\text{adv}}^{\text{adler}}(\boldsymbol{\theta}, \phi)$  was described in (5), and discriminator loss  $-\mathcal{L}_{\text{adv}}^{\text{adler}}(\boldsymbol{\theta}, \phi) + \alpha_1 \mathcal{L}_{\text{gp}}(\phi) + \alpha_2 \mathcal{L}_{\text{drift}}(\phi)$  with the values of  $\alpha_1 = 10$ ,  $\alpha_2 = 0.001$ , and  $\beta_{\text{adv}} = 1$ , as in the original paper.

Ohayon et al.’s cGAN [23] uses generator loss  $\beta_{\text{adv}}\mathcal{L}_{\text{adv}}(\theta, \phi) + \mathcal{L}_{2, P_{\text{train}}}(\theta)$ , where  $\mathcal{L}_{2, P_{\text{train}}}(\theta)$  was described in (14), and discriminator loss  $-\mathcal{L}_{\text{adv}}(\theta, \phi) + \alpha_1\mathcal{L}_{\text{gp}}(\phi) + \alpha_2\mathcal{L}_{\text{drift}}(\phi)$  with the values  $\alpha_1 = 10$ ,  $\alpha_2 = 0.001$ , and  $\beta_{\text{adv}} = 1\text{e-}5$ . We modify  $\beta_{\text{adv}}$  to re-balance the loss due to an increased magnitude of our discriminator’s outputs.

All three cGANs used the same generator and discriminator architectures (detailed below), except that Adler and Öktem’s discriminator used extra input channels to facilitate the 3-input loss  $\mathcal{L}_{\text{adv}}^{\text{adler}}(\theta, \phi)$  from (5).

### I.1.2 cGAN generator architecture

For our MRI experiments, we take inspiration from the UNet architecture from [36], using it as the basis for the cGAN generators. The primary input  $y$  is concatenated with the code vector  $z$  and fed through the UNet. The network consists of 4 pooling layers with 128 initial channels. However, instead of pooling, we opt to use convolutions with kernels of size  $3 \times 3$ , “same” padding, and a stride of 2 when downsampling. Likewise, we upsample using transpose convolutions, again with kernels of size  $3 \times 3$ , “same” padding, and a stride of 2. All other convolutions utilize kernels of size  $3 \times 3$ , “same” padding, and a stride of 1.

Within each encoder and decoder layer we include a residual block, the architecture of which can be found in [8]. We use instance-norm for all normalization layers and parametric ReLUs as our activation functions, in which the network learns the optimal “negative slope.” Finally, we include 5 residual blocks at the base of the UNet, in between the encoder and decoder. This is done in an effort to artificially increase the depth of the network and is inspired by [54]. Our generator has 86 734 334 trainable parameters.

### I.1.3 cGAN discriminator architecture

Our discriminator is a standard CNN with 6 layers and 1 fully-connected layer. In the first 3 layers, we use convolutions with kernels of size  $3 \times 3$ , “same” padding. We reduce spatial resolution with average pooling, using  $2 \times 2$  kernels with a stride of 2. We use batch-norm as our normalization layer and leaky ReLUs with a “negative-slope” of 0.2 as our activation functions. The network outputs an estimated Wasserstein score for the whole image.

### I.1.4 E2E-VarNet

For the Sriram et al.’s E2E-VarNet [37], we use the same training procedure and hyperparameters outlined in [19] other than replacing the sampling pattern with the GRO undersampling mask. As in [19], we use the SENSE-based coil-combined image as the ground truth instead of the RSS image.

### I.1.5 Langevin approach

For Jalal et al.’s MRI approach [19], we do not modify the original implementation from [38] other than replacing the default sampling pattern with the GRO undersampling mask. We generated 32 samples for 72 different test images using a batch-size of 4, which took roughly 6 days. These samples were generated on a server with 4 NVIDIA V100 GPUs, each with 32 GB of memory. We used 4 samples per batch (and recorded the time to generate 4 samples in Table 1) because the code from [38] is written to generate one sample per GPU.

## I.2 Inpainting

### I.2.1 Our cGAN

For our generator and discriminator, we use the CoModGAN networks from [9]. Unlike CoModGAN, however, we train our cGAN with  $\mathcal{L}_{1, \text{SD}, P_{\text{train}}}$  regularization and we do not use MBSD at the discriminator. We use the same general training and testing procedure described in Section 4.2, but with  $\beta_{\text{adv}} = 5\text{e-}3$ ,  $n_{\text{batch}} = 100$ , and 110 epochs of cGAN training. Running PyTorch on a server with 4 Tesla A100 GPUs, each with 82 GB of memory, the training takes approximately 2 days.

Table J.1: The mean and covariance components of CFID, along with the total CFID, for the generative models in the MRI and inpainting experiments. For the MRI experiment, CFID<sup>1</sup> used 72 test samples and  $P = 32$ , CFID<sup>2</sup> used 2 376 test samples and  $P = 8$ , and CFID<sup>3</sup> used all 14 576 samples and  $P = 1$ . For the inpainting experiment, CFID<sup>1</sup> used 1 000 test images and  $P = 32$ , CFID<sup>2</sup> used 3 000 test and validation images and  $P = 8$ , and CFID<sup>3</sup> used all 30 000 images and  $P = 1$ .

Model	CFID <sup>1</sup> <sub>mean</sub> ↓	CFID <sup>1</sup> <sub>cov</sub> ↓	CFID <sup>1</sup> ↓	CFID <sup>2</sup> <sub>mean</sub> ↓	CFID <sup>2</sup> <sub>cov</sub> ↓	CFID <sup>2</sup> ↓	CFID <sup>3</sup> <sub>mean</sub> ↓	CFID <sup>3</sup> <sub>cov</sub> ↓	CFID <sup>3</sup> ↓
<i>R = 4 MRI</i>									
Langevin (Jalal [19])	1.89	3.40	5.29	-	-	-	-	-	-
cGAN (Adler [8])	3.12	3.27	6.39	2.79	1.48	4.27	2.71	1.10	3.82
cGAN (Ohayon [23])	1.94	2.12	4.06	2.27	1.00	3.27	2.29	0.66	2.95
cGAN (Ours)	<b>0.98</b>	<b>2.12</b>	<b>3.10</b>	<b>0.86</b>	<b>0.68</b>	<b>1.54</b>	<b>0.86</b>	<b>0.43</b>	<b>1.29</b>
<i>R = 8 MRI</i>									
Langevin (Jalal [19])	2.61	4.73	7.34	-	-	-	-	-	-
cGAN (Adler [8])	5.00	5.10	10.10	4.16	2.14	6.30	4.09	1.63	5.72
cGAN (Ohayon [23])	2.73	3.31	6.04	3.07	1.52	4.59	3.30	0.97	4.27
cGAN (Ours)	<b>1.55</b>	<b>3.32</b>	<b>4.87</b>	<b>1.24</b>	<b>0.99</b>	<b>2.23</b>	<b>1.17</b>	<b>0.62</b>	<b>1.79</b>
<i>Inpainting</i>									
Score SDE (Song [20])	0.97	<b>38.69</b>	<b>39.66</b>	-	-	-	0.90	<b>4.21</b>	5.11
CoModGAN (Zhao [9])	0.42	41.21	41.63	0.35	25.39	25.74	0.32	4.98	5.29
cGAN (Ours)	<b>0.32</b>	39.41	39.73	<b>0.25</b>	<b>22.32</b>	<b>22.58</b>	<b>0.24</b>	4.45	<b>4.69</b>

## I.2.2 CoModGAN

We use the PyTorch implementation of CoModGAN from [44] and train the model to inpaint a  $128 \times 128$  centered square on  $256 \times 256$  CelebA-HQ images. The total training time on a server with 4 NVIDIA A100 GPUs, each with 82 GB of memory, is roughly 2 days.

## I.2.3 Score-based SDE

For the inpainting experiment in Section 4.3, we compare against Song et al.’s more recent SDE technique [20], for which we use the publicly available pretrained weights, the suggested settings for the  $256 \times 256$  CelebA-HQ dataset, and the code from the official PyTorch implementation [45]. We generate 32 samples for all 1 000 images in our test set, using a batch-size of 20 and generating 32 samples for each batch element concurrently. The total generation time on a server with 4 NVIDIA A100 GPUs, each with 82 GB of memory, is roughly 9 days.

# J Additional experimental results

## J.1 CFID decomposition into mean and covariance components

In this section, we investigate the small-sample bias effects of CFID, which have been previously noted in [22]. To do this, we write the CFID from (F.1) as a sum of two terms: a term that quantifies the conditional-mean error and a term that quantifies the conditional-covariance error:

$$\text{CFID} = \text{CFID}_{\text{mean}} + \text{CFID}_{\text{cov}} \quad (\text{J.1})$$

$$\text{CFID}_{\text{mean}} \triangleq \mathbb{E}_y \{ \|\mu_{\mathbf{x}|y} - \mu_{\hat{\mathbf{x}}|y}\|_2^2 \} \quad (\text{J.2})$$

$$\text{CFID}_{\text{cov}} \triangleq \text{tr} [\Sigma_{\mathbf{xx}|y} + \Sigma_{\mathbf{xx}|y} - 2(\Sigma_{\mathbf{xx}|y}^{1/2} \Sigma_{\mathbf{xx}|y} \Sigma_{\mathbf{xx}|y}^{1/2})^{1/2}]. \quad (\text{J.3})$$

To verify that (J.3) quantifies the error in  $\Sigma_{\mathbf{xx}|y}$ , notice that (J.3) equals zero when  $\Sigma_{\mathbf{xx}|y} = \Sigma_{\mathbf{xx}|y}$  and is otherwise positive (by Cauchy Schwarz).

In Table J.1, we report CFID<sub>mean</sub> and CFID<sub>cov</sub> for the MRI and inpainting experiments, in addition to the total CFID (also shown in Tables 1 and 4). As before, we computed CFID on three test sets for each experiment, which contained 72, 2 376, and 14 576 samples respectively for MRI, and 1000, 3000, and 30 000 samples respectively for inpainting. Due to the slow sample-generation time of the Langevin/score-based methods [19, 20], we did not have the computational resources to evaluate them on all datasets, and that’s why certain table entries are blank.

For both MRI experiments, Table J.1 shows our method outperforming the competing methods in both the mean and covariance components of CFID (and thus the total CFID) for all sample sizes.

And, in the inpainting experiment, Table J.1 shows our method outperforming CoModGAN in both the mean and covariance components (and thus the total CFID) for all sample sizes.

For the inpainting experiment, Table J.1 shows our method outperforming the score-based approach in total CFID on the 3000- and 30 000-sample test sets but not on the 1000-sample test set. However, we now argue that the 1000-sample inpainting experiment is heavily affected by small-sample bias, and therefore untrustworthy. Looking at the mean component of CFID (i.e.,  $\text{CFID}_{\text{mean}}^1$ ,  $\text{CFID}_{\text{mean}}^2$ , and  $\text{CFID}_{\text{mean}}^3$ ) across the inpainting experiments, we see that the values are relatively small and stable with sample size. But looking at the covariance component of CFID (i.e.,  $\text{CFID}_{\text{cov}}^1$ ,  $\text{CFID}_{\text{cov}}^2$ , and  $\text{CFID}_{\text{cov}}^3$ ) across the inpainting experiments, we see that the values are large and heavily dependent on sample size. For the 1000-sample inpainting experiment, the total CFID is dominated by the covariance component and thus strongly affected by small-sample bias. Consequently, for the 1000-sample inpainting experiment, the total CFID is not trustworthy.

## K Additional reconstruction plots

### K.1 $R = 4$ MRI Reconstruction

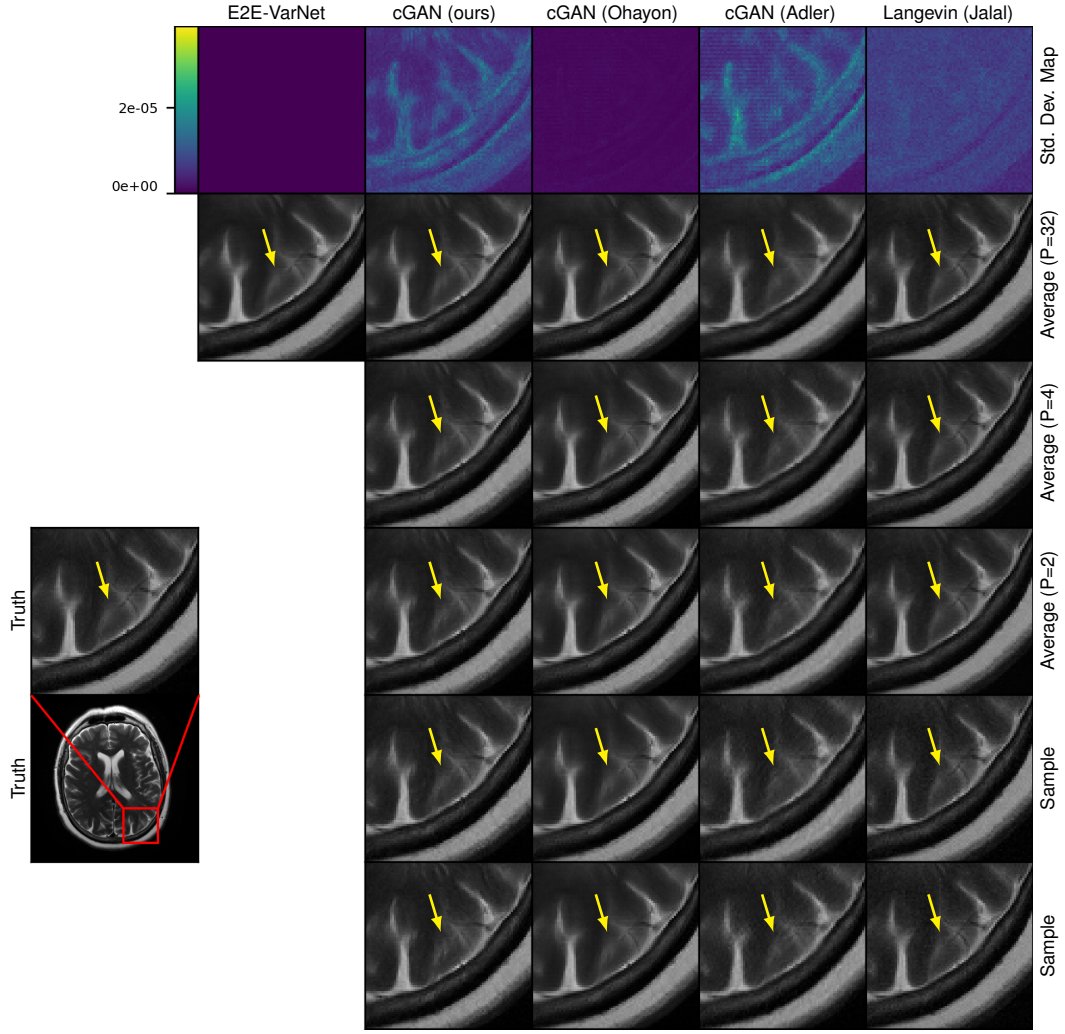


Figure K.1: Example  $R = 4$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{x}_{(P)}$  with  $P = 32$ , Row three:  $\hat{x}_{(P)}$  with  $P = 4$ , Row four:  $\hat{x}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

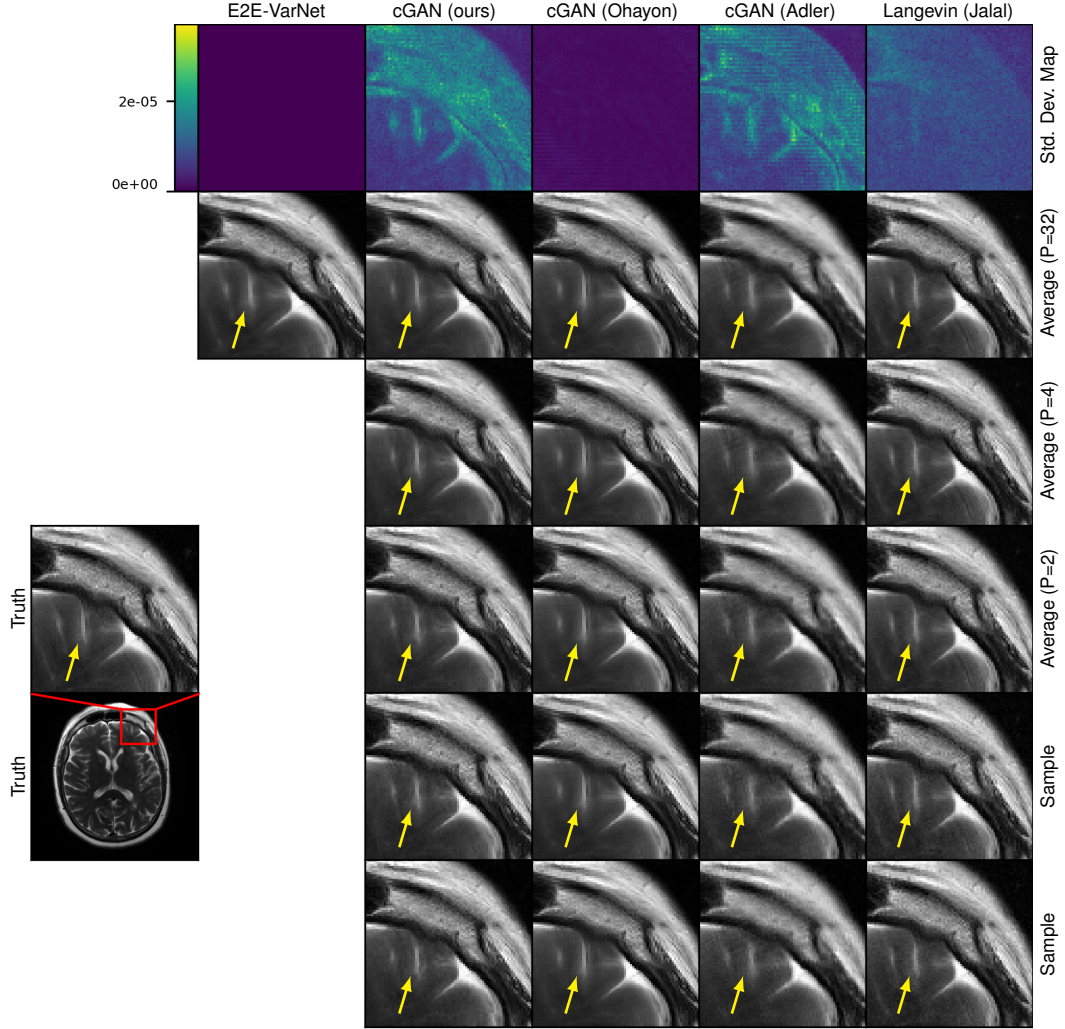


Figure K.2: Example  $R = 4$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.



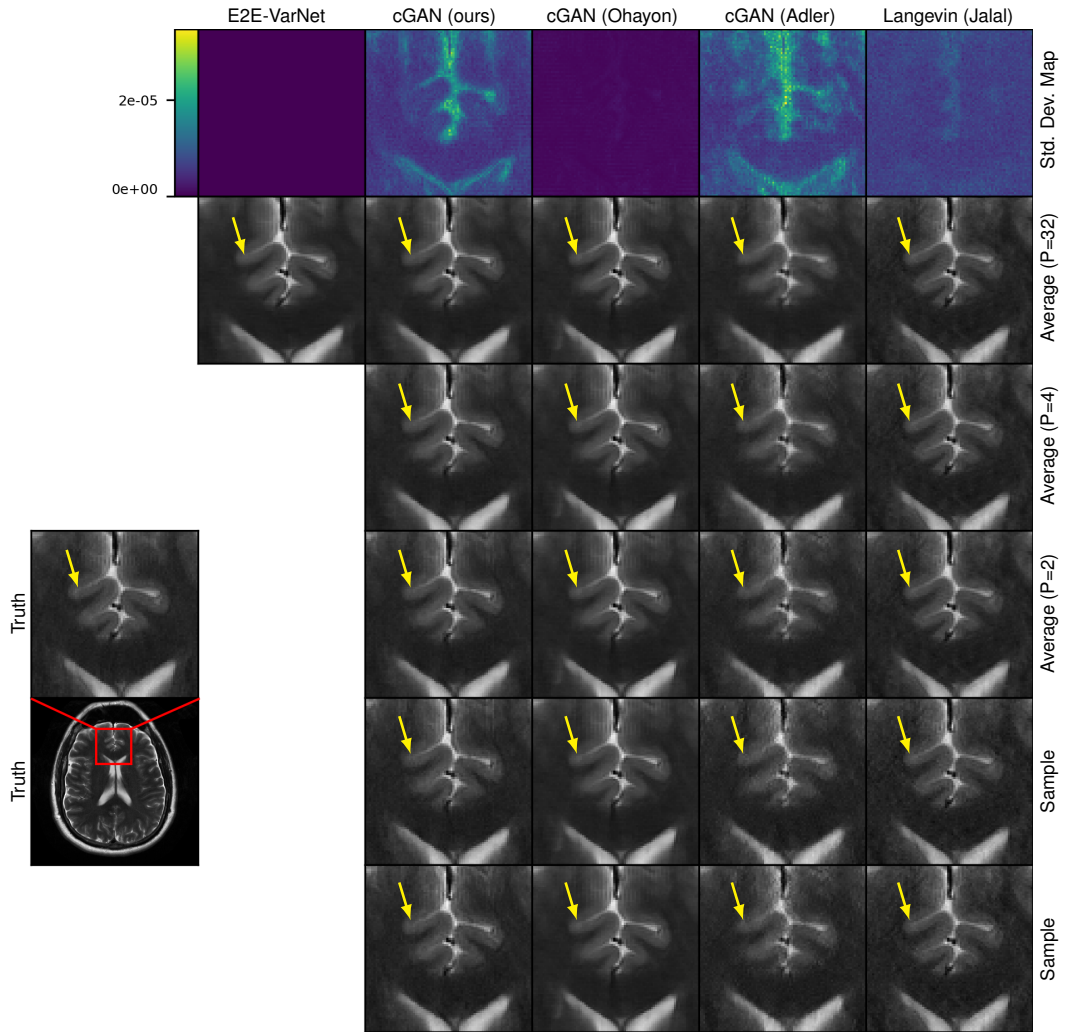


Figure K.3: Example  $R = 4$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

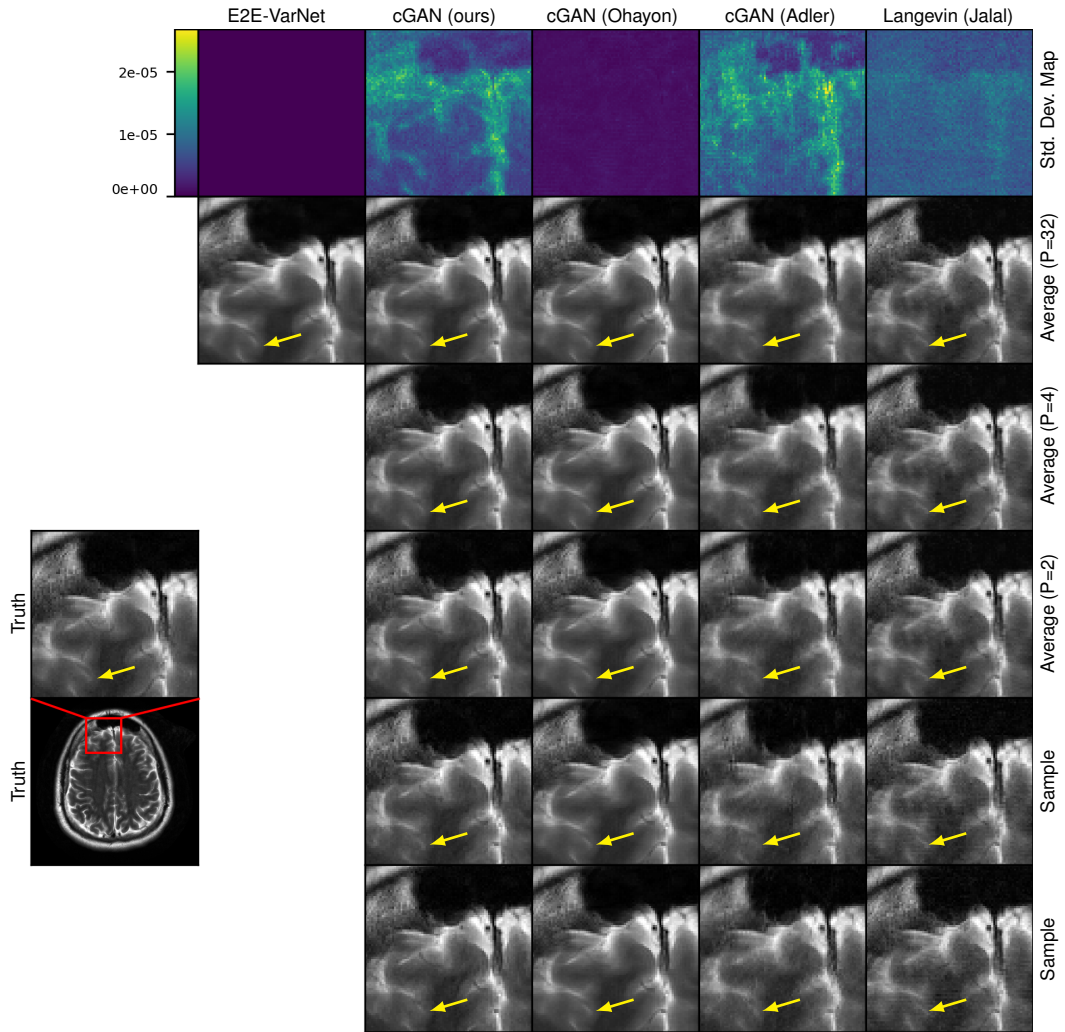


Figure K.4: Example  $R = 4$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

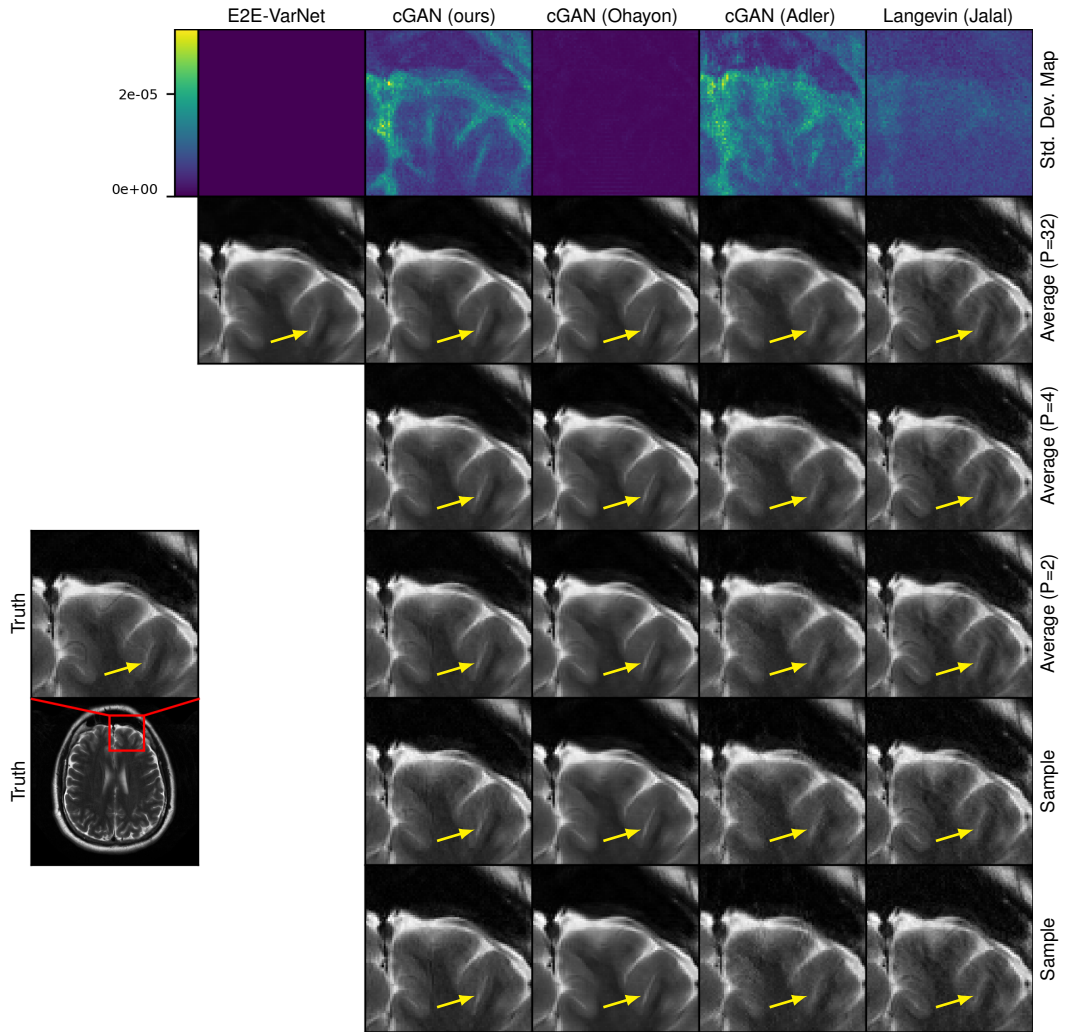


Figure K.5: Example  $R = 4$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

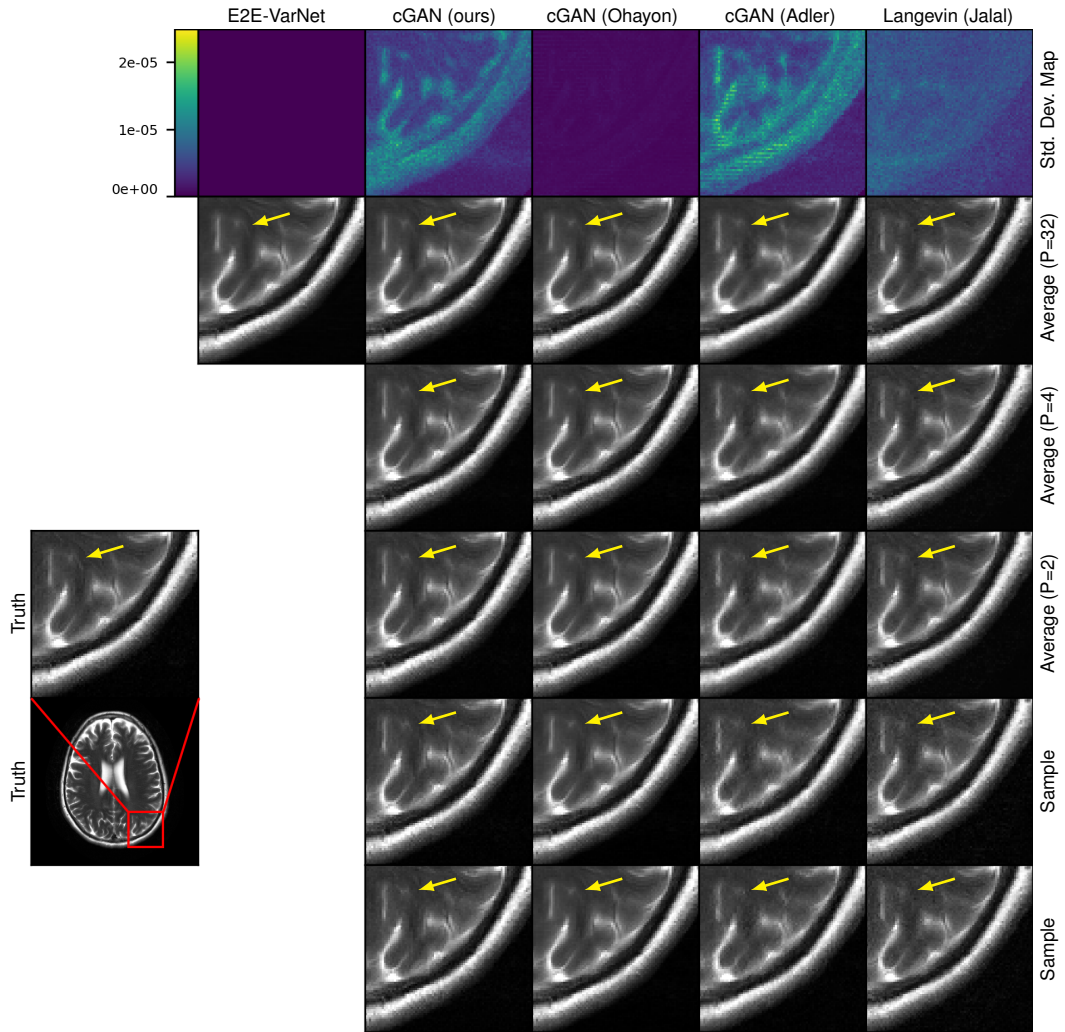


Figure K.6: Example  $R = 4$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

## K.2 $R = 8$ MRI Reconstruction

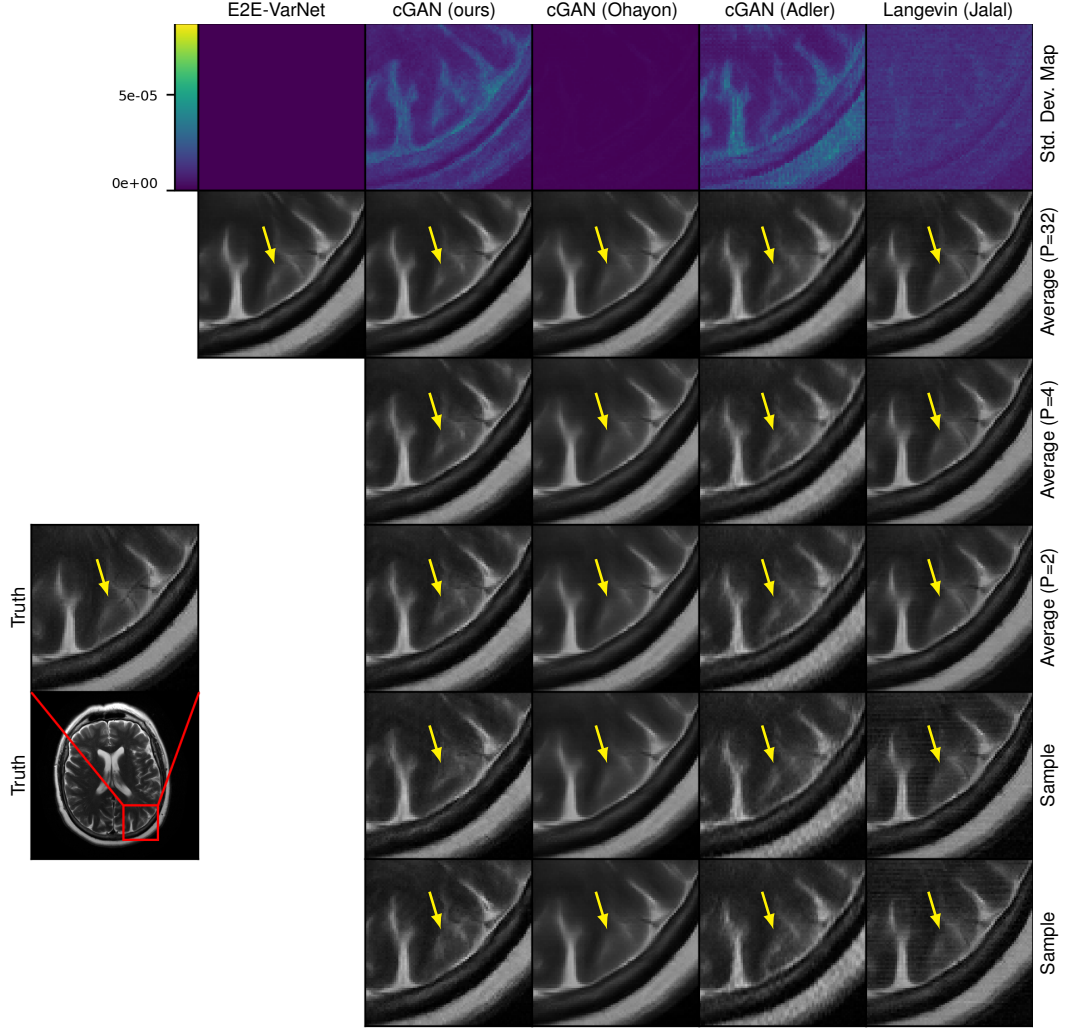


Figure K.7: Example  $R = 8$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.



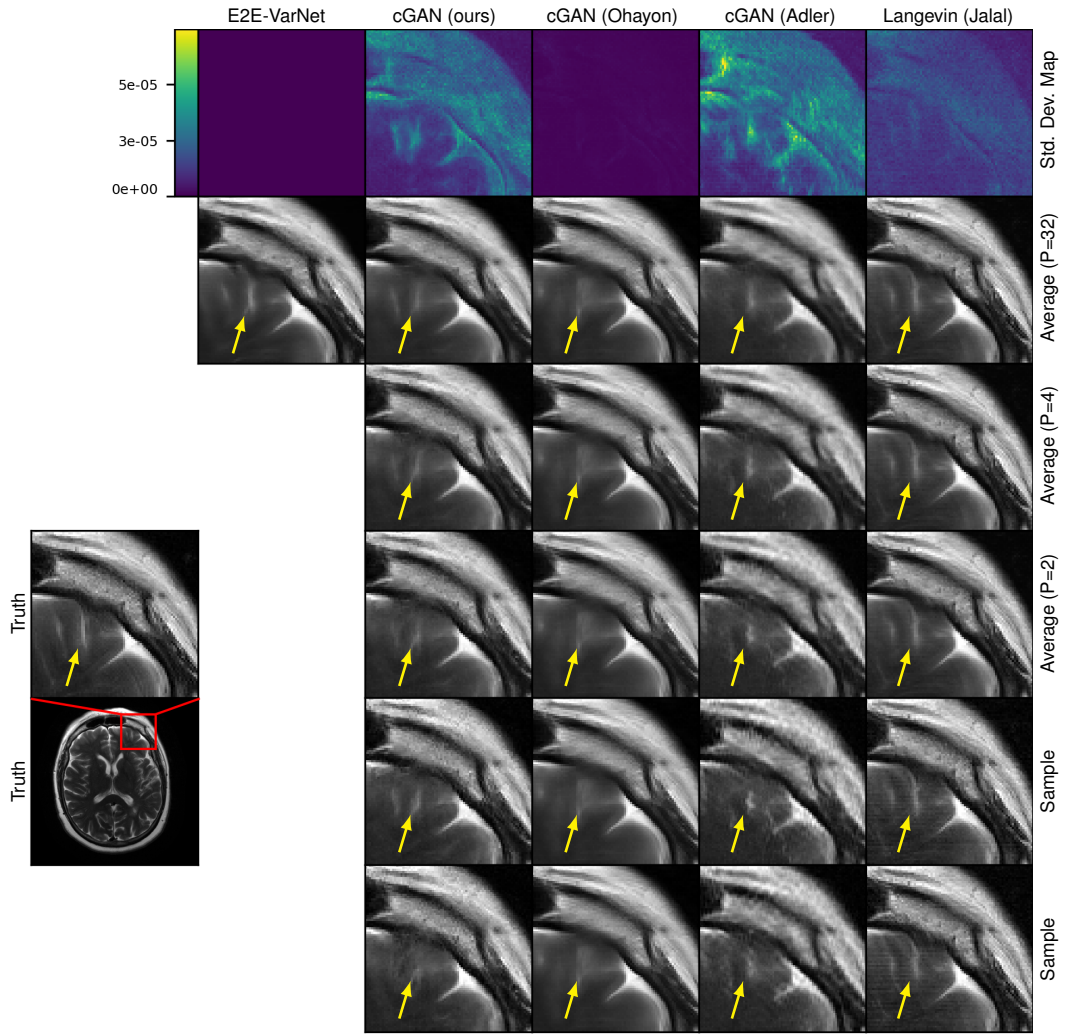


Figure K.8: Example  $R = 8$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

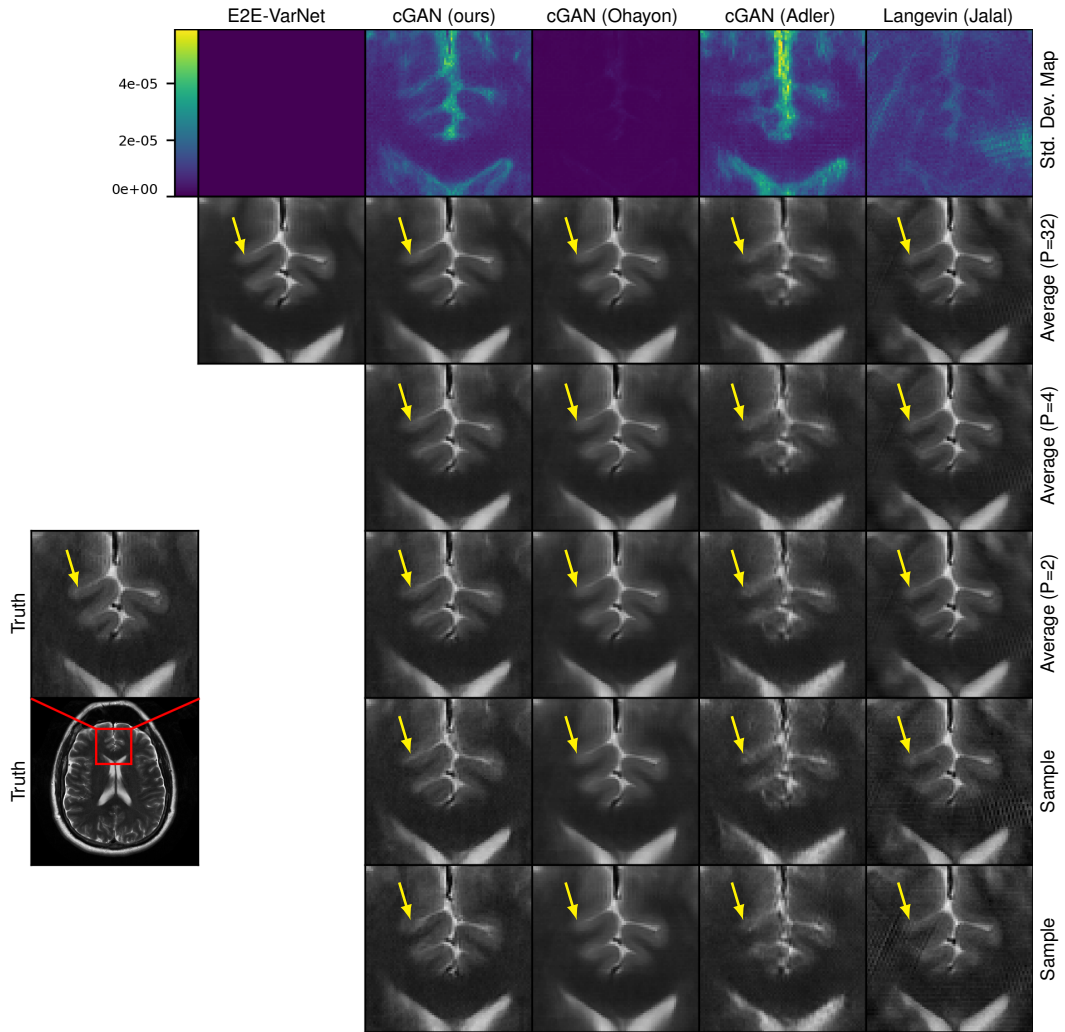


Figure K.9: Example  $R = 8$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

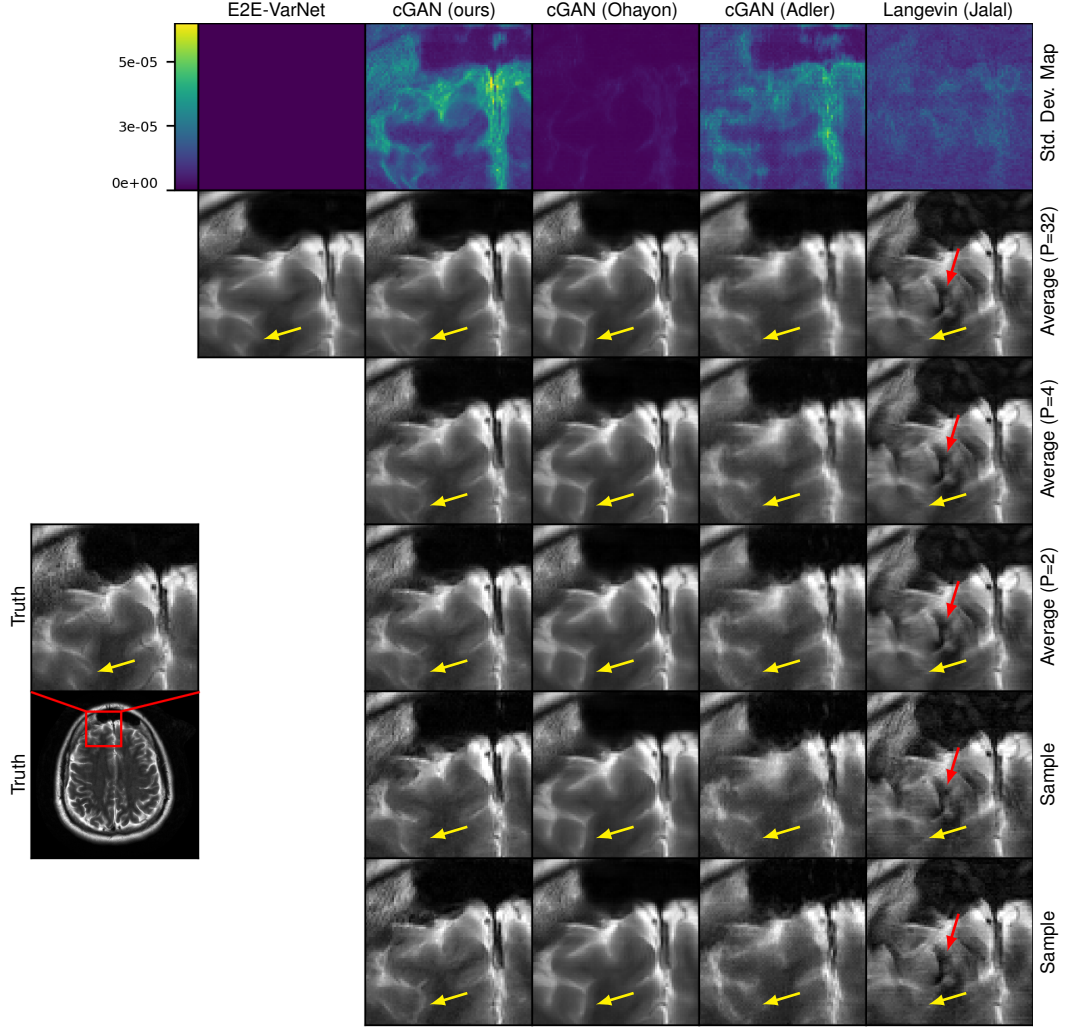


Figure K.10: Example  $R = 8$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The yellow arrows indicate regions of meaningful variation across posterior samples. The red arrows show visible artifacts in the Langevin recovery.



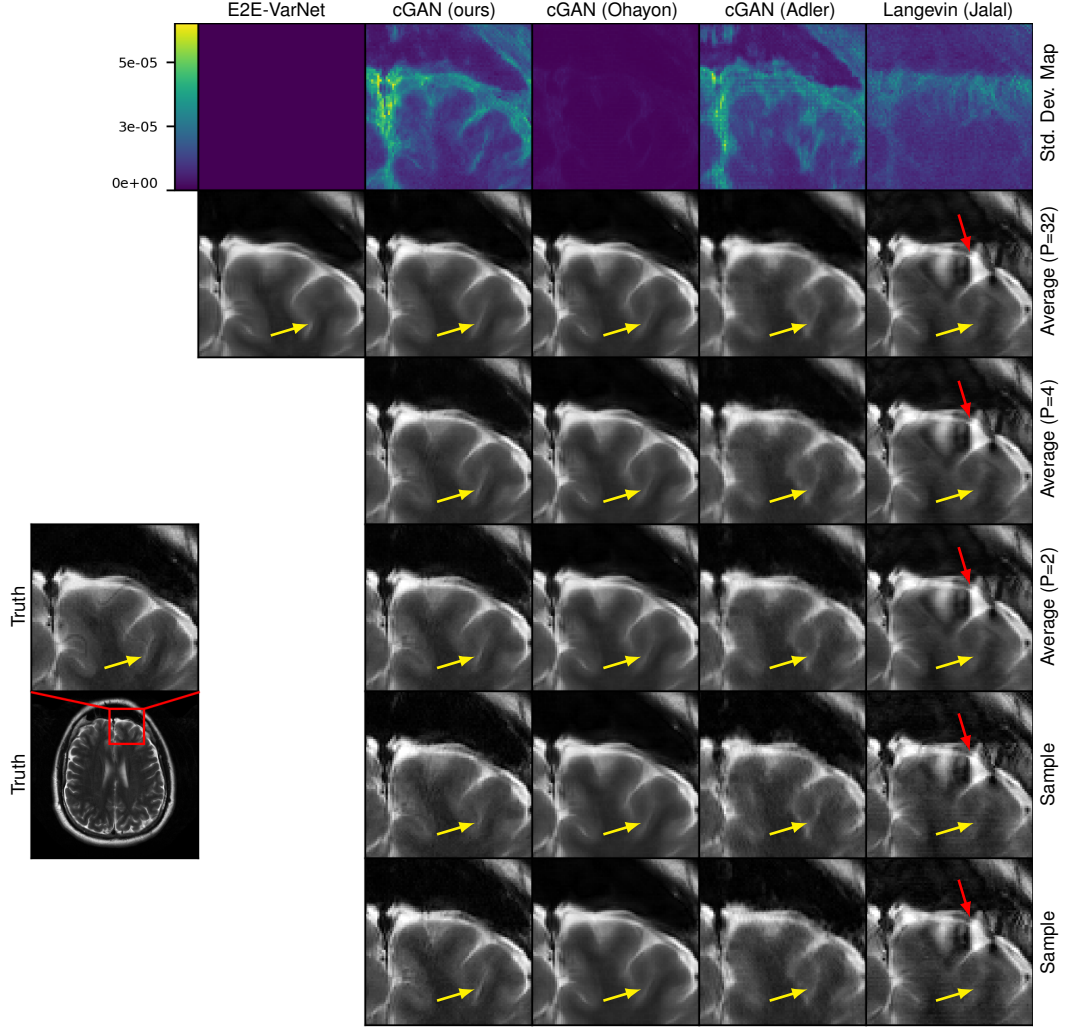


Figure K.11: Example  $R = 8$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The yellow arrows indicate regions of meaningful variation across posterior samples. The red arrows show visible artifacts in the Langevin recovery.

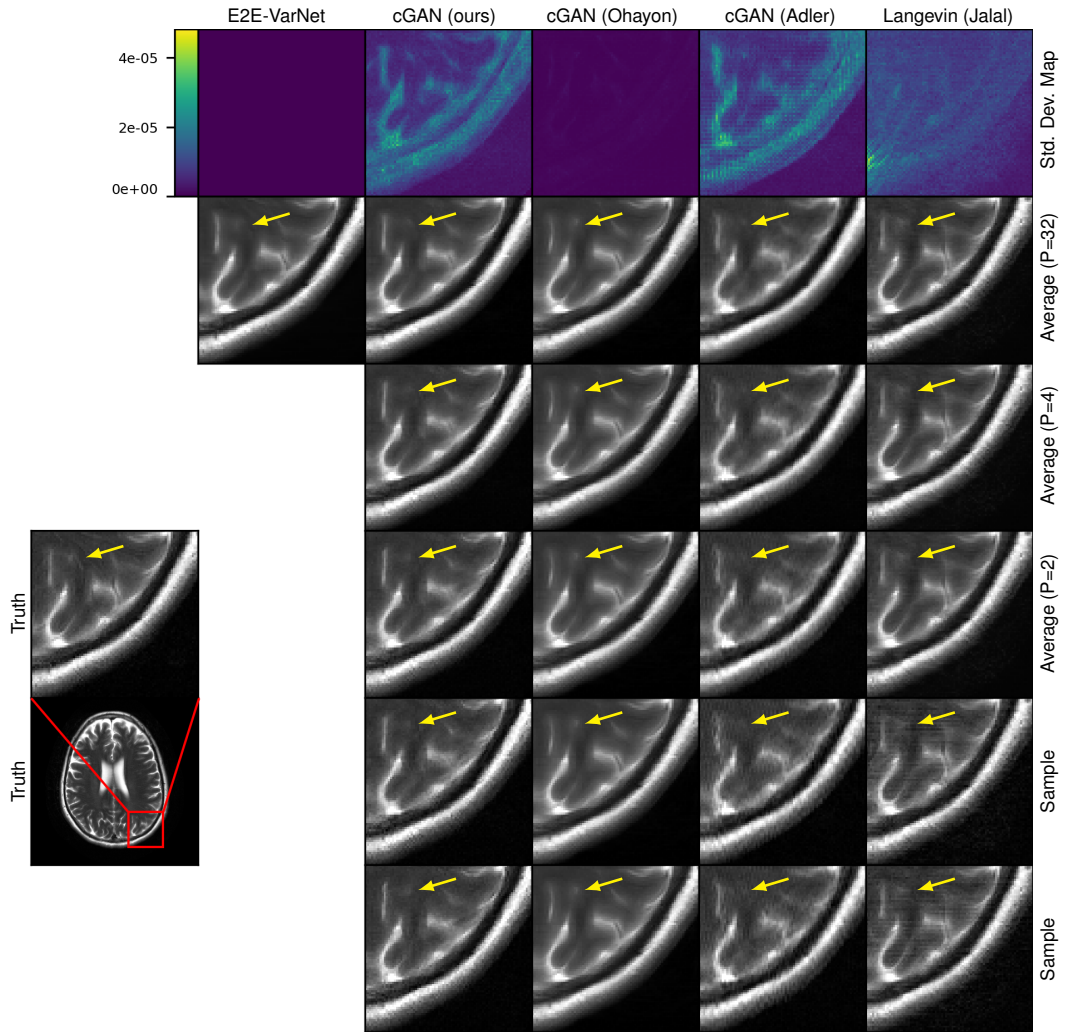


Figure K.12: Example  $R = 8$  MRI reconstruction. Row one: pixel-wise SD with  $P = 32$ , Row two:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 32$ , Row three:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 4$ , Row four:  $\hat{\mathbf{x}}_{(P)}$  with  $P = 2$ , Rows five and six: posterior samples. The arrows indicate regions of meaningful variation across posterior samples.

### K.3 Inpainting



Figure K.13: Example of inpainting a  $128 \times 128$  square on a  $256 \times 256$  resolution CelebA-HQ image.

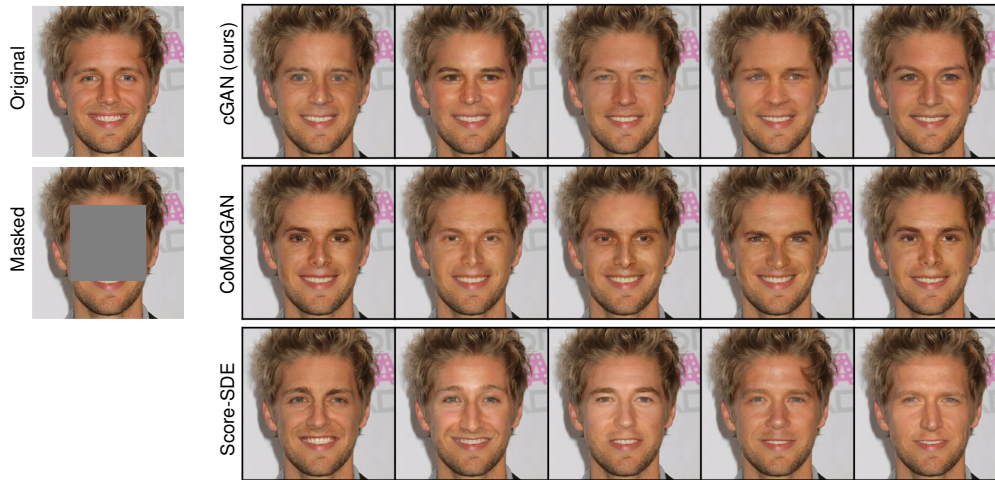


Figure K.14: Example of inpainting a  $128 \times 128$  square on a  $256 \times 256$  resolution CelebA-HQ image.

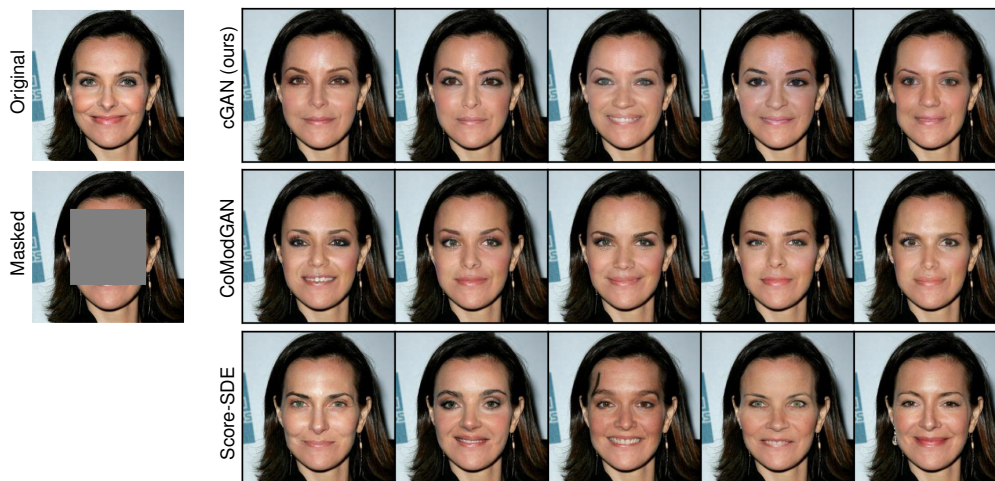


Figure K.15: Example of inpainting a  $128 \times 128$  square on a  $256 \times 256$  resolution CelebA-HQ image.



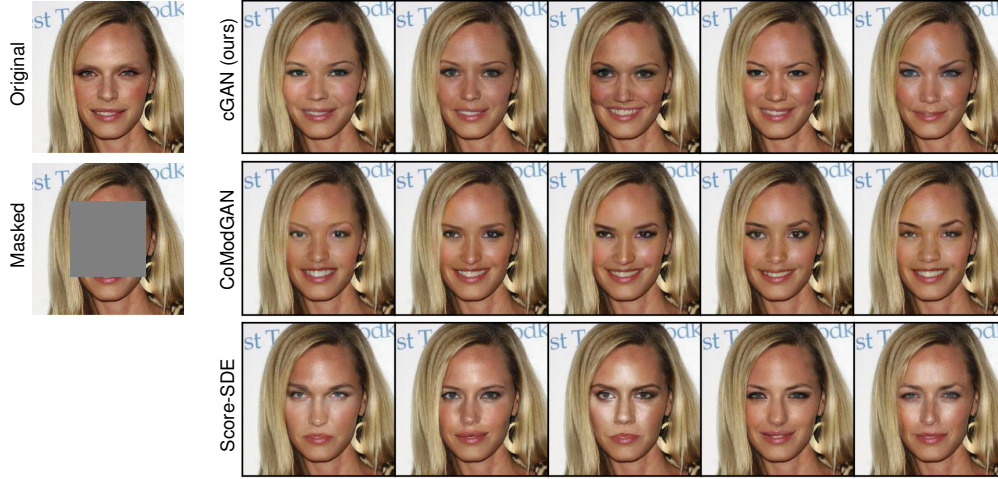


Figure K.16: Example of inpainting a  $128 \times 128$  square on a  $256 \times 256$  resolution CelebA-HQ image.



Figure K.17: Example of inpainting a  $128 \times 128$  square on a  $256 \times 256$  resolution CelebA-HQ image.

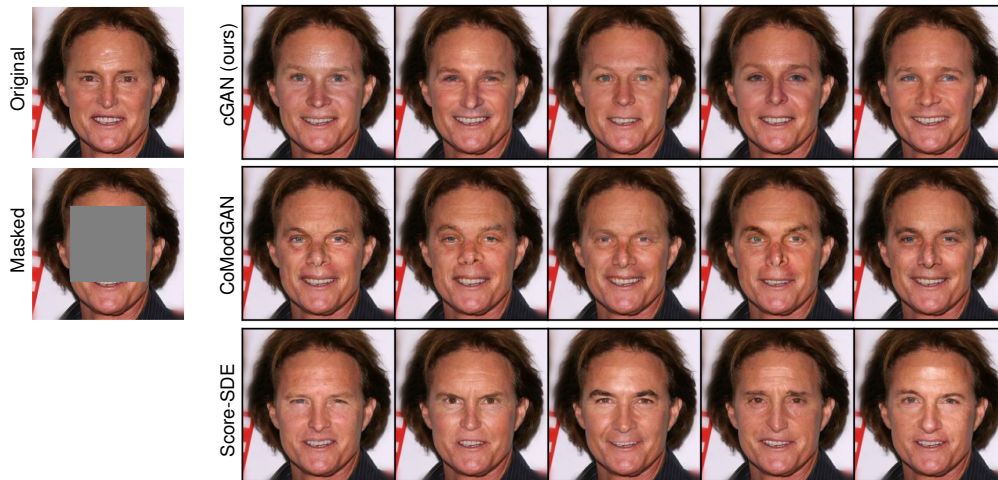


Figure K.18: Example of inpainting a  $128 \times 128$  square on a  $256 \times 256$  resolution CelebA-HQ image.