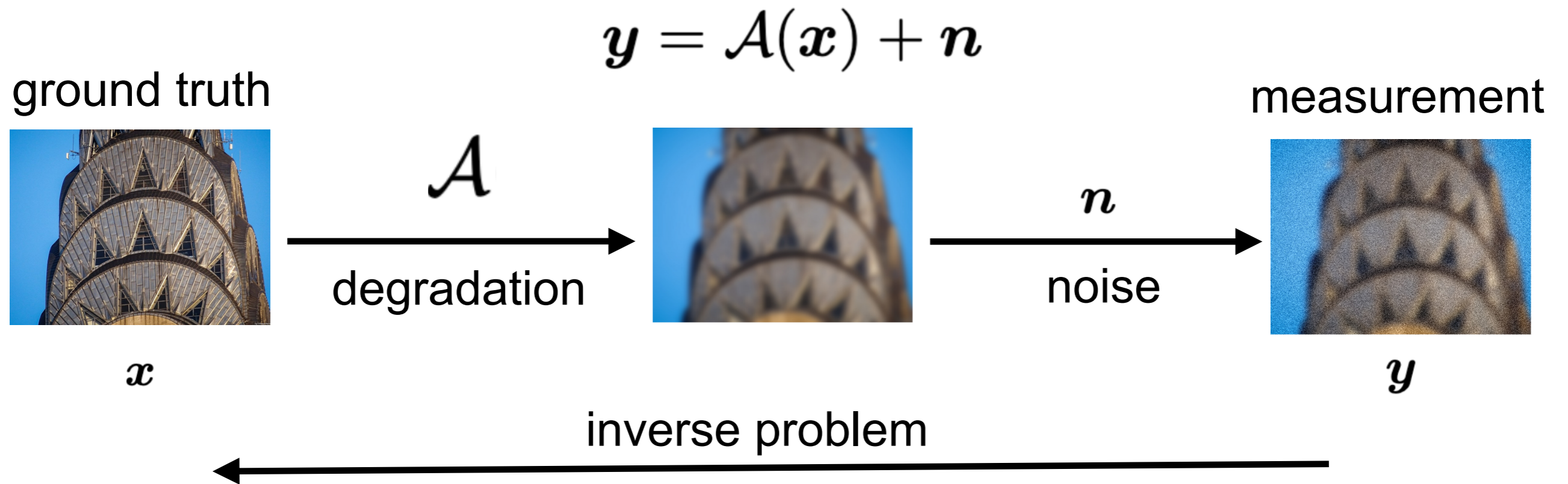


# ADAPT AND DIFFUSE: SAMPLE-ADAPTIVE RECONSTRUCTION VIA LATENT DIFFUSION MODELS

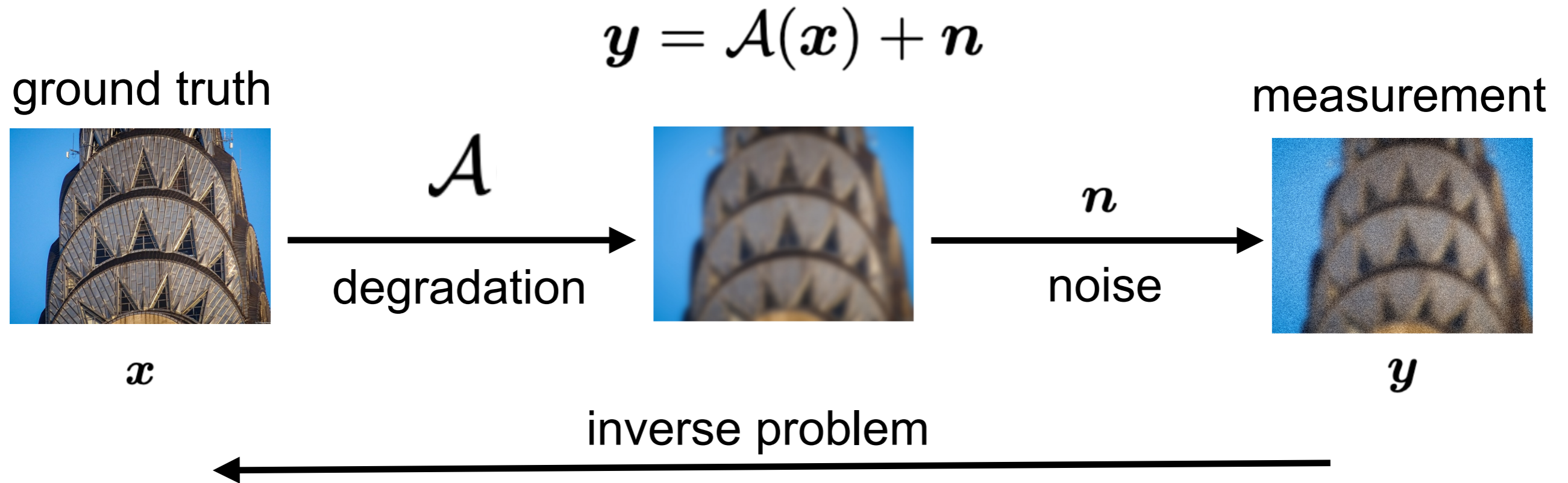
Zalan Fabian\*, Berk Tinaz\*, Mahdi Soltanolkotabi



# Inverse Problems



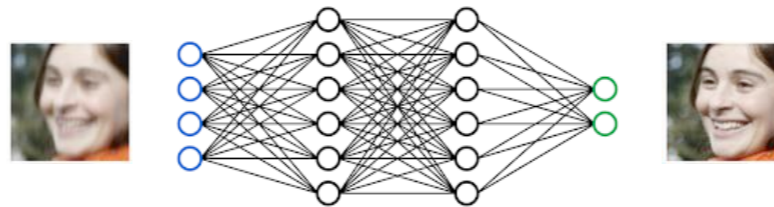
# Inverse Problems



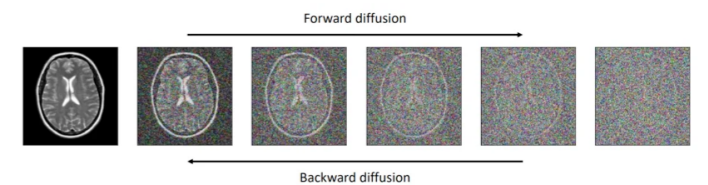
- Solving inverse problems

$$\hat{x} = \arg \min_x \|\mathcal{A}(x) - y\|^2 + \mathcal{R}(x)$$

Classical approaches



End-to-end deep learning



Diffusion solvers

# Reconstruction difficulty

- **Sample-by-sample variation in reconstruction difficulty**

Degradation level



# Reconstruction difficulty

- **Sample-by-sample variation in reconstruction difficulty**

Degradation level



Target



Noise level



Target



# Reconstruction difficulty

- Sample-by-sample variation in reconstruction **difficulty**

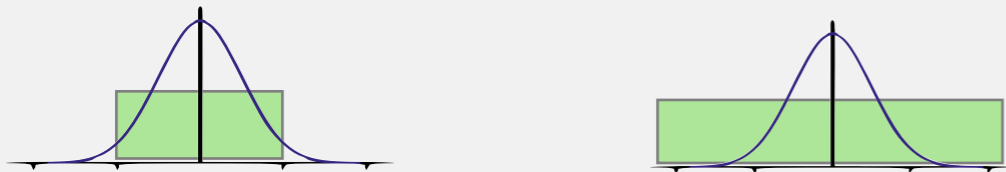
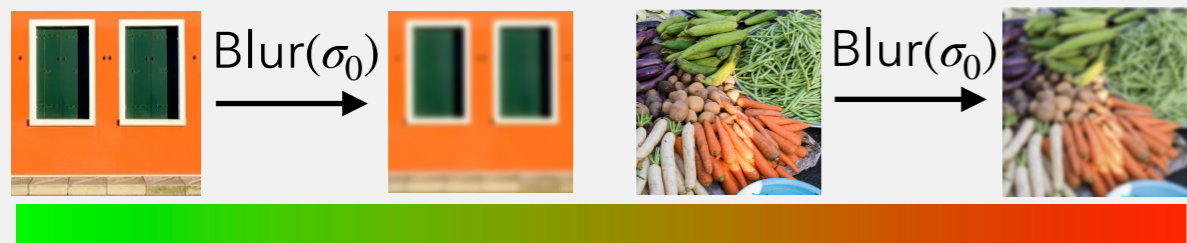
Degradation level



Noise level



Sample + degradation



# Reconstruction difficulty

- Sample-by-sample variation in reconstruction **difficulty**

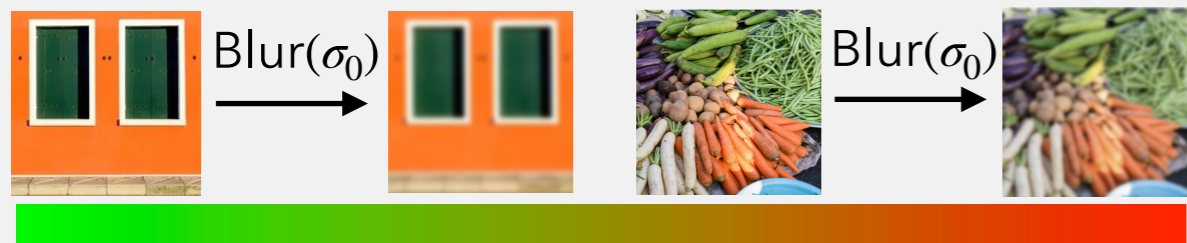
Degradation level



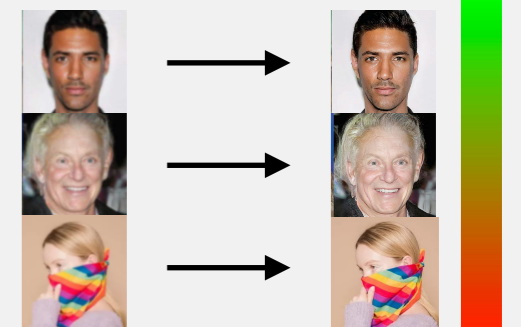
Noise level



Sample + degradation



Test time distribution shift

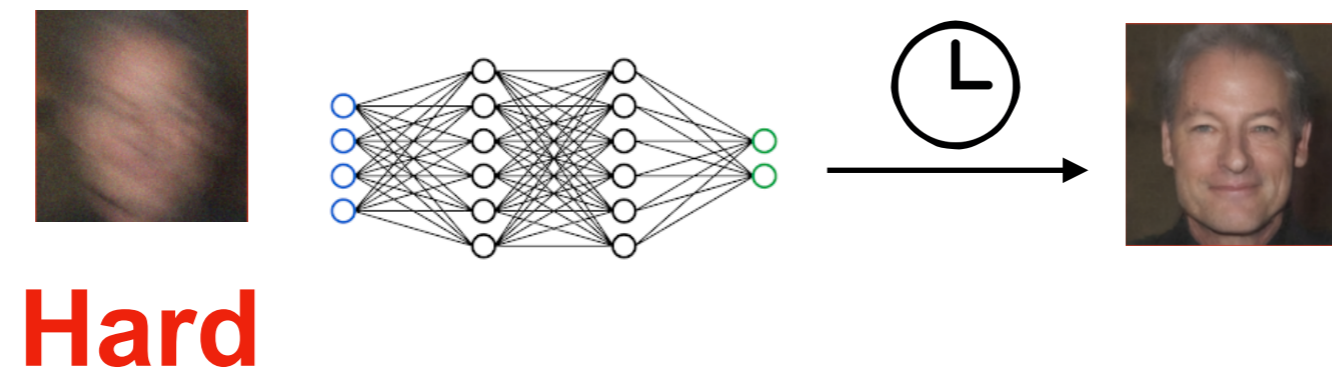
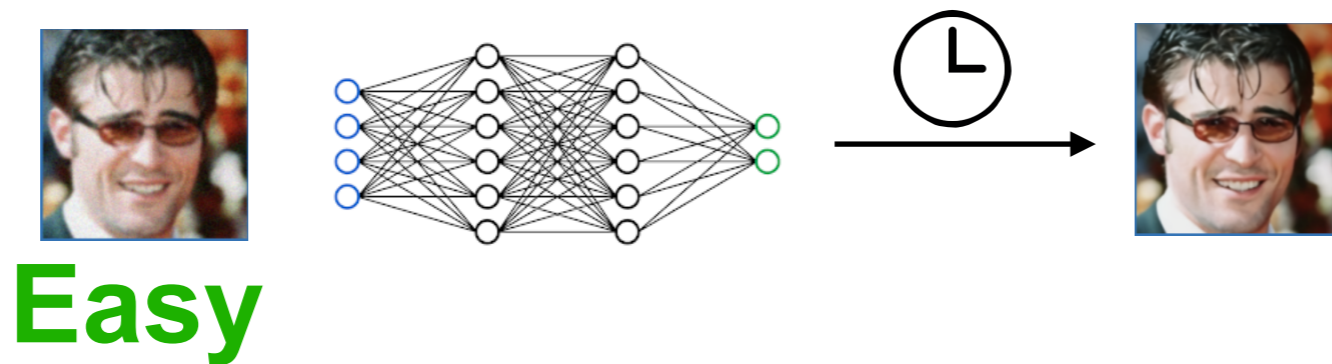


Training set

Test set

# Sample-adaptive compute

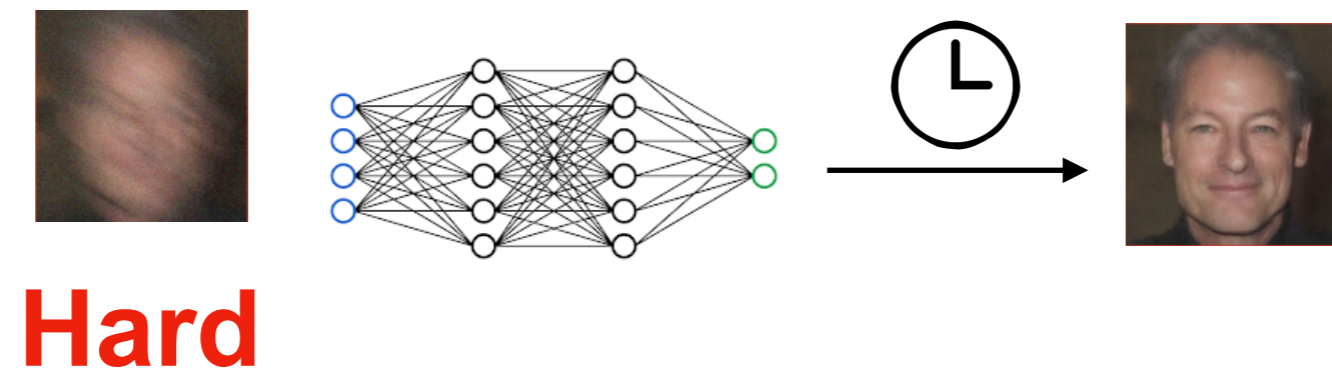
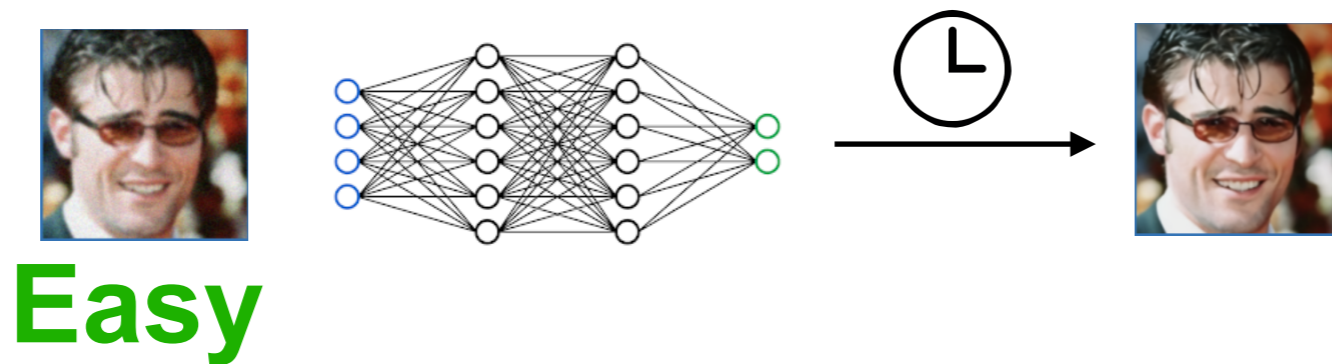
- Expending the same amount of resources to reconstruct any sample (easy or hard) is potentially wasteful.





# Sample-adaptive compute

- Expending the same amount of resources to reconstruct any sample (easy or hard) is potentially wasteful.



**Idea:** **adapt** the compute allocation based on the difficulty of the problem on a **sample-by-sample** basis in test time!

**Adapt...**

# Quantifying difficulty

## Image space



- arbitrarily high perturbation
- reconstruction is trivial

$$y = c \cdot x, \quad c \in \mathbb{R}^+$$

# Quantifying difficulty

## Image space



$$y = c \cdot x, \quad c \in \mathbb{R}^+$$

- arbitrarily high perturbation
- reconstruction is trivial



$$y = x + n, \quad n \sim \mathcal{N}(0, \sigma^2 I)$$

- small perturbation
- reconstruction is challenging

# Quantifying difficulty

## Image space



$$y = c \cdot x, \quad c \in \mathbb{R}^+$$

- arbitrarily high perturbation
- reconstruction is trivial



$$y = x + n, \quad n \sim \mathcal{N}(0, \sigma^2 I)$$

- small perturbation
- reconstruction is challenging

## Autoencoder latents

- compressed representation of relevant information in image

# Quantifying difficulty

## Image space



$$y = c \cdot x, \quad c \in \mathbb{R}^+$$

- arbitrarily high perturbation
- reconstruction is trivial



$$y = x + n, \quad n \sim \mathcal{N}(0, \sigma^2 I)$$

- small perturbation
- reconstruction is challenging

## Autoencoder latents

- compressed representation of relevant information in image
- natural space to quantify loss of information due to corruption

# Quantifying difficulty

## Image space



$$y = c \cdot x, \quad c \in \mathbb{R}^+$$

- arbitrarily high perturbation
- reconstruction is trivial



$$y = x + n, \quad n \sim \mathcal{N}(0, \sigma^2 I)$$

- small perturbation
- reconstruction is challenging

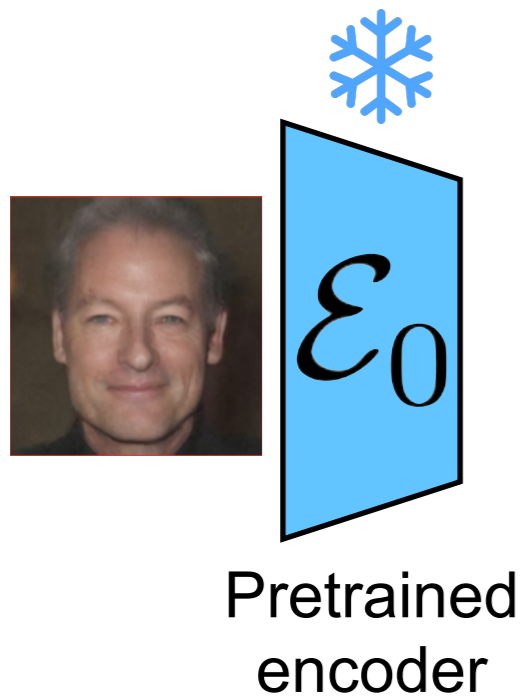
## Autoencoder latents

- compressed representation of relevant information in image
- natural space to quantify loss of information due to corruption

**Idea:** quantify severity of degradation **in the latent space** of an autoencoder

# Severity Encoding

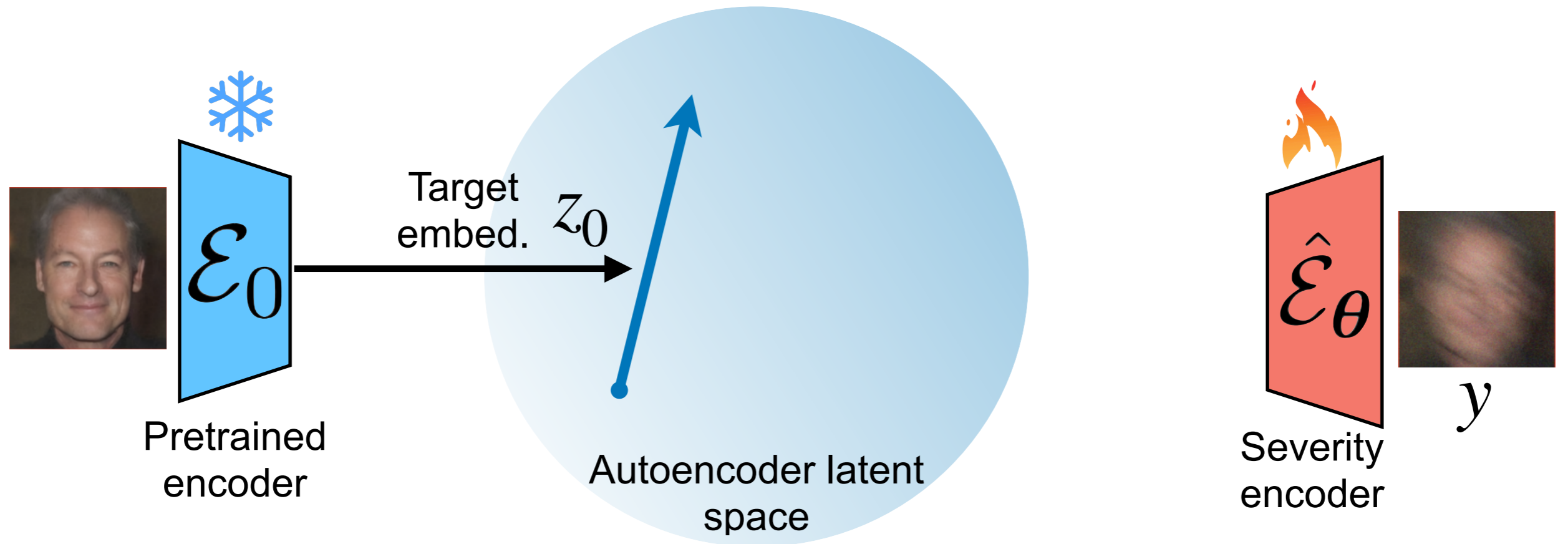
- Estimate degradation severity given corrupted image





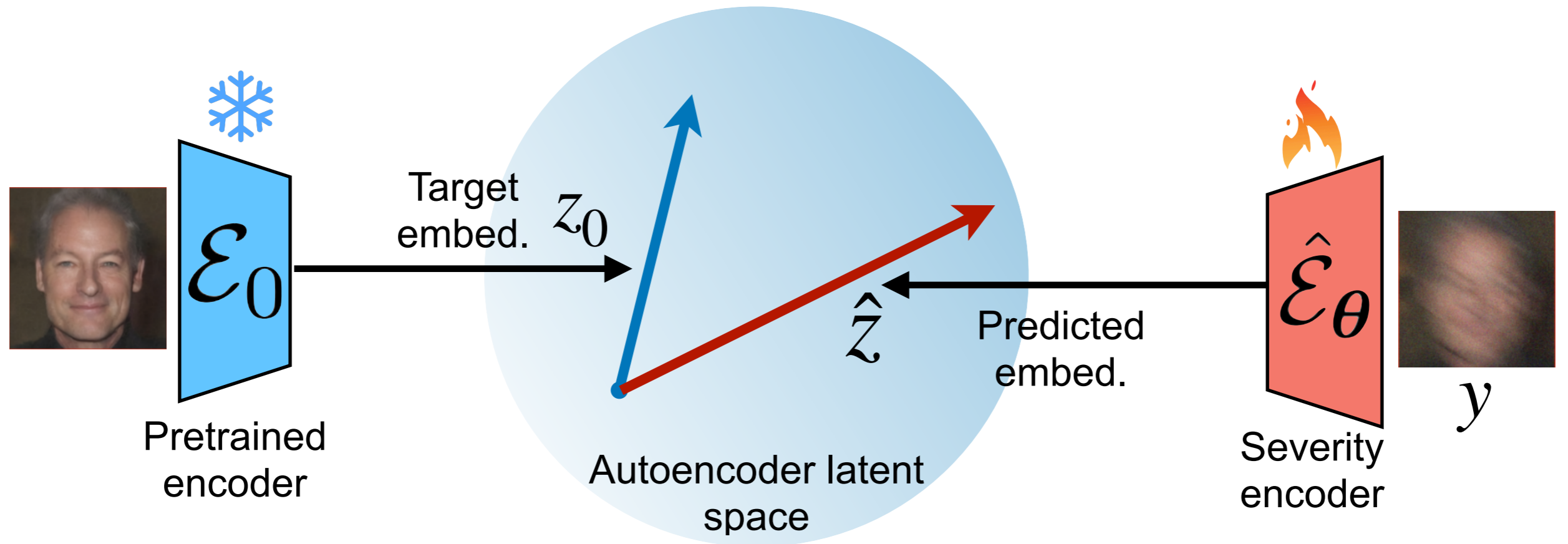
# Severity Encoding

- Estimate degradation severity given corrupted image



# Severity Encoding

- Estimate degradation severity given corrupted image

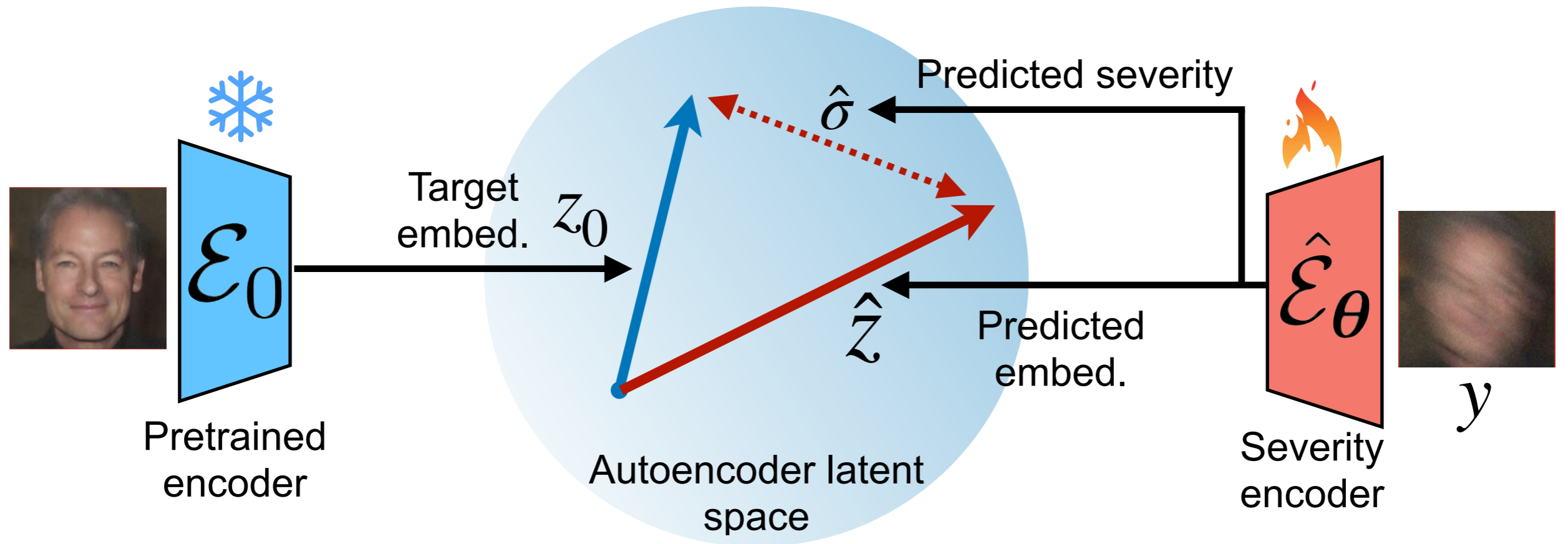


## Objectives:

1. Predict latent of clean image

# Severity Encoding

- Estimate degradation severity given corrupted image

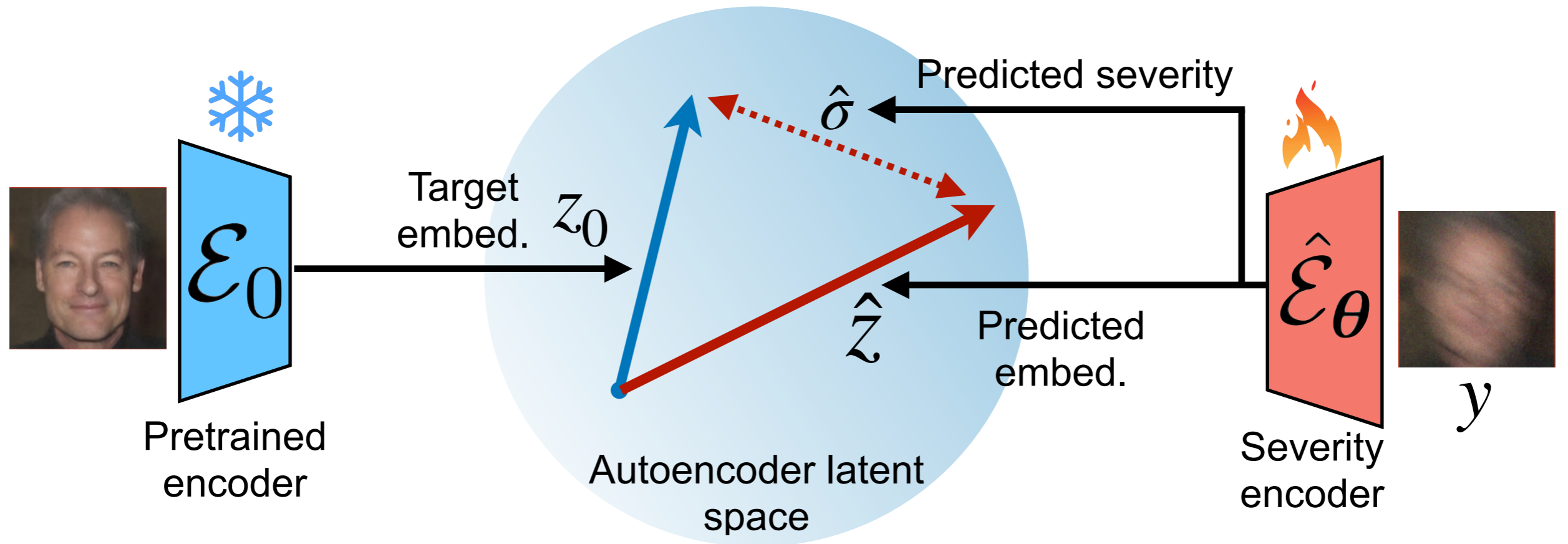


## Objectives:

1. Predict latent of clean image
2. Estimate magnitude of error

# Severity Encoding

- Estimate degradation severity given corrupted image



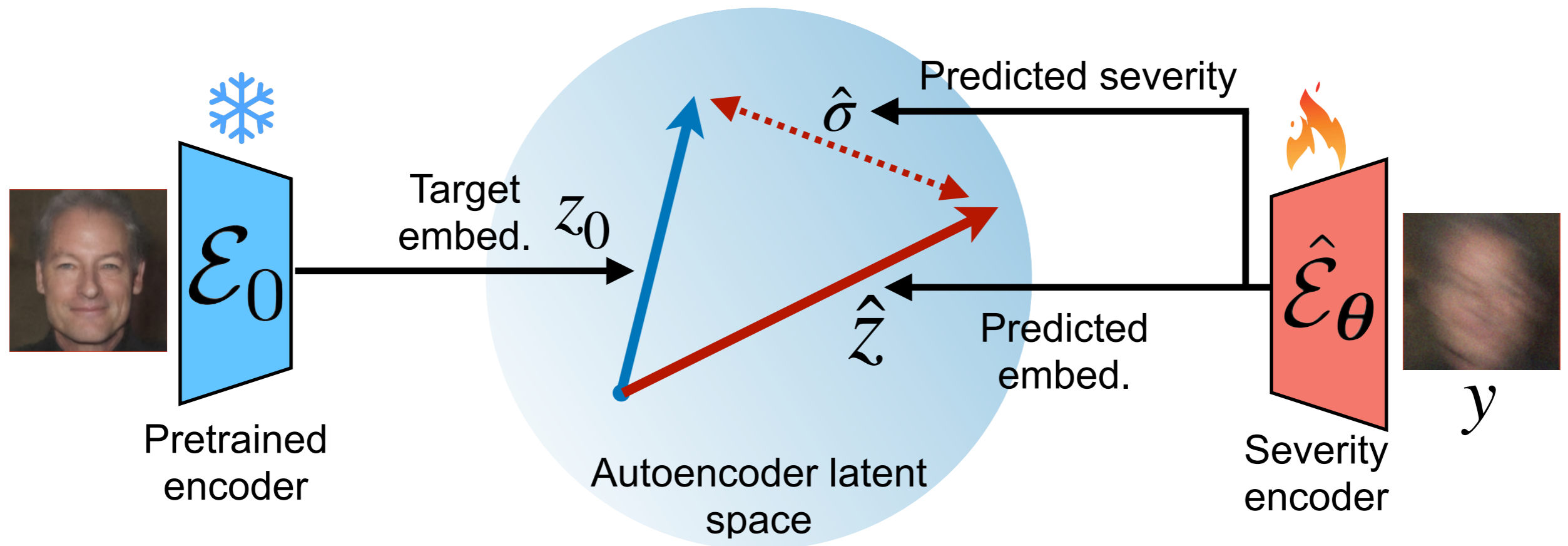
## Objectives:

1. Predict latent of clean image
2. Estimate magnitude of error

Intuition: the more degraded the input, the larger the prediction error

# Severity Encoding

- Estimate degradation severity given corrupted image



## Objectives:

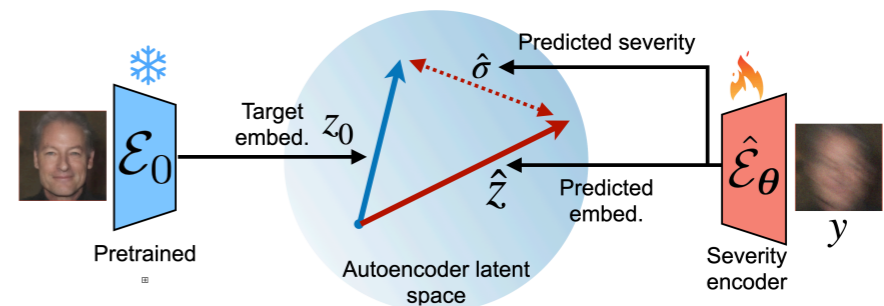
1. Predict latent of clean image
2. Estimate magnitude of error

Intuition: the more degraded the input, the larger the prediction error

We leverage **latent prediction error as a proxy** for degradation severity!

# Training objective

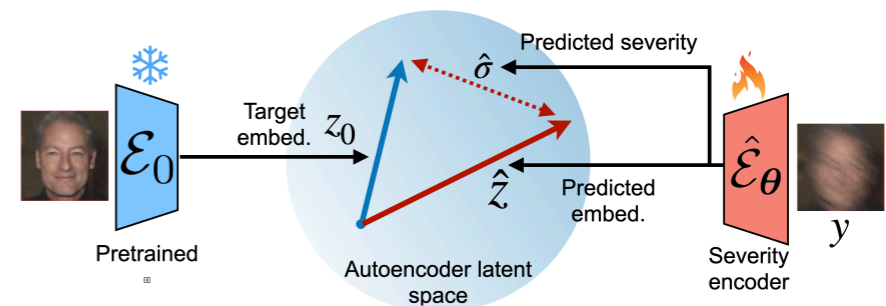
- We minimize the following loss:



$$\min_{\theta} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0), \mathbf{y} \sim \mathcal{N}(\mathcal{A}(\mathbf{x}_0), \sigma_y^2 \mathbf{I})} \left[ \underbrace{\left\| \mathbf{z}_0 - \hat{\mathbf{z}}(\mathbf{y}; \theta) \right\|^2}_{\text{reconstruction}} + \lambda_\sigma \underbrace{\left\| \bar{\sigma}(\mathbf{y}; \mathbf{z}_0) - \hat{\sigma}(\mathbf{y}; \theta) \right\|^2}_{\text{error prediction}} \right]$$

# Training objective

- We minimize the following loss:



$$\min_{\theta} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0), \mathbf{y} \sim \mathcal{N}(\mathcal{A}(\mathbf{x}_0), \sigma_y^2 \mathbf{I})} \left[ \underbrace{\left\| \mathbf{z}_0 - \hat{\mathbf{z}}(\mathbf{y}; \theta) \right\|^2}_{\text{reconstruction}} + \lambda_{\sigma} \underbrace{\left\| \bar{\sigma}(\mathbf{y}; \mathbf{z}_0) - \hat{\sigma}(\mathbf{y}; \theta) \right\|^2}_{\text{error prediction}} \right]$$

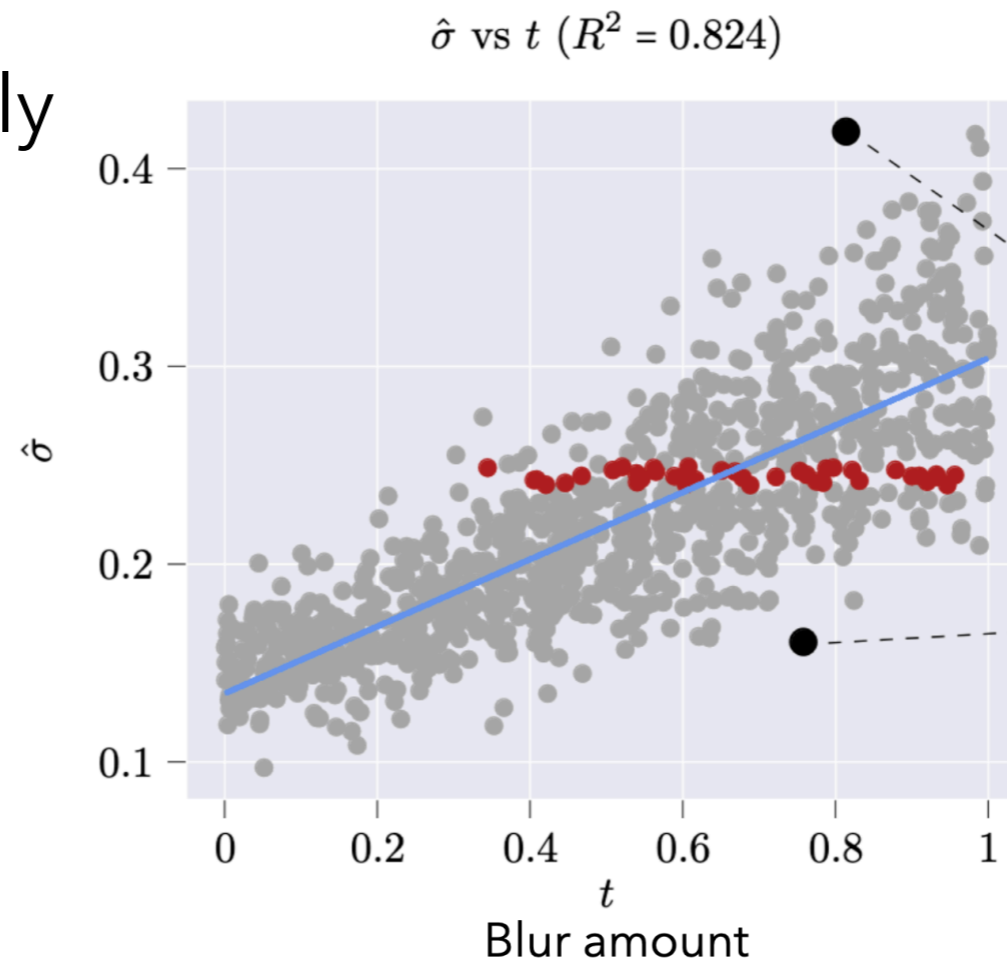
Assumption: prediction error is zero-mean i.i.d. Gaussian:

$$\mathbf{e}(\mathbf{y}) = \hat{\mathbf{z}} - \mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, \sigma_*^2(\mathbf{y}) \mathbf{I})$$

$$\bar{\sigma}^2(\mathbf{y}, \mathbf{z}_0) = \frac{1}{d-1} \sum_{i=1}^d (\mathbf{e}^{(i)} - \frac{1}{d} \sum_{j=1}^d \mathbf{e}^{(j)})^2$$

# Severity experiments

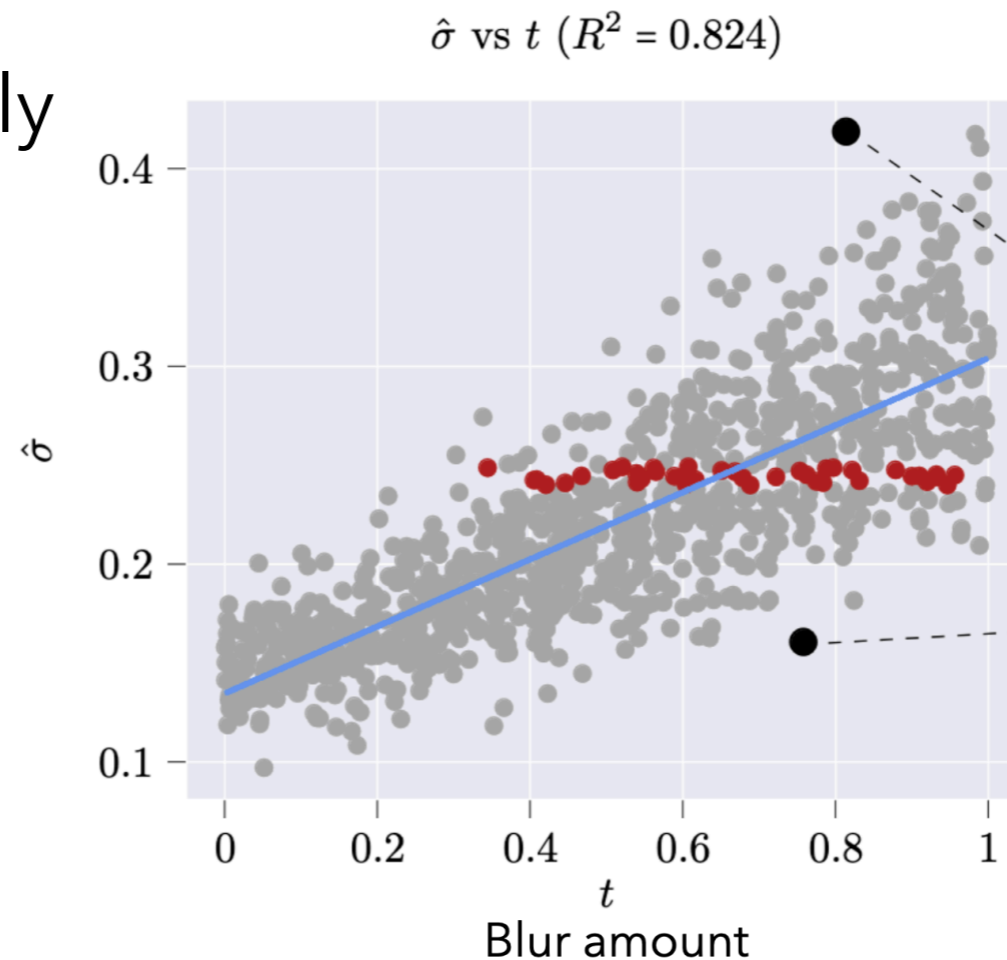
- Predicted severity strongly correlates with blur level





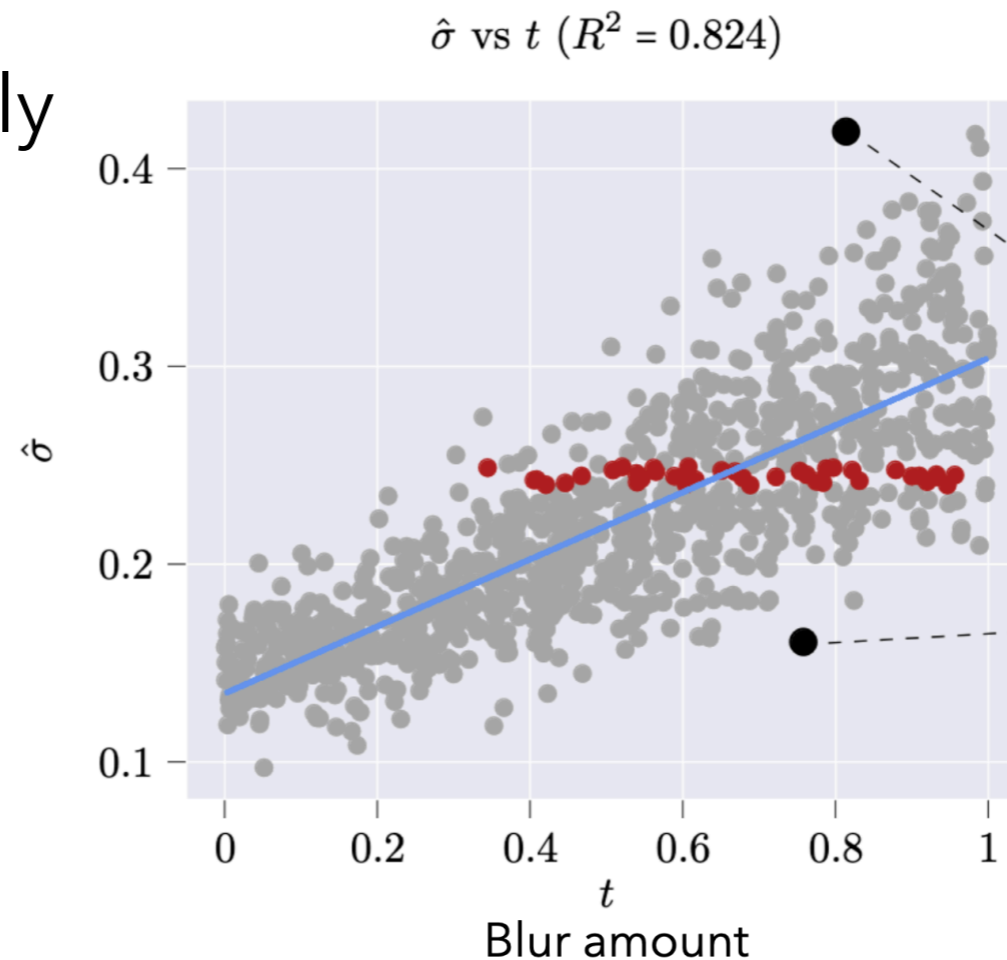
# Severity experiments

- Predicted severity strongly correlates with blur level
- Outliers indicate the presence of additional contributing factors

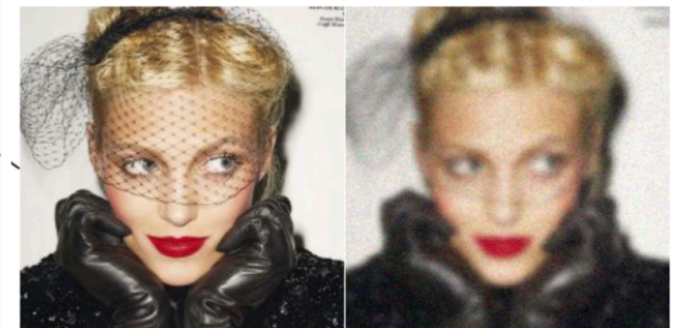


# Severity experiments

- Predicted severity strongly correlates with blur level
- Outliers indicate the presence of additional contributing factors

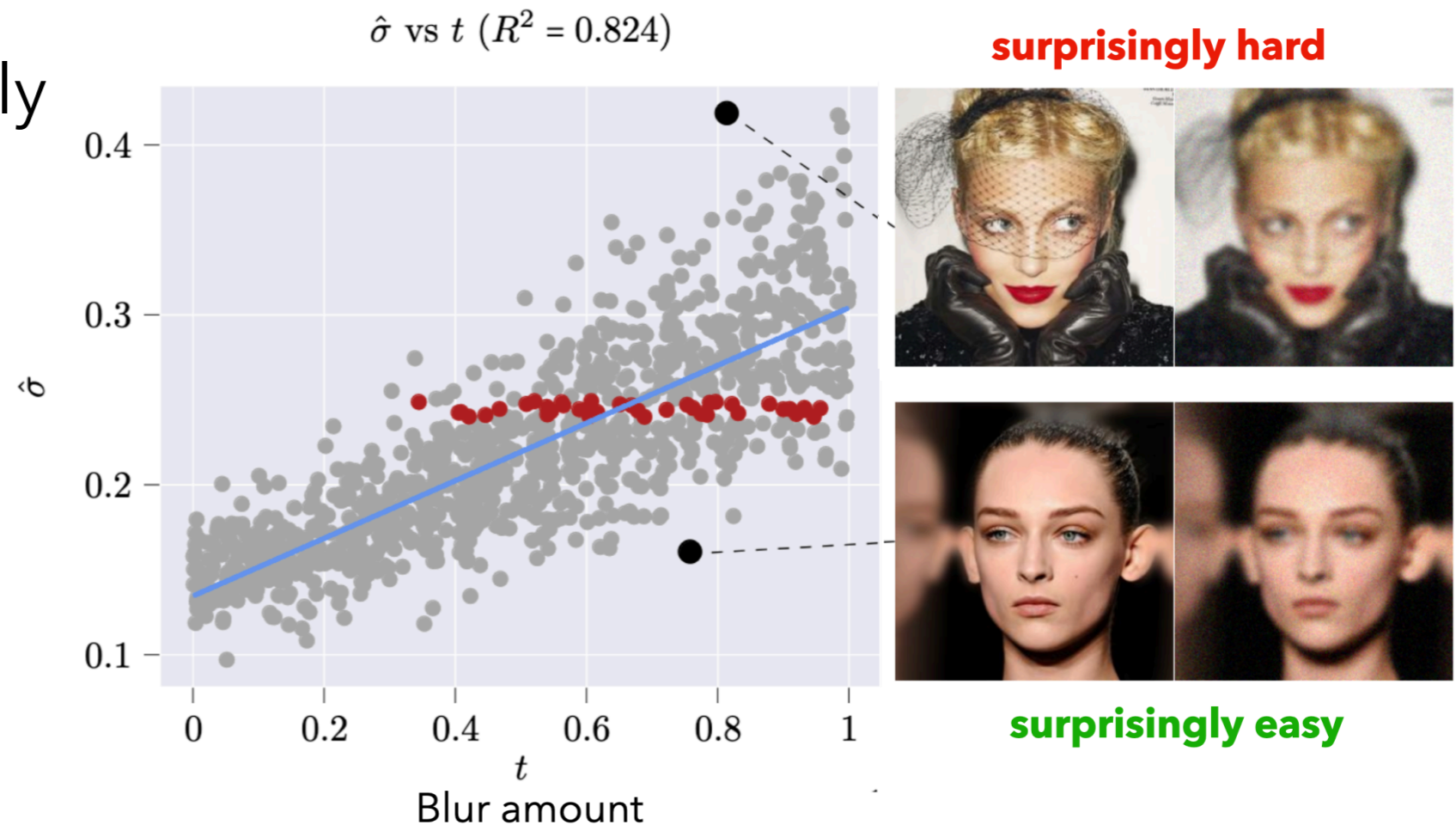


**surprisingly hard**



# Severity experiments

- Predicted severity strongly correlates with blur level
- Outliers indicate the presence of additional contributing factors



... and diffuse

# Diffusion-based Inverse Problem Solving

- Bayesian framework

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{n}$$

*solve inverse problem = sample from posterior*

$$p_{\theta}(\mathbf{x}|\mathbf{y}) \propto p_{\theta}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$$

$$\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}|\mathbf{y}) = \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y}|\mathbf{x})$$

# Diffusion-based Inverse Problem Solving




- Bayesian framework

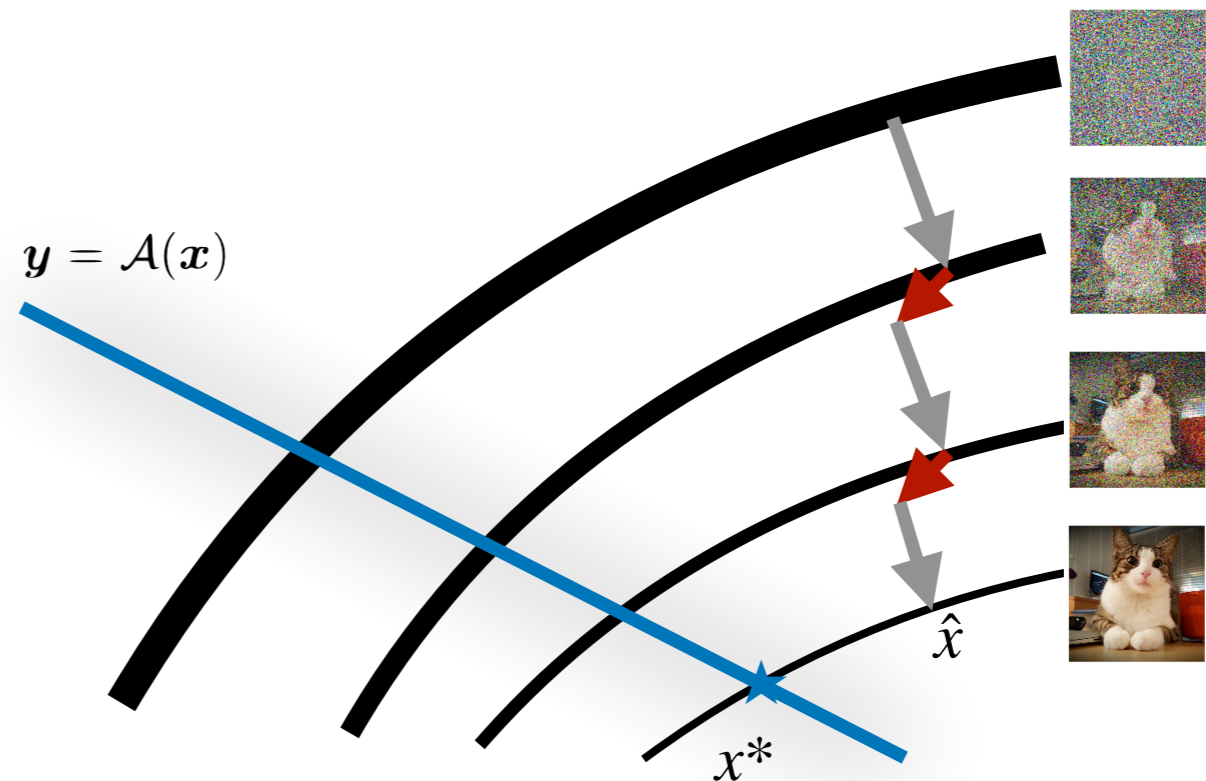
$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{n}$$

*solve inverse problem = sample from posterior*

$$p_{\theta}(\mathbf{x}|\mathbf{y}) \propto p_{\theta}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$$

$$\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}|\mathbf{y}) = \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y}|\mathbf{x})$$

Data-consistency update   
Reverse diffusion   
Noisy image manifolds 



# Flash-Diffusion



Hard example

Reconstruction in 200 steps



Easy example

Reconstruction in 50 steps

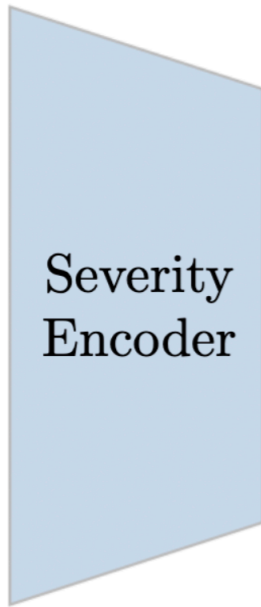


# Flash-Diffusion

Reconstruction in 200 steps



Hard example



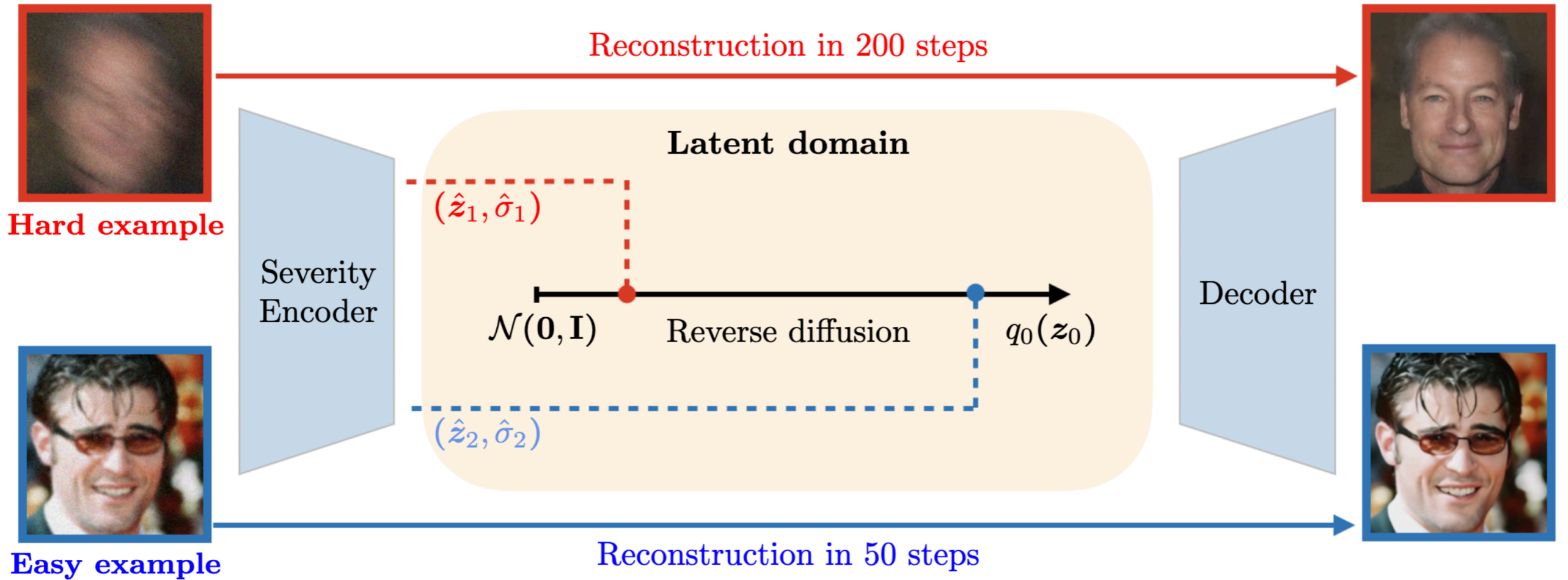
Easy example

Reconstruction in 50 steps

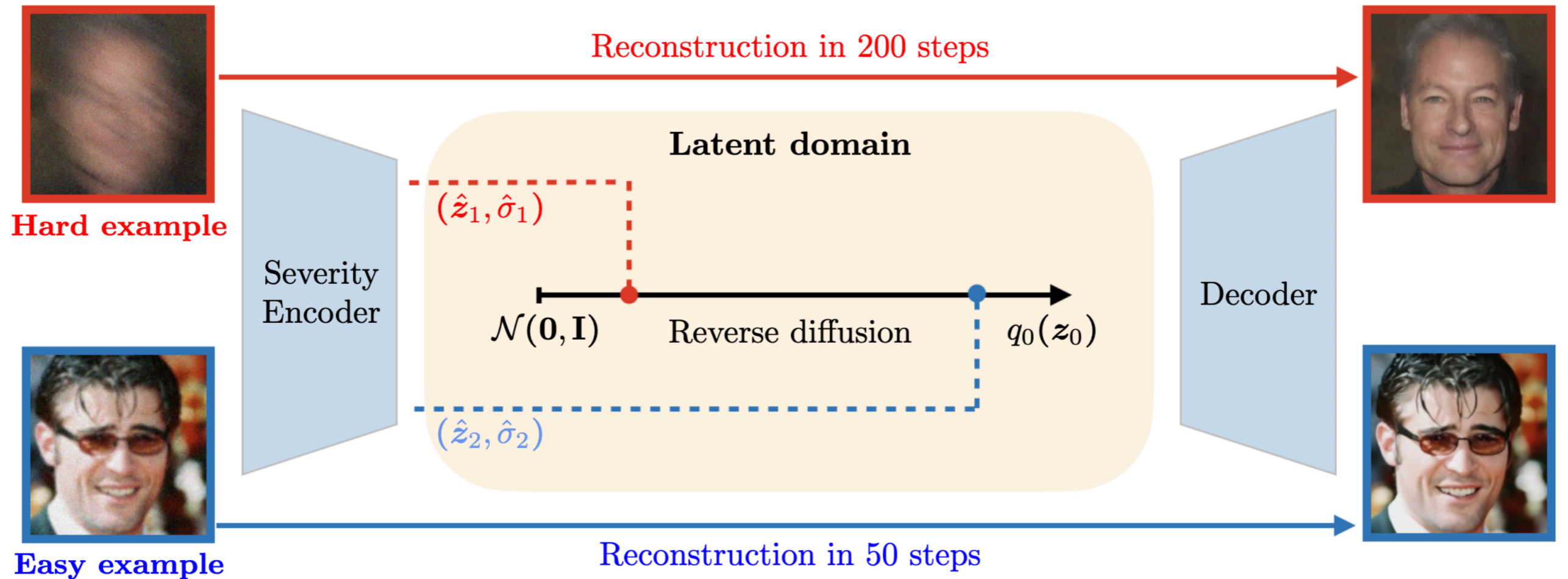




# Flash-Diffusion



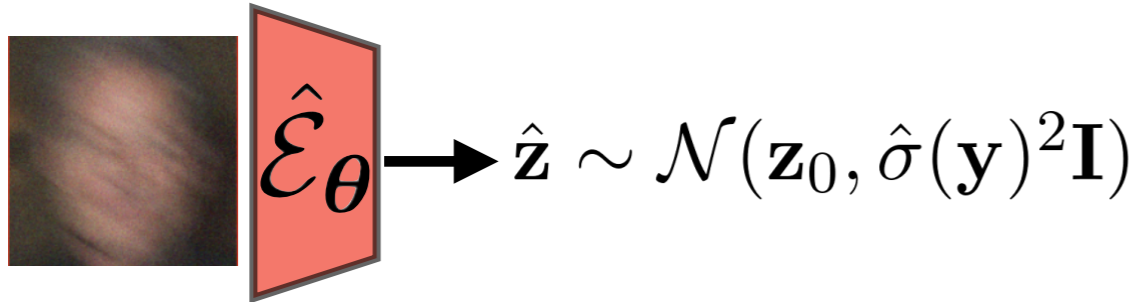
# Flash-Diffusion



Flash-Diffusion **acts as a wrapper** around *any* baseline latent diffusion solver, imbuing it with sample-adaptivity.

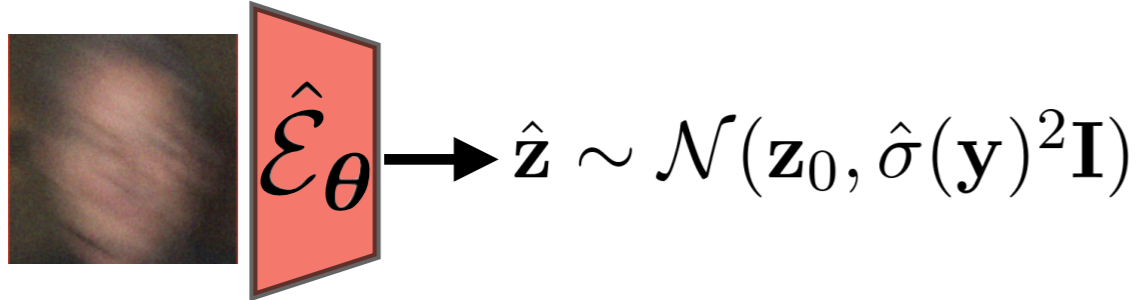
# Finding optimal starting time

Severity encoder



# Finding optimal starting time

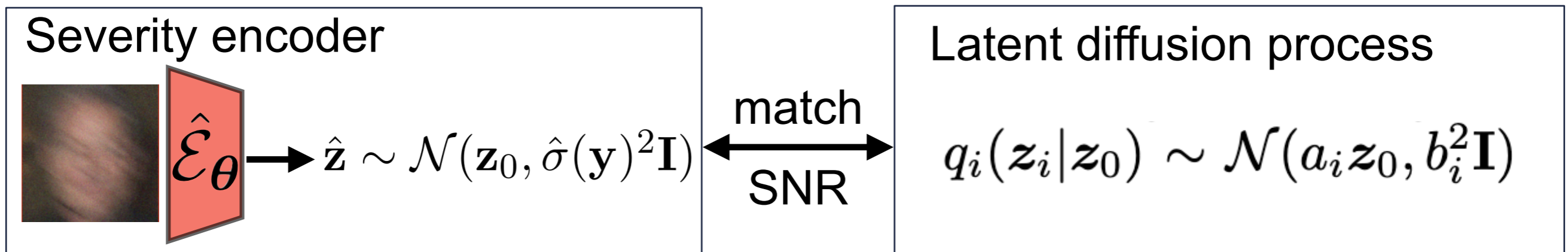
Severity encoder



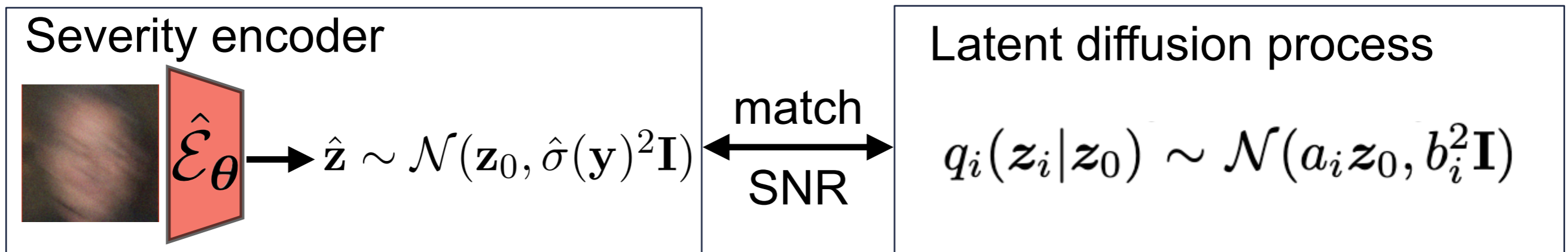
Latent diffusion process

$$q_i(\mathbf{z}_i | \mathbf{z}_0) \sim \mathcal{N}(a_i \mathbf{z}_0, b_i^2 \mathbf{I})$$

# Finding optimal starting time



# Finding optimal starting time



Adaptive starting time: 
$$i_{start}(\mathbf{y}) = \arg \min_{i \in [1, 2, \dots, N]} \left| \frac{1}{\hat{\sigma}(\mathbf{y})^2} - \frac{a_i^2}{b_i^2} \right|$$

# Experiments

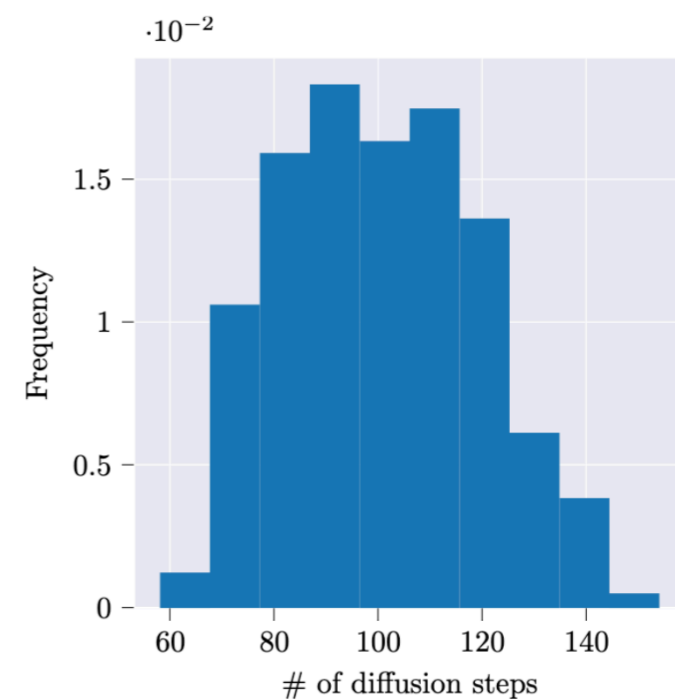
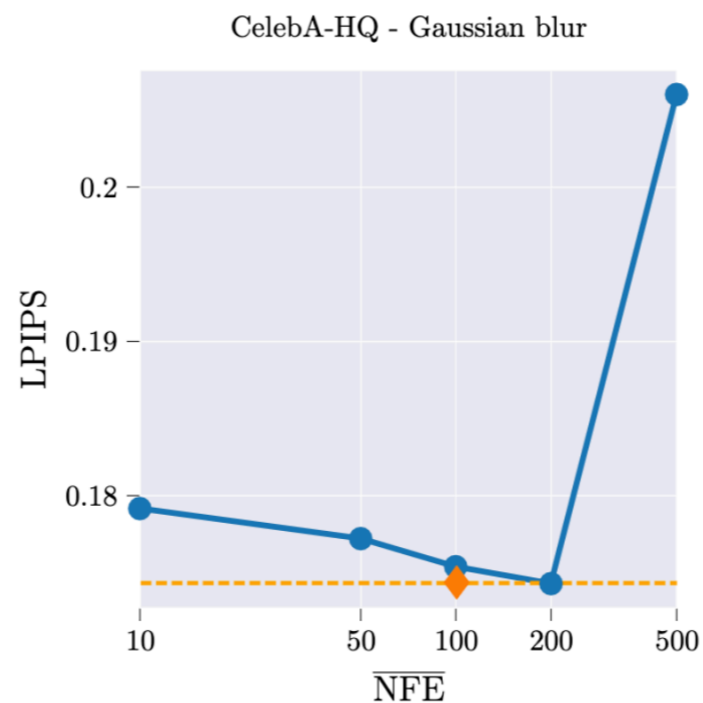
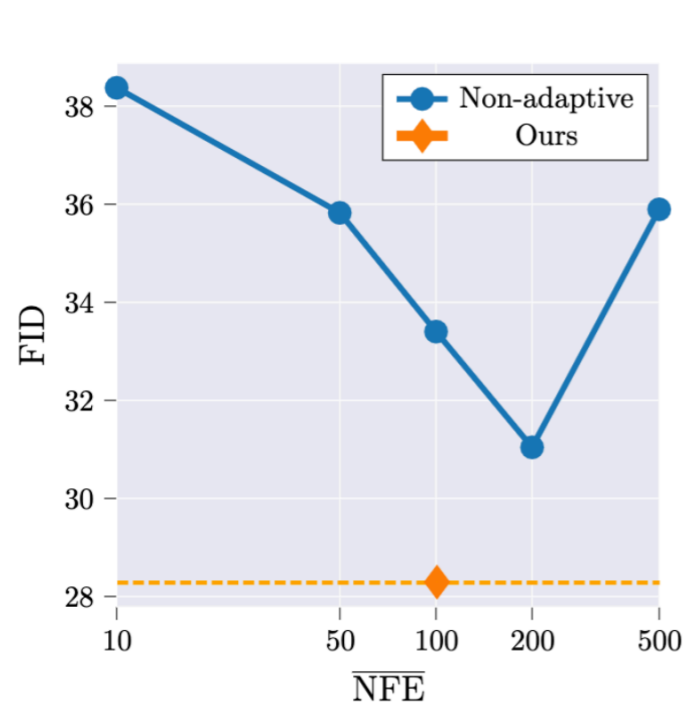
- Comparison with baseline solvers

FFHQ Method	Gaussian Deblurring (Varying)					Gaussian Deblurring (Fixed)					Nonlinear Deblurring					Random Inpainting				
	PSNR(↑)	SSIM(↑)	LPIPS(↓)	FID(↓)	NFE	PSNR(↑)	SSIM(↑)	LPIPS(↓)	FID(↓)	NFE	PSNR(↑)	SSIM(↑)	LPIPS(↓)	FID(↓)	NFE	PSNR(↑)	SSIM(↑)	LPIPS(↓)	FID(↓)	NFE
Latent-DPS	23.69	0.6418	0.3579	87.26	1000	22.88	0.6136	0.3690	89.38	1000	22.07	0.5974	0.3814	90.89	1000	23.96	0.6566	0.3666	93.65	1000
Flash(Latent-DPS)	<u>29.17</u>	<u>0.8182</u>	<u>0.2240</u>	<u>55.57</u>	100.3	<u>27.44</u>	<u>0.7691</u>	<u>0.2823</u>	<u>80.44</u>	127.7	<u>27.17</u>	<u>0.7659</u>	<u>0.2695</u>	<u>69.78</u>	136.1	<u>29.21</u>	<u>0.8414</u>	<b>0.1945</b>	<b>53.95</b>	104.7
PSLD (Rout et al., 2024)	25.06	0.6769	0.3194	79.79	1000	23.72	0.6183	0.3324	88.45	1000	-	-	-	-	-	24.94	0.6617	0.3672	85.64	1000
Flash(PSLD)	<u>29.26</u>	<u>0.8205</u>	<b>0.2203</b>	<b>53.27</b>	100.3	<u>27.44</u>	<u>0.7657</u>	<b>0.2797</b>	<b>65.35</b>	127.7	-	-	-	-	-	<u>27.06</u>	<u>0.8018</u>	<u>0.2185</u>	<u>55.12</u>	104.7
GML-DPS (Rout et al., 2024)	24.98	0.6884	0.3471	100.27	1000	24.01	0.6574	0.3621	102.80	1000	23.00	0.6426	0.3812	108.79	1000	25.20	0.7044	0.3527	103.3	1000
Flash(GML-DPS)	<u>29.21</u>	<u>0.8276</u>	<u>0.2274</u>	<u>69.16</u>	100.3	<u>27.47</u>	<u>0.7699</u>	<u>0.2816</u>	<u>69.81</u>	127.7	<u>27.11</u>	<u>0.7640</u>	<u>0.2756</u>	<u>81.93</u>	136.1	<u>28.95</u>	<u>0.8437</u>	<u>0.1957</u>	<u>59.39</u>	104.7
ReSample (Song et al., 2023)	28.77	0.8219	0.2587	81.96	500	27.62	0.7789	0.3148	102.47	500	26.61	0.7318	0.2838	68.57	500	27.51	0.7892	0.2460	63.39	500
Flash(ReSample)	<u>29.07</u>	<u>0.8330</u>	<u>0.2383</u>	<u>74.76</u>	49.9	<u>27.77</u>	<u>0.7845</u>	<u>0.3092</u>	<u>100.84</u>	63.6	<u>26.88</u>	<u>0.7660</u>	<b>0.2667</b>	<b>64.57</b>	67.8	<u>28.13</u>	<u>0.8260</u>	<u>0.2045</u>	<u>56.67</u>	52.1
AE	29.46	0.8358	0.2671	89.29	-	27.69	0.7820	0.3396	110.56	-	27.17	0.7786	0.3364	111.24	-	29.23	0.8432	0.2515	85.87	-
SwinIR (Liang et al., 2021)	<b>30.69</b>	<b>0.8583</b>	0.2409	87.61	-	<b>28.41</b>	<b>0.8021</b>	0.3091	108.49	-	<b>27.60</b>	<b>0.7928</b>	0.3093	99.56	-	<b>30.08</b>	<b>0.8654</b>	0.2223	78.32	-
DPS (Chung et al., 2022a)	28.34	0.7791	0.2465	81.70	1000	25.49	0.6829	0.3035	97.89	1000	22.77	0.6191	0.3601	109.58	1000	28.30	0.8049	0.2451	82.78	1000

Flash-Diffusion **accelerates the baseline** solver by a factor of up to 10x on average and greatly **improves reconstruction quality**.

# Experiments

- **Adaptivity**

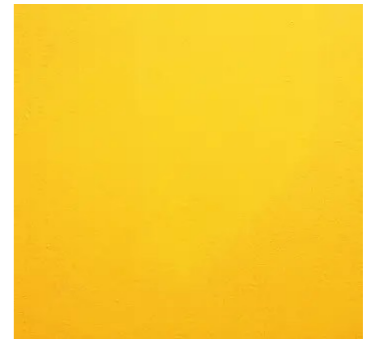
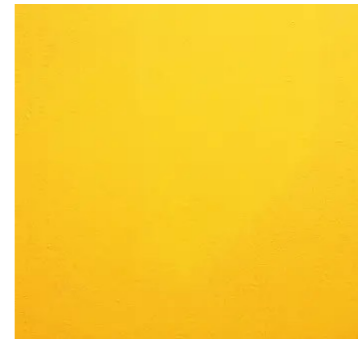


FlashDiffusion achieves best perceptual quality compared to any non-adaptive starting time.



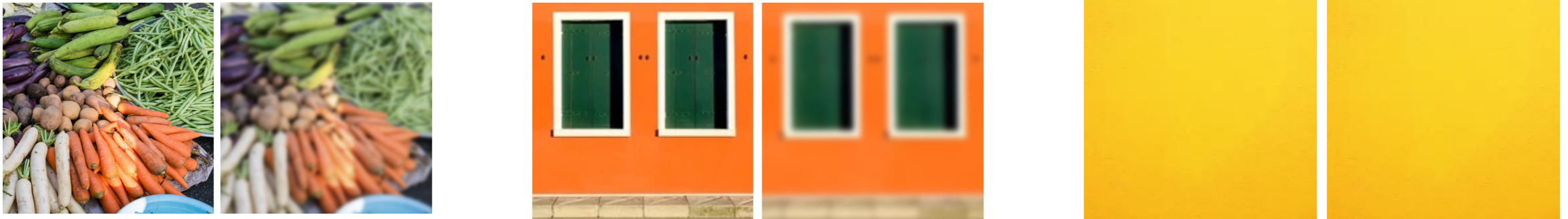
# Conclusion

1. Difficulty of a reconstruction problem may vary greatly on a **sample-by-sample** basis.

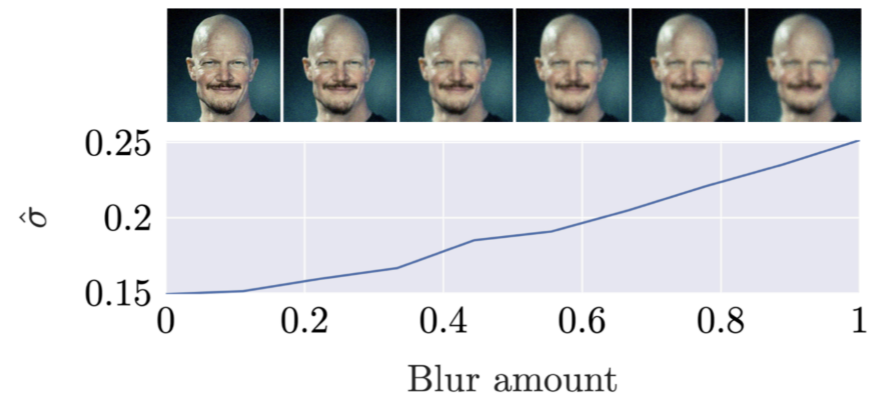


# Conclusion

1. Difficulty of a reconstruction problem may vary greatly on a **sample-by-sample** basis.

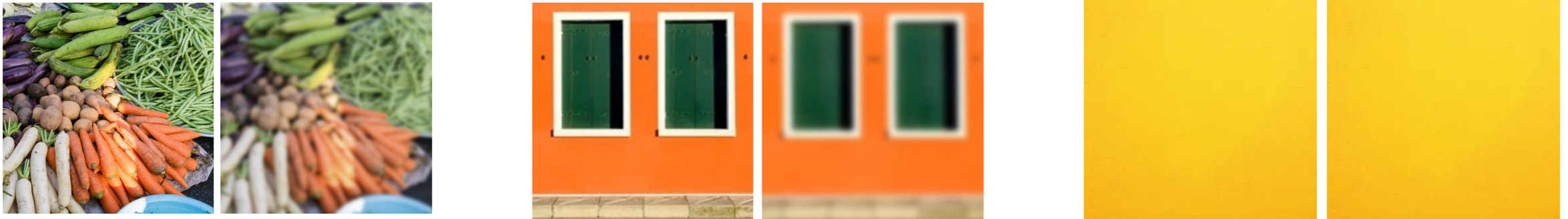


2. We propose **Severity Encoding** to estimate degradation severity in test-time.

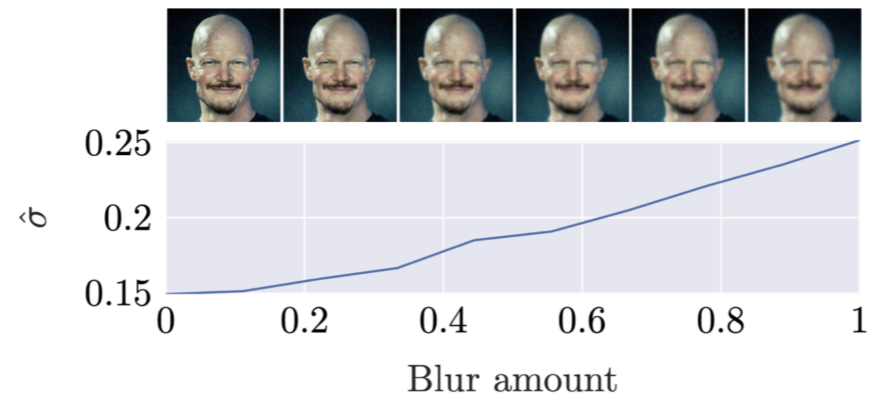


# Conclusion

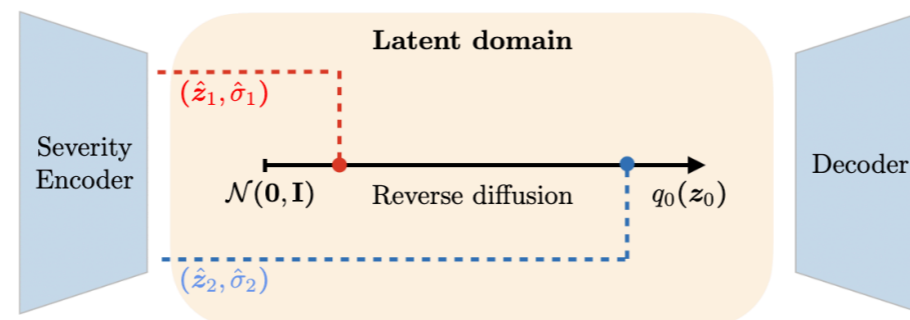
1. Difficulty of a reconstruction problem may vary greatly on a **sample-by-sample** basis.



2. We propose **Severity Encoding** to estimate degradation severity in test-time.



3. Flash-Diffusion **automatically scales reconstruction effort** with degradation severity via latent diffusion.



# Thank you for your attention!

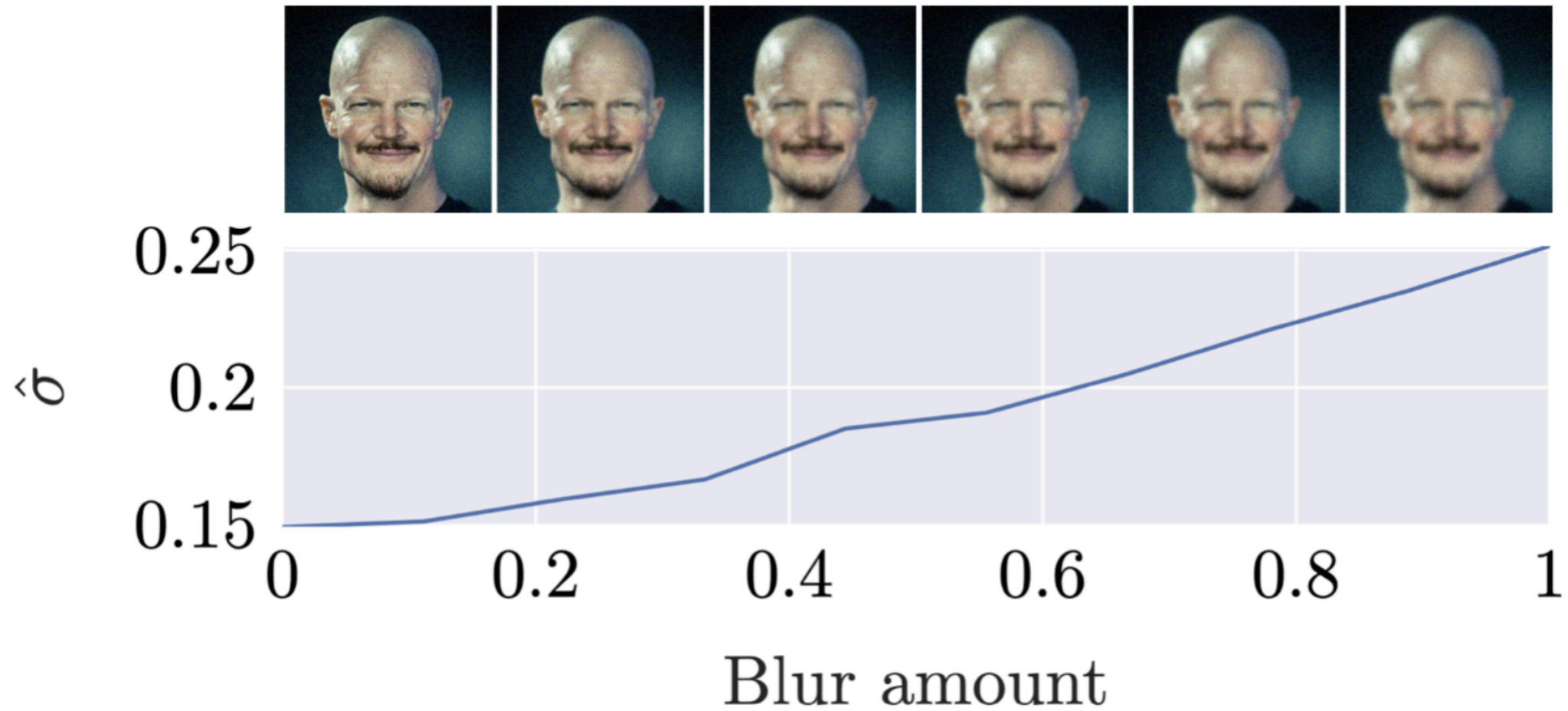


Paper



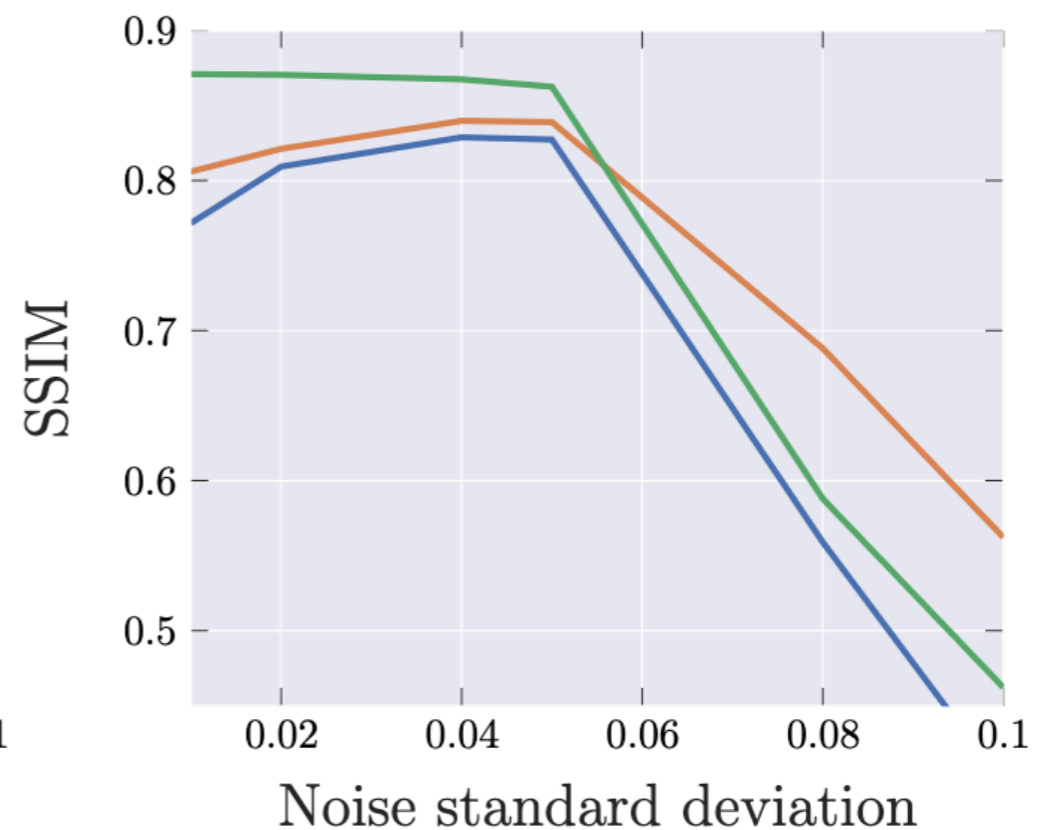
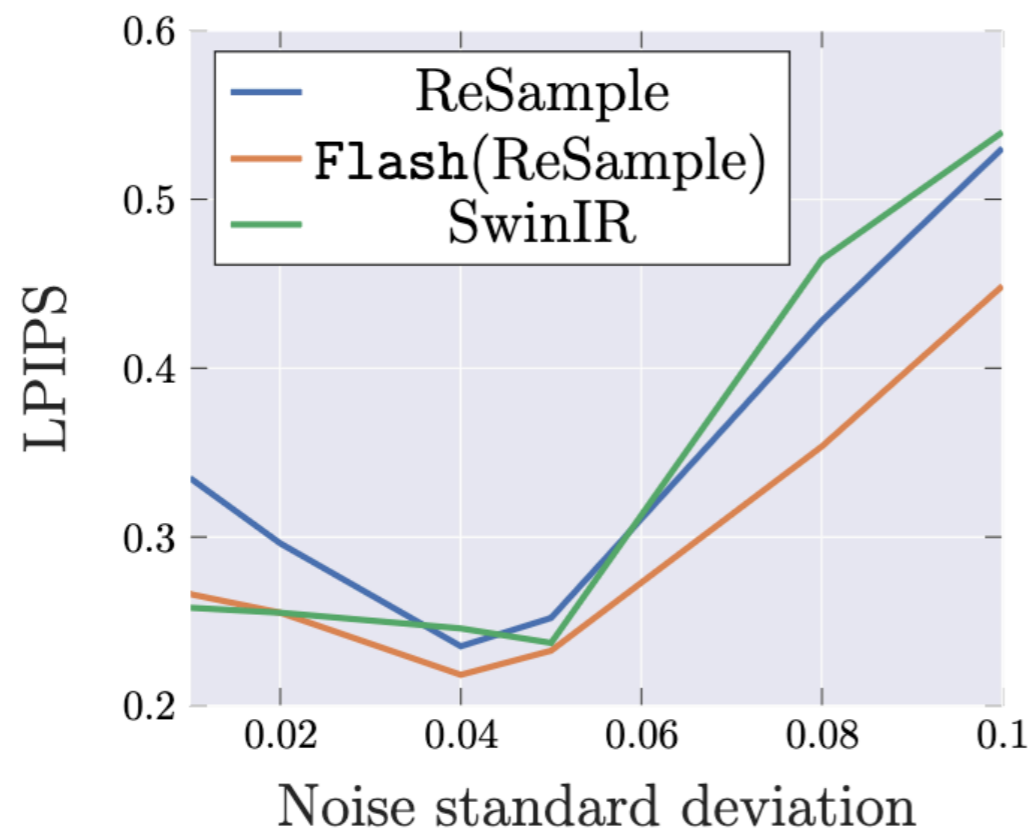
Code

# Severity experiments



# Experiments

- **Robustness**



Flash-Diffusion performance degrades more gracefully than the baseline solver's performance.