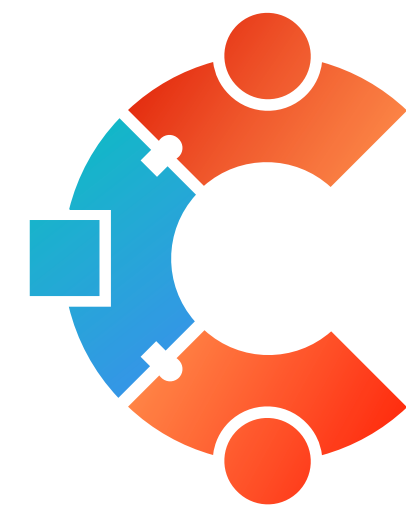


# Selective Amnesia:

## *A Continual Learning Approach to Forgetting in Deep Generative Models*

*NeurIPS 2023 Spotlight*

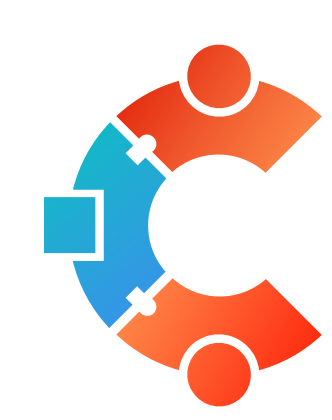
**Alvin Heng and Harold Soh**



Collaborative,  
**L**earning, and  
Adaptive  
**R**obots



School of Computing



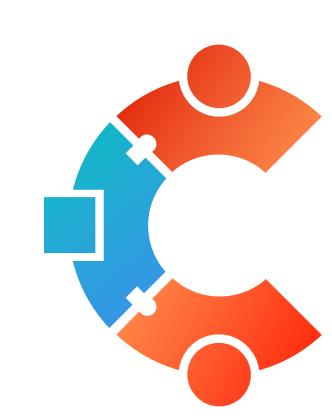
# Deep Generative Models are Awesome!

- Deep generative models can now generate photorealistic images (DALL-E 3, Stable Diffusion XL etc)



Image source: <https://stability.ai/blog/sdxl-09-stable-diffusion>





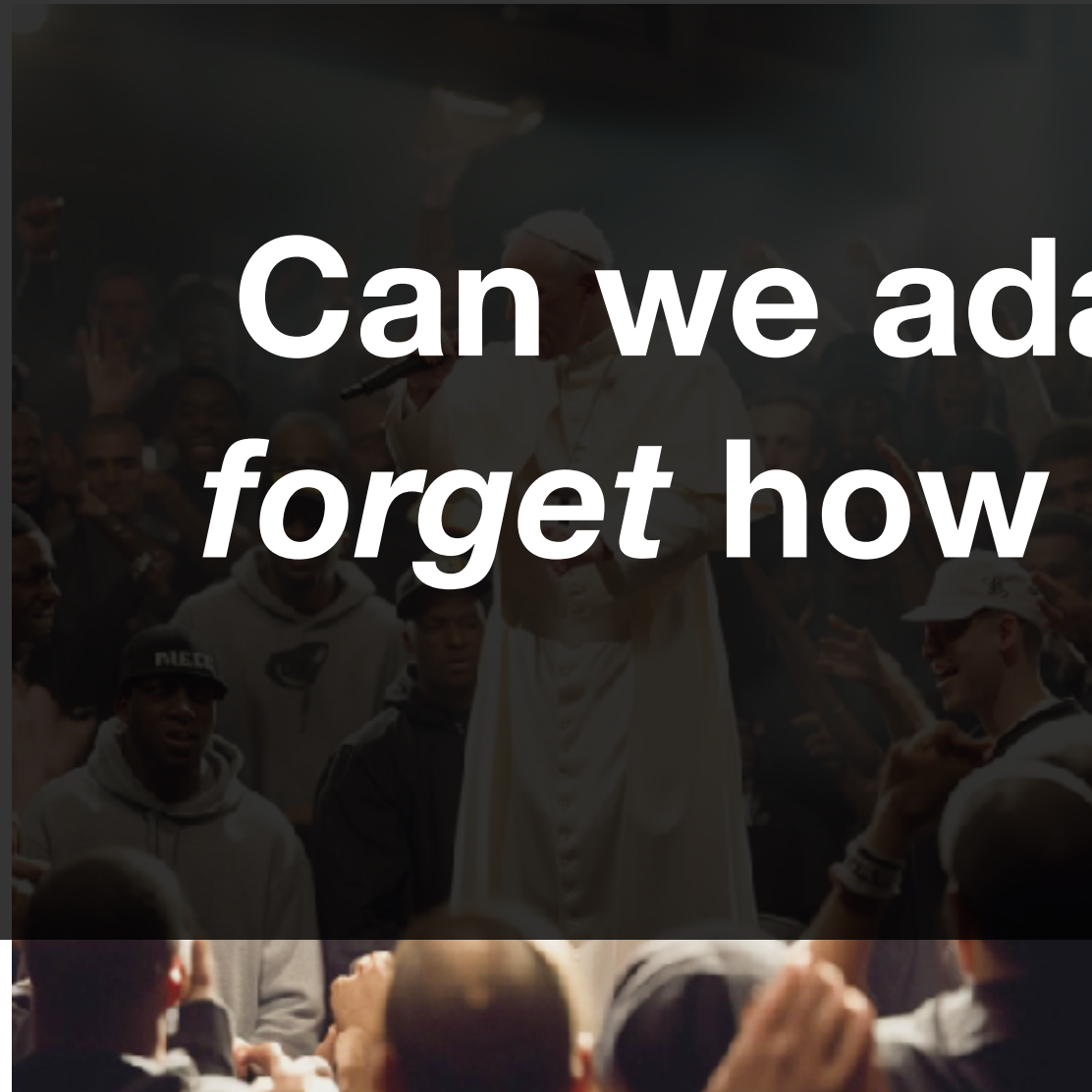
# But...

- Also able to generate inappropriate or harmful images

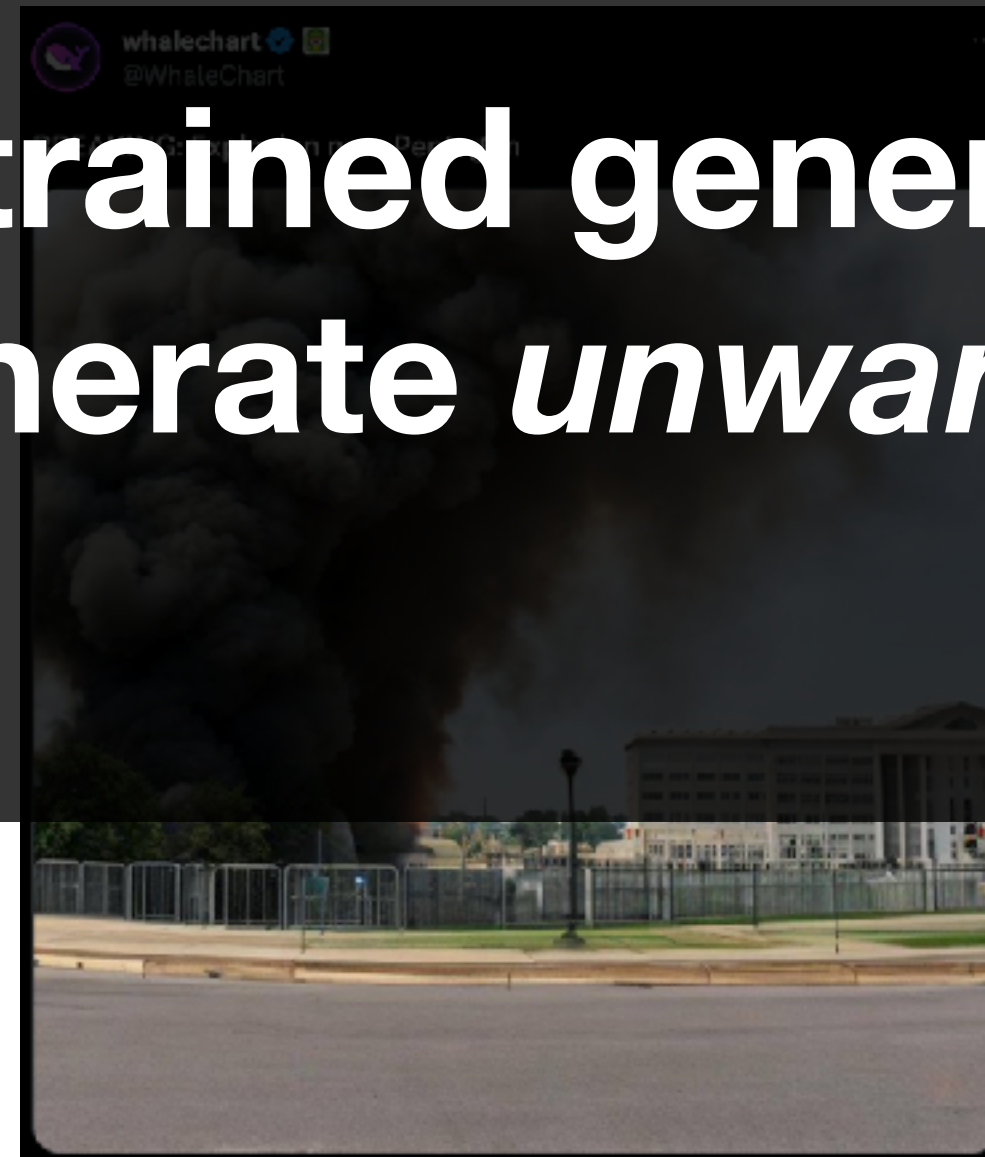
**Fake Pentagon explosion photo goes viral:  
How to spot an AI image**

*A picture claiming to show an explosion near the Pentagon raises concerns about AI's ability to produce misinformation.*

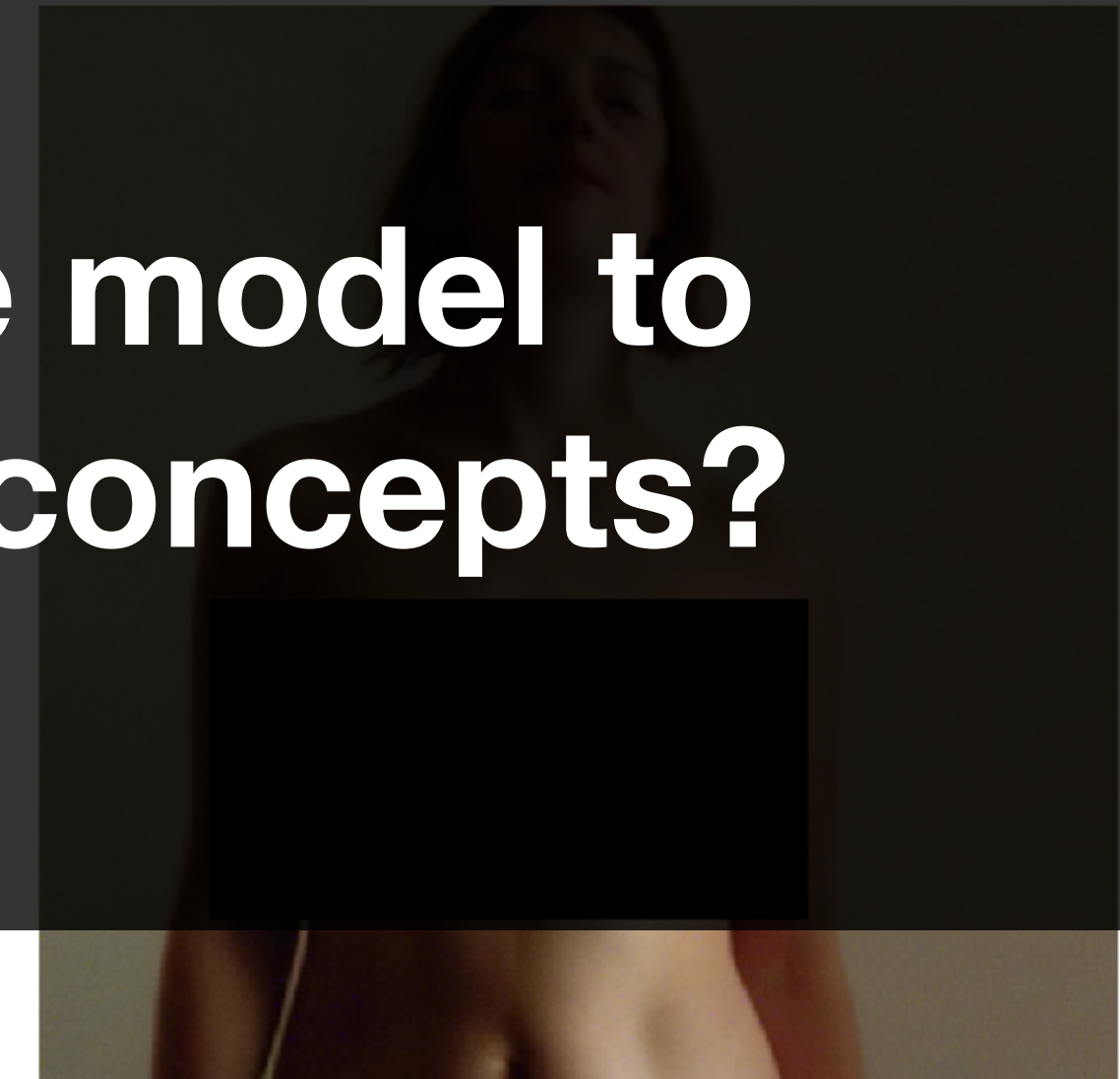
**Can we adapt a trained generative model to  
*forget* how to generate *unwanted* concepts?**



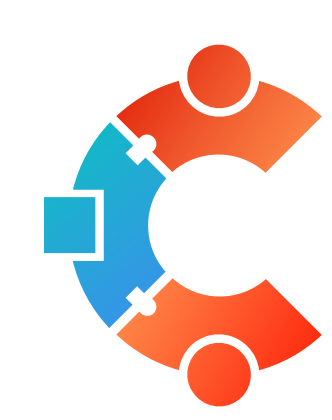
Deepfakes



Fake news

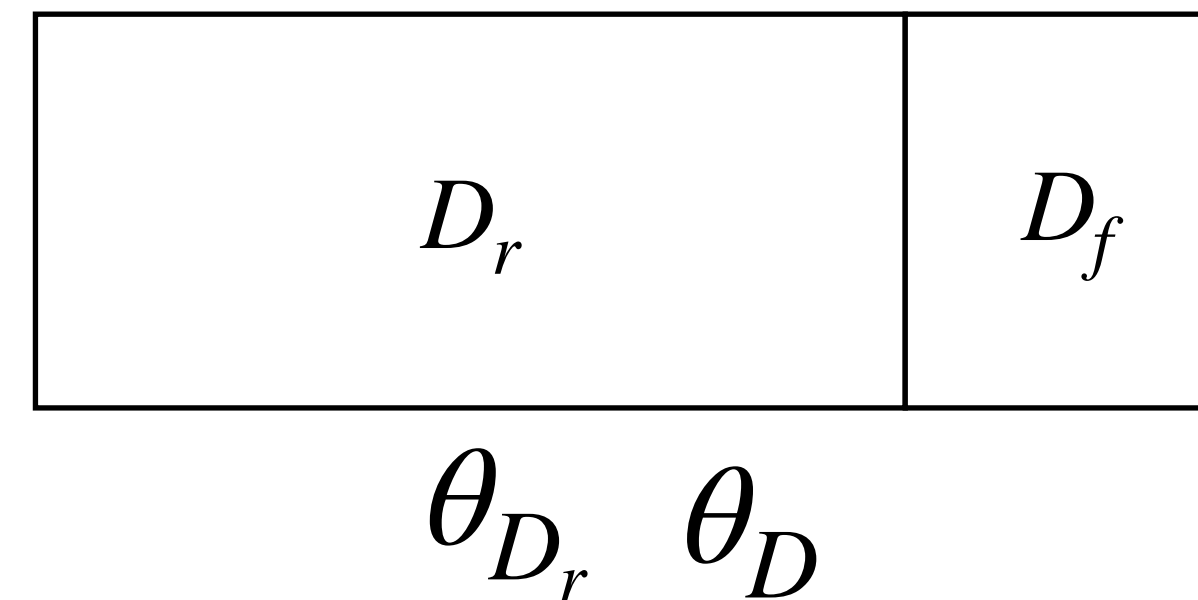
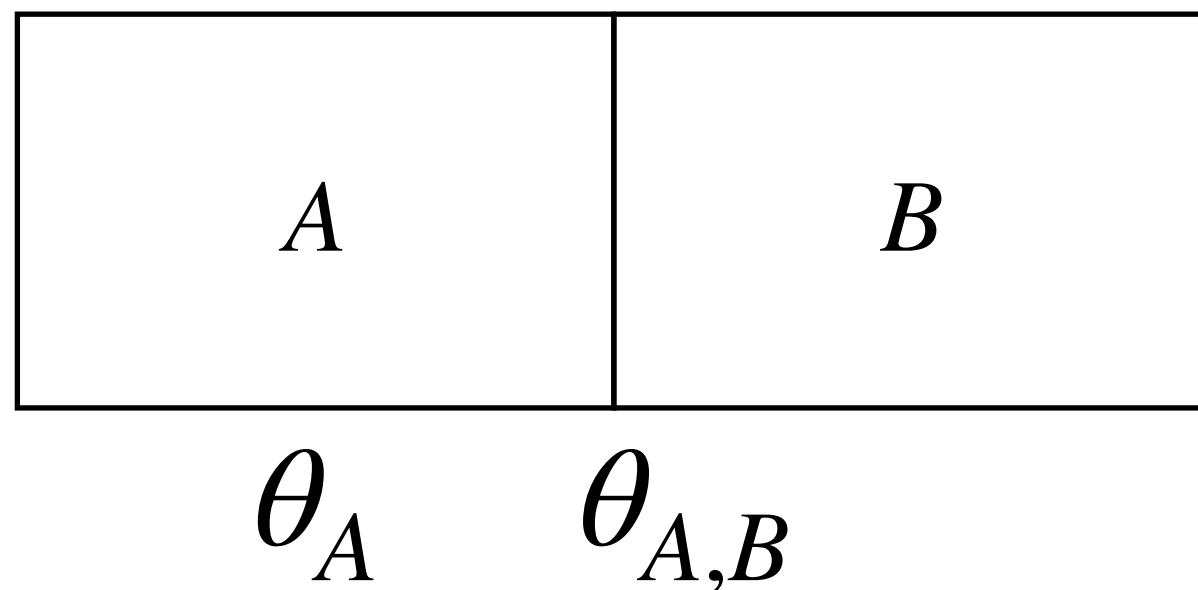


Nudity

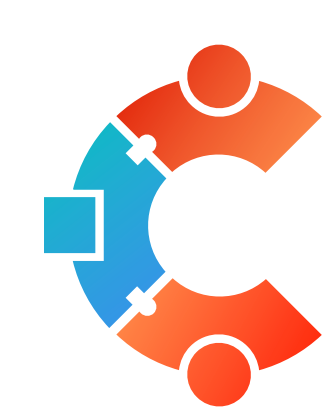


# Continual Learning For Forgetting

- Continual learning traditionally about *preventing* forgetting.
- Learn dataset B without forgetting dataset A,  $\theta_A \rightarrow \theta_{A,B}$
- **This work:** forget  $D_f$  while remembering  $D_r$ ,  $\theta_D \rightarrow \theta_{D_r}$







# Selective Amnesia

- Bayesian approach to maximize the posterior  $\log p(\theta | D_r)$
- Focus on conditional variational generative models with only access to likelihood lower bounds

$$\mathcal{L} = \underbrace{-\mathbb{E}_{p(\mathbf{x}|\mathbf{c})p_f(\mathbf{c})} [\log p(\mathbf{x}|\theta, \mathbf{c})]}_{\text{Want to maximize}} - \underbrace{\lambda \sum_i \frac{F_i}{2} (\theta_i - \theta_i^*)^2}_{\text{Elastic Weight Consolidation (Laplace approx.)}} + \underbrace{\mathbb{E}_{p(\mathbf{x}|\mathbf{c})p_r(\mathbf{c})} [\log p(\mathbf{x}|\theta, \mathbf{c})]}_{\text{Maximize likelihood of remembering set}}$$

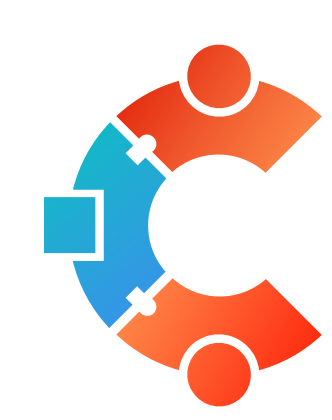
Want to maximize

Minimise likelihood of dataset to forget

Elastic Weight Consolidation (Laplace approx.)

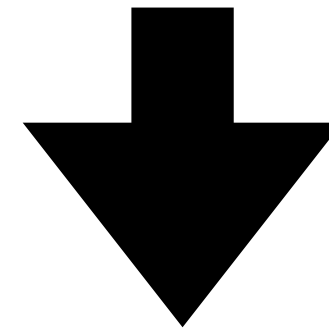
Maximize likelihood of remembering set

**Minimizing the ELBO does not guarantee to reduce the likelihood!**



# Surrogate Distribution

$$\mathcal{L} = -\mathbb{E}_{p(\mathbf{x}|\mathbf{c})p_f(\mathbf{c})} [\log p(\mathbf{x}|\theta, \mathbf{c})] - \lambda \sum_i \frac{F_i}{2} (\theta_i - \theta_i^*)^2 + \mathbb{E}_{p(\mathbf{x}|\mathbf{c})p_r(\mathbf{c})} [\log p(\mathbf{x}|\theta, \mathbf{c})]$$



$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}|\mathbf{c})p_f(\mathbf{c})} [\log p(\mathbf{x}|\theta, \mathbf{c})] - \lambda \sum_i \frac{F_i}{2} (\theta_i - \theta_i^*)^2 + \mathbb{E}_{p(\mathbf{x}|\mathbf{c})p_r(\mathbf{c})} [\log p(\mathbf{x}|\theta, \mathbf{c})]$$

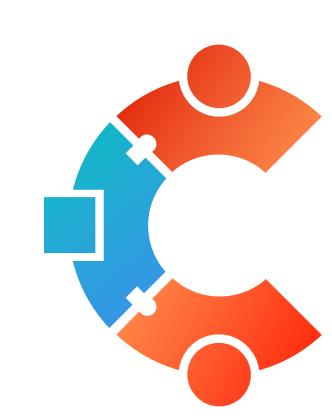
Want to maximize

Maximize likelihood by *remapping* the concept to some  $q(x|\mathbf{c}_f)$

Elastic Weight Consolidation (Laplace approx)

Maximize likelihood of remembering set





# Surrogate Distribution

Suppose: Want to forget nudity for SFW applications.

Let  $c =$  “naked, nude, erotic, sexual” etc.

**Minimize**

$$\mathbb{E}_{p(x|c)p_f(c)} \log p(x|\theta, c)$$

$p(x|c)$



**Maximize**

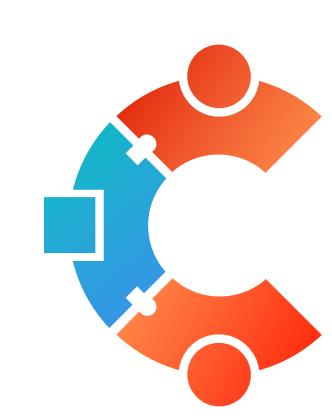
$$\mathbb{E}_{q(x|c)p_f(c)} \log p(x|\theta, c)$$

$q(x|c)$



In other words:  
“Generate clothed persons  
when prompted for nudity”

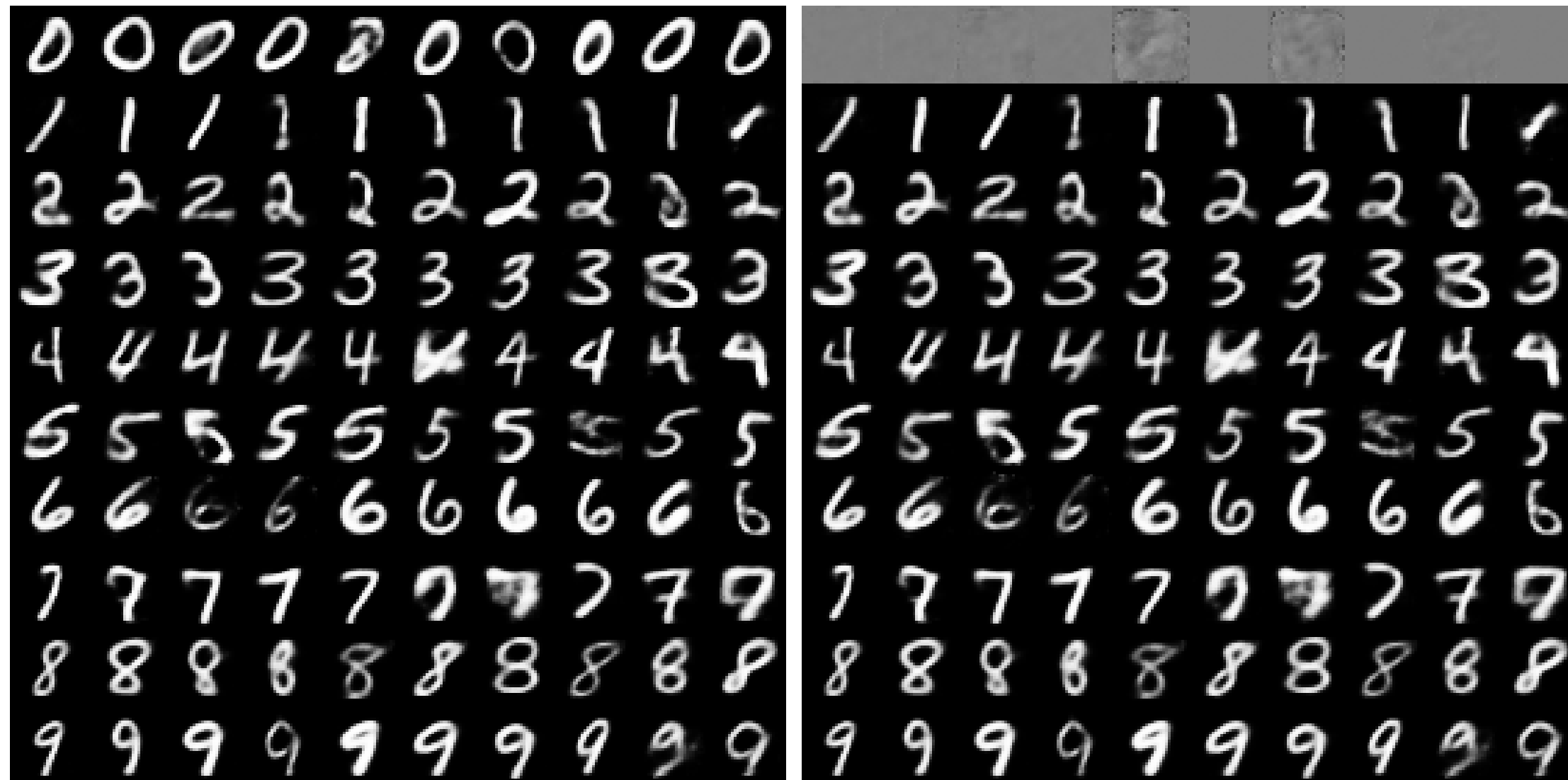




# Qualitative Results

Original

$\lambda = 100$

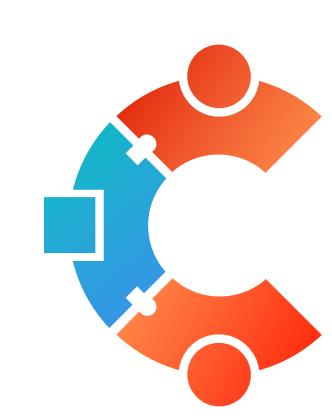


Original

$\lambda = 10$







# Qualitative Results

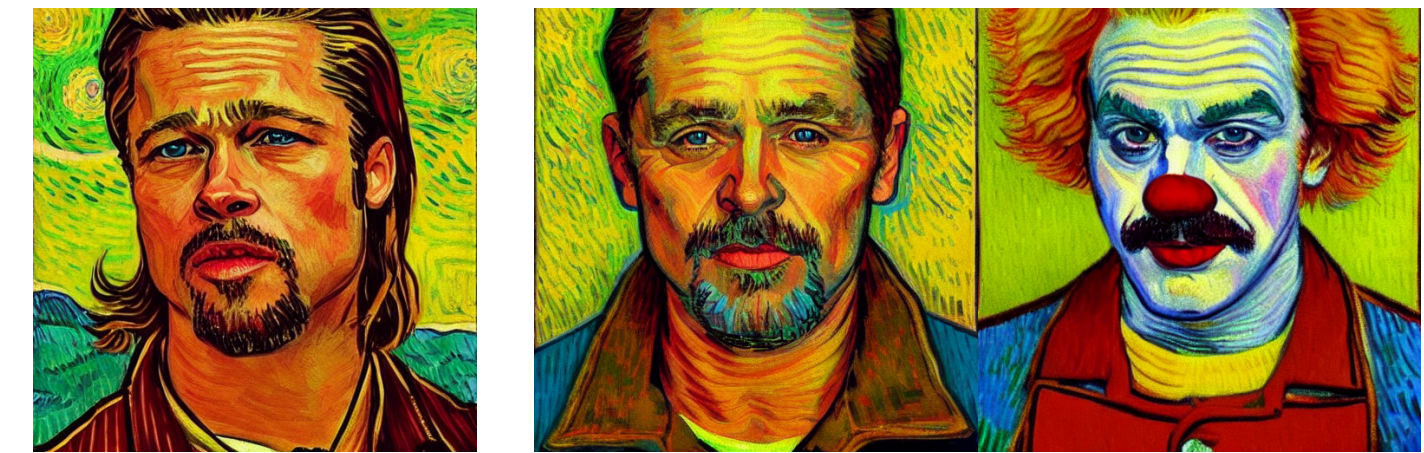
[...] in a grocery store



[...] petting a corgi



painting of [...] in Van Gogh style



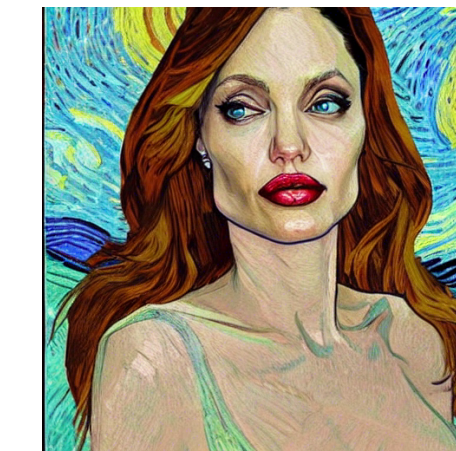
SD v1.4

SA (Ours)



SD v1.4

SA (Ours)



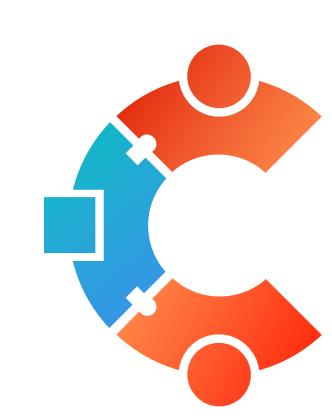
SD v1.4

SA (Ours)

Brad Pitt:  $q(x \mid c_f)$  “middle aged man” or “male clown”

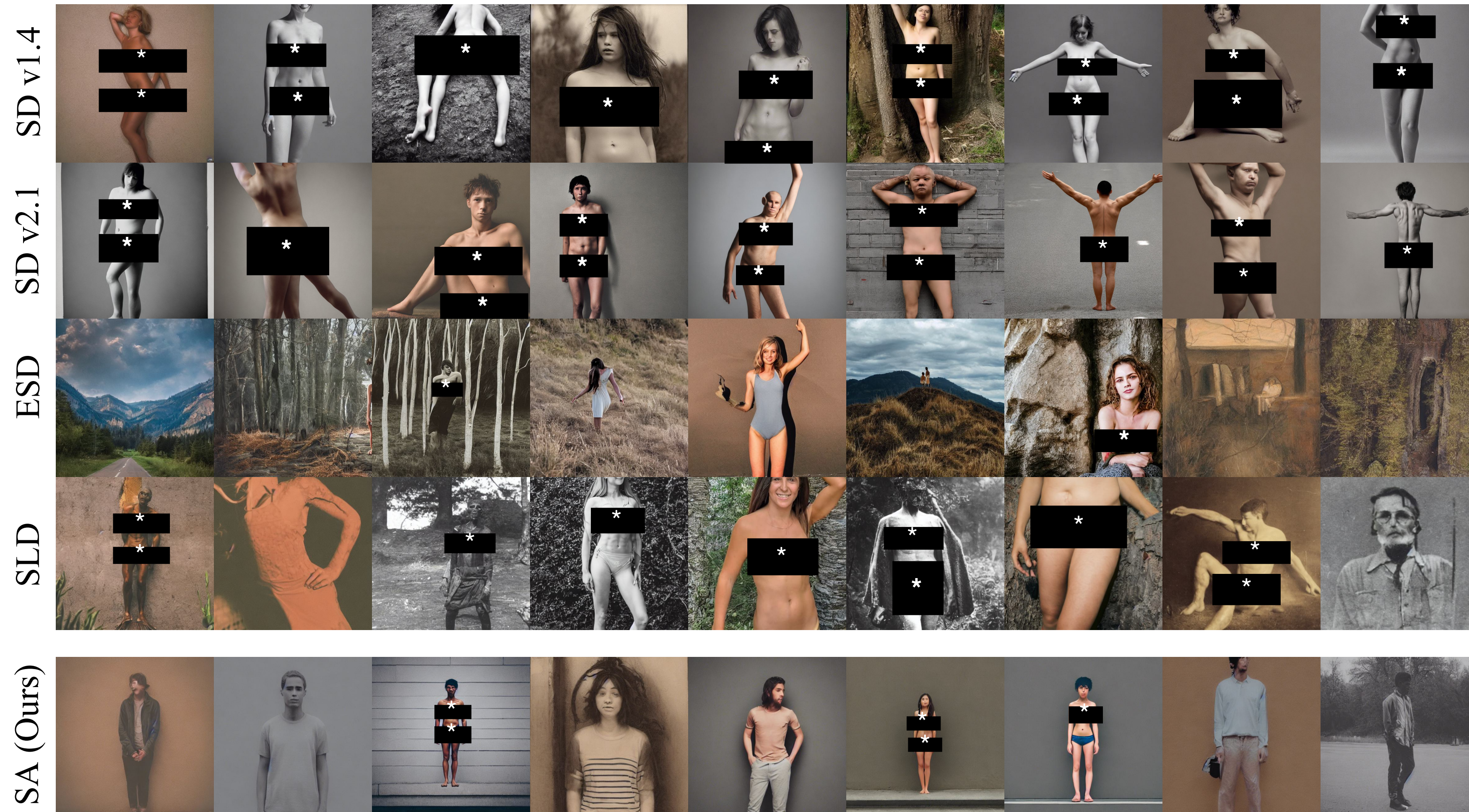
Angelina Jolie:  $q(x \mid c_f)$  “middle aged woman” or “female clown”



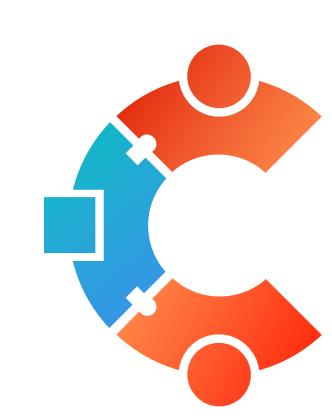


# Qualitative Results

“A photo of a naked person”







# For More Information



[alvin.heng@u.nus.edu](mailto:alvin.heng@u.nus.edu) or [harold@comp.nus.edu.sg](mailto:harold@comp.nus.edu.sg)