

Enhancing Privacy-Utility Trade-offs to Mitigate Memorization in Diffusion Models

Chen Chen, Daochang Liu, Mubarak Shah, Chang Xu

CVPR 2025



THE UNIVERSITY OF
SYDNEY



THE UNIVERSITY OF
WESTERN
AUSTRALIA



Memorization in Diffusion Models

- Stable Diffusion can memorize training images, leading them to reproduce
 - Entire images (global memorization)
 - Parts of images (local memorization)
- This has sparked concerns about the
 - Originality of the generated images
 - Privacy issues

Global Memorization

Real



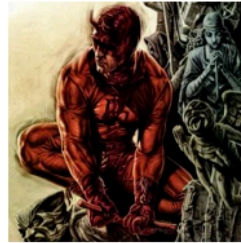
Stable Diffusion Generations



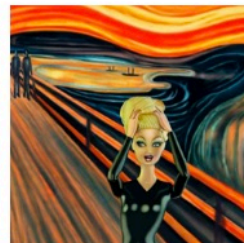
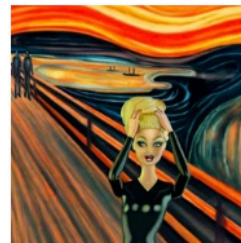
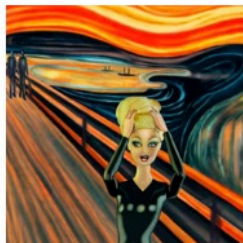
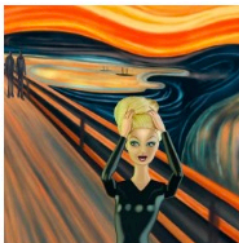
PRSS Generations (Ours)



“<i>I Am Chris Farley</i> Documentary Releases First Trailer”



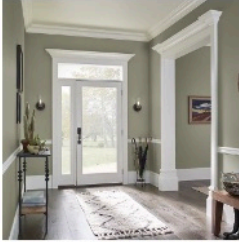
“As Punisher Joins <i>Daredevil</i> Season Two, Who Will the New Villain Be?”



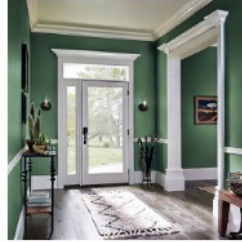
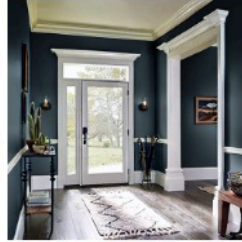
“If Barbie Were The Face of The World's Most Famous Paintings”

Local Memorization

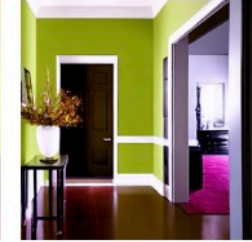
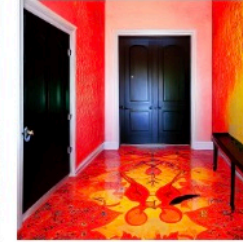
Real



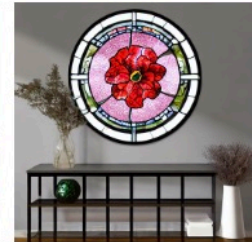
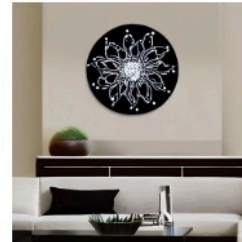
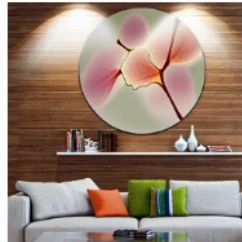
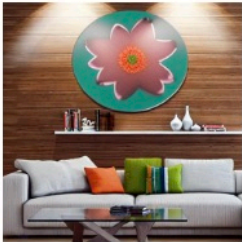
Stable Diffusion Generations



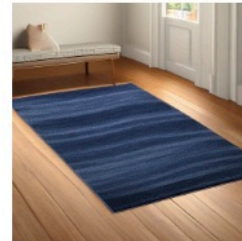
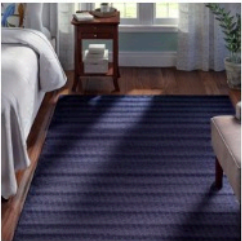
PRSS Generations (Ours)



“Foyer painted in HABANERO”



“Designart Canada White Stained Glass Floral Design 29-in Round Metal Wall Art”



“Falmouth Navy Blue Area Rug by Andover Mills”

Research Gaps

- Existing mitigation strategies
 - Require re-training, searching over billions of images, or
 - Suffer from sub-optimal privacy-utility trade-offs

Research Gaps

Real



Stable Diffusion Generations



Baseline Generations



PRSS Generations (Ours)



“Living in the Light with Ann Graham Lotz”



“Mothers influence on her young hippo”



“Long-Lost F. Scott Fitzgerald Story Rediscovered and Published, 76 Years Later”

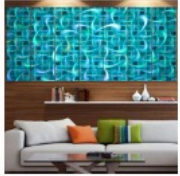


“<i>The Colbert Report</i> Gets End Date”

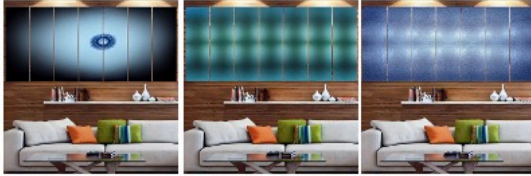
Global Memorization Examples

Research Gaps

Real



Stable Diffusion Generations



Baseline Generations



PRSS Generations (Ours)



“Designart Blue Fractal Abstract Illustration Abstract Canvas Wall Art - 7 Panels”



“Meditation Floor Pillow”



“Designart Circled Blue Psychedelic Texture Abstract Art On Canvas - 7 Panels”



“3D Black & White Skull King Design Luggage Covers 007”

Local Memorization Examples

Trade-off Analysis of Existing Strategies

- Classifier-free guidance (CFG):

$$\hat{e} \leftarrow \epsilon_{\theta}(x_t, e_{\phi}) + s(\epsilon_{\theta}(x_t, e_p) - \epsilon_{\theta}(x_t, e_{\phi})),$$

- The one-step memorization detection strategy:

$$m_t = \|\epsilon_{\theta}(x_t, e_p) - \epsilon_{\theta}(x_t, e_{\phi})\|_2,$$

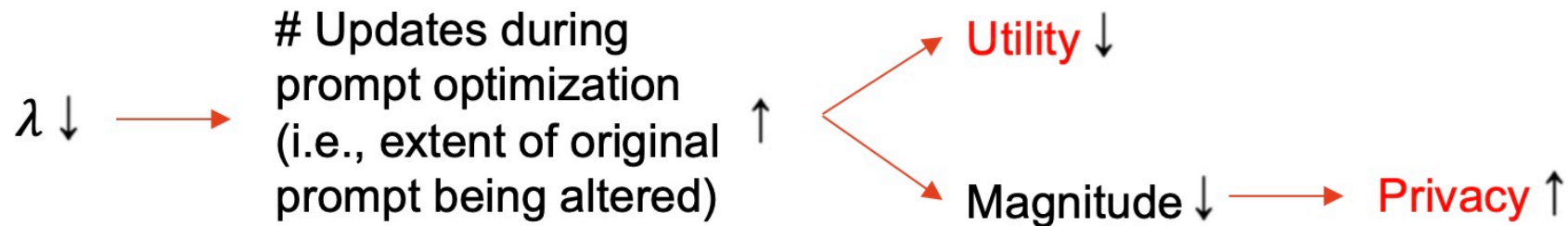
- The corresponding one-step memorization mitigation strategy:

$$\begin{aligned} \hat{e} \leftarrow & [\epsilon_{\theta}(x_t, e_{\phi}) + s(\epsilon_{\theta}(x_t, e_p) - \epsilon_{\theta}(x_t, e_{\phi}))] \mathbb{1}_{\{m_{T-1} < \lambda\}} \\ & + [\epsilon_{\theta}(x_t, e_{\phi}) + s(\epsilon_{\theta}(x_t, e^*) - \epsilon_{\theta}(x_t, e_{\phi}))] \mathbb{1}_{\{m_{T-1} > \lambda\}}, \end{aligned}$$

Trade-off Analysis of Existing Strategies

$$\hat{e} \leftarrow [\epsilon_{\theta}(x_t, e_{\phi}) + s(\epsilon_{\theta}(x_t, e_p) - \epsilon_{\theta}(x_t, e_{\phi}))]\mathbb{1}_{\{m_{T-1} < \lambda\}} \\ + [\epsilon_{\theta}(x_t, e_{\phi}) + s(\epsilon_{\theta}(x_t, e^*) - \epsilon_{\theta}(x_t, e_{\phi}))]\mathbb{1}_{\{m_{T-1} > \lambda\}},$$

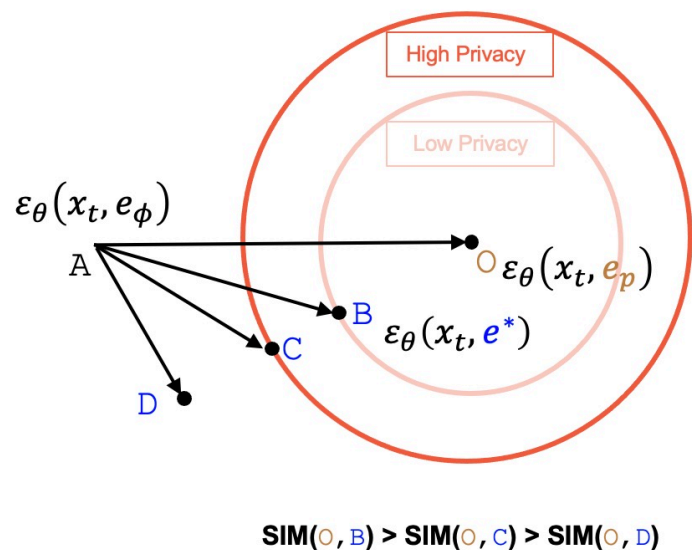
- Currently, the privacy-utility trade-off is governed solely by how extensively the prompt is modified:



- Is it possible to reduce memorization more effectively using a method that achieves a better privacy-utility trade-off?

Trade-off Analysis of Existing Strategies

a) Privacy-Utility Trade-Off in Baseline



$$\hat{\epsilon} \leftarrow [\underbrace{\epsilon_{\theta}(x_t, e_{\phi}) + s(\epsilon_{\theta}(x_t, e_p) - \epsilon_{\theta}(x_t, e_{\phi}))}_{\text{Text-conditional prediction}}] \mathbb{1}_{\{m_{T-1} < \lambda\}} + [\underbrace{\epsilon_{\theta}(x_t, e_{\phi}) + s(\epsilon_{\theta}(x_t, e^*) - \epsilon_{\theta}(x_t, e_{\phi}))}_{\text{Unconditional prediction}}] \mathbb{1}_{\{m_{T-1} > \lambda\}},$$



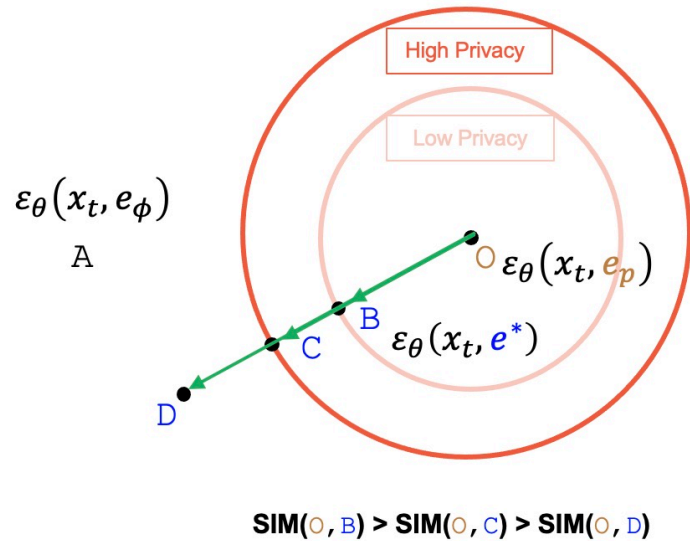
Insights

- We observe that both inputs in the standard CFG framework (the text-conditional and unconditional predictions) are suboptimal for mitigating memorization.
 - Text-conditional prediction
 - Unconditional prediction
- Thus, we propose two meticulously crafted strategies to refine each of them.

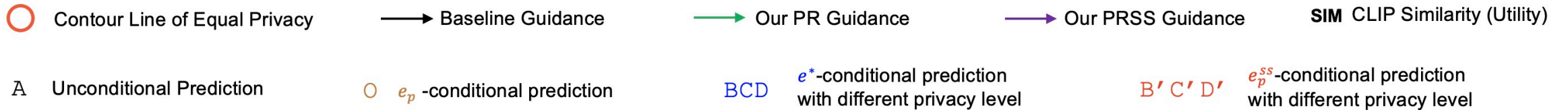
Method: Prompt Re-anchoring (PR)

- PR recognizes e_p as a valuable anchor point for replacing the e_\emptyset in CFG.

b) Enhancing Privacy with Prompt Re-Anchoring (PR)

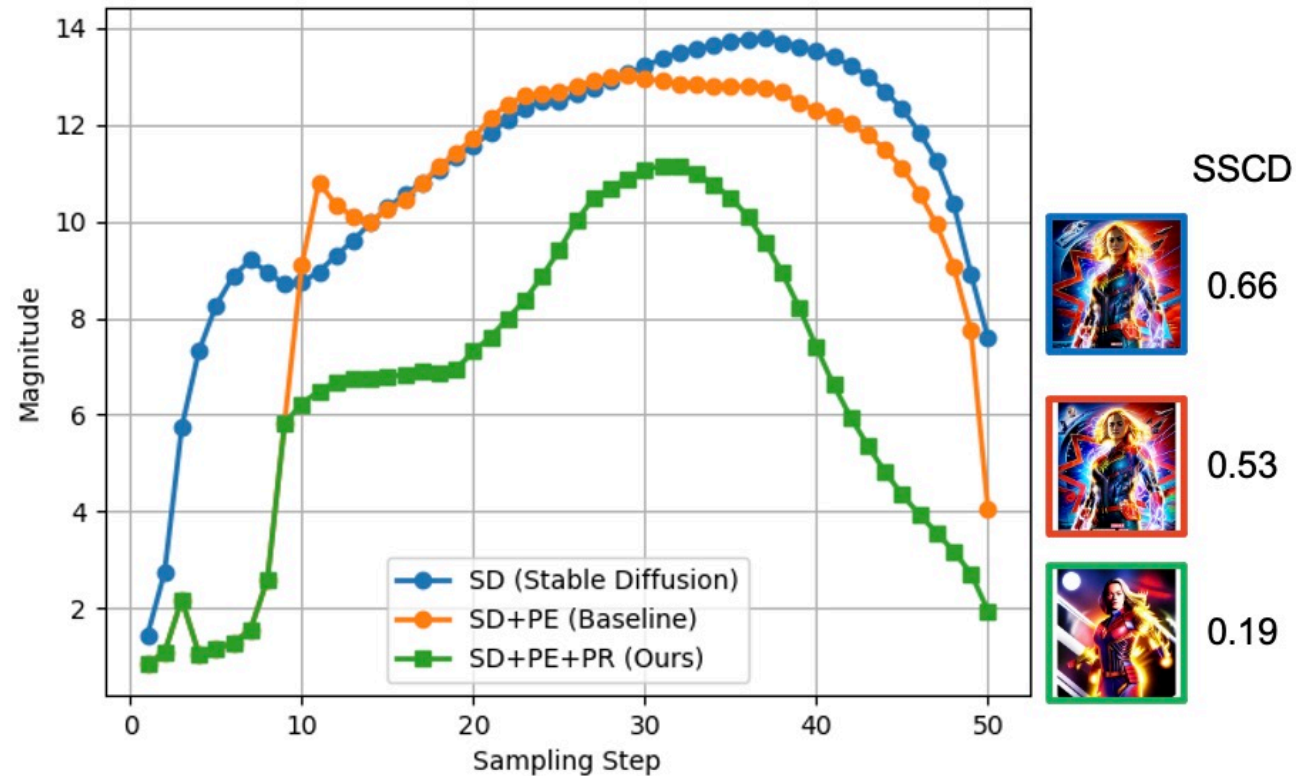


$$\hat{\epsilon} \leftarrow [\underbrace{\epsilon_\theta(x_t, e_\phi) + s(\epsilon_\theta(x_t, e_p) - \epsilon_\theta(x_t, e_\phi))}_{\text{Text-conditional prediction}}] \mathbb{1}_{\{m_{T-1} < \lambda\}} + [\underbrace{\epsilon_\theta(x_t, e_p) + s(\epsilon_\theta(x_t, e^*) - \epsilon_\theta(x_t, e_p))}_{\text{Unconditional prediction}}] \mathbb{1}_{\{m_{T-1} > \lambda\}}.$$



Method: Prompt Re-anchoring (PR)

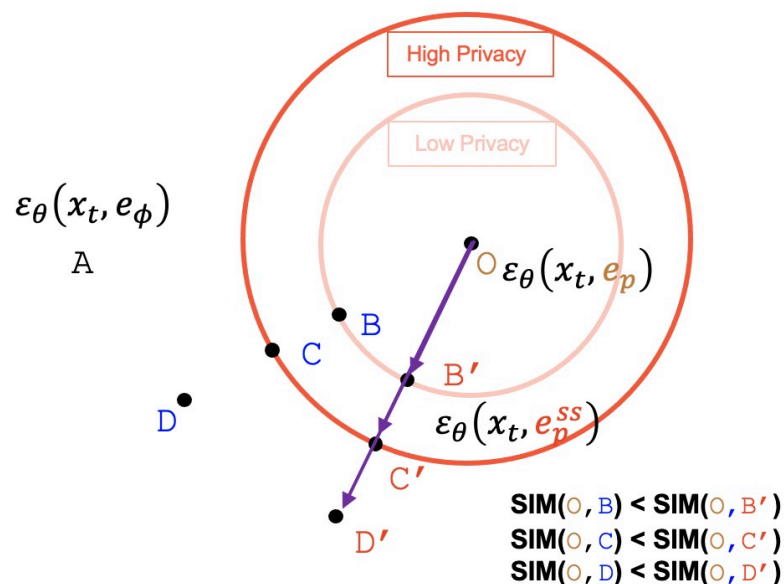
- Analysing PR



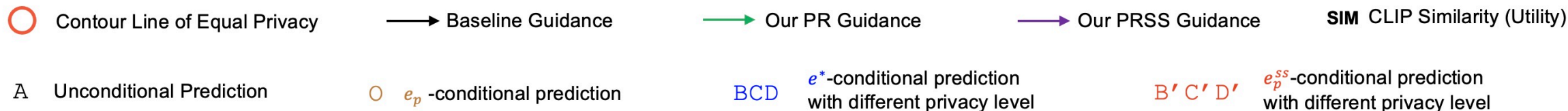
Method: Semantic Search (SS)

- Instead of e^* , SS uses an LLM to find less memorized e_p^{ss} that is semantically similar to e_p :

c) Enhancing Utility with Semantic Search (SS)



$$\hat{\epsilon} \leftarrow [\underbrace{\epsilon_{\theta}(x_t, e_p) + s(\epsilon_{\theta}(x_t, e_p) - \epsilon_{\theta}(x_t, e_{\phi}))}_{\text{Text-conditional prediction}}] \mathbb{1}_{\{m_{T-1} < \lambda\}} + [\underbrace{\epsilon_{\theta}(x_t, e_p) + s(\epsilon_{\theta}(x_t, e_p^{ss}) - \epsilon_{\theta}(x_t, e_p))}_{\text{Unconditional prediction}}] \mathbb{1}_{\{m_{T-1} > \lambda\}}.$$



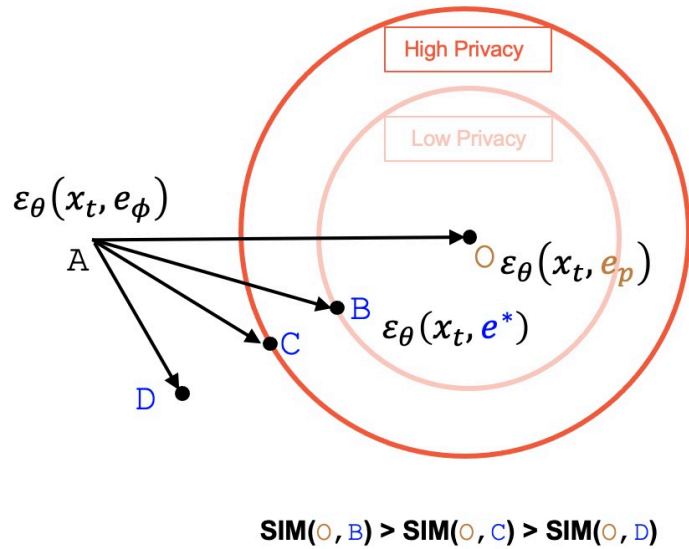
Method: Semantic Search (SS)

- Analysing SS

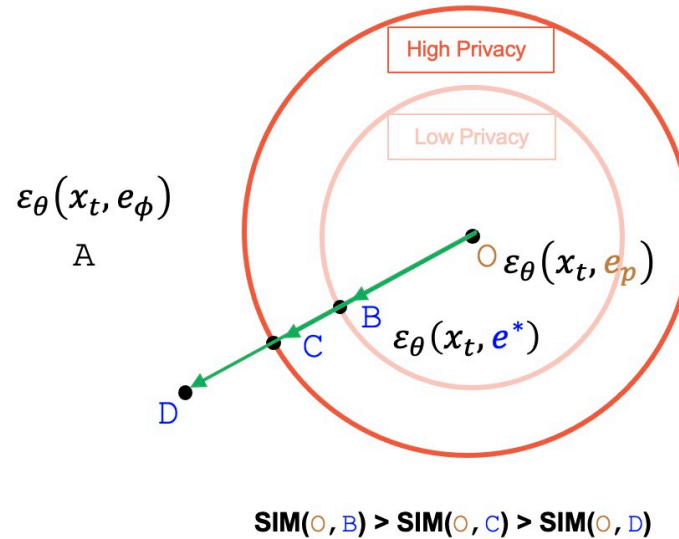
	Stable Diffusion (SD)	w/ Prompt Eng. Baseline (PE)	w/ Semantic Search (SS)
			
Memorization	Yes	No	No
Privacy m_{T-1}	Poor (7.38)	Good (1.15)	Good (0.78)
Utility $CLIP_{txt}$	High (100.00)	Low (2.32)	High (93.82)
$CLIP_{img}$	High (35.92)	Low (13.25)	High (25.20)

Overview of Our Improvements via PRSS

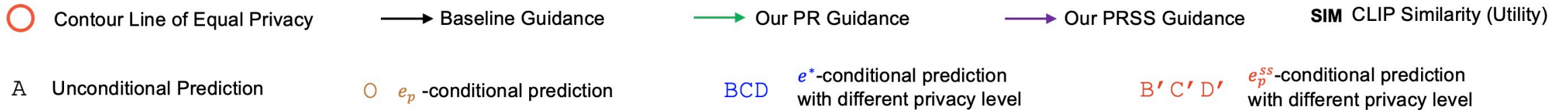
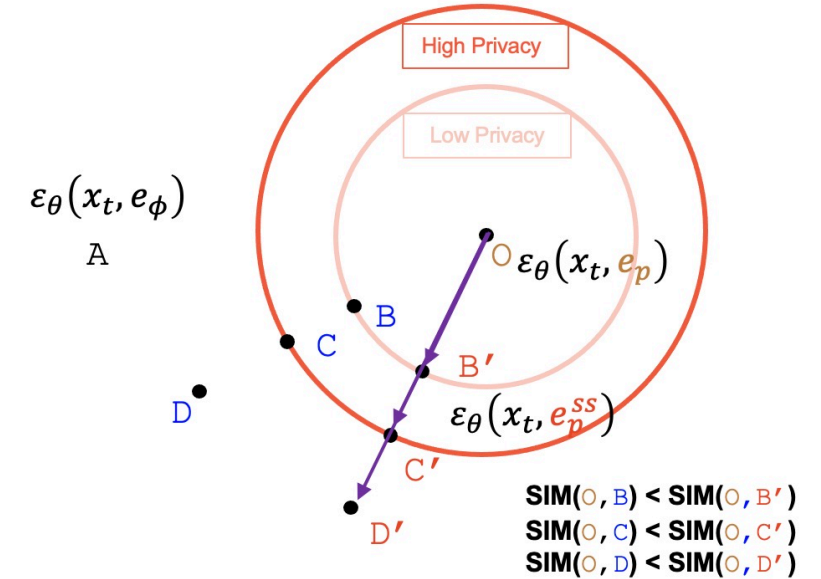
a) Privacy-Utility Trade-Off in Baseline



b) Enhancing Privacy with Prompt Re-Anchoring (PR)

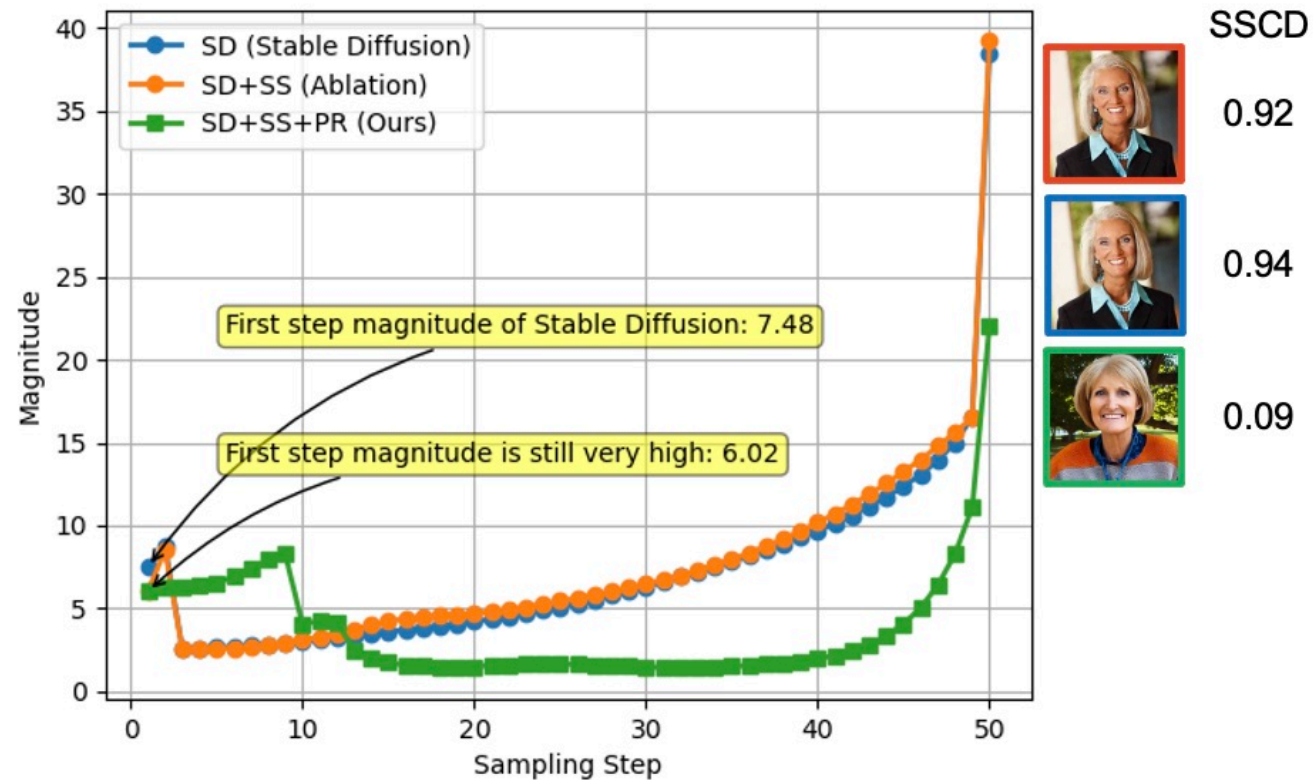


c) Enhancing Utility with Semantic Search (SS)

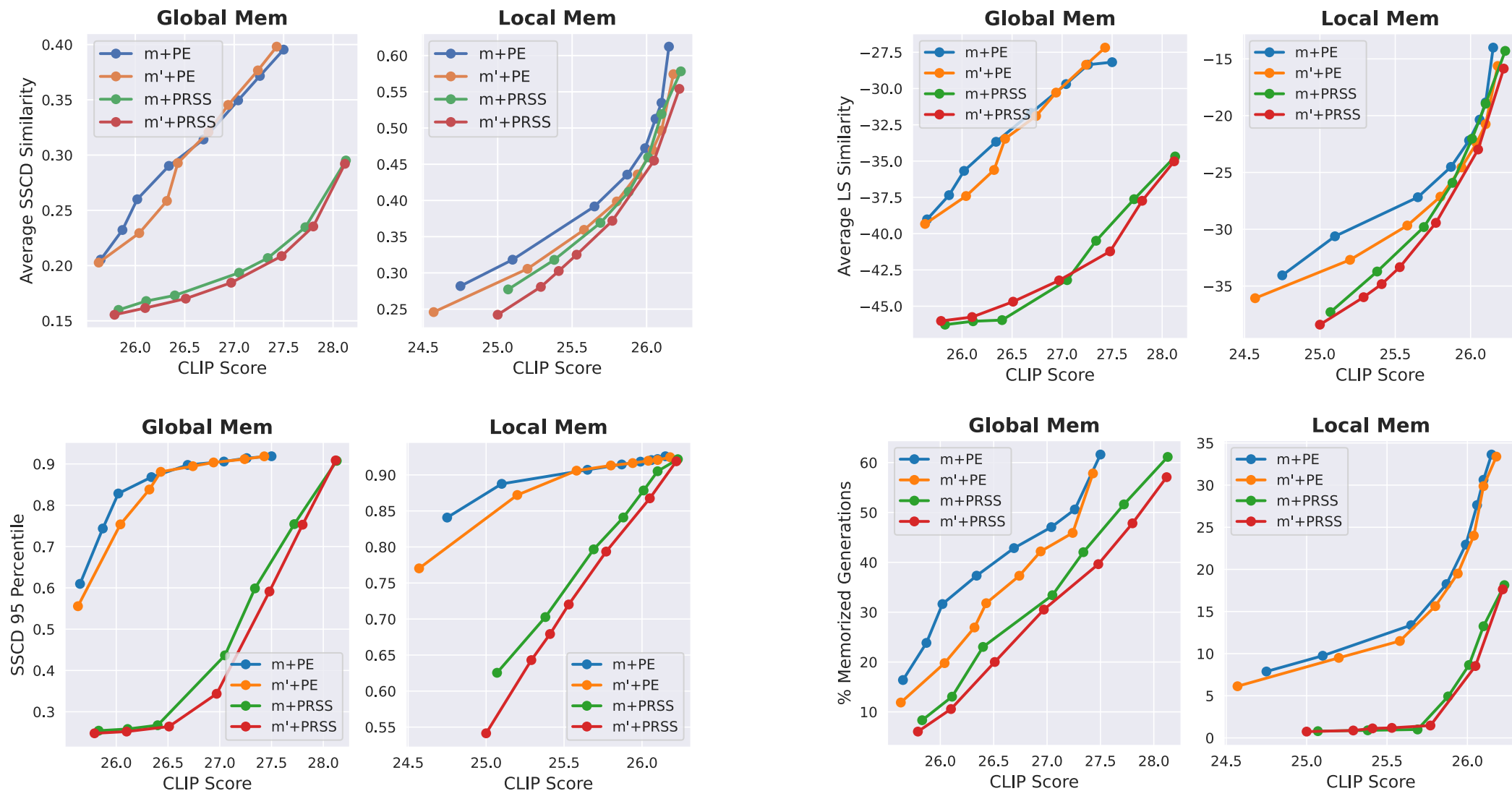


Synergy Effect of the Two Strategies (PR and SS)

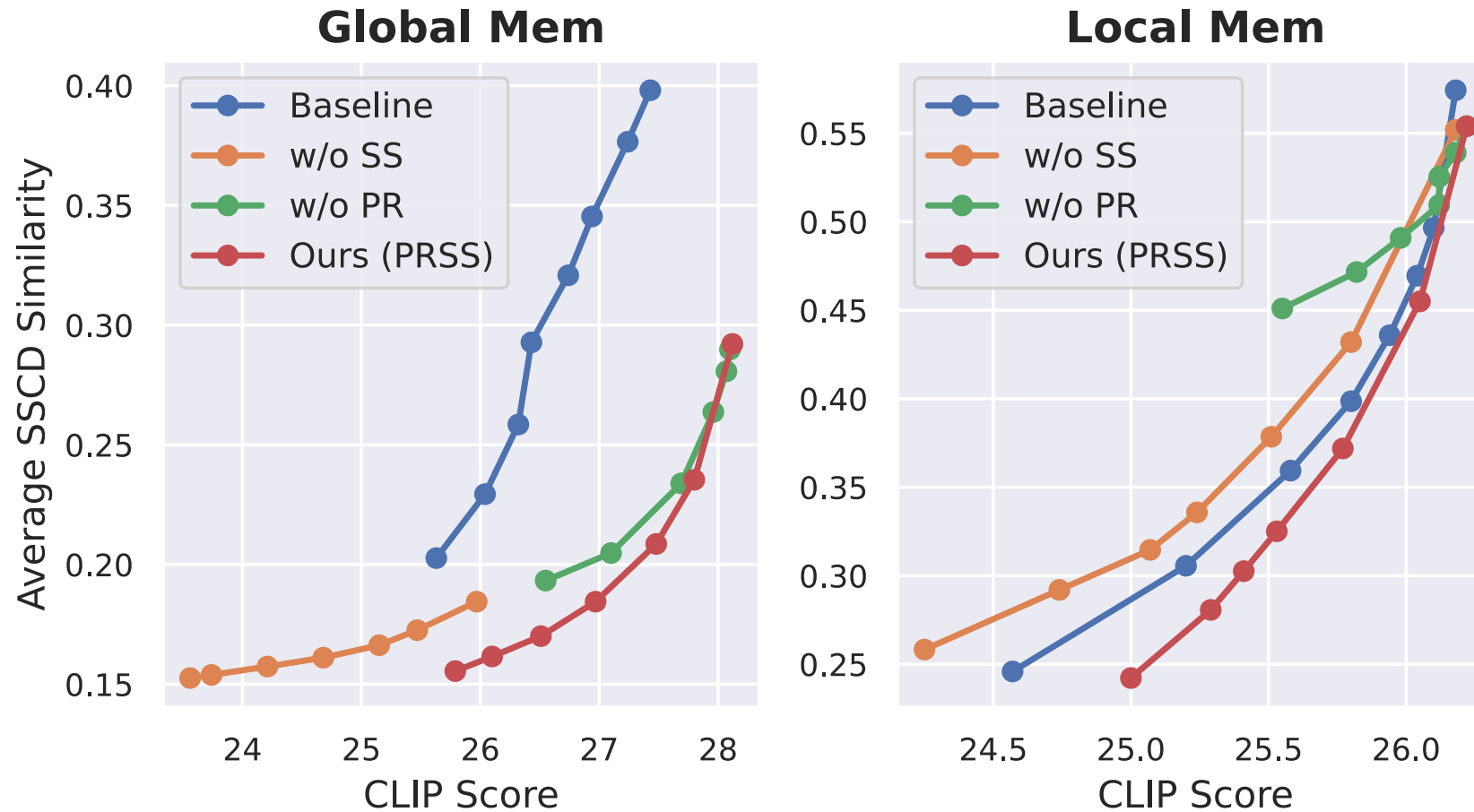
- Analysing PRSS



Quantitative Results



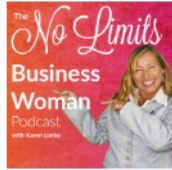
Ablation Studies



Qualitative Results

- Global memorization mitigation

Real



Stable Diffusion Generations



Baseline Generations



PRSS Generations (Ours)



"The No Limits Business Woman Podcast"



"Living in the Light with Ann Graham Lotz"



"Mothers influence on her young hippo"



"Long-Lost F. Scott Fitzgerald Story Rediscovered and Published, 76 Years Later"

Qualitative Results

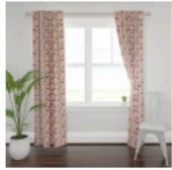
- Global memorization mitigation



Qualitative Results

- Local memorization mitigation

Real



Stable Diffusion Generations



Baseline Generations



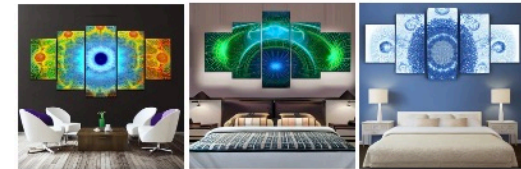
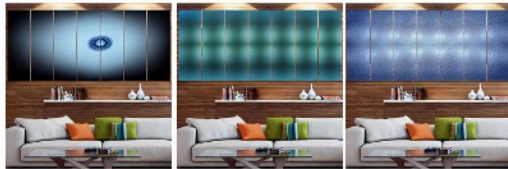
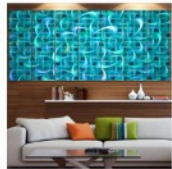
PRSS Generations (Ours)



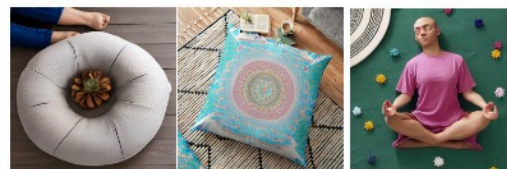
"Plymouth Curtain Panel featuring Madelyn - White Botanical Floral Large Scale by heatherdutton"



"Aero 51-204710 51 Series 15x10 Wheel, Spun, 5 on 4-3/4 BP, 1 Inch BS"



"Designart Blue Fractal Abstract Illustration Abstract Canvas Wall Art - 7 Panels"



"Meditation Floor Pillow"

Qualitative Results

- Local memorization mitigation

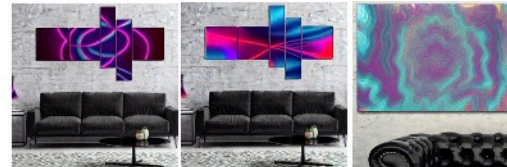
Real



Stable Diffusion Generations



Baseline Generations



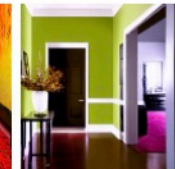
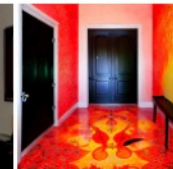
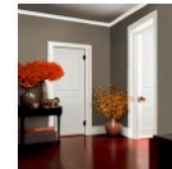
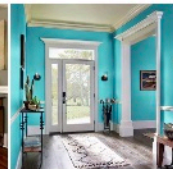
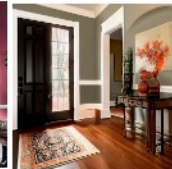
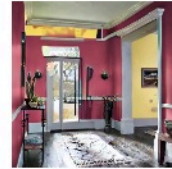
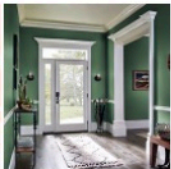
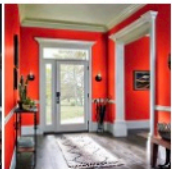
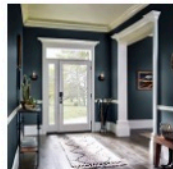
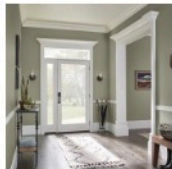
PRSS Generations (Ours)



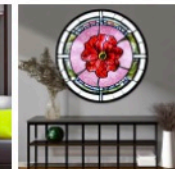
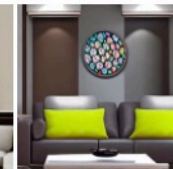
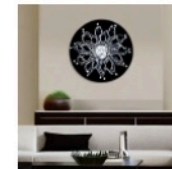
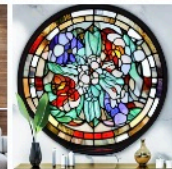
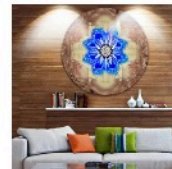
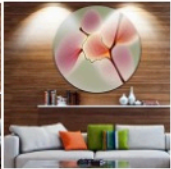
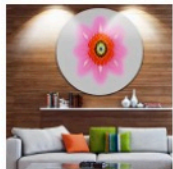
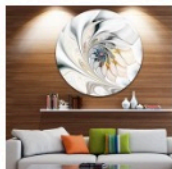
“Designart Circled Blue Psychedelic Texture Abstract Art On Canvas - 7 Panels”



“3D Black & White Skull King Design Luggage Covers 007”



“Foyer painted in HABANERO”



“Designart Canada White Stained Glass Floral Design 29-in Round Metal Wall Art”

Conclusion

- This paper introduces *PRSS*, which consists of *Prompt Re-anchoring (PR)* & *Semantic Search (SS)* as deliberately designed strategies to mitigate memorization in diffusion models.
- Results show that PRSS enhances existing methods' privacy-utility trade-off at different privacy levels.
- Implementationally, PRSS is simple and efficient, necessitating only adjustments to the existing CFG equation without re-training, fine-tuning, or searching over the training data.

Limitation and Future Work

- PRSS relies on the accuracy of a pre-defined detection strategy to trigger the mitigation mechanism.
- Thus, improving the detection mechanism's accuracy in future works would orthogonally bolster the effectiveness of PRSS.