

Mitigating Semantic Collapsing in Generative Personalization with Test-time Embedding Adjustment



Anh Bui¹, Thuy-Trang Vu¹, Trung Le¹, Junae Kim², Tamas Abraham², Rollin Omari², Amar Kaur², Dinh Phung¹
¹ Monash University, ² DSTG

Qualitative Results



OminiControl \uparrow vs OminiControl with \downarrow
Prompt: "V* wearing glasses"



Quantitative Results

Method	CLIP _T ^p \uparrow	CLIP _T ^f \uparrow	CLIP-I \uparrow	DINO-I \uparrow	VLM-P \uparrow	VLM-I \uparrow
CC101 - Pet Dog						
ES	18.54	26.02	61.33	43.71	64.25	74.00
ES+TEA	18.72 (+0.18)	26.11 (+0.09)	64.56 (+3.23)	48.32 (+4.61)	66.50 (+2.25)	77.25 (+3.25)
CC101 - Plushie Teddybear						
ES	20.48	26.80	81.64	49.08	78.00	80.25
ES+TEA	20.61 (+0.13)	27.3 (+0.50)	82.84 (+1.20)	51.17 (+2.09)	80.25 (+2.25)	81.50 (+1.25)
Subject - Clock						
OC	18.11	23.90	81.37	32.41	67.50	62.25
OC+TEA	18.78 (+0.67)	23.98 (+0.08)	83.10 (+1.73)	34.48 (+2.07)	71.75 (+4.25)	64.50 (+2.25)
Subject - Oranges						
OC	21.49	27.62	70.43	30.33	68.50	53.00
OC+TEA	21.60 (+0.11)	27.70 (+0.08)	71.90 (+1.47)	31.64 (+1.31)	70.00 (+1.50)	55.50 (+2.50)
Subject - Penguin						
OC	20.30	31.61	78.58	45.59	86.25	83.25
OC+TEA	20.33 (+0.03)	32.02 (+0.41)	80.64 (+2.06)	49.37 (+3.78)	90.50 (+4.25)	86.75 (+3.50)
Relationship - A <Carved by> B						
RV	25.64	27.74	N/A	N/A	N/A	N/A
RV+TEA	27.84 (+2.20)	30.17 (+2.43)	N/A	N/A	N/A	N/A
Relationship - A <Inside> B						
RV	24.97	27.87	N/A	N/A	N/A	N/A
RV+TEA	25.15 (+0.18)	28.40 (+0.53)	N/A	N/A	N/A	N/A
Relationship - A <Painted on> B						
RV	23.98	30.07	N/A	N/A	N/A	N/A
RV+TEA	24.38 (+0.40)	30.35 (+0.28)	N/A	N/A	N/A	N/A

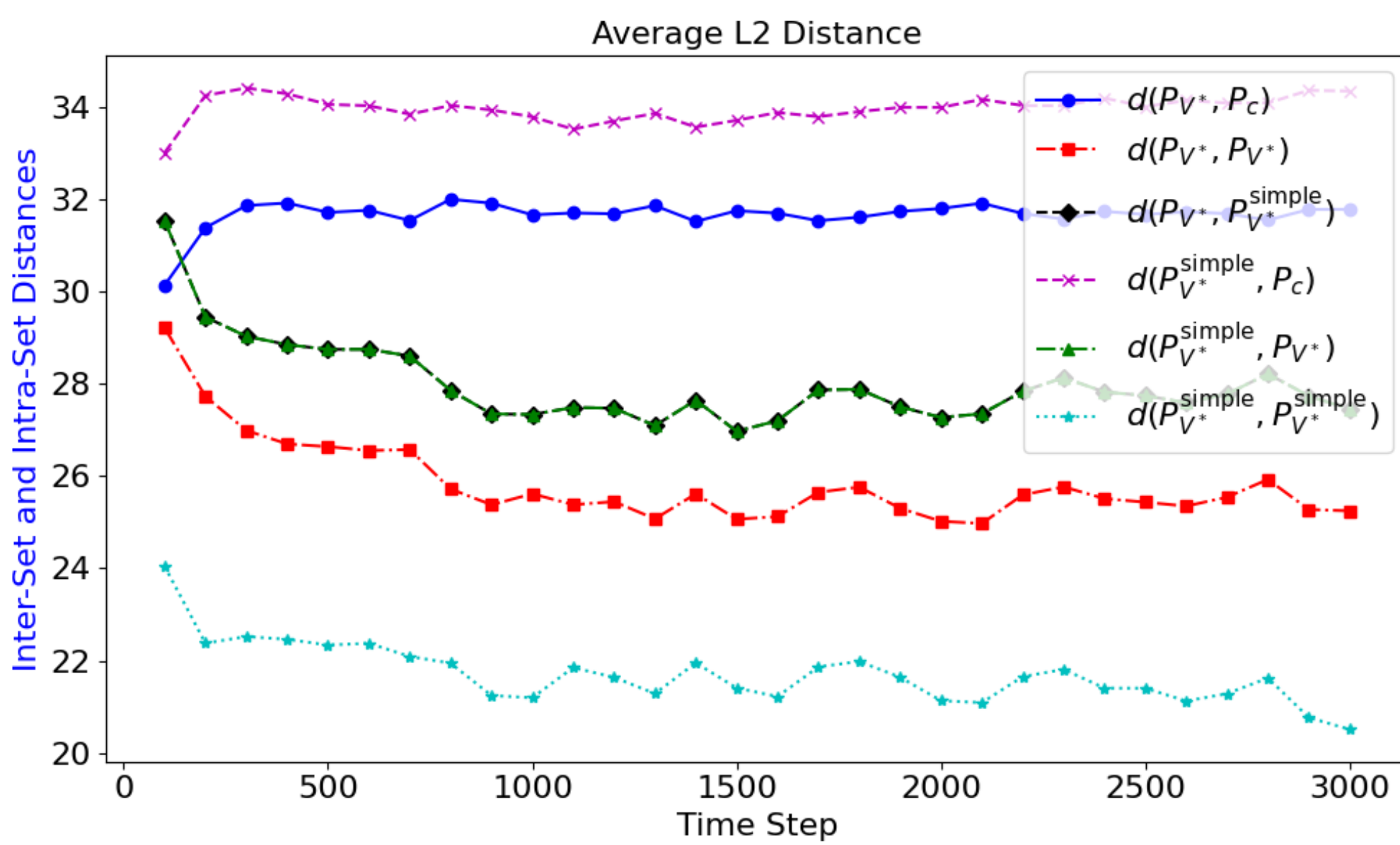
The "Semantic Collapsing Problem"

Well known phenomenon \downarrow but don't understand Why yet

$$G(\lfloor p, V^* \rfloor) \rightarrow G(V^*)$$

Finding 1: the phenomenon \uparrow because of the SCP \downarrow

$$\tau(\lfloor p, V^* \rfloor) \rightarrow \tau(V^*)$$

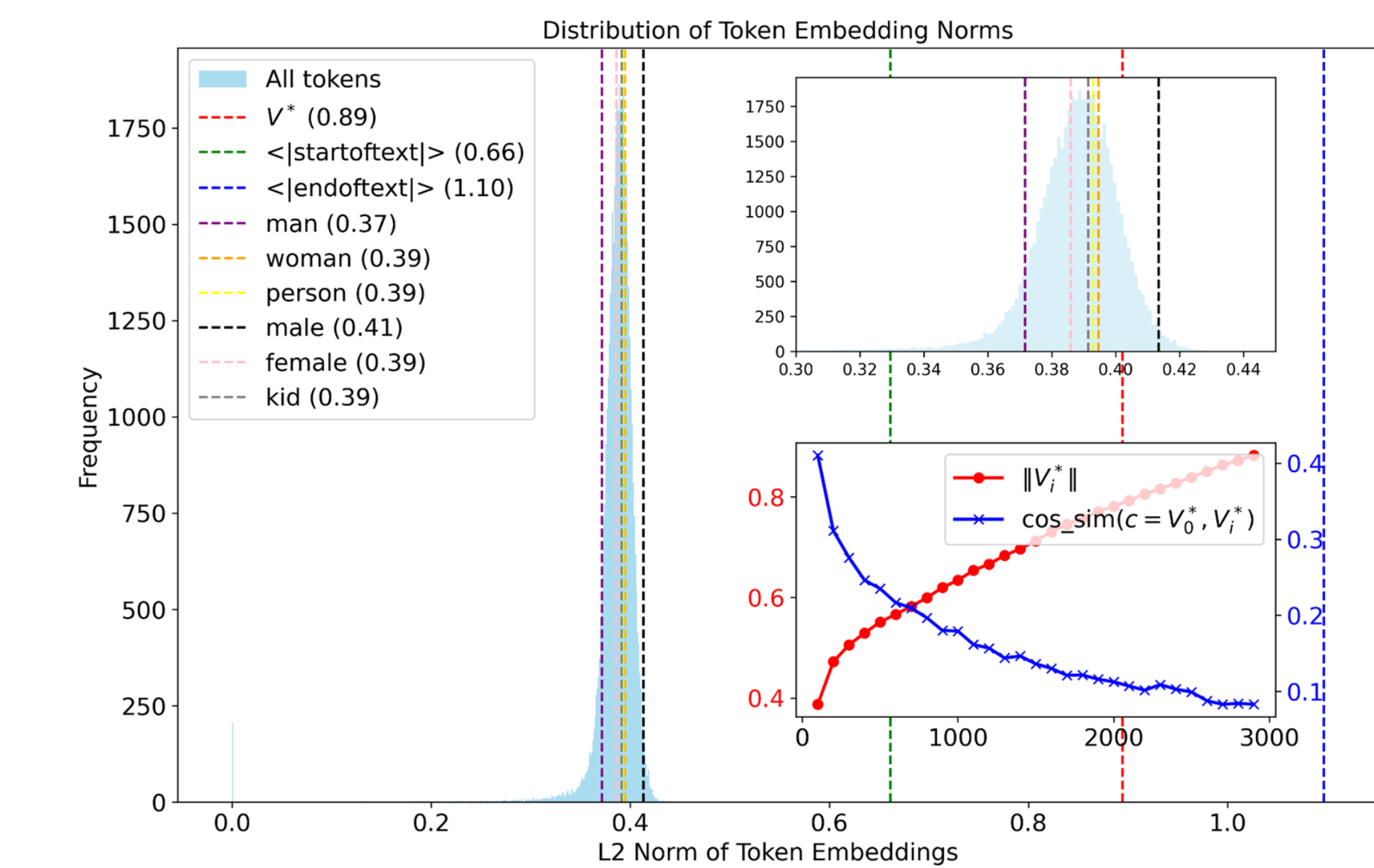


Finding 2: Root cause of SCP because of the unconstrained optimization in personalization process, i.e., TI or DB

$$TI: \min_{V^*} \mathbb{E}_{x,p,\epsilon,t} \left[\|\epsilon - \epsilon_\theta(x_{t,\epsilon}, t, \lfloor p, V^* \rfloor)\|_2^2 \right]$$

$$DB: \min_{\theta, V^*} \mathbb{E}_{x,p,\epsilon,t} \left[\|\epsilon - \epsilon_\theta(x_{t,\epsilon}, t, \lfloor p, V^* \rfloor)\|_2^2 + \lambda \|\epsilon' - \epsilon_\theta(x'_{t,\epsilon'}, t', p^{DB})\|_2^2 \right]$$

Embedding of V* can be arbitrary large and shifted from c



Test-time Embedding Adjustment

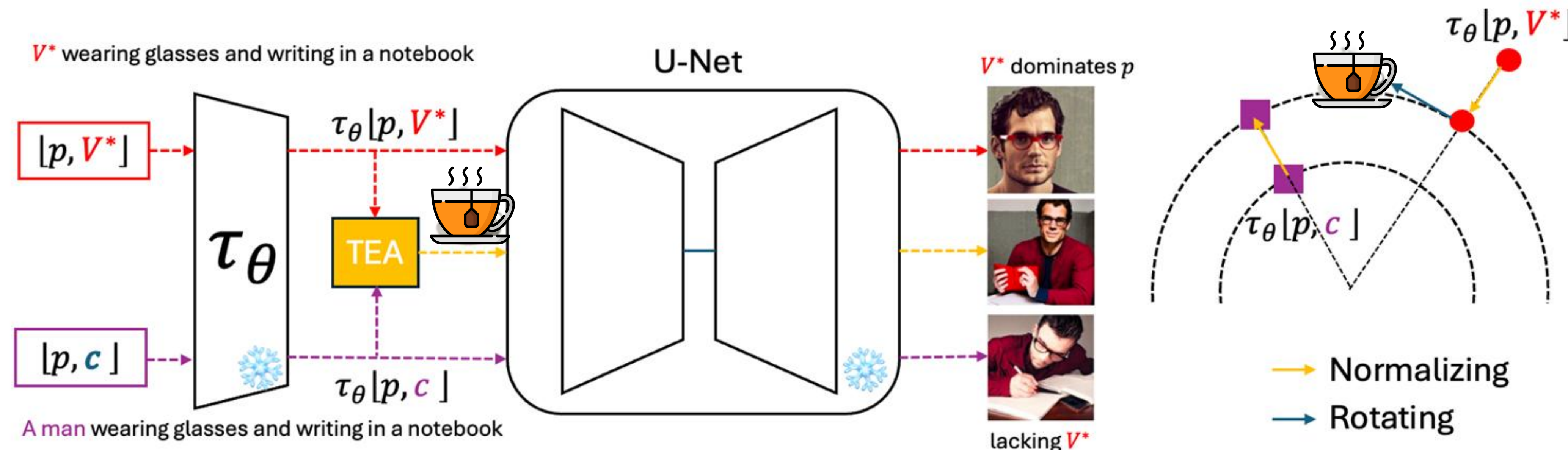
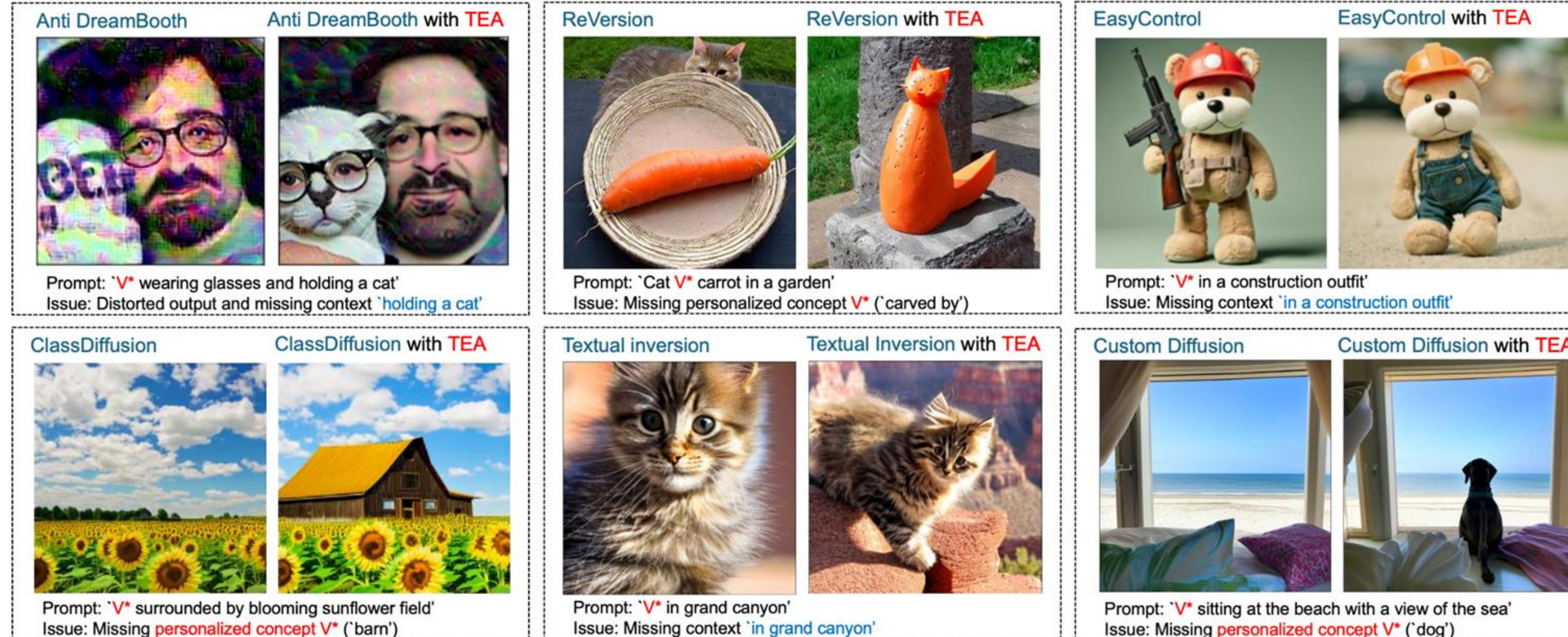


$$\hat{\tau}(\lfloor p, V^* \rfloor) = \frac{\sin((1-\alpha)\theta)}{\sin(\theta)} \hat{\tau}(\lfloor p, V^* \rfloor) + \frac{\sin(\alpha\theta)}{\sin(\theta)} \hat{\tau}(\lfloor p, c \rfloor)$$

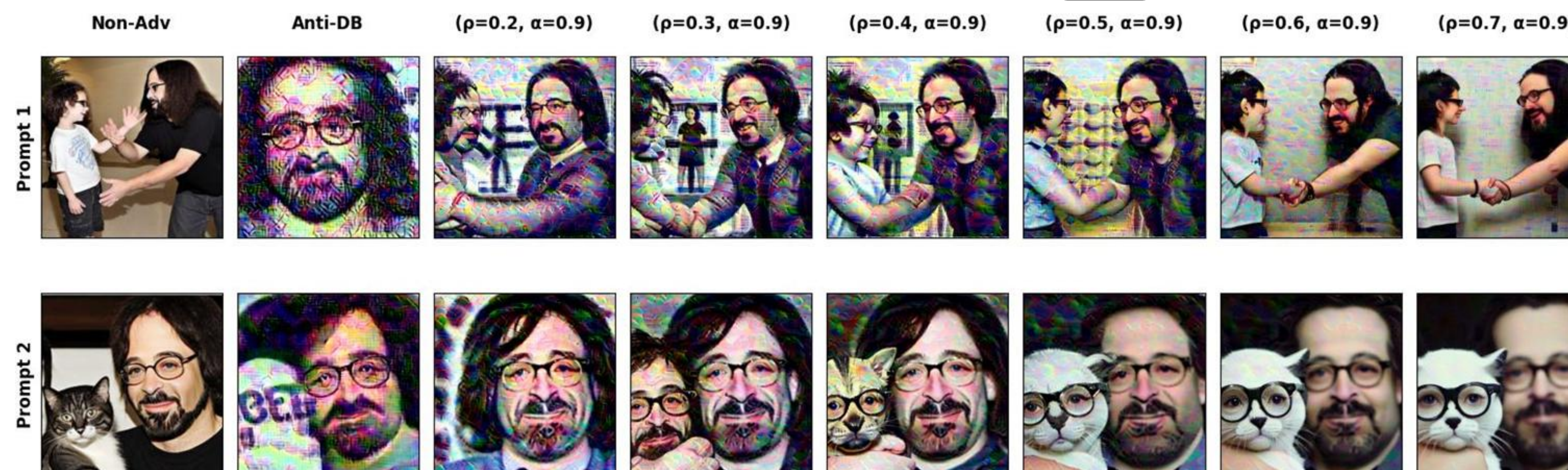
$$\hat{\tau}(\lfloor p, V^* \rfloor) = \beta \frac{\tau(\lfloor p, V^* \rfloor)}{\|\tau(\lfloor p, V^* \rfloor)\|}$$

$$\hat{\tau}(\lfloor p, c \rfloor) = \beta \frac{\tau(\lfloor p, c \rfloor)}{\|\tau(\lfloor p, c \rfloor)\|}$$

$$\theta = \arccos(\hat{\tau}(\lfloor p, c \rfloor), \hat{\tau}(\lfloor p, V^* \rfloor))$$



Anti-Personalization with



Non-Adv: Standard Personalization \rightarrow Attackers can generate inappropriate images of anyone
 Anti-DB: Anti-Personalization \rightarrow Defenders protect users' images by adding invisible mask \rightarrow Attackers fail
 Underlying mechanism: Adversarial mask makes the SCP worse \rightarrow V* over-domination \rightarrow generating distortion
 Anti-DB with TEA \rightarrow TEA can mitigate SCP \rightarrow reducing concept domination \rightarrow generate clean images

MORE INSIGHTFUL FINDINGS
tuananhbui89.github.io

