



LUSD: Localized Update Score Distillation for Text-Guided Image Editing

Worameth Chinchuthakun*1,2 Pitchaporn Rewatbowornwong¹

Tossaporn Saengja*1,3 Nontawat Tritrong¹ Pramook Khungurn⁴ Supasorn Suwajanakorn¹

³ Faculty of Medicine Siriraj Hospital VISTEC ² Siam Commercial Bank ⁴ pixiv Inc.

*Equal contributions



Problem & Background

- Image editing: modify the input image to match the prompt while preserving some elements of the input image.
- Recent diffusion-based approaches:
- Supervised: fine-tune models on synthetic data

SCB 6 DIXIV

Unsupervised: invert diffusion process, score distillation

Challenges

1. Supervised methods

- often biased, synthetic training data 🌉
- leads to poor generalization

2. Inversion-based

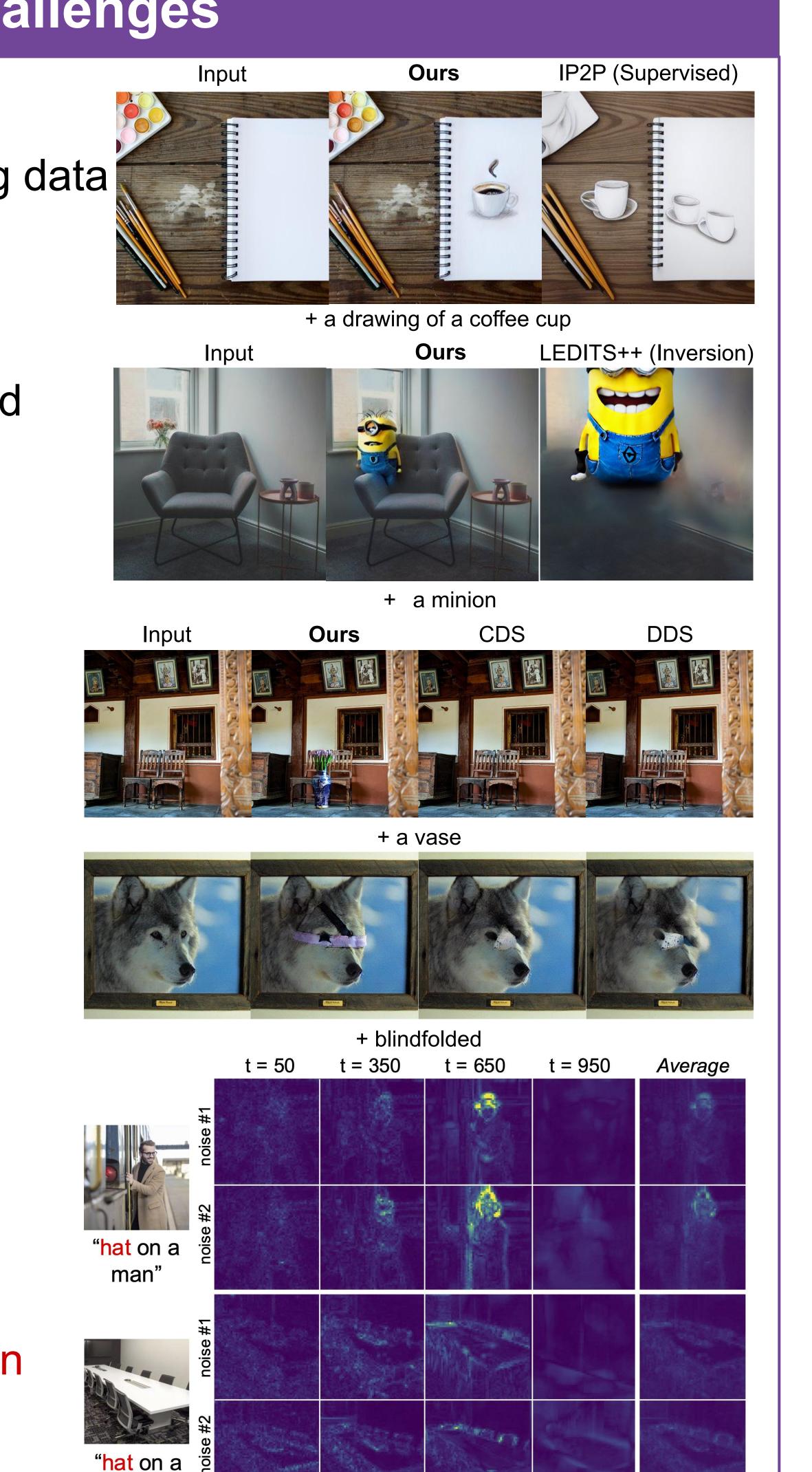
- relies solely on implicit, derived mask from attention features
- lacks explicit background preservation constraints

3. Score distillation

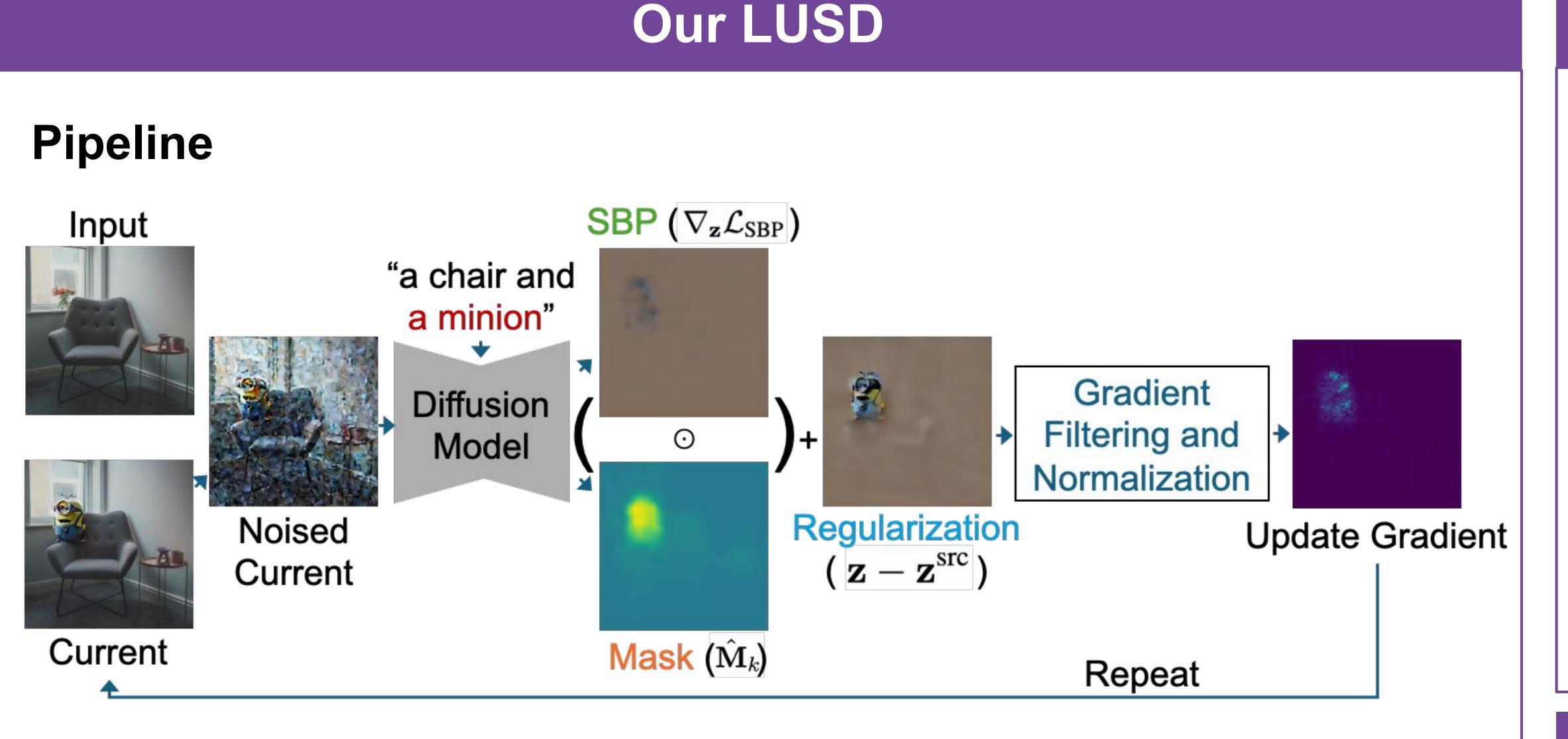
- highly sensitive to prompts, inputs, noise seeds
- conflicting gradient updates
- unstable optimization
- leads to poor object insertion

Our ideas

- stabilizes gradient updates in score distillation
- uses implicit mask to progressively focus gradient updates on the relevant region
- actively filters out counterproductive gradients



Gradient variation



Key ideas

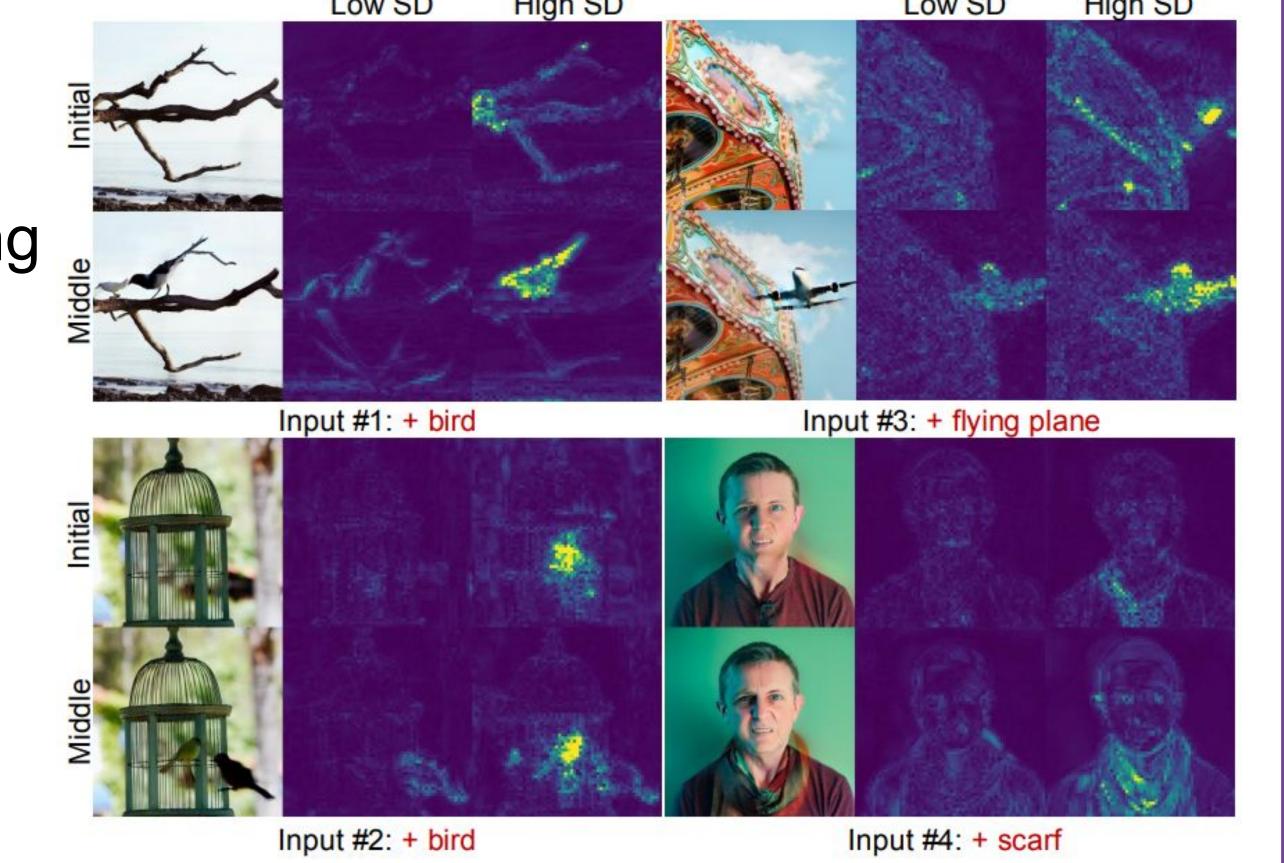
1. An implicit editing mask derived from attention features to reduce spatial variation of gradient updates

Score distillation formulation (SBP) [McAllister et al. 2024] $\nabla_{\mathbf{z}} \mathcal{L}_{\text{SBP-reg}} = (1 - \lambda)(\hat{\mathbf{M}}_k \odot \nabla_{\mathbf{z}} \mathcal{L}_{\text{SBP}}) + \lambda(\mathbf{z} - \mathbf{z}^{\text{src}})$ Background nask from cross- &

self-attention features

2. Gradient Filtering and Normalization to filter out "counterproductive" gradients using thresholding when the magnitude standard deviation is low.

$$\nabla_{\mathbf{z}} \mathcal{L}_{\text{LUSD}} = \gamma \frac{\nabla_{\mathbf{z}} \mathcal{L}_{\text{SBP-reg}}}{\text{SD}(\nabla_{\mathbf{z}} \mathcal{L}_{\text{SBP-reg}})}$$



Results – Quantitative

Better image editing measured by a user study and 4 metrics: one (CLIP-T) for prompt fidelity and three (CLIP-AUC, L1*, CLIP-I*) for background preservation

User preference (lower = ours wins)					Method	Time (mins)	CLIP-T↑	CLIP-AUC	↑ L1 * ↓	CLIP-I* ↑
Method	Background Prompt Quality Overall				Instruction-guided methods					
	Dackground	Trompt	Quanty		IP2P [6]	0.06	0.275	0.053	0.029	0.180
IP2P [6]	33.5%	40.0%	36.5%	36.0%	HIVE [51]	0.13	0.272	0.040	0.024	0.189
HIVE [51]	47.0%	40.5%	45.0%	39.0%	Global description-guided methods					
LEDITS++ [5]	35.5%	33.0%	37.0%	35.0%	LEDITS++ [5]	0.13	0.279	0.067	0.022	0.182
DDS [16]	43.5%	37.0%	38.0%	38.5%	DDS [16]	0.22	0.277	0.048	0.017	0.195
CDS [30]	44.5%	40.0%	43.0%	42.0%	CDS [30]	0.62	0.272	0.034	0.016	0.197
	88.90 Pen - Nagarangan		989909999 - 103109799		SBP [26]	0.30	0.285	0.068	0.024	0.174
SBP [26]	38.3%	40.2%	38.5%	42.3%	Ours	1.79	0.287	0.074	0.015	0.192

Results – Qualitative

