

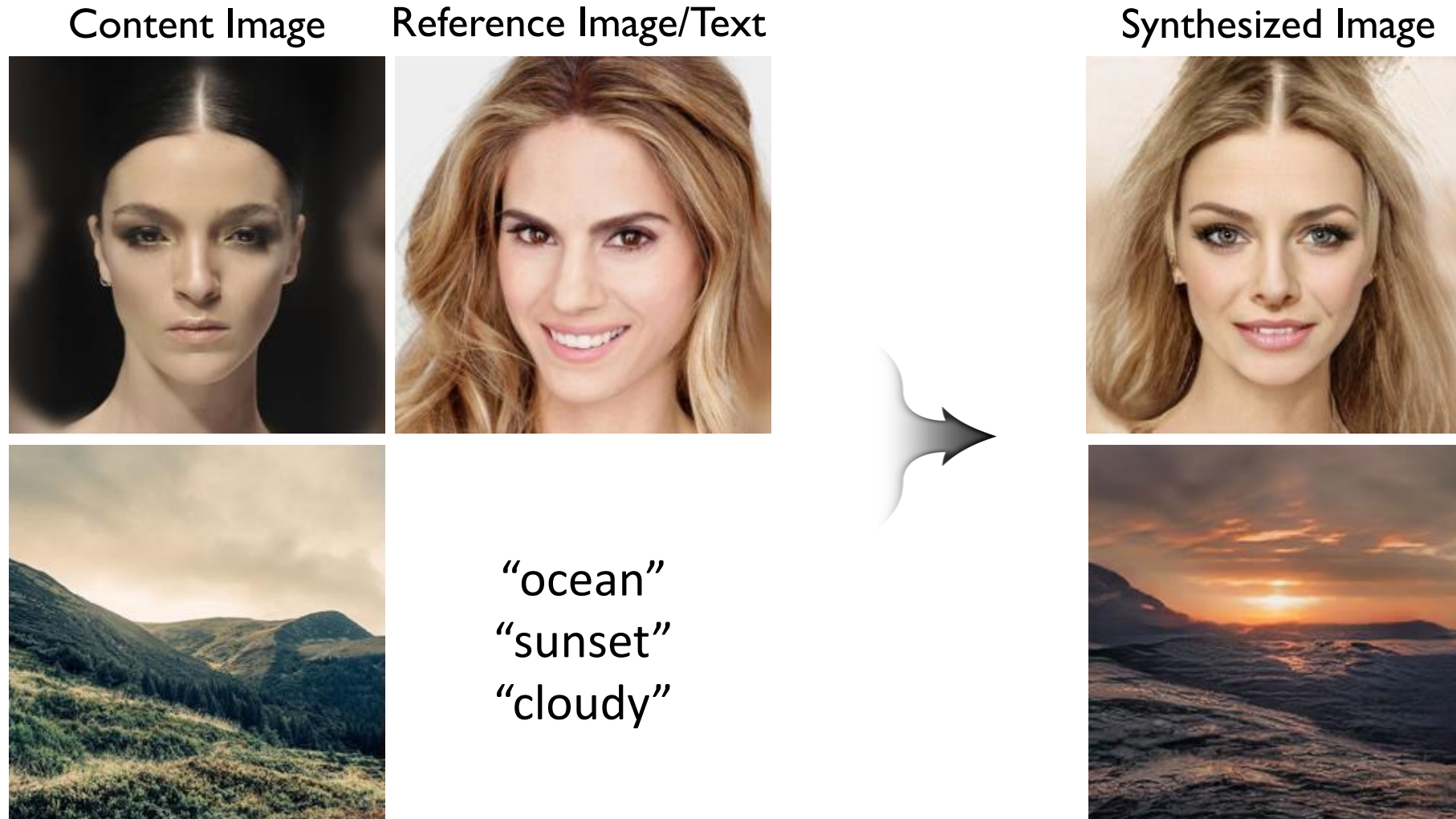
# LANIT: Language-Driven Image-to-Image Translation for Unlabeled Data

## CVPR 2023

Jihye Park\*, Sunwoo Kim\*, Soohyun Kim\*,  
Seokju Cho, Jaejun Yoo, Youngjung Uh, Seungryong Kim  
Korea University



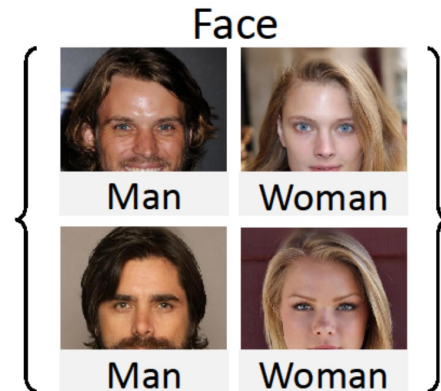
# What is Image-to-Image Translation



# Motivation and Problem Formulation

## Per-sample-level

- Each image is annotated corresponding classes through manual labeling

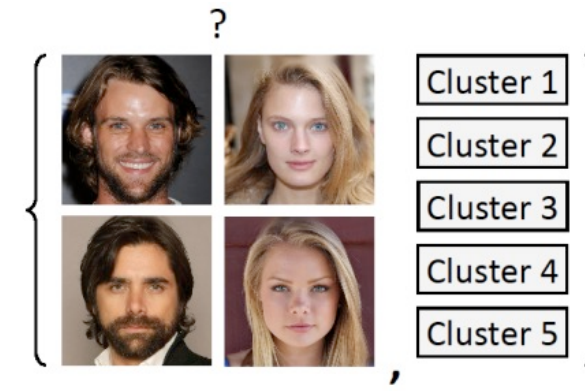


✗ *Intensive labeling resource*

*StarGANv2 (CVPR'20)*

## Unsupervised

- Pseudo labels from clustering or SSL of unsupervised manner



✓ *No ground-truth labels*

*TUNIT (ICCV'21)*

*Styla-aware Discriminator. (CVPR'22)*

# Motivation and Problem Formulation



Clusters do not have semantic meaning

Designated candidate domains have semantic meanings

# Motivation and Problem Formulation

## Per-sample-level

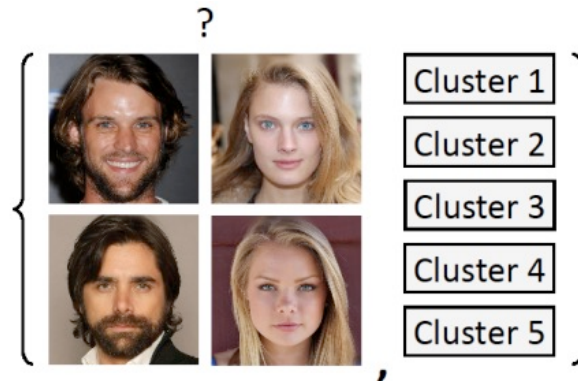
- Each image is annotated corresponding classes through manual labeling



- ✗ *Intensive labeling resource*
- ✓ *Applicable: we can accurately target the attributes to synthesize*

## Unsupervised

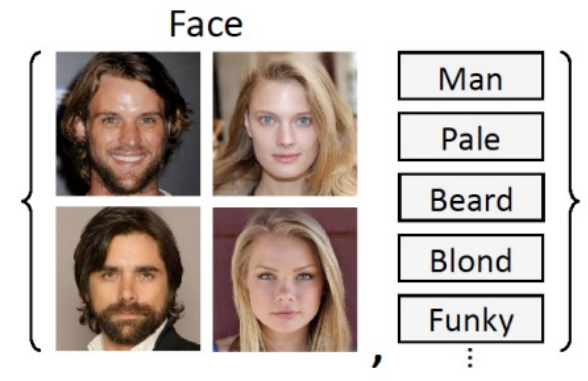
- Pseudo labels from clustering or SSL of unsupervised manner



- ✓ *No ground-truth labels*
- ✗ *Not applicable: we can't accurately target the attributes to synthesize*

## Dataset-level (Ours)

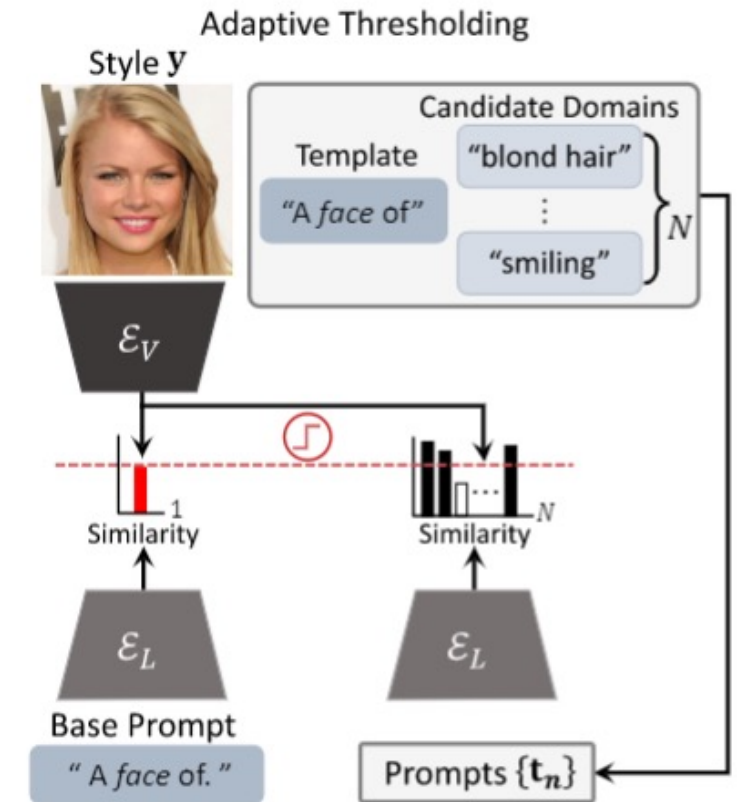
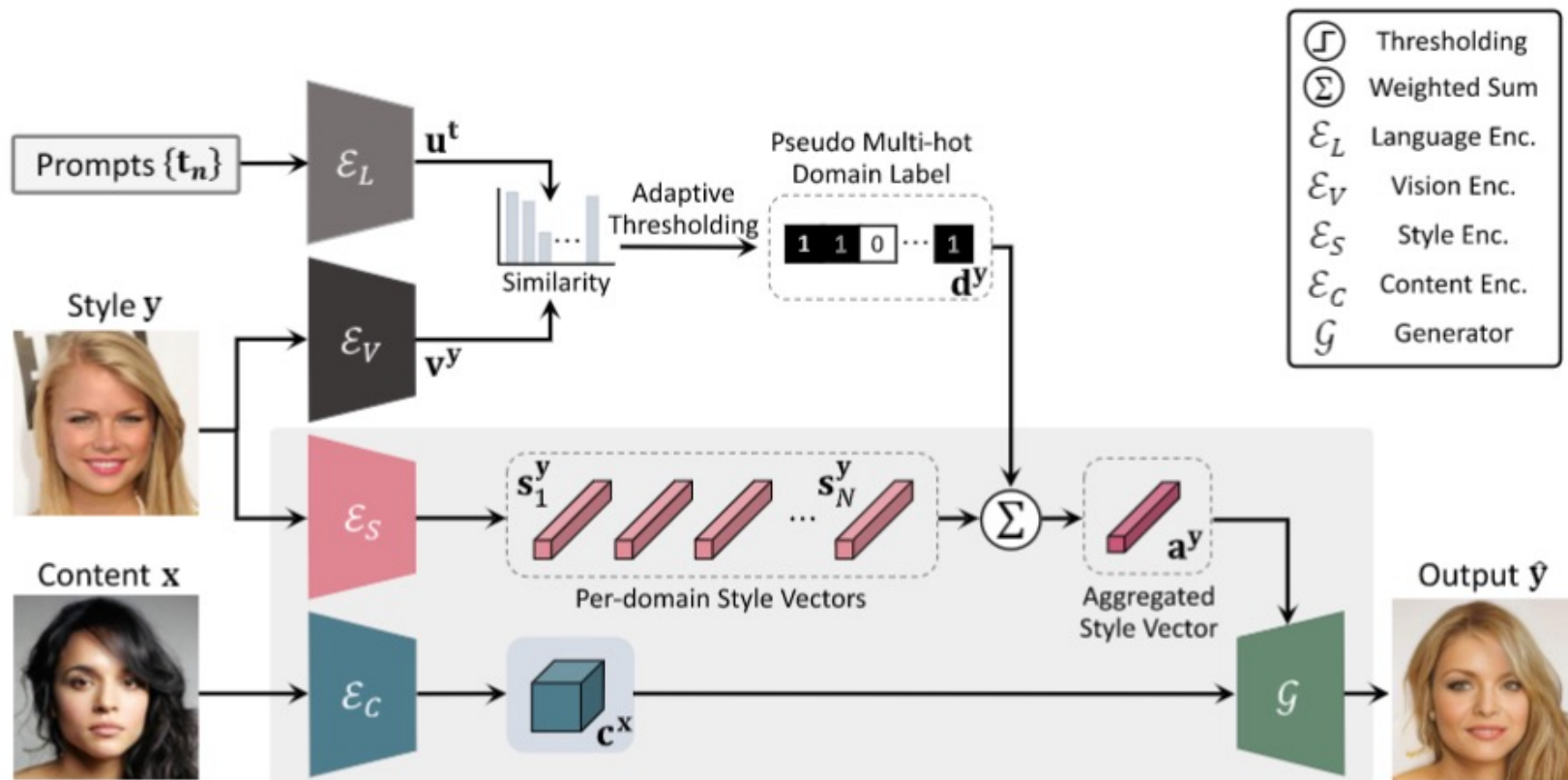
- Pseudo labels from **Weak human's supervision** and **large VL model**



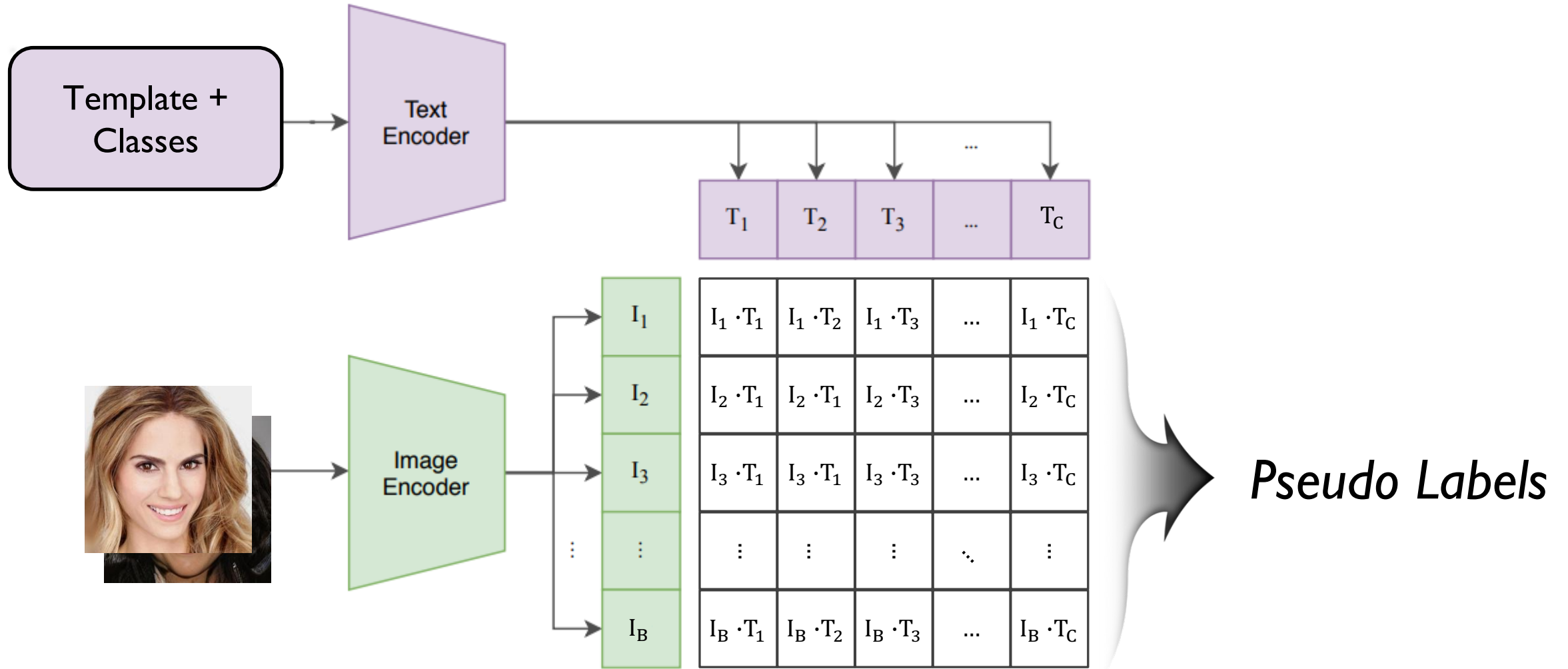
- ✓ *No ground-truth labels*
- ✓ *Applicable: we can accurately target the attributes to synthesize*



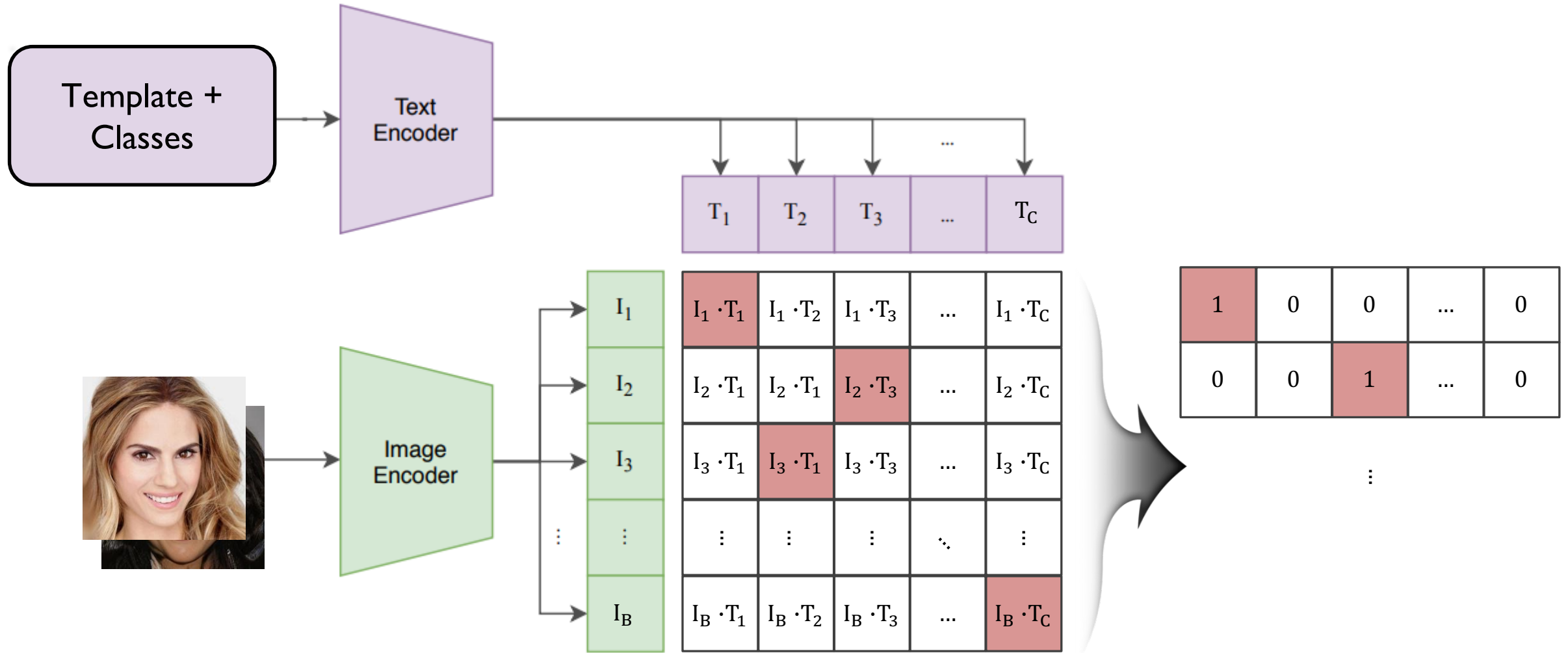
# Proposed Network



# How to make pseudo labels?

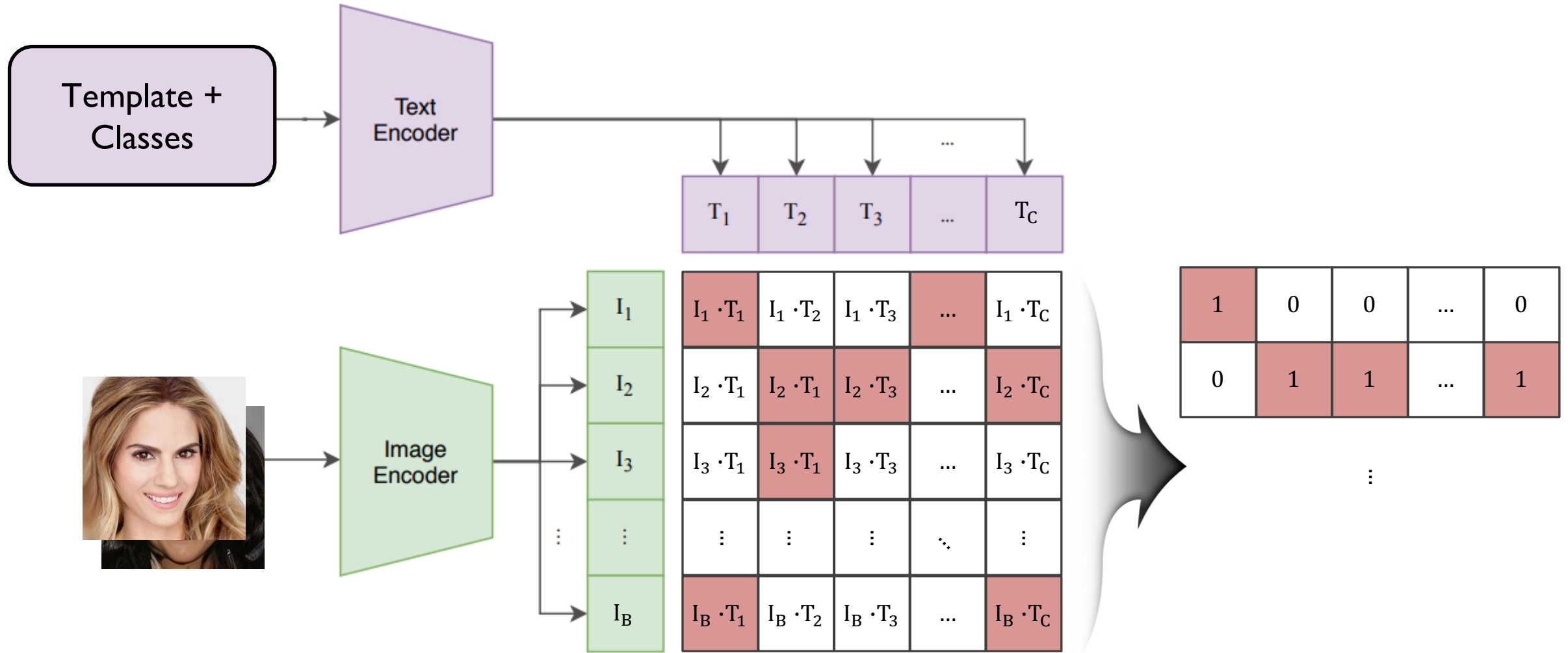


# How to make pseudo labels?: Topk (k=1)



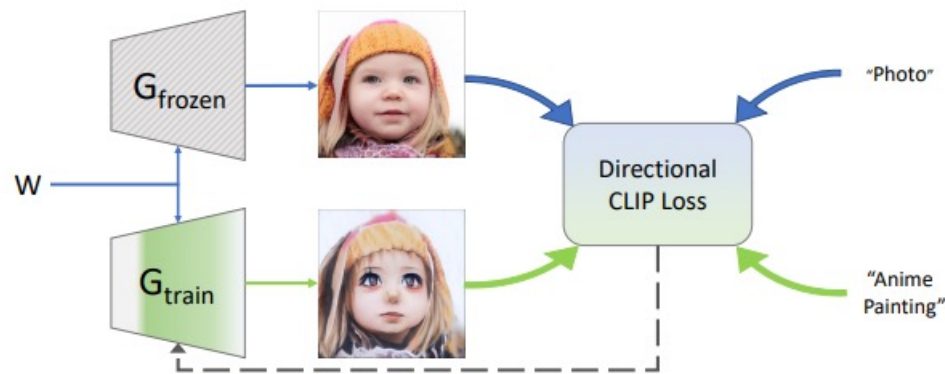


# How to make pseudo labels?: Thres ( $\tau=0.2$ )



# How to make pseudo labels?: Adaptive Thres

## Direction cliploss



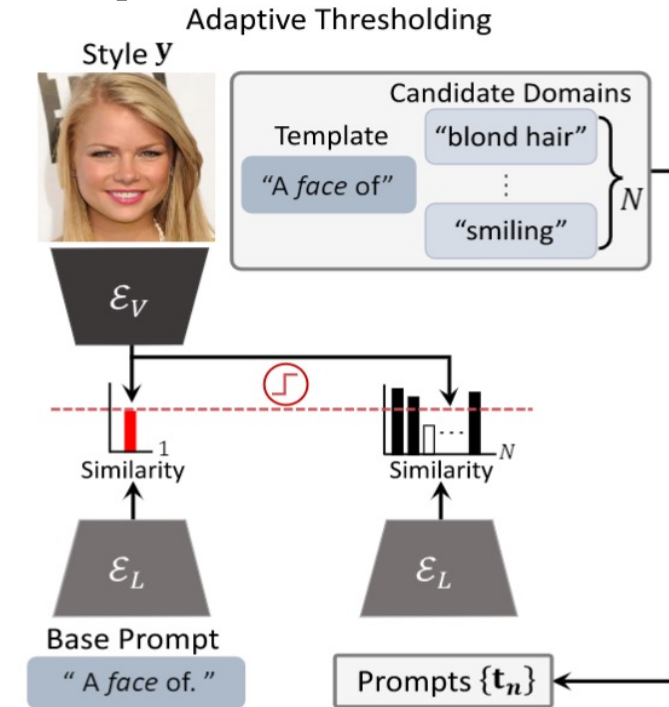
$$\Delta T = E_T(t_{target}) - E_T(t_{source}),$$

$$\Delta I = E_I(G_{train}(w)) - E_I(G_{frozen}(w)),$$

$$\mathcal{L}_{direction} = 1 - \frac{\Delta I \cdot \Delta T}{|\Delta I| |\Delta T|}.$$

StyleGAN-NADA (SIGGRAPH'22)

## Adaptive Thresholding

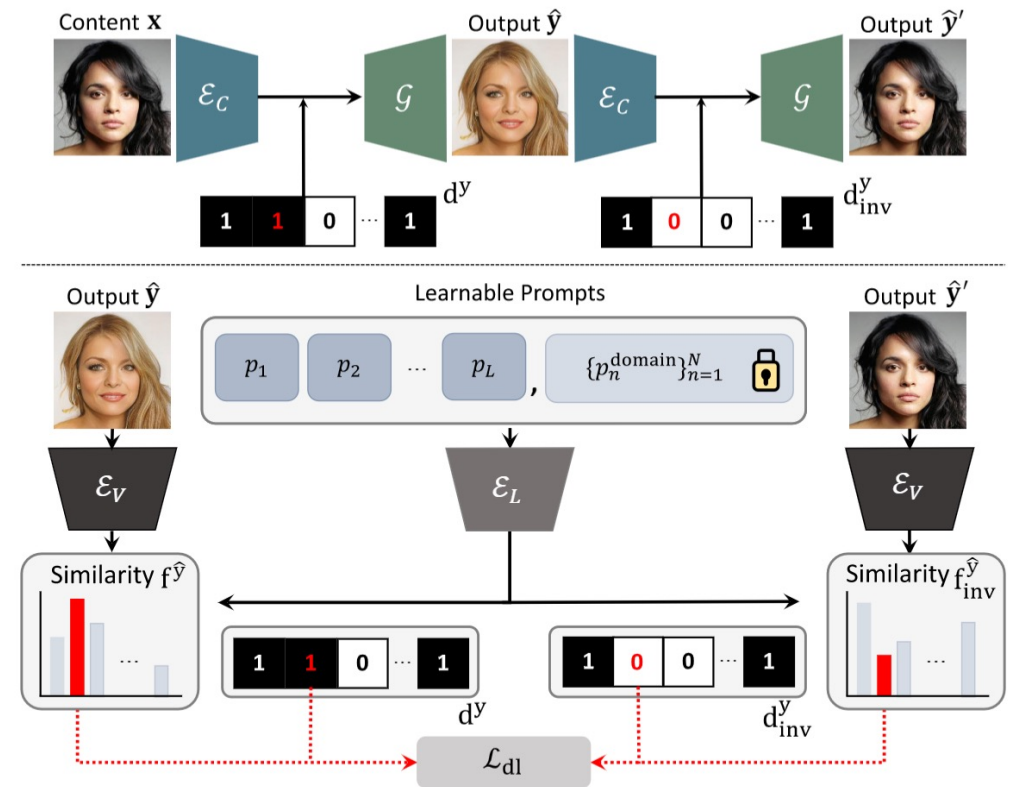


# How to improve performance of CLIP?

## Template Augmentation

[  
 ['a face photo with.'],  
 ['a face photo of the.'],  
 ['the face photo of the.'],  
 ['a good face photo of the.'],  
 ["high quality face photo of."],  
 ["a face image of."],  
 ["the face image of."],  
 ["high quality face image of."],  
 ["a high quality face image of."],  
 ]

## Prompt learning



# Loss Functions

---

Adversarial Loss( $\mathcal{L}_{adv}$ )

$$= \mathbb{E}_{\mathbf{x}, \mathbf{y}} \sum_{n=1}^N [\log \mathcal{D}_n(\mathbf{y})d_n^{\mathbf{y}} + \log(1 - \mathcal{D}_n(\mathcal{G}(\mathbf{x}, \mathbf{a}^{\mathbf{y}}))d_n^{\mathbf{y}})]$$

n'th Discriminator:  $\mathcal{D}_n(\cdot)$

Binary Cross Entropy:  $\mathcal{H}(\cdot, \cdot)$

Domain Regularization Loss( $\mathcal{L}_{dl}$ )

$$\mathcal{L}_{dl} = \mathcal{H}(d_n^{\mathbf{y}}, f_n^{\hat{\mathbf{y}}}) + \mathcal{H}(d_{inv, n}^{\mathbf{y}}(n), f_{inv, n}^{\hat{\mathbf{y}}}).$$

Cycle-Consistency Loss( $\mathcal{L}_{cyc}$ )

$$\mathcal{L}_{cyc} = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\|\mathbf{x} - \mathcal{G}(\mathbf{c}^{\hat{\mathbf{y}}}, \mathbf{a}^{\mathbf{x}})\|_1]$$

Style Reconstruction Loss( $\mathcal{L}_{sty}$ )

$$\mathcal{L}_{sty} = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\|\mathbf{s}^{\mathbf{y}} - \mathcal{E}_S(\hat{\mathbf{y}})\|_1]$$

Overall Objective( $\mathcal{L}_{total}$ )

$$\mathcal{L}_{total} = \lambda_{adv}\mathcal{L}_{adv} + \lambda_{dl}\mathcal{L}_{dl} + \lambda_{cyc}\mathcal{L}_{cyc} + \lambda_{sty}\mathcal{L}_{sty}$$

# Quantitative Results

*Our proposed techniques improve the performance of CLIP in our framework !*

N	AnimalFaces-10 [6]					CelebA-HQ [7]				
	Top-1	Top-3	Baseline	TextAug	Prompt learning	Top-1	Top-3	Baseline	TextAug	Prompt learning
4	0.762	0.672	0.678	0.796	<b>0.832</b>	0.372	0.421	0.423	0.435	<b>0.481</b>
7	0.903	0.701	0.688	0.862	<b>0.893</b>	0.423	0.613	0.610	0.631	<b>0.652</b>
10	0.956	0.723	0.693	0.835	<b>0.880</b>	0.355	0.562	0.610	0.638	<b>0.670</b>
13	0.826	0.654	0.606	0.785	<b>0.801</b>	0.293	0.533	0.591	0.612	<b>0.639</b>
16	0.753	0.630	0.601	0.753	<b>0.783</b>	0.300	0.522	0.562	0.613	<b>0.641</b>

Table 4. F1 score on the varying the number of domains.



# Quantitative Results

## Comparison to other works

Method	CelebA-HQ [41]		AnimalFaces-10 [35]		Food-10 [7]	
	mFID	D&C	mFID	D&C	mFID	D&C
StarGAN2 [10] (sup.)	32.16	1.22 / <b>0.446</b>	<b>33.67</b>	<b>1.54 / 0.91</b>	65.03	1.09 / 0.76
Smoothing [40] (sup.)	35.93	<b>1.25</b> / 0.431	38.93	0.97 / 0.75	61.13	0.96 / 0.68
TUNIT [3] (unsup.)	61.29	0.24 / 0.13	47.70	1.04 / 0.81	52.2	1.08 / <b>0.88</b>
Kim <i>et al.</i> [29] (unsup.)	41.33	0.60 / 0.241	36.83	1.06 / 0.82	<b>49.34</b>	1.06 / 0.80
LANIT	<b>27.96</b>	0.91 / 0.34	34.11	1.46 / 0.89	49.50	<b>1.24</b> / 0.86

## Ablation study

( get candidate dominas form dictionary )

Datasets	Default		Dictionary	
	mFID	D&C	mFID	D&C
CelebA-HQ [15]	27.96	0.91/0.34	28.34	0.77/0.23
AnimalFaces-10 [35]	34.11	1.46/0.89	40.48	1.01/0.78
Food-10 [7]	48.08	1.24/0.86	49.50	1.17/0.81

## Ablation study( #of domain to train )

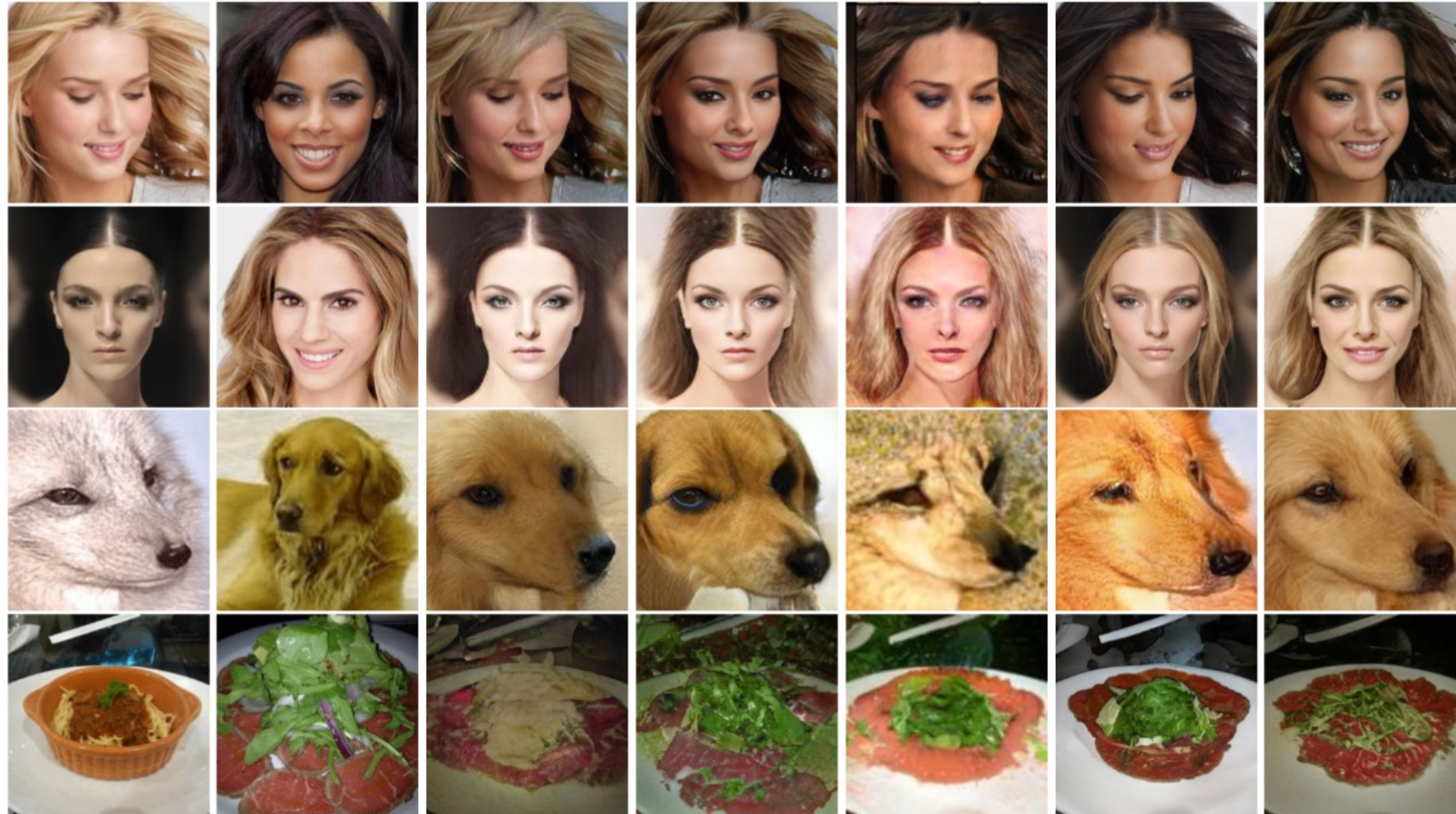
N	Method	AnimalFaces-10 [35]		CelebA-HQ [41]	
		mFID	D&C	mFID	D&C
4	TUNIT	77.7	0.88 / 0.74	61.5	0.24 / 0.12
	LANIT	71.6	1.35 / 0.46	49.3	0.33 / 0.14
7	TUNIT	62.7	1.02 / 0.73	54.7	0.33 / 0.16
	LANIT	49.9	1.47 / 0.66	43.2	0.44 / 0.19
10	TUNIT	47.7	1.04 / 0.81	61.3	0.24 / 0.13
	LANIT	<b>34.1</b>	1.46 / <b>0.89</b>	27.96	<b>0.91 / 0.34</b>
13	TUNIT	56.8	0.99 / 0.72	98.9	0.08 / 0.03
	LANIT	30.13	1.43 / 0.85	34.8	0.58 / 0.21
16	TUNIT	54.1	1.09 / 0.78	127.7	0.04 / 0.02
	LANIT	35.8	<b>1.49</b> / 0.82	<b>27.92</b>	0.76 / 0.23

## Ablation study( components )

Method	CelebA-HQ [41]		
	mFID	D&C	F1.
A LANIT (Top-1)	49.65	0.56 / 0.32	0.347
B LANIT (Top-3)	41.68	0.68 / 0.30	0.564
C Baseline	33.79	0.70 / 0.26	0.613
D + DL Loss	29.21	0.86 / 0.32	0.632
E + Prompt Learning	<b>27.96</b>	<b>0.91 / 0.34</b>	<b>0.718</b>



# Qualitative Results: Comparison



(a) Content

(b) Style

(c) StarGANv2 [10]

(d) Smoothing [40]

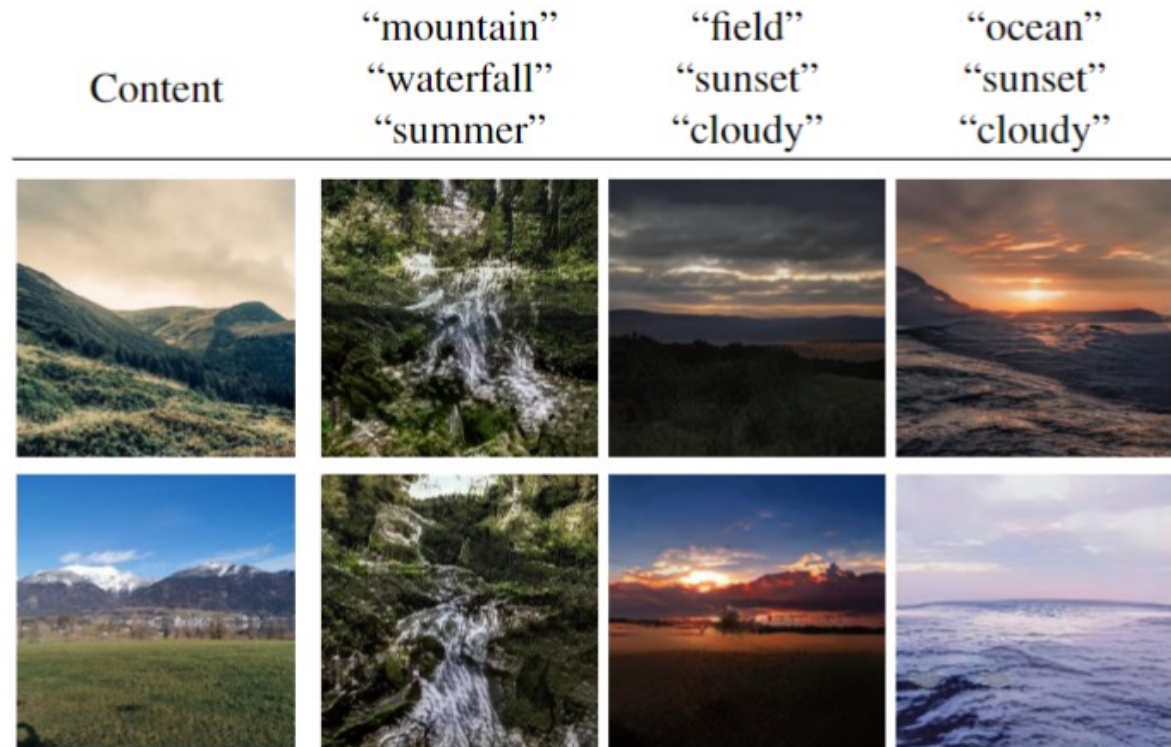
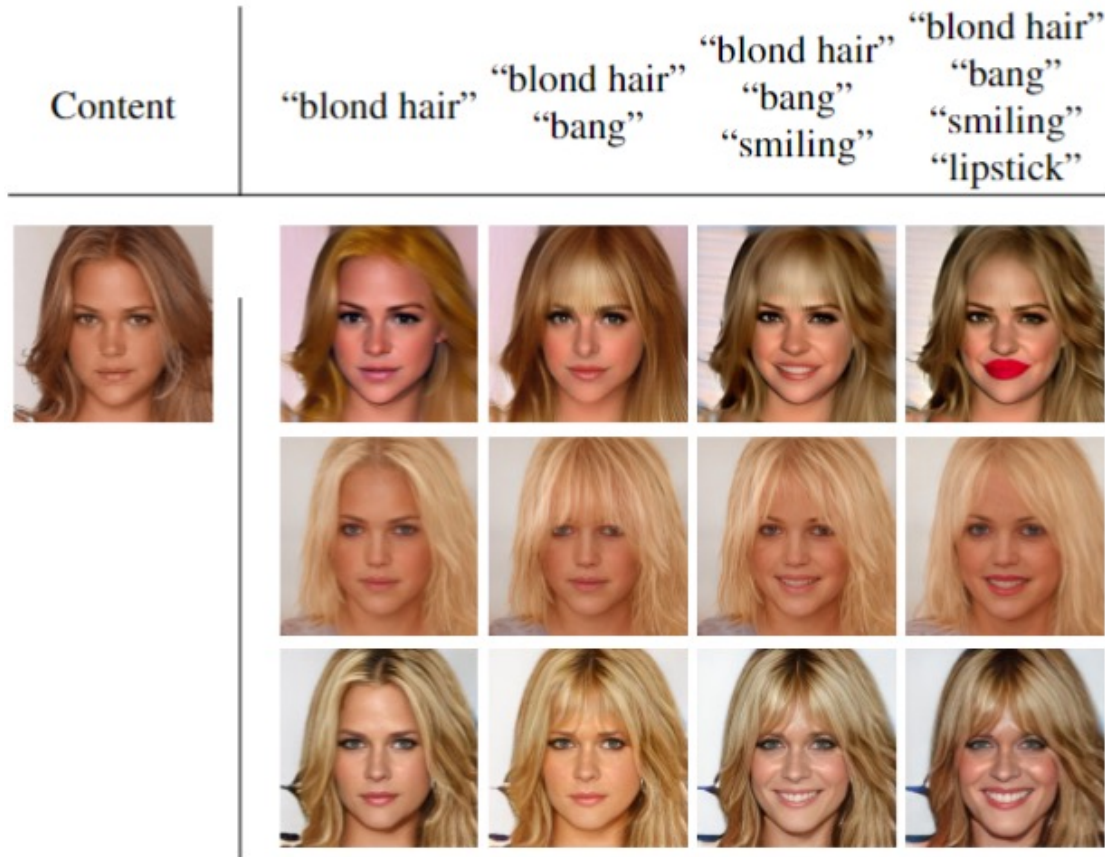
(e) TUNIT [3]

(f) Kim *et al.* [29]

(g) Ours



# Latent Guided Manipulation



# Reference Guided Manipulation



((a)) Content

((b)) Style

((c)) TUNIT [1]

((d)) Kim *et al.* [5]

((e)) Ours(K=1)

((f)) Ours(K=2)

((g)) Ours(K=3)



# Qualitative Results on Various Datasets

## AnimalFaces-10



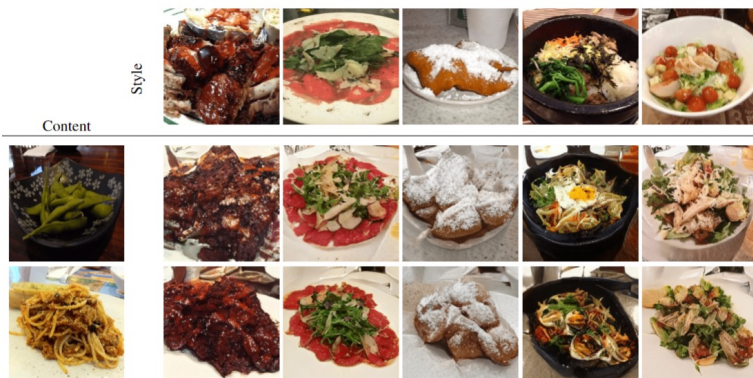
## MetFace



## LSUN-Church



## Food-10



## Anime



## LSUN-Car





---

Thank you!

