

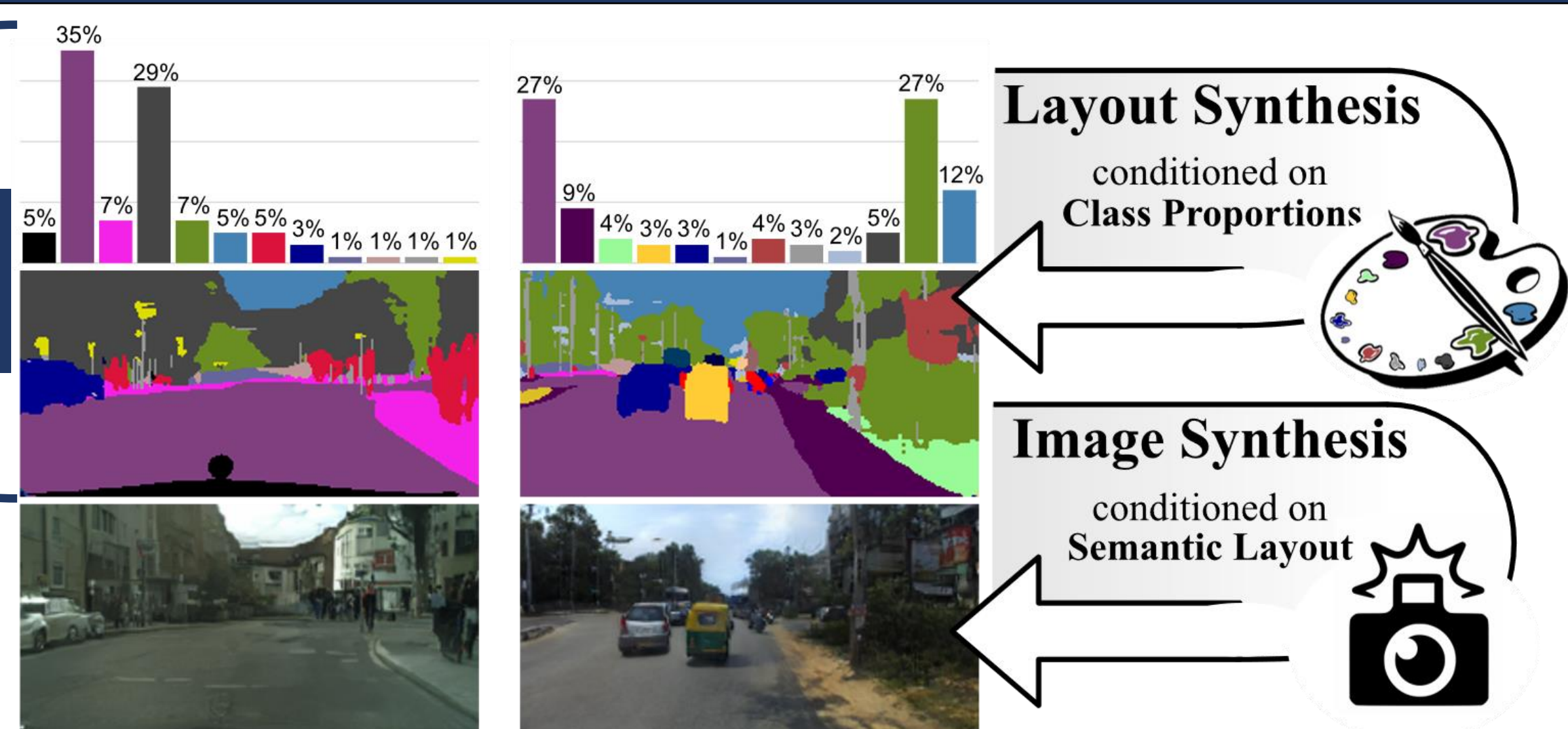
# Semantic Palette: Guiding Scene Generation with Class Proportions

Guillaume Le Moing<sup>1,†</sup>, Tuan-Hung Vu<sup>2</sup>, Himalaya Jain<sup>2</sup>, Patrick Pérez<sup>2</sup>, Matthieu Cord<sup>2</sup>

<sup>1</sup>Inria, École normale supérieure, CNRS, PSL Research University, France <sup>2</sup>Valeo.ai <sup>†</sup>work performed during an internship at Valeo

## Introduction

Our focus



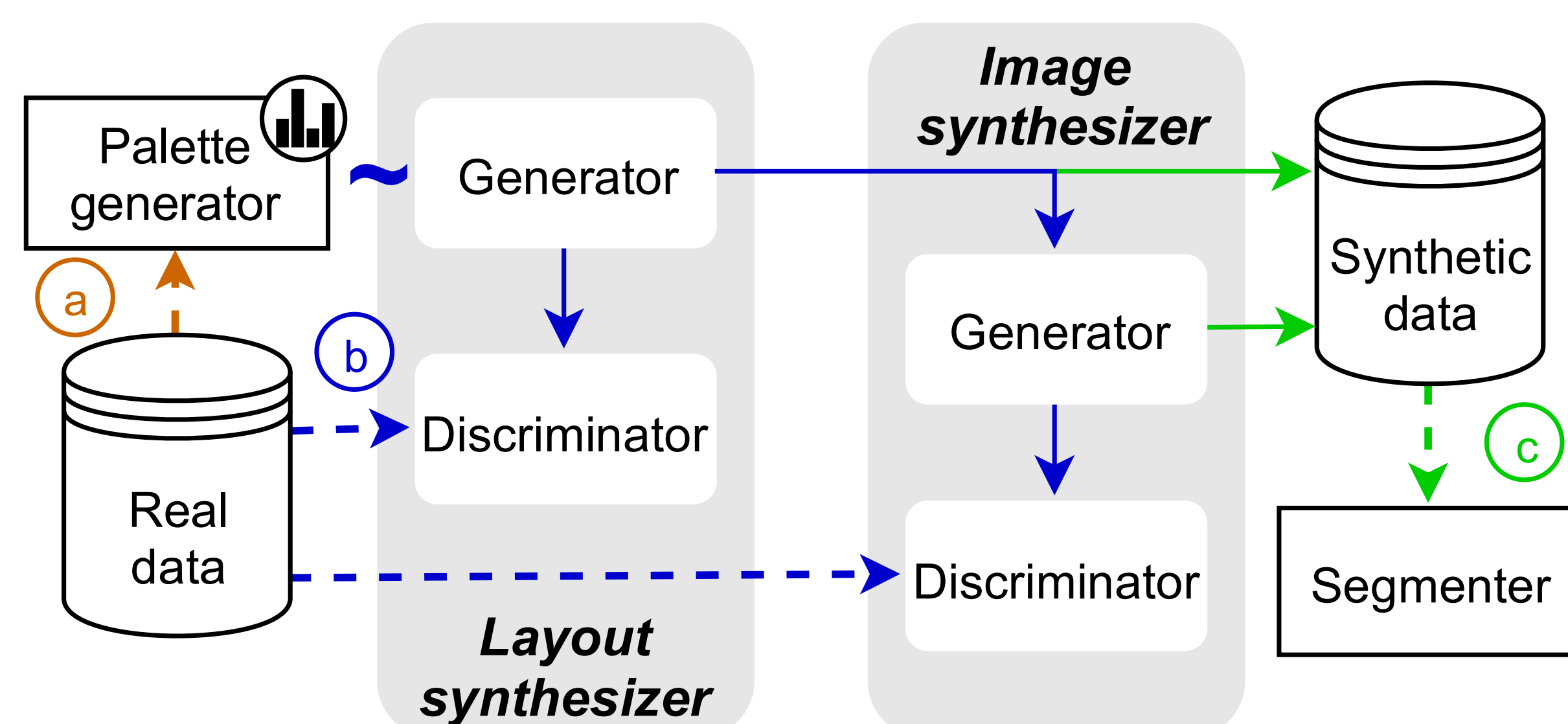
Towards controllable quality scene generation, we propose:

1. the task of layout synthesis conditioned on class proportions,
2. the combination of a layout generator and an image generator to produce images of greater practical use and quality,
3. the extension of our framework to partial editing of scenes.

## Method

### Overall framework

- Learn **palettes** (i.e., semantic class histograms) distribution.
- Train layout/image synthesizers individually, then end-to-end.
- Use synthetic pairs for downstream tasks (e.g., segmentation)



## Results

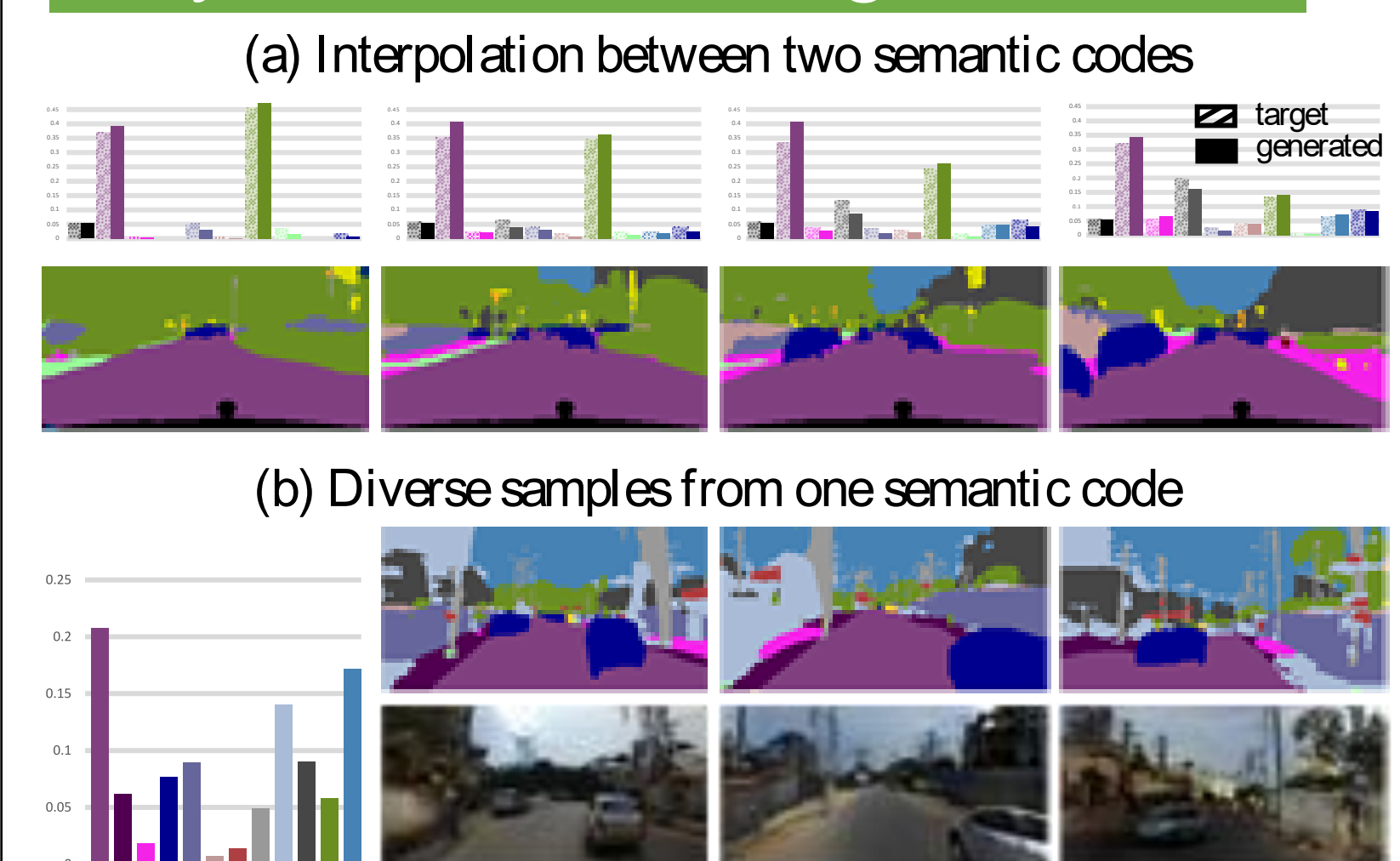
### Metrics

**KL**: KL-divergence between syn. and target palette  
**FID**: Fréchet inception distance in image space  
**FSD**: Fréchet inception distance in palette space  
**GAN-train**: train on synthetic, test on real  
**GAN-test**: train on real, test on synthetic

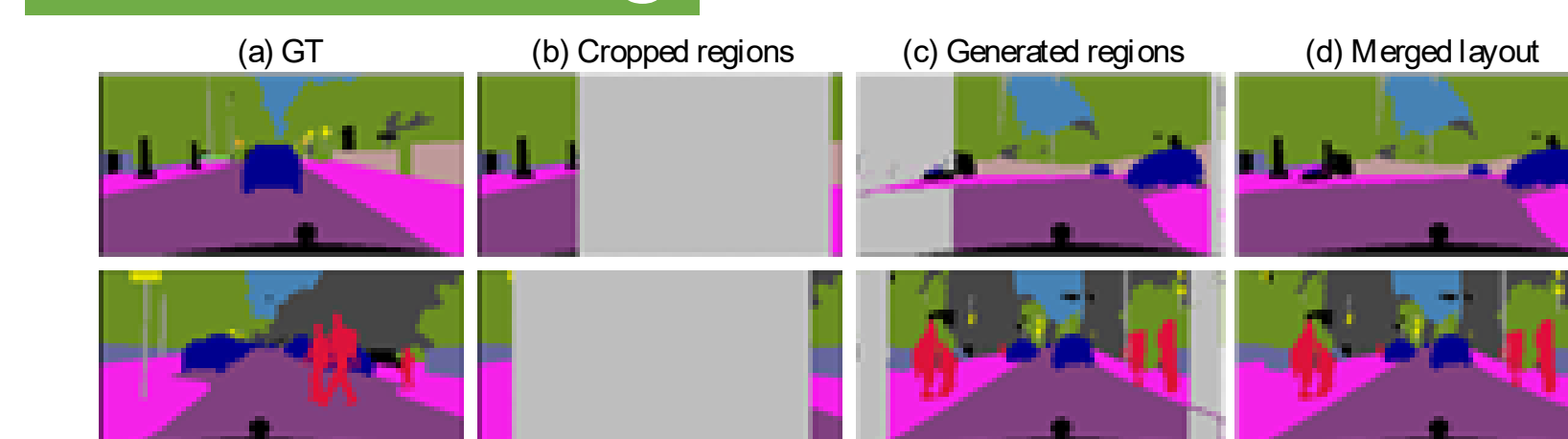
### Palette vs. Conditional GANs

Method	Layout KL ↓	Image FID ↓	GAN-test		GAN-train	
			mIoU <sup>*</sup>	mIoU	mIoU <sup>*</sup>	mIoU
Baseline 1	1.17	69.2	33.7	42.8	29.6	38.5
Baseline 2	0.32	69.0	35.3	46.9	30.2	39.4
Sem. Palette	<b>0.07</b>	60.7	34.6	45.7	30.6	40.1
Sem. Palette <sub>e2e</sub>	0.08	<b>51.0</b>	<b>36.8</b>	<b>48.6</b>	<b>33.3</b>	<b>44.5</b>
Oracle	-	28.2	-	-	36.9	48.1

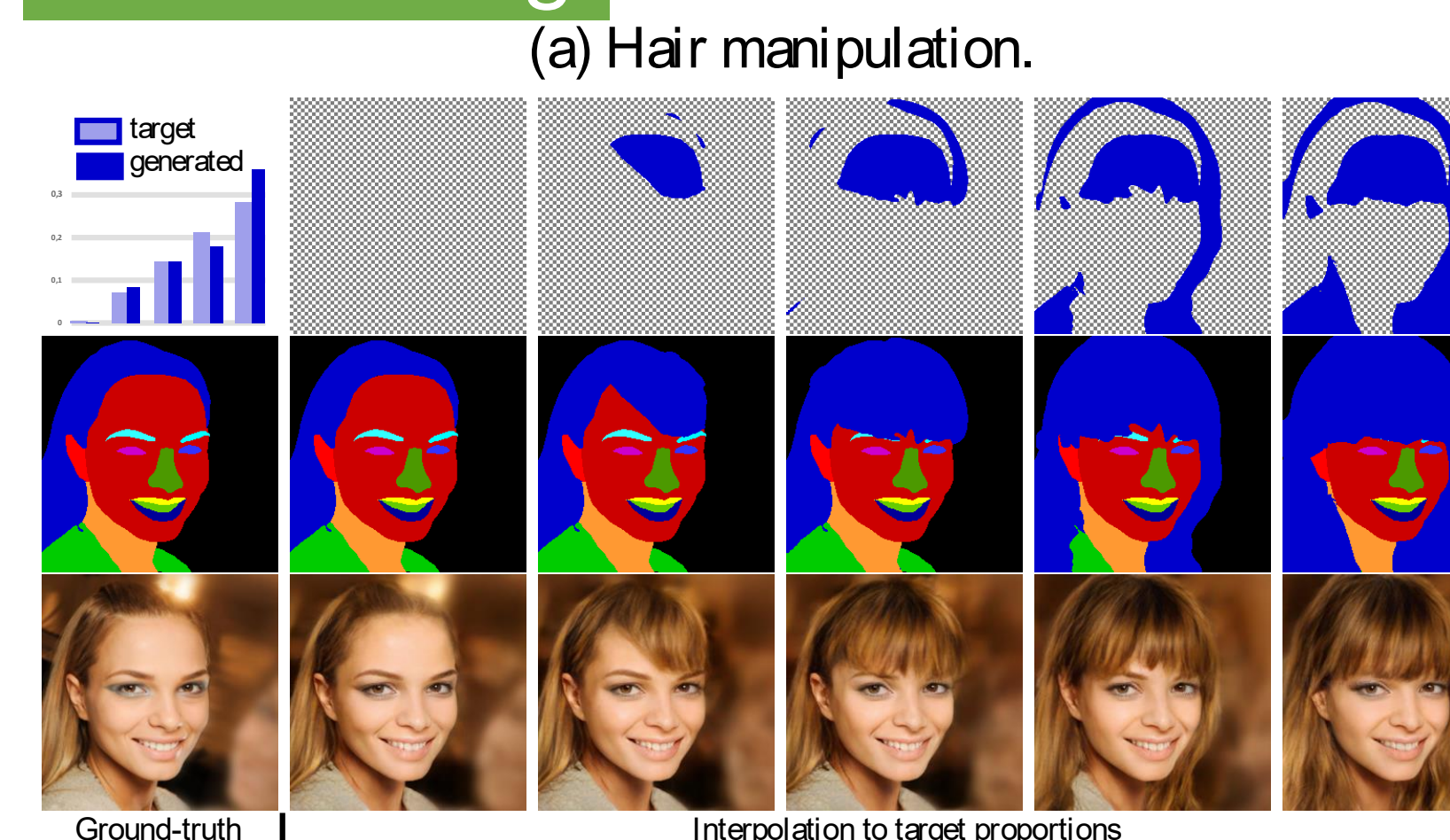
### Layout-and-scene generation



### Partial editing



### Face editing



### Palette vs. Unconditional GANs

Method	Layout FSD ↓	Image FID ↓	GAN-test		GAN-train	
			mIoU <sup>*</sup>	mIoU	mIoU <sup>*</sup>	mIoU
PCGAN	63.8	85.7	30.4	39.0	28.2	35.7
SB-GAN	63.8	71.0	31.8	41.2	28.8	37.2
Sem. Palette	<b>25.3</b>	<b>60.7</b>	<b>34.6</b>	<b>45.7</b>	<b>30.6</b>	<b>40.1</b>
SB-GAN <sub>e2e</sub>	20.4	61.8	34.5	44.7	29.6	37.0
Sem. Palette <sub>e2e</sub>	<b>11.8</b>	<b>51.0</b>	<b>36.8</b>	<b>48.6</b>	<b>33.3</b>	<b>44.5</b>
Oracle	-	28.2	-	-	36.9	48.1

### Data augmentation

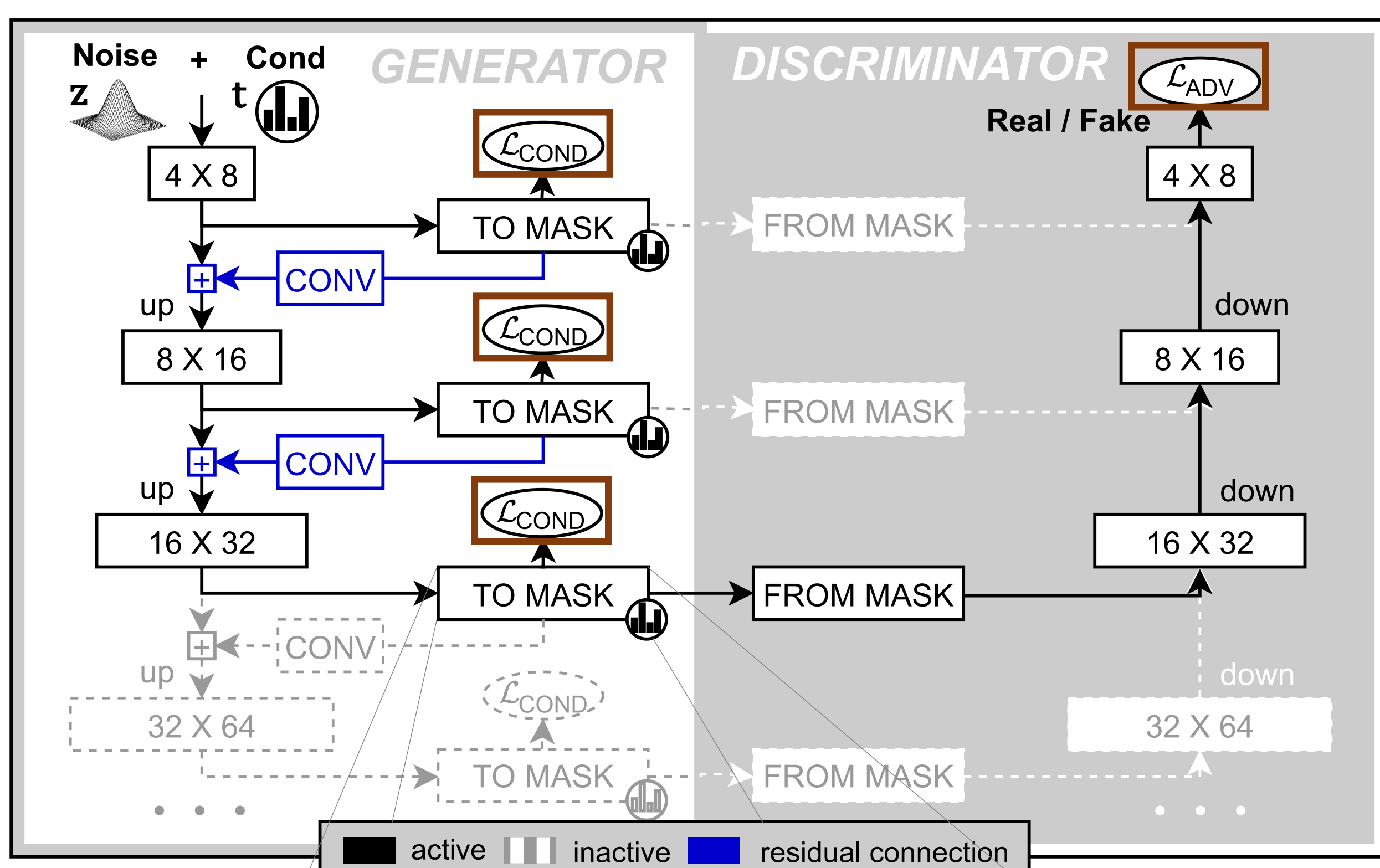
Data	Method	(a) Cityscapes		(b) IDD	
		mIoU <sup>*</sup>	mIoU	mIoU <sup>*</sup>	mIoU
Real	Baseline	36.9	48.1	33.8	43.8
Real + Semi-Syn	CauGAN	37.2 <sub>±0.3</sub>	48.2 <sub>±0.1</sub>	33.6 <sub>±0.2</sub>	43.5 <sub>±0.3</sub>
Real + Syn	SB-GAN	34.6 <sub>±2.3</sub>	45.5 <sub>±2.6</sub>	33.5 <sub>±0.3</sub>	43.4 <sub>±0.4</sub>
	Sem. Palette	38.0 <sub>±1.1</sub>	49.4 <sub>±1.3</sub>	33.8 <sub>±0.7</sub>	43.8 <sub>±0.9</sub>
	Sem. Palette (DA)	38.6 <sub>±1.7</sub>	51.6 <sub>±3.5</sub>	34.5 <sub>±0.7</sub>	44.7 <sub>±0.9</sub>
	Sem. Palette (Part.)	<b>40.7<sub>±3.8</sub></b>	<b>51.9<sub>±3.8</sub></b>	<b>35.6<sub>±1.8</sub></b>	<b>46.1<sub>±2.3</sub></b>
	Sem. Palette (Part. + DA)	<b>40.7<sub>±3.8</sub></b>	<b>52.6<sub>±4.5</sub></b>	<b>35.3<sub>±1.5</sub></b>	<b>45.8<sub>±2.0</sub></b>

## Conclusion

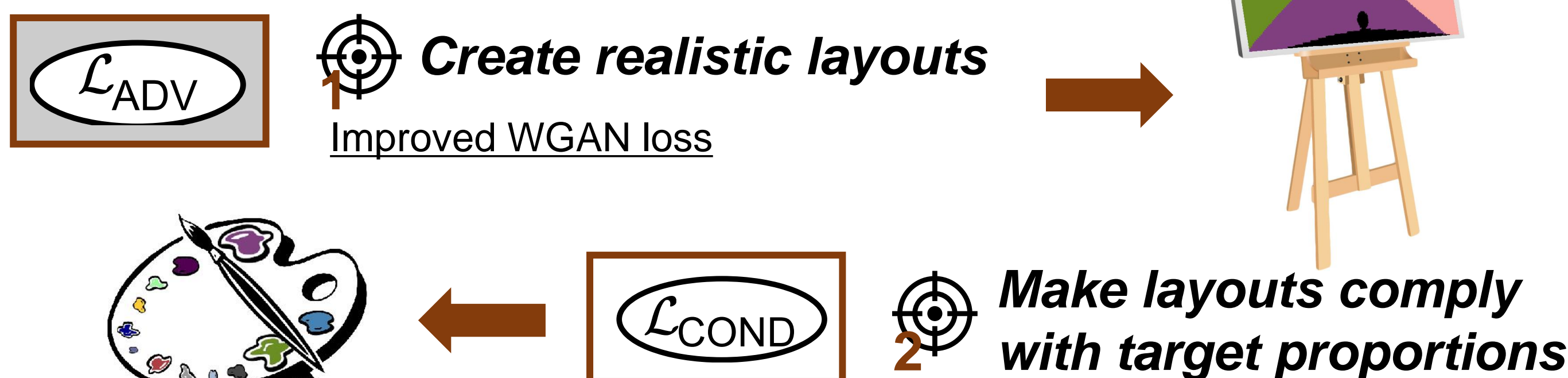
We have proposed **Semantic Palette**, a new framework for scene generation, and editing, guided by semantic proportions.

- Our novel architecture effectively accommodates class proportions while proposing plausible layouts, which then translate into realistic images.
- Semantic Palette better captures the distribution of real layouts / images than unconditional layout-and-scene GANs.
- Semantic Palette better follows target proportions and produces higher quality layouts than conditional baselines.
- Partial editing is an efficient data-augmentation strategy and opens up interesting applications like face editing.

## Conditional layout synthesis



## Learning objectives



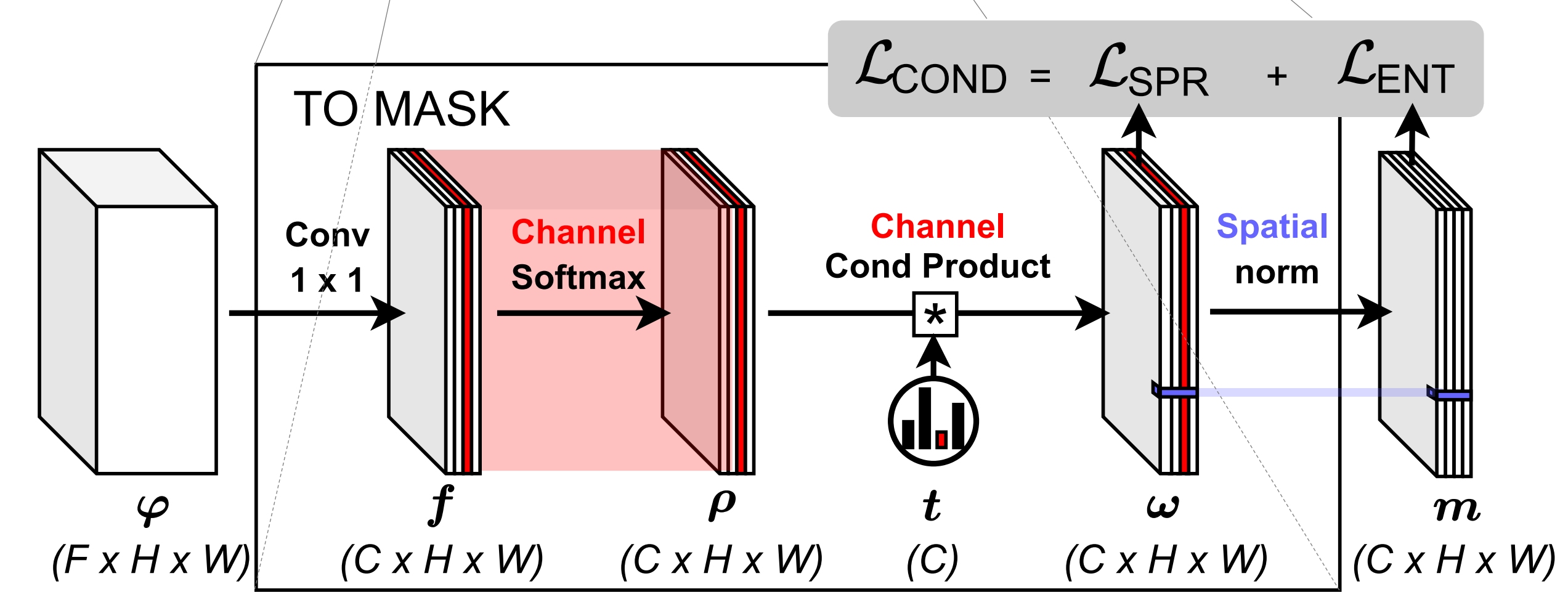
**Spread loss**: favor an even spatial semantic coverage

$$\mathcal{L}_{SPR} = \frac{1}{HW} \mathbb{E}_{(z,t)} \left[ \sum_{(i,j) \in \Omega} s_{i,j} \right], \quad s_{i,j} = \left( 1 - HW \sum_{c \in [1,C]} \omega_{c,i,j} \right)^2$$

**Entropy loss**: favor a peaky class distribution at each pixel

$$\mathcal{L}_{ENT} = \frac{1}{HW} \mathbb{E}_{(z,t)} \left[ \sum_{(i,j) \in \Omega} e_{i,j} \right], \quad e_{i,j} = - \sum_{c \in [1,C]} m_{c,i,j} \ln(m_{c,i,j})$$

## Semantically-assisted activation



## Partial editing of real layouts

