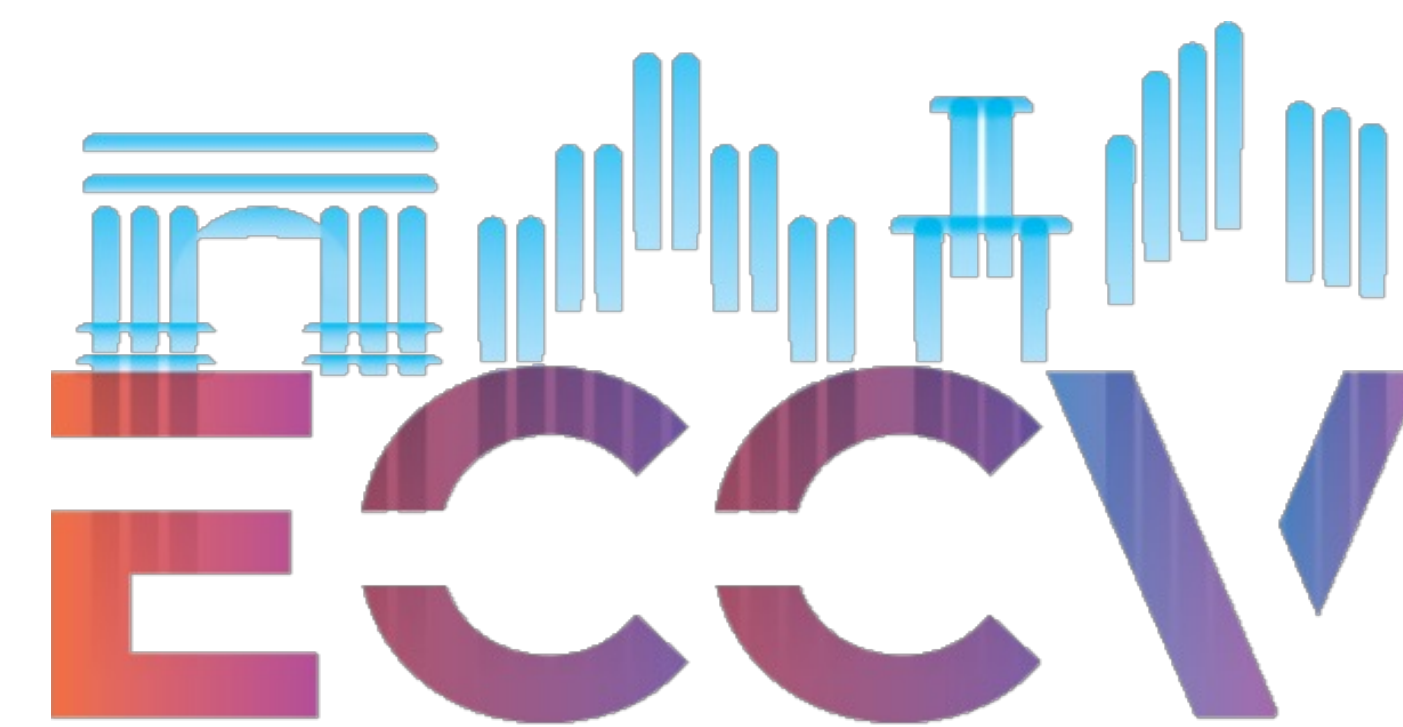




LayoutDETR: Detection Transformer Is a Good Multimodal Layout Designer

Ning Yu Chia-Chih Chen Zeyuan Chen Rui Meng
Gang Wu Paul Josel Juan Carlos Niebles Caiming Xiong
Salesforce Research

<https://ningyu1991.github.io/projects/LayoutDETR.html>



Motivations

- Graphic designs are at the foundation of communication between marketers and target audience.
- Graphic designs require designers' thoughtful understanding of multimodal inputs:
 - Background images
 - Multiple foreground texts
 - Multiple foreground product images
- Graphic designs require reasonable and aesthetically appealing compositions.
- Manual graphic designs are skill-demanding, time-consuming, and not scalable.

Goal

- Automate the multimodal layout design process by learning a layout generator.

Contributions

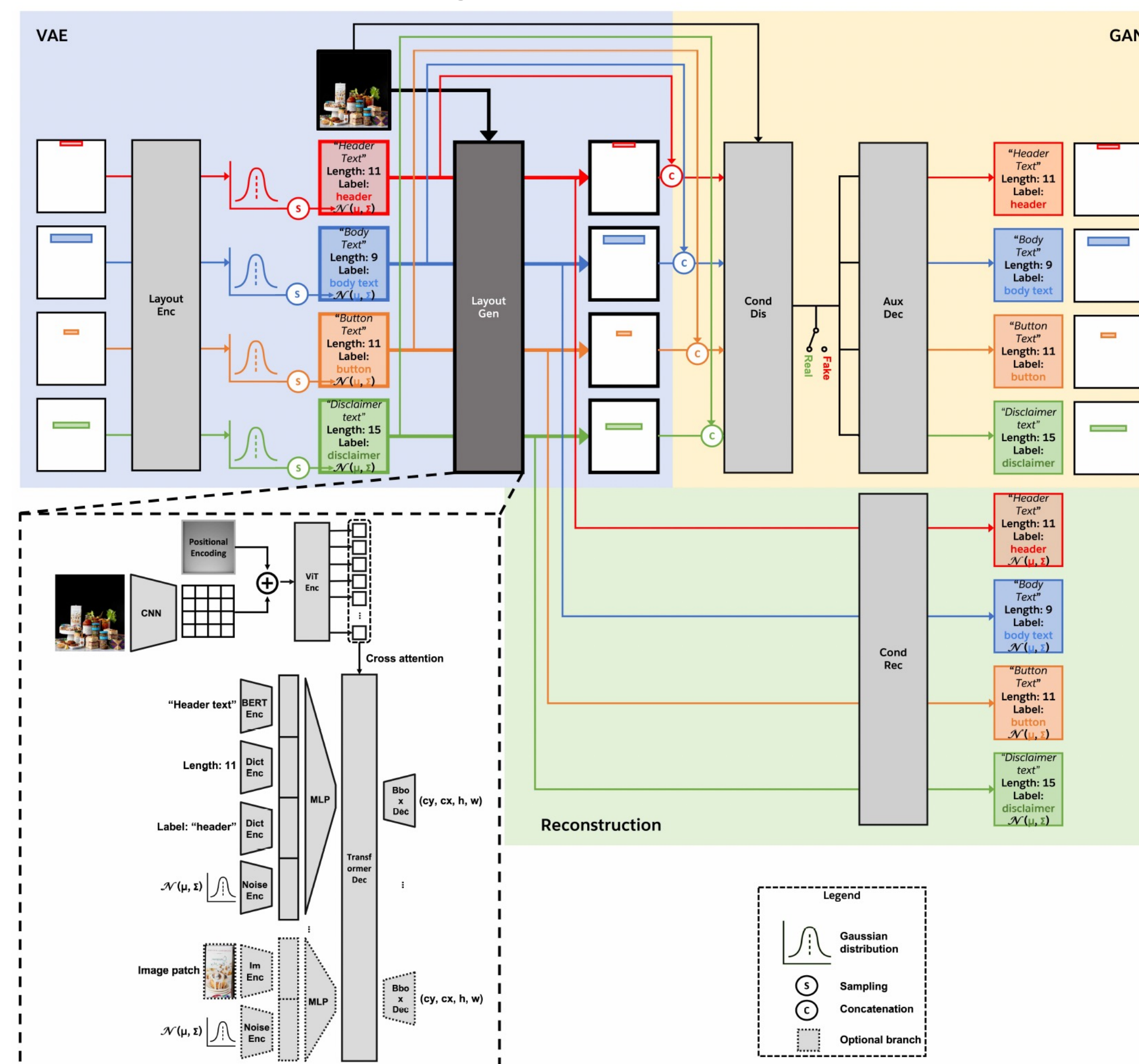
- Method:** We bridge layout generation and visual detection into one framework that solves multimodal graphic layout design.
- Dataset:** A large-scale multimodal ad banner dataset with 7,196 samples.
- State-of-the-art performance** in six evaluation metrics, which measure the realism, accuracy, and regularity of generated layouts
- Graphical system and user study:** Scales up layout generation and facilitates user studies. Users prefer our designs by significant margins.

Pipeline



Framework

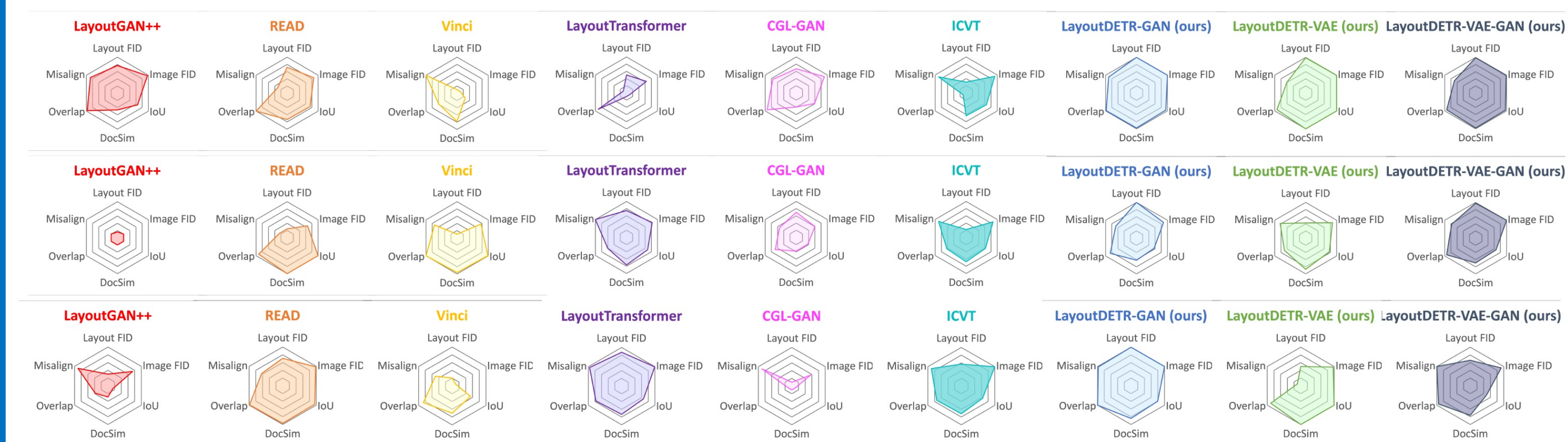
- GAN+VAE
- Unconditional/condition discriminators
- Auxiliar reconstructor
- DETR-based multimodal architecture
- BERT text encoder; ViT image encoder



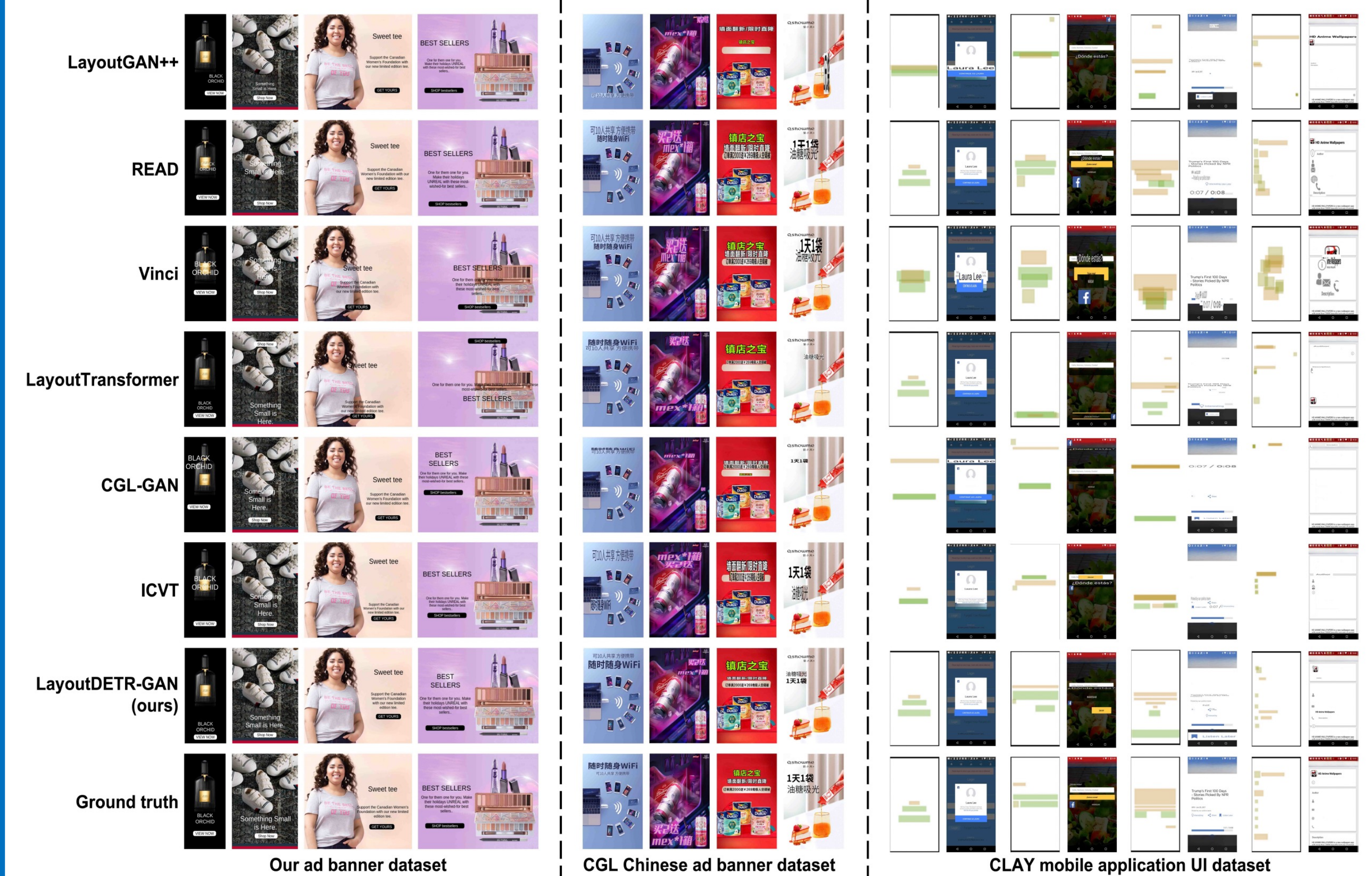
Ablation study

Method	Realism				Accuracy		Regularity		
	Layout FID ↓	Layout KID ($\times 10^{-3}$) ↓	Image FID ↓	Image KID ($\times 10^{-5}$) ↓	IoU ↑	DocSim ↑	Overlap ↓	Misalign ($\times 10^{-2}$) ↓	
Random layout on real bg	58.21 \pm 4.04	525.93 \pm 45.08	51.01 \pm 0.41	582.47 \pm 7.53	—	—	—	—	
Conditional LayoutGAN++ + Aux. Dec. (Eq. 4-7)	11.33 \pm 0.10	44.77 \pm 0.36	36.06 \pm 0.02	115.16 \pm 3.37	0.111 \pm 0.001	0.121 \pm 0.001	0.374 \pm 0.006	2.194 \pm 0.058	
+ Gen. Rec. (Eq. 11)	4.25 \pm 0.01	16.62 \pm 0.05	28.40 \pm 0.06	58.5 \pm 1.45	0.163 \pm 0.002	0.130 \pm 0.001	0.104 \pm 0.003	0.759 \pm 0.021	
+ Uncond. Dis. D^u	3.27 \pm 0.01	11.80 \pm 0.04	29.56 \pm 0.06	11.29 \pm 0.20	0.186 \pm 0.002	0.148 \pm 0.001	0.125 \pm 0.003	0.853 \pm 0.016	
+ gIoU loss (Eq. 10)	3.70 \pm 0.05	16.23 \pm 0.08	29.21 \pm 0.08	25.09 \pm 0.02	0.177 \pm 0.002	0.140 \pm 0.001	0.103 \pm 0.003	0.681 \pm 0.011	
+ Overlap & Misalign loss = LayoutDETR-GAN (ours)	3.23 \pm 0.01	11.60 \pm 0.02	28.20 \pm 0.04	10.51 \pm 0.09	0.182 \pm 0.002	0.138 \pm 0.001	0.106 \pm 0.003	0.721 \pm 0.011	
- Text length embeddings	3.24 \pm 0.01	9.25 \pm 0.05	28.65 \pm 0.03	11.42 \pm 0.35	0.191 \pm 0.002	0.144 \pm 0.001	0.117 \pm 0.003	0.807 \pm 0.012	
- Text class embeddings	25.17 \pm 0.54	171.88 \pm 5.17	29.25 \pm 0.25	139.16 \pm 4.44	0.166 \pm 0.002	0.132 \pm 0.001	0.110 \pm 0.001	0.000 \pm 0.000	

Quantitative comparisons



Quantitative comparisons



User study

Method	READ	Vinci	LayoutTransformer	CGL-GAN	ICVT	LayoutDETR-GAN (ours)
LayoutGAN++	49.8% $p=0.4$	45.6% $p=3e-3$	44.4% $p=3e-4$	53.9% $p=0.01$	47.1% $p=0.04$	55.7% $p=2e-4$
READ	—	45.1% $p=1e-3$	44.5% $p=3e-4$	53.8% $p=0.01$	53.0% $p=0.04$	54.2% $p=5e-3$
Vinci	—	—	51.7% $p=0.2$	55.8% $p=2e-4$	56.9% $p=1e-5$	62.6% $p=3e-15$
LayoutTransformer	—	—	—	57.1% $p=8e-6$	56.0% $p=2e-4$	63.5% $p=2e-17$
CGL-GAN	—	—	—	—	48.9% $p=0.2$	54.7% $p=3e-3$
ICVT	—	—	—	—	—	55.4% $p=6e-4$