

# Unifying Layout Generation with a Decoupled Diffusion Model

Mude Hui<sup>1\*</sup> Zhizheng Zhang<sup>2</sup> Xiaoyi Zhang<sup>2</sup> Wenxuan Xie<sup>2</sup> Yuwang Wang<sup>3</sup> Yan Lu<sup>2</sup>  
<sup>1</sup>Xi'an Jiaotong University <sup>2</sup>Microsoft Research Asia <sup>3</sup>Tsinghua University

Poster ID: TUE-AM-185



# Summary

- We present that various layout generation subtasks can be comprehensively unified with a single diffusion model.
- We propose the ***Layout Diffusion Generative Model (LDGM)***, which allows parallel decoupled diffusion processes for different attributes and a joint denoising process for generation with sufficient global message passing and context exploitation. It conforms to the characteristics of layouts and achieves high generation qualities.
- Extensive qualitative and quantitative experiment results demonstrate that our proposed scheme outperforms existing layout generation models in terms of the functionality and performance on different benchmark datasets.

# Background

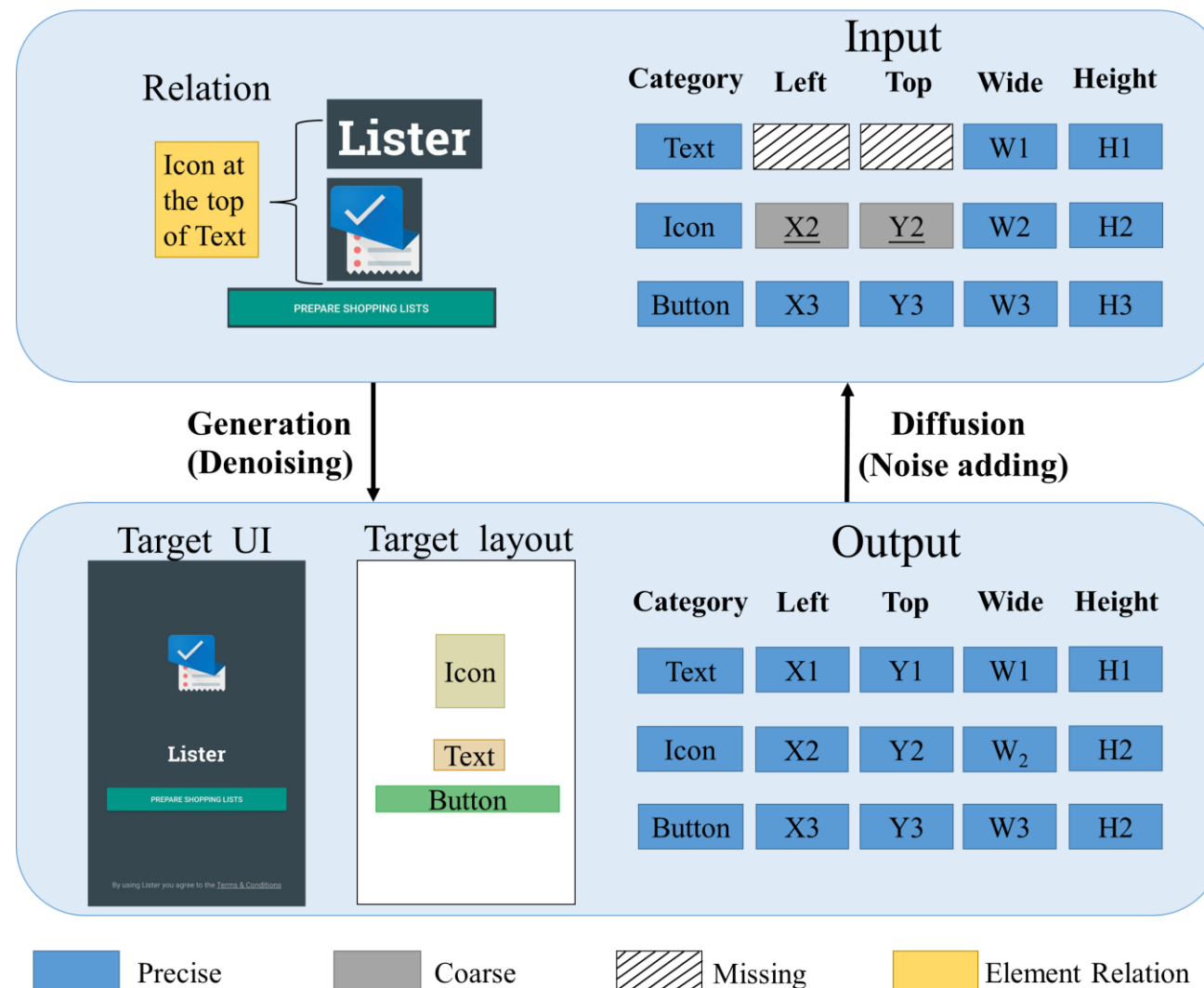
## Manual Layout Designs:

- Time-consuming
- Requiring expertise in design



## AI-based Layout Generation:

- Diverse demands (versatility)
- Aesthetics & practicality



# Generic Settings

**U-Gen:** unconditional generation

**Gen-T:** conditioned on types

**Gen-TS:** conditioned on types & sizes

**Gen-TR:** conditioned on types & relations

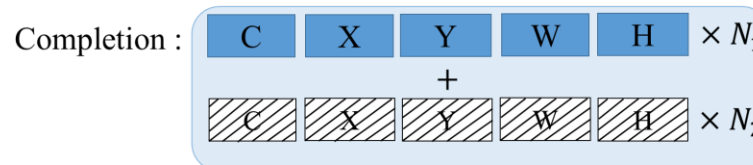
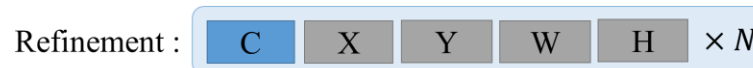
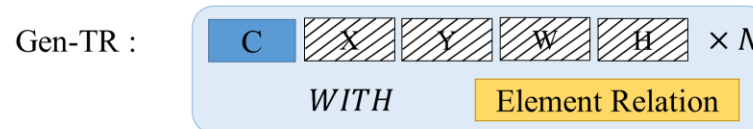
**Refinement:** update coarse attributes

**Completion:** generate missing attributes

**P:** Precise (attributes)

**M:** Missing (attributes)

**C:** Coarse (attributes)




Typical subtasks defined in previous works

 Precise

 Missing

 Missing / Precise

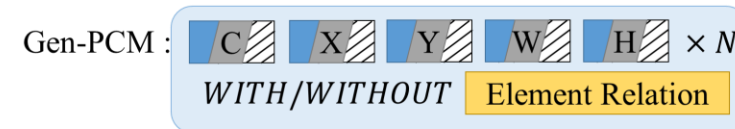
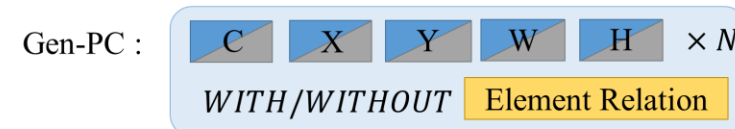
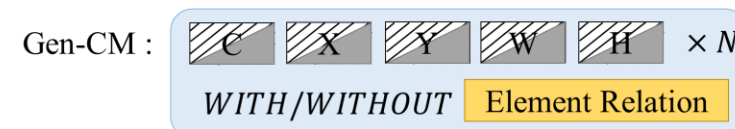
 Precise / Coarse

 Coarse

 Relation

 Missing / Coarse

 Missing / Precise / Coarse



General layout generation task settings

# Method

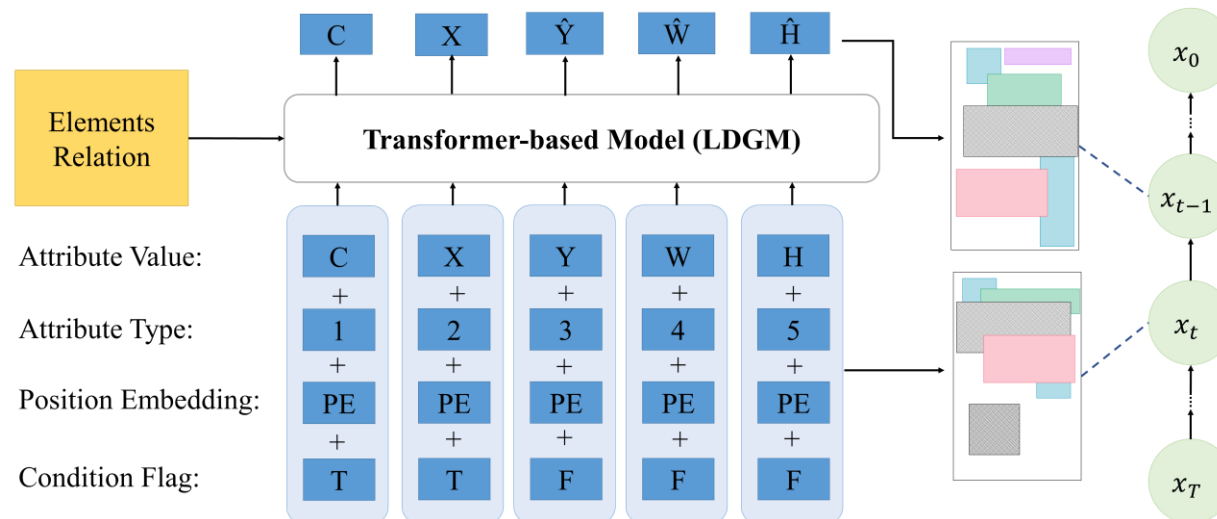
- Unification with Diffusion Modelling

The process from a completed layout to fully corruption. → A diffusion process.

- **Layout formulation:**

$$l = [c_1, x_1, y_1, w_1, h_1, c_2, x_2, y_2, \dots, h_N, \mathcal{E}]$$

- **Framework:**



# Method

## □ Decoupled Diffusion (Training)

---

**Algorithm 1** Training of the LDGM

---

**Require:** Transition matrices  $\{Q_t^c, Q_t^p, Q_t^s\}$ , initial network parameters  $\theta$ , loss weight  $\lambda$ , and learning rate  $\eta$ .

```
1: repeat
2:    $l \leftarrow$  sample a layout from the training set
3:    $timsteps = \text{zeros}(\text{len}(l))$   $\triangleright$  Record t of attributes.
4:    $\hat{l} = \text{RandSelect}(l)$   $\triangleright$  Select attributes for diffusion.
5:    $\hat{l} = [C, P, S]$   $\triangleright$  Group  $\hat{l}$  upon the semantics.
6:   for  $g$  in  $[C, P, S]$  do
7:     sample  $t \sim \text{Uniform}(\{1, \dots, T\})$ 
8:     for  $x$  in  $g$  do
9:        $timsteps[x.index] = t$ 
10:       $x = x_t \leftarrow$  sample from  $q(x_t|x_0)$   $\triangleright$  Eqn. 2
11:    end for
12:  end for
13:   $\mathcal{L}_x = \begin{cases} \lambda \mathcal{L}_{rec}, & \text{if } timsteps[x.index] = 0 \\ \mathcal{L}_0, & \text{if } timsteps[x.index] = 1 \\ \mathcal{L}_{t-1}, & \text{otherwise} \end{cases}$ 
14:   $\mathcal{L} = \sum_{x \in l} \mathcal{L}_x$ 
15:   $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$   $\triangleright$  Update network parameters.
16: until converged
```

---

Different attributes have their own semantics.



**Core idea:** “*Decouple First Diffusion Then*”.

- **For category  $c$** , we adopt noises of a uniform distribution for its diffusion.
- **For position  $(x, y)$  and size  $(w, h)$** , we adopt discretized Gaussian noises for their diffusion.

# Method

## □ A Joint Denoising Process (Generation Inference)

---

### Algorithm 2 Inference of the LDGM

---

**Require:** Initial layout  $l_T$ , condition flags, and maximum denoising steps  $T$ .

```
1:  $l_T \leftarrow$  tokenize  $l_T$  with condition flags
2:  $l_T^m \leftarrow$  GetMiss( $l_T$ )  $\triangleright$  Get missing attributes from  $l_T$ .
3:  $N_m \leftarrow$  len( $l_T^m$ )
4:  $k \leftarrow \lceil N_m/T \rceil$ 
5: for  $t = T, \dots, 1$  do
6:    $p_\theta(l_{t-1}|l_t) = LDGM(l_t)$ 
7:    $l_{t-1}, p_{t-1} \leftarrow$  sample from  $p_\theta(l_{t-1}|l_t)$ 
8:   if  $N_m > 0$  then
9:      $l_{t-1}^m, p_{t-1}^m \leftarrow$  GetMiss( $l_{t-1}, p_{t-1}$ )
10:     $l_{t-1}^m \leftarrow$  Top- $k$ Keep( $l_{t-1}^m, p_{t-1}^m$ )
11:     $N_m \leftarrow N_m - k$ 
12:   end if
13: end for
14: return  $l_0$ 
```

---

**Handling different attributes of Precise/Missing/Coarse statuses all in one.**

- **GetMiss()** refers to an operation of splitting the missing attributes from the entire attribute set.
- **Top- $k$ Keep()** refers to an operation of preserving the predicted results of missing attributes with top- $k$  high confidences and re-mark the remaining ones as absorbing status until all missing attributes are predicted.



# Comparison with SOTAs

Subtasks	Methods	Magazine				Rico				PubLayNet			
		MaxIoU $\uparrow$	FID $\downarrow$	Align. $\downarrow$	Overlap $\downarrow$	MaxIoU $\uparrow$	FID $\downarrow$	Align. $\downarrow$	Overlap $\downarrow$	MaxIoU $\uparrow$	FID $\downarrow$	Align. $\downarrow$	Overlap $\downarrow$
U-Gen	LayoutTrans. [7]	0.18	47.84	0.59	47.98	0.46	46.64	0.66	64.10	0.32	49.72	0.37	36.63
	BLT [13]	0.20	44.91	0.55	55.56	0.51	33.81	0.59	67.33	0.34	48.24	0.27	42.79
	UniLayout [9]	0.31	36.61	0.49	<b>44.50</b>	<b>0.62</b>	26.68	0.40	59.26	0.33	32.29	<b>0.22</b>	22.19
	LDGM (Ours)	<b>0.38</b>	<b>32.73</b>	<b>0.47</b>	46.43	<b>0.62</b>	<b>26.06</b>	<b>0.36</b>	<b>56.35</b>	<b>0.46</b>	<b>25.94</b>	0.25	<b>19.83</b>
Gen-T	LayoutGAN++ [12]	0.26	36.35	0.54	58.44	0.46	34.43	0.58	59.85	0.36	30.48	0.19	32.80
	BLT [13]	0.22	48.26	0.69	64.01	0.44	39.64	0.57	56.83	0.37	44.86	0.21	38.21
	UniLayout [9]	0.32	28.37	0.51	53.56	0.55	18.06	0.48	57.92	0.41	27.34	0.20	20.98
	LDGM (Ours)	<b>0.36</b>	<b>24.67</b>	<b>0.45</b>	<b>45.11</b>	<b>0.58</b>	<b>16.64</b>	<b>0.39</b>	<b>55.87</b>	<b>0.44</b>	<b>20.69</b>	<b>0.15</b>	<b>16.88</b>
Gen-TS	BLT [13]	0.33	22.72	0.59	61.94	0.51	42.88	0.46	57.74	0.40	24.32	0.16	31.06
	UniLayout [9]	0.35	19.35	0.58	56.43	0.55	20.42	0.49	58.72	0.43	27.47	<b>0.16</b>	23.82
	LDGM (Ours)	<b>0.37</b>	<b>17.65</b>	<b>0.45</b>	<b>44.25</b>	<b>0.62</b>	<b>12.59</b>	<b>0.35</b>	<b>55.92</b>	<b>0.47</b>	<b>19.02</b>	<b>0.16</b>	<b>10.09</b>
Gen-TR	CLG-LO [12]	0.27	33.88	0.59	59.43	0.38	38.89	0.54	<b>56.51</b>	0.38	31.87	0.21	34.39
	UniLayout [9]	0.36	<b>19.24</b>	0.54	49.61	0.57	26.38	0.46	66.93	<b>0.46</b>	27.73	0.17	27.35
	LDGM (Ours)	<b>0.39</b>	20.58	<b>0.48</b>	<b>47.27</b>	<b>0.61</b>	<b>16.98</b>	<b>0.39</b>	58.75	0.44	<b>19.54</b>	<b>0.16</b>	<b>21.28</b>
Refinement	RUIE [24]	0.24	44.27	0.64	54.26	0.46	36.70	0.57	64.13	0.32	41.72	0.49	35.74
	UniLayout [9]	0.33	19.78	0.49	49.02	0.56	24.41	0.42	56.04	0.44	22.34	0.11	27.23
	LDGM (Ours)	<b>0.39</b>	<b>14.95</b>	<b>0.42</b>	<b>37.22</b>	<b>0.62</b>	<b>13.19</b>	<b>0.33</b>	<b>52.17</b>	<b>0.48</b>	<b>15.28</b>	<b>0.10</b>	<b>13.05</b>
Completion	LayoutTrans. [7]	0.17	39.36	0.67	55.32	0.46	36.15	0.66	67.10	0.32	41.72	0.37	39.81
	UniLayout [9]	0.23	28.78	0.52	46.43	0.59	25.18	0.45	55.99	0.41	32.04	0.19	22.90
	LDGM (Ours)	<b>0.38</b>	<b>24.35</b>	<b>0.49</b>	<b>39.26</b>	<b>0.60</b>	<b>16.42</b>	<b>0.36</b>	<b>53.15</b>	<b>0.44</b>	<b>25.31</b>	<b>0.10</b>	<b>19.45</b>
Gen-PM		0.38	27.33	0.47	39.02	0.58	21.64	0.38	56.56	0.46	23.58	0.10	14.11
Gen-CM	LDGM (Ours)	0.37	28.74	0.51	43.25	0.57	26.15	0.38	57.74	0.44	24.94	0.11	16.26
Gen-PC		0.37	22.56	0.47	42.95	0.60	18.13	0.36	53.67	0.50	16.42	0.09	12.51
Gen-PCM		0.37	24.45	0.49	44.41	0.59	21.59	0.40	54.77	0.42	25.76	0.14	19.68
GT	-	0.41	9.89	0.43	34.27	0.66	7.05	0.26	49.86	0.64	9.38	0.008	5.18

Generation tasks supported by previous tasks.

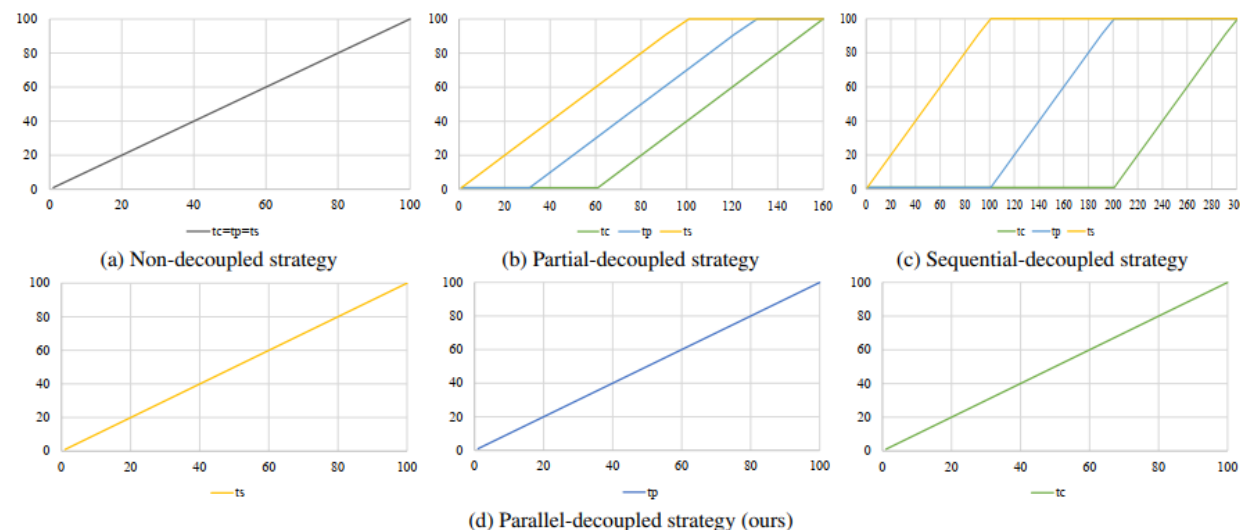
More generic generation tasks supported by ours.



# Ablation Studies

➤ Effectiveness of our proposed decoupled corruption strategy:

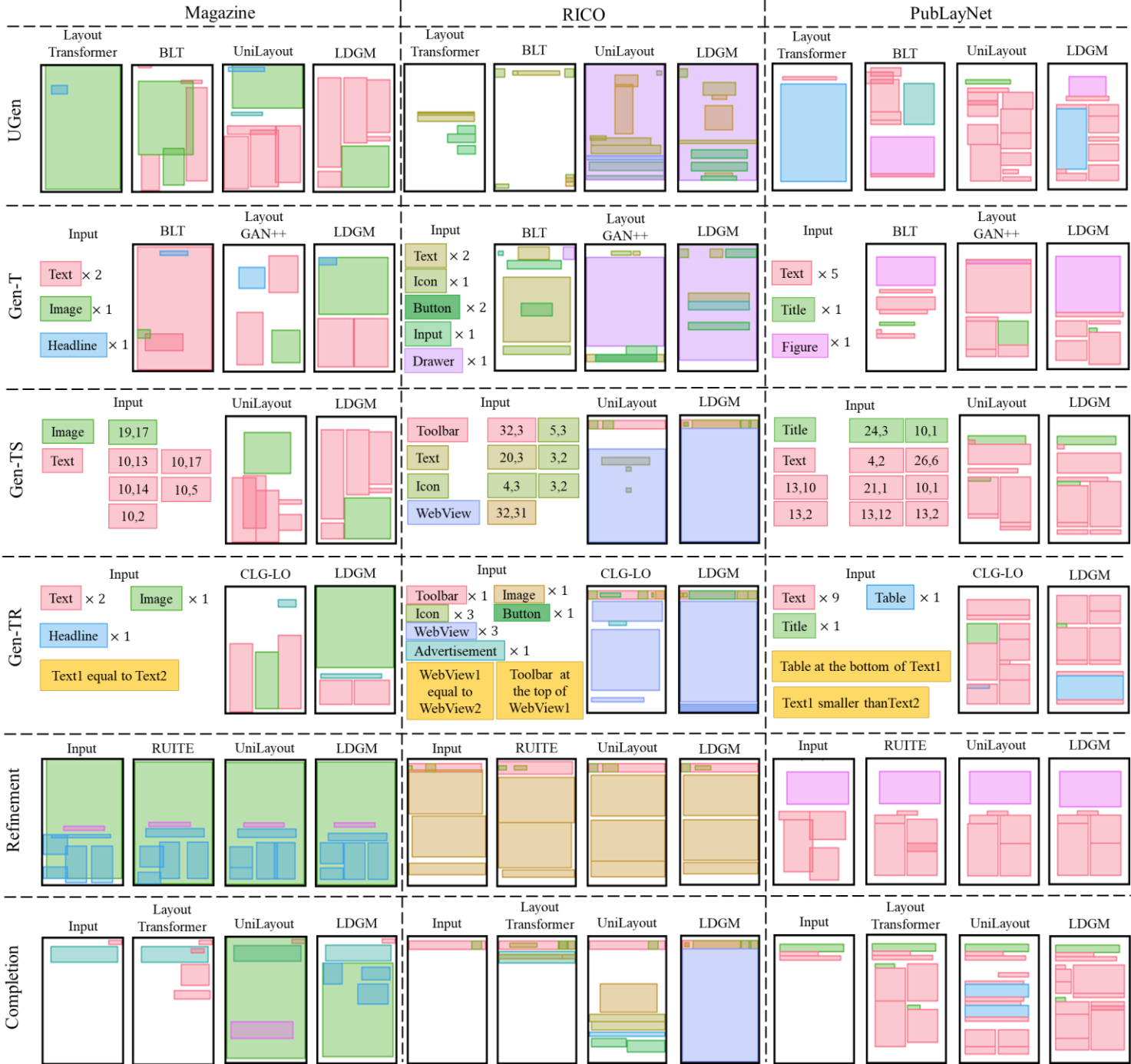
Model	MaxIoU $\uparrow$	FID $\downarrow$	Align. $\downarrow$	Overlap $\downarrow$
Non-decoupled	0.56	29.24	0.43	60.04
Partial	0.57	27.71	0.48	<b>54.24</b>
Sequential	0.56	26.69	0.43	57.17
Parallel (Ours)	<b>0.59</b>	<b>21.59</b>	<b>0.40</b>	54.77



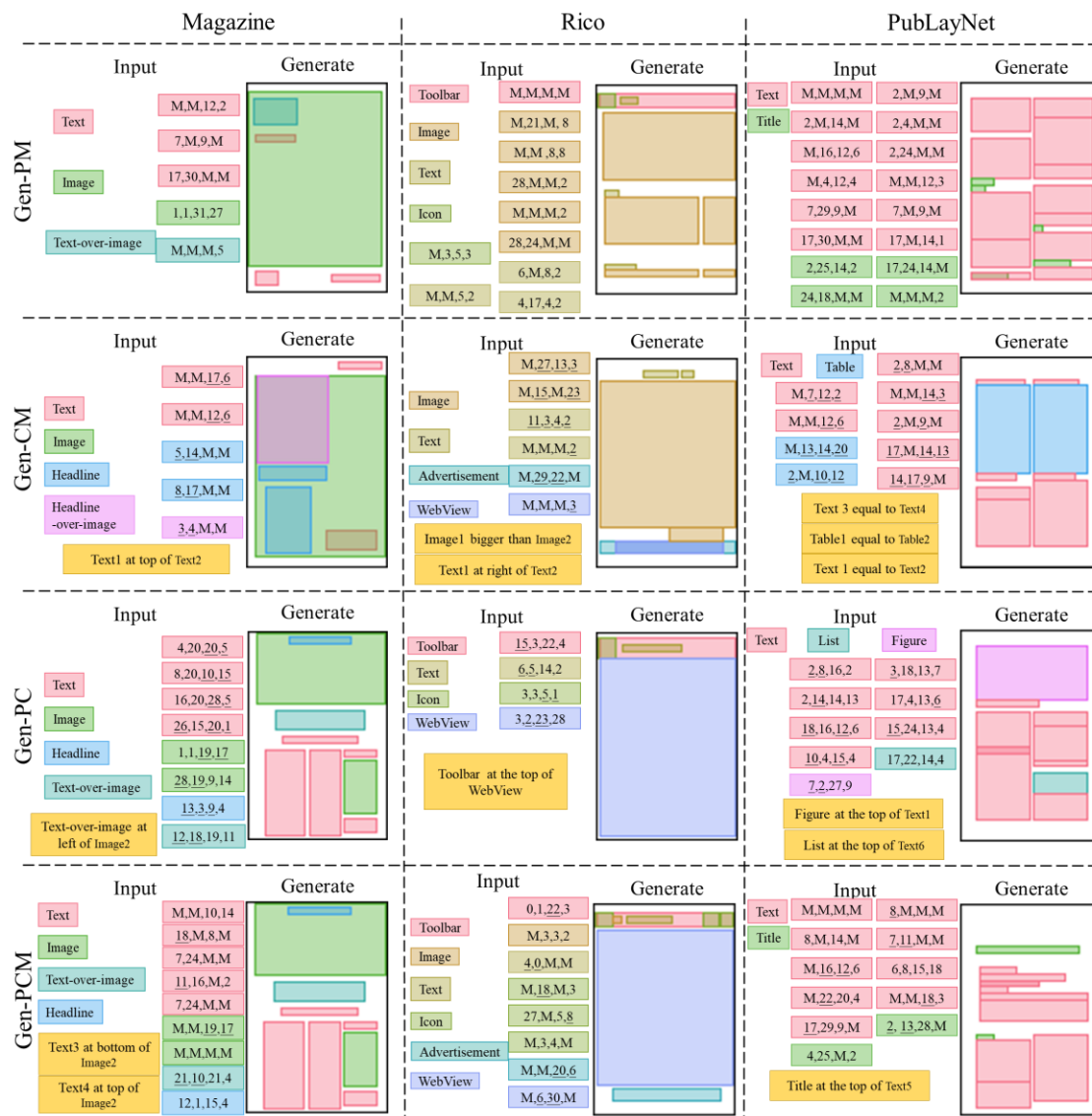
➤ Effectiveness of our proposed inference strategy:

Model	MaxIoU $\uparrow$	FID $\downarrow$	Align. $\downarrow$	Overlap $\downarrow$
AutoReg	<b>0.60</b>	23.16	0.42	56.87
Non-AutoReg	0.57	25.14	0.44	58.63
Ours	0.59	<b>21.59</b>	<b>0.40</b>	<b>54.77</b>

# Visualization



# Visualization



# Thank You!