

# LayoutFlow: Flow Matching for Layout Generation

Julian Jorge Andrade Guerreiro<sup>1</sup>, Naoto Inoue<sup>2</sup>, Kento Masui<sup>2</sup>,  
Mayu Otani<sup>2</sup>, and Hideki Nakayama<sup>1</sup>

<sup>1</sup> The University of Tokyo, Japan

{guerreiro,nakayama}@nlab.ci.i.u-tokyo.ac.jp

<sup>2</sup> CyberAgent, Japan

{inoue\_naoto,masui\_kento,otani\_mayu}@cyberagent.co.jp

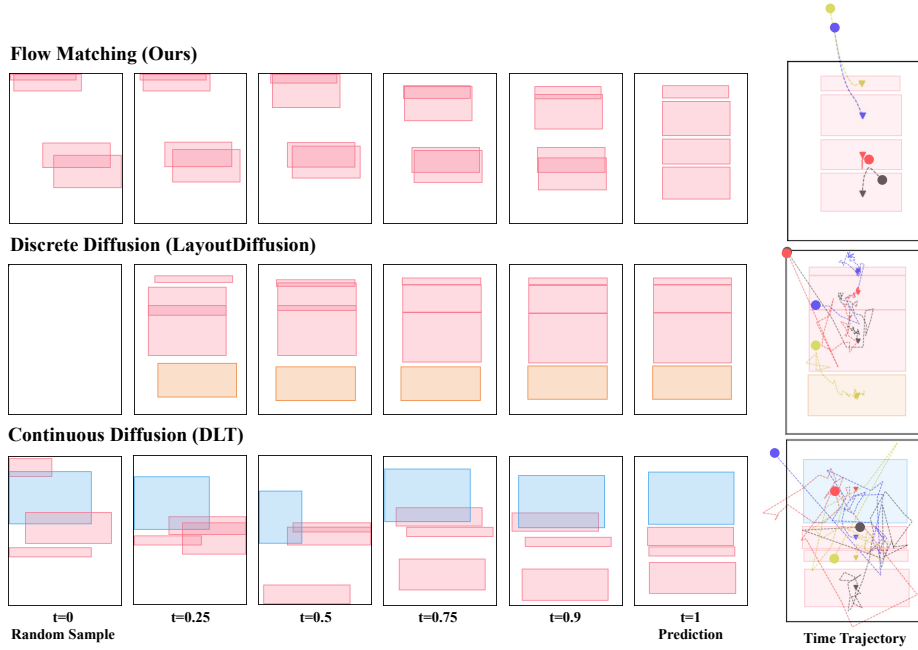
**Abstract.** Finding a suitable layout represents a crucial task for diverse applications in graphic design. Motivated by simpler and smoother sampling trajectories, we explore the use of Flow Matching as an alternative to current diffusion-based layout generation models. Specifically, we propose LayoutFlow, an efficient flow-based model capable of generating high-quality layouts. Instead of progressively denoising the elements of a noisy layout, our method learns to gradually move, or flow, the elements of an initial sample until it reaches its final prediction. In addition, we employ a conditioning scheme that allows us to handle various generation tasks with varying degrees of conditioning with a single model. Empirically, LayoutFlow performs on par with state-of-the-art models while being significantly faster.

**Keywords:** Layout Generation · Flow Matching · Generative Models

## 1 Introduction

Layout design describes the process of arranging various elements, such as images, text, or other components, on a page or screen. Finding an appropriate layout forms a crucial part of creating documents, user interfaces, graphic designs, and other compositions since the arrangement of different elements can substantially impact how the intended message or purpose is communicated. Over the years, several approaches have been explored to automate layout generation in a data-driven manner using machine learning methods. Nonetheless, current layout generation models still leave room for improvement in terms of layout quality and sampling speed.

Motivated by this observation, we propose LayoutFlow, a generative framework based on Flow Matching [1, 32, 33, 44] that can produce high-quality layouts while requiring less time than previous methods. Flow Matching has recently been introduced as a powerful generative framework and demonstrated strong performance on various tasks, including image generation [34]. Intuitively, flow-based models try to learn a *flow* that moves samples from a base distribution to a target distribution defined by the training data. In this work, we investigate



**Fig. 1: Comparison of different layout generation trajectories.** Given a randomly initialized layout at time  $t = 0$  with fixed element sizes, we visualize different states of the generation process until the final layout at  $t = 1$ . In addition, we overlay the trajectory, which can be interpreted as the movement of the elements over time, on top of the final layout. A circle marks the location of the initial sample and a triangle marks the final location. Flow Matching produces smooth and directed paths, whereas both diffusion models slowly converge to the final prediction under noisy trajectories with a long path length. As a result, flow-based models require fewer evaluation steps than diffusion, leading to faster sampling.

how to apply Flow Matching for layout generation and show that flows offer a more intuitive generation process from a geometrical interpretation compared to previous models.

Most recent layout generation models [4, 5, 8, 15, 20, 22, 28, 48] have been based on the diffusion framework [18] which has become the standard for various generative tasks [39]. While diffusion models are commonly described as models that learn to gradually remove random noise that was added to the training data, the generation process can also be interpreted in analogy to flow-based models as moving samples from a Gaussian base distribution to a data distribution. In the diffusion case, however, additional randomness is introduced along the way by adding noise. Mathematically, this corresponds to a Stochastic Differential Equation (SDE), whereas a flow is defined by a simpler Ordinary Differential Equation (ODE), which we describe in more detail in Sec. 3.

In the context of layout generation, the differences between diffusion-based and flow-based models become apparent when visualizing the trajectories created by the generation process, as illustrated in Fig. 1. The trajectories show how the sample from the base distribution is moved by the flow or backward diffusion process until it reaches the final layout prediction. While our LayoutFlow model produces smooth trajectories with short path lengths, diffusion-based models follow a noisy trajectory that constantly changes directions due to the added noise. As a result, diffusion models tend to require more sampling steps, *i.e.*, longer inference time, than flow-based models. Altogether, we argue that learning how to move the initial sample straight toward the final prediction is much more natural than trying to do so under additional noise. Overall, our contributions can be summarized as follows:

1. Motivated by its geometrical interpretation, we apply Flow Matching to the task of layout generation and propose a novel model called LayoutFlow, which can handle various tasks utilizing a conditioning mechanism.
2. On top of generating simpler trajectories, flow-based models generally also offer more choices than diffusion models, such as different prior distributions or training trajectories. We, therefore, extensively explore this additional flexibility offered by Flow Matching in the context of layout generation.
3. Empirically, we demonstrate that our proposed flow-based model significantly outperforms previous diffusion-based layout generation models of similar size and performs on par compared to a significantly larger model. In either case, our model greatly speeds up inference, requiring only a fraction of the time previous models need to generate samples.

## 2 Related Work

### 2.1 Layout Generation

Over the years, layout generation has been investigated using various methods. Early works started exploring generating layouts by minimizing an energy function based on pre-defined constraints [36, 37]. With the advent of generative machine learning techniques, however, research shifted to a more data-centric learning approach and started to use Variational Autoencoders (VAE) [2, 13, 27, 47], and Generative Adversarial Networks (GAN) [25, 29, 30]. While earlier works explored a variety of different architectures, such as Graph Convolutional Networks [27], Long-Short Term Memories (LSTM) [13, 24] or Transformers [14, 45, 47], most recent research has converged to using a Transformer-based architecture due to its flexibility and strong performance.

**Diffusion-based Models.** As diffusion models have proven to generate more diverse data and are generally easy to train, especially compared to GANs, most recent research has shifted towards exploring the diffusion process for layout generation tasks. Since layouts are represented by categorical data defining types of elements and continuous numerical values describing the element position and

size, applying standard diffusion models is not straightforward. As a result, different approaches have incorporated a diffusion loss into their training process, which can be divided into continuous [18] and discrete diffusion models [3]. For the discrete case, the continuous coordinates are quantized and interpreted as different states along with the categorical element type [15, 20, 22, 48]. In this discrete scenario, the forward diffusion process is modeled as a random walk between discrete states that, in addition to the quantized states, also include a mask state, which removes elements from the canvas. While LayoutDM [22] and LDGM [20] rely on the mask-and-replace strategy proposed in [12], LayoutDiffusion [48] introduces a new mild forward process that is closer to the continuous process while still increasingly masking elements over time.

On the other hand, continuous diffusion models have been less common as they have not attained the same performance as discrete models. In particular, continuous models have struggled to produce layouts with well-aligned elements. The first continuous layout generation model based on diffusion was proposed by Chai *et al.* [4], but is restricted to generating continuous element coordinates given categorical type data as a condition. Cheng *et al.* [8] introduced a latent diffusion model for layout design focusing on user constraints by embedding the layout into continuous representation, on which the diffusion process is performed. The first continuous diffusion-based layout generation model acting on the layout representation and able to handle unconditional generation was proposed by Chen *et al.* [5]. Their work addresses the alignment issue observed for continuous diffusion by introducing a specific regularization loss. Lastly, Levi *et al.* [28] proposed DLT, which applies a continuous diffusion process on the element coordinates and a discrete process on the element types.

**Large Language Models.** Motivated by the strong generalization capabilities of Large Language Models (LLM), the most recent works [11, 31, 43] have investigated the layout generation abilities of LLMs. While LLMs provide some interesting applications, such as zero-shot synthesis, they typically consist of billions of parameters, resulting in long inference times. Moreover, current models are limited to conditional layout generation. In contrast, our proposed approach aims to offer a lightweight model for unconditional and conditional layout generation.

## 2.2 Flow Matching and Diffusion Models

Flow-based models are based on learning a mapping between samples from a simpler distribution to a data distribution. Over the years, there have been various methods, such as Normalizing Flows or Continuous Normalizing Flows, which have used neural networks to estimate such a mapping, which can subsequently be used to sample from the learned distribution [26]. Until recently, the biggest issue with flow-based models has been the need to backpropagate through an ODE during training. However, recent models have proposed training algorithms that only require explicitly solving an ODE during inference [1, 32, 33, 35, 44] demonstrating impressive results and outperforming diffusion models on the image generation task [34]. Our proposed model builds on this improvement by

utilizing Flow Matching for layout generation while offering an improved geometrical interpretation.

Diffusion models, which were first proposed by Sohl-Dickstein *et al.* [40] and later popularized by Ho *et al.* [18], can be formulated similarly to flow-based models. In fact, the term Flow Matching was inspired by the similar Score Matching loss [21] used by diffusion models. However, flows are characterized by an ODE, whereas the diffusion process can be formulated as an SDE with an additional stochastic component described by Brownian Motion [42]. As a result, diffusion models are restricted to a Gaussian as a base distribution, while flow-based models allow for more flexibility, which we also explore in our experiments. Notably, Song *et al.* [41] proposed DDIM, an alternative sampling method that makes it possible to sample using an ODE formulation on a model trained with an SDE forward process, effectively leading to smoother trajectories. In the context of layout generation, however, models have been using the SDE formulation following DDPM [18]. Our argument for introducing Flow Matching instead of just using DDIM for smoother trajectories lies in its increased flexibility and simplicity. Moreover, if the goal is to sample using an ODE, it seems more natural to also use an ODE during training instead of the more complex SDE used in diffusion.

### 3 Preliminary: Flow Matching

In this section, we first present the mathematical definition of a flow and then introduce how to train a neural network to estimate a flow via Flow Matching. The term *flow* in Flow Matching refers to a mapping between samples of two distributions. In general, flow-based models aim to estimate the flow between a known source probability distribution  $p_0(x)$ , *e.g.*, a Gaussian distribution, and the typically more complex data distribution  $p_1(x)$ . Given a data point  $x \in \mathbb{R}^d$ , the flow  $\phi : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  can be defined by the following Ordinary Differential Equation (ODE)

$$\frac{d}{dt}\phi_t(x) = v_t(\phi_t(x)), \quad \phi_0(x) = x_0, \quad (1)$$

where the time-dependent vector field  $v_t : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  is said to construct the flow [32] with the boundary condition  $\phi_0(x) = x_0$ . In other words, the flow  $\phi_t(x)$  describes how an initial sample  $x_0$  is transported over time and is given as the solution to the ODE in Eq. (1). As the flow transports samples with respect to time, the associated probability distribution is transformed according to the change of variables formula and creates a probability density path

$$p_t = [\phi_t]_* p_0 = p_0(\phi_t^{-1}(x)) \det \left[ \frac{d\phi_t^{-1}}{dt}(x) \right] \quad (2)$$

with  $\phi_t^{-1}$  denoting the inverse of the flow and  $*$  representing the push-forward operator. If the flow constructed by the vector field  $v_t$  satisfies Eq. (2), it is

said that  $v_t$  generates the probability path  $p_t$ . This is the case if  $v_t$  follows the continuity equation, which holds for our considerations in this paper.

Instead of estimating the flow directly, Flow Matching trains a neural network  $u_\theta(t, x)$  with weights  $\theta$  to match the vector field  $v_t$  and then solves the ODE for a given boundary condition  $x_0$  to obtain the flow. The Flow Matching loss function is defined as the expectation over uniformly sampled timesteps and samples along the probability density path as follows

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1), x \sim p_t(x)} \|u_\theta(t, x) - v_t(x)\|^2. \quad (3)$$

However, this loss alone does not allow us to train the model as we usually do not have access to  $p_t$  nor  $v_t$ . As an extension, Lipman *et al.* [32] have proposed a Conditional Flow Matching objective and is given as

$$\mathcal{L}_{CFM}(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1), z \sim q(z), x \sim p_t(x|z)} \|u_\theta(t, x) - v_t(x|z)\|^2, \quad (4)$$

with  $v_t(x|z)$  denoting a conditional vector field generating the conditional probability path  $p_t(x|z)$  and  $q(z)$  being a distribution over a latent conditioning variable  $z$ . By introducing the conditional formulation of the Flow Matching objective, training a neural network becomes possible by choosing and training with a conditional vector field and its associated probability path. For an overview of the various choices and more details on Conditional Flow Matching, we recommend [32, 44].

## 4 Conditional Flow Matching for Layout Generation

### 4.1 Data Representation

A layout typically consists of various elements, such as text, tables, or images placed on a canvas. Mathematically, we can describe a layout  $\mathcal{A}$  as a set of  $N$  elements  $\mathbf{e}^k$  with  $k \in [1, N]$  in the following way

$$\mathcal{A} = \{\mathbf{e}^k = (\mathbf{g}^k, \mathbf{a}^k) | k \in [1, N]\}, \quad (5)$$

where we split each element into its geometrical information  $\mathbf{g}^k = (c_x^k, c_y^k, w^k, h^k)$  described by the bounding box center  $(c_x^k, c_y^k)$ , width  $w^k$  and height  $h^k$ , and its category information  $\mathbf{a}^k$  describing the element type, such as text or image. In our considerations, the elements on the canvas do not possess an inherent order, which is why the canvas is considered a set of elements. Layout generation can be divided into various tasks depending on the conditioning assumptions. For example, in the unconditional generation setting, it is assumed that no prior knowledge is available, while for the refinement task, a complete but coarse layout is given. Since Flow-Matching requires the input data to be continuous, we convert the categorical information  $\mathbf{a}^k$  into a continuous embedding, denoted as  $\tilde{\mathbf{a}}^k$ , using  $B$  Analog Bits [7]. To learn the flow, we define a data sample  $\mathbf{x} \in \mathbb{R}^{(4+B) \cdot N_{max}}$  as the concatenation of all the elements  $\mathbf{g}^k$  and  $\tilde{\mathbf{a}}^k$  in a single layout. If the number of elements is smaller than  $N_{max}$ , the maximum number of elements contained in a single layout across the entire dataset, we pad the remaining dimensions.

**Algorithm 1: Training**

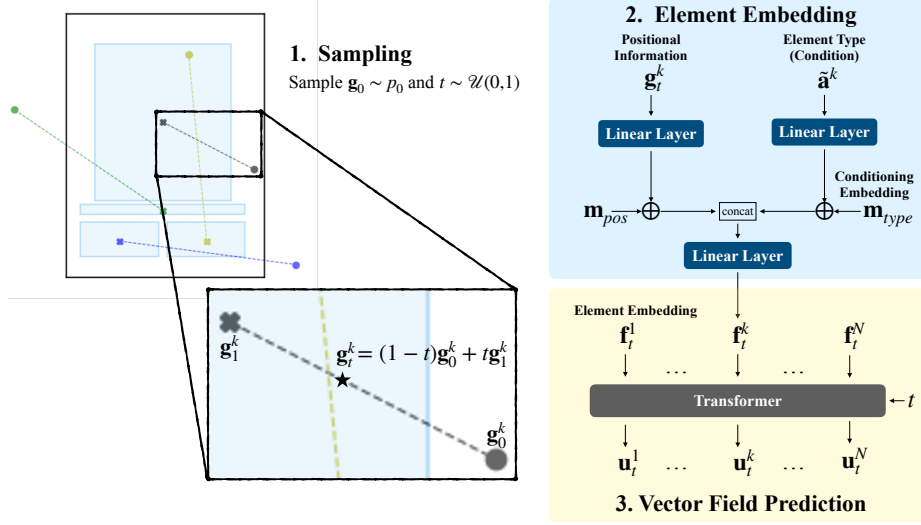

---

```

 $\mathbf{x}_0 \sim p_0(\mathbf{x}), \mathbf{x}_1 \sim p_{data}(\mathbf{x}), t \sim \mathcal{U}(0,1)$       // Sampling from distributions
 $\mathbf{x}_t \leftarrow (1-t)\mathbf{x}_0 + t\mathbf{x}_1$                           // Linear interpolation
 $\mathbf{v}_t \leftarrow \mathbf{x}_1 - \mathbf{x}_0$                             // Calculate the constant vector field
 $\mathbf{u}_t \leftarrow f_\theta(\mathbf{x}_t, t)$                         // Get vector field prediction from network
 $\mathcal{L}(\theta) = \|\mathbf{u}_t - \mathbf{v}_t\|^2 + \lambda \|\mathbf{u}_\theta^g(t, \mathbf{x}_t) - (\mathbf{g}_1 - \mathbf{g}_0)\|_1$   // Compute loss

```

---



**Fig.2: Overview of the training procedure of LayoutFlow for the type-conditioned scenario.** First, we sample an initial layout from a base distribution and a time  $t$ . Then, an intermediate sample  $\mathbf{g}_t$  is calculated by linearly interpolating between the initial sample and the ground truth layout. Each intermediate element is embedded jointly with the given element condition  $\tilde{\mathbf{a}}^k$ . Lastly, the Transformer architecture takes all the element embeddings to predict a vector field.

## 4.2 Training

To generate layouts, we want to find a way to sample from a data distribution  $p_{data}(\mathbf{x})$ , whereas we only have access to a limited amount of samples  $\mathbf{x}$  through our training data. More specifically, our goal is to train a neural network that allows us to sample from the underlying data distribution and guide the sampling based on certain conditional constraints. To that end, we propose LayoutFlow, a layout generation model based on Flow Matching.

An overview of the training procedure for the type-conditioned scenario is illustrated in Fig. 2 and summarized in Algorithm 1 for the unconditional case. In the first step, we randomly sample a layout  $\mathbf{x}_0$  from a prior distribution  $p_0$ , *e.g.*, a Gaussian distribution. This sample will typically not be well-aligned, and some elements might not be located on the canvas. In addition, we also randomly select a time  $t$ , which determines the intermediate sample  $\mathbf{x}_t$  on the flow trajec-

tory used during training. For LayoutFlow, we train on a simple linear trajectory following [32, 44]. Therefore, the intermediate sample is simply obtained by linearly interpolating between the ground truth layout  $\mathbf{x}_1$  and the initial sample  $\mathbf{x}_0$ . The same trajectory is not used again during training since  $\mathbf{x}_0$  and  $t$  change with every iteration.

The network  $f_\theta$  is trained to output the conditional vector field, corresponding to the derivative of the trajectory as defined in Eq. (1). This derivative is simply a constant for the linear trajectory case, specifically the difference between  $\mathbf{x}_0$  and  $\mathbf{x}_1$ . Intuitively, the network learns to output a direction pointing toward a data sample, which can be used during sampling to move the initial data samples toward a better prediction. While we train our network using linear trajectories, this does not mean that the model will also produce straight trajectories during sampling. As mentioned in Eq. (4), the linear trajectory, which is derived by the conditional vector field in Eq. (4), acts as a proxy to learn the actual vector field in Eq. (3) defined by Flow Matching.

In our experiments, we found that using only the Mean Squared Error (MSE) as a loss function tends to lead to poorly aligned layouts, which results in suboptimal perceptual quality. This observation is similar to the issue described in [5]. We hypothesize that this can be attributed to the fact that the MSE penalizes all mistakes within the same distance equally, regardless of direction, whereas the alignment between elements is usually only measured along the horizontal and vertical direction. As a result, the minimization objective does not fully align with the perceptual quality. To tackle this issue, we add an  $L_1$  regularization on the geometrical part of the output  $u_\theta^{\mathbf{g}}$  as the  $L_1$  loss encourages sparsity of the error vector, which can also be interpreted as minimizing the error along the axis dimensions. Therefore, our final training loss can be written as

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{CFM}}(\theta) + \lambda \mathcal{L}_1(\theta) = \mathcal{L}_{\text{CFM}}(\theta) + \lambda \|u_\theta^{\mathbf{g}}(t, \mathbf{x}_t) - (\mathbf{g}_1 - \mathbf{g}_0)\|_1, \quad (6)$$

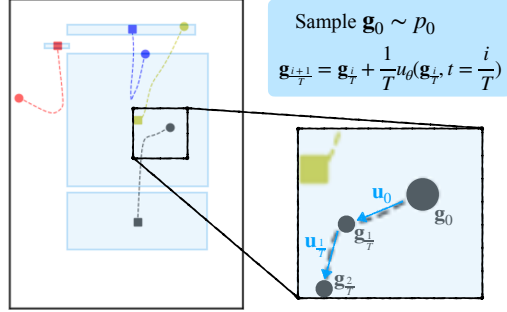
where  $\lambda$  is a hyperparameter.

### 4.3 Sampling

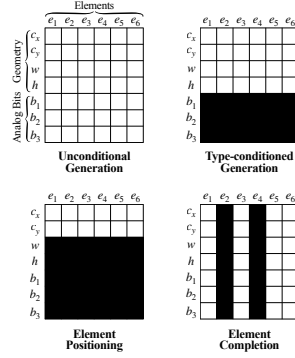
After training our neural network to predict a vector field, we can generate new layouts by sampling from our initial distribution and solving the ODE describing the flow as defined in Eq. (1). In practice, this can be done using any numerical ODE solver. Figure 3 illustrates the sampling process using the Euler method [10] to solve the ODE. For a fixed number of evaluation steps  $T$ , the initial sample  $x_0$  is moved along the direction predicted by the network in an autoregressive manner. Therefore, the backward process can be denoted as

$$\mathbf{x}_{\frac{i+1}{T}} = \mathbf{x}_{\frac{i}{T}} + \frac{1}{T} u_\theta \left( \frac{i}{T}, \mathbf{x}_{\frac{i}{T}} \right), \quad (7)$$

with  $i \in \mathbb{N}$  and  $i \in [0, T-1]$  indicating how far along the sample has been moved.



**Fig. 3: Overview of the type-conditioned inference process.** The initial sample is autoregressively moved in the predicted direction.



**Fig. 4: Condition masks.** Black indicates parts given as conditions.

#### 4.4 Network Architecture

Here, we describe the actual implementation of the network  $f_\theta$  to efficiently handle the layout. We first embed the elements of the sampled intermediate layout  $\mathbf{x}_t$  and then use a Transformer to predict the conditional vector field as depicted in Fig. 2. First, we separately embed the geometrical information  $\mathbf{g}_t^k$  and the element type  $\tilde{\mathbf{a}}_t^k$  using a linear layer. A conditioning mask is embedded and added to the geometrical and type embedding to inform the network about which parts of the input act as a condition. The example in Fig. 2 illustrates this for the type-conditioned case, where the network is informed through the masks  $\mathbf{m}_{type}$  and  $\mathbf{m}_{pos}$  that  $\tilde{\mathbf{a}}_t^k$  is taken from  $\mathbf{x}_1$  while only the geometrical information  $\mathbf{g}_t^k$  was taken from the training trajectory. In a subsequent step, we concatenate both embeddings and pass them through another linear layer to obtain a high-dimensional embedding  $\mathbf{f}_t^k$  representing an element of the intermediate layout  $\mathbf{x}_t$ . Next, all element embeddings  $\mathbf{f}_t^k$  are processed by a transformer encoder network together with the temporal information  $t$  incorporated in each layer. Finally, the output of the Transformer is projected back to the same dimension as the input layout using a linear layer representing the vector field.

#### 4.5 Conditioning Mechanism

As layout generation involves various tasks with different assumptions about available knowledge, we employ a conditioning mechanism to perform these tasks with just a single model. Similarly to [5, 20, 28], we mark the conditioning input by including the information in the element embedding while distinguishing between the four different scenarios depicted in Fig. 4. During training, we simply randomly sample each mask with the same probability. Alternatively, it is also possible to only train an unconditional model and impose the conditions during sampling through a guidance strategy, similar to [22, 48]. However, we empiri-

cally found that including conditioning during training works better. For a more detailed analysis, we refer to the supplementary material.

## 5 Experiments

To compare LayoutFlow with other methods, we first perform a series of experiments on various layout generation tasks. In addition, we conduct an ablation study on our proposed method.

### 5.1 Experimental Settings

**Datasets.** We evaluate the performance of our network on the RICO [9] and PubLayNet [49] datasets, following previous methods. RICO consists of over 66k User Interface (UI) layouts with 25 element categories, while PubLayNet comprises over 360k document layouts annotated with 5 different element types. We train our network using the dataset split described in [23,48], which discards layouts containing more than 20 elements.

**Implementation Details.** Similar to LayoutDM and DLT, we employ a 4-layer Transformer network with 8 attention heads and a model dimension of 512 for a fair comparison. The model is trained using the AdamW optimizer with a learning rate of 0.0005. In addition to the MSE Flow Matching loss, we regularize the training by calculating the  $L_1$  loss between the geometrical information as described in Eq. (6) with  $\lambda = 0.2$ . For inference, we obtain the flow using a NeuralODE solver [6] for 100 steps with the Euler method.

**Evaluation Metrics.** Following previous work [4, 5, 28, 48], we evaluate our method as well as comparable approaches using *Frechet Inception Distance (FID)* [16], *Maximum Intersection over Union (mIoU)* [25], *Alignment* [27], and *Overlap* [29]. *FID* and *mIoU* inherently evaluate both fidelity and diversity [17] of the generated results. For the *FID* calculation, we use the same network with identical weights as LayoutDiffusion [48] to measure the similarity between the generated layouts and the original dataset based on the feature space. The *mIoU* metric compares sets of layouts and measures their similarity by optimally matching generated and real layouts that maximize the average *IoU*. On the other hand, *Alignment* and *Overlap* capture the fidelity property. Nonetheless, both metrics should be judged with respect to a reference dataset, where the best outputs lead to similar values.

**Comparison with Existing Approaches.** To validate the performance of our proposed method, we compare it to state-of-the-art techniques. Since there is no standard way of splitting the datasets, we need to retrain the models that do not follow the dataset split used in LayoutDiffusion [48] for a fair comparison. This fact limits us to only compare to methods that have made their code publicly available. Therefore, we compare our method against the two discrete diffusion models LayoutDM [4] and LayoutDiffusion [48] as well as the discrete-continuous diffusion model DLT [28]. Where applicable, we additionally compare

**Table 1: Quantitative results for various layout generation tasks on the RICO and PubLayNet datasets.** The two best results are highlighted in **bold** and underlined. The  $\rightarrow$  symbol indicates best results are the ones closest to the validation data. Models marked with \* have been retrained.

Task	Model	RICO				PubLayNet			
		FID↓	Ali→	Ove→	mIoU↑	FID↓	Ali→	Ove→	mIoU↑
Un-Gen	LayoutTransformer	24.32	0.037	0.542	0.587	30.05	0.067	0.005	0.359
	LayoutFormer++	20.20	<u>0.051</u>	0.546	<b>0.634</b>	47.08	0.228	<u>0.001</u>	0.401
	LayoutDM*	4.43	0.143	0.584	0.582	36.85	0.180	0.132	0.382
	DLT*	13.02	0.271	0.571	0.566	12.70	0.117	0.036	<b>0.431</b>
	LayoutDiffusion	<u>2.49</u>	<b>0.069</b>	<u>0.502</u>	<u>0.620</u>	<b>8.63</b>	<u>0.065</u>	<b>0.003</b>	0.417
	LayoutFlow (ours)	<b>2.37</b>	0.150	<b>0.498</b>	0.570	<u>8.87</u>	<b>0.057</b>	0.009	<u>0.424</u>
Gen-Type	NDN-none	13.76	0.560	0.550	<u>0.350</u>	35.67	0.350	0.170	0.310
	LayoutFormer++	2.48	<u>0.124</u>	0.537	<b>0.377</b>	10.15	<b>0.025</b>	<u>0.009</u>	0.333
	LayoutDM*	2.39	0.222	0.598	0.341	39.12	0.267	0.139	0.348
	DLT*	6.64	0.303	0.616	0.326	7.09	0.097	0.040	<u>0.349</u>
	LayoutDiffusion	<u>1.56</u>	<b>0.124</b>	<b>0.491</b>	0.345	<u>3.73</u>	<u>0.029</u>	<b>0.005</b>	0.343
	LayoutFlow (ours)	<b>1.48</b>	0.176	<u>0.517</u>	0.322	<b>3.66</b>	0.037	0.011	<b>0.350</b>
Gen-TypeSize	LayoutDM*	<u>1.76</u>	<b>0.175</b>	<u>0.606</u>	0.424	29.91	0.246	0.160	<u>0.436</u>
	DLT*	6.27	0.332	0.609	0.424	<u>5.35</u>	<u>0.130</u>	<u>0.053</u>	0.426
	LayoutFlow (ours)	<b>1.03</b>	<u>0.283</u>	<b>0.523</b>	<b>0.470</b>	<b>1.26</b>	<b>0.041</b>	<b>0.031</b>	<b>0.454</b>
<b>Validation Data</b>		2.10	0.093	0.466	0.658	8.10	0.022	0.003	0.434

LayoutFlow to the non-diffusion-based models LayoutTransformer [14], LayoutFormer++ [23], NDN-none [27] and RUIE [38].

**Tasks.** We test LayoutFlow and existing approaches on various common unconditional and conditional layout generation tasks. *Un-Gen* describes the layout generation task without any constraints. For the conditional scenario, we consider *Gen-Type* as assuming the element types as given, while *Gen-TypeSize* assumes element types and sizes are known. Lastly, we consider the *Completion* task, which takes in a partial layout that is missing up to 20% of its elements, and the *Refinement* task following [48] with a standard deviation of 0.01.

## 5.2 Evaluation

**Quantitative Evaluation.** In Tab. 10, we report the results of our experiments comparing LayoutFlow with existing approaches for different layout generation tasks. In terms of *FID*, our model outperforms state-of-the-art methods across all three tasks, except for a close second place on the PubLayNet *Un-Gen* task, proving its strong capabilities to model the underlying sample distribution. In particular, the *FID* score of 2.37 produced by LayoutFlow on the unconditional generation for RICO almost matches the score obtained when comparing the validation with the test dataset. Note that for the conditional tasks, the *FID* score can become even lower than the validation dataset since the generated

**Table 2: Quantitative results for the completion and refinement tasks on the RICO and PubLayNet datasets.** The two best results are highlighted in **bold** and underlined. Models marked with \* have been retrained.

Task	Model	RICO				PubLayNet			
		FID↓	Ali→	Ove→	mIoU↑	FID↓	Ali→	Ove→	mIoU↑
Completion	LayoutDM*	6.80	<b>0.054</b>	0.630	0.678	25.02	0.169	0.107	0.678
	LayoutFlow (ours)	<b>1.51</b>	0.150	<b>0.474</b>	<b>0.741</b>	<b>1.10</b>	<b>0.054</b>	<b>0.127</b>	<b>0.746</b>
Refinement	RUITE	7.93	0.177	0.492	0.658	7.89	0.073	0.038	0.637
	LayoutFormer++	3.67	<u>0.141</u>	0.503	0.656	2.94	0.042	0.013	0.642
	LayoutDM*	2.91	0.143	0.575	0.437	48.61	0.286	0.173	0.372
	LayoutDiffusion	<b>0.55</b>	<b>0.102</b>	<b>0.469</b>	<b>0.719</b>	<u>2.05</u>	<u>0.035</u>	<u>0.008</u>	<u>0.660</u>
	LayoutFlow (ours)	<u>0.77</u>	0.152	<u>0.455</u>	<u>0.700</u>	<b>0.18</b>	<b>0.020</b>	<b>0.005</b>	<b>0.723</b>
	<b>Validation Data</b>	2.10	0.093	0.466	0.658	8.10	0.022	0.003	0.434

images share more similarities through conditioning. Regarding the geometrical metrics *Alignment*, *Overlap*, and *mIoU*, our method consistently outperforms LayoutDM and DLT while usually producing values close to LayoutDiffusion, a significantly larger model with 85M parameters compared to around 15M used by LayoutFlow. A comparison between the quality and the sampling speed is illustrated in Fig. 8, clearly showing that our method closes the gap between inference speed and performance. While non-diffusion models outperform diffusion models on some geometrical metrics, the drastically higher *FID* score implies that the generated results lack diversity.

On the completion task, shown in Tab. 2, LayoutFlow is able to clearly outperform LayoutDM, the only diffusion model that can handle that task, except for falling behind on *Alignment* for RICO. On the refinement task, LayoutFlow significantly improves the noisy layouts on PubLayNet, while being the second-best model for RICO, underlining its broad capabilities.

**Qualitative Evaluation.** We provide some qualitative examples in Fig. 6 for the unconditional Fig. 7 as well as for the conditional generation task on the RICO dataset. More samples for the other tasks and PubLayNet can be found in the supplementary material. Overall, LayoutFlow shows a strong performance across all tasks, producing visually pleasing results that resemble real layouts.

**Limitations.** We discuss the limitations of our proposed methods in the supplementary material.

### 5.3 Ablation Study

In order to justify our design choices, we perform an ablation study on unconditional layout generation using the RICO dataset.

**Diffusion vs. Flow.** We show that employing Flow Matching as the generative model is essential by also training the model used for LayoutFlow and only substituting Flow Matching with diffusion. We test two methods for sampling

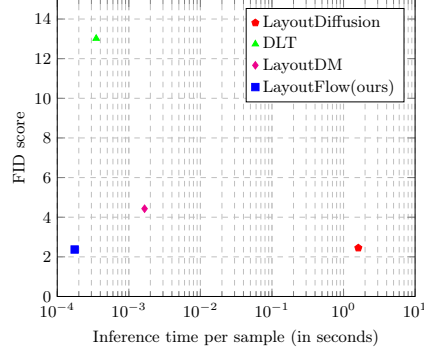


Fig. 5: Quality-Speed Comparison

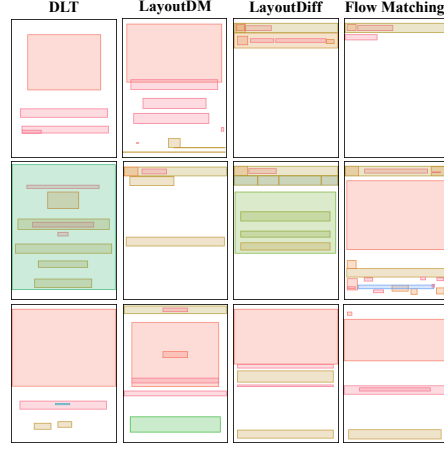


Fig. 6: Examples for Un-Gen

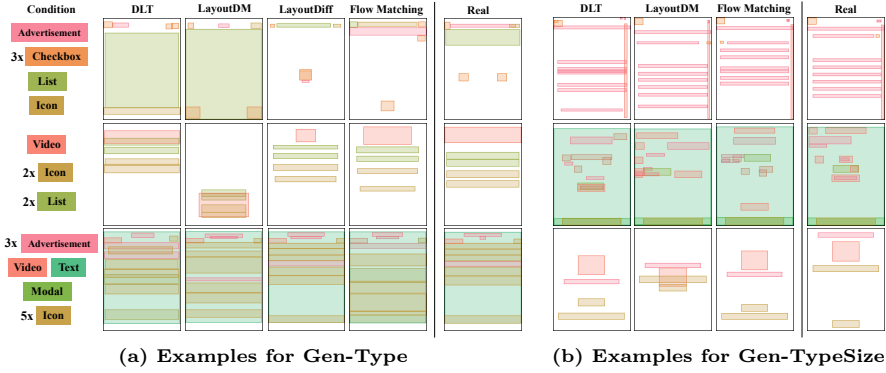


Fig. 7: Conditional Generation Examples

from the diffusion-based model, i.e., DDPM and DDIM. As shown in Tab. 3, neither approach is able to reach the same performance as LayoutFlow, and both seem to be particularly lacking in terms of *Overlap*. Interestingly, DDIM outperforms DDPM by a large margin, which might be related to DDIM sampling based on an ODE instead of an SDE.

**Regularization Loss.** We validate the efficacy of our regularization loss proposed in Eq. (6) and the choice of  $\lambda$  in Tab. 4. It can be clearly observed that a stronger regularization using the  $L1$ -loss provides better alignment but increases the *FID*. Therefore, we considered the trade-off across all metrics and chose  $\lambda = 0.2$  for LayoutFlow.

**Initial Distribution.** Since flow-based models allow for arbitrary initial distributions, we explore two alternatives. First, we try a uniform distribution to ensure all initial samples are placed on the canvas. In addition, as Analog

**Table 3: Ablation study on the choice of training approaches using the same architecture**

Approach	FID↓	Ali→	Ove→	mIoU↑
Diff.(DDPM)	34.89	<b>0.128</b>	0.335	0.116
Diff.(DDIM)	3.18	0.181	0.561	0.548
Flow	<b>2.37</b>	0.150	<b>0.498</b>	<b>0.570</b>

**Table 5: Ablation study on different initial distributions**

Distribution	FID↓	Ali→	Ove→	mIoU↑
Gaussian	<b>2.37</b>	0.150	0.498	0.570
Uniform	2.61	<b>0.115</b>	<b>0.481</b>	<b>0.584</b>
Mixture	2.52	0.160	0.494	0.562

**Table 4: Ablation study on different loss functions**

Loss	FID↓	Ali→	Ove→	mIoU↑
$\mathcal{L}_{\text{CFM}}$	2.27	0.194	0.507	0.570
$\mathcal{L}_{\text{CFM}} + 0.1\mathcal{L}_1$	<b>2.20</b>	0.155	<b>0.495</b>	0.559
$\mathcal{L}_{\text{CFM}} + 0.2\mathcal{L}_1$	2.37	0.150	0.498	0.570
$\mathcal{L}_{\text{CFM}} + 0.3\mathcal{L}_1$	2.40	0.144	0.498	0.581
$\mathcal{L}_{\text{CFM}} + 0.4\mathcal{L}_1$	2.60	<b>0.134</b>	0.498	0.579
$L_1$ -loss only	26.64	0.159	0.625	<b>0.586</b>

**Table 6: Ablation study on different training trajectories**

Trajectory	FID↓	Ali→	Ove→	mIoU↑
Linear	2.37	<b>0.150</b>	0.498	<b>0.570</b>
Sine/Cosine	<b>2.33</b>	0.172	0.533	0.557
Sine	2.48	0.152	<b>0.480</b>	0.565

Bits have a smaller sample space, we tested using a Gaussian for the geometrical elements and a uniform distribution for the Analog Bits. Altogether, the results in Tab. 5 indicate that a Gaussian initial distribution provides the best overall performance measured by the *FID* score, whereas a uniform distribution provides stronger results in terms of geometrical metrics, such as *Alignment* and *Overlap*.

**Training Trajectories.** As long as the conditional vector field fulfills the criteria presented in Sec. 3, it can be freely chosen for training. In addition to the linear training trajectory proposed in [32,44], we also explore the sine/cosine interpolation introduced in [1] and a sine-based interpolation, as shown in Tab. 6. The choice of training trajectories only slightly affects the performance.

## 6 Conclusion

In this paper, we explored the application of Flow Matching for layout generation and, as a result, proposed LayoutFlow, a flow-based model that is able to handle various layout generation tasks. Our model is able to significantly speed up inference compared to other diffusion-based models while providing state-of-the-art performance. Even though we tested a diverse set of design options, there still remains room for further exploration into applying Flow Matching to layout generation, for example, different training trajectories, conditioning methods, or initial distributions. Furthermore, LayoutFlow might also be extended to generate layouts considering the content in the elements, similar to related work [19,46]. Overall, we demonstrate that Flow Matching provides a highly flexible and powerful tool for layout generation that offers a natural geometrical interpretation.

## Acknowledgments

This work was supported by JSPS KAKENHI Grant Number JP23K28139.

## References

1. Albergo, M.S., Vanden-Eijnden, E.: Building normalizing flows with stochastic interpolants. In: ICLR (2023)
2. Arroyo, D.M., Postels, J., Tombari, F.: Variational transformer networks for layout generation. In: CVPR (2021)
3. Austin, J., Johnson, D.D., Ho, J., Tarlow, D., Van Den Berg, R.: Structured denoising diffusion models in discrete state-spaces. In: NeurIPS (2021)
4. Chai, S., Zhuang, L., Yan, F.: Layoutdm: Transformer-based diffusion model for layout generation. In: CVPR (2023)
5. Chen, J., Zhang, R., Zhou, Y., Chen, C.: Towards aligned layout generation via diffusion model with aesthetic constraints. In: ICLR (2024)
6. Chen, R.T., Rubanova, Y., Bettencourt, J., Duvenaud, D.K.: Neural ordinary differential equations. In: NeurIPS (2018)
7. Chen, T., ZHANG, R., Hinton, G.: Analog bits: Generating discrete data using diffusion models with self-conditioning. In: ICLR (2023)
8. Cheng, C.Y., Huang, F., Li, G., Li, Y.: Play: Parametrically conditioned layout generation using latent diffusion. In: ICML (2023)
9. Deka, B., Huang, Z., Franzen, C., Hibsichman, J., Afergan, D., Li, Y., Nichols, J., Kumar, R.: Rico: A mobile app dataset for building data-driven design applications. In: UIST (2017)
10. Euler, L.: Institutiones calculi integralis, vol. 4. Academia Imperialis Scientiarum (1794)
11. Feng, W., Zhu, W., Fu, T.J., Jampani, V., Akula, A.R., He, X., Basu, S., Wang, X.E., Wang, W.Y.: LayoutGPT: Compositional visual planning and generation with large language models. In: NeurIPS (2023)
12. Gu, S., Chen, D., Bao, J., Wen, F., Zhang, B., Chen, D., Yuan, L., Guo, B.: Vector quantized diffusion model for text-to-image synthesis. In: CVPR (2022)
13. Guo, S., Jin, Z., Sun, F., Li, J., Li, Z., Shi, Y., Cao, N.: Vinci: an intelligent graphic design system for generating advertising posters. In: CHI (2021)
14. Gupta, K., Lazarow, J., Achille, A., Davis, L.S., Mahadevan, V., Shrivastava, A.: Layouttransformer: Layout generation and completion with self-attention. In: ICCV (2021)
15. He, L., Lu, Y., Corring, J., Florencio, D., Zhang, C.: Diffusion-based document layout generation. In: ICDAR (2023)
16. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: NeurIPS (2017)
17. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: NeurIPS (2017)
18. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: NeurIPS (2020)
19. Hsu, H.Y., He, X., Peng, Y., Kong, H., Zhang, Q.: Posterlayout: A new benchmark and approach for content-aware visual-textual presentation layout. In: CVPR (2023)

20. Hui, M., Zhang, Z., Zhang, X., Xie, W., Wang, Y., Lu, Y.: Unifying layout generation with a decoupled diffusion model. In: CVPR (2023)
21. Hyvärinen, A., Dayan, P.: Estimation of non-normalized statistical models by score matching. JMLR **6**(4) (2005)
22. Inoue, N., Kikuchi, K., Simo-Serra, E., Otani, M., Yamaguchi, K.: Layoutdm: Discrete diffusion model for controllable layout generation. In: CVPR (2023)
23. Jiang, Z., Guo, J., Sun, S., Deng, H., Wu, Z., Mijovic, V., Yang, Z.J., Lou, J.G., Zhang, D.: Layoutformer++: Conditional graphic layout generation via constraint serialization and decoding space restriction. In: CVPR (2023)
24. Jyothi, A.A., Durand, T., He, J., Sigal, L., Mori, G.: Layoutvae: Stochastic scene layout generation from a label set. In: ICCV (2019)
25. Kikuchi, K., Simo-Serra, E., Otani, M., Yamaguchi, K.: Constrained graphic layout generation via latent optimization. In: ACM MM (2021)
26. Kobayez, I., Prince, S.J., Brubaker, M.A.: Normalizing flows: An introduction and review of current methods. IEEE TPAMI **43**(11) (2020)
27. Lee, H.Y., Jiang, L., Essa, I., Le, P.B., Gong, H., Yang, M.H., Yang, W.: Neural design network: Graphic layout generation with constraints. In: ECCV (2020)
28. Levi, E., Brosh, E., Mykhailych, M., Perez, M.: Dlt: Conditioned layout generation with joint discrete-continuous diffusion layout transformer. In: ICCV (2023)
29. Li, J., Yang, J., Hertzmann, A., Zhang, J., Xu, T.: Layoutgan: Generating graphic layouts with wireframe discriminators. In: ICLR (2018)
30. Li, J., Yang, J., Zhang, J., Liu, C., Wang, C., Xu, T.: Attribute-conditioned layout gan for automatic graphic design. IEEE TVCG **27**(10) (2020)
31. Lin, J., Guo, J., Sun, S., Yang, Z.J., Lou, J.G., Zhang, D.: Layoutprompter: Awaken the design ability of large language models. In: NeurIPS (2023)
32. Lipman, Y., Chen, R.T.Q., Ben-Hamu, H., Nickel, M., Le, M.: Flow matching for generative modeling. In: ICLR (2023)
33. Liu, X., Gong, C., et al.: Flow straight and fast: Learning to generate and transfer data with rectified flow. In: ICLR (2022)
34. Ma, N., Goldstein, M., Albergo, M.S., Boffi, N.M., Vanden-Eijnden, E., Xie, S.: Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. arXiv preprint arXiv:2401.08740 (2024)
35. Neklyudov, K., Severo, D., Makhzani, A.: Action matching: A variational method for learning stochastic dynamics from samples. In: ICML (2023)
36. O'Donovan, P., Agarwala, A., Hertzmann, A.: Designscape: Design with interactive layout suggestions. In: CHI (2015)
37. O'Donovan, P., Agarwala, A., Hertzmann, A.: Learning layouts for single-pagegraphic designs. IEEE TVCG **20**(8) (2014)
38. Rahman, S., Sermuga Pandian, V.P., Jarke, M.: Ruite: Refining ui layout aesthetics using transformer encoder. In: 26th International Conference on Intelligent User Interfaces-Companion (2021)
39. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR (2022)
40. Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: ICML (2015)
41. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. In: ICLR (2020)
42. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Ermon, S., Poole, B.: Score-based generative modeling through stochastic differential equations. In: ICLR (2021)
43. Tang, Z., Wu, C., Li, J., Duan, N.: LayoutNUWA: Revealing the hidden layout expertise of large language models. In: ICLR (2024)

44. Tong, A., Malkin, N., Huguet, G., Zhang, Y., Rector-Brooks, J., Fatras, K., Wolf, G., Bengio, Y.: Improving and generalizing flow-based generative models with mini-batch optimal transport. In: ICMLW (2023)
45. Xie, Y., Huang, D., Wang, J., Lin, C.Y.: Canvasemb: Learning layout representation with large-scale pre-training for graphic design. In: ACM MM (2021)
46. Xu, C., Zhou, M., Ge, T., Jiang, Y., Xu, W.: Unsupervised domain adaption with pixel-level discriminator for image-aware layout generation. In: CVPR (2023)
47. Yamaguchi, K.: Canvasvae: Learning to generate vector graphic documents. In: ICCV (2021)
48. Zhang, J., Guo, J., Sun, S., Lou, J.G., Zhang, D.: Layoutdiffusion: Improving graphic layout generation by discrete diffusion probabilistic models. In: ICCV (2023)
49. Zhong, X., Tang, J., Yepes, A.J.: Publaynet: largest dataset ever for document layout analysis. In: ICDAR (2019)

## Supplementary Material

In this document, we provide additional explanations and evaluations to supplement the contents of the main paper. Specifically, we present further discussions on training trajectories, conditioning, and the limitations of our method, as well as additional quantitative and qualitative results.

### A Training Trajectories

The conditional flow, or training trajectory,  $\phi_t(\mathbf{x}) = \mathbf{x}_t$  can be chosen freely as long as the conditional vector field follows the continuity equation and  $\phi_0(\mathbf{x}) = \mathbf{x}_0$  and  $\phi_1(\mathbf{x}) = \mathbf{x}_1$ . While LayoutFlow is trained using a linear trajectory following [32, 44], we also tested the sine/cosine interpolation introduced in [1] and a sine-based interpolation. More details are shown in Tab. 7. Our main motivation for the sine-based interpolation was to preserve the direction of the conditional vector field but apply a different time schedule to the trajectory. In the linear scenario, the conditional vector field remains independent of the time and always points towards  $x_1$ . On the other hand, the sine/cosine trajectory changes its direction over time, only gradually leading to  $x_1$  in a circular trajectory. The sine-based interpolation maintains the direction of the linear trajectory but weighs it depending on the time. At an early time of the trajectory, the magnitude of the direction is large, whereas towards the end, the derivative starts to slow down. In analogy to time importance sampling in diffusion models [28], we hypothesized that this might help focus on the later stages of the flow process during training, where more detailed features become important. However, based on the paper’s ablation results, the cosine-based interpolation does not seem to affect the performance. Nonetheless, this scheduling property of conditional fields might be further explored with other functions of the general form:

$$v_t = \kappa(t)(\mathbf{x}_1 - \mathbf{x}_0), \quad (8)$$

while still fulfilling the conditions mentioned above. In the cosine case of  $\kappa(t) = \cos(\frac{\pi}{2}t)$ , the schedule might have been too close to the linear schedule to make a noticeable change, and different choices might provide better results.

### B Inference Speed Analysis

Since previous layout generation approaches rely on different architectures, we provide some insights on how the architectural components, *i.e.*, the model size measured by the numbers of trainable parameters and the number of tokens

**Table 7: Overview of different training trajectories and conditional vector fields.** Note that the conditional vector field is the time-derivative of the training trajectory.

Name	Training Trajectory $\mathbf{x}_t$	Conditional Vector Field $v_t$
Linear	$(1 - t)\mathbf{x}_0 + t\mathbf{x}_1$	$\mathbf{x}_1 - \mathbf{x}_0$
Sine/Cosine	$\cos(\frac{\pi}{2}t)\mathbf{x}_0 + \sin(\frac{\pi}{2}t)\mathbf{x}_1$	$\cos(\frac{\pi}{2}t)\mathbf{x}_1 - \sin(\frac{\pi}{2}t)\mathbf{x}_0$
Sine	$(1 - \sin(\frac{\pi}{2}t))\mathbf{x}_0 + \sin(\frac{\pi}{2}t)\mathbf{x}_1$	$\cos(\frac{\pi}{2}t)(\mathbf{x}_1 - \mathbf{x}_0)$

needed to represent a single layout, influence the inference speed in Table 8. In addition, we also report the number of steps each method requires to generate a layout that is tied to the choice of the generative model. Overall, it can be seen that the speed-up for LayoutFlow can be attributed to architectural improvements as well as changing from a diffusion model to a flow-based model due to the reduced number of generation steps required.

**Table 8: Overview of different factors contributing to inference speed**

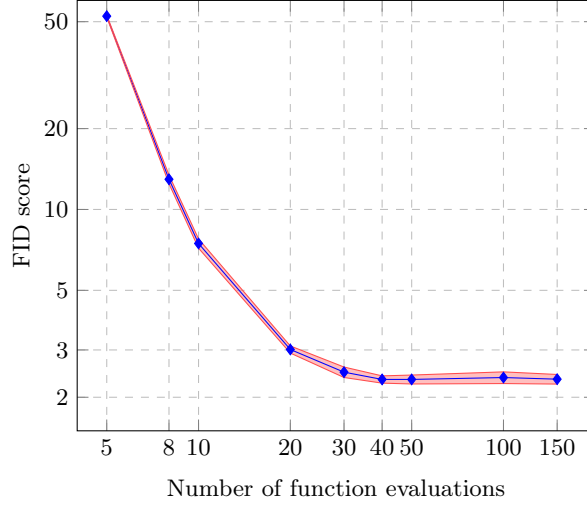
	Model Size	Token per Layout	Generation Steps	Inference Speed
LayoutDiffusion	86M	>100	160	1600.00ms
LayoutDM	12.4M	100	100	16.60ms
DLT	9.0M	20	100	3.50ms
LayoutFlow (ours)	12.7M	20	50	<b>1.75ms</b>

## C Quality Speed Trade-off

In Fig. 8, we compare the number of function evaluations of our model, corresponding to the number of steps taken to solve the ODE, to the performance on the RICO dataset. While only 5 function evaluations result in a bad score of an *FID* of roughly 50, just three more evaluations can bring down the *FID* to 13. The improvement observed by increasing the number of steps plateaus at around 40. In comparison, diffusion models typically require at least 100 steps to perform well. Overall, LayoutFlow offers a strong quality-speed trade-off, where slightly reducing the number of steps, *e.g.*, from 30 to 20, only has a small impact on performance.

## D Conditioning Analysis

There are several ways to introduce conditions to a generative task. For layout generation, given conditions can differ significantly depending on the task. For



**Fig. 8: Quality-Speed Trade-Off.** We investigate the relationship between the step size employed by the Euler method, which translates to the number of function evaluations, and the *FID* score using LayoutFlow trained on the RICO dataset.

LayoutFlow, we employ a masking strategy during training that indicates which part of the input serves as a condition, similar to [5, 20, 28]. Intuitively, the masking mechanism can be considered a switch that allows training multiple flow-based task-specific models with a single shared architecture in just one model. Alternatively, it is possible only to train an unconditional model and impose the conditions during sampling through a guidance strategy as done in [22, 48]. This strategy only requires training the model on the unconditional task and inserting the condition during inference, effectively altering the trajectory.

To condition the trajectory during inference, we employ a mask  $\mathbf{m}$  that is zero for the condition dimension, *i.e.*, all ones in the case of unconditional generation. Essentially, we use the mask only to update the direction of the condition dimension to  $\mathbf{u}'$  and can be described as

$$\mathbf{u}'_{\frac{k+1}{T}} = \mathbf{m} \odot \mathbf{u}_\theta(\mathbf{x}_{\frac{k}{T}} + \frac{1}{T} \mathbf{u}'_{\frac{k}{T}}) + (1 - \mathbf{m}) \odot (\mathbf{x}_c - \mathbf{x}_{\frac{k}{T}}),$$

where  $\mathbf{x}_c$  denotes the conditions and  $\odot$  represents the Hadamard product. For the conditional task, we update the direction only while  $k < 0.8T$  as we find it to work better empirically. Table 9 compares both conditioning methods and illustrates that a training-based approach significantly outperforms conditioning during inference. While we tried several other methods to introduce the conditions during inference, there might be more sophisticated conditioning methods that perform better and might be a potential future research direction.

**Table 9: Quantitative comparison of training LayoutFlow on multiple tasks using a masking strategy and inserting conditions to an unconditional LayoutFlow model by modifying the trajectory during inference.** The best results are highlighted in **bold**.

Task	Condition	RICO			
		FID↓	Ali→	Ove→	mIoU↑
Un-Gen	None	<b>2.28</b>	0.155	0.511	<b>0.610</b>
	Mask	2.37	<b>0.150</b>	<b>0.498</b>	0.570
Gen-Type	Trajectory	33.38	0.202	0.635	0.160
	Mask	<b>1.48</b>	<b>0.176</b>	<b>0.517</b>	<b>0.322</b>
Gen-TypeSize	Trajectory	174.66	0.862	0.728	0.181
	Mask	<b>1.03</b>	<b>0.283</b>	<b>0.523</b>	<b>0.470</b>
Completion	Trajectory	7.96	0.193	0.520	0.657
	Mask	<b>1.51</b>	<b>0.150</b>	<b>0.474</b>	<b>0.741</b>
<b>Validation Dataset</b>		2.10	0.093	0.466	0.658

## E Results Using Different Data Split

Due to the different data splits used by previous methods, comparing the numbers reported by other papers has become difficult. The FID values depend on a separately trained network, making comparing even more difficult as the same weights and network are required. Our experiments used the same settings as in [23, 48]. As a result, we had to retrain LayoutDM [22] and DLT [28] with that specific data split. To ensure that LayoutFlow works robustly across different dataset settings, we also trained a model on the data split used in LayoutDM [4], which includes up to 25 elements per layout as opposed to only 20 elements in the other setting. Additionally, this allows us to compare with a concurrent work called LACE [5], which uses a continuous diffusion model.

## F Completion Task

The completion task describes generating a complete layout given an incomplete layout. In this scenario, the given incomplete layout acts as a strong condition. There are different ways to define the completion task, *e.g.*, Inoue *et al.* assume that up to 20% of the layout is already given. While this scenario is closer to unconditional generation as it still offers a high degree of flexibility, we also looked into completing almost finished layouts containing more than 80% of their elements. To show that LayoutFlow can handle either scenario, we separately trained models with the respective conditional mask and present the results in Tab. 11 compared to LayoutDM and show the qualitative results in Appendix H.

**Table 10: Quantitative results on a different dataset split [22] for the PubLayNet dataset.** The best result is highlighted in **bold**. The  $\rightarrow$  symbol indicates best results are the ones closest to the validation data.

Task	Model	PubLayNet			
		FID↓	Ali→	Ove→	mIoU↑
Un-Gen	LACE	8.45	0.141	0.075	-
	LayoutDM	13.90	0.195	0.134	-
	LayoutFlow (ours)	<b>7.48</b>	<b>0.066</b>	<b>0.015</b>	0.423
Gen-Type	LACE	5.14	0.046	0.018	<b>0.383</b>
	LayoutDM	7.95	0.106	0.164	0.310
	LayoutFlow (ours)	<b>3.58</b>	<b>0.046</b>	<b>0.018</b>	0.349
Gen-TypeSize	LACE	2.66	0.061	<b>0.034</b>	0.418
	LayoutDM	4.25	0.119	0.189	0.381
	LayoutFlow (ours)	<b>0.80</b>	<b>0.052</b>	0.036	<b>0.443</b>
<b>Validation Data</b>		6.25	0.021	0.003	0.438

**Table 11: Quantitative results for the completion task with different percentages of the missing elements on the RICO and PubLayNet datasets.** The best results are highlighted in **bold**. Models marked with \* have been retrained.

Task	Model	RICO				PubLayNet			
		FID↓	Ali→	Ove→	mIoU↑	FID↓	Ali→	Ove→	mIoU↑
Completion (20%)	LayoutDM*	6.80	<b>0.054</b>	0.630	0.678	25.02	0.169	0.107	0.678
	LayoutFlow (ours)	<b>1.51</b>	0.150	<b>0.474</b>	<b>0.741</b>	<b>1.10</b>	<b>0.054</b>	<b>0.127</b>	<b>0.746</b>
Completion (80%)	LayoutDM*	5.21	<b>0.094</b>	0.658	0.574	28.73	0.223	0.146	0.385
	LayoutFlow (ours)	<b>3.59</b>	0.182	<b>0.605</b>	<b>0.628</b>	<b>4.02</b>	<b>0.050</b>	<b>0.024</b>	<b>0.445</b>
<b>Validation Dataset</b>		2.10	0.093	0.466	0.658	8.10	0.022	0.003	0.434

## G Limitations and Broader Impact

While our work aims to improve the state of layout generation further and help make design tasks easier, there are also some potential negative effects. Easier and better layout generation might be misused to build spam websites or large amounts of content or potentially be used for scams, particularly phishing attempts. The potential harm does not come from enabling such actions but derives from the scale that automation allows. Nonetheless, we believe the democratization of technologies like ours outweighs these risks.

Our proposed model, LayoutFlow, has shown remarkable results. However, some areas can still be improved in the future. Alignment remains one of the largest challenges, especially for models working in the continuous data space. While LayoutFlow utilizes an L1 regularization loss to improve the alignment, the issue is not fully resolved. For the future, it is crucial to find a better loss function that can reflect the perceptual error between layouts rather than relying

too heavily on metrics that do not necessarily correlate with important design aspects such as alignment.

## **H More Qualitative Results**

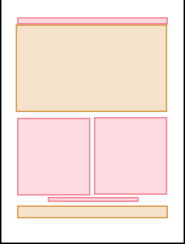
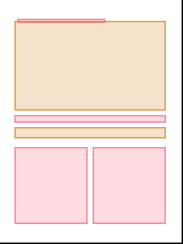
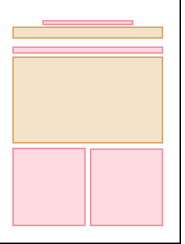
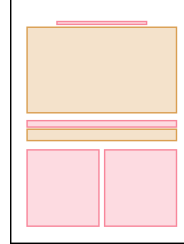
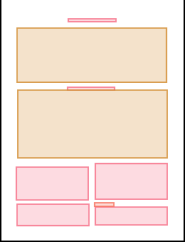
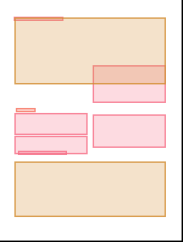
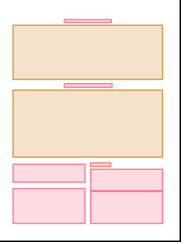
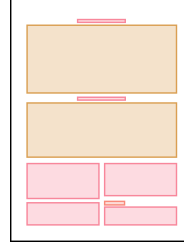
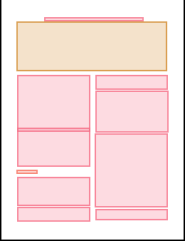
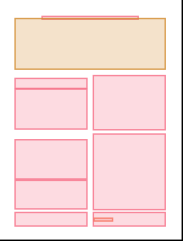
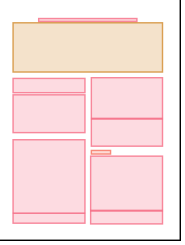
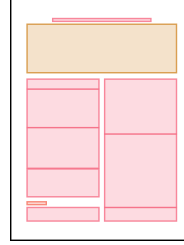
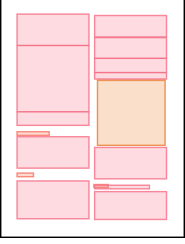
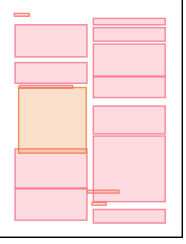
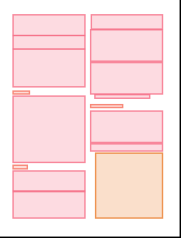
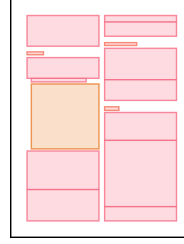
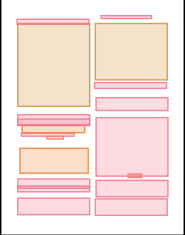
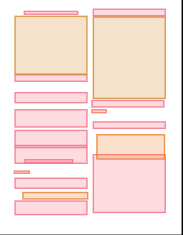
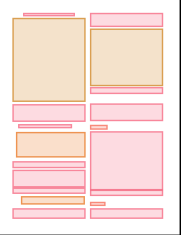
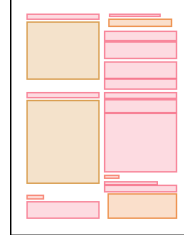
Due to limited space in the main paper, we present more qualitative results for each task and dataset in this section.

DLT	LayoutDM	LayoutDiffusion	LayoutFlow

**Table 12:** Unconditional Generation on PubLayNet

Input		DLT	LayoutDM	LayoutDiff.	LayoutFlow	Ground Truth
4x	Text					
	2x Table					
6x	Text					
	2x Table					
	Title					
7x	Text					
	Title					
	Figure					
8x	Text					
	Title					
	List					
9x	Text					
	Title					
	Table					

Table 13: Category Conditional Generation on PubLayNet

DLT	LayoutDM	LayoutFlow	Ground Truth
			
			
			
			
			

**Table 14:** Size Conditional Generation on PubLayNet

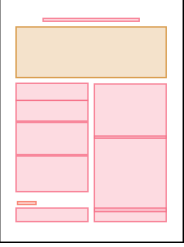

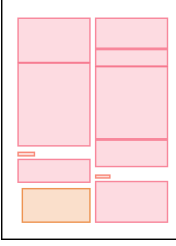
Input	LayoutDM	LayoutFlow	Ground Truth
			
			
			
			
			

Table 15: Element Completion (0-20%) on PubLayNet

Input	LayoutDM	LayoutFlow	Ground Truth
			
			
			
			
			

**Table 16:** Element Completion (80-100%) on PubLayNet

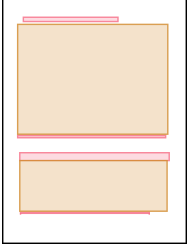
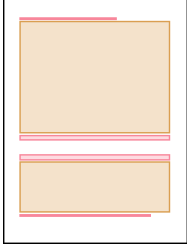
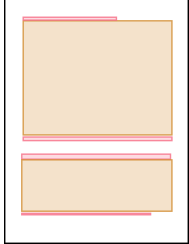
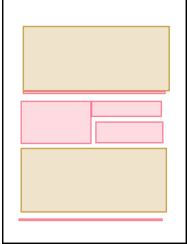
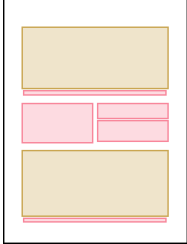
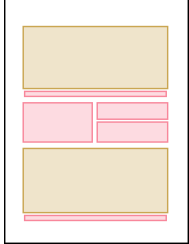
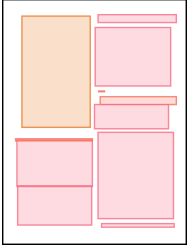
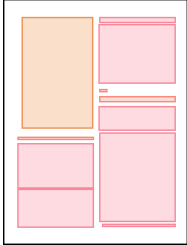
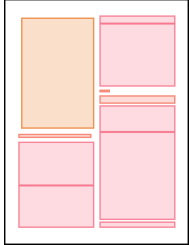
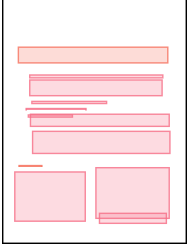
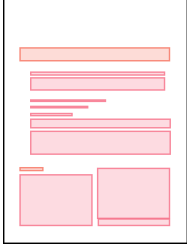
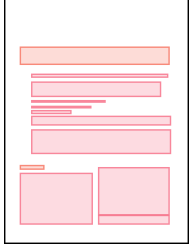
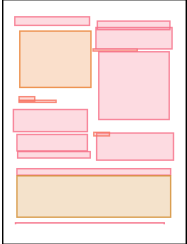
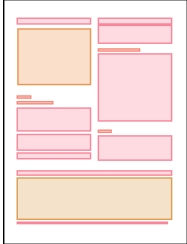
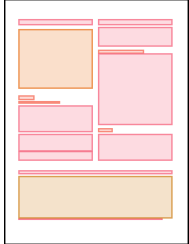
Input	LayoutFlow	Ground Truth
		
		
		
		
		

Table 17: Refinement on PubLayNet

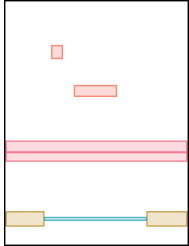
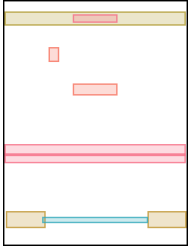
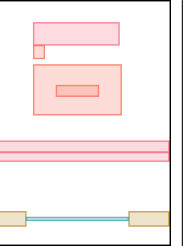
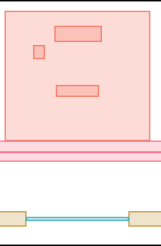
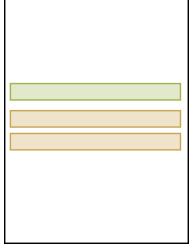
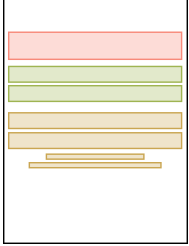

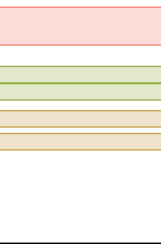
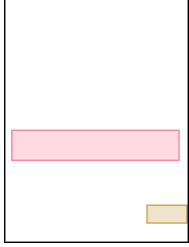
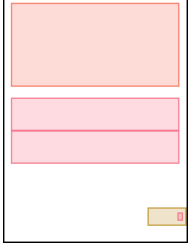
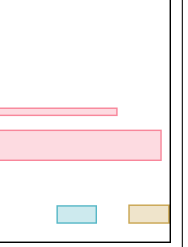
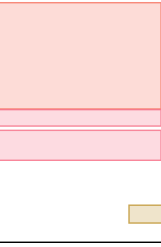


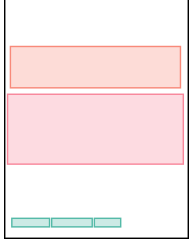
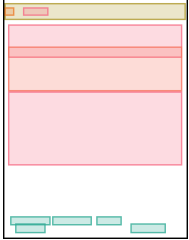
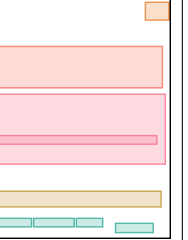
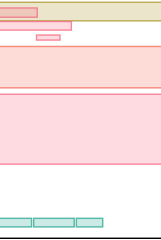
Input	LayoutDM	LayoutFlow	Ground Truth
			
			
			
			
			

Table 18: Element Completion (0-20%) on RICO

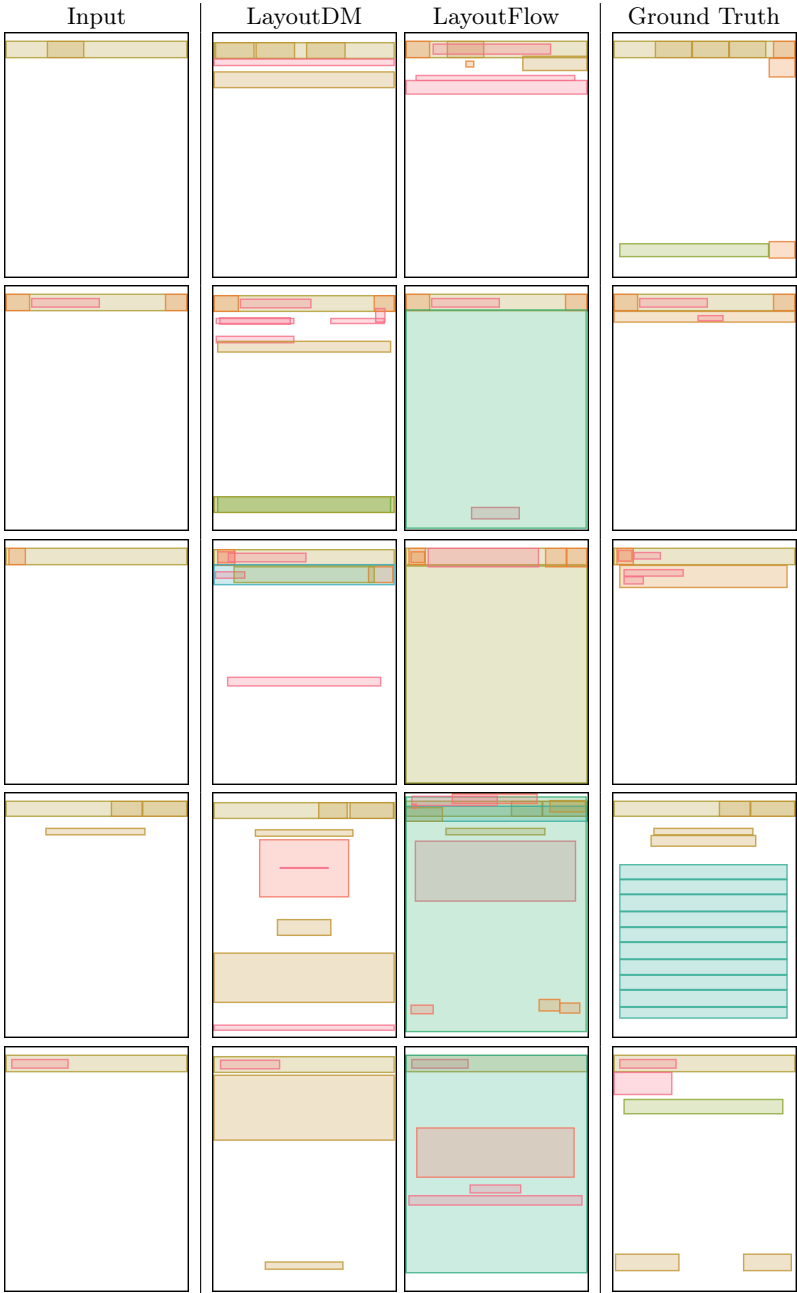


Table 19: Element Completion (80-100%) on RICO

Input	LayoutDM	LayoutDiffusion	LayoutFlow	Ground Truth

Table 20: Refinement on RICO