Edit Flows: Flow Matching with Edit Operations

arXiv 2025

Marton Havasi¹, Brian Karrer¹, Itai Gat¹, Ricky T. Q. Chen¹ ¹FAIR at Meta

2025.07.04

Seungwoo Yoo

KAIST

Non-autoregressive models have become the standard across high-dimensional modalities, including images and videos.



Esser et al., Stable Diffusion 3, ICML 2024

Wan 2.1, Team Wan, arXiv 2025

Recent efforts have aimed to extend these advances to discrete code and text generation, as exemplified by LLaDA [Nie *et al.*, 2025] and Mercury [Inception Labs, 2025].



Mercury, Inception Labs

However, these approaches do not support token insertion or deletion—operations that are fundamental to sequence generation and editing.

Discrete Flow Model

Discrete sequence trajectory: x_t



<u>Generative Flows on Discrete State-Spaces, Campbell et al., ICML 2024</u>

This work presents **Edit Flows**, an extension of Discrete Flow Matching [Gat *et al.*, 2024] that enables the training of discrete flow models generating variable-length sequences.

```
Insert-only Insert + Delete + Substitute

def is_prime(n: int) -> bool:
    def is_prime(n: int) -> bool:
        if n <= 1:
            return True
        for i in range(2, n):
            if i % n == 0:
                return True
        return True
        return False</pre>
```

Consider a continuous-time Markov chain (CTMC) over a discrete state space \mathcal{X} , a Markov process that generates trajectories $(X_t)_{t \in [0,1]}$ characterized by:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + hu_t(x | x_t) + o(h)$$

where u_t denotes a **rate**-the infinitesimal transition probabilities between states.

Consider a continuous-time Markov chain (CTMC) over a discrete state space \mathcal{X} , a Markov process that generates trajectories $(X_t)_{t \in [0,1]}$ characterized by:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + hu_t(x | x_t) + o(h)$$

where u_t denotes a **rate**-the infinitesimal transition probabilities between states.

Sampling from a CTMC is done by iteratively applying the above formula, which often starts from an initial sample $x_0 \sim p(x_0)$ and ends at $x_1 \sim q(x_1)$.



(a) Flow(b) Diffusion(c) Jump(d) CTMCLipman et al., Flow Matching Guide and Code, NeurIPS 2024 Tutorial

The rate $u_t(x|x_t)$ is the infinitesimal probabilities of transitioning from a state x_t to any other state x at time t, which must satisfy the **rate condition**

$$u_t(x|x_t) \ge 0$$
 for all $x \ne x_t$, $\sum_x u_t(x|x_t) = 0$

for the probability mass in the previous slide to be valid (*i.e.*, sum to one).

Note that this implies that

$$u_t(x_t|x_t) = -\sum_{x \neq x_t} u_t(x|x_t)$$

In CTMC, state transitions occur as discrete jumps over the course of continuous time:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + hu_t(x | x_t) + o(h)$$



Example adapted from <u>a talk by Andrew Campbell and Jason Yim</u>.

In CTMC, state transitions occur as discrete jumps over the course of continuous time:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + hu_t(x | x_t) + o(h)$$



Example adapted from <u>a talk by Andrew Campbell and Jason Yim</u>.

In CTMC, state transitions occur as discrete jumps over the course of continuous time:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + hu_t(x | x_t) + o(h)$$



Example adapted from <u>a talk by Andrew Campbell and Jason Yim</u>.

In CTMC, state transitions occur as discrete jumps over the course of continuous time:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + hu_t(x | x_t) + o(h)$$



Example adapted from <u>a talk by Andrew Campbell and Jason Yim</u>.

In CTMC, state transitions occur as discrete jumps over the course of continuous time:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + \frac{hu_t(x | x_t)}{hu_t(x | x_t)} + o(h)$$



Example adapted from <u>a talk by Andrew Campbell and Jason Yim</u>.

A rate u_t is said to generate a **probability path** p_t if the time marginals of the associated CTMC are samples from p_t , *i.e.*, $X_t \sim p_t$.

Remark: The probability path is shown in continuous state space for visualization purposes only.

Lipman et al., Flow Matching Guide and Code, NeurIPS 2024 Tutorial

A rate u_t is said to generate a **probability path** p_t if the time marginals of the associated CTMC are samples from p_t , *i.e.*, $X_t \sim p_t$.

Specifically, such a probability path satisfies the Kolmogorov Forward Equation (KFE):

$$\frac{\partial}{\partial t}p_t(x) = \sum_{y} u_t(x|y)p_t(y) = \sum_{y \neq x} u_t(x|y)p_t(y) - \sum_{y \neq x} u_t(y|x)p_t(x)$$
"The rate of change of a state's probability equals the net probability flux at that state."

Lipman et al., Flow Matching Guide and Code, NeurIPS 2024 Tutorial

Consider a discrete space over sequences of fixed length $N, \mathcal{X} = \mathcal{T}^N$ where $\mathcal{T} = \{1, 2, \dots, M\}$ is a vocabulary of size M.

The goal of DFM is to learn the marginal rate $u_t(x|x_t)$ and simulate transitions of a CTMC:

$$\mathbb{P}(X_{t+h} = x | X_t = x_t) = \delta_{x_t}(x) + hu_t(x | x_t) + o(h)$$

that generates a probability path p_t interpolating an easy-to-sample distribution and the data distribution.

Lipman et al., Flow Matching Guide and Code, NeurIPS 2024 Tutorial

However, directly regressing the marginal rate is intractable, as it involves averaging over the data distribution. Instead, we can use the **conditional rate** as the regression target:

$$u_t(x|x_t, x_0, x_1)$$

which generates a **conditional probability path** $p_t(x|x_0, x_1)$ with $p_0(x|x_0, x_1) = \delta_{x_0}(x)$ and $p_1(x|x_0, x_1) = \delta_{x_1}(x)$.

However, directly regressing the marginal rate is intractable, as it involves averaging over the data distribution. Instead, we can use the **conditional rate** as the regression target:

$$u_t(x|x_t, x_0, x_1)$$

which generates a **conditional probability path** $p_t(x|x_0, x_1)$ with $p_0(x|x_0, x_1) = \delta_{x_0}(x)$ and $p_1(x|x_0, x_1) = \delta_{x_1}(x)$.

 (x_0, x_1) is a sample from a **coupling distribution** $\pi(x_0, x_1)$, *i.e.*, $(x_0, x_1) \sim \pi(x_0, x_1)$, which, in its simplest form, is modeled as the product of two independent densities:

$$\pi(x_0, x_1) = p(x_0)q(x_1)$$

The marginal probability path and marginal rate can be obtained via simple averaging:

$$p_t(x) = \sum_{x_0, x_1} p_t(x|x_0, x_1) \pi(x_0, x_1)$$
$$u_t(x|x_t) = \mathbb{E}_{p_t(x_0, x_1|x_t)} u_t(x|x_t, x_0, x_1)$$
where $p_0(x) = p(x)$ and $p_1(x) = q(x)$.

One popular target conditional probability path $p_t(x|x_0, x_1)$ is the **factorized token-wise probability path** [Gat *et al.*, NeurIPS 2024]:

$$p_t(x^i | x_0^i, x_1^i) = (1 - \kappa_t) \delta_{x_0^i}(x^i) + \kappa_t \delta_{x_1^i}(x^i)$$
$$u_t(x^i | x_t^i, x_0^i, x_1^i) = \frac{\dot{\kappa}_t}{1 - \kappa_t} \left(\delta_{x_1^i}(x^i) - \delta_{x_t^i}(x^i) \right)$$

where κ_t is a scheduler satisfying $\kappa_0 = 0$ and $\kappa_1 = 1$.

One popular target conditional probability path $p_t(x|x_0, x_1)$ is the **factorized token-wise probability path** [Gat *et al.*, NeurIPS 2024]:

$$p_t(x^i | x_0^i, x_1^i) = (1 - \kappa_t) \delta_{x_0^i}(x^i) + \kappa_t \delta_{x_1^i}(x^i)$$
$$u_t(x^i | x_t^i, x_0^i, x_1^i) = \frac{\dot{\kappa}_t}{1 - \kappa_t} \left(\delta_{x_1^i}(x^i) - \delta_{x_t^i}(x^i) \right)$$

where κ_t is a scheduler satisfying $\kappa_0 = 0$ and $\kappa_1 = 1$.

The marginal distributions of (corrupted) sequences are modeled as:

$$p_t(x|x_0, x_1) = \prod_{i=1}^{N} p_t(x^i|x_0^i, x_1^i)$$

that are characterized by the conditional rate:

$$u_t(x|x_t, x_0, x_1) = \sum_i \delta_{x_t}(x^{\neg i}) u_t(x^i|x_t^i, x_0^i, x_1^i) \quad \text{with } \delta_{x_t}(x^{\neg i}) = \prod_{j \neq i} \delta_{x_t^j}(x^j)$$

However, DFM has difficulty generating sequences of varying lengths, as illustrated in the example below.

However, DFM has difficulty generating sequences of varying lengths, as illustrated in the example below.

However, DFM has difficulty generating sequences of varying lengths, as illustrated in the example below.

Where to put a new token?

However, DFM has difficulty generating sequences of varying lengths, as illustrated in the example below.

Where to put a new token?

Which token to put?

However, DFM has difficulty generating sequences of varying lengths, as illustrated in the example below.

However, DFM has difficulty generating sequences of varying lengths, as illustrated in the example below.

In this work, we extend the DFM framework by allowing two sequences $x, x_t \in \bigcup_{n=0}^{N} \mathcal{T}^n$ differ by one *edit operation*, which is one of *insertion*, *deletion*, and *substitution*:

- $ins(x, i, a) = (x^1, ..., x^i, a, x^{i+1}, ..., x^{n(x)})$
- $del(x,i) = (x^1, ..., x^{i-1}, x^{i+1}, ..., x^{n(x)})$
- $\operatorname{sub}(x, i, a) = (x^1, \dots, x^{i-1}, a, x^{i+1}, \dots, x^{n(x)})$

Edit Flow sampling process.

Since insertions, deletions, and substitutions produce mutually exclusive outcomes, we can define separate rates for each operation:

- $u_t^{\theta}(\operatorname{ins}(x, i, a)|x) = \lambda_{t,i}^{\operatorname{ins}}(x)Q_{t,i}^{\operatorname{ins}}(a|x)$
- $u_t^{\theta}(\operatorname{del}(x,i)|x) = \lambda_{t,i}^{\operatorname{del}}(x)$

for $i \in \{1, ..., n(x)\}$ for $i \in \{1, ..., n(x)\}$ for $i \in \{1, ..., n(x)\}$

• $u_t^{\theta}(\operatorname{sub}(x, i, a)|x) = \lambda_{t,i}^{\operatorname{sub}}(x)Q_{t,i}^{\operatorname{sub}}(a|x)$

Since insertions, deletions, and substitutions produce mutually exclusive outcomes, we can define separate rates for each operation:

- $u_t^{\theta}(\operatorname{ins}(x, i, a)|x) = \lambda_{t,i}^{\operatorname{ins}}(x)Q_{t,i}^{\operatorname{ins}}(a|x)$
- $u_t^{\theta}(\operatorname{del}(x,i)|x) = \lambda_{t,i}^{\operatorname{del}}(x)$
- $u_t^{\theta}(\operatorname{sub}(x, i, a)|x) = \lambda_{t,i}^{\operatorname{sub}}(x)Q_{t,i}^{\operatorname{sub}}(a|x)$

for $i \in \{1, ..., n(x)\}$ for $i \in \{1, ..., n(x)\}$ for $i \in \{1, ..., n(x)\}$

The total rates of editing operations "Which operation should be applied to which token in *x*?"

Since insertions, deletions, and substitutions produce mutually exclusive outcomes, we can define separate rates for each operation:

- $u_t^{\theta}(\operatorname{ins}(x, i, a)|x) = \lambda_{t,i}^{\operatorname{ins}}(x)Q_{t,i}^{\operatorname{ins}}(a|x)$
- $u_t^{\theta}(\operatorname{del}(x,i)|x) = \lambda_{t,i}^{\operatorname{del}}(x)$
- $u_t^{\theta}(\operatorname{sub}(x, i, a)|x) = \lambda_{t,i}^{\operatorname{sub}}(x)Q_{t,i}^{\operatorname{sub}}(a|x)$

for $i \in \{1, ..., n(x)\}$ for $i \in \{1, ..., n(x)\}$ for $i \in \{1, ..., n(x)\}$

Normalized distributions over token values $a \in \mathcal{T}$ "Which value should be added or used to replace the selected token?"

Indeed, Edit Flow is a generalization of existing techniques, each of which can be obtained by restricting the rates.

- **DFM**: A special case that permits only substitutions, while disallowing token insertions and deletions (*i.e.*, $\lambda_{t,i}^{\text{ins}} = 0$ and $\lambda_{t,i}^{\text{del}} = 0$) \rightarrow **Lacks support for varying lengths**;
- AR: A special case that permits insertions to occur at the rightmost location (*i.e.*, $\lambda_{t,n(x)}^{\text{ins}} \neq 0 \rightarrow \text{Incapable of making corrections}$.

Edit Flows are powerful. But how can we train them?

The main challenge lies in defining the conditional probability path $p_t(x|x_0, x_1)$ and its rate $u_t(x|x_0, x_1)$ that our model can regress due to **the enormous number of possible transitions** that take one sequence to another:

K $arepsilon$ ITTEN	ΚΙΤΤΕΝεεεεεε	ΚΙΤΤΕΝεεεεεε
$\downarrow \downarrow$	$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$	$\downarrow \downarrow $
SMITTEN	S M I T T Ε Ν εεεεε	εεεεε Ε Μ Ι Τ Τ Ε Ν
Optimal	Sub-optimal	Least Optimal

The key idea for simplifying the problem is to **augment the space** $\mathcal{X} = \bigcup_{n=0}^{N} \mathcal{T}^{n}$ by introducing an auxiliary token ε and define an auxiliary Markov process via **alignment**.

Edit operations can be recovered as tuples $(a \rightarrow b)$ with $a, b \in \mathcal{T} \cup \{\varepsilon\}$.

Insertion: $a = \varepsilon$

Edit operations can be recovered as tuples $(a \rightarrow b)$ with $a, b \in \mathcal{T} \cup \{\varepsilon\}$.

Deletion: $b = \varepsilon$

Edit operations can be recovered as tuples $(a \rightarrow b)$ with $a, b \in \mathcal{T} \cup \{\varepsilon\}$.

Substitution: $a \neq \varepsilon$ and $b \neq \varepsilon$

Formally, let

- $\mathcal{X} = \bigcup_{n=0}^{N} \mathcal{T}^{n}$: The space of normal (raw) sequences of different lengths;
- $\mathcal{Z} = (\mathcal{T} \cup \{\varepsilon\})^N$: The space of aligned sequences (*i.e.*, padded with ε 's);
- $f_{\rm rm-blanks}$ ($\mathcal{Z} \to \mathcal{X}$): An operation that strips away all ε tokens;
 - $K_{\underline{\varepsilon}}$ ITTEN \rightarrow KITTEN
 - KITTEN<u> $\epsilon\epsilon\epsilon$ </u> \rightarrow KITTEN
 - $K_{\underline{\varepsilon}}I_{\underline{\varepsilon}}T_{\underline{\varepsilon}}T_{\underline{\varepsilon}}E_{\underline{\varepsilon}}N \rightarrow KITTEN$

Given samples from the source $x_0 \sim p(x)$ and target $x_1 \sim q(x)$ in \mathcal{X} , one can construct aligned sequences z_0 and z_1 by:

- randomly padding the sequences;
- solving for the optimal alignment that corresponds to the minimal edit distance.

This yields a coupling $\pi(z_0, z_1)$ over the auxiliary variables satisfying

$$p(x) = \sum_{z_0} \sum_{z_1} \pi(z_0, z_1) \delta_{f_{rm-blanks}(z_0)}(x)$$
$$q(x) = \sum_{z_0} \sum_{z_1} \pi(z_0, z_1) \delta_{f_{rm-blanks}(z_1)}(x)$$

Furthermore, for $z_0, z_1 \sim \pi$, we define a conditional probability path over the augmented space $\mathcal{X} \times \mathcal{Z}$:

$$p_t(x, z | x_0, z_0, x_1, z_1) = p_t(x, z | z_0, z_1) = p_t(z | z_0, z_1) \delta_{f_{\text{rm-blanks}(z)}}(x)$$

where $p_t(z|z_0, z_1)$ is a token-wise mixture probability path:

$$p_t(z|z_0, z_1) = \prod_{i=1}^N p_t(z^i|z_0^i, z_1^i)$$

The conditional rate $u_t(x, z | x_t, z_t, z_0, z_1)$ corresponding to the probability path is:

$$u_{t}(x, z | x_{t}, z_{t}, z_{0}, z_{1}) = \delta_{f_{\text{rm-blanks}(z)}}(x) \sum_{i=1}^{N} \frac{\dot{\kappa}_{t}}{1 - \kappa_{t}} \left(\delta_{z_{1}^{i}}(z^{i}) - \delta_{z_{t}^{i}}(z^{i}) \right) \delta_{z_{t}}(z^{\neg i})$$

which is analogous to the conditional rate discussed in DFM:

$$u_t(x|x_t, x_0, x_1) = \sum_i \frac{\dot{\kappa}_t}{1 - \kappa_t} \left(\delta_{x_1^i}(x^i) - \delta_{x_t^i}(x^i) \right) \delta_{x_t}(x^{\neg i})$$

By marginalizing the conditional rate $u_t(x, z | x_t, z_t, z_0, z_1)$, we obtain the unconditional rate of a CTMC over the original state space \mathcal{X} :

$$u_t(x|x_t) = \sum_{z} \mathbb{E}_{p_t(z_0, z_1, z_t|x_t)} u_t(x, z|x_t, z_t, z_0, z_1)$$

It is proved that this unconditional rate can be learned by minimizing the training loss of form:

$$\mathcal{L}(\theta) = \mathbb{E}_{\substack{\pi(z_0, z_1) \\ t, p_t(x_t, z_t | z_0, z_1)}} \left[\sum_{x \neq x_t} u_t^{\theta}(x | x_t) - \sum_{i=1}^N \mathbf{1}_{[z_1^i \neq z_t^i]} \frac{\dot{\kappa}_t}{1 - \kappa_t} \log u_t^{\theta}(x(z_t, i, z_1^i) | x_t) \right]$$

where $x(z_t, i, z_1^i) = f_{\text{rm-blanks}}(z_t^1, \dots, z_t^{i-1}, z_1^i, z_t^{i+1}, \dots, z_t^N)$ is an output of one of the possible edit operations (insertion, deletion, and substitution).

It is proved that this unconditional rate can be learned by minimizing the training loss of form:

$$\mathcal{L}(\theta) = \mathbb{E}_{\substack{\pi(z_0, z_1) \\ t, p_t(x_t, z_t | z_0, z_1)}} \left[\sum_{\substack{x \neq x_t}} u_t^{\theta}(x | x_t) - \sum_{i=1}^N \mathbf{1}_{[z_1^i \neq z_t^i]} \frac{\dot{\kappa}_t}{1 - \kappa_t} \log u_t^{\theta}(x(z_t, i, z_1^i) | x_t) \right]$$

Regularize spurious editing operations

It is proved that this unconditional rate can be learned by minimizing the training loss of form:

$$\mathcal{L}(\theta) = \mathbb{E}_{\substack{\pi(z_0, z_1) \\ t, p_t(x_t, z_t | z_0, z_1)}} \left[\sum_{x \neq x_t} u_t^{\theta}(x | x_t) - \sum_{i=1}^N \mathbf{1}_{[z_1^i \neq z_t^i]} \frac{\dot{\kappa}_t}{1 - \kappa_t} \log u_t^{\theta}(x(z_t, i, z_1^i) | x_t) \right]$$

Weighted cross-entropy over edit operations for $x_t \rightarrow x_1$

It is proved that this unconditional rate can be learned by minimizing the training loss of form:

$$\mathcal{L}(\theta) = \mathbb{E}_{\substack{\pi(z_0, z_1) \\ t, p_t(x_t, z_t | z_0, z_1)}} \left[\sum_{x \neq x_t} u_t^{\theta}(x | x_t) - \sum_{i=1}^N \mathbf{1}_{[z_1^i \neq z_t^i]} \frac{\dot{\kappa}_t}{1 - \kappa_t} \log u_t^{\theta}(x(z_t, i, z_1^i) | x_t) \right]$$

where $x(z_t, i, z_1^i) = f_{rm-blanks}(z_t^1, ..., z_t^{i-1}, z_1^i, z_t^{i+1}, ..., z_t^N)$ is an output of one of the possible edit operations (insertion, deletion, and substitution).

Theorem 3.1 (Flow Matching with Auxiliary Processes). Let $u_t(x, z | x_t, z_t)$ be a rate over the augmented space of $\mathcal{X} \times \mathcal{Z}$ that generates $p_t(x, z)$, then

$$u_t(x|x_t) \triangleq \sum_z \mathbb{E}_{p_t(z_t|x_t)} u_t(x, z|x_t, z_t) \qquad generates \qquad p_t(x) \triangleq \sum_z p_t(x, z), \tag{17}$$

and furthermore, for any Bregman divergence $D_{\phi}(a,b) = \phi(a) - \phi(b) - \langle a - b, \frac{d}{db}\phi(b) \rangle$ defined by a convex function ϕ , we have that

$$\frac{\mathrm{d}}{\mathrm{d}\theta} \mathbb{E}_{x_t, z_t \sim p_t(x, z)} D_\phi \Big(\sum_z u_t(\cdot, z | x_t, z_t), u_t^\theta(\cdot | x_t) \Big) = \frac{\mathrm{d}}{\mathrm{d}\theta} \mathbb{E}_{x_t \sim p_t(x)} D_\phi \left(u_t(\cdot | x_t), u_t^\theta(\cdot | x_t) \right).$$
(18)

Sampling Edit Flows

To sample from the model, we transport a source sample $x_0 \sim p$ to time t = 1, by simulating the CTMC defined with the learned rate $u_t^{\theta}(x|x_t)$. The exact probability of an edit operation characterized by the rate $\lambda_{t,i}$ occurring within (t, t + h) is:

 $1 - e^{-\int_{t}^{t+h} \lambda_{t,i}(X_t) dt} \approx h \lambda_{t,i}(X_t)$ (Taylor Expansion)

Sampling Edit Flows

As with previous discrete diffusion/flow models, each next states are sampled **independently**. The sampling proceeds as follows:

- 1. For each position *i*, sample whether to insert with probability $h\lambda_{t,i}^{ins}(X_t)$ and whether to delete or substitute with probability $h\left(\lambda_{t,i}^{del}(X_t) + \lambda_{t,i}^{sub}(X_t)\right)$;
- 2. If insertion or substitution occurs at token *i*, sample the new token value from $Q_{t,i}^{ins}(a|X_t)$ or $Q_{t,i}^{sub}(a|X_t)$;
- 3. Move to the next timestep: $t \leftarrow t + h$.

Sampling Edit Flows

Following [Nisonoff *et al.*, 2024], the authors apply classifier-free guidance (CFG) [Ho and Salimans, 2021] during sampling, by adjusting the predicted rates:

$$\tilde{u}_t(x|x_t,c) = \tilde{\lambda}_{t,i}(x_t,c)\tilde{Q}_t(a|x_t,c)$$

where $\tilde{\lambda}_t(x_t, c) = \lambda_{t,i}(x_t|c)^{1+w}\lambda_{t,i}(x_t)^{-w}$ and w is guidance strength.

Edit Flow is mainly tested on text generation tasks, including:

- 1. Image-to-Text Generation (280M models);
- 2. Text Generation (1.3B models);
- 3. Code Generation (1.3B models).

As baselines, the authors consider SotA AR/non-AR models:

- 1. Llama 2 [Touvron et al., 2023]
- 2. Masked DFM [Gat et al., 2024]

These baselines are compared against three variants of Edit Flow:

- **1.** Default Model: Starts from $p = \delta_{\phi}$ (an empty string) and inserts/delete tokens;
- 2. Uniform X₀ + Edit Flow: $X_0 = (X^1, ..., X^{100})$ where $X^i \sim p_{emp}$ using empirical training token distribution;
- **3.** Localized Edit Flow: A variant of edit flow using a localized propagation process.

B.1 Localized propagation paths

Edit Flows leverage an underlying conditional probability path $p_t(z|z_0, z_1)$ and associated rates $u_t(z|z_t)$, so far given by the factorized token-wise mixture. Let us further generalize this probability path and associated rate to be non-factorized, applying auxiliary variables again. We first re-express this probability path through an auxiliary boolean variable $m \in \{false, true\}^N$:

$$p_t(z|z_0, z_1) = \sum_{\mathsf{m}} p_t(\mathsf{m}|z_0, z_1) p_t(z|\mathsf{m}, z_0, z_1),$$

where $p_t(z|\mathsf{m}, z_0, z_1) = \prod_{i=1}^N \mathbb{1}_{[\neg \mathsf{m}^i]} \delta_{z_0^i}(z^i) + \mathbb{1}_{[\mathsf{m}^i]} \delta_{z_1^i}(z^i),$

(37)

(38)

All baselines and variants of Edit Flow share the architecture from Llama [Grattafiori *et al.*, 2024; Touvron *et al.*, 2023], with 280M and 1.3B parameters.

- All models are trained using the same compute budget;
- The maximum sequence length during training is set to 1024 tokens;
- AR baselines use causal attention, while non-AR methods full self-attention;
- For Edit Flow, FlexAttention [Dong *et al.*, 2024] is used to handle varying lengths.

For image captioning tasks, all methods are trained on the MS COCO dataset [Lin *et al*., 2014] and Image Captioning 3M dataset.

Method	MS	S COCO		Image Captioning 3M			
	METEOR	CIDEr	SPICE	ROUGE-L	CIDEr	SPICE	
$\overline{\rm VLP^{\dagger}}$ (Zhou et al., 2020)	28.4	117.7	21.3	24.3	77.5	16.5	
$ClipCap^{\dagger}$ (Mokady et al., 2021)	27.1	108.3	20.1	26.7	87.2	18.5	
Autoregressive	25.7	95.5	19.6	25.2	85.8	17.8	
Mask DFM	25.3	95.6	19.2	27.4	96.2	20.3	
Edit Flow (Ours)	27.4	108.1	21.1	29.0	101.9	21.7	
Localized Edit Flow (Ours)	27.4	105.1	22.1	28.3	99.7	20.8	

For text benchmarks, all methods are trained on the DCLM baseline 1.0 [Li *et al.*, 2024] dataset. Edit Flow bridges the gap between non-AR- and AR-based methods.

Method	HellaSwag	ARC-E	ARC-C	PIQA	OBQA	WinoGrande
Autoregressive	49.5	71.0	36.3	76.0	30.4	62.1
Mask DFM	38.3	55.4	27.8	65.3	22.6	52.3
Edit Flow (CFG applied to u_t) (Ours)	49.0	63.1	33.0	68.8	28.6	53.6
Edit Flow (CFG applied to \mathcal{L}) (Ours)	54.5	61.0	34.0	65.0	37.2	54.3

For code generation, all methods are trained on the CodeLlama datamix [Roziere *et al.*, 2023]. *Oracle Length* denotes the case where GT code length is known during generation.

Method		HumanEval		HumanEval+		MBPP	
		Pass@1	Pass@10	Pass@1	Pass@10	Pass@1	Pass@10
	Autoregressive (Gat et al., 2024)	14.3	21.3			17.0	34.3
	Autoregressive [†]	17.0	34.7	14.0	28.6	25.6	45.4
Non-AR	Mask DFM (Gat et al., 2024)	6.7	13.4			6.7	20.6
	Mask DFM (Oracle Length) (Gat et al., 2024)	11.6	18.3			13.1	28.4
	$\mathrm{Mask} \ \mathrm{DFM}^\dagger$	9.1	17.6	7.9	13.4	6.2	25.0
	Uniform X_0 + Edit Flow (Ours)	9.7	24.3	9.7	19.5	9.4	33.4
	Edit Flow (Ours)	12.8	24.3	10.4	20.7	10.0	36.4
	Localized Edit Flow (Ours)	14.0	22.6	10.4	18.9	14.8	34.0

Conclusion

To summarize, this paper presents:

- Edit Flow, a novel extension of Discrete Flow Matching [Gat *et al.*, 2024], which enables variable-length generation via edit operations-insertion, deletion, and substitution;
- A novel CTMC with conditional rates defined over an augmented state space, enabling a tractable training objective for Edit Flow;
- Extensive evaluations in image captioning and text/code generation demonstrate that non-autoregressive models can serve as potential replacements for current AR models.

Edit Flows: Flow Matching with Edit Operations

```
Insert-only Insert + Delete + Substitute

def is_prime(n: int) -> bool:
    def is_prime(n: int) -> bool:
        if n <= 1:
            return True
        for i in range(2, n):
            if i % n == 0:
                return True
        reture
        return True
        retu
```

Seungwoo Yoo KAIST <u>https://dvelopery0115.github.io</u>