

A Continuous-Time Markov Chain Framework for Insertion Language Models

Dhruvesh Patel*

Benjamin Rozonoyer*

Soumitra Das*

Tahira Naseem†

Tim G. J. Rudner‡§

Andrew McCallum*

AISTATS 2026

UMassAmherst

UMass Amherst*



IBM Research†



UNIVERSITY OF
TORONTO

University of
Toronto‡



Vijil§

Why insertions?

ARMs hard-code a left-to-right ordering assumption. Not suitable for some tasks.

ARM	<input type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length
	<input checked="" type="checkbox"/> Probabilistically Principled	<input type="checkbox"/> Multi-token

The quick brown fox jumps over the **the**

The quick brown fox jumps over the lazy **lazy**

The quick brown fox jumps over the lazy dog **dog**

Why insertions?

ARMs hard-code a left-to-right ordering assumption. Not suitable for some tasks.

MDMs relax generation order. Variable-length generation is non-trivial.

ARM	<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length
	<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Multi-token
The quick brown fox jumps over the		
The quick brown fox jumps over the lazy		
The quick brown fox jumps over the lazy dog		
MDM	<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length
	<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Multi-token
The quick [M] fox jumps over the [M] [M] [M]		
The quick brown fox jumps over the lazy [M] [M]		
The quick brown fox jumps over the lazy dog [P] [P]		

Why insertions?

ARMs hard-code a left-to-right ordering assumption. Not suitable for some tasks.

MDMs relax generation order. Variable-length generation is non-trivial.

ILMs have the right flexibility, but existing objectives and sampling rules are ad hoc.

ARM	<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Arbitrary length
	<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Multi-token
The quick brown fox jumps over the		
The quick brown fox jumps over the lazy		
The quick brown fox jumps over the lazy dog		
MDM	<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Multi-token
	<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length
The quick [M] fox jumps over the [M] [M] [M]		
The quick brown fox jumps over the lazy [M] [M]		
The quick brown fox jumps over the lazy dog [P] [P]		
ILM	<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Multi-token
	<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length
The quick fox over the lazy dog		
The quick fox jumps over the lazy dog		
The quick brown fox jumps over the lazy dog		

Principled Insertion-based Generation

ARMs hard-code a left-to-right ordering assumption. Not suitable for some tasks.

MDMs relax generation order. Variable-length generation is non-trivial.

ILMs have the right flexibility, but existing objectives and sampling rules are ad hoc.

DILMs have diffusion-style denoising objective and sampling to combine all the benefits.

ARM

<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Arbitrary length
<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Multi-token

The quick brown fox jumps over the
The quick brown fox jumps over the lazy
The quick brown fox jumps over the lazy dog

MDM

<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Multi-token
<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length

The quick [M] fox jumps over the [M] [M] [M]
The quick brown fox jumps over the lazy [M] [M]
The quick brown fox jumps over the lazy dog [P] [P]

ILM

<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Multi-token
<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length

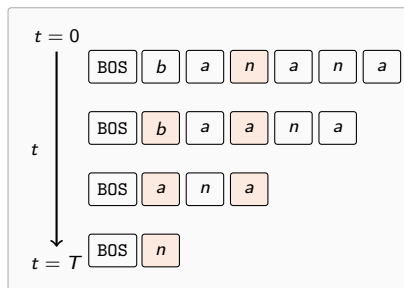
The quick fox over the lazy dog
The quick fox jumps over the lazy dog
The quick brown fox jumps over the lazy dog

DILM

<input checked="" type="checkbox"/> Probabilistically Principled	<input checked="" type="checkbox"/> Multi-token
<input checked="" type="checkbox"/> Flexible generation order	<input checked="" type="checkbox"/> Arbitrary length

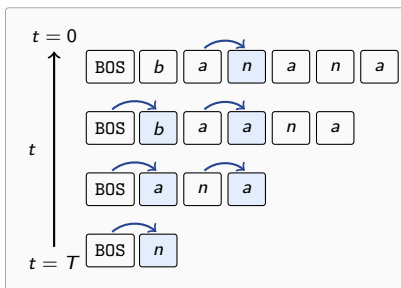
The quick fox over the dog
The quick fox jumps over the dog
The quick brown fox jumps over the lazy dog

A CTMC framework for Insertion Language Models



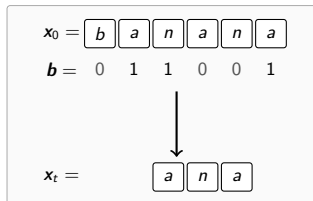
State space: $\mathbb{X} = \bigcup_{n=0}^{\ell_{\max}} \{n\} \times \mathbb{V}^n$, so a state contains both length and tokens.

Forward/noising: $\mathbf{x}_0 \rightarrow \mathbf{x}_t$, $t \in [0, 1]$, where \mathbf{x}_t is a random subsequence of \mathbf{x}_0 obtained by deleting tokens.



Reverse/generative: $\hat{R}_t^\theta(\mathbf{x}, \mathbf{y})$ directly parameterizes insertion rates, giving a CTMC that inserts tokens in available gaps. This yields a diffusion-style denoising objective from the CTMC path likelihood. **Key benefit:** sampling quality can trade off with the number of reverse-time steps.

Noising: Random Deletion



The noising process deletes tokens from x_0 to form a subsequence x_t .

1. Sample $t \sim \text{Uniform}[0, T]$.
2. **Joint:** $d \sim \text{Poisson}(\bar{\sigma}_{0,t})$.
3. **Independent:**
 $d \sim \text{Binomial}(n, 1 - \exp(-\bar{\sigma}_{0,t}))$.
4. Set $d \leftarrow \min(d, n)$.
5. Uniformly at random select d indices from $\{1, \dots, n\}$ to drop; let $\mathbf{b} \in \mathcal{B}_{n, n-d}$ enumerate the complementary kept indices.
6. Return $((n-d, \dot{x}[\mathbf{b}]), t)$.

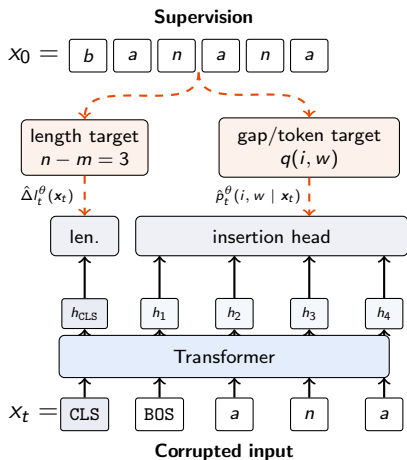
Joint noising

$$\begin{aligned}
 & p_{t|s}^{\text{joint}}((m, \dot{y}) | (n, \dot{x})) \\
 &= \begin{cases} \frac{\exp(-\bar{\sigma}_{s,t})(\bar{\sigma}_{s,t})^{n-m} m!}{n!} \sum_{\mathbf{b} \in \mathcal{B}_{n,m}} \delta_{\dot{x}[\mathbf{b}]}(\dot{y}), & 0 < m \leq n, \\ 1 - \sum_{r=1}^n \frac{\exp(-\bar{\sigma}_{s,t})(\bar{\sigma}_{s,t})^{n-r} r!}{n!}, & m = 0. \end{cases}
 \end{aligned}$$

Independent noising

$$\begin{aligned}
 & p_{t|s}^{\text{ind}}((m, \dot{y}) | (n, \dot{x})) \\
 &= \rho_{s,t}^m (1 - \rho_{s,t})^{n-m} \sum_{\mathbf{b} \in \mathcal{B}_{n,m}} \delta_{\dot{x}[\mathbf{b}]}(\dot{y}), \quad \rho_{s,t} = \exp(-\bar{\sigma}_{s,t}).
 \end{aligned}$$

Denoising: Joint Denoising with single-token Insertion (DILM-S)



DILM-S

Supervision. Delete tokens from x_0 using mask \mathbf{b} , giving $x_t = x_0[\mathbf{b}]$. The deleted tokens define the length target $n - m$ and the reverse insertion target

$$q(i, w) = \frac{\sum_{k \in S_i(\mathbf{b})} \delta_w(x_0^k)}{n - m}.$$

Predictions. The transformer encodes $\tilde{x}_t = (\text{BOS}, x_t^1, \dots, x_t^m)$. The length head predicts $\hat{\Delta}_t^\theta(x_t) = \mathbb{E}[l | x_t]$. The insertion head predicts one joint softmax

$$\hat{p}_t^\theta(i, w | x_t) \propto \exp s^\theta(i, w | x_t).$$

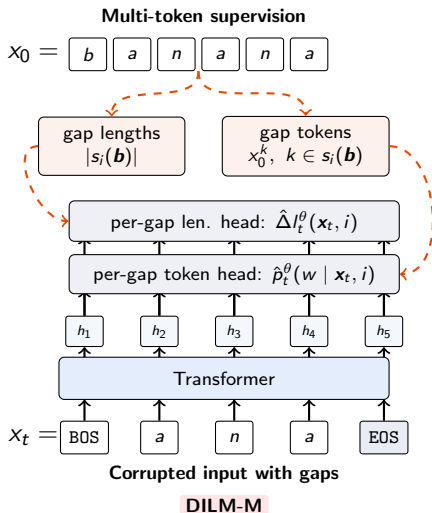
Denosing: Independent Denoising with multi-token Insertion (DILM-M)

Supervision. Independent deletion noise gives a corrupted subsequence $x_t = x_0[\mathbf{b}]$. For each gap i , the deleted indices $s_i(\mathbf{b})$ define a gap-length target $|s_i(\mathbf{b})|$ and token targets x_0^k , $k \in s_i(\mathbf{b})$.

$$x_0^k, \quad k \in s_i(\mathbf{b}).$$

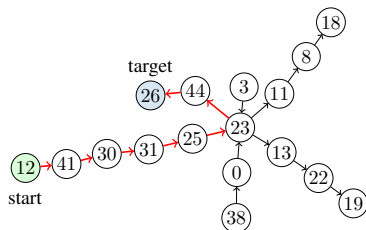
Predictions. The transformer gives one representation per gap. The length head predicts $\hat{\Delta} l_t^\theta(x_t, i) = \mathbb{E}[l \mid x_t, i]$. The token head predicts a separate vocabulary softmax for each gap,

$$\hat{p}_t^\theta(w \mid x_t, i) \propto \exp s^\theta(w \mid x_t, i).$$



Empirical Evaluation: Planning Task

12 41 13 22 ... 11 8 12 26 <s> 12 41 ... 44 26



Given the edges and the start and the target node, the goal is to predict the path from the start to the target node.

Model	Easy	Med.	Hard
ARM	32.3	75.0	23.0
MDM	100.0	36.5	21.0
ILM	100.0	100.0	99.1
DILM-S	100.0	100.0	99.3
DILM-M	100.0	98.6	95.1

ARM: Due to the presence of easy edges, a left-to-right learner like ARM cannot predict the correct edge at the junction

MDM: An arbitrary order learner using absolute positions needs to predict the length of the path to work backwards

DILM: Insertion based models can overcome both these limitations

Empirical Evaluation: Language Modeling

Setup: 12-layer Transformer with 768 hidden dimensions and 12 attention heads, trained from scratch with AdamW on 8 A100 GPUs

LM1B: short news sequences tokenized with the BERT tokenizer and padded to length 128

OpenWebText: longer web-text sequences tokenized with the GPT-2 tokenizer and padded to length 1024

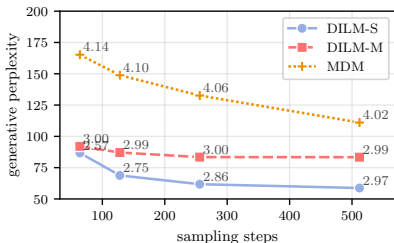
Evaluation: per-token NLL of 1000 unconditional samples under a pre-trained LM

Data	Model	NLL ↓
LM1B	ILM	4.67
	DILM-S	<u>4.65</u>
	DILM-M	4.76
OWT	MDM	3.96
	ILM	4.65
	DILM-S	3.74
	DILM-M	<u>3.92</u>

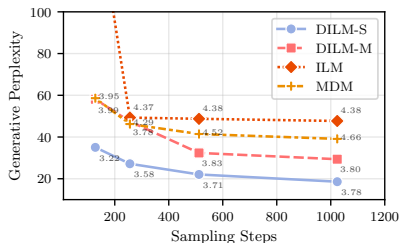
DILM-S and DILM-M improve sample quality as reverse steps increase, giving a smooth quality/-compute trade-off.

Empirical Evaluation: Sampling Trade-off

LM1B



OpenWebText



Generative perplexity improves as the number of sampling steps increases; DILM-S and DILM-M exhibit smoother quality/compute trade-offs than non-diffusion Insertion Language Model (ILM).

CTMC framework: The CTMC framework gives a principled diffusion-style training and sampling methods for insertion language models.

Using this framework, we derive two types of noising and respective generative processes: DILM-S and DILM-M.

Sampling flexibility: We gain sampling flexibility by using the CTMC framework to control the number of sampling steps while keeping the benefits of insertion-based generation.

Future work: Adapting pre-trained autoregressive language models into diffusion-based insertion language models remains an active direction for future work.



Thank you