A Practical Guide to

# Diffusion Language Models

Himanshu Sahni
06/08/2025

# 1. Why?

# Language Models

The

quick

brown

fox

jumped

over

the

lazy

# Language Models

The

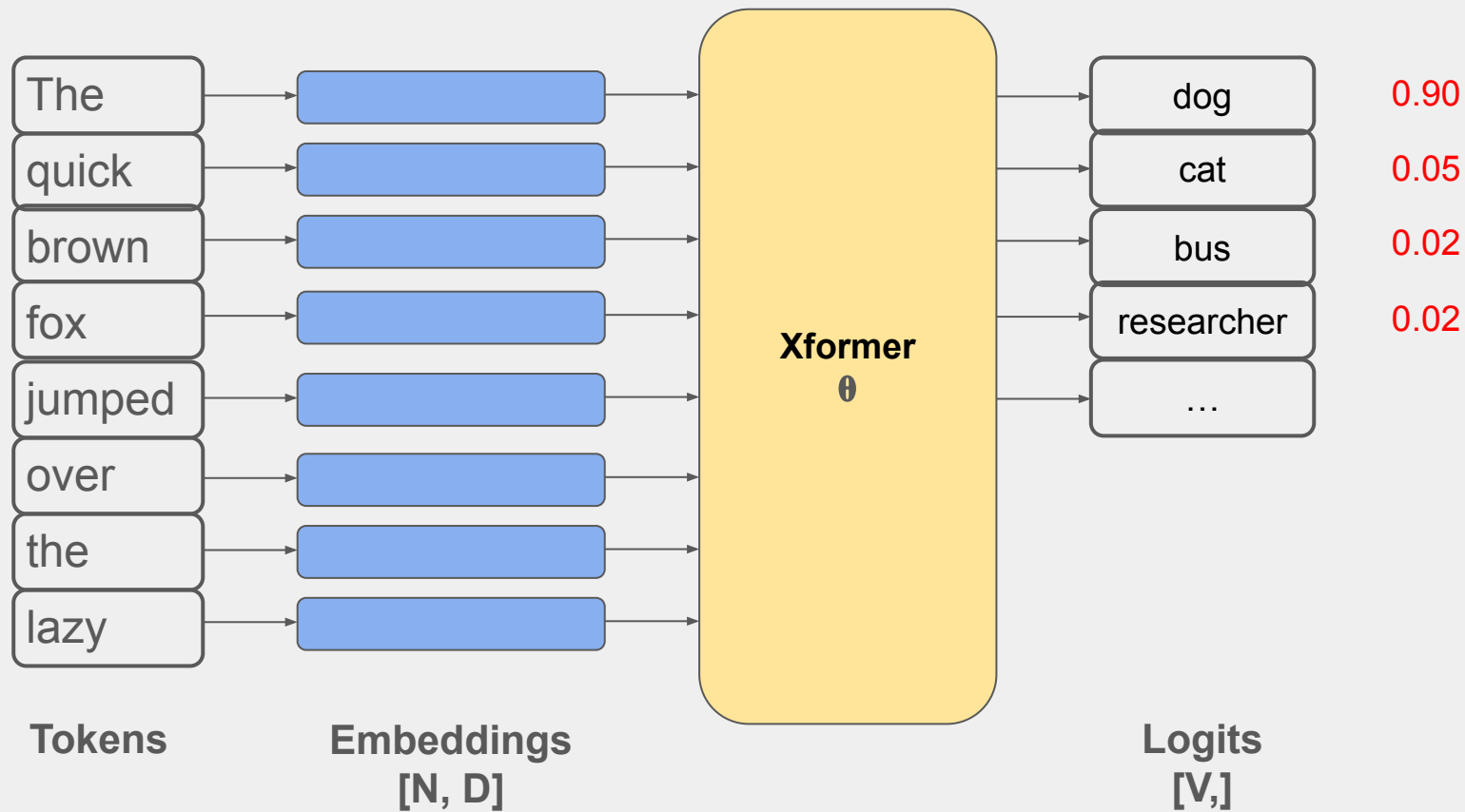quick

brown

fox

jumped

over
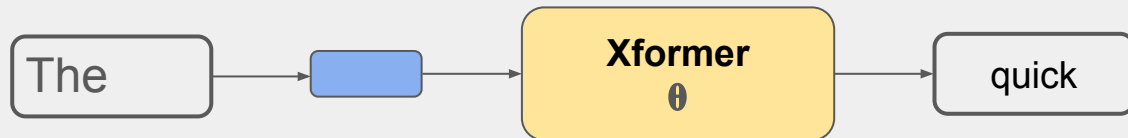
the

lazy

**Tokens**

# Language Models

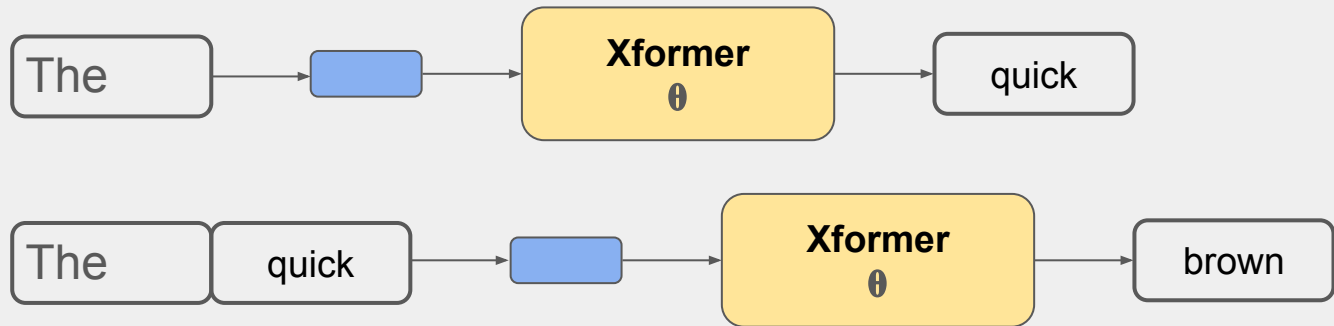| | |
|---|---|
| The | ▬ |
| quick | ▬ |
| brown | ▬ |
| fox | ▬ |
| jumped | ▬ |
| over | ▬ |
| the | ▬ |
| lazy | ▬ |

**Tokens**

**Embeddings
[N, D]**

# Language Models

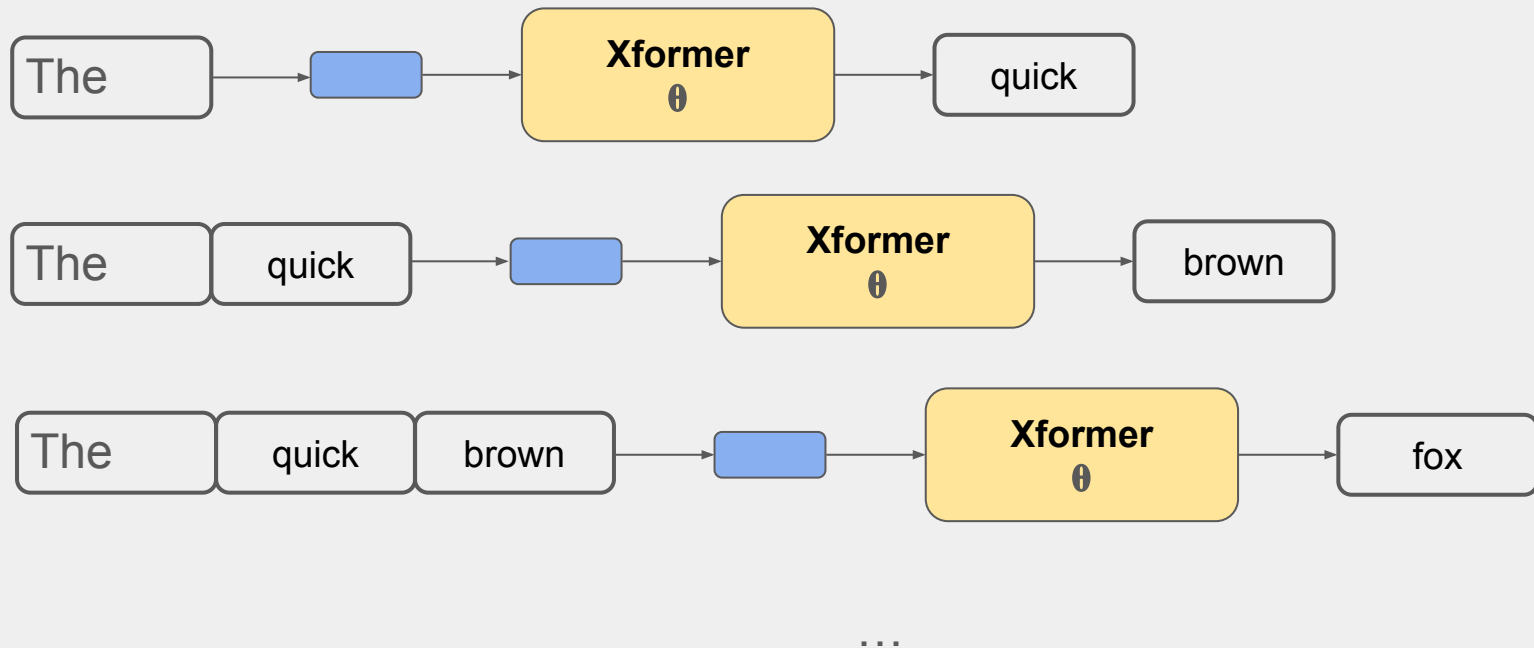# Why Diffusion Language Models?

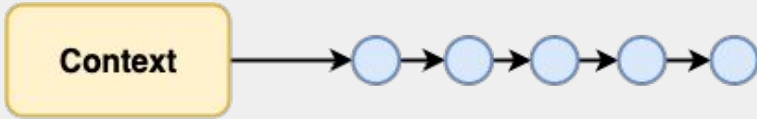# Why Diffusion Language Models?

# Why Diffusion Language Models?

# Diffusion vs. Autoregressive Language Models



1 token at a time :(

n tokens at a time :)

# Diffusion vs. Autoregressive Language Models

|                   | Autoregression | Diffusion |
|-------------------|:--------------:|:---------:|
| High Quality      | ✓              | ✓         |
| Arbitrary Length  | ✓              | ✓         |
| KV Caching        | ✓              | ✓         |
| Parallel Decoding | ✗              | ✓         |

# Diffusion vs. Autoregressive Language Models

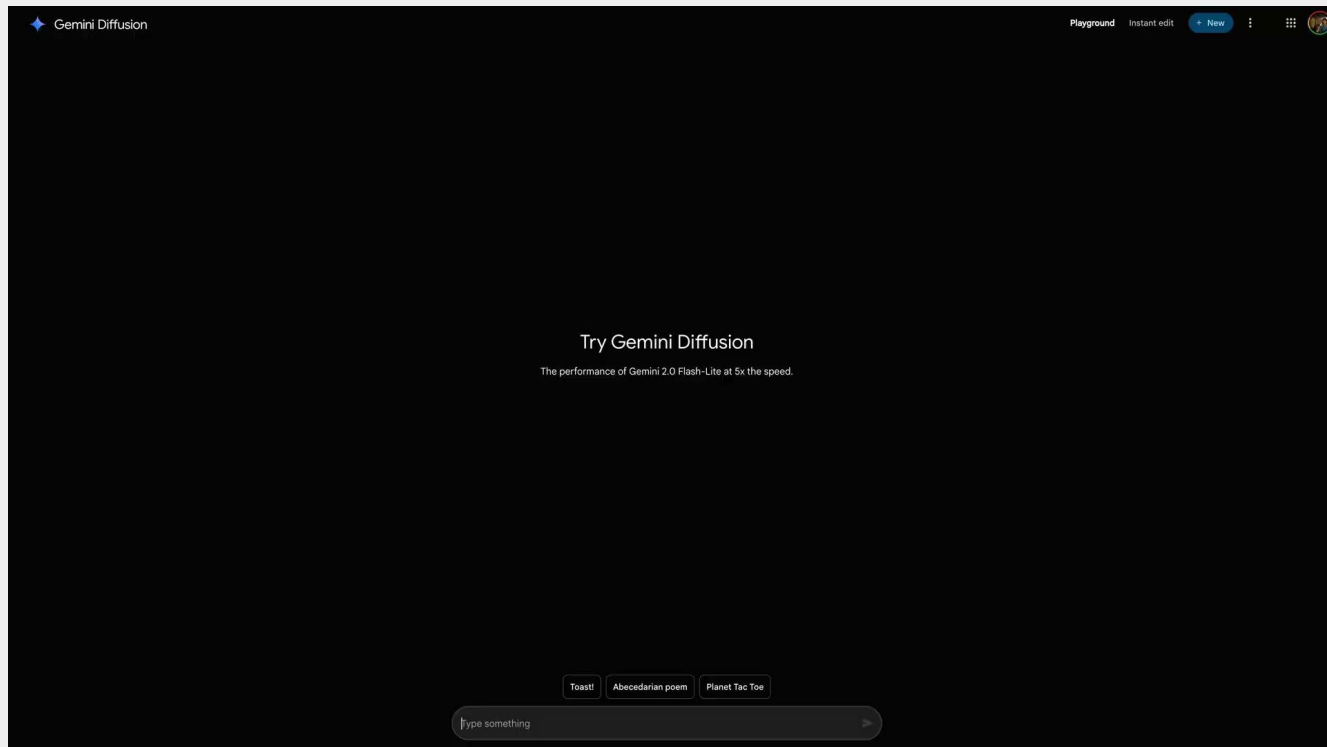|  | Autoregression | Diffusion |
|---|:---:|:---:|
| High Quality | ✓ | ✓ |
| Arbitrary Length | ✓ | ✓ |
| KV Caching | ✓ | ✓ |
| Parallel Decoding | ✗ | ✓ |
| Self-correction | ✗ | ✓ |
| Unify multimodality | ✗ | ✓ |

# 2. The hype

# The hype



# Gemini Diffusion

Our state-of-the-art, experimental text diffusion model

Join the waitlist ›

# The hype

# The hype



**Gemini Diffusion** · Playground · Instant edit · + New

## Try Gemini Diffusion

The performance of Gemini 2.0 Flash-Lite at 5x the speed.

Drawing · Fractal · Toast!

Type something

# The hype

## Benchmarks

Gemini Diffusion's external benchmark performance is comparable to much larger models, whilst also being faster.

| Benchmark | GEMINI DIFFUSION | GEMINI 2.0 FLASH-LITE |
|---|---|---|
| Code<br>LiveCodeBench (v6) | 30.9% | 28.5% |
| Code<br>BigCodeBench | 45.4% | 45.8% |
| Code<br>LBPP (v2) | 56.8% | 56.0% |
| Code<br>SWE-Bench Verified* | 22.9% | 28.5% |
| Code<br>HumanEval | 89.6% | 90.2% |
| Code<br>MBPP | 76.0% | 75.8% |
| Science<br>GPQA Diamond | 40.4% | 56.5% |
| Mathematics<br>AIME 2025 | 23.3% | 20.0% |
| Reasoning<br>BIG-Bench Extra Hard | 15.0% | 21.0% |
| Multilingual<br>Global MMLU (Lite) | 69.1% | 79.0% |

# The hype

# The hype

Nie et al., "Large Language Diffusion Models," arXiv:2502.09992, 2025.

# 3. How?

# 3. How?
# (Part 1)

# Image Diffusion



Forward Process

Reverse Process

# Latent Diffusion



[H,W,C]

VAE

Latent $\mathbf{x}_o$
[D,]

+

**Noise** $\sim \mathcal{N}(0, \boldsymbol{\varepsilon}_t)$
[D,]

Noised Latent $\mathbf{x}_t$
[D,]

$$x_t = \sqrt{\bar{\alpha}_t}\, x_0 + \sqrt{1 - \bar{\alpha}_t}\, \varepsilon_t$$

# Latent Diffusion



Noised Latent $\mathbf{x}_t$
[D,]

**Xformer**
$\boldsymbol{\theta}$

$\boldsymbol{\varepsilon_\theta}(x_t, t)$

# Latent Diffusion

**Network prediction**

**Actual Noise**

Train 🔥

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{x_0, \varepsilon, t} \left[ \left\| \varepsilon_\theta(x_t, t) - \varepsilon \right\|^2 \right]$$
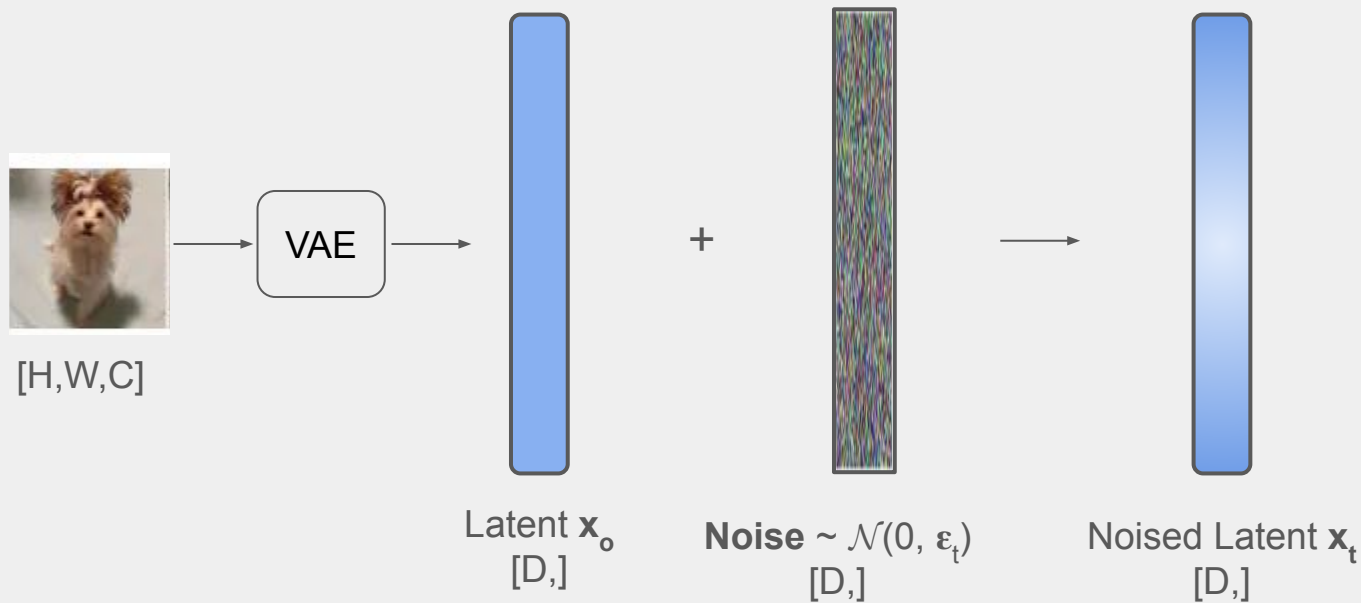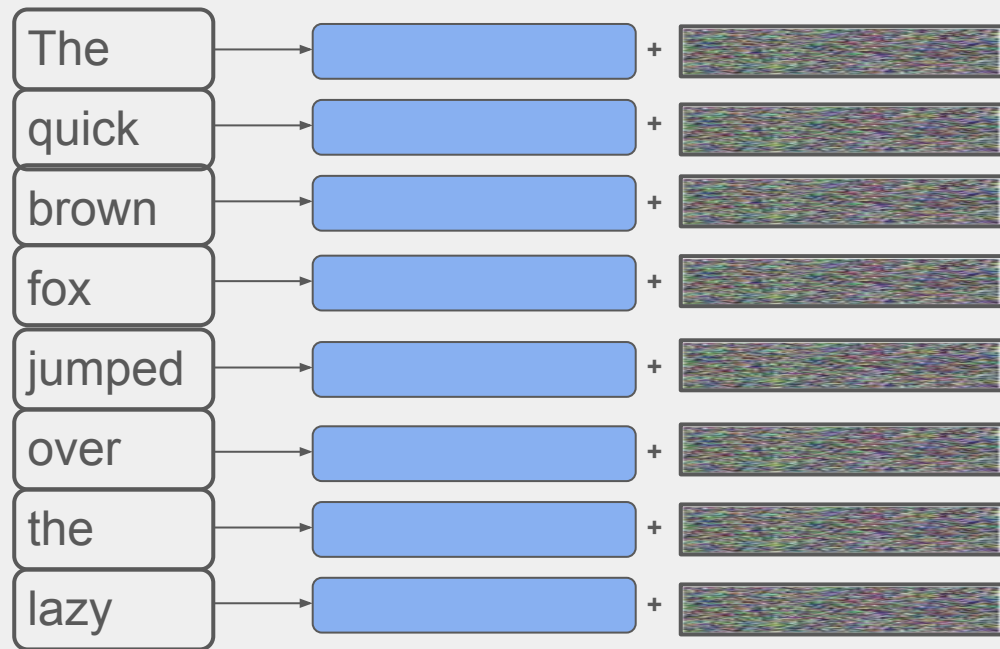
Sample 🌊

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_\theta(x_t, t) \right) + \beta_t \varepsilon$$

# Latent Diffusion



VAE

[H,W,C]

Latent $\mathbf{x_o}$
[D,]

**Noise** $\sim \mathcal{N}(0, \varepsilon_t)$
[D,]

Noised Latent $\mathbf{x_t}$
[D,]

$$x_t = \sqrt{\bar{\alpha}_t}\, x_0 + \sqrt{1 - \bar{\alpha}_t}\, \varepsilon_t$$

# Latent Diffusion Language Models

| Tokens | | Embeddings | Noise |
|--------|--|-----------|-------|
| The | → | | + |
| quick | → | | + |
| brown | → | | + |
| fox | → | | + |
| jumped | → | | + |
| over | → | | + |
| the | → | | + |
| lazy | → | | + |

**Tokens**

**Embeddings
[N, D]**

**Noise
[N, D]**

# Latent Diffusion Language Models

| | | | |
|---|---|---|---|
| The | ⬜ | + | 🟫 |
| quick | ⬜ | + | 🟫 |
| brown | ⬜ | + | 🟫 |
| fox | ⬜ | + | 🟫 |
| jumped | ⬜ | + | 🟫 |
| over | ⬜ | + | 🟫 |
| the | ⬜ | + | 🟫 |
| lazy | ⬜ | + | 🟫 |

**Xformer**
θ

**Tokens**  **Embeddings**  **Noise**
**[N, D]**  **[N, D]**

# Latent Diffusion Language Models
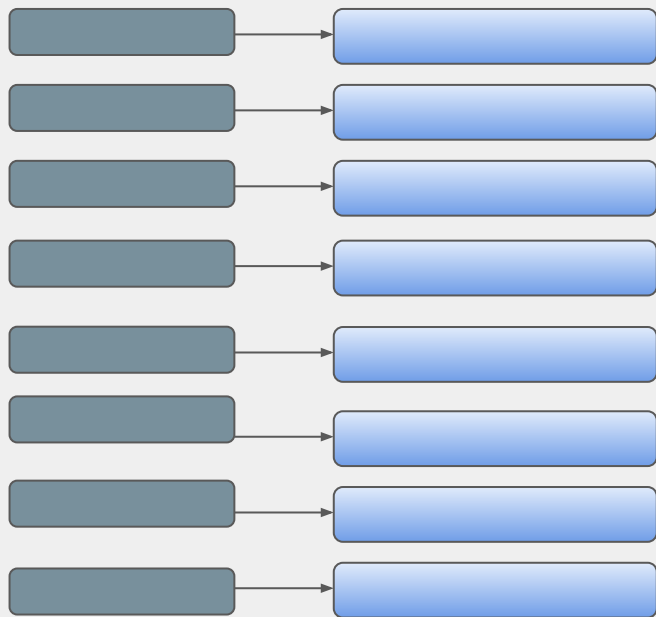
# Latent Diffusion Language Models



**Logits
[N,V]**

For each logit, compute:

$$\mathbb{E}_{p(x_0|x,t)}[x_0] = \sum_{i=1}^{V} p(x_0 = e_i|x,t) \cdot e_i$$
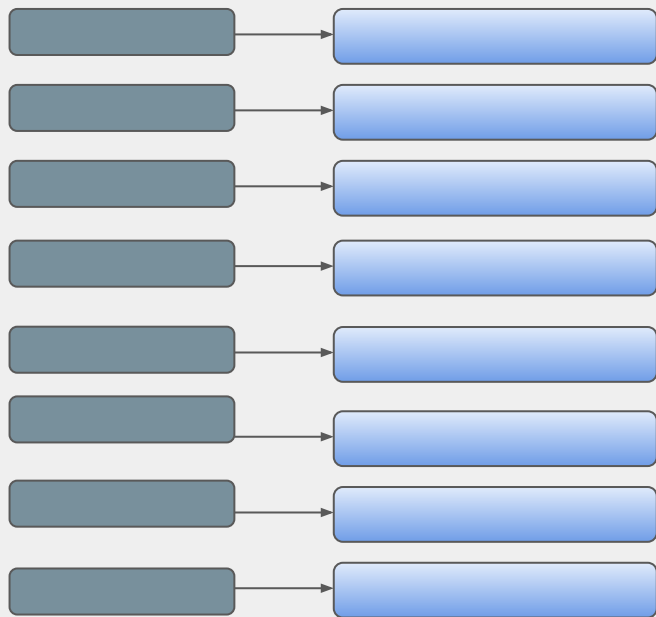
# Latent Diffusion Language Models

$$\mathbb{E}_{p(x_0|x,t)}[x_0] = \sum_{i=1}^{V} p(x_0 = e_i|x,t) \cdot e_i$$

Let's say we're good at predicting these

**Logits [N,V]**

**Predicted Embeddings [N, D]**

# Latent Diffusion Language Models
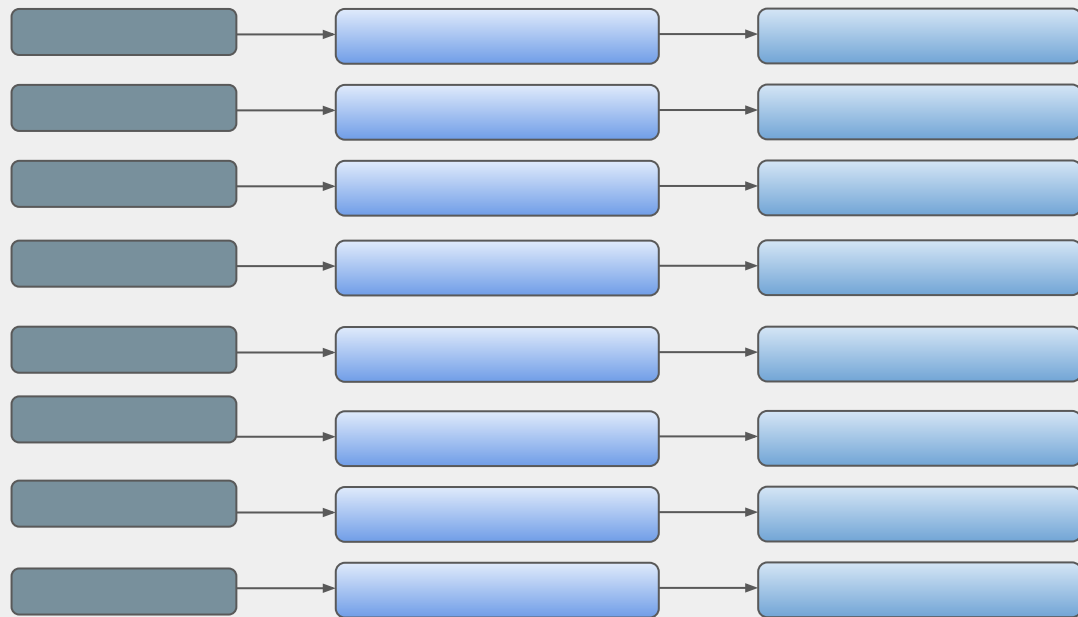


**Logits
[N,V]**

**Predicted
Embeddings
[N, D]**

$$\nabla_x \log p_t(x) = s(x, t)$$

$$s(x, t) \approx s(x, t \mid x_0) = \frac{1}{t^2}(x_0 - x)$$

$$\hat{s}(\mathbf{x}, t) = \frac{\mathbb{E}_{p(x_0 \mid \mathbf{x}, t)}[\mathbf{x}_0] - \mathbf{x}}{t^2}$$

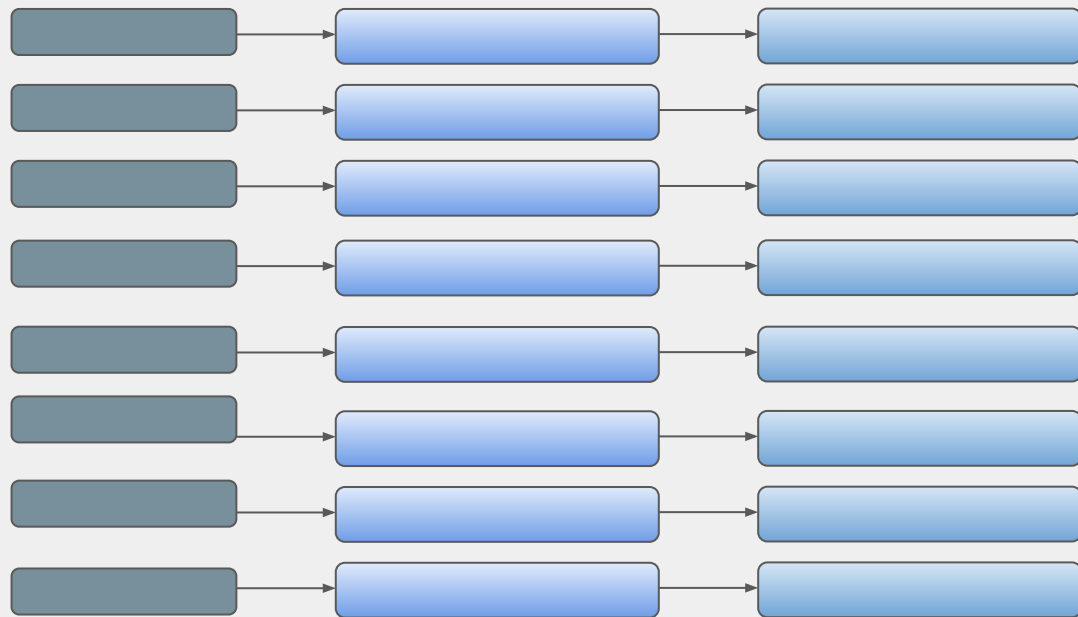Dieleman et al., "Continuous Diffusion for Categorical Data," arXiv:2211.15089, 2022.

# Latent Diffusion Language Models

$$\hat{\mathbf{s}}(\mathbf{x}, t) = \frac{\mathbb{E}_{p(x_0|\mathbf{x},t)}[\mathbf{x}_0] - \mathbf{x}}{t^2}$$



**Logits
[N,V]**

**Predicted
Embeddings
[N, D]**

$\nabla_x \log p_t(\mathbf{x})$
**[N, D]**

# Latent Diffusion Language Models



**Logits
[N,V]**

**Predicted
Embeddings
[N, D]**

$\nabla_x \log p_t(\mathbf{x})$
**[N, D]**

$$x_{t-\Delta t} = x_t + t \cdot \nabla_x \log p_t(x) \cdot \Delta t$$

T. Karras et al. "Elucidating the design space of diffusion-based generative models". arXiv:2206.00364, 2022.

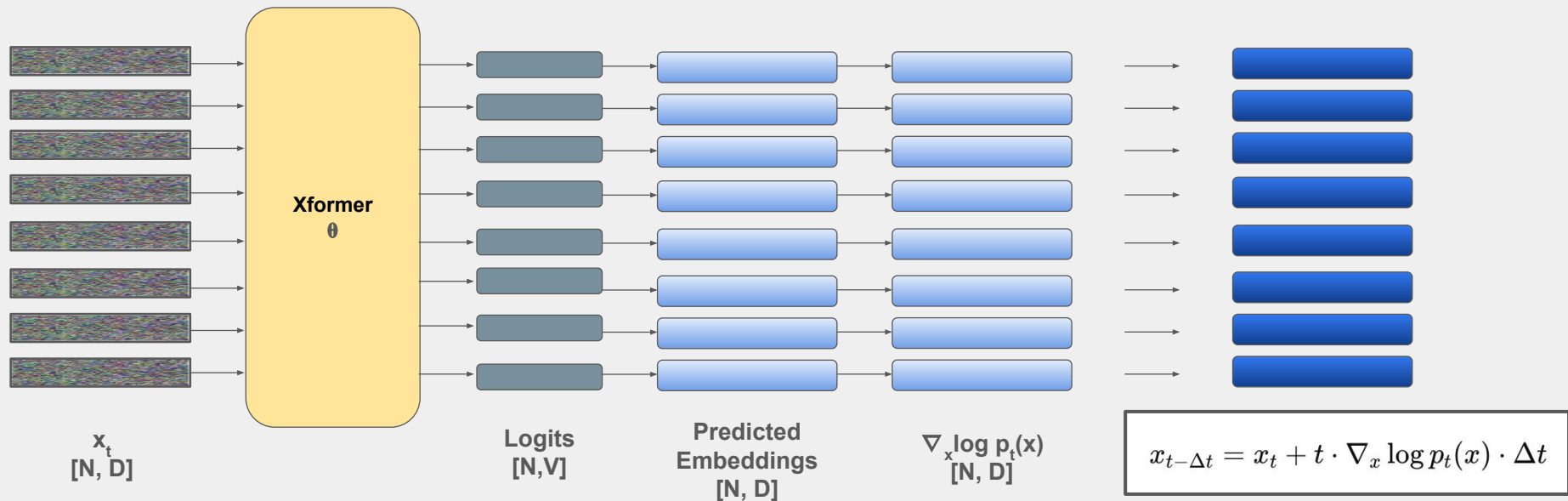# Latent Diffusion Language Models - Sampling



$x_t$
[N, D]

Xformer
θ

Logits
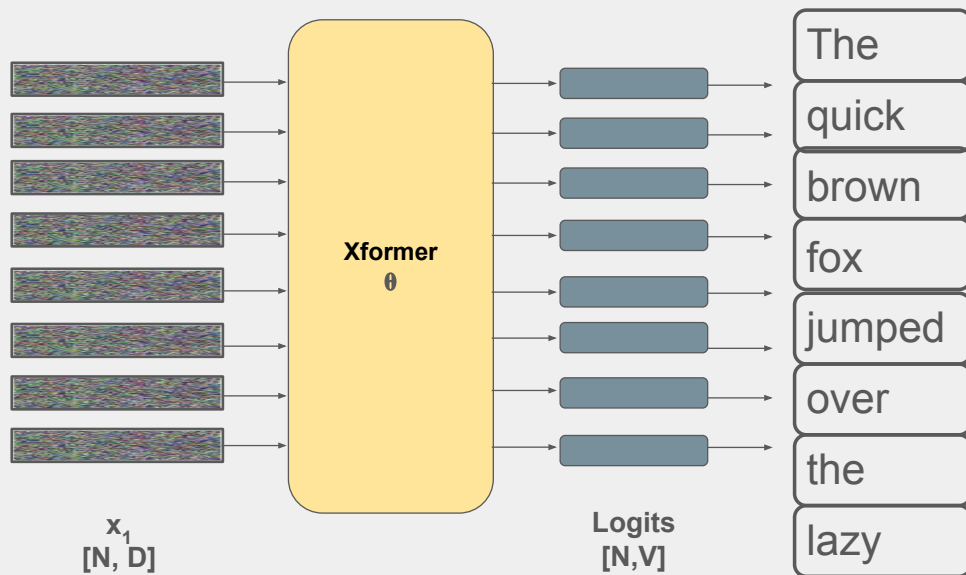[N,V]

Predicted
Embeddings
[N, D]

$\nabla_x \log p_t(x)$
[N, D]

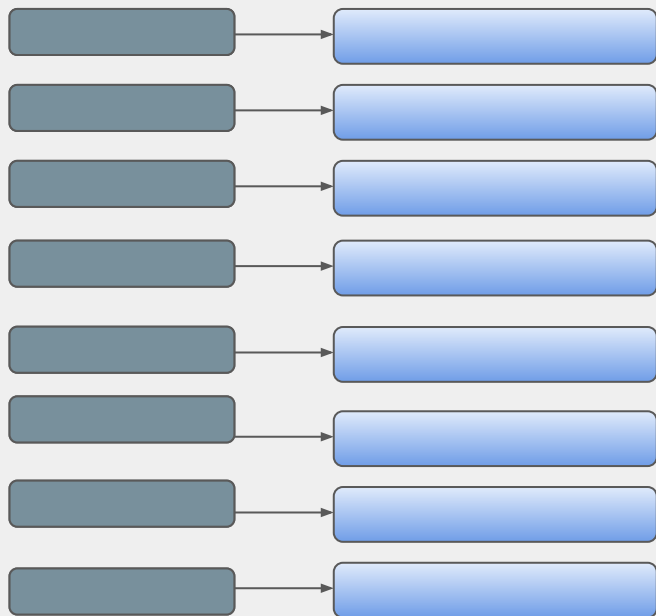$$x_{t-\Delta t} = x_t + t \cdot \nabla_x \log p_t(x) \cdot \Delta t$$

# Latent Diffusion Language Models - Sampling

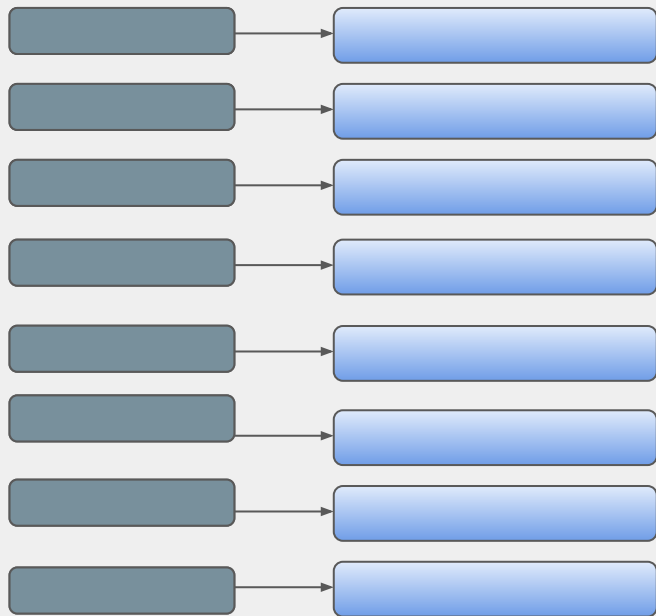# Latent Diffusion Language Models - Training



Logits
[N,V]

Predicted
Embeddings
[N, D]

How do we get good at predicting $x_0$?

# Latent Diffusion Language Models - Training



**Logits
[N,V]**

**Predicted
Embeddings
[N, D]**

$$\left\| \mathbb{E}_{p_\theta(x_0|x,t)}[x_0] - x_0 \right\|^2$$

# Latent Diffusion Language Models - Training

$$\mathcal{L}_{\mathrm{mse}} = \left\| \sum_i p_\theta(i \mid x, t) \cdot e_i - e_{x_0} \right\|^2$$

Dieleman et al., "Continuous Diffusion for Categorical Data," arXiv:2211.15089, 2022.

# Latent Diffusion Language Models - Training

Turns out, we can also train the embeddings!

$$\mathcal{L}_{\mathrm{mse}} = \left\| \sum_i p_\theta(i \mid x, t) \cdot e_i - e_{x_0} \right\|^2$$

Dieleman et al., "Continuous Diffusion for Categorical Data," arXiv:2211.15089, 2022.

# Latent Diffusion Language Models - Training

Turns out, we can also train the embeddings!

$$\mathcal{L}_{\mathrm{mse}} = \left\| \sum_i p_\theta(i \mid x, t) \cdot e_i - e_{x_0} \right\|^2$$

Not so fast!

Dieleman et al., "Continuous Diffusion for Categorical Data," arXiv:2211.15089, 2022.

# Latent Diffusion Language Models - Training

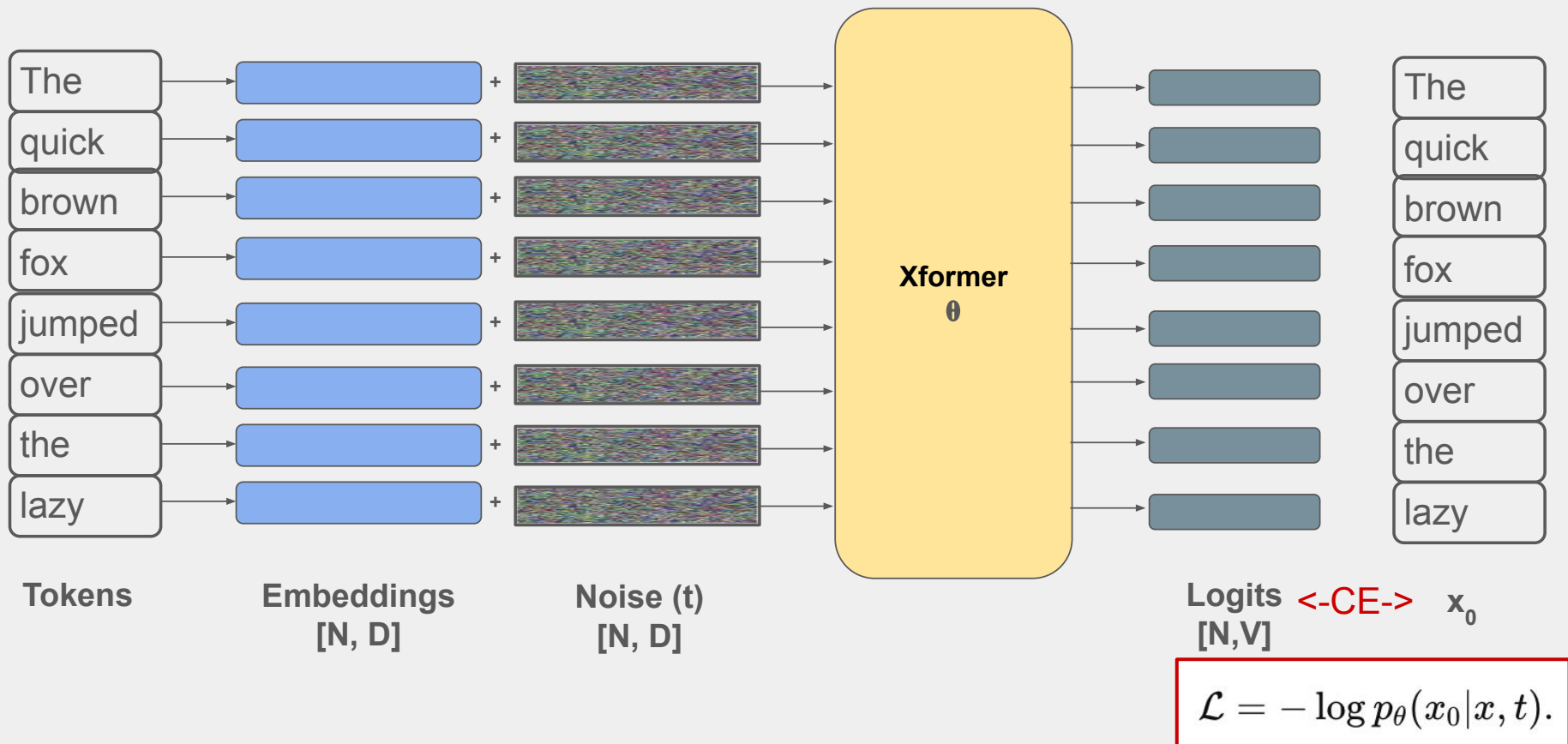Turns out, we can also train the embeddings!

$$\mathcal{L}_{\mathrm{mse}} = \left\| \sum_i p_\theta(i \mid x, t) \cdot e_i - e_{x_0} \right\|^2$$
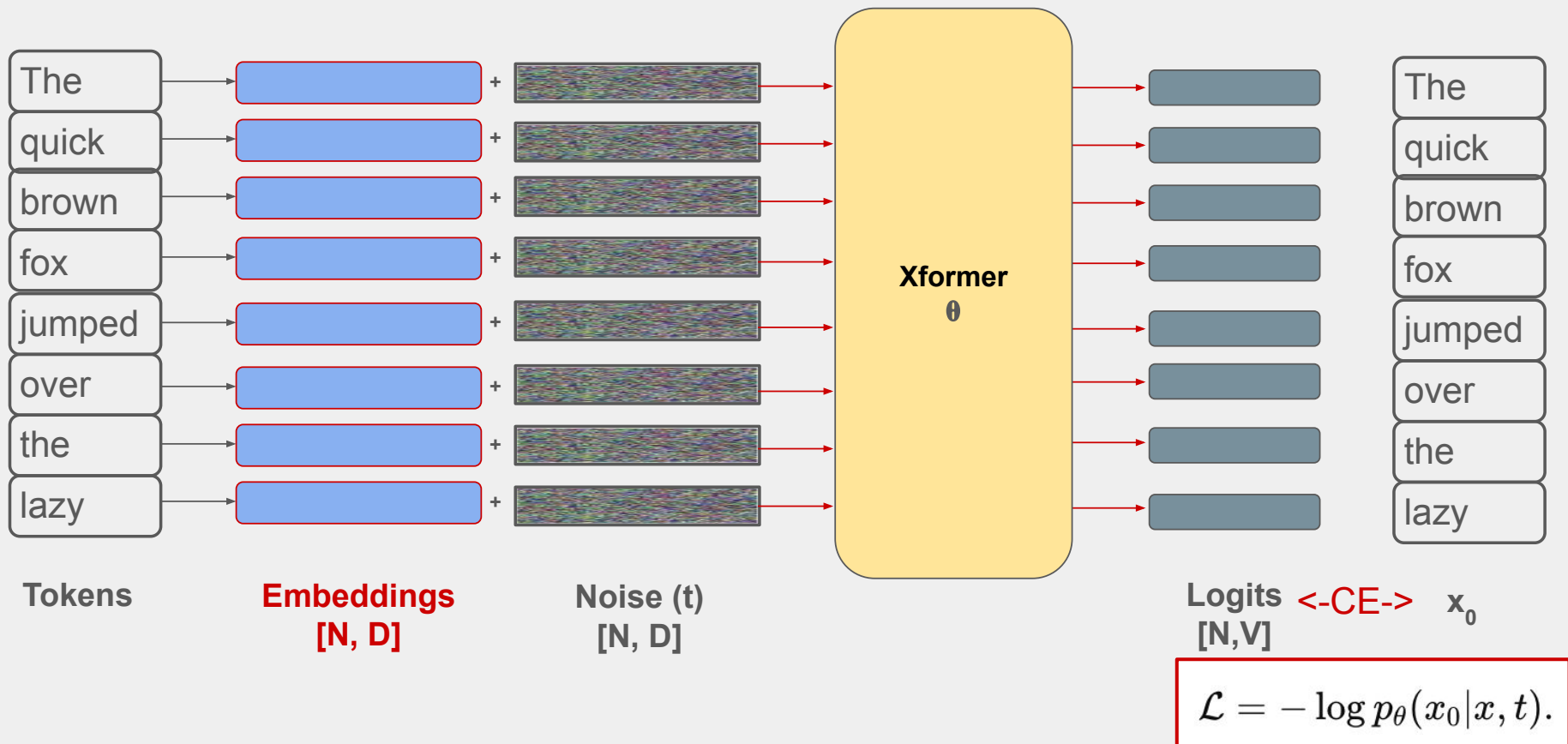
Not so fast!

The model can just learn to make all the embeddings

the same!

Dieleman et al., "Continuous Diffusion for Categorical Data," arXiv:2211.15089, 2022.

# Latent Diffusion Language Models - Training

Turns out, we can also train the embeddings!

**Not so fast!**

The model can just learn to make all the embeddings the same!

Also, MSE is not a good loss for MSE data.

$$\mathcal{L}_{\mathrm{mse}} = \left\| \sum_i p_\theta(i \mid x, t) \cdot e_i - e_{x_0} \right\|^2$$

Dieleman et al., "Continuous Diffusion for Categorical Data," arXiv:2211.15089, 2022.

# Latent Diffusion Language Models - Training



| Tokens | Embeddings [N, D] | Noise (t) [N, D] | Xformer θ | Logits [N,V] <-CE-> $x_0$ |
|---|---|---|---|---|

$$\mathcal{L} = -\log p_\theta(x_0 | x, t).$$

# Latent Diffusion Language Models - Training



**Tokens**

**Embeddings [N, D]**

**Noise (t) [N, D]**

**Xformer θ**

**Logits [N,V]** <-CE-> $x_0$

$$\mathcal{L} = -\log p_\theta(x_0|x, t).$$

# Latent Diffusion Language Models - Training

Turns out, we can also train the embeddings!

Not so fast!

# Latent Diffusion Language Models - Training

Turns out, we can also train the embeddings!

Not so fast!

The model can just learn to push embeddings to the

extreme to make them easy to predict from noise!!*

# 3. How?
# (Part Deux)

# Latent Diffusion Language Models - Training

The model can just learn to push embeddings to the extreme to make them easy to predict from noise!!

- At scale, ^^ turns out to be a huge problem!
- The higher dimensional latent you use, the bigger the problem!
- How to learn embeddings suited for diffusion - high frequency to low frequency with gaussian noise?

ooof!

# Discrete Diffusion Language Models

| Continuous noising process | Discrete noising process |
|---|---|
| Embed tokens to latent space | Operate on tokens directly |
| Gaussian Noise added to embeddings | Flip tokens to other tokens in the vocabulary or MASK token |
| Uncertainty maintained until last step | Commit to tokens at each step |
| Difficulty scaling beyond 1B | Straightforward to scale to 8B |

# Discrete noising process

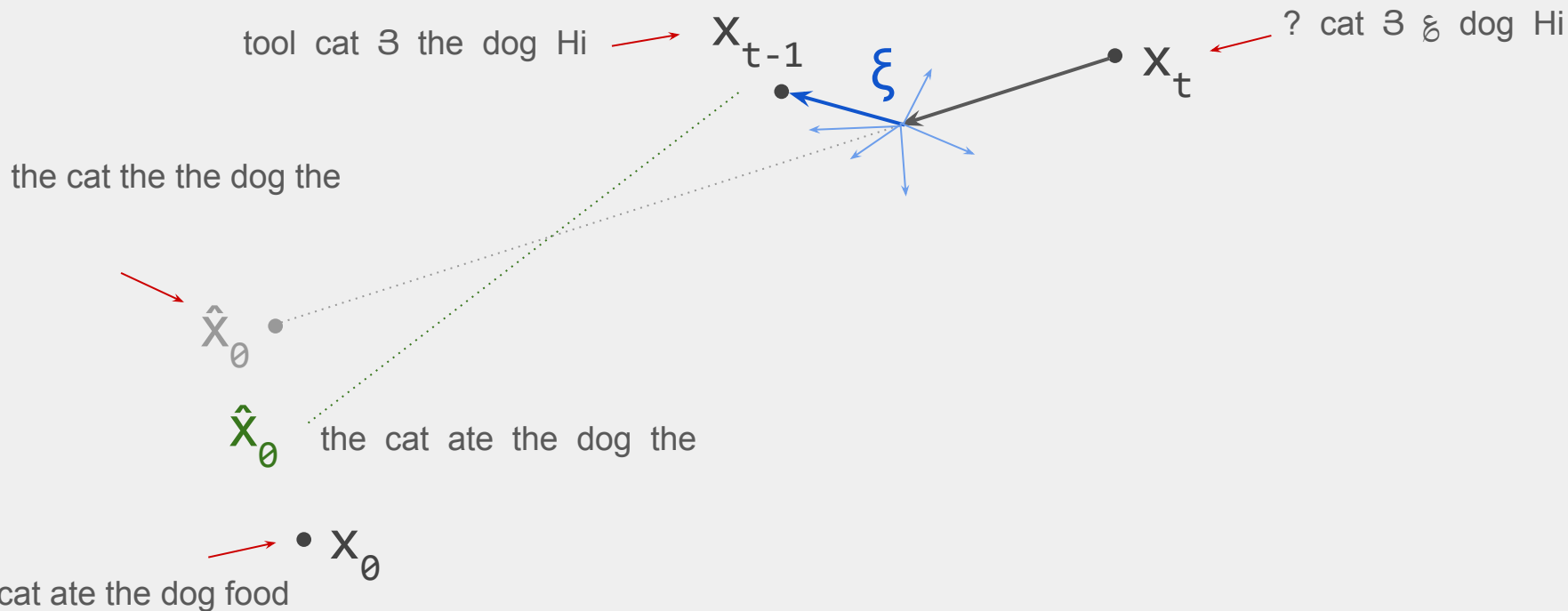| | | | | | | |
|---|---|---|---|---|---|---|
| The | cat | ate | the | dog | food | t = 0 |
| The | cat | ate | ৪ | dog | food | t = 0.1 |
| The | 吹 | ate | ৪ | dog | Hi | t = 0.2 |
| | | | … | | | |
| ? | 吹 | 3 | ৪ | ফ | Hi | t = 1.0 |

Forward Process

# Discrete noising process

The  [M]  ate  [M]  dog  [M]  $\xrightarrow{\quad p_\theta(x_t,\ t)\quad}$  The  cat  ate  the  dog  food

$$-\mathbb{E}_{t,p_0,r_0,r_t}\left[\frac{1}{t}\sum_{i=1}^{L'}\mathbf{1}[r_t^i=\mathbf{M}]\log p_\theta(r_0^i|p_0,r_t)\right],$$

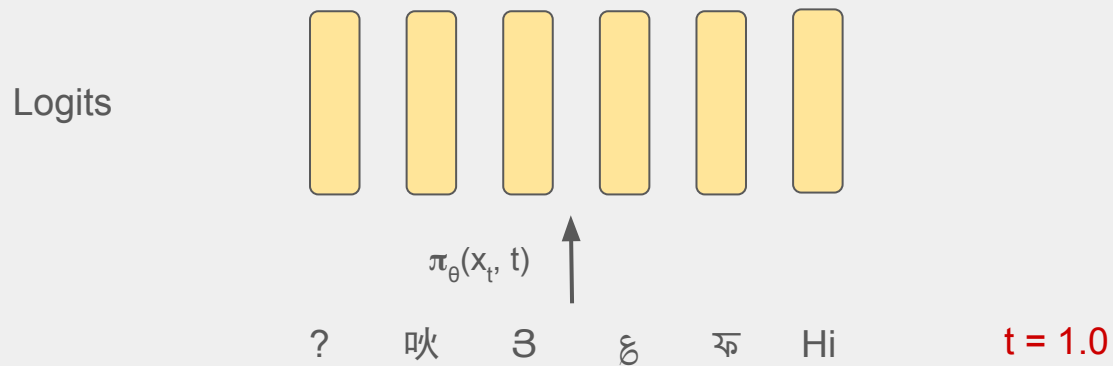Nie et al., "Large Language Diffusion Models," arXiv:2502.09992, 2025.

# Discrete diffusion for text - Inference



tool  cat  Ƨ  the  dog  Hi  →  $X_{t-1}$

? cat  Ƨ  ℅  dog  Hi  ←

$\xi$

$X_t$

the cat the the dog the

$\hat{X}_0$

the cat ate the dog food  →  $X_0$

# Discrete diffusion for text - Inference



tool cat Ɛ the dog Hi → $X_{t-1}$

? cat Ɛ ୫ dog Hi ←

$\xi$

$X_t$

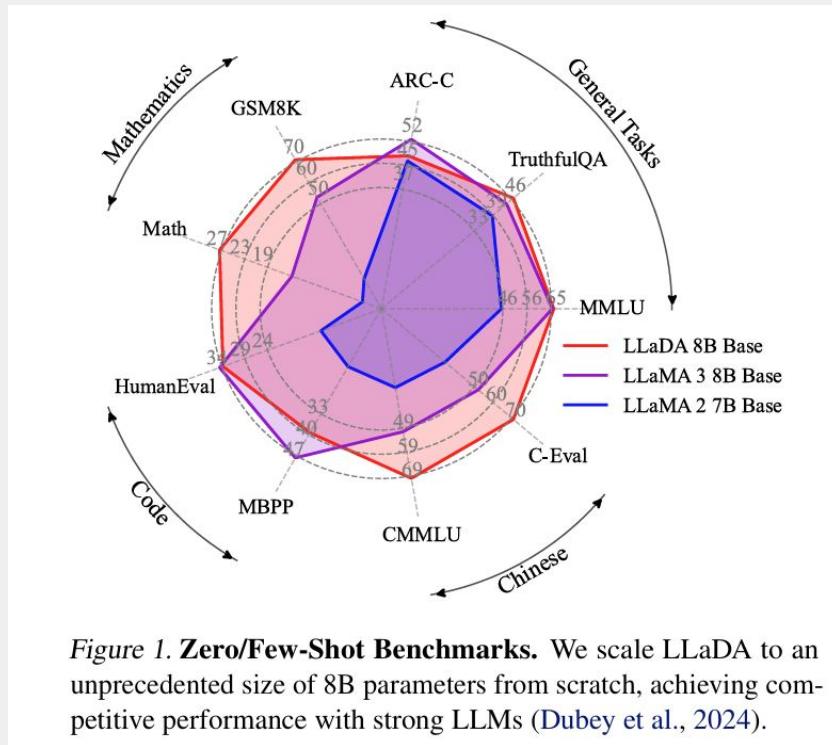the cat the the dog the

$\hat{X}_0$

$\hat{X}_0$ the cat ate the dog the

$X_0$

the cat ate the dog food

# Discrete diffusion for text - Sampler

Logits

$\pi_\theta(x_t, t)$

? 吹 Ʒ ৬ ফ Hi

t = 1.0

# Discrete diffusion for text - Sampler

# Discrete diffusion for text - Sampler
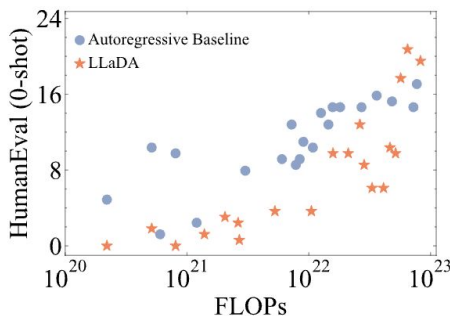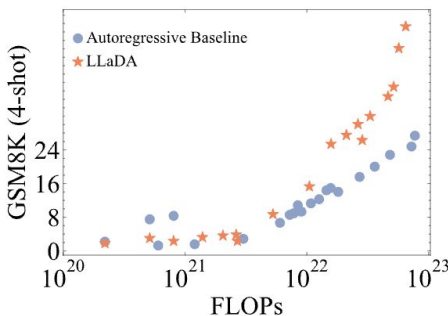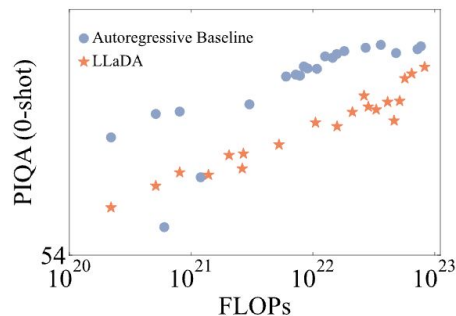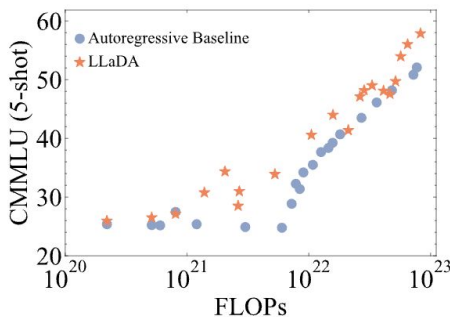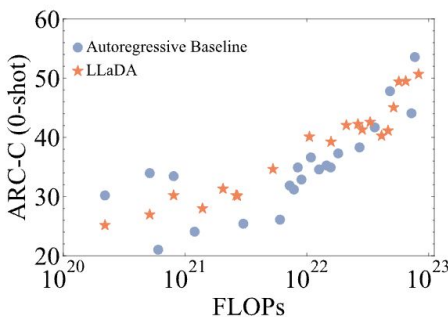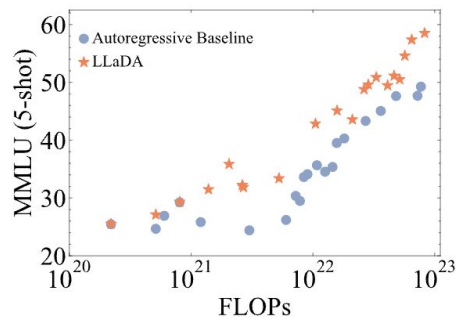
# Discrete diffusion for text - Sampler

Many design choices and hypers

- Sampling temperature - $\mathbf{T}(t_n)$

- Acceptance probability profile - $\boldsymbol{p_a}(t_n, t_{n-1})$

- Renoising probability profile - $\boldsymbol{p_r}(t_{n-1})$

- Number of denoising steps - $\mathbf{N}$

# Discrete Diffusion Language Models are SOTA



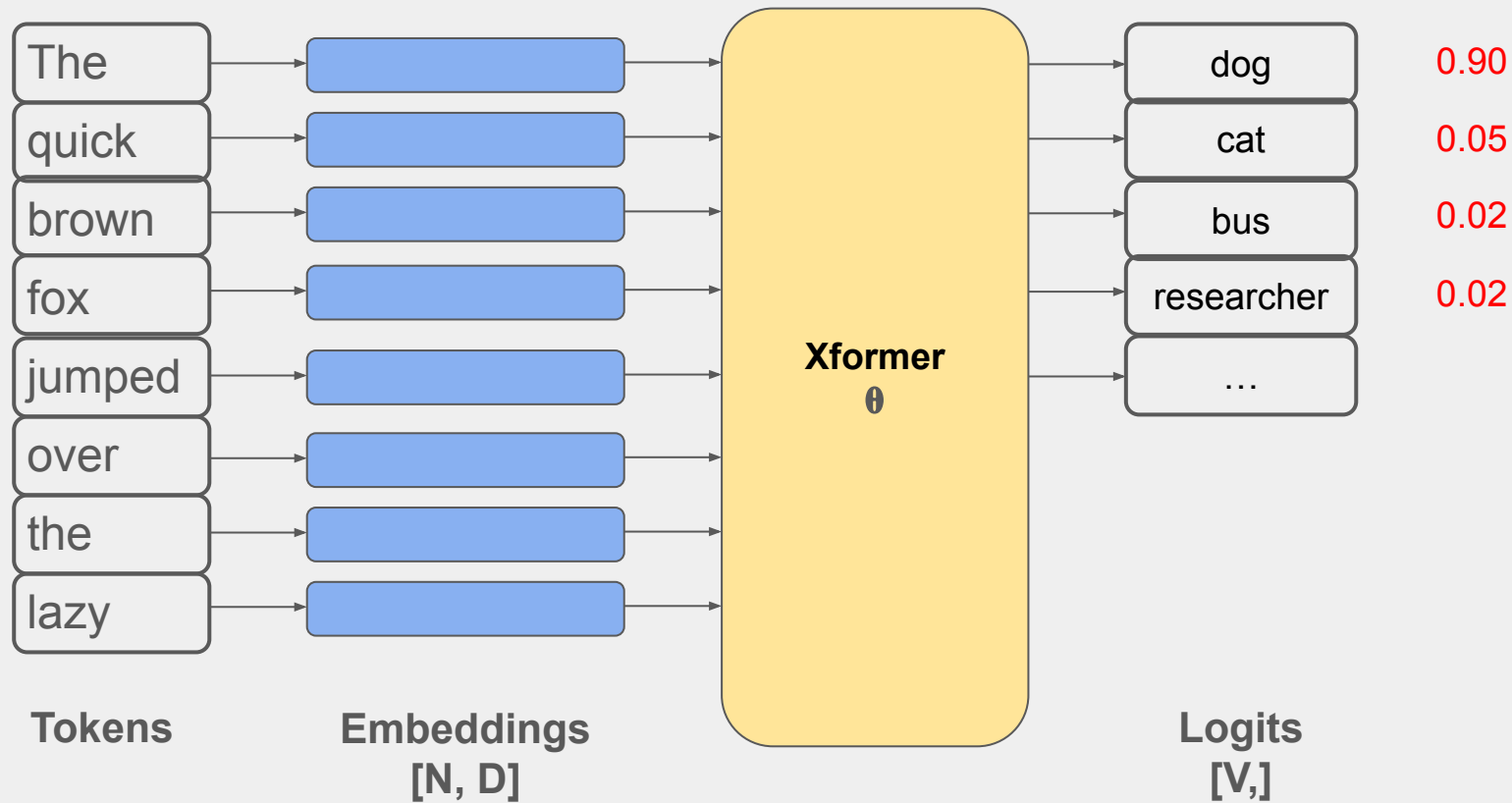Figure 1. **Zero/Few-Shot Benchmarks.** We scale LLaDA to an unprecedented size of 8B parameters from scratch, achieving competitive performance with strong LLMs (Dubey et al., 2024).

Nie et al., "Large Language Diffusion Models," arXiv:2502.09992, 2025.

# Discrete Diffusion Language Models are SOTA



Nie et al., "Large Language Diffusion Models," arXiv:2502.09992, 2025.

# 3. How?
# (Part tres)

# Language Models

| Tokens | Embeddings [N, D] | | Logits [V,] |
|---|---|---|---|
| The | | | dog — 0.90 |
| quick | | | cat — 0.05 |
| brown | | | bus — 0.02 |
| fox | Xformer θ | | researcher — 0.02 |
| jumped | | | … |
| over | | | |
| the | | | |
| lazy | | | |

# Language Models - Efficient Training

# Language Models - Efficient Training



Tokens

Logits [N,V] <-CE-> Next Tokens

# Language Models - Efficient Training

| | The | qui.. | bro.. | fox | jum. | over | the | lazy |
|---|---|---|---|---|---|---|---|---|
| The | | | | | | | | |
| quick | | | | | | | | |
| brown | | | | | | | | |
| fox | | | | | | | | |
| jumped | | | | | | | | |
| over | | | | | | | | |
| the | | | | | | | | |
| lazy | | | | | | | | |

**Tokens**

No for loops!

# Discrete diffusion for text - Training

| Context | Clean Sequence |
|---------|----------------|

| Context | Noised Sequence |
|---------|-----------------|

# Discrete diffusion for text

| Context | Clean Sequence |
|---------|----------------|

| Context | Noised Sequence |
|---------|-----------------|

[16384,] tokens

Full attention over all tokens can be very expensive!

# Discrete diffusion for text - Training



$$\pi_\theta(x_t, t)$$

# Discrete diffusion for text - Training

for block in blocks:

Block 1 | Context | Block | Remaining Sequence |

Block 2 | Context | Clean | Block | Remaining Sequence |

Block 3 | Context | Clean | Clean | Block | Remaining Sequence |

...

# Discrete diffusion for text - Inference



Block 1 — Context | Noisy Canvas ↻ DNS

Block 2 — Context | Noisy Canvas ↻ DNS

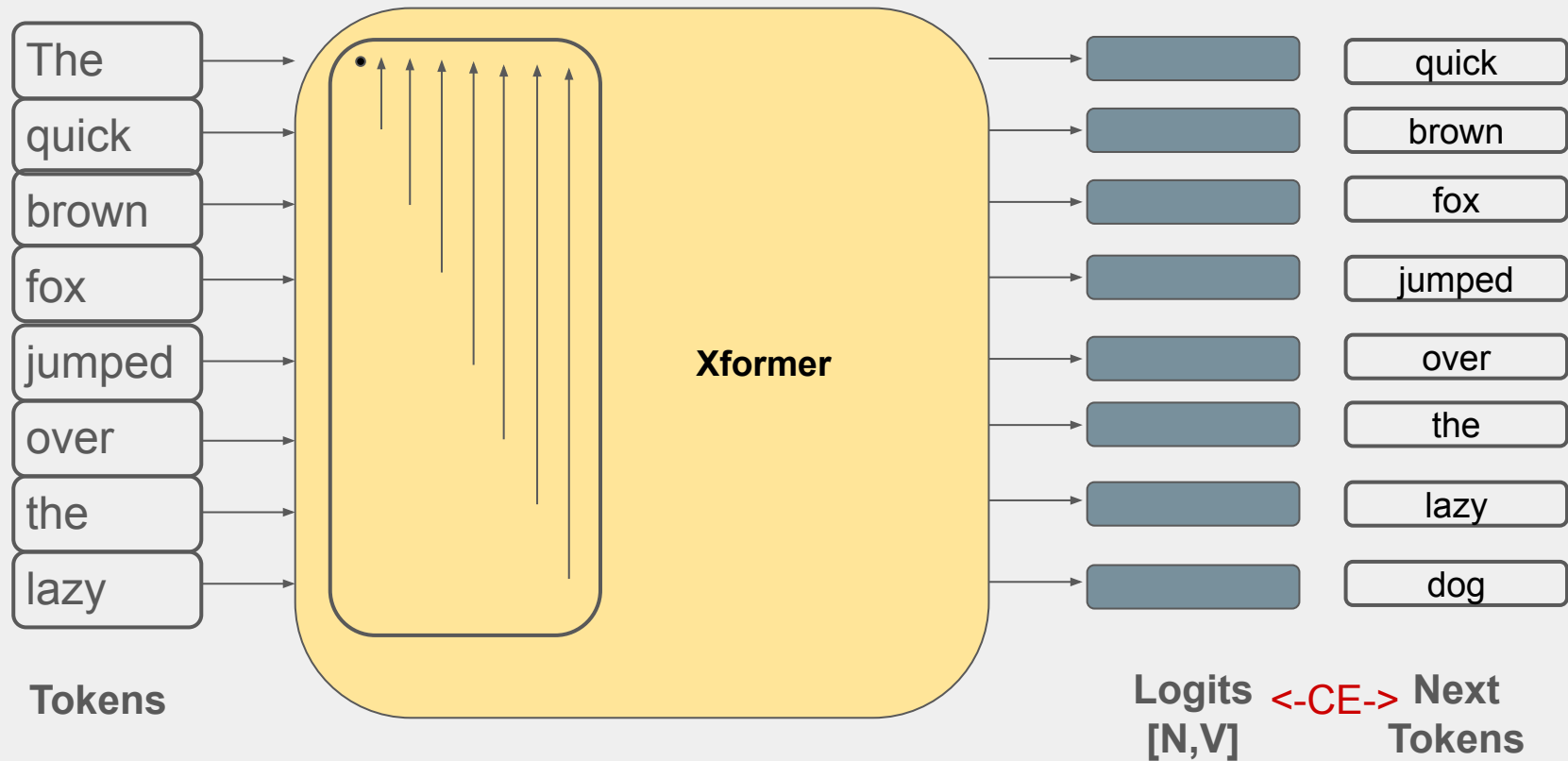Block 3 — Context | Noisy Canvas ↻ DNS

...

Enables KV caching between blocks by using block causal attention pattern!

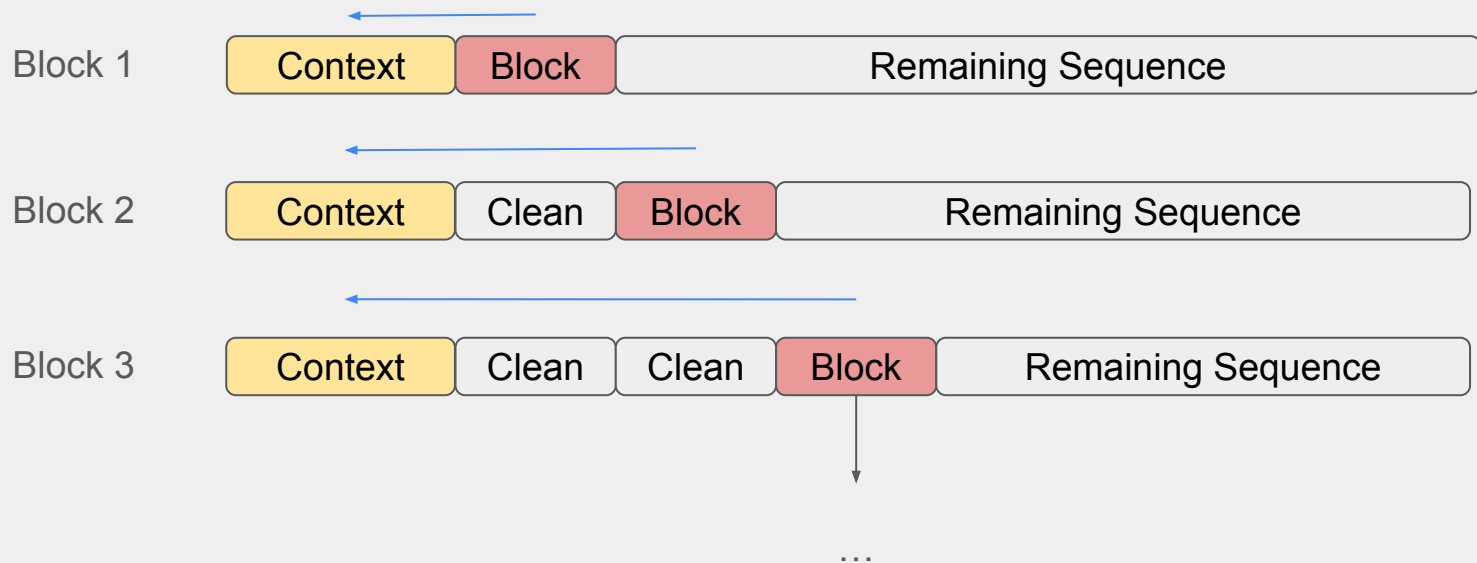# Discrete diffusion for text - Training

Can we do better?

# Language Models - Efficient Training

# Discrete diffusion for text - Training

for block in blocks:

Block 1     | Context | Block | Remaining Sequence |

Block 2     | Context | Clean | Block | Remaining Sequence |

Block 3     | Context | Clean | Clean | Block | Remaining Sequence |

…

# Discrete diffusion for text - Efficient Training

*Concatenate all noised blocks with the clean sequence!*

| Block 1 | Block 2 | Block 3 | clean sequence |
|---------|---------|---------|----------------|

Arriola et al., "Block Diffusion," arXiv:2503.09573, 2025.

# Discrete diffusion for text - Efficient Training

| Block 1 | Block 2 | Block 3 | clean sequence |



Figure 3: Example of Specialized Attention Mask

Arriola et al., "Block Diffusion," arXiv:2503.09573, 2025.

# Discrete diffusion for text - Efficient Training



Arriola et al., "Block Diffusion," arXiv:2503.09573, 2025.

# Discrete diffusion for text - Inference



Figure 4. **A Conceptual Overview of the Semi-autoregressive Sampling.**

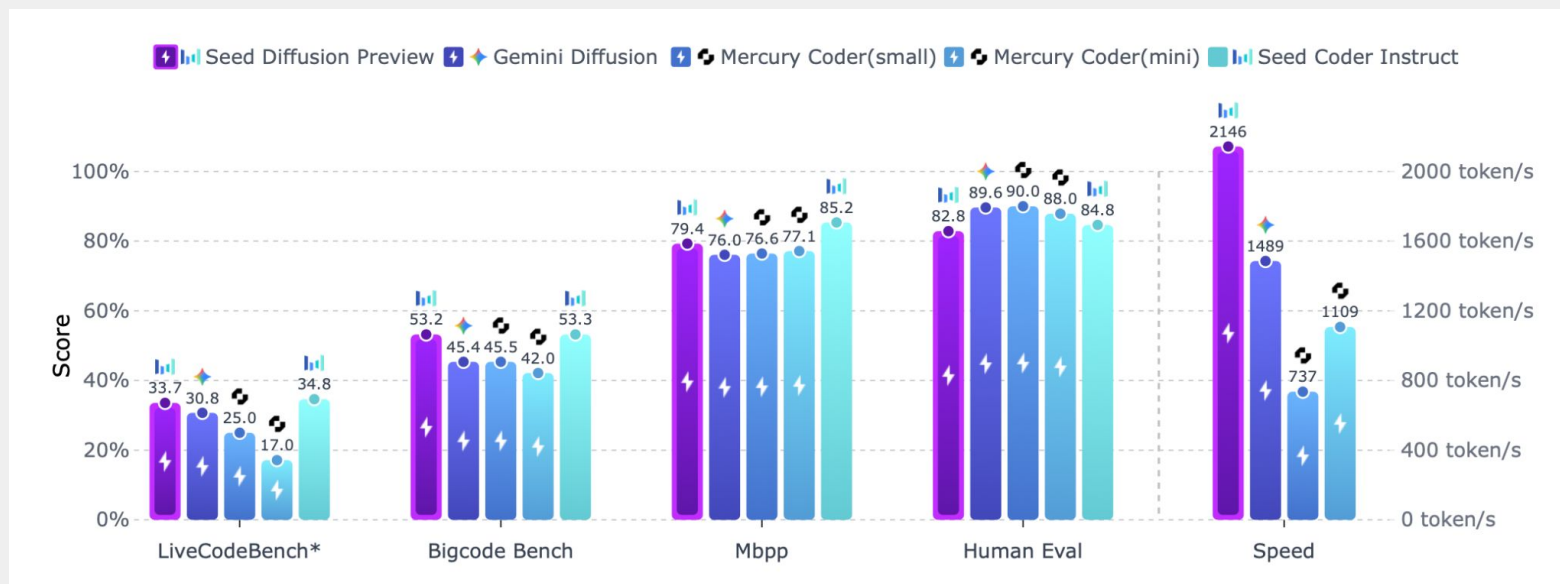Nie et al., "Large Language Diffusion Models," arXiv:2502.09992, 2025.

# 4. What's Next?

# Diffusion Language Models

Recent, fast Moving Field

- Dieleman et al., "Continuous Diffusion for Categorical Data - Dec. 2022
- Nie et al., "Large Language Diffusion Models - Feb 2025
- Gong et al., "Scaling Diffusion Language Models" - April 2025
- Arriola et al., "Block Diffusion" - May 2025
- Google Deepmind, "Gemini Diffusion" - May 2025
- ByteDance, "Seed Diffusion" - August 2025

# Diffusion Language Models



ByteDance, "Seed Diffusion", arXiv: 2508.02193v1

# Diffusion Language Models

Lots of ideas in the air:

- Edit-based Noise Process
- On-policy Diffusion Learning
- Distillation
- Adaptive Computation
- Combining Continuous and Discrete
- Post-training and RL

    …

# Diffusion Language Models

Twitter: @sahnihim

Gmail: himsahni@google.com