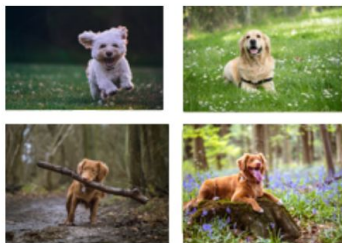


Discrete Diffusion Models

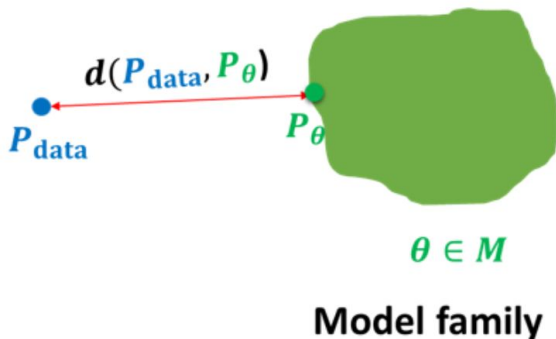


The Task of Generative Modeling

Suppose we are given a training set of examples, e.g., images of dogs



$$x_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$

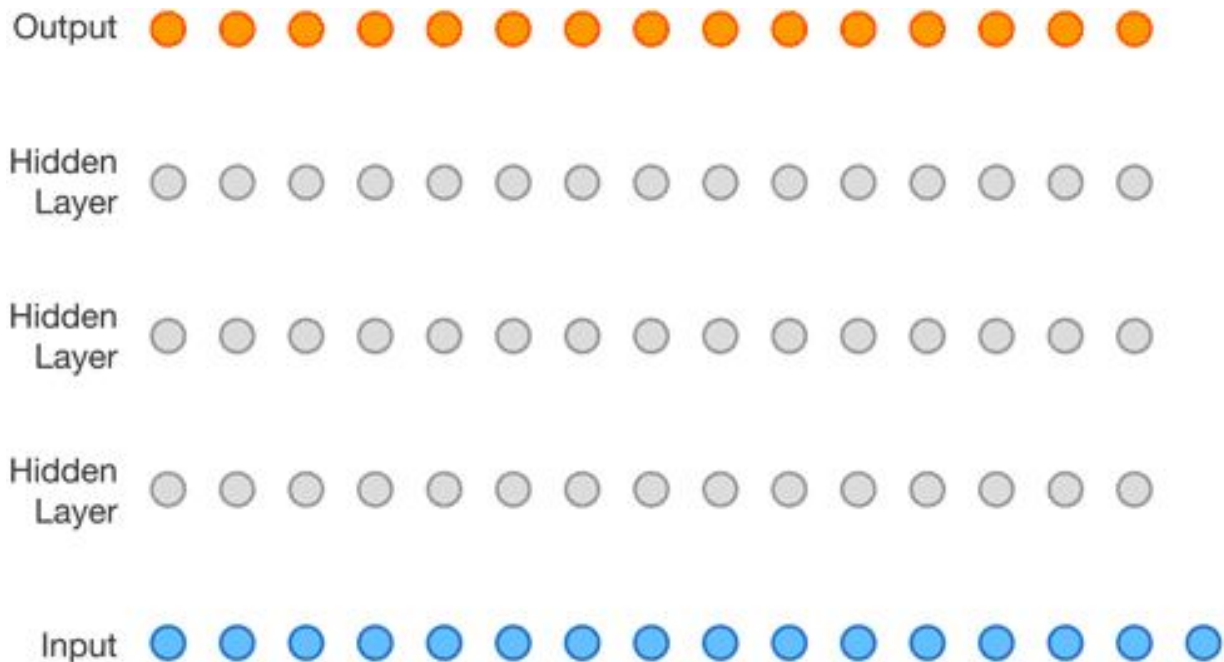


Our Goal: define a probability distribution $p(x)$ over images x .

- 1 **Modeling:** Define a family of distributions p_{θ} parameterized by θ .
- 2 **Learning:** Search for model parameters θ based on training data \mathcal{D} .
- 3 **Inference:** Generate data x by sampling from $p_{\theta}(x)$.

Last Lecture: Autoregressive Models (Overview)

Autoregressive models generate data sequentially by predicting the next element from previously generated ones.



Last Lecture: Autoregressive Models (Theory)

- Without loss of generality, we can use chain rule for factorization

$$p(x_1, \dots, x_{784}) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \cdots p(x_n | x_1, \dots, x_{n-1})$$

- We parameterize the conditionals using **logistic regression**:

$$p(x_1, \dots, x_{784}) = p(x_1; \theta) p_{\text{logit}}(x_2 | x_1; \theta) p_{\text{logit}}(x_3 | x_1, x_2; \theta) \cdots p_{\text{logit}}(x_n | x_1, \dots, x_{n-1}; \theta)$$

Given training data $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$, **maximum likelihood learning** solves:

$$\max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \sum_{i=1}^n \log p_{\text{neural}}(x_i | \mathbf{x}_{<i}; \theta)$$

Last Lecture: Autoregressive Models (Pros & Cons)

Autoregression

There

- ✓ Exact Likelihood Evaluation & Training
- ✓ Exact Sampling
- ✓ Generates Sequences of Arbitrary Length

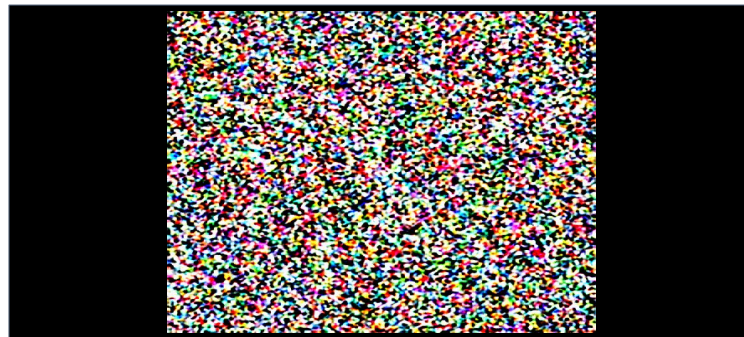
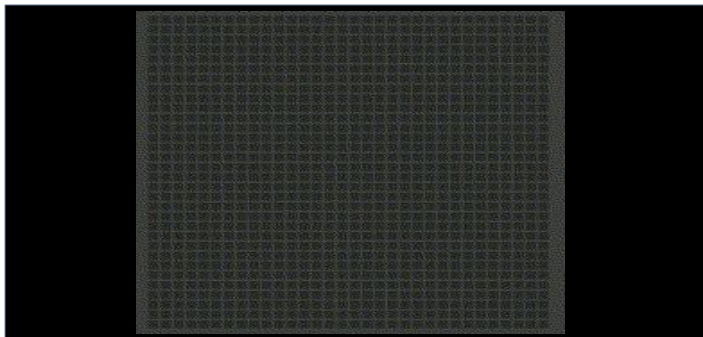
- ✗ Sequential (Slow) Generation
- ✗ No Latent Representations
- ✗ Compounding Errors At Generation

Today's Generative Modeling Landscape

Autoregression

Diffusion

Images



Text



```
]? trunk })  
oka(()ns  
.....  
(((Abstractreported-but',//  
.....  
((.(.(ns(handlerewis(___ originates.getDayROSS.  
).....cameractic" norms prosper taxing]));  
Models.  
<
```

This talk

Diffusion vs. Autoregressive Language Models

Autoregression

There

- ✗ Sequential Generation
- ✗ Previous Tokens Can't Be Edited
- ✗ Fixed Number of Generation Steps
- ✗ Causal Context

Diffusion

Hirsch
account

- ✓ Parallel (fast) Generation
- ✓ Can Fix Previous Mistakes
- ✓ Generate in More/Fewer Steps
- ✓ Bidirectional Context

Diffusion vs. Autoregressive Language Models

Autoregression

There



Diffusion

Hirsch
account

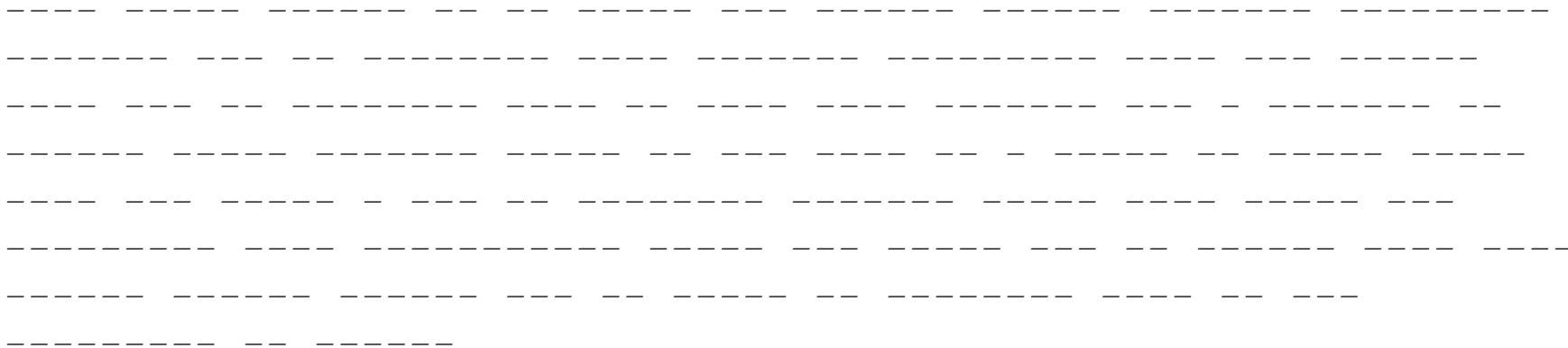
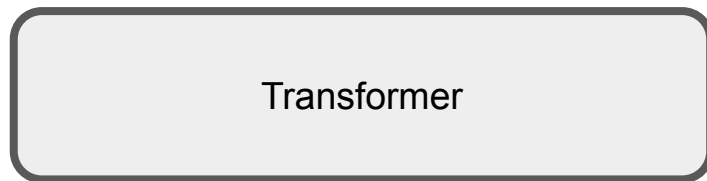
Today's hardware in parallel, but LLM inference algorithms are sequential.
Diffusion is the right language modeling approach for today's hardware.

This Talk: Simple Masked Diffusion Language Models

$$x \sim p_{\theta}(x)$$

Many years later, as he faced the firing squad, Colonel Aureliano Buendía was to remember that distant afternoon when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the bank of a river of clear water that ran along a bed of polished stones, which were white and enormous, like prehistoric eggs. The world was so recent that many things lacked names, and in order to indicate them it was necessary to point.

Sampling from a Masked Language Model



Transformer

----- squad, -----
----- to remember ----- when -----
----- to -----
twenty adobe ----- of ----- water
----- along ----- polished ----- which -----
----- prehistoric ----- many
----- in -----

Transformer

Many years later, as he faced --- ----- squad, -----
Buendía was to remember that distant ----- when his -----
took --- to discover ---- At that ---- ----- was a village of
twenty adobe houses, built -- the bank of - river of clear water
that ran along - --- of polished stones, which were ----- and
----- like prehistoric eggs. --- ----- --- so recent ---- many
things lacked names, and in ----- -- ----- them it was
----- to -----

Transformer

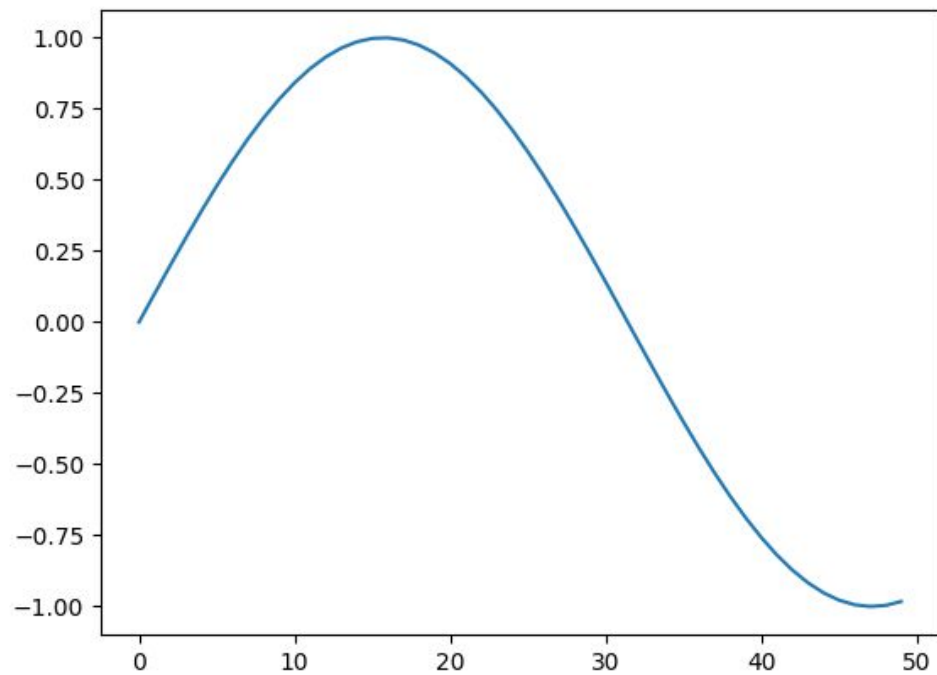
Many years later, as he faced the firing squad, Colonel Aureliano Buendía was to remember that distant ----- when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the bank of - river of clear water that ran along a bed of polished stones, which were white and enormous, like prehistoric eggs. The world --- so recent that many things lacked names, and in order -- indicate them it was ----- to point.

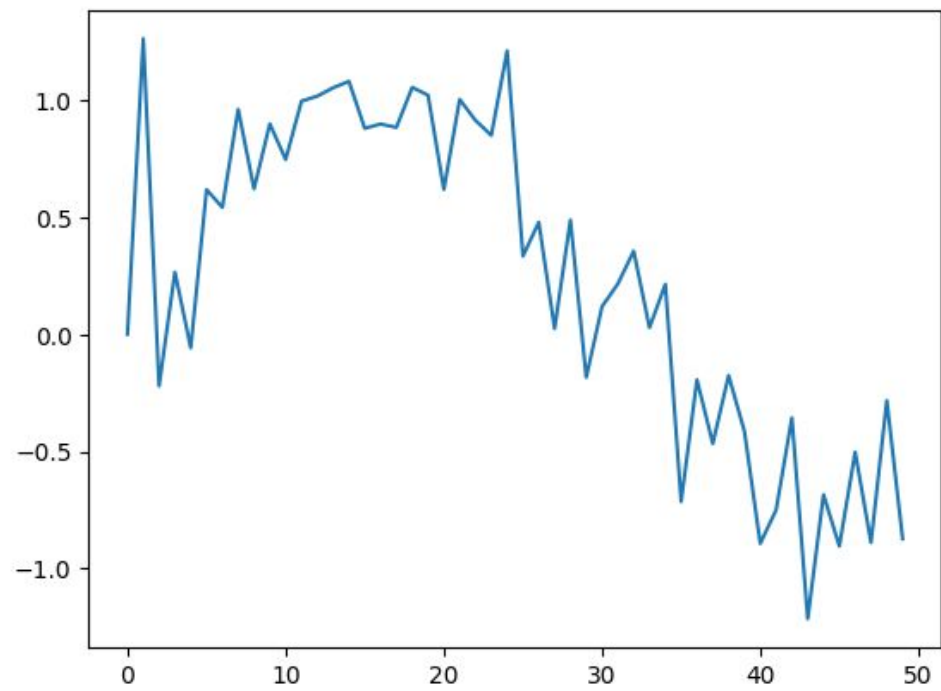
$$x \sim p_{\theta}(x)$$

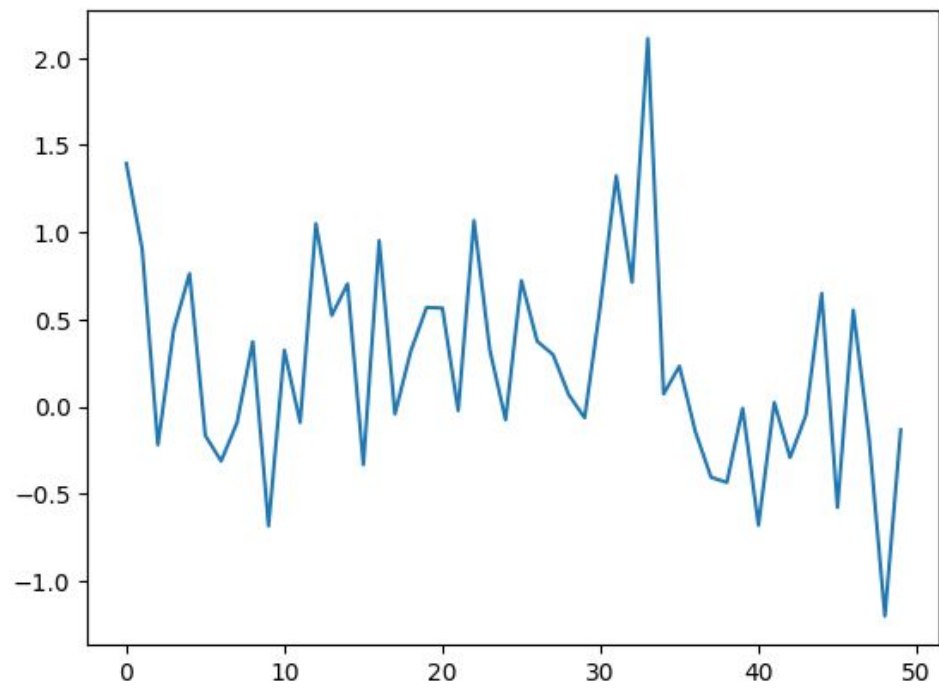
Many years later, as he faced the firing squad, Colonel Aureliano Buendía was to remember that distant afternoon when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the bank of a river of clear water that ran along a bed of polished stones, which were white and enormous, like prehistoric eggs. The world was so recent that many things lacked names, and in order to indicate them it was necessary to point.

Diffusion Background

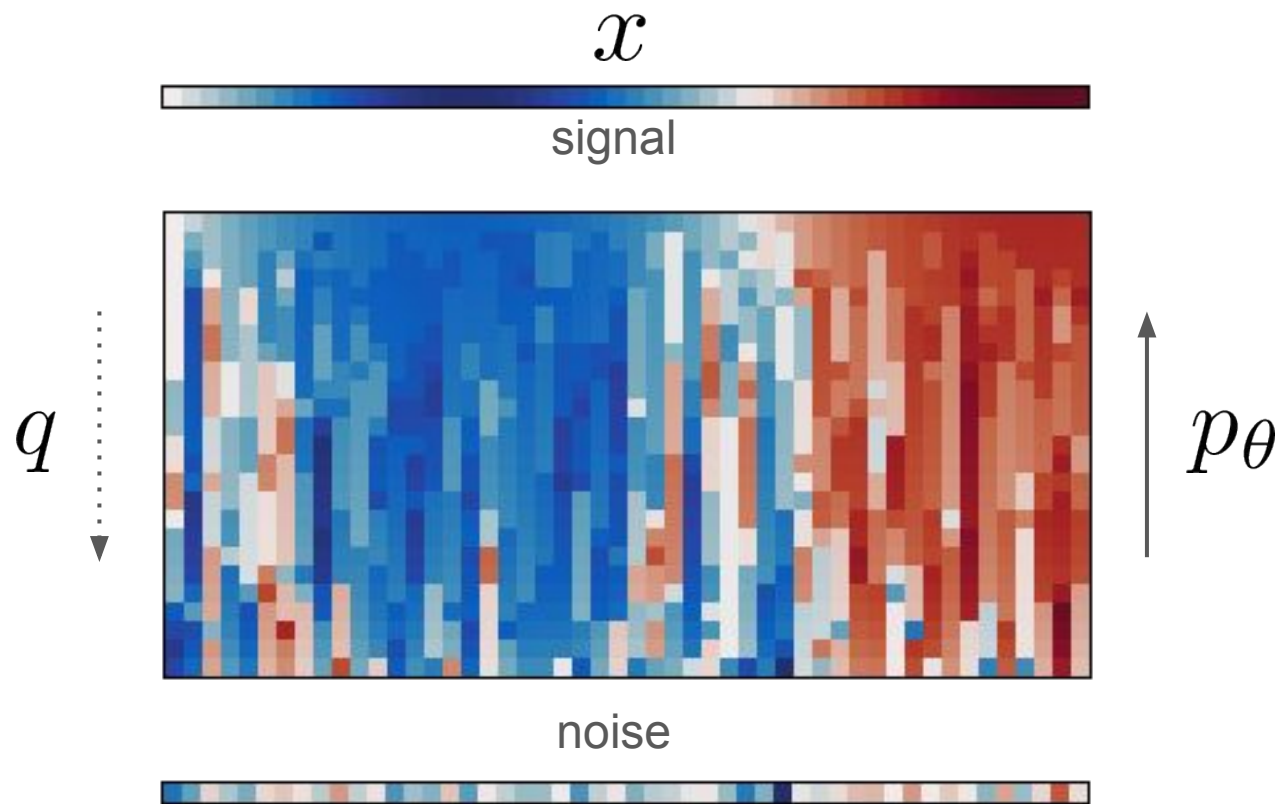
$$x \sim p_{\theta}(x)$$



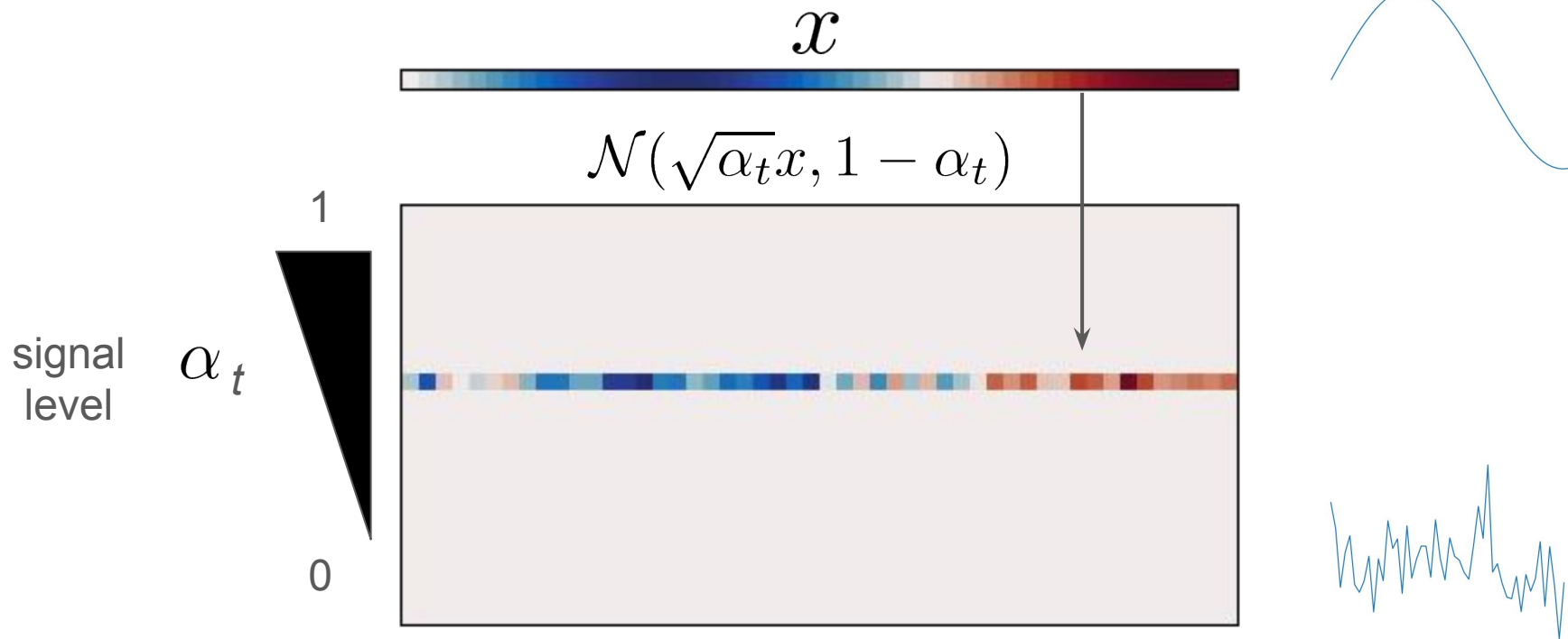




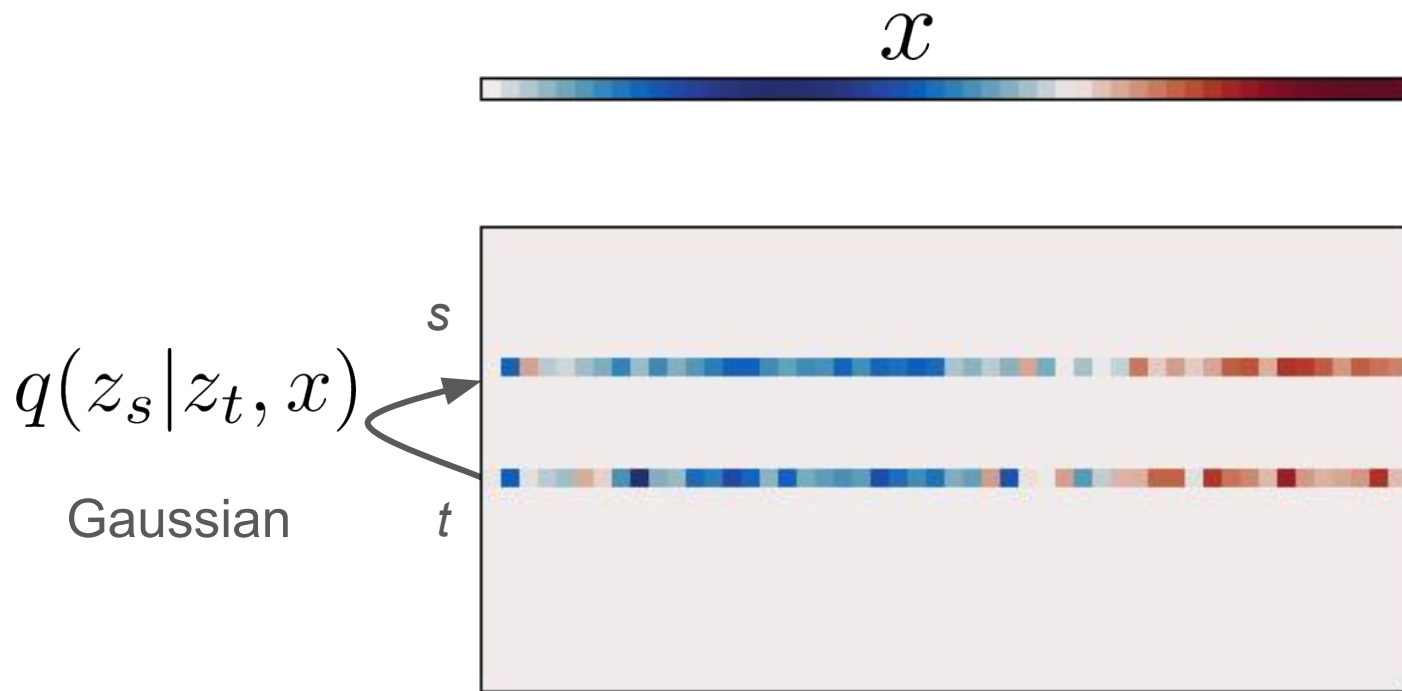
Gaussian Diffusion



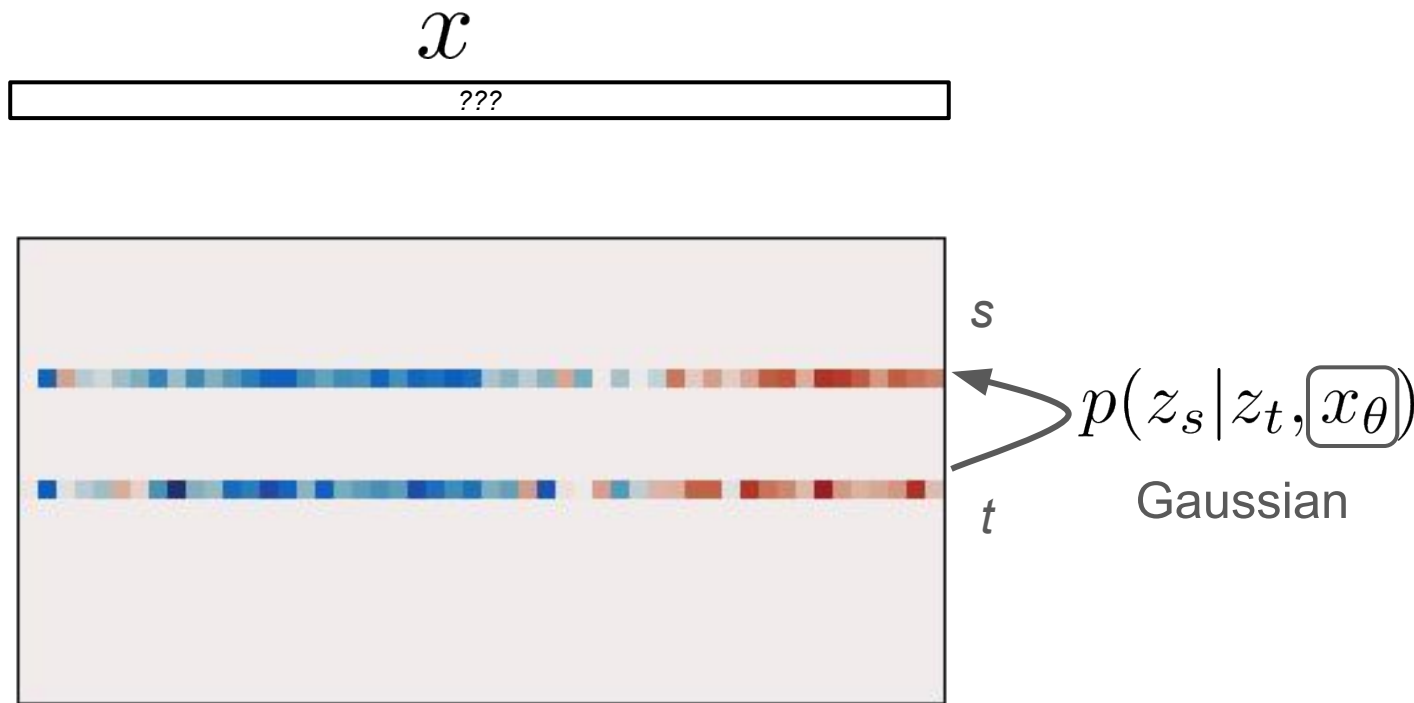
Noising Process: Gaussian Noise



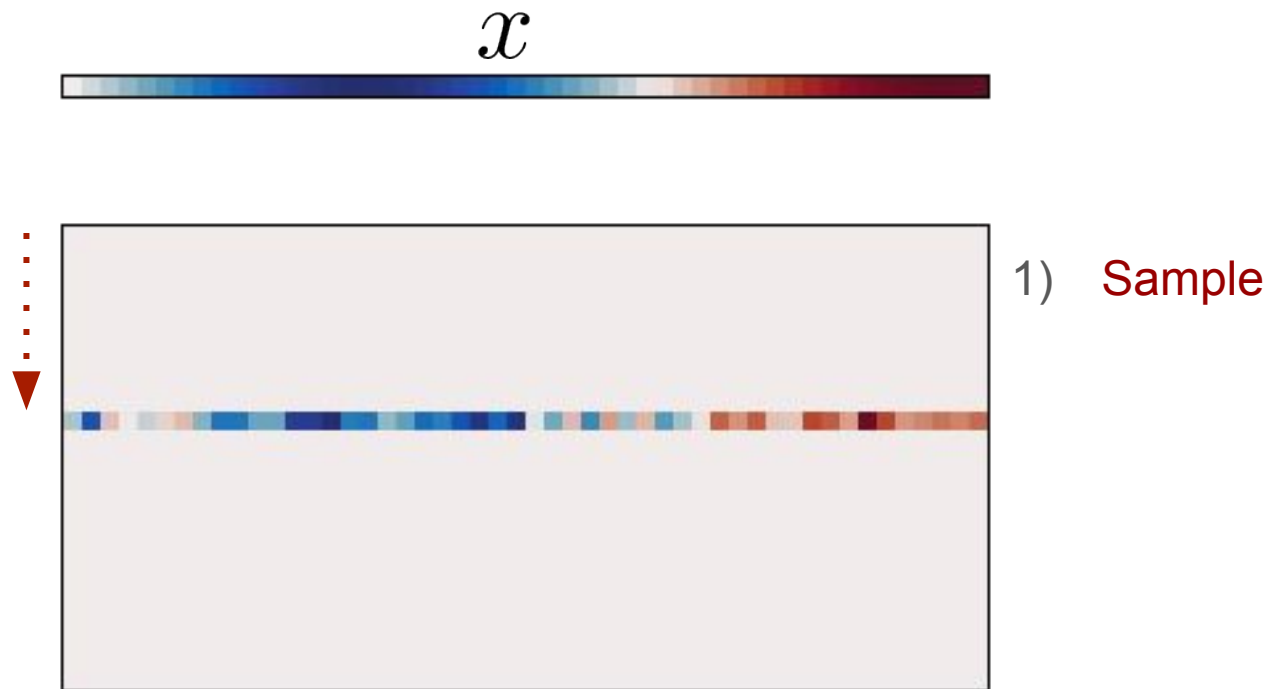
Gaussian Noise Implies Gaussian Posterior



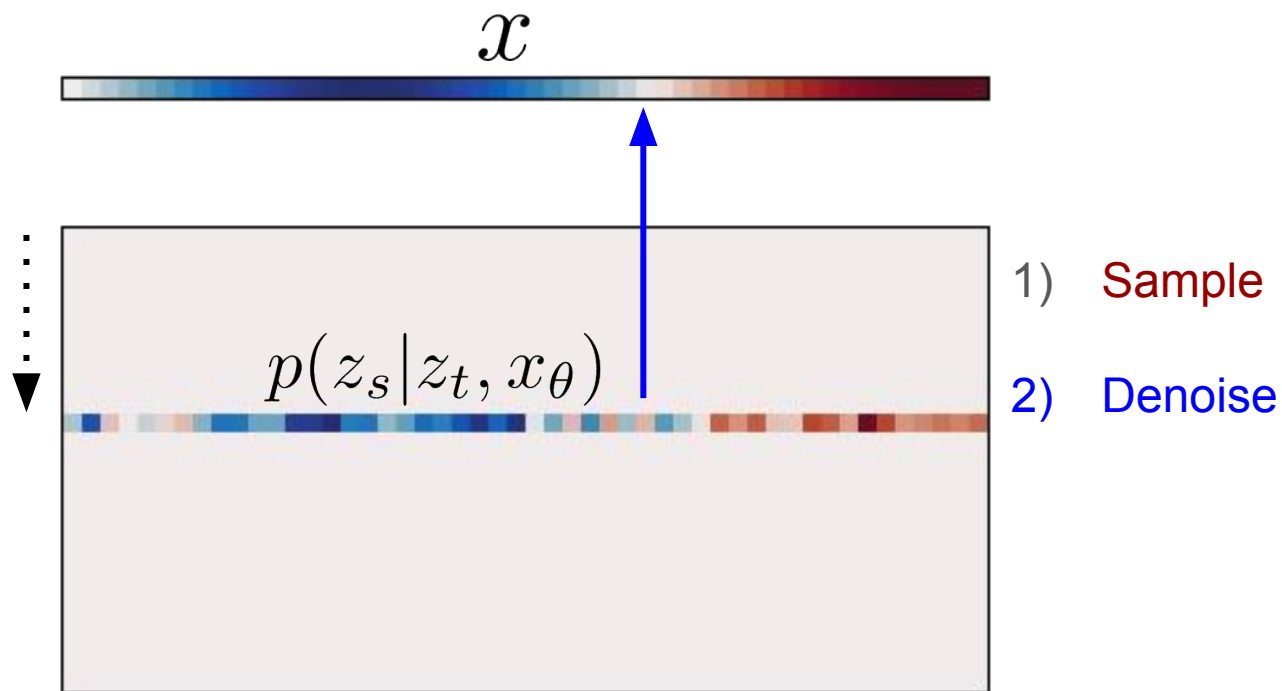
Reverse Process: Denoising



Learning to Denoise



Learning to Denoise

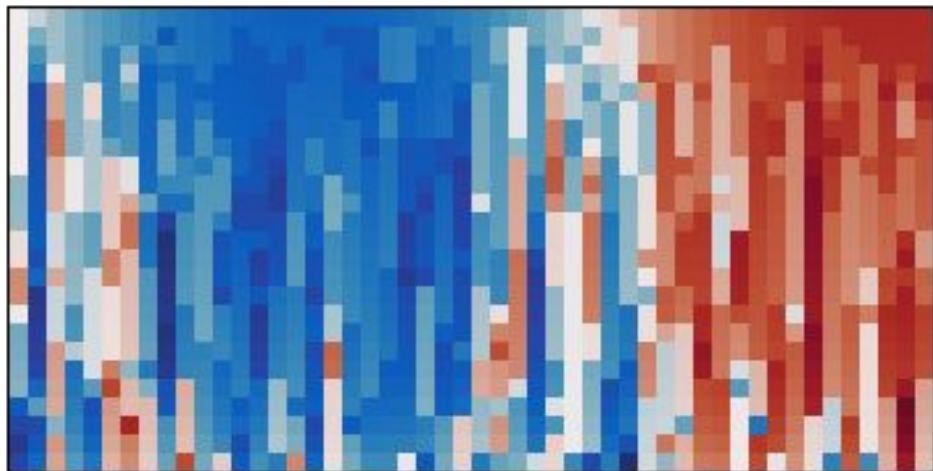


Generation: Iterative Denoising

x



signal



p_θ

noise



Simple Discrete Masking Diffusion

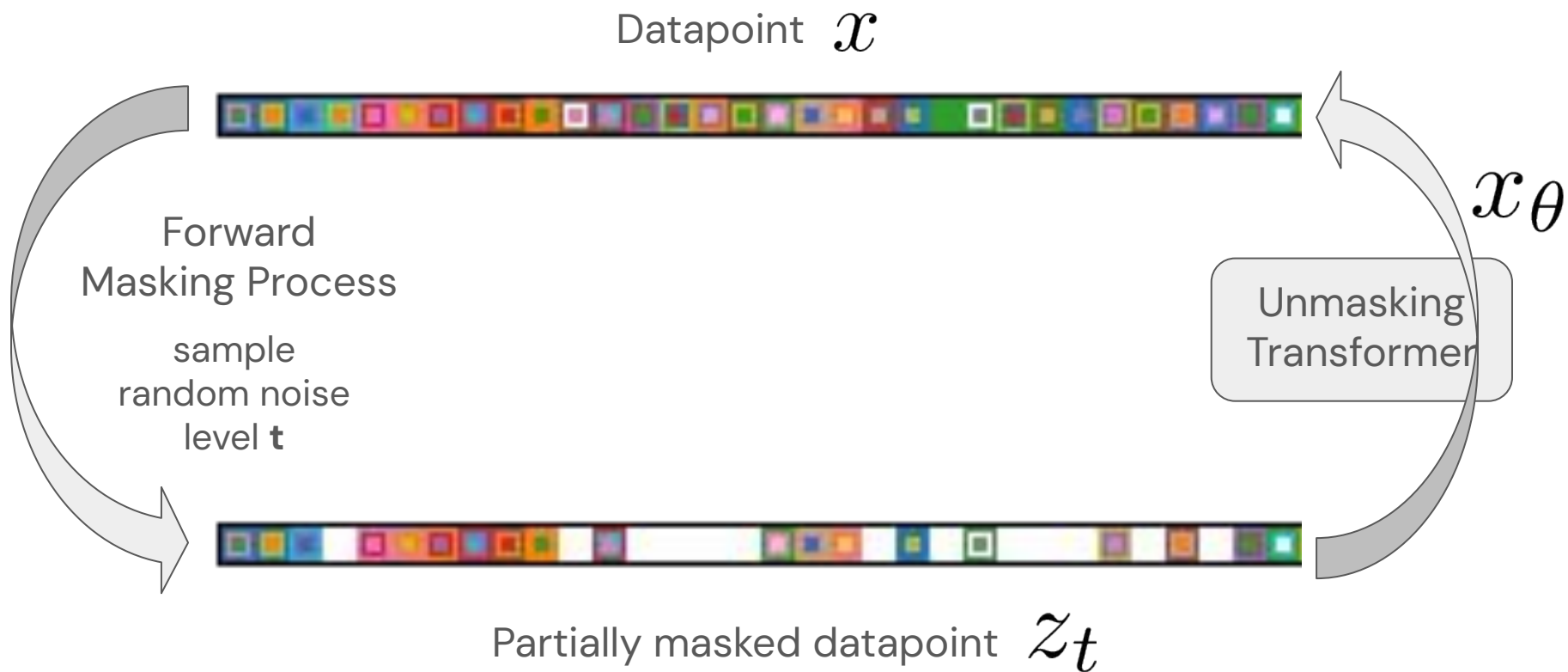
Many years later, as he faced the firing squad, Colonel Aureliano Buendía was to remember that distant afternoon when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the bank of a river of clear water that ran along a bed of polished stones, which were white and enormous, like prehistoric eggs. The world was so recent that many things lacked names, and in order to indicate them it was necessary to point.



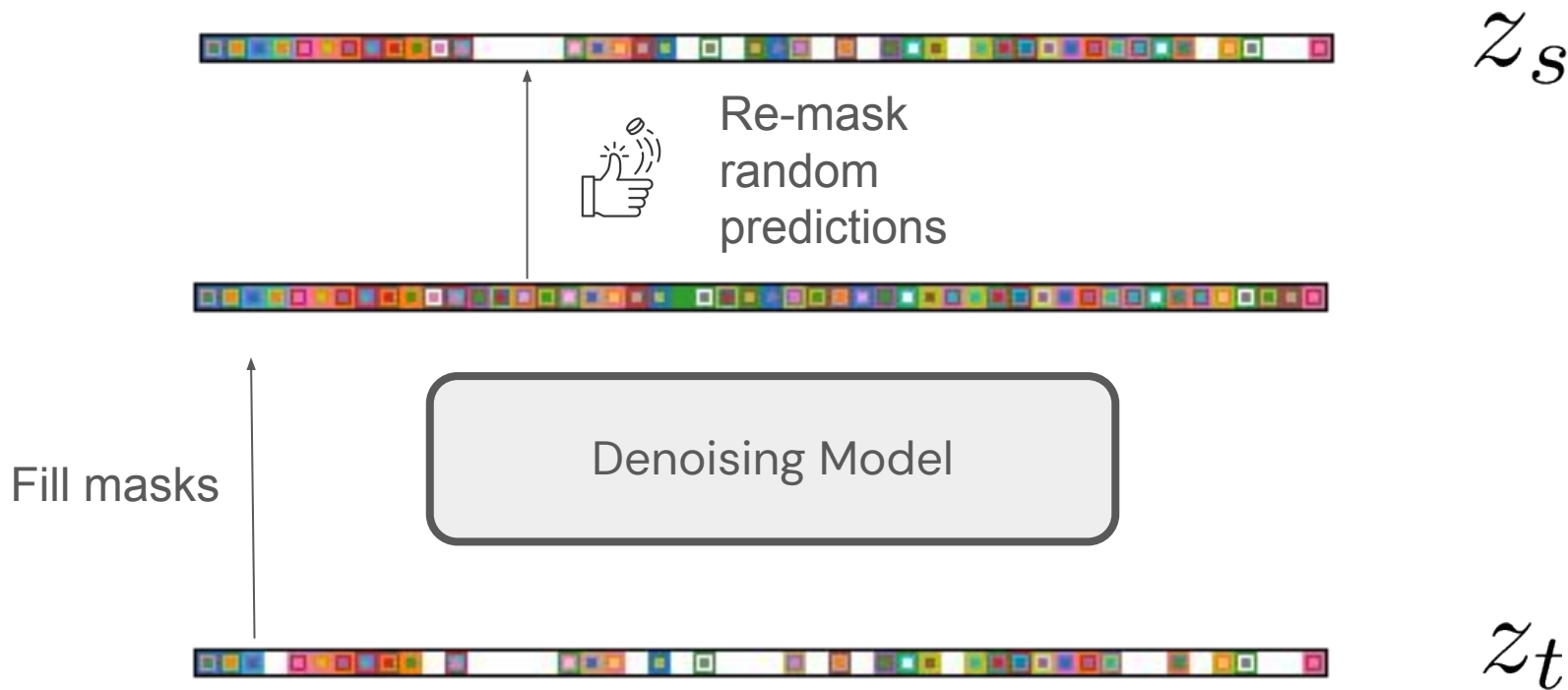
Many years later, -- he faced --- firing squad, ----- Aureliano Buendía was to remember that ----- afternoon when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the ---- of a river of clear water that ran ----- a bed of polished stones, which were ----- --- enormous, like prehistoric eggs. The ----- was so recent that many things lacked names, and in order to indicate them -- was necessary to point.



Model Parameterization & Training: Unmasking Transformer



One Step of Generation



Multiple Steps of Generation

*Many years later, as he faced squad,
Buendía was to remember that distant when his
took to discover At that was a village of
twenty adobe houses, built the bank of river of clear water
that ran along of polished stones, which were and
like prehistoric eggs. so recent many
things lacked names, and in them it was
to*



Multiple Steps of Generation

Many years later, as he faced the firing squad, Colonel Aureliano Buendía was to remember that distant when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the bank of river of clear water that ran along a bed of polished stones, which were white and enormous, like prehistoric eggs. The world so recent that many things lacked names, and in order indicate them it was to point.



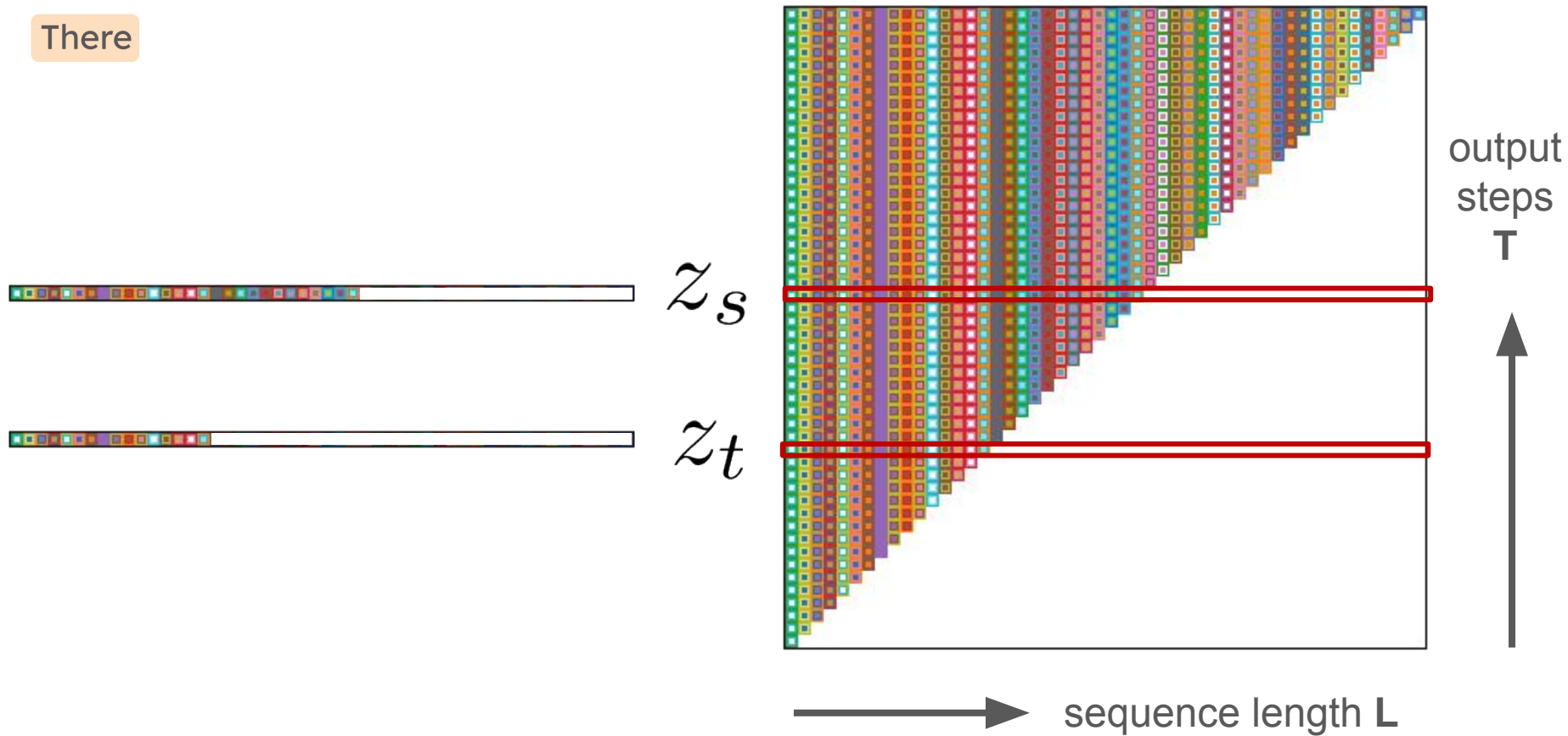
Multiple Steps of Generation

Many years later, as he faced the firing squad, Colonel Aureliano Buendía was to remember that distant afternoon when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the bank of a river of clear water that ran along a bed of polished stones, which were white and enormous, like prehistoric eggs. The world was so recent that many things lacked names, and in order to indicate them it was necessary to point.

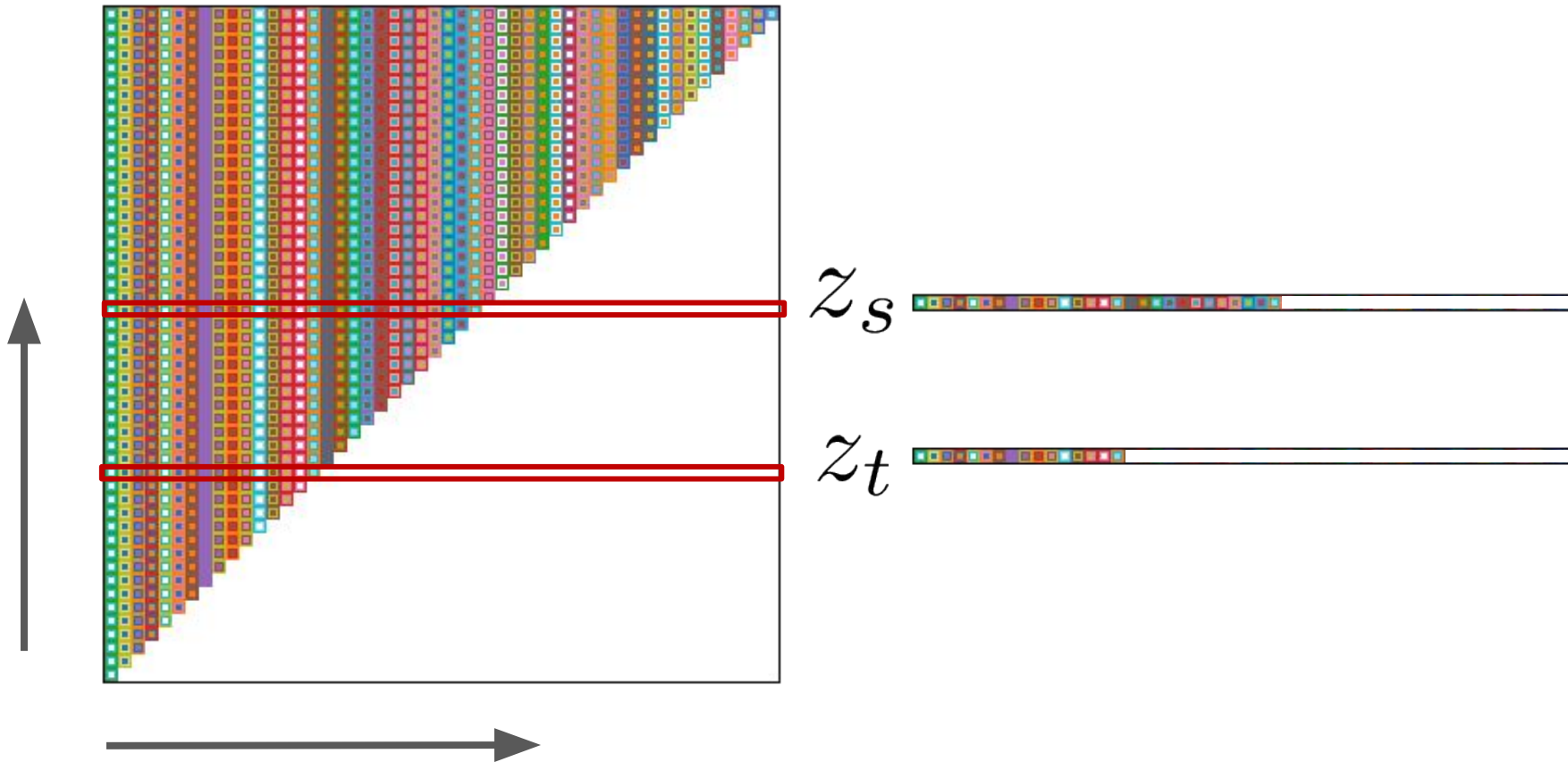


Example: Autoregressive Unmasking

There



Example: Autoregressive Unmasking



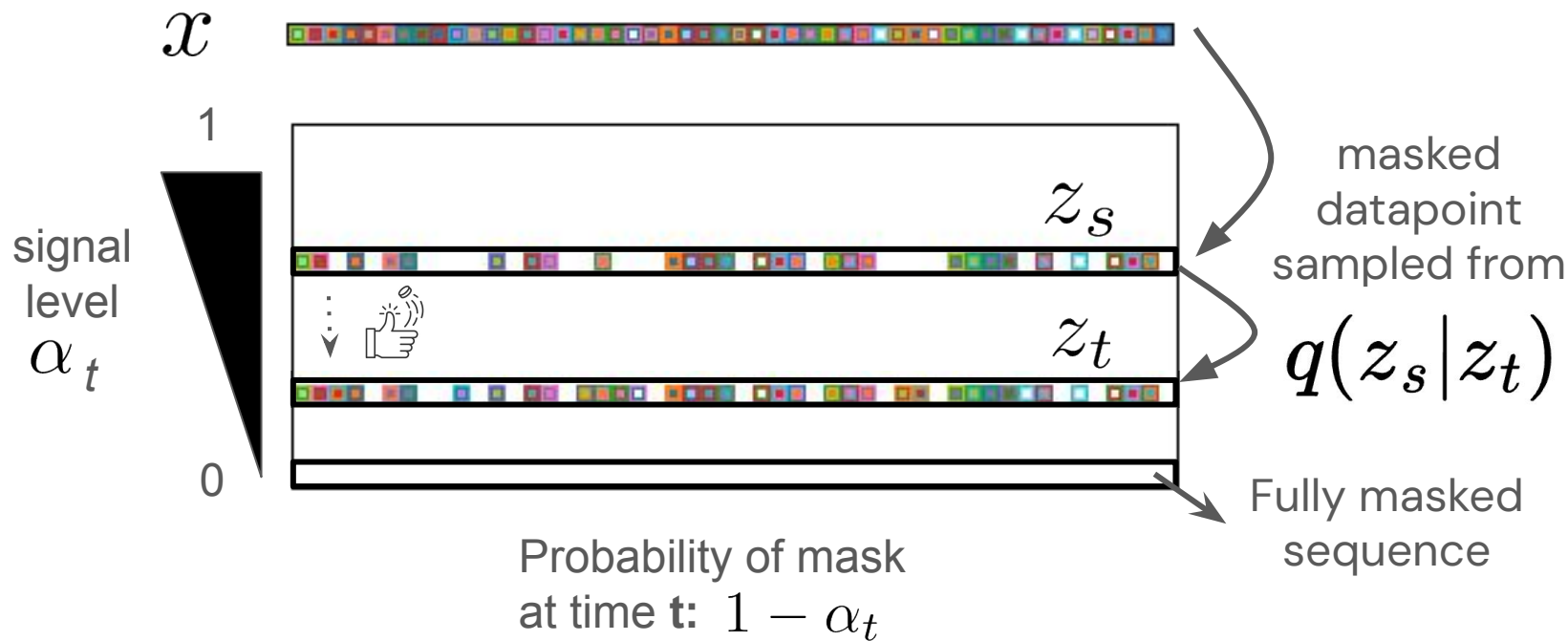
Modeling: Defining a Generative Model p

Hirsch
account

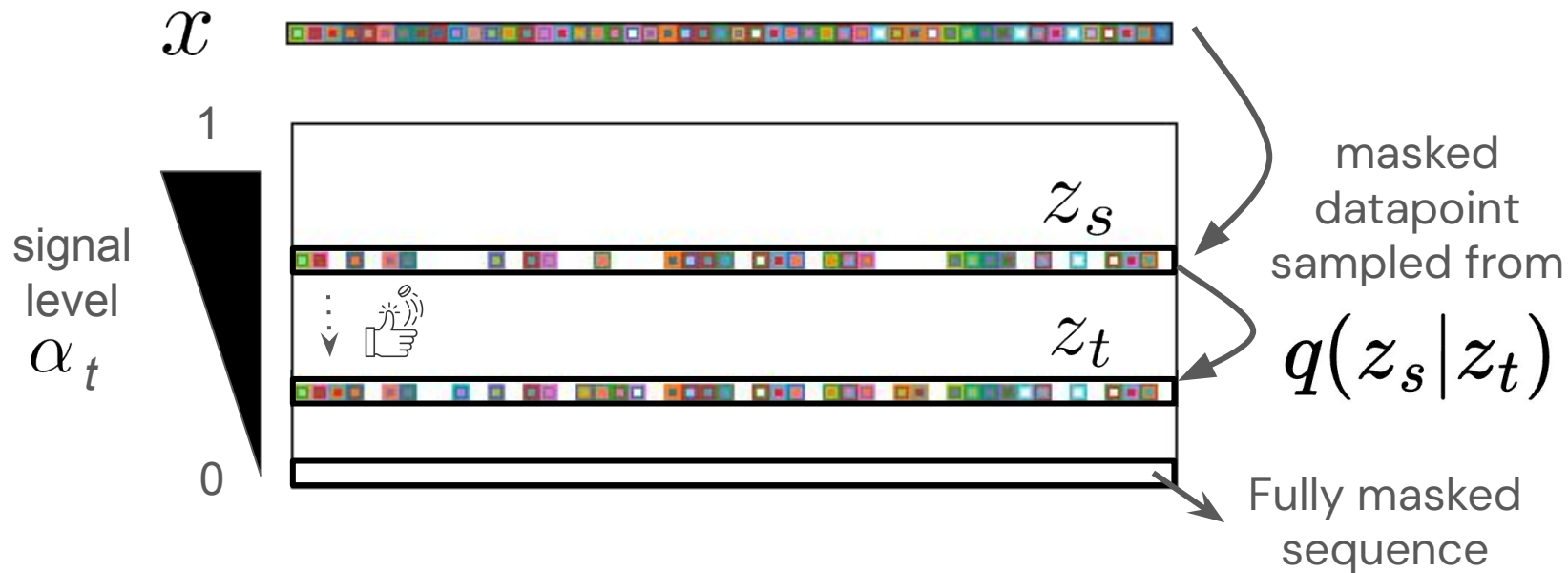
$$x, z_0, \dots, z_T \sim p(x|z_0) \prod_{t=1}^T p(z_{t-1}|z_t)$$



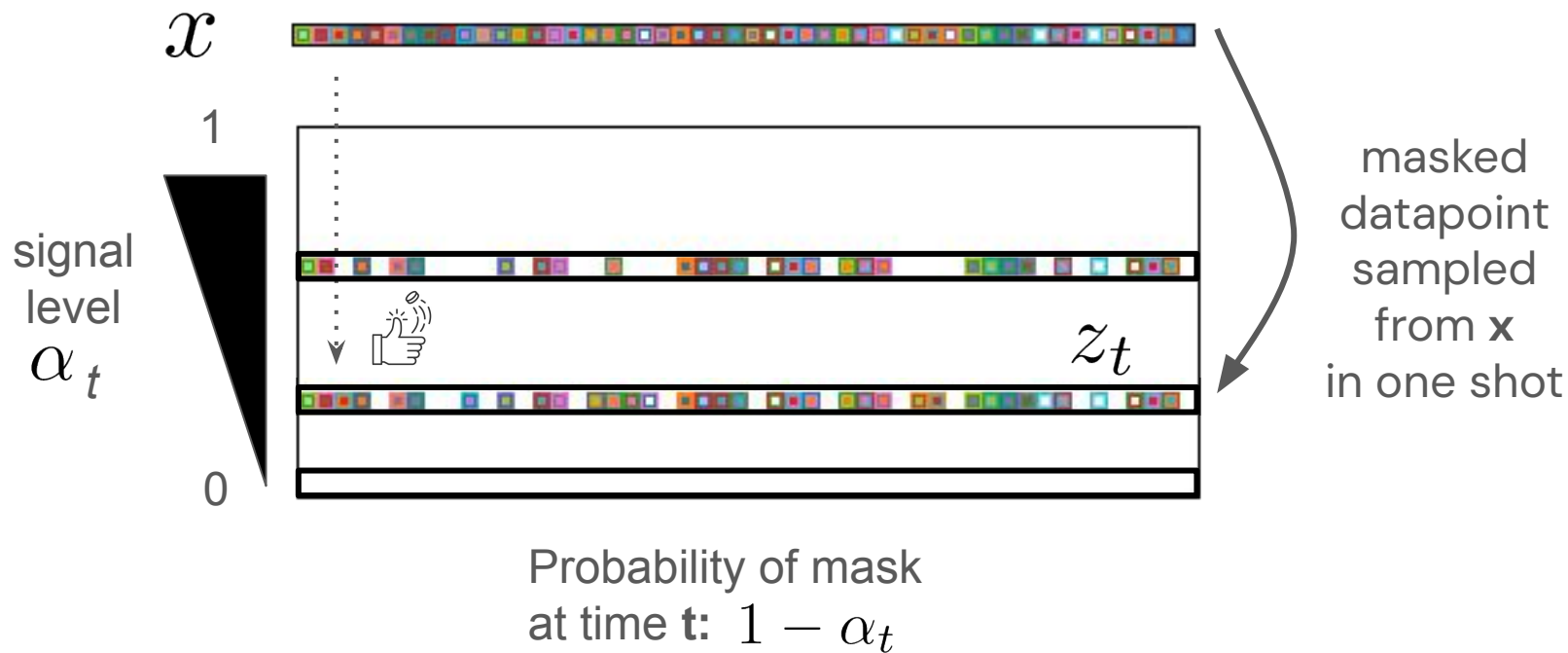
How Do We Train p ? To Reverse Noise Process q



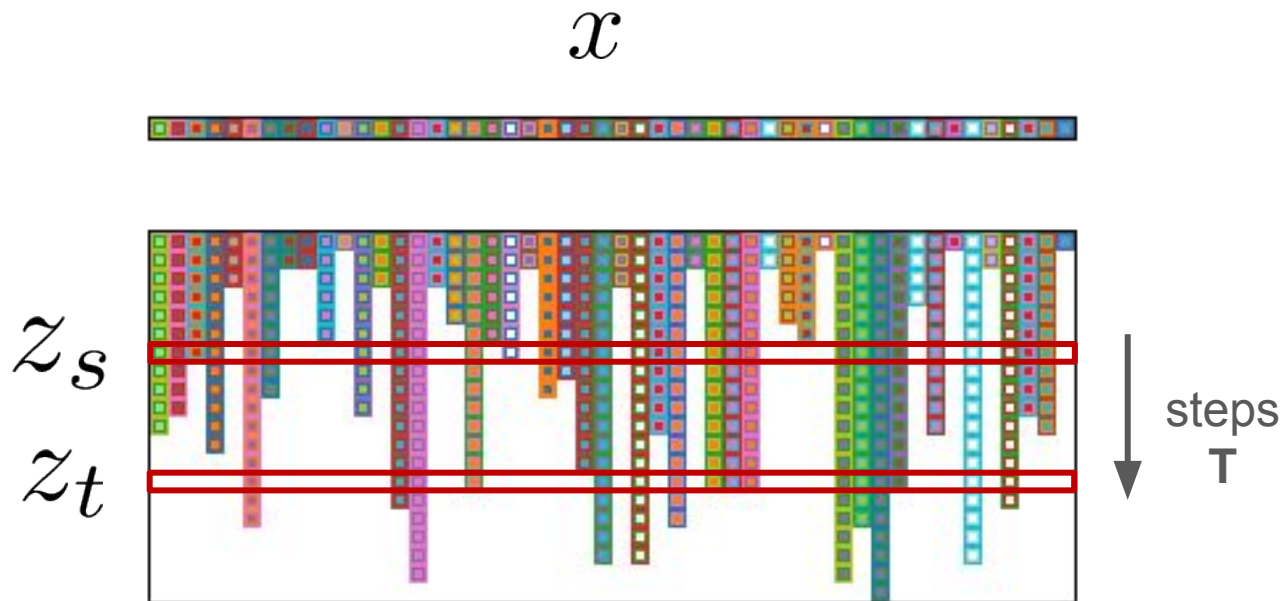
How Do We Train p ? To Reverse Noise Process q



How Do We Train p ? To Reverse Noise Process q



How Do We Train p ? To Reverse Noise Process q



The noise process defines a path from clean x to masked z 's.

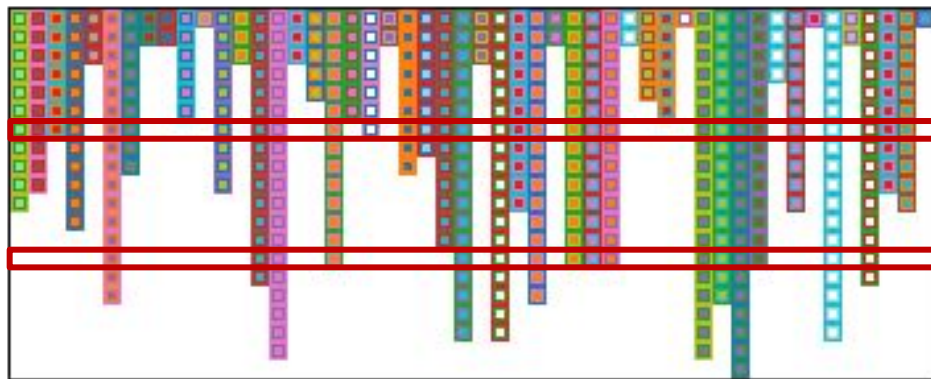
How Do We Train p ? To Reverse Noise Process q

x



$q(z_s | z_t, x)$

Optimal
unmasking
process



$p(z_s | z_t, x_\theta)$

Learned
denoising
model



The noise process defines a path from clean x to masked z 's.

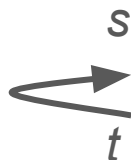
How Do We Train p ? To Reverse Noise Process q

x



$$q(z_s | z_t, x)$$

Unmasking
posterior



$$p(z_s | z_t, x_\theta)$$

Denoising
model

Learning \mathbf{p} : A Variational Inference Perspective

The classical approach to learning a latent variable \mathbf{p} is to maximize the ELBO:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})]$$

The (negative) ELBO has (in general) the following form:

$$\mathcal{L}_{\text{NELBO}} = \mathbb{E}_q \left[\mathcal{L}_{\text{reconstr}} + \underbrace{\mathbb{E}_t \text{KL}(q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) || p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t))}_{\mathcal{L}_{\text{diffusion}}} + \mathcal{L}_{\text{prior}} \right]$$

Learning \mathbf{p} : Simplifying the ELBO

The (negative) ELBO for masked diffusion has (in general) the following form:

$$\mathcal{L}_{\text{NELBO}} = \mathbb{E}_q \left[\mathcal{L}_{\text{reconstr}} + \underbrace{\mathbb{E}_t \text{KL}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) || p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t))}_{\mathcal{L}_{\text{diffusion}}} + \mathcal{L}_{\text{prior}} \right]$$

Naively, an unsimplified discrete diffusion loss equals the following:

$$\mathcal{L}_{\text{diffusion}} = \mathbb{E} \left[\frac{\alpha_s - \alpha_t}{1 - \alpha_t} \log \frac{\alpha_t \mathbf{x}_\theta(\mathbf{z}_t)^\top \mathbf{m} + (1 - \alpha_t)}{(1 - \alpha_t) \mathbf{x}_\theta(\mathbf{z}_t)^\top \mathbf{x}} + \frac{1 - \alpha_s}{1 - \alpha_t} \log \frac{(1 - \alpha_s)(\alpha_t \mathbf{x}_\theta(\mathbf{z}_t)^\top \mathbf{m} + (1 - \alpha_t))}{(1 - \alpha_t)(\alpha_s \mathbf{x}_\theta(\mathbf{z}_t)^\top \mathbf{m} + (1 - \alpha_s))} \right] \mathbf{z}^\top \mathbf{m}$$

Learning \mathbf{p} : Simplifying the ELBO

The (negative) ELBO for masked diffusion has (in general) the following form:

$$\mathcal{L}_{\text{NELBO}} = \mathbb{E}_q \left[\mathcal{L}_{\text{reconstr}} + \underbrace{\mathbb{E}_t \text{KL}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) || p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t))}_{\mathcal{L}_{\text{diffusion}}} + \mathcal{L}_{\text{prior}} \right]$$

Our work uses properties of asking to obtain the following simplified form:

$$\mathcal{L}_{\text{diffusion}} = \mathbb{E}_{t, \mathbf{z}_t} \frac{\alpha_{t-1} - \alpha_t}{1 - \alpha_t} \log p_\theta(\mathbf{x} | \mathbf{z}_t)$$

Learning \mathbf{p} : Simplifying the ELBO

The (negative) ELBO for masked diffusion has (in general) the following form:

$$\mathcal{L}_{\text{NELBO}} = \mathbb{E}_q \left[\mathcal{L}_{\text{reconstr}} + \underbrace{\mathbb{E}_t \text{KL}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) || p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t))}_{\mathcal{L}_{\text{diffusion}}} + \mathcal{L}_{\text{prior}} \right]$$

If we take the number of time steps \mathbf{T} to infinity, we get the following bound:

$$\mathcal{L}_{\text{NELBO}} = -\mathbb{E}_q \int_{\alpha=0}^{\alpha=1} \frac{\log p_\theta(\mathbf{x} | \mathbf{z}_\alpha)}{1 - \alpha} d\alpha$$

Learning \mathbf{p} : Simplifying the ELBO

The (negative) ELBO for masked diffusion has (in general) the following form:

$$\mathcal{L}_{\text{NELBO}} = \mathbb{E}_q \left[\mathcal{L}_{\text{reconstr}} + \underbrace{\mathbb{E}_t \text{KL}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) || p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t))}_{\mathcal{L}_{\text{diffusion}}} + \mathcal{L}_{\text{prior}} \right]$$

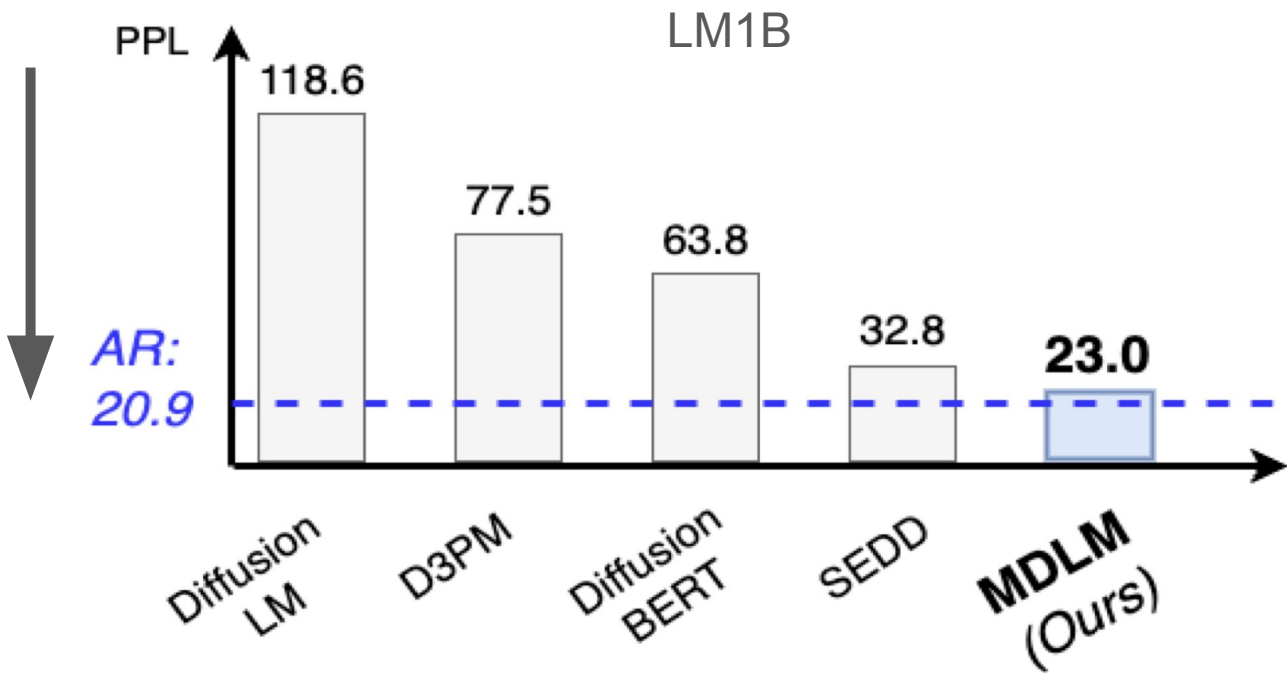
If we take the number of time steps \mathbf{T} to infinity, we get the following bound:

$$\mathcal{L}_{\text{NELBO}} = -\mathbb{E}_q \int_{\alpha=0}^{\alpha=1} \frac{\log p_\theta(\mathbf{x} | \mathbf{z}_\alpha)}{1 - \alpha} d\alpha$$

Observations

- The model is trained as a BERT model with a randomized masking rate
- However, the model admits principled sampling algorithms: generative BERT!
- This BERT defines probabilistic model p with an ELBO for evaluation

Experiments: Masked Diffusion Helps Close Gap To AR



Austin et al., "Structured Denoising Diffusion Models in Discrete State-Spaces." NeurIPS 2021
Li et al. "Diffusion-LM Improves Controllable Text Generation", NeurIPS 2022
Lou et al. "Discrete Diffusion Modeling by Estimating the Ratios of the Data Distribution", ICML 2024

Features of Discrete Diffusion

Parallel Generation

Autoregression

There

- ✗ Sequential
- ✗ Memory bound
- ✗ One token/step

Diffusion

Hirsch
account

- ✓ Parallel
- ✓ Compute bound
- ✓ Multiple tokens/step

Diffusion accelerates inference, producing multiple tokens/step.



Variable-Length Generation

Autoregression

There are three possible categories of the typical rate of return for companies that share...

- ✗ Sequential**
- ✗ Causal context only**
- ✓ Higher quality
- ✓ Arbitrary-length
- ✓ KV caching

Block Diffusion

On

- ✓ Parallel
- ✓ Semi-global context
- ✓ Higher quality
- ✓ Arbitrary-length
- ✓ KV caching

Diffusion

Hirsch will need to take account data that's to be released by tomorrow.

- ✓ Parallel
- ✓ Global context
- ✗ Lower quality**
- ✗ Fixed-length**
- ✗ No KV caching**

Iterative Refinement

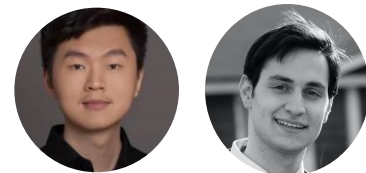


Autoregressive models:
once a token is produced,
it cannot be changed

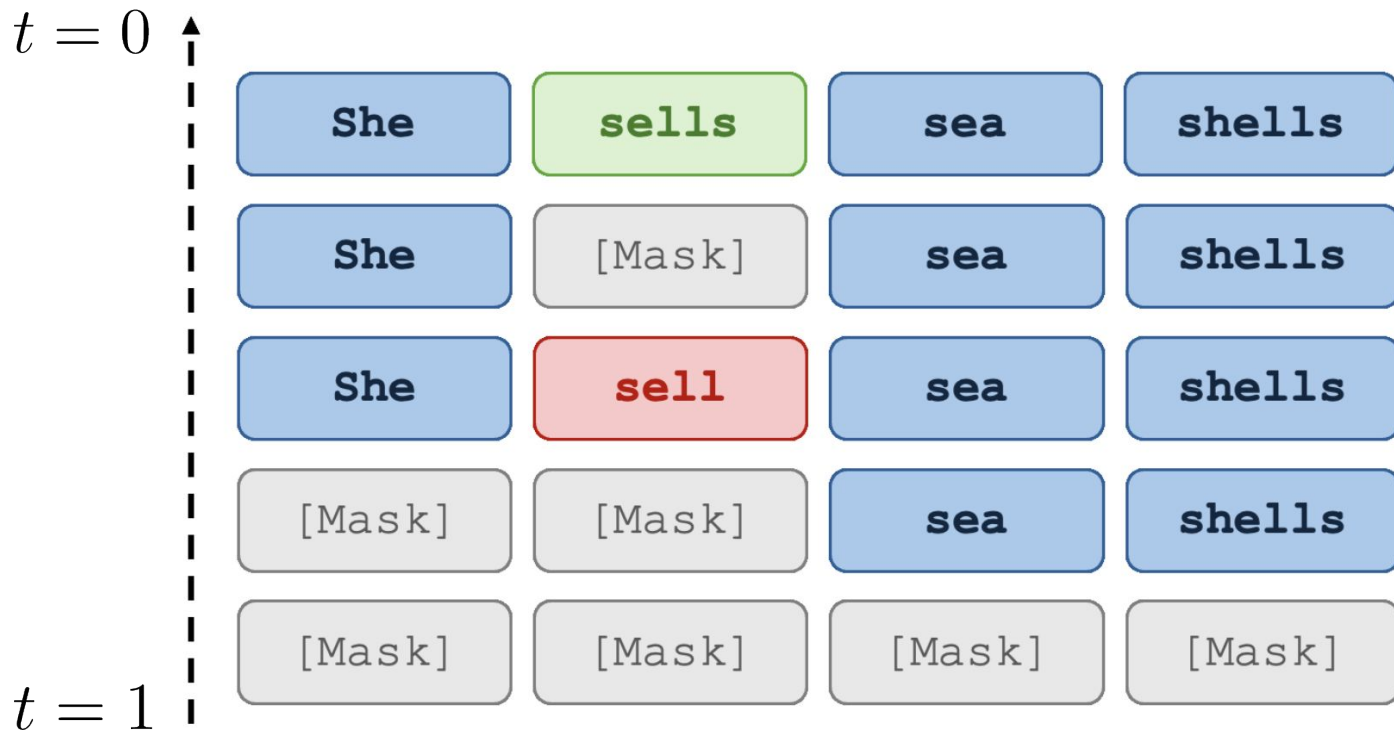


Diffusion models: can
perform **iterative error
correction**

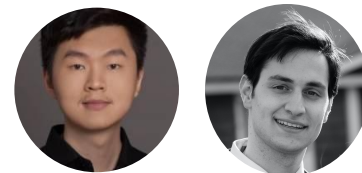




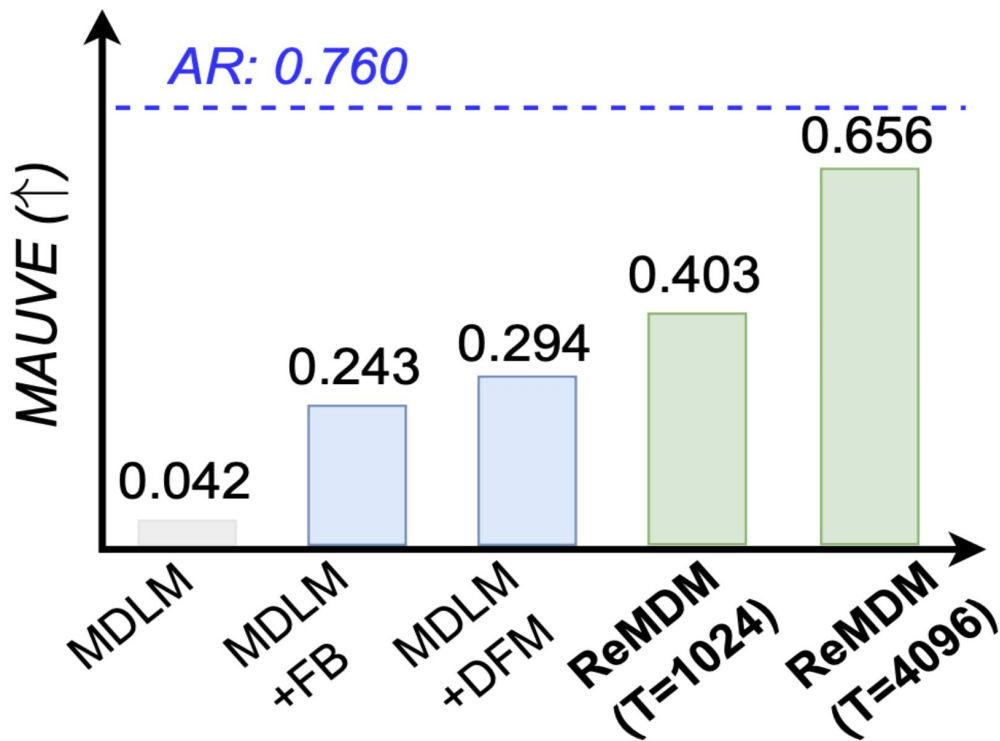
Iterative Refinement: ReMDM Sampler



Model produces
“sell” because it’s
a valid completion,
but later corrects
itself



Inference-Time Scaling

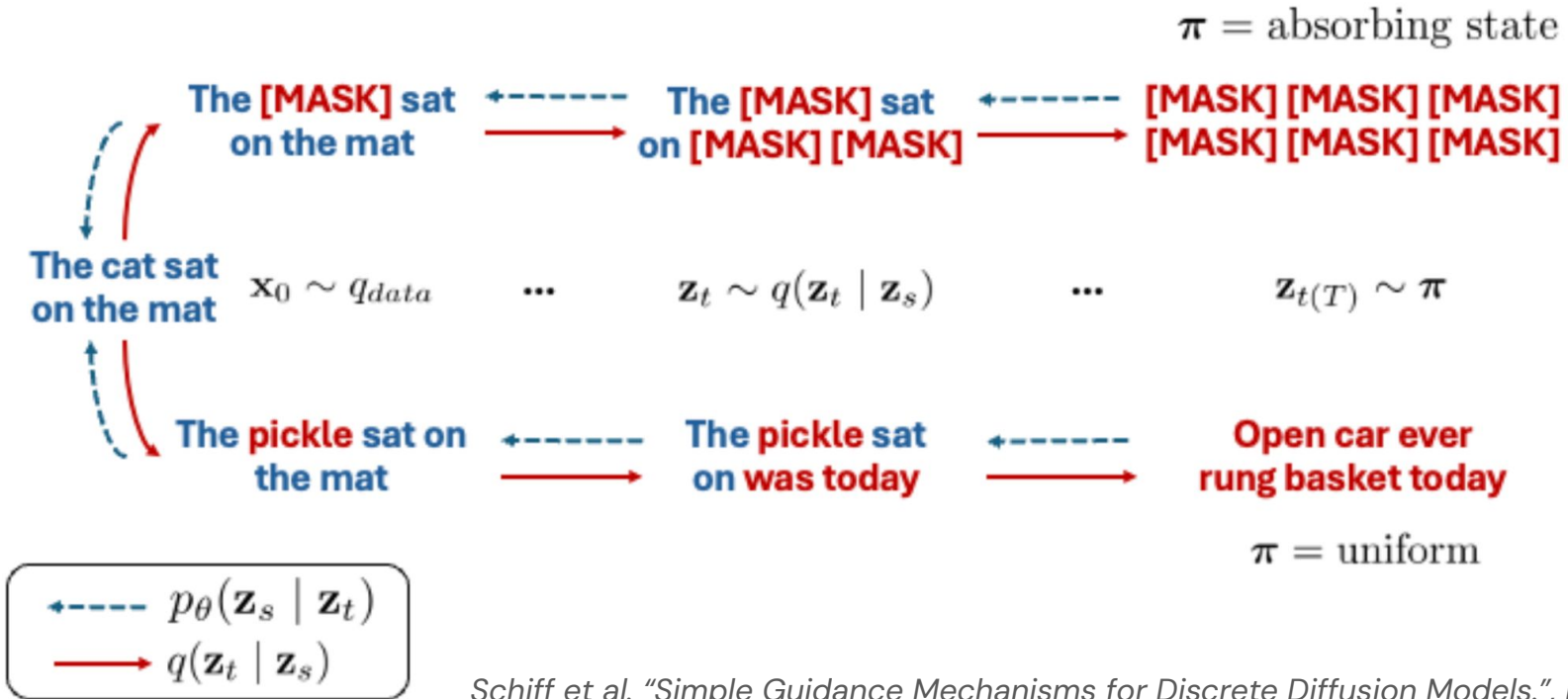


Diffusion modeling with inference-time scaling further closes the gap in sample quality to AR.



Uniform Noise Diffusion Language Models

An alternative to masked diffusion is to randomly swap words.





Improved Controllable & Guided Generation

X Autoregressive models:
make “**local**” predictions

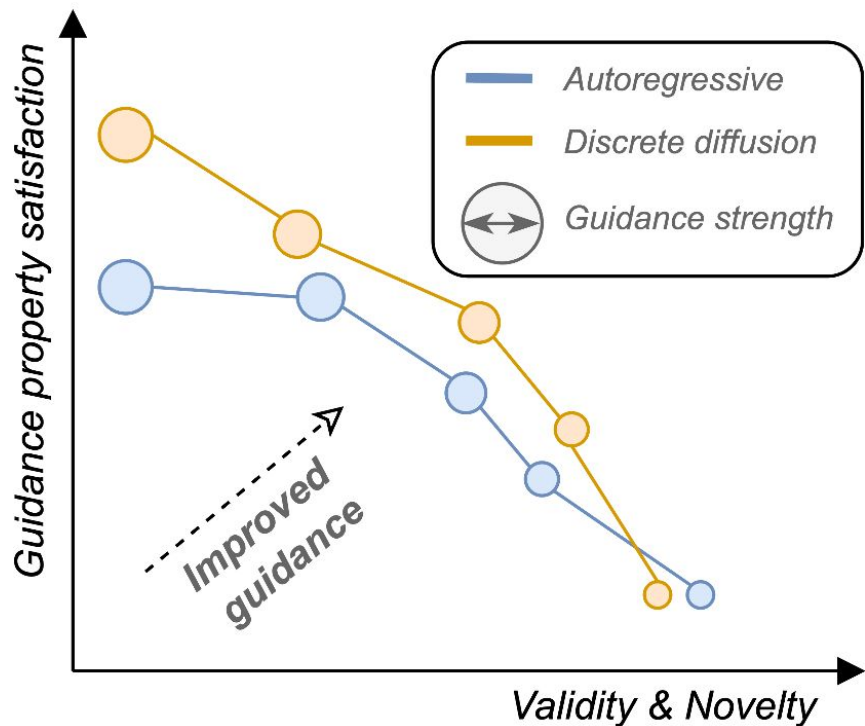
✓ Diffusion models: make
“**global**” (and **iterative**)
refinements



An ink sketch style illustration of a small hedgehog holding a piece of watermelon with its tiny paws, taking little bites with its eyes closed in delight.



Improved Controllable & Guided Generation



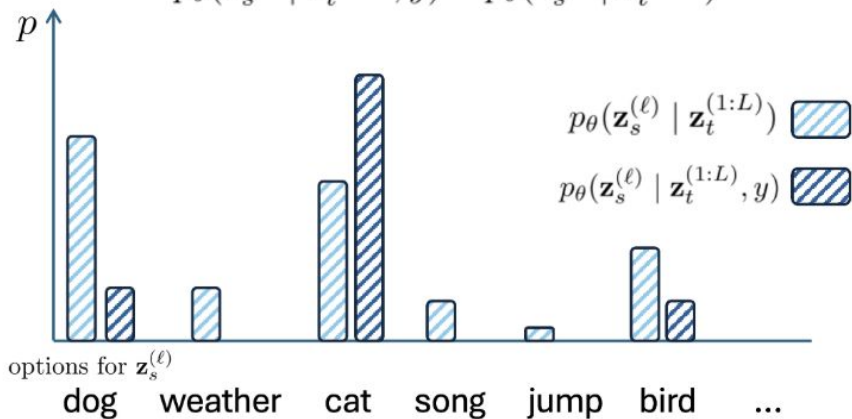
Consider generating \mathbf{x} to maximize property \mathbf{y} (e.g., binding affinity of drug).

Diffusion yields better Pareto frontier between realistic \mathbf{x} and high \mathbf{y} .

Suppose you need to update this word: **The pickle sat on the mat**

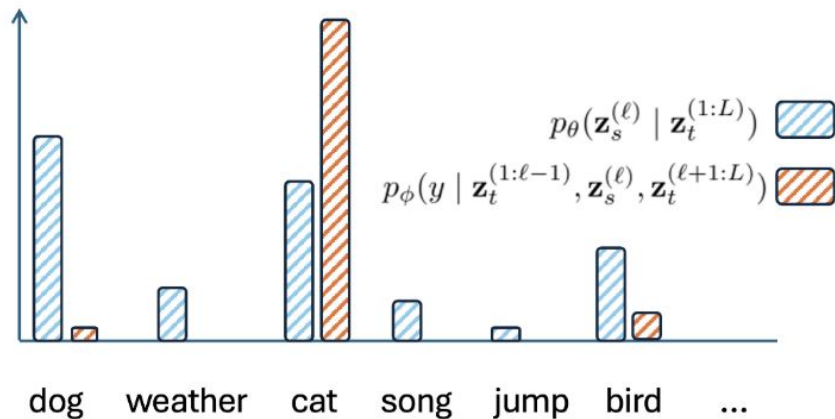
Discrete Classifier-Free Guidance

$$p_{\theta}^{\gamma}(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y) \propto p_{\theta}(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y)^{\gamma} \cdot p_{\theta}(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)})^{1-\gamma}$$



Discrete Classifier-Based Guidance

$$p_{\theta}^{\gamma}(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y) \propto p_{\phi}(y \mid \mathbf{z}_t^{(1:\ell-1)}, \mathbf{z}_s^{(\ell)}, \mathbf{z}_t^{(\ell+1:L)})^{\gamma} \cdot p_{\theta}(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)})$$



Classifier-Based and Classifier-Free Guidance

Suppose we have a classifier $p(y|\mathbf{z}_t)$ of property y from noisy \mathbf{z}_t .

$$\underbrace{\log p(\mathbf{z}_s|\mathbf{z}_t, y)}_{\text{conditional distribution}} = \underbrace{\log p(y|\mathbf{z}_t, \mathbf{z}_s)}_{\text{classifier}} + \underbrace{\log p(\mathbf{z}_s|\mathbf{z}_t)}_{\text{unconditional model}} + c$$

Classifier-Based and Classifier-Free Guidance

Suppose we have a classifier $p(y|z_t)$ of property y from noisy z_t .

$$\underbrace{\log p(z_s|z_t, y)}_{\text{conditional distribution}} = \underbrace{\log p(y|z_t, z_s)}_{\text{classifier}} + \underbrace{\log p(z_s|z_t)}_{\text{unconditional model}} + c$$

Discrete Classifier-Based Guidance (Schiff et al., 2024)

We can use this classifier to guide generation, and sample from

$$\underbrace{\log p^{(\gamma)}(z_s|z_t, y)}_{\text{new unnormalized distribution}} = \gamma \underbrace{\log p(y|z_t, z_s)}_{\text{guidance term}} + \underbrace{\log p(z_s|z_t)}_{\text{original diffusion model}}$$

where $\gamma > 0$ is a guidance strength parameter.

This extends to sequences if $p(z_s|z_t)$ is fully-factored across tokens.

Suppose we have a conditional model $p(\mathbf{z}_s|\mathbf{z}_t, y)$ and an unconditional $p(\mathbf{z}_s|\mathbf{z}_t)$. Recall the CBG factorization is:

$$\log p^{(\gamma)}(\mathbf{z}_s|\mathbf{z}_t, y) = \gamma \cdot \underbrace{\log p(y|\mathbf{z}_t, \mathbf{z}_s)}_{\text{apply Bayes rule}} + \log p(\mathbf{z}_s|\mathbf{z}_t) + c$$

By Bayes rule, $\log p(y|\mathbf{z}_t, \mathbf{z}_s) = \log p(\mathbf{z}_s|\mathbf{z}_t, y) - \log p(\mathbf{z}_s|\mathbf{z}_t) + c$.

Suppose we have a conditional model $p(\mathbf{z}_s|\mathbf{z}_t, y)$ and an unconditional $p(\mathbf{z}_s|\mathbf{z}_t)$. Recall the CBG factorization is:

$$\log p^{(\gamma)}(\mathbf{z}_s|\mathbf{z}_t, y) = \gamma \cdot \underbrace{\log p(y|\mathbf{z}_t, \mathbf{z}_s)}_{\text{apply Bayes rule}} + \log p(\mathbf{z}_s|\mathbf{z}_t) + c$$

By Bayes rule, $\log p(y|\mathbf{z}_t, \mathbf{z}_s) = \log p(\mathbf{z}_s|\mathbf{z}_t, y) - \log p(\mathbf{z}_s|\mathbf{z}_t) + c$.

Discrete Classifier-Free Guidance (Schiff et al., 2024)

We can use this classifier to guide generation, and sample from

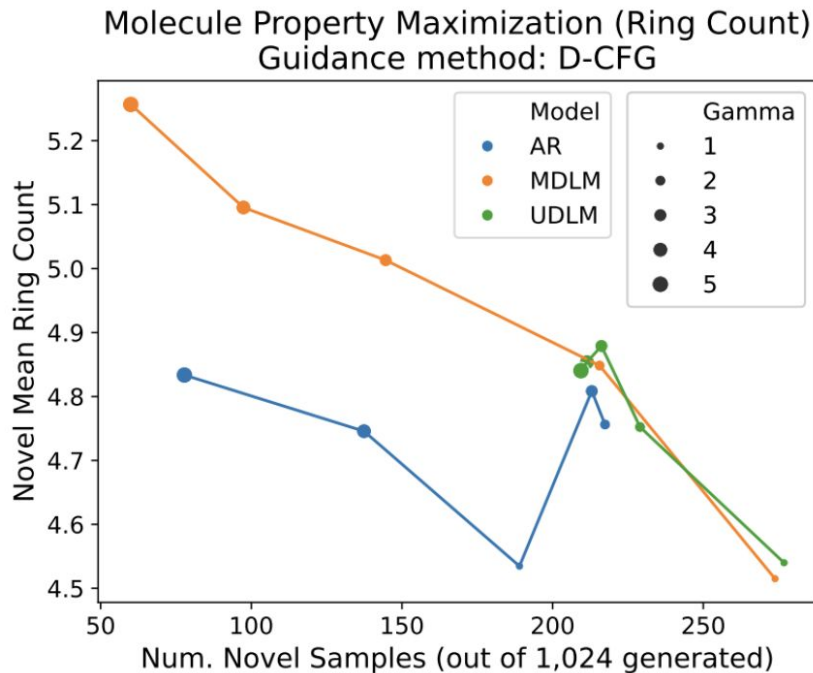
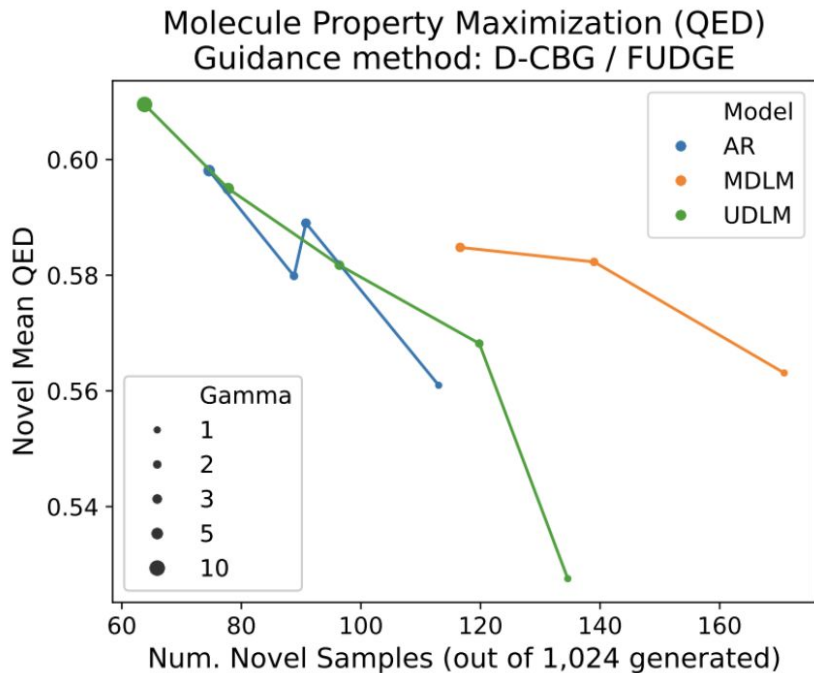
$$\underbrace{\log p^{(\gamma)}(\mathbf{z}_s|\mathbf{z}_t, y)}_{\text{new unnormalized distribution}} = \gamma \underbrace{\log p(\mathbf{z}_s|\mathbf{z}_t, y)}_{\text{conditional model}} + (1-\gamma) \underbrace{\log p(\mathbf{z}_s|\mathbf{z}_t)}_{\text{unconditional model}}$$

where $\gamma > 0$ is a guidance strength parameter.

In practice, we parameterize $p(\mathbf{z}_s|\mathbf{z}_t, y), p(\mathbf{z}_s|\mathbf{z}_t)$ with same model.

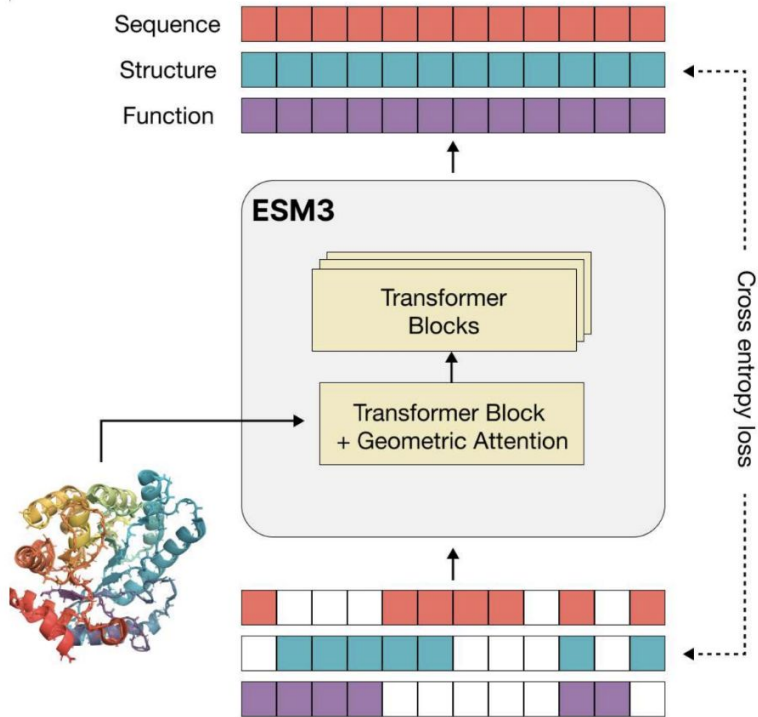


Improved Controllable & Guided Generation



Large Diffusion Language Models

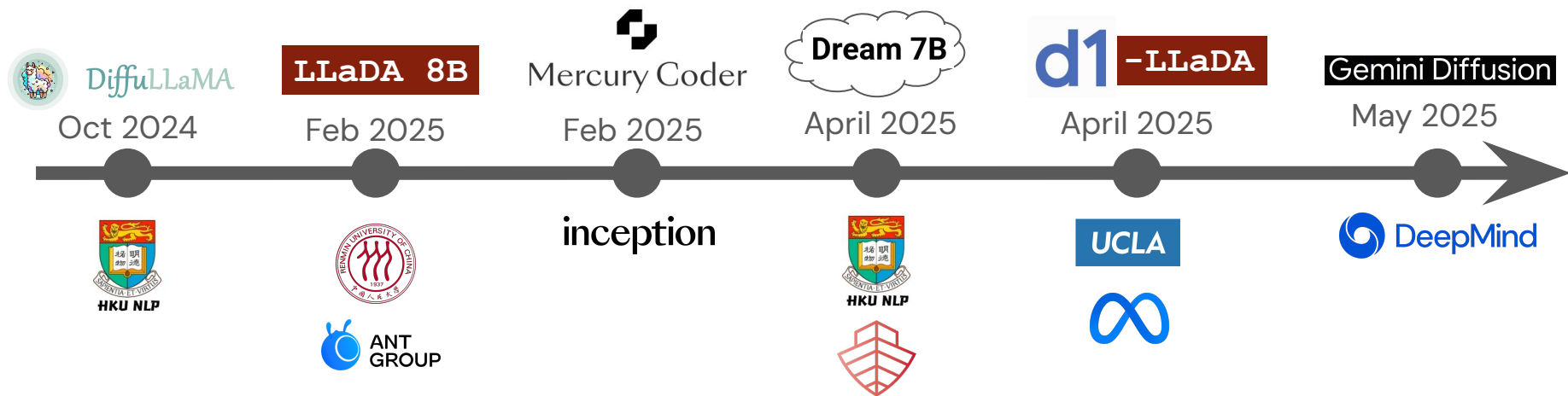
Protein Language Models are Discrete Diffusion Models



Recent ESM3 100B protein language models use masked diffusion for training and sampling.

Protein LMs can be used to accelerating folding

Large Diffusion Language Models



Diffusion LLMs Exhibit Favorable Scaling Laws Relative to AR

- Both uniform and masked diffusion LLMs exhibit similar scaling behavior to AR LLMs
- However, the optimal model sizes of dLLMs are larger

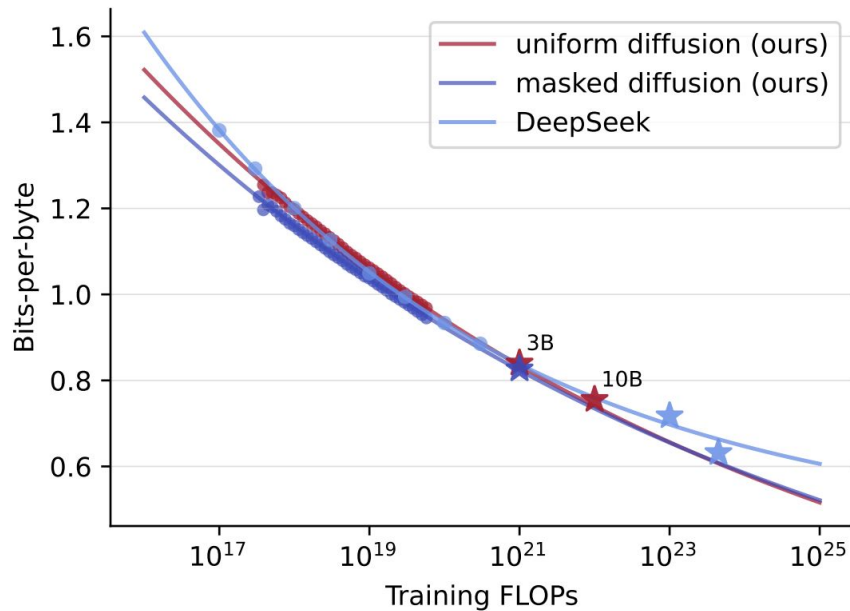
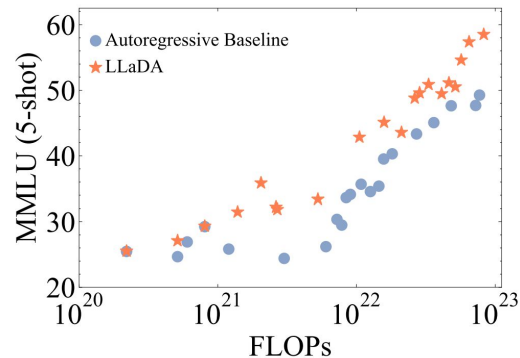
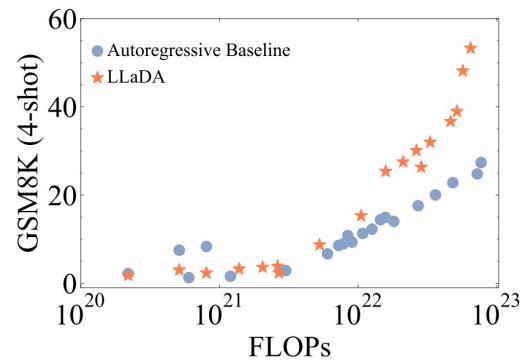
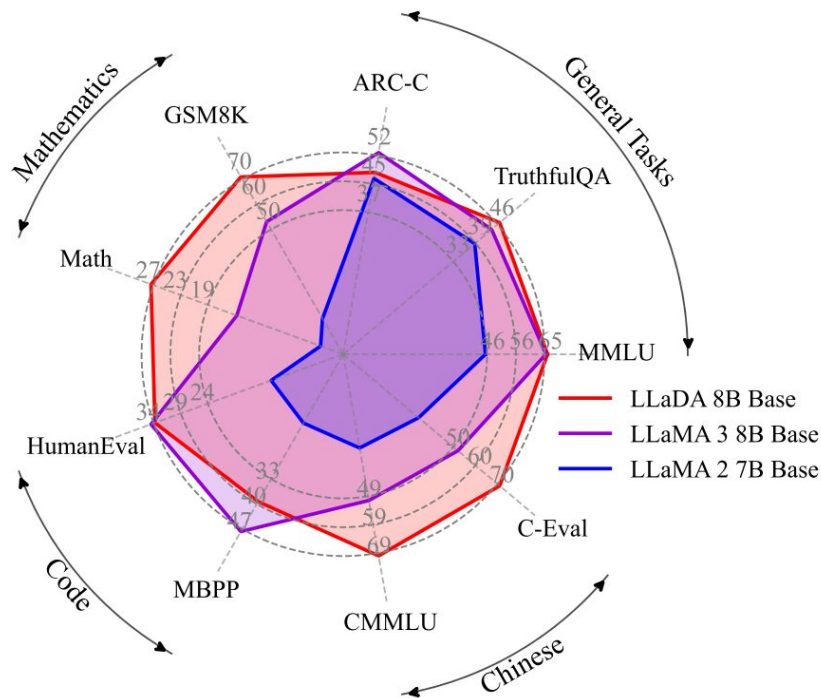


Figure 1: Our proposed scaling laws extrapolate well to 3B and 10B models trained on up to $50\times$ larger compute budgets and suggest that DLMs can be competitive with ALMs at scale, even in compute-bound training settings.

LLaDA: An 8B Open-Weights Diffusion LLM





Iterations

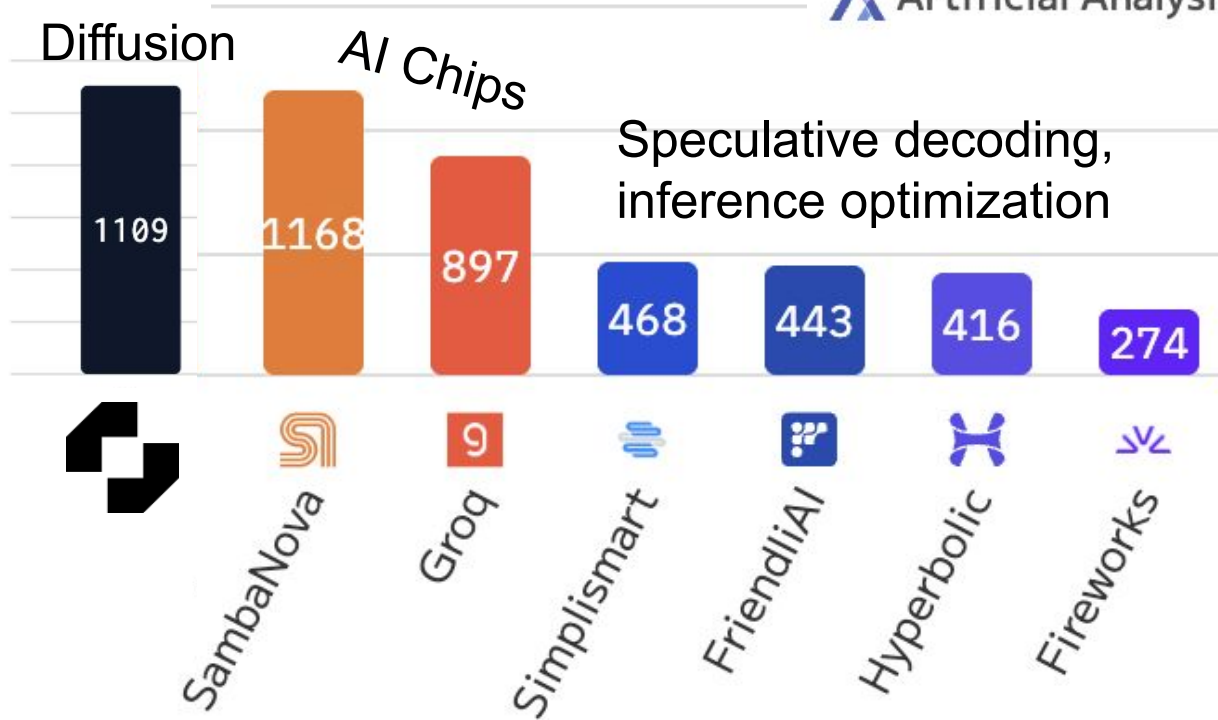
0

AUTOREGRESSIVE LLM
LEFT-TO-RIGHT GENERATION

Iterations

0

INCEPTION DIFFUSION LLM
COARSE-TO-FINE GENERATION



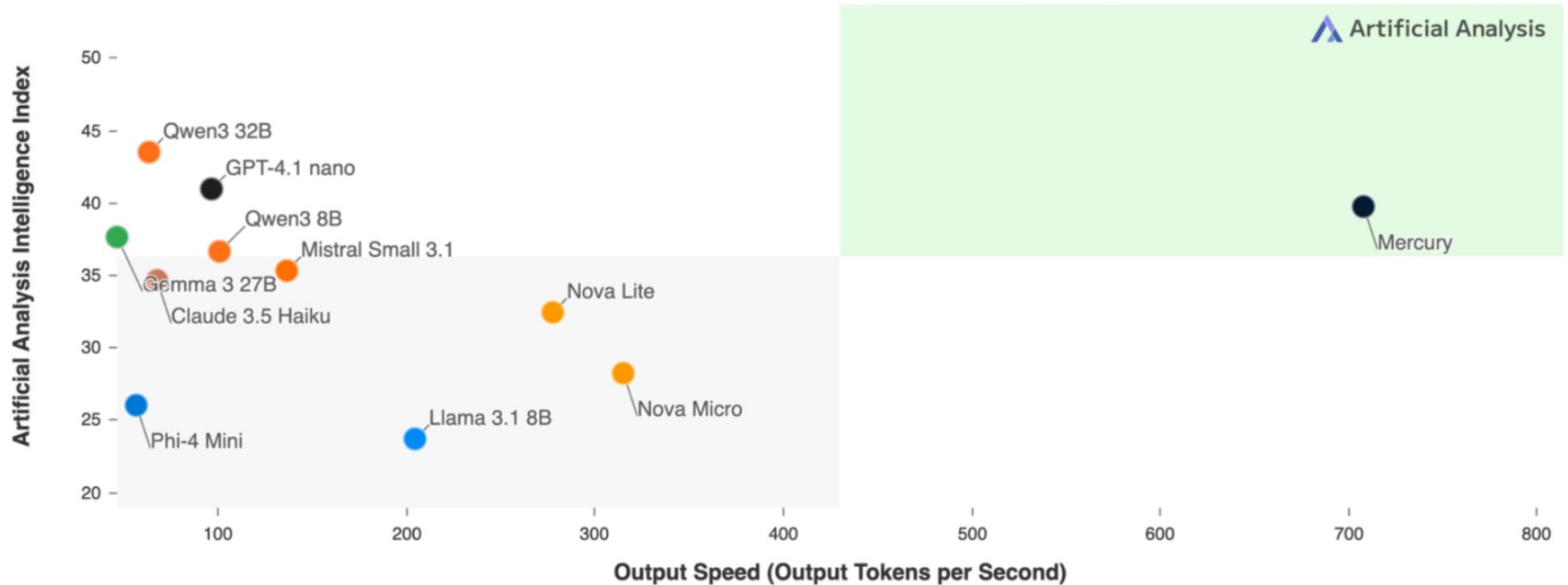
Throughput of Llama 8b vs. Mercury Mini (higher quality) in tokens/second across inference providers

Intelligence vs. Output Speed: Small, Non-reasoning models

Artificial Analysis Intelligence Index; Output Speed: Output Tokens per Second

Most attractive quadrant

- GPT-4.1 nano
- Gemma 3 27B
- Claude 3.5 Haiku
- Mistral Small 3.1
- Nova Lite
- Nova Micro
- Phi-4 Mini
- Mercury
- Qwen3 32B
- Qwen3 8B
- Llama 3.1 8B
- Gemini 2.5 Flash



Conclusion

- Masked diffusion is emerging as an alternative to autoregressive models.
- Advantages of diffusion include:
 - **Fast parallel generation:** much greater tok/sec at inference time
 - **Built-in error correction:** improves quality via iterative refinement with bidirectional context
 - **Controllability:** easier to guide model because of error correction
 - **New form of inference-time scaling:** can denoise more or less at inference time
 - **Multimodality:** one paradigm for all modalities
- Large-scale diffusion language models are already here!

Simple Diffusion Language Models

