

Diffusion Models are Evolutionary Algorithms

Yanbo Zhang¹, Benedikt Hartl^{1,2}, Hananel Hazan¹, Michael Levin^{1,3}

¹Allen Discovery Center at Tufts University

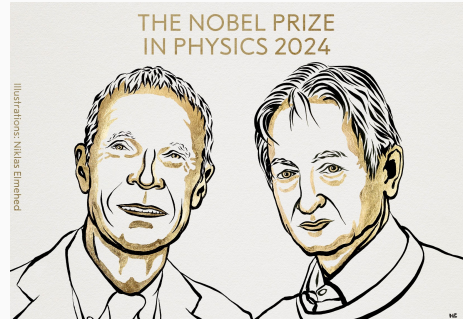
²Institute for Theoretical Physics, TU Wien, Austria

³Wyss Institute for Biologically Inspired Engineering at Harvard University

Intelligence and Learning

What is Intelligence?

- Ilya: Intelligence = compression
- Geoffrey Hinton: Intelligence = learning
 - Adaptive, reasoning, ...
 - self-organize



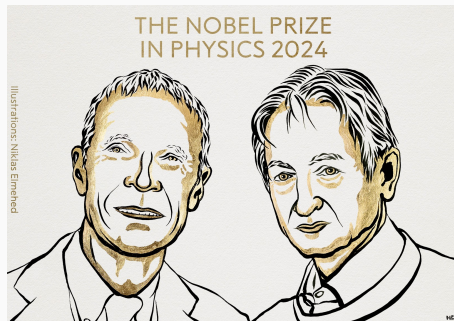
Intelligence of Evolution

- Adapting
- Immune systems, neurons, ...
- Evolution itself can be seen as a kind of intelligence

Intelligence and Learning

What is Intelligence?

- Ilya: Intelligence = compression
- Geoffrey Hinton: Intelligence = learning
 - Adaptive, reasoning, ...
 - self-organize



Evolution of Intelligence

- Machine Intelligence: diffusion model
- Iterative: denoise, optimize, adding noise

Diffusion Model

Quick review of Diffusion Models:

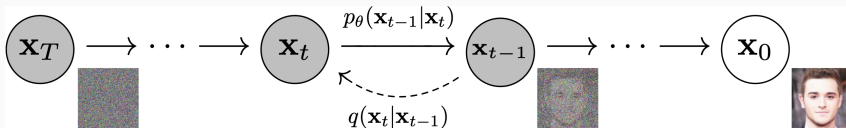


Figure 1: Example of Diffusion Models on image generation.

Diffusion and Training

Train a model to predict the added noise given \mathbf{x}_t [1, 2].

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\epsilon,$$

and

$$\mathcal{L} = \mathbb{E}_{\mathbf{x}_0 \sim p_{\text{data}}, \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_\theta(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\epsilon) - \epsilon\|^2$$

Diffusion Model

Quick review of Diffusion Models:

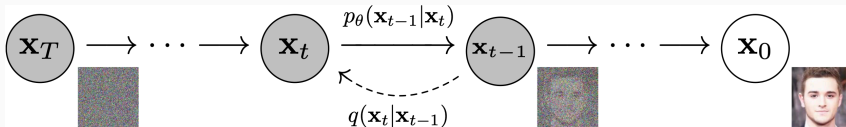


Figure 1: Example of Diffusion Models on image generation.

Denosing

Use the trained model to denoise from Gaussian distribution [2]:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \epsilon_\theta(\mathbf{x}_t, t) + \sigma_t \mathbf{w}$$

Similarity between evolution and diffusion

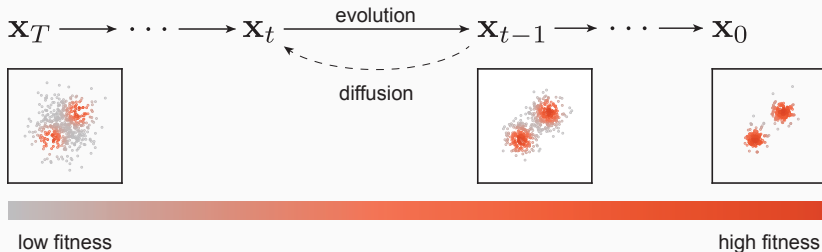
Diffusion and Evolution are both:

- directed: natural selection v.s. denoise
- randomness: mutation v.s. noise term
- optimization process

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\hat{\mathbf{x}}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2}\hat{\mathbf{e}} + \sigma_t\mathbf{w}$$

Key idea

diffusion as reversed evolution, and denoising as evolution.



Connecting Diffusion to Evolution

Three problems need to be solved:

Mapping fitness to probability density

Higher fitness $f(x)$ corresponds to higher density $p(x)$, which requires a mapping function $g : \mathbb{R} \rightarrow \mathbb{R}^+$:

$$p(\mathbf{x}) = g[f(\mathbf{x})].$$

Connecting Diffusion to Evolution

Three problems need to be solved:

Mapping fitness to probability density

Higher fitness $f(x)$ corresponds to higher density $p(x)$, which requires a mapping function $g : \mathbb{R} \rightarrow \mathbb{R}^+$:

$$p(\mathbf{x}) = g[f(\mathbf{x})].$$

Denoise Model[†]

Diffusion Models have $\epsilon(x_t, t)$ to predict noise, evolutionary algorithm also need a predictive model:

$$p(\epsilon|x_t), \text{ etc.}$$

Connecting Diffusion to Evolution

Three problems need to be solved:

Mapping fitness to probability density

Higher fitness $f(x)$ corresponds to higher density $p(x)$, which requires a mapping function $g : \mathbb{R} \rightarrow \mathbb{R}^+$:

$$p(\mathbf{x}) = g[f(\mathbf{x})].$$

Denoise Model[†]

Diffusion Models have $\epsilon(x_t, t)$ to predict noise, evolutionary algorithm also need a predictive model:

$$p(\epsilon|x_t), etc.$$

Determine iteration method

Directly use iteration method from diffusion models.

Diffusion models are doing early prediction

Most common understanding: diffusion models are trained to predict added noise: $p(\epsilon|x_t)$.

It is hard to analyze in this perspective. However, they are also predicting the origin directly:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\epsilon \iff \hat{\mathbf{x}}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \alpha_t}\epsilon}{\sqrt{\alpha_t}}$$

Hence,

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\hat{\mathbf{x}}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \frac{\mathbf{x}_t - \sqrt{\alpha_t}\hat{\mathbf{x}}_0}{\sqrt{1 - \alpha_t}} + \sigma_t \mathbf{w}$$

Early prediction $p(x_0 = x|x_t)$

Diffusion models are making **early predictions**, then move toward it by small steps. The model is also predicting $p(x_0 = x|x_t)$.

Estimating high-fitness parameter

Apply Bayesian equation on early prediction:

$$p(\mathbf{x}_0 = \mathbf{x} | \mathbf{x}_t) = \frac{\overbrace{p(\mathbf{x}_t | \mathbf{x}_0 = \mathbf{x})}^{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}, 1 - \alpha_t)} \overbrace{p(\mathbf{x}_0 = \mathbf{x})}^{g[f(\mathbf{x})]}}{p(\mathbf{x}_t)}$$

Diffusion models are trained with MSE loss \rightarrow We need average over \mathbf{x} :

$$\hat{\mathbf{x}}_0(\mathbf{x}_t, \alpha, t) = \sum_{\mathbf{x} \sim p_{\text{eval}}(\mathbf{x})} p(\mathbf{x}_0 = \mathbf{x} | \mathbf{x}_t) \mathbf{x} = \frac{1}{Z} \sum_{\mathbf{x} \in \mathbf{X}_t} g[f(\mathbf{x})] \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}, 1 - \alpha_t) \mathbf{x},$$

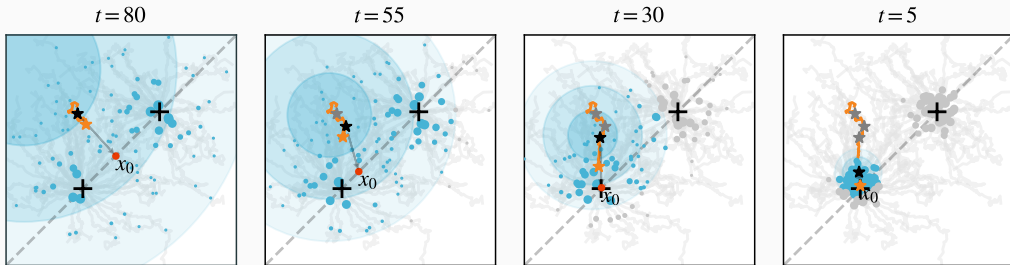
Take this back to diffusion models sampler, we got **Diffusion Evolution!**

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \hat{\mathbf{x}}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \frac{\mathbf{x}_t - \sqrt{\alpha_t} \hat{\mathbf{x}}_0}{\sqrt{1 - \alpha_t}} + \sigma_t \mathbf{w}$$

Diffusion Evolution Algorithm

An example:

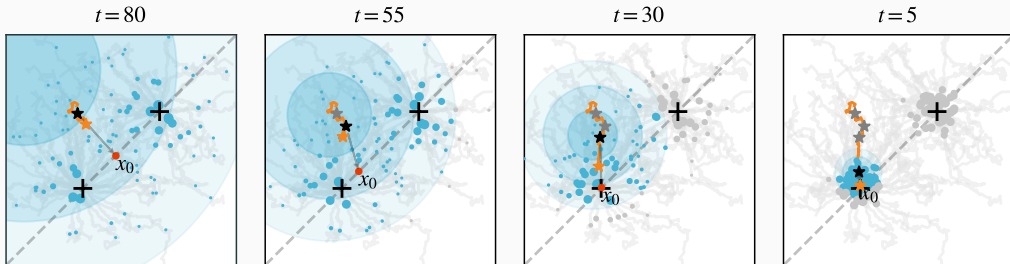
- 1 Random initialized population $\{\mathbf{x}_T^{(i)}\} \sim \mathcal{N}(0, I)$
- 2 Evaluate fitness and compute density $g[f(\mathbf{x})]$
- 3 Estimate $\hat{\mathbf{x}}_0^{(i)}$ for each individual
- 4 Move toward the $\hat{\mathbf{x}}_0^{(i)}$ plus mutation noise $\sigma_t \mathbf{w}$



Selection, Mutation, and Reproductive Isolation

$$\mathbf{x}_{t-1} = \underbrace{\sqrt{\alpha_{t-1}}\hat{\mathbf{x}}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \frac{\mathbf{x}_t - \sqrt{\alpha_t}\hat{\mathbf{x}}_0}{\sqrt{1 - \alpha_t}}}_{\text{directed evolution}} + \underbrace{\sigma_t \mathbf{w}}_{\text{undirected mutation}}$$

$$\hat{\mathbf{x}}_0(\mathbf{x}_t, \alpha, t) = \frac{1}{Z} \sum_{\mathbf{x} \in \mathbf{X}_t} \underbrace{g[f(\mathbf{x})]}_{\text{selection}} \underbrace{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}, 1 - \alpha_t)}_{\text{reproductive isolation}} \mathbf{x}$$



Mapping Between Diffusion and Evolution

Diffusion		Evolution
MES loss	\leftrightarrow	Weighted average on $\hat{\mathbf{x}}_0$
$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}$	\leftrightarrow	Gaussian weight as “reproductive isolation”
DDPM/DDIM sampling	\leftrightarrow	Evolution iteration

Table 1: Diffusion models can be decomposed into three parts, associate to three parts in evolutionary algorithm respectively.

Experiments: parallel between
diffusion and evolution

Latent Space Diffusion Evolution

Original Diffusion Evolution struggles on high-dimensional space:

$$\hat{x}_0(x_t, \alpha, t) = \frac{1}{Z} \sum_{x \in X_t} g[f(x)] \underbrace{\mathcal{N}(x_t; \sqrt{\alpha_t}x, 1 - \alpha_t)}_{\text{become more local at high-dimension}} x$$

Latent Space Diffusion Evolution

Evolve at lower-dimensional space may resolve this problem.



Diffusion Evolution on high-dimensional space

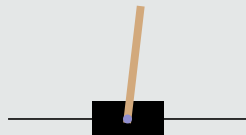
We train neural networks to solve Cart-pole system by Latent Space Diffusion Evolution.

Cart-pole system

Using the cart's survival time as the reward, the cart ends when:

- If cart exceed certain range;
- If pole falls down

Four inputs $(x, \theta, \dot{x}, \dot{\theta})$ and two possible control signals (move left and right).



Diffusion Evolution on high-dimensional space

We train neural networks to solve Cart-pole system by Latent Space Diffusion Evolution.

Neural network controller

- **Simple version:** Three layers, (4, 8, 2), with 58 neurons;
- **Complex version:** Four layers, (4, 128, 128, 2), with 17,410 neurons.

Both use ReLU activation.

Diffusion Evolution on high-dimensional space

We train neural networks to solve Cart-pole system by Latent Space Diffusion Evolution.

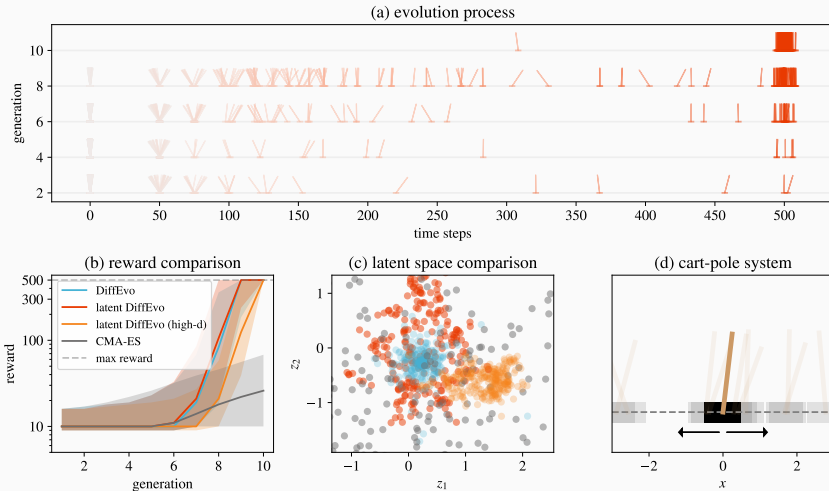


Table 2: QD-scores

	Diffusion Evolution	Latent Diffusion Evolution	CMA-ES	PEPG	Open-ES	MAP-Elite
Rosenbrock	<u>35.4</u> (1.82)	<u>23.4</u> (9.78)	1.00 (0.00)	1.25 (0.43)	0.73 (0.12)	42.0 (0.36)
Beale	<u>20.0</u> (0.94)	<u>13.0</u> (4.50)	1.00 (0.00)	2.00 (0.00)	1.84 (0.72)	23.7 (0.48)
Himmelblau	<u>15.8</u> (1.18)	<u>11.4</u> (3.99)	1.00 (0.00)	1.00 (0.00)	1.00 (0.00)	18.1 (0.42)
Ackley	<u>28.3</u> (1.35)	<u>16.7</u> (9.07)	13.4 (4.53)	1.13 (0.34)	2.14 (0.63)	33.0 (0.53)
Rastrigin ²	<u>35.2</u> (2.18)	<u>17.8</u> (9.36)	18.9 (6.24)	2.30 (0.64)	3.99 (0.00)	45.2 (0.93)
Rastrigin ⁴	<u>10.6</u> (0.50)	33.1 (6.66)	6.03 (2.30)	0.67 (0.12)	1.01 (0.00)	<u>13.5</u> (0.22)
Rastrigin ³²	<u>0.96</u> (0.04)	73.4 (2.10)	0.12 (0.03)	0.06 (0.01)	0.07 (0.01)	<u>1.20</u> (0.02)
Rastrigin ²⁵⁶	0.11 (0.01)	70.2 (0.62)	0.01 (0.00)	0.01 (0.00)	0.00 (0.00)	<u>0.14</u> (0.00)

Source code available:

<https://github.com/Zhangyanbo/diffusion-evolution>

ICLR 2025:

<https://openreview.net/forum?id=xVefsBbG20>

email: Yanbo.Zhang@tufts.edu



Yanbo
Zhang





Benedikt
Hartl



Hananel
Hazan



Michael
Levin

-  J. Ho, A. Jain, and P. Abbeel.
Denoising diffusion probabilistic models.
Advances in neural information processing systems, 33:6840–6851, 2020.
-  J. Song, C. Meng, and S. Ermon.
Denoising diffusion implicit models.
arXiv preprint arXiv:2010.02502, 2020.

Appendix I: Algorithm

Require: Population size N , parameter dimension D , fitness function f , density mapping function g , total evolution steps T , diffusion schedule α and noise schedule σ .

- 1: $[\mathbf{x}_T^{(1)}, \mathbf{x}_T^{(2)}, \dots, \mathbf{x}_T^{(N)}] \leftarrow \mathcal{N}(0, I^{N \times D})$ ▷ Initialize population
- 2: **for** $t \in [T, T - 1, \dots, 2]$ **do**
- 3: $\forall i \in [1, N] : Q_i \leftarrow g[f(\mathbf{x}_t^{(i)})]$ ▷ Fitness are cached to avoid repeated evaluations
- 4: **for** $i \in [1, 2, \dots, N]$ **do**
- 5: $Z \leftarrow \sum_{j=1}^N Q_j \mathcal{N}(\mathbf{x}_t^{(i)}; \sqrt{\alpha_t} \mathbf{x}_t^{(j)}, 1 - \alpha_t)$
- 6: $\hat{\mathbf{x}}_0 \leftarrow \frac{1}{Z} \sum_{j=1}^N Q_j \mathcal{N}(\mathbf{x}_t^{(i)}; \sqrt{\alpha_t} \mathbf{x}_t^{(j)}, 1 - \alpha_t) \mathbf{x}_t^{(j)}$
- 7: $\mathbf{w} \leftarrow \mathcal{N}(0, I^D)$
- 8: $\mathbf{x}_{t-1}^{(i)} \leftarrow \sqrt{\alpha_{t-1}} \hat{\mathbf{x}}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \frac{\mathbf{x}_t^{(i)} - \sqrt{\alpha_t} \hat{\mathbf{x}}_0}{\sqrt{1 - \alpha_t}} + \sigma_t \mathbf{w}$
- 9: **end for**
- 10: **end for**