

Discretized Approximate Ancestral Sampling

Jona Ballé
NYU Tandon School of Engineering

ISIT Learn to Compress & Compress to
Learn Workshop
25 June 2025



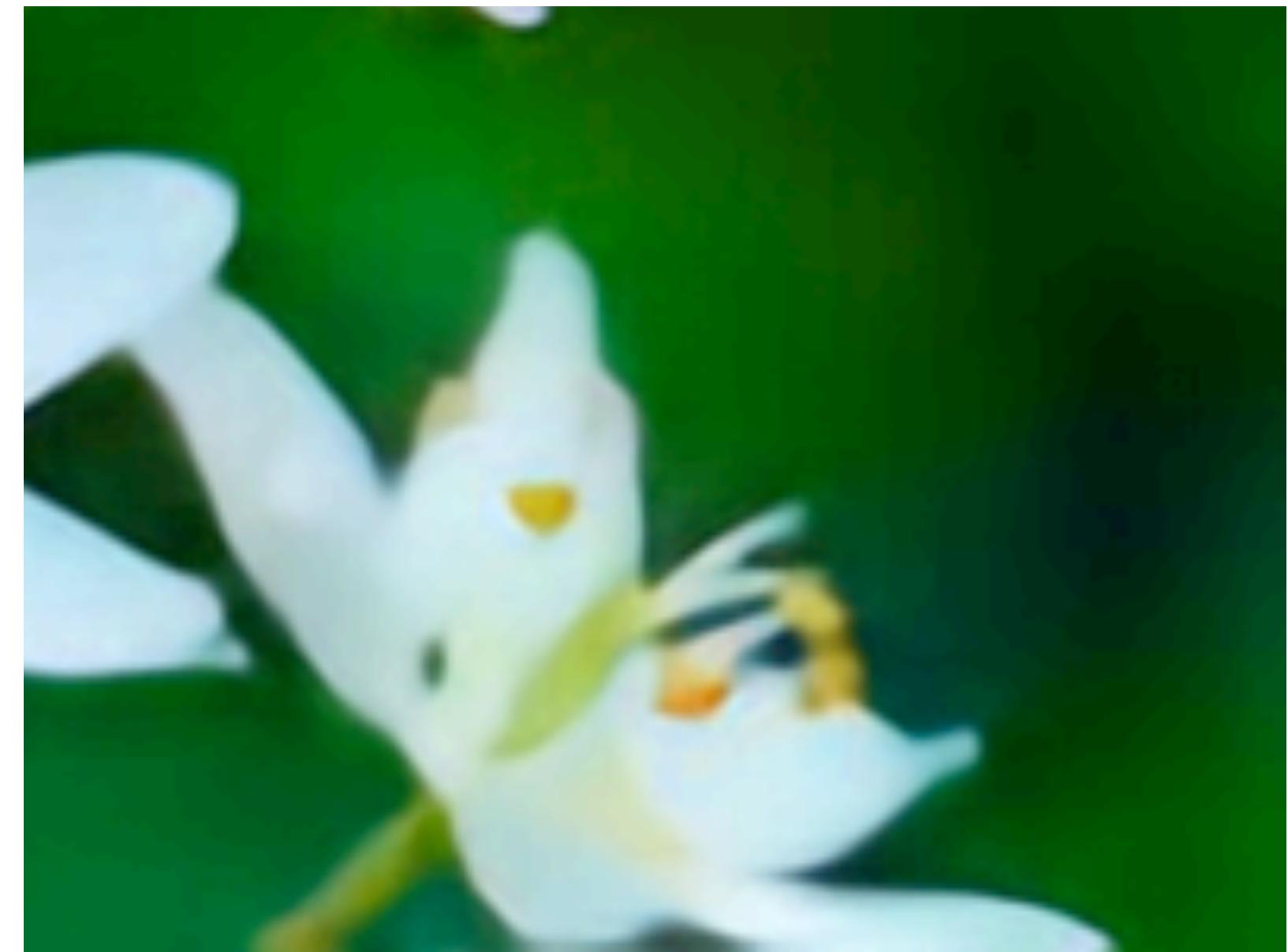
Alfredo De la Fuente
Google



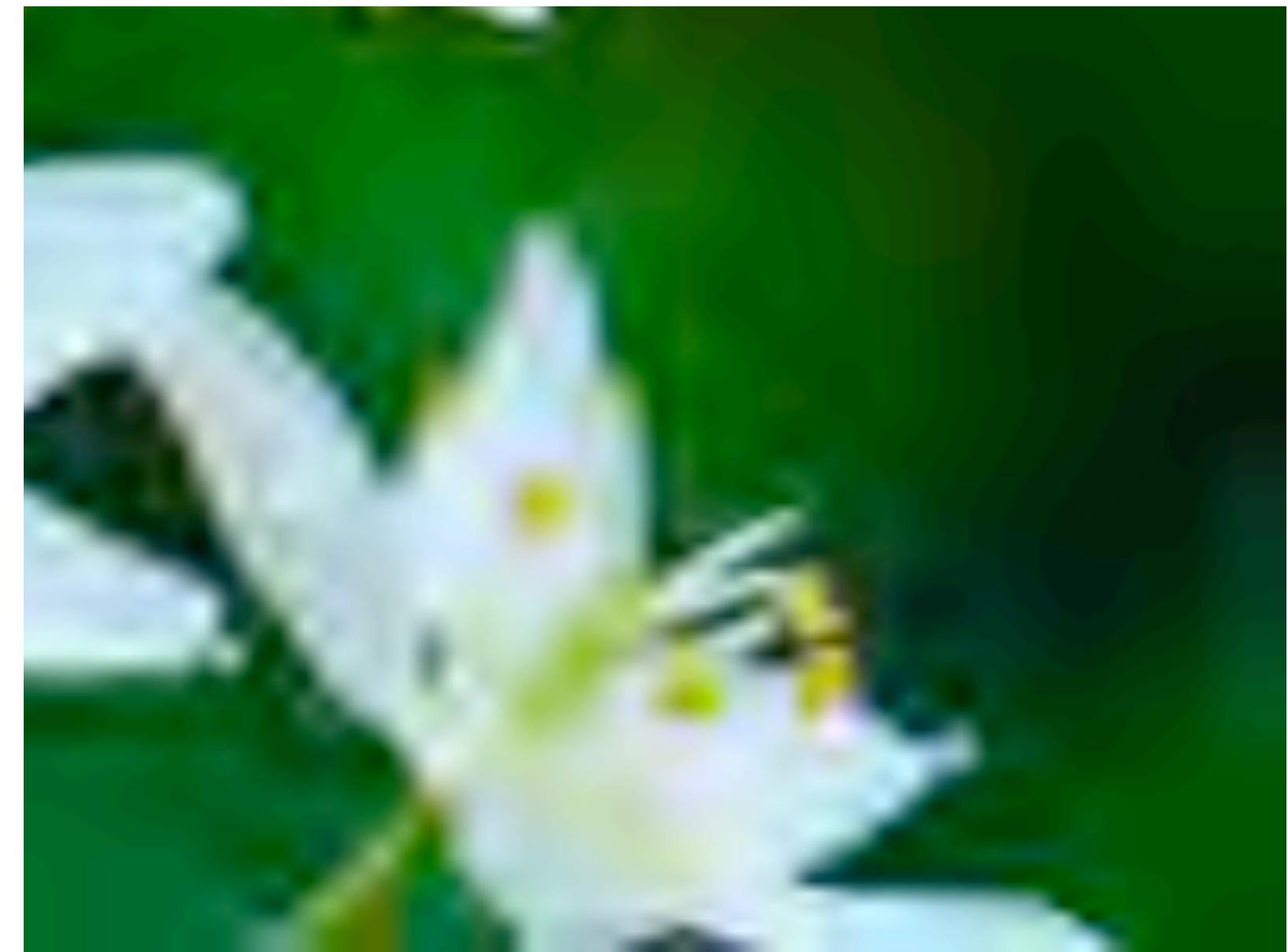
Saurabh Singh
Google DeepMind



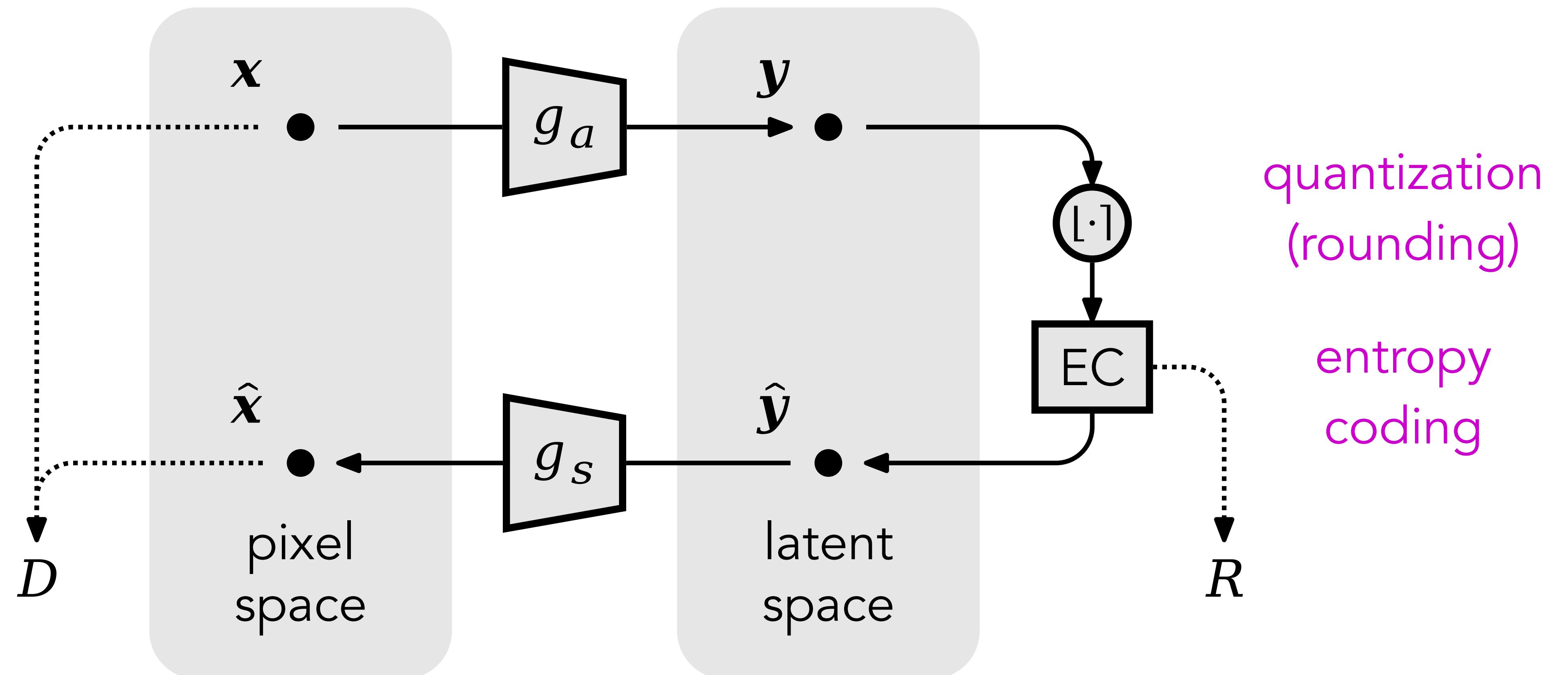
original NTC (2017)



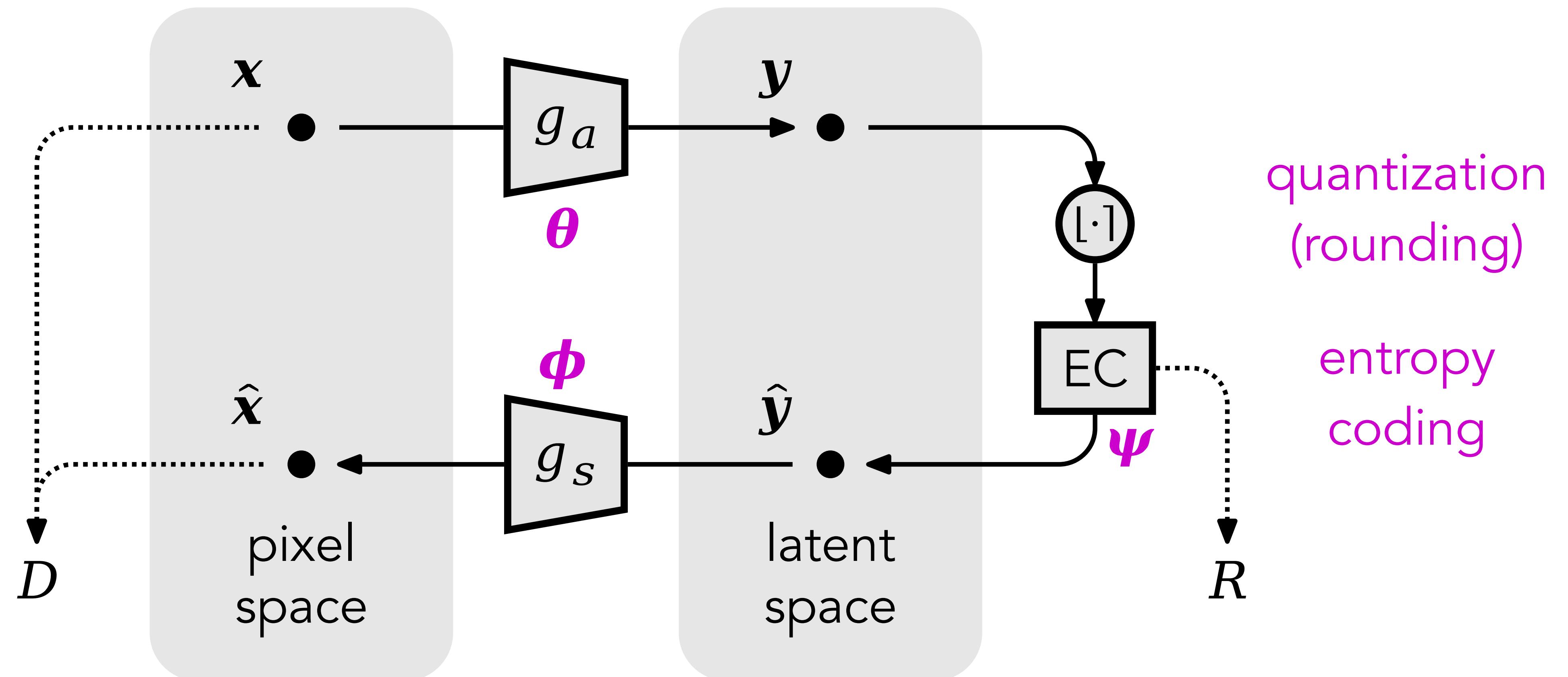
JPEG JPEG 2000



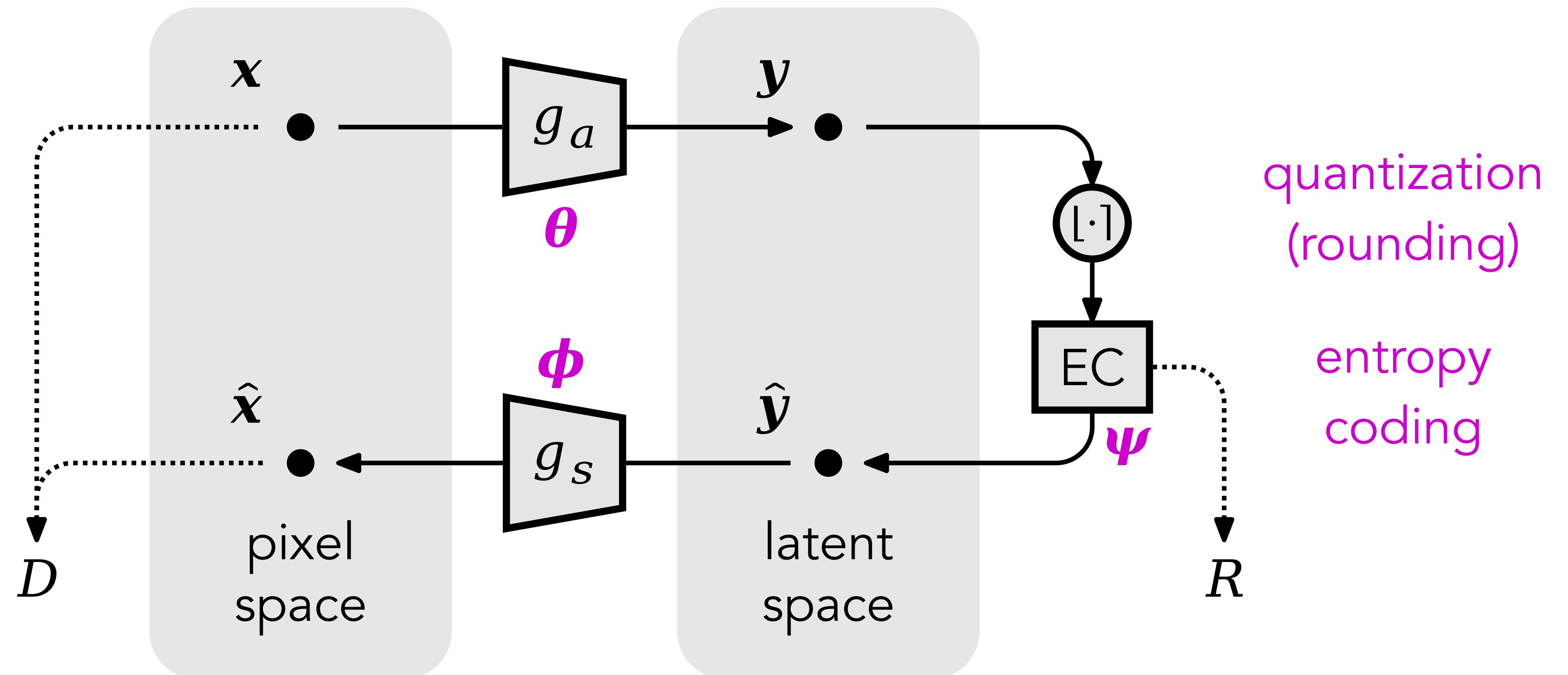
Nonlinear transform coding (end-to-end optimized)



Nonlinear transform coding (end-to-end optimized)

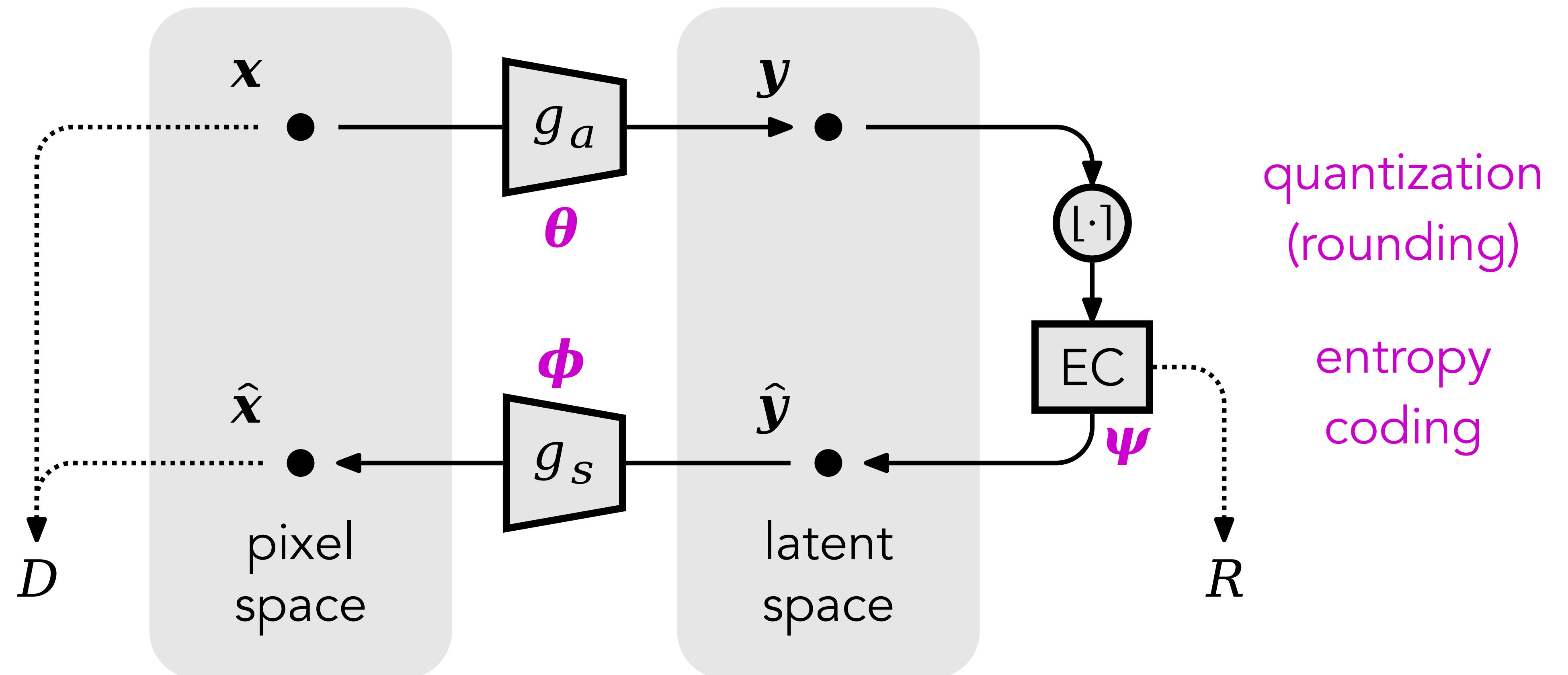


Nonlinear transform coding (end-to-end optimized)



$$L(\theta, \phi, \psi) = \underbrace{\mathbb{E}_{\mathbf{x}}[-\log_2 p_{\hat{\mathbf{y}}}(\hat{\mathbf{y}})]}_R + \lambda \underbrace{\mathbb{E}_{\mathbf{x}}[\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2]}_D$$

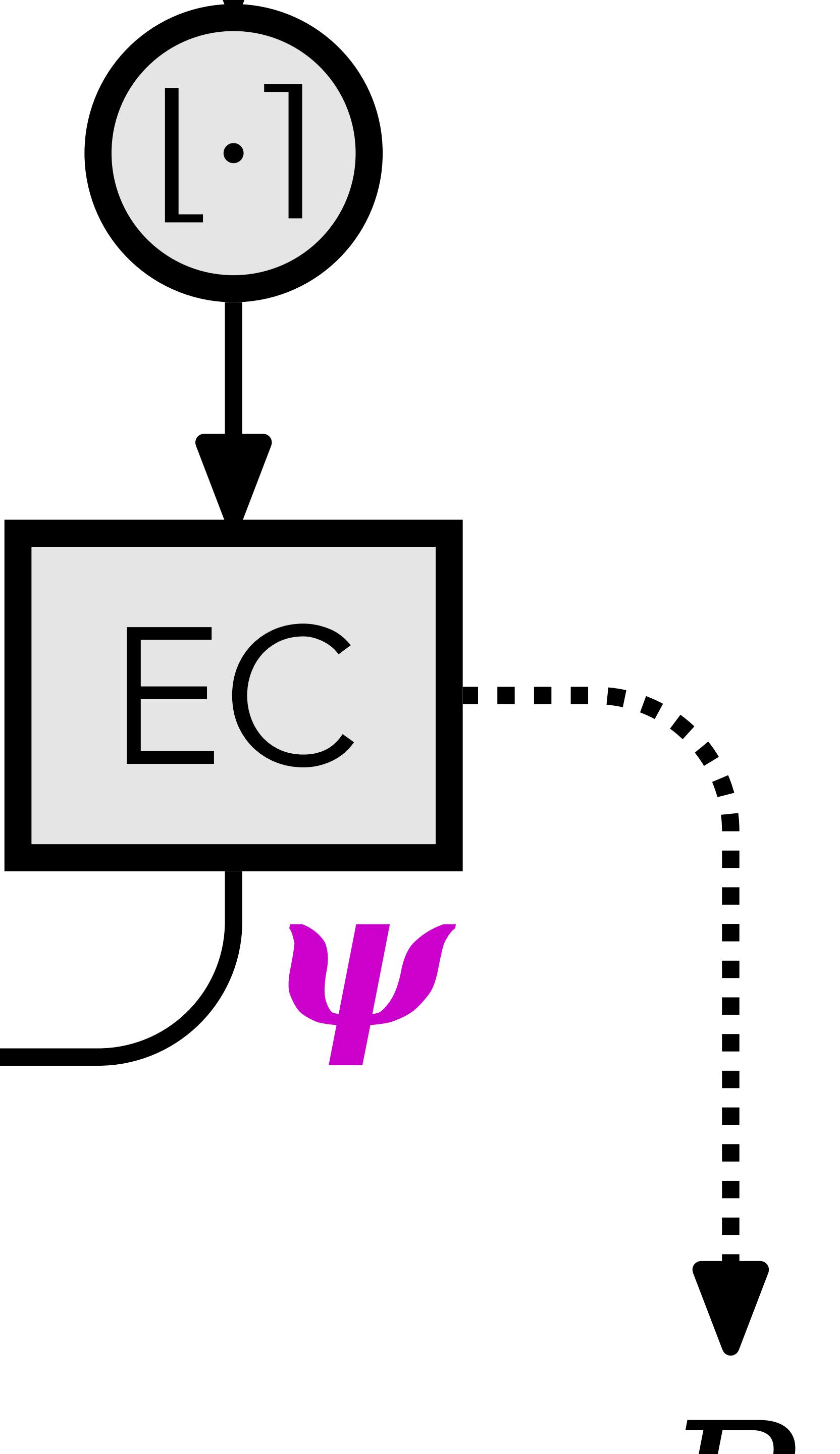
Nonlinear transform coding (end-to-end optimized)



$$L(\theta, \phi, \psi) = \underbrace{\mathbb{E}_{\mathbf{x}}[-\log_2 p_{\hat{\mathbf{y}}}(\hat{\mathbf{y}})]}_R + \lambda \underbrace{\mathbb{E}_{\mathbf{x}}[\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2]}_D$$

\hat{y}

latent



quarrelz
(round
entro
codi

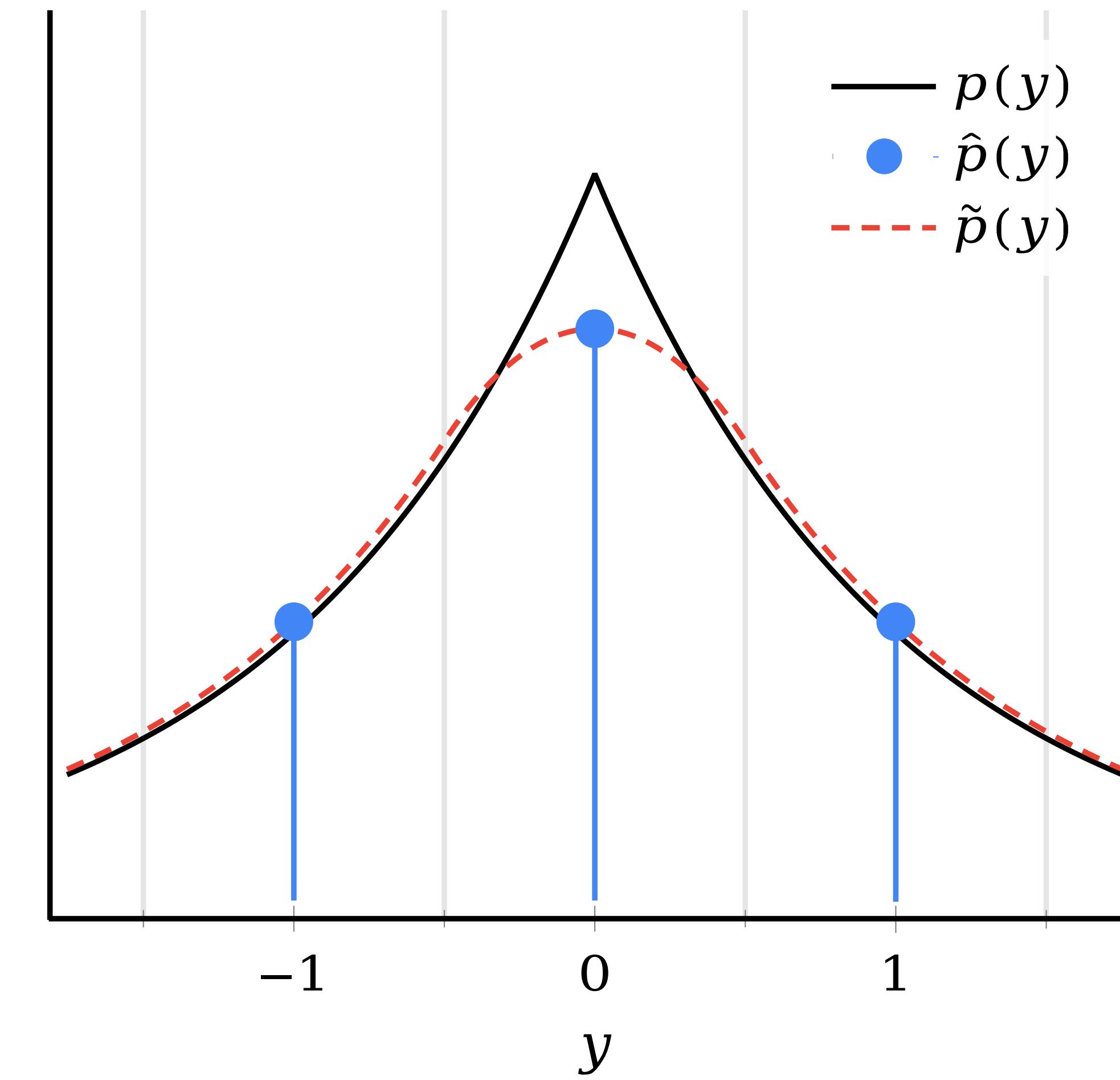
Quantization replaced by uniform noise for training

Quantized y described by probability mass function at integer locations.

“Noisy” y described by probability density function.

At integer locations, they coincide by construction.

Estimate PDF during training, then use PMF for compression.

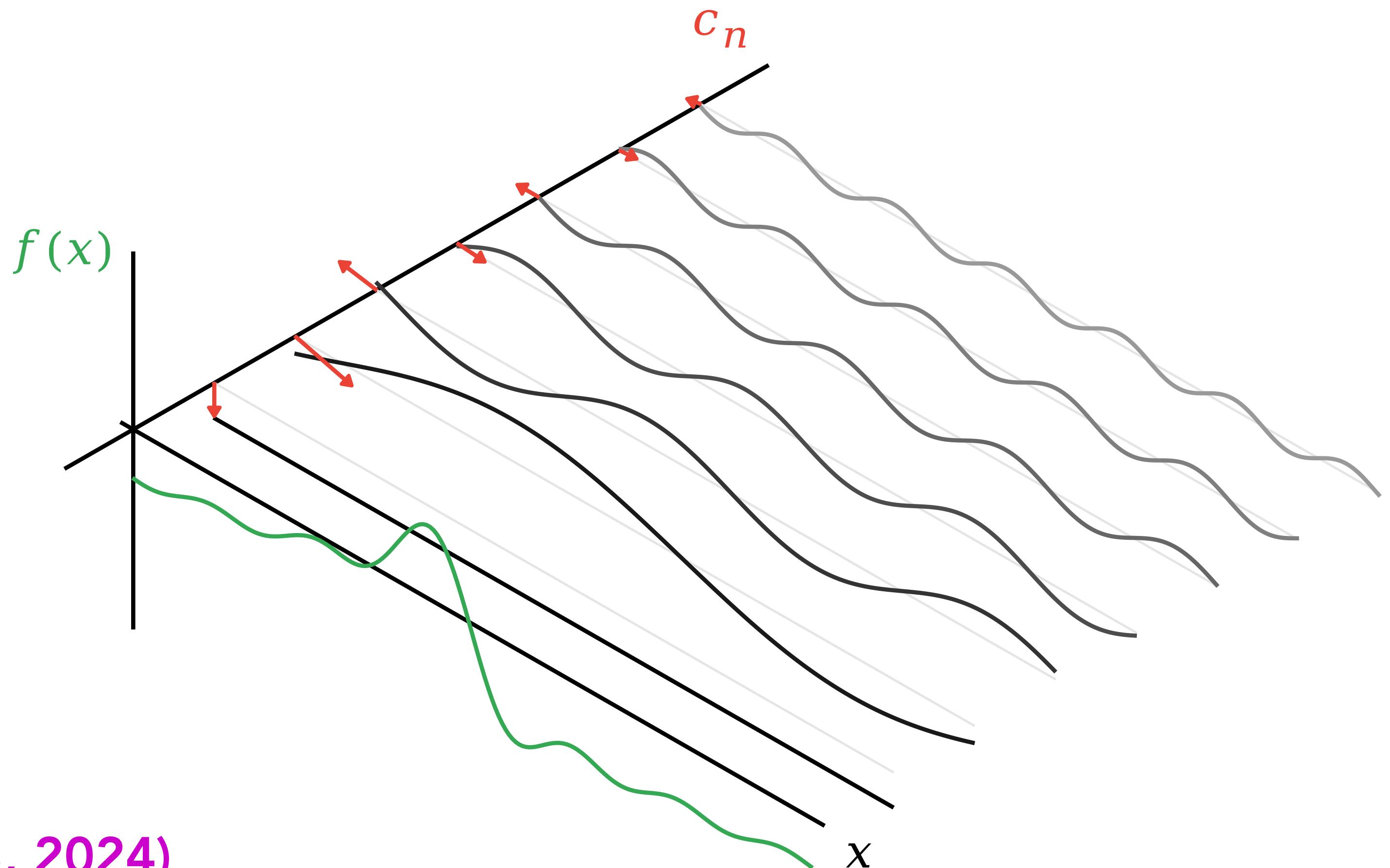


Fourier Basis Density Model

- Parametric densities are inflexible.
- Mixture models get stuck in local minima.
- The cumulative density model we proposed earlier is slow to train.
- **Idea:** represent a density by a truncated Fourier series.

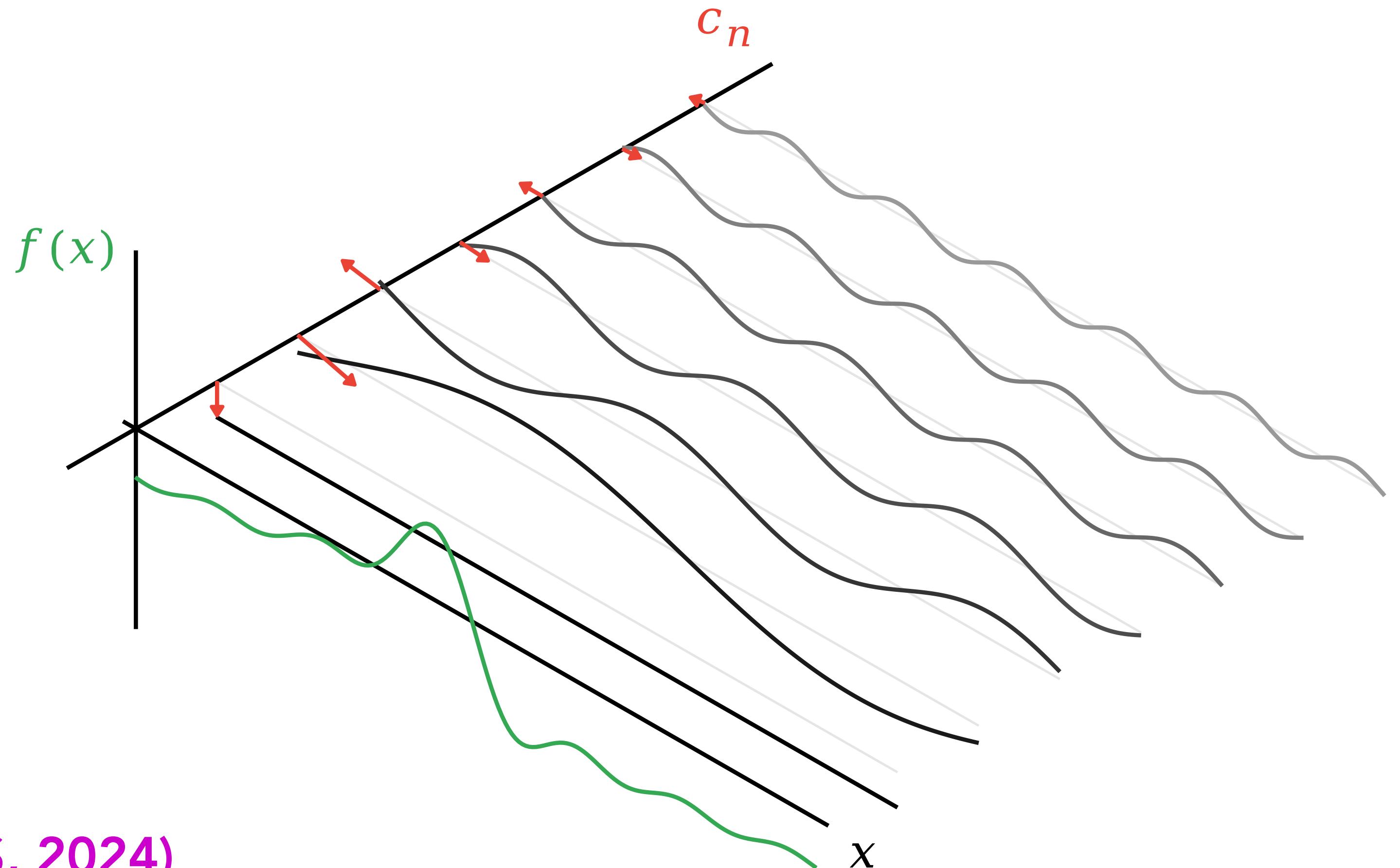
1. Let's be real!

$$f(x; \{c_n\}_{-N}^N) = \sum_{n=-N}^N c_n \exp(in\pi x)$$



1. Let's be real!

$$f(x; \{c_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$



2. Don't be negative!

$$f(x; \{c_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

2. Don't be negative!

$$f(x; \{c_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

Theorems: Herglotz and Wiener–Khinchin

2. Don't be negative!

$$f(x; \{c_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

Theorems: Herglotz and Wiener–Khinchin

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

2. Don't be negative!

$$f(x; \{c_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

Theorems: Herglotz and Wiener–Khinchin

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

$$\begin{bmatrix} d & e & f & g \\ c & d & e & f \\ b & c & d & e \\ a & b & c & d \end{bmatrix}$$

2. Don't be negative!

$$f(x; \{c_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

Theorems: Herglotz and Wiener–Khinchin

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

$$\begin{bmatrix} d & e & f & g \\ c & d & e & f \\ b & c & d & e \\ a & b & c & d \end{bmatrix} \quad (C = A^\top A)$$

2. Don't be negative!

$$f(x; \{a_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

Theorems: Herglotz and Wiener–Khinchin

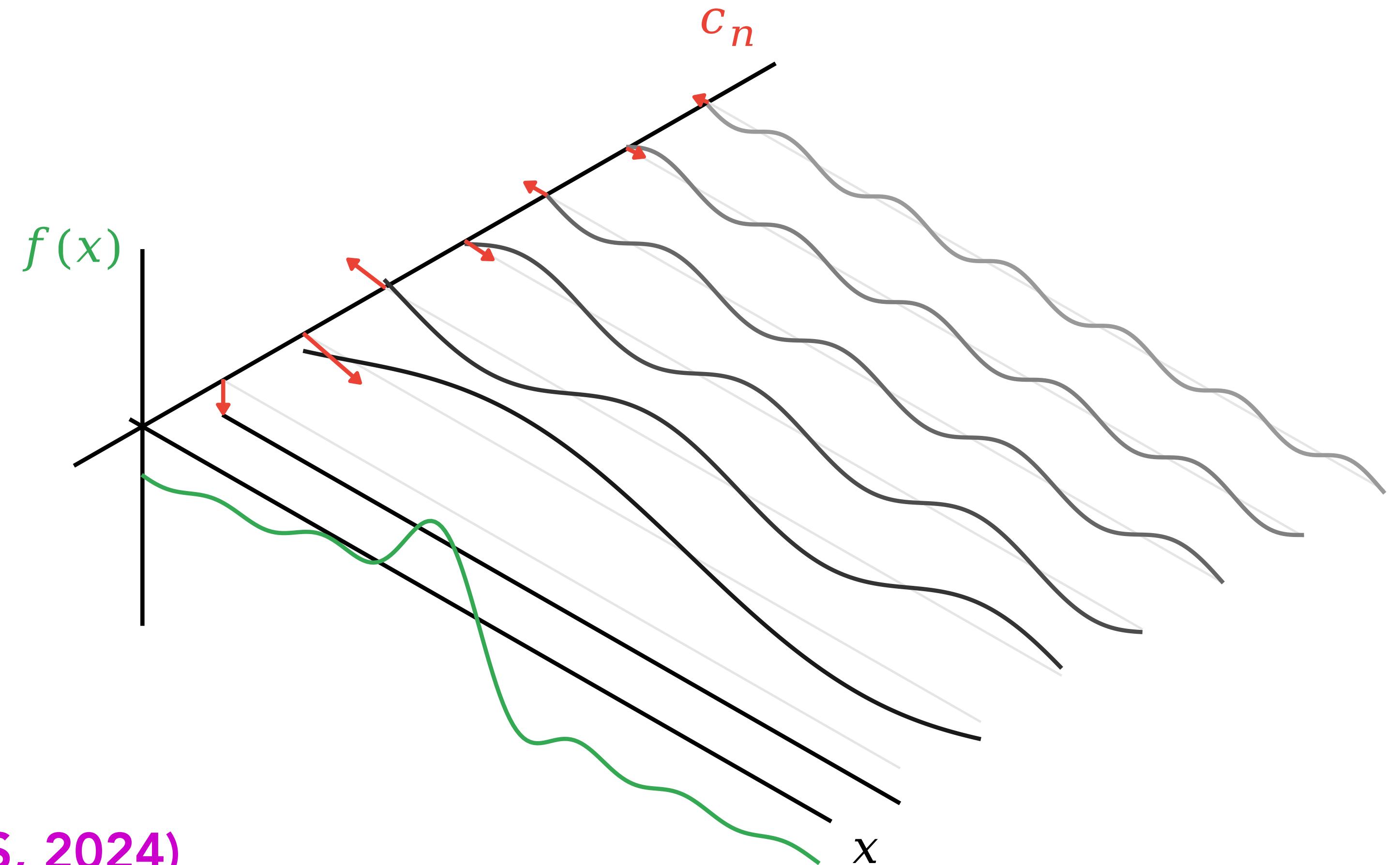
$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

$$\begin{bmatrix} d & e & f & g \\ c & d & e & f \\ b & c & d & e \\ a & b & c & d \end{bmatrix} \quad (C = A^\top A)$$

2. Don't be negative!

$$f(x; \{a_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

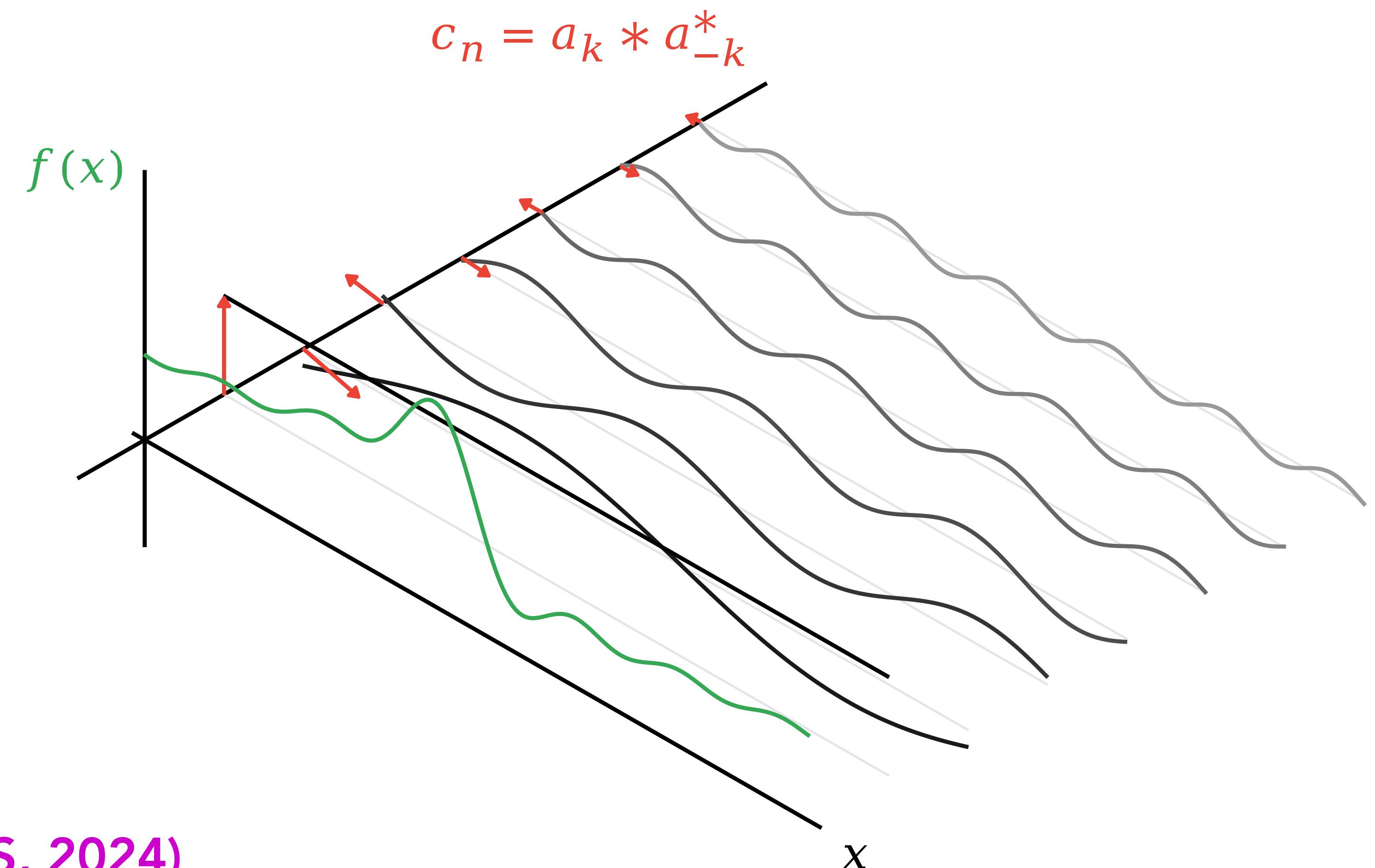
$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$



2. Don't be negative!

$$f(x; \{a_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

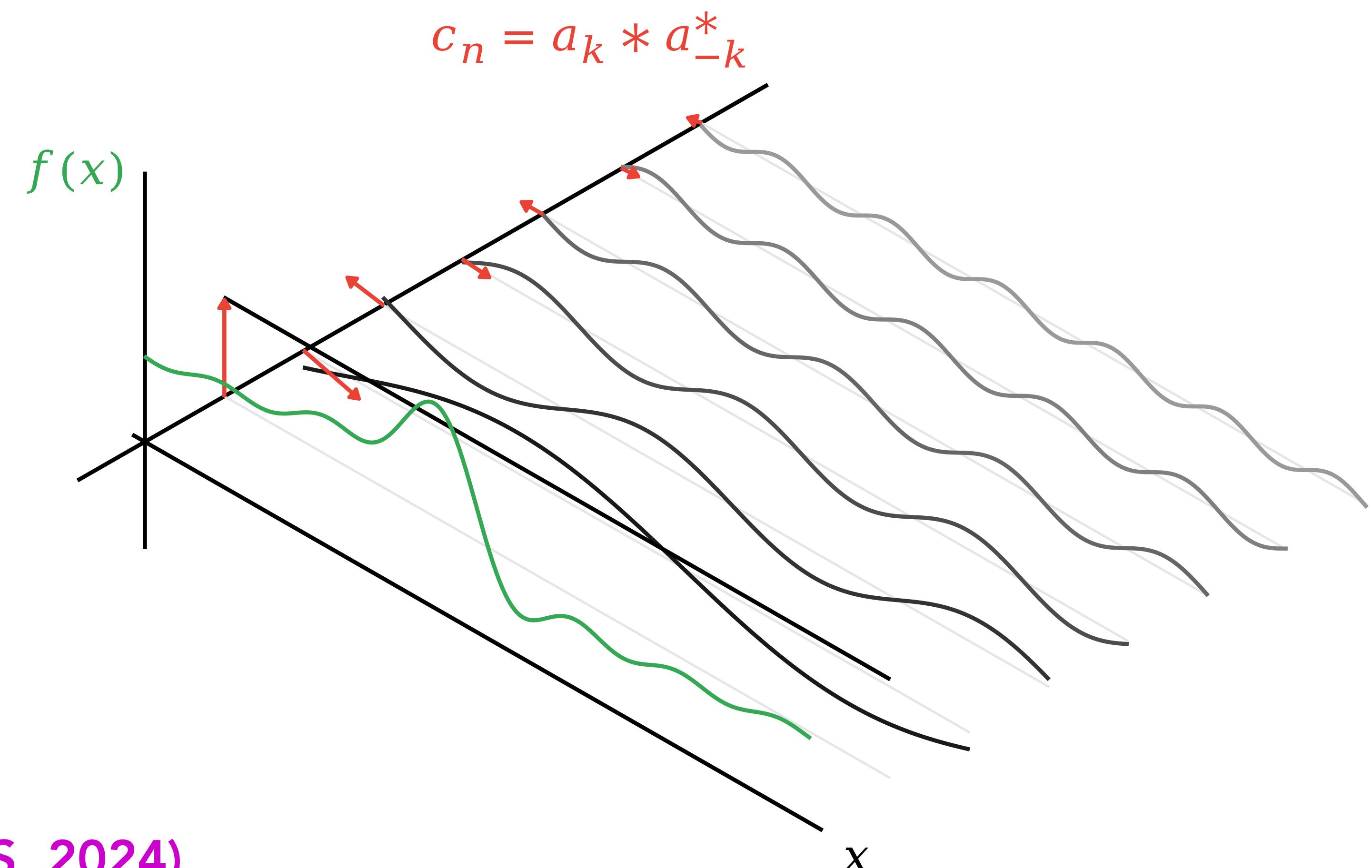
$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$



3. Normalize it!

$$f(x; \{a_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

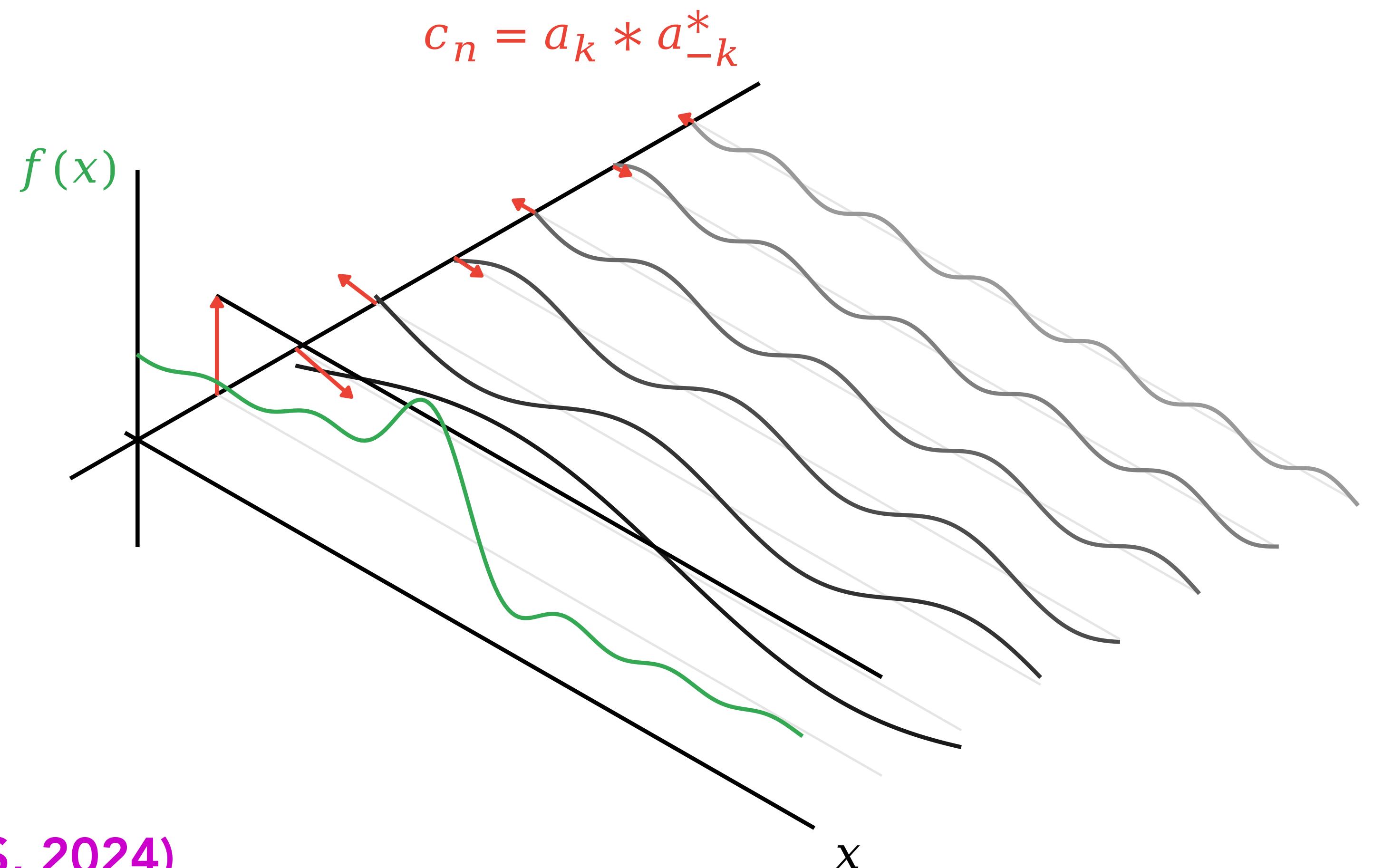


3. Normalize it!

$$f(x; \{a_n\}_0^N) = c_0 + 2 \sum_{n=1}^N \Re\{c_n \exp(in\pi x)\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

$$\int_{-1}^1 f(x; \{a_n\}_0^N) = 2c_0$$

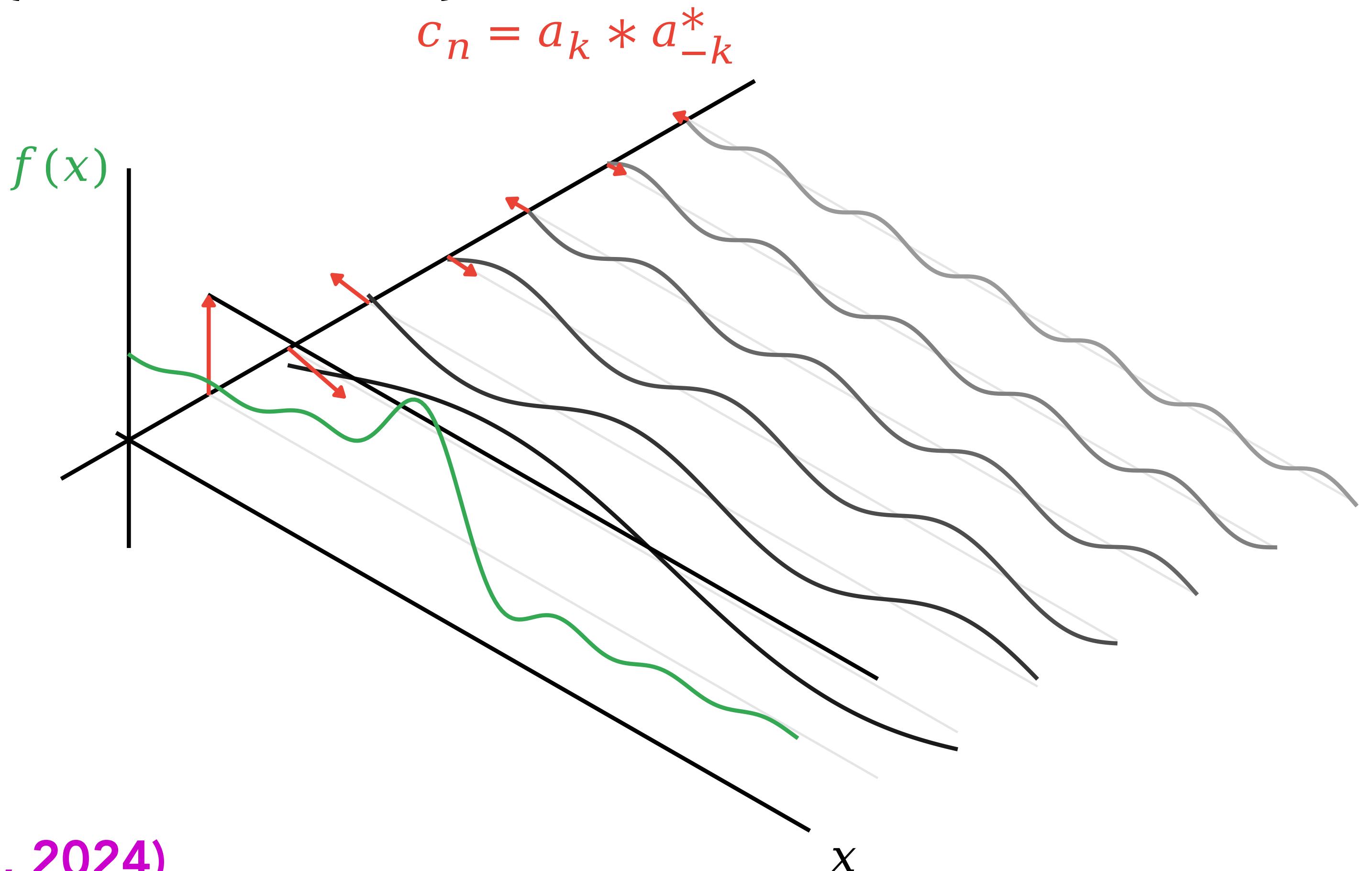


3. Normalize it!

$$p(x; \{a_n\}_0^N) = \frac{1}{2} + \sum_{n=1}^N \Re \left\{ \frac{c_n}{c_0} \exp(in\pi x) \right\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

$$\int_{-1}^1 f(x; \{a_n\}_0^N) dx = 2c_0$$

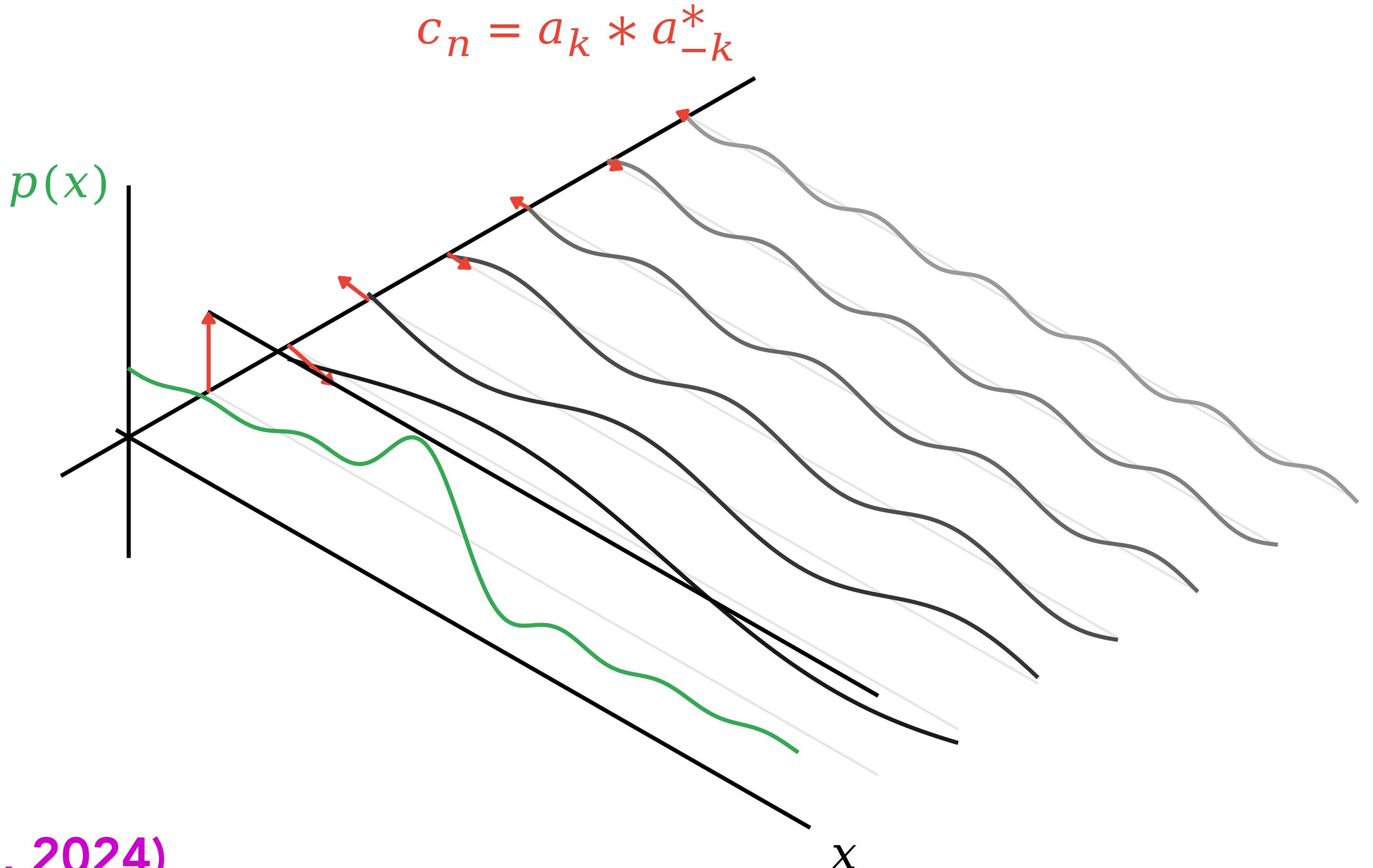


3. Normalize it!

$$p(x; \{a_n\}_0^N) = \frac{1}{2} + \sum_{n=1}^N \Re \left\{ \frac{c_n}{c_0} \exp(in\pi x) \right\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

$$\int_{-1}^1 f(x; \{a_n\}_0^N) dx = 2c_0$$



4. Expand support to \mathbb{R}

$$p(x; \{a_n\}_0^N) = \frac{1}{2} + \sum_{n=1}^N \Re \left\{ \frac{c_n}{c_0} \exp(in\pi x) \right\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

4. Expand support to \mathbb{R}

$$p(x; \{a_n\}_0^N) = \frac{1}{2} + \sum_{n=1}^N \Re \left\{ \frac{c_n}{c_0} \exp(in\pi x) \right\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

Change of variables:

$$g : (-1, 1) \rightarrow \mathbb{R}$$

4. Expand support to \mathbb{R}

$$p(x; \{a_n\}_0^N) = \frac{1}{2} + \sum_{n=1}^N \Re \left\{ \frac{c_n}{c_0} \exp(in\pi x) \right\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

Change of variables:

$$g(x; s, t) = s \cdot \tanh^{-1}(x) + t$$

$$g : (-1, 1) \rightarrow \mathbb{R}$$

$$g^{-1}(x; s, t) = \tanh\left(\frac{x - t}{s}\right)$$

4. Expand support to \mathbb{R}

$$p(x; \{a_n\}_0^N) = \frac{1}{2} + \sum_{n=1}^N \Re \left\{ \frac{c_n}{c_0} \exp(in\pi x) \right\}$$

$$c_n = \sum_{k=0}^{N-n} a_k a_{k+n}^*$$

Change of variables:

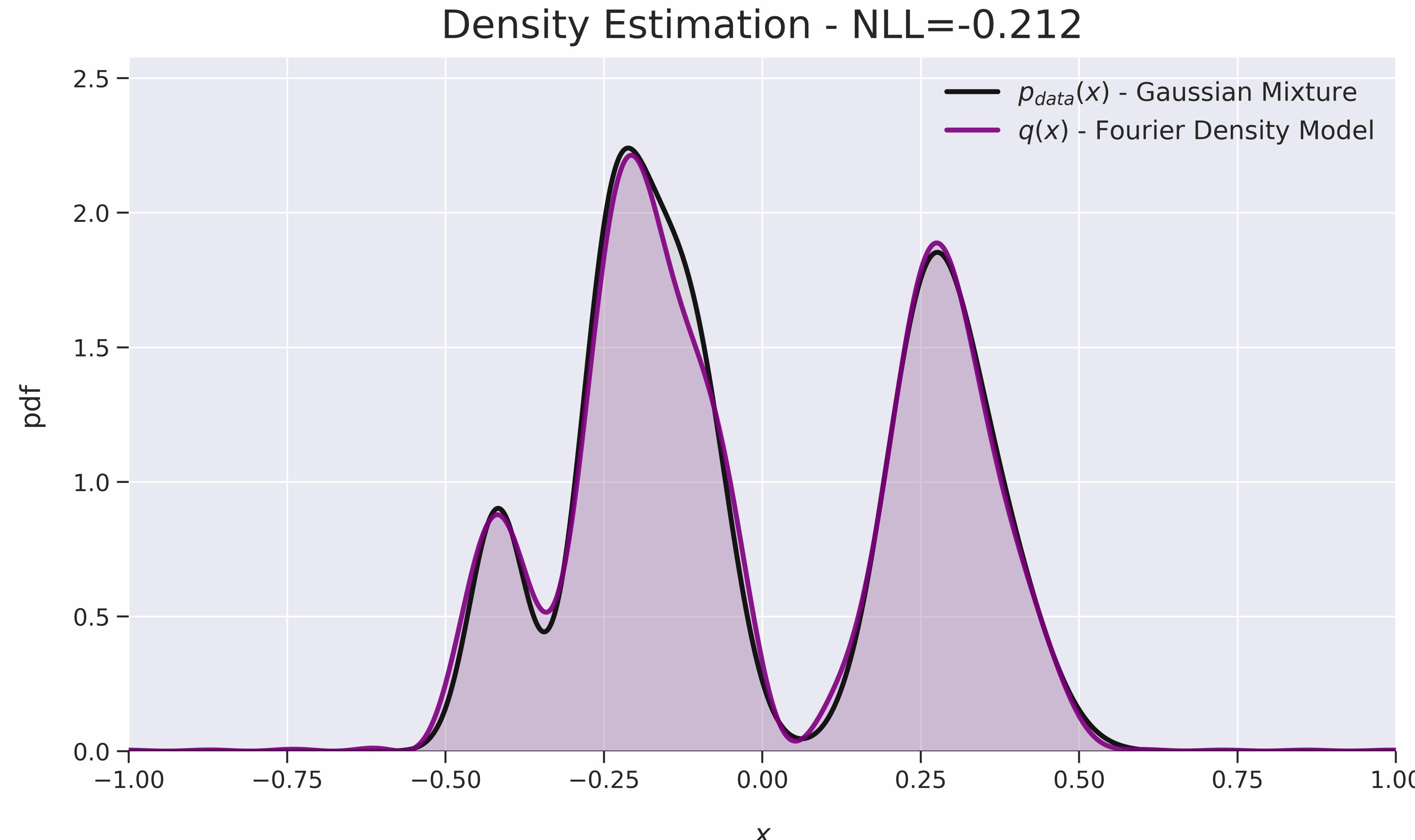
$$g(x; s, t) = s \cdot \tanh^{-1}(x) + t$$

$$g : (-1, 1) \rightarrow \mathbb{R}$$

$$g^{-1}(x; s, t) = \tanh\left(\frac{x - t}{s}\right)$$

$$q(x; \{a_n\}_0^N, s, t) = p(g^{-1}(x; s, t); \{a_n\}_0^N) \cdot (g^{-1})'(x; s, t)$$

Density modeling: training example



How to draw samples from this density?

- **Inverse transform sampling:** Draw a sample U from a uniform distribution, then compute $P^{-1}(U)$, where P is the cumulative distribution function.
 - Not feasible since P^{-1} is not accessible in closed form.
- **Rejection sampling:** Define an envelope distribution $e(x)$ such that $M e(x) \geq p(x)$ everywhere. Then draw a sample X from $e(x)$, and reject samples based on the ratio $p(X) / (M e(X))$.
 - Difficult to determine $e(x)$ and M such that envelope is tight; if not tight, low acceptance rate makes it slow.

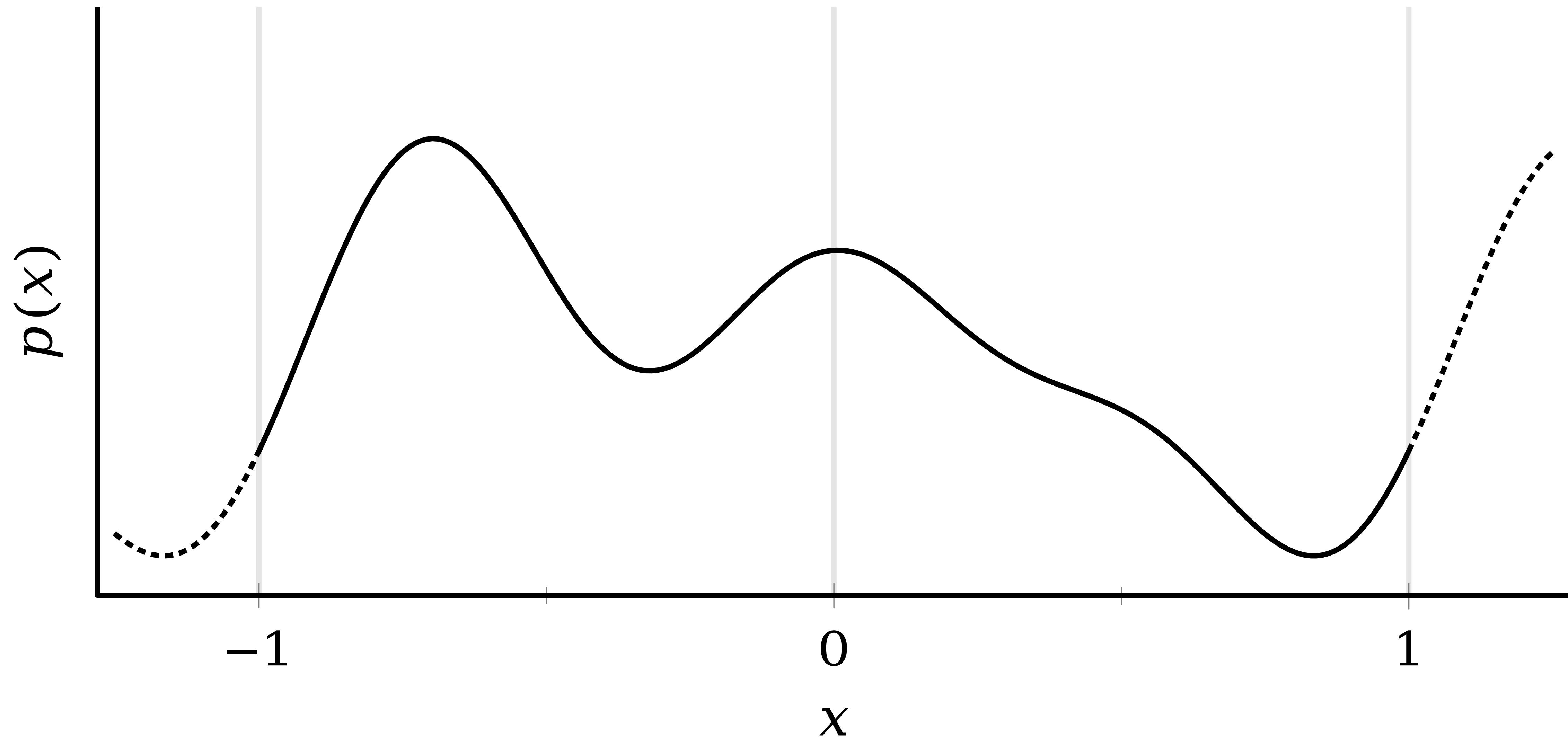
Speaking of sampling...

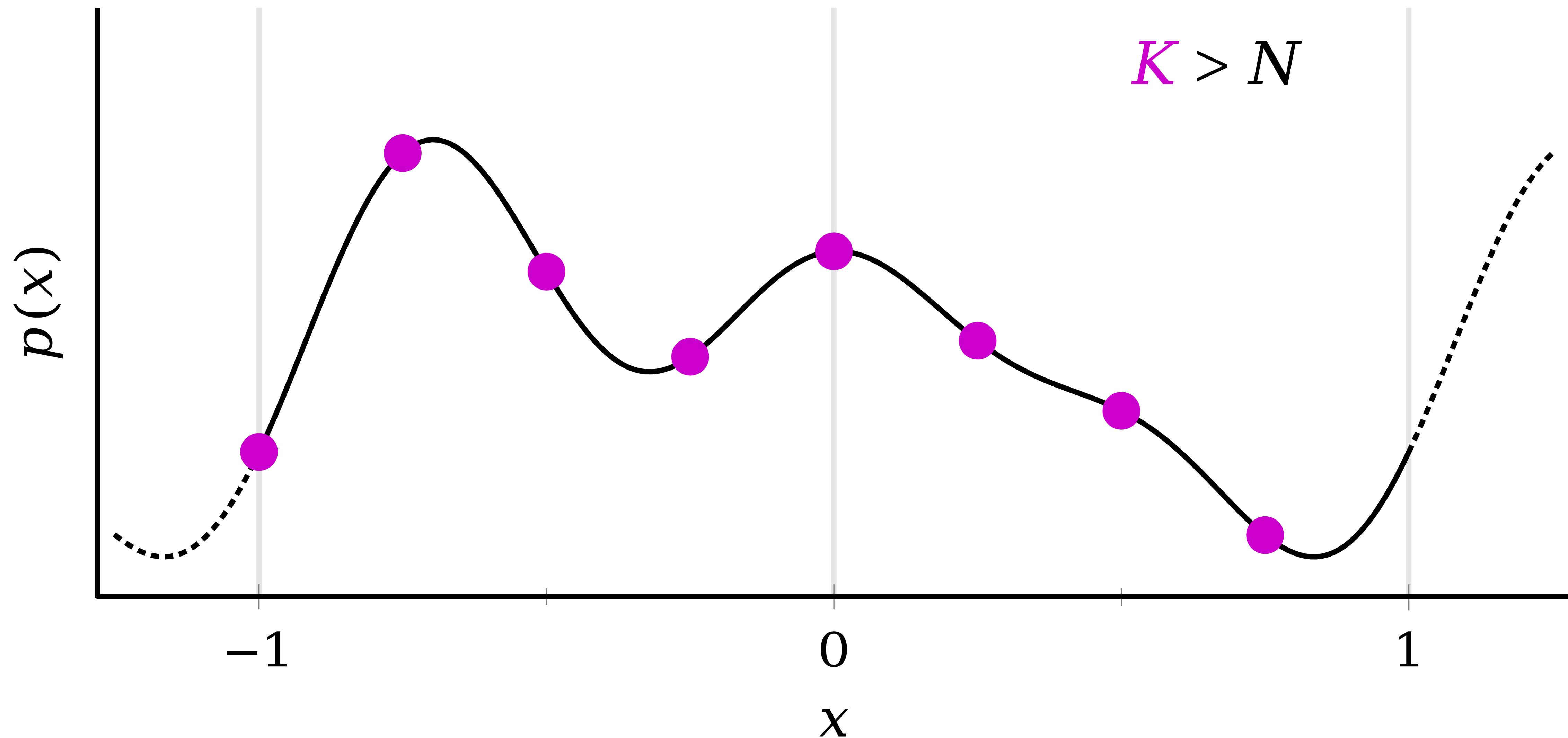
- In **Digital Signal Processing** (DSP), band-limited signals are **sampled** above the Nyquist rate, so they can be reconstructed with linear filtering.
- Here, the **density** is band-limited, instead of a signal.

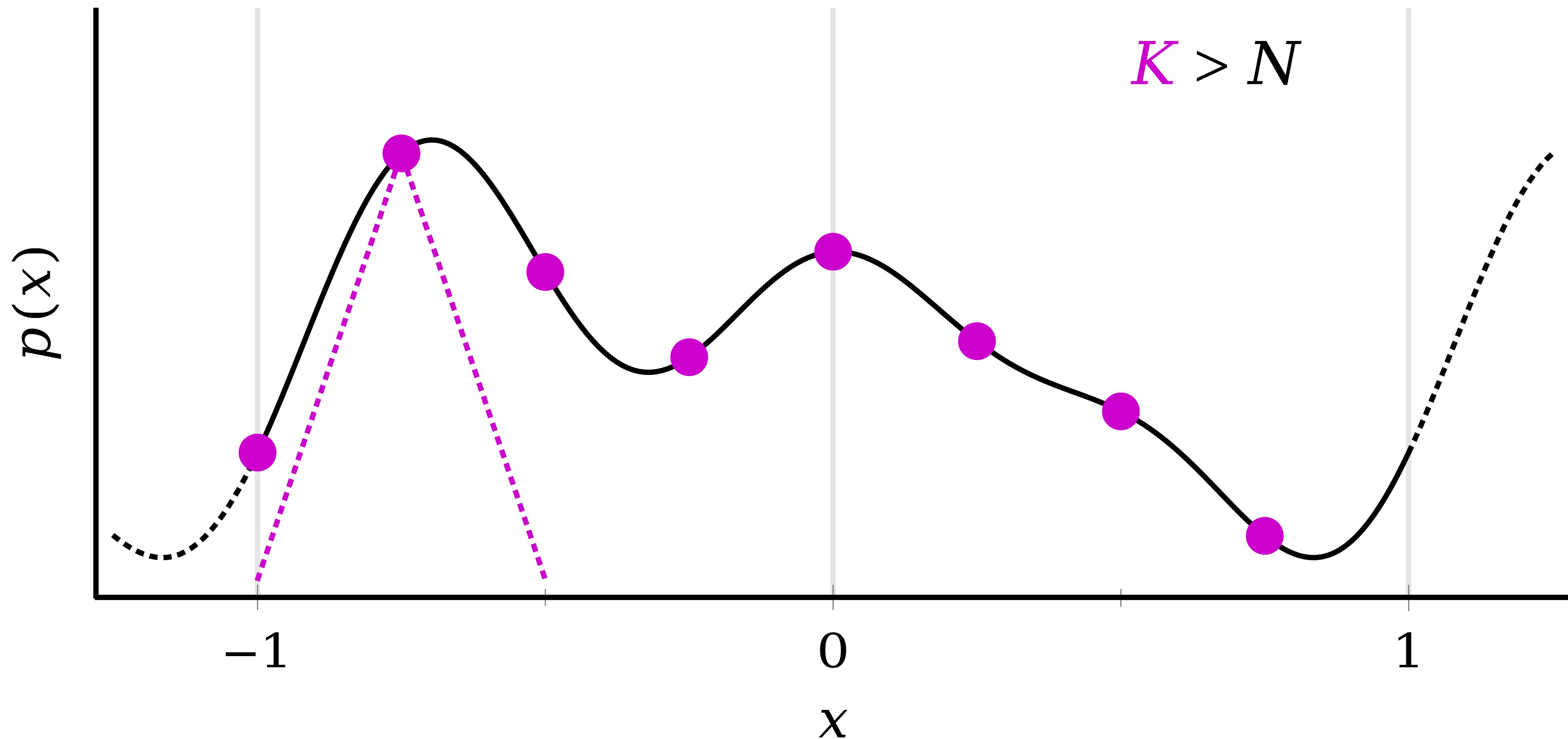
Speaking of sampling...

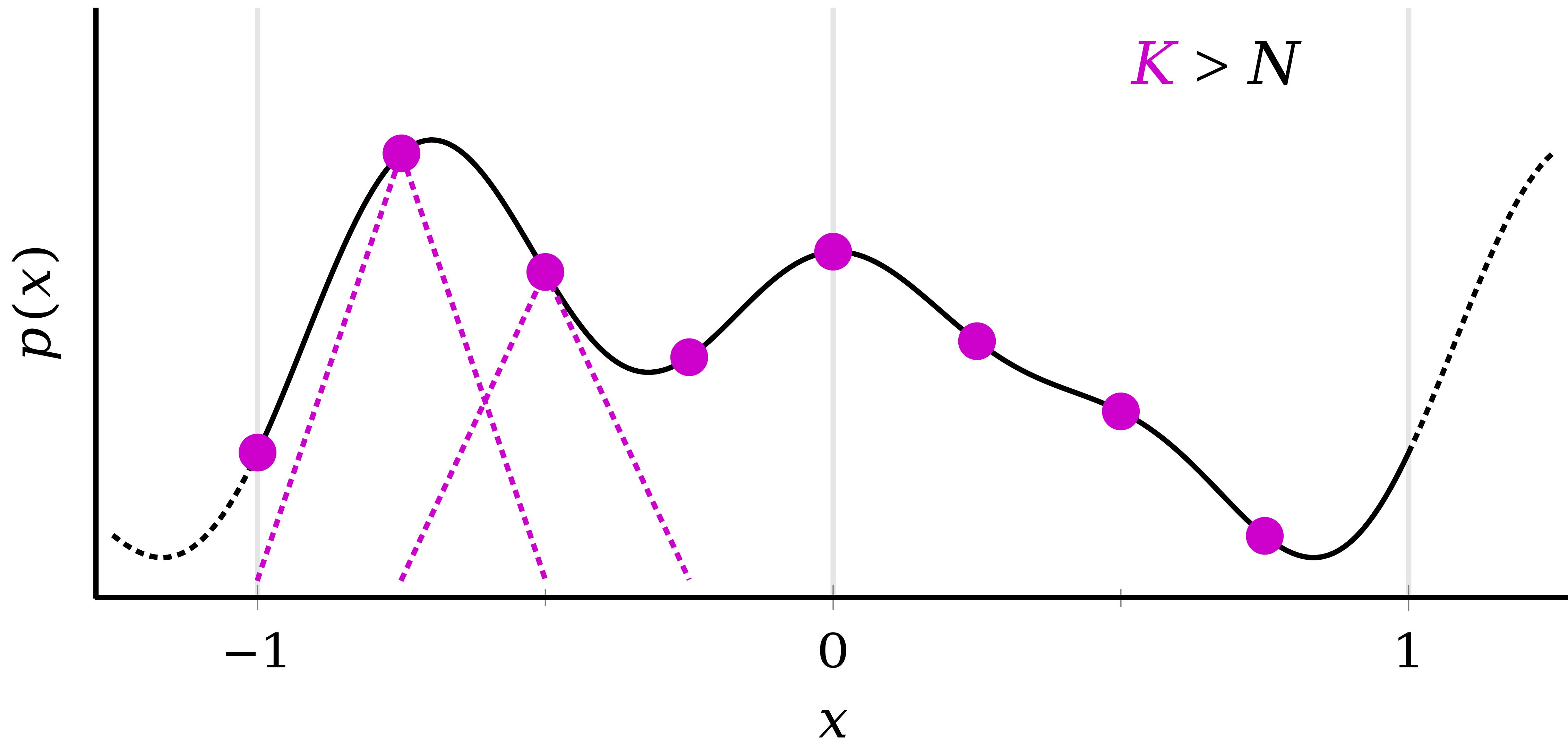
- In **Digital Signal Processing** (DSP), band-limited signals are **sampled** above the Nyquist rate, so they can be reconstructed with linear filtering.
- Here, the **density** is band-limited, instead of a signal.

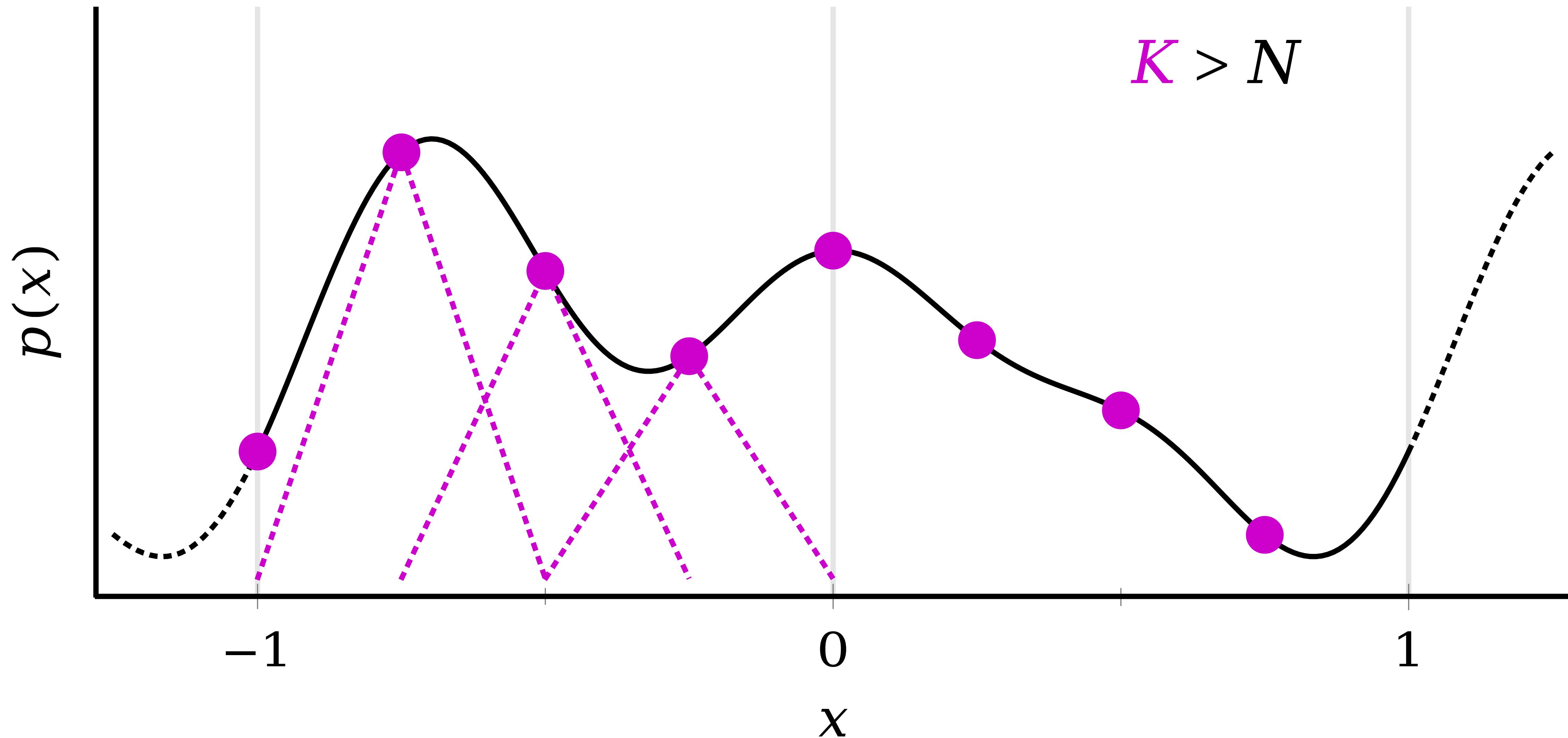


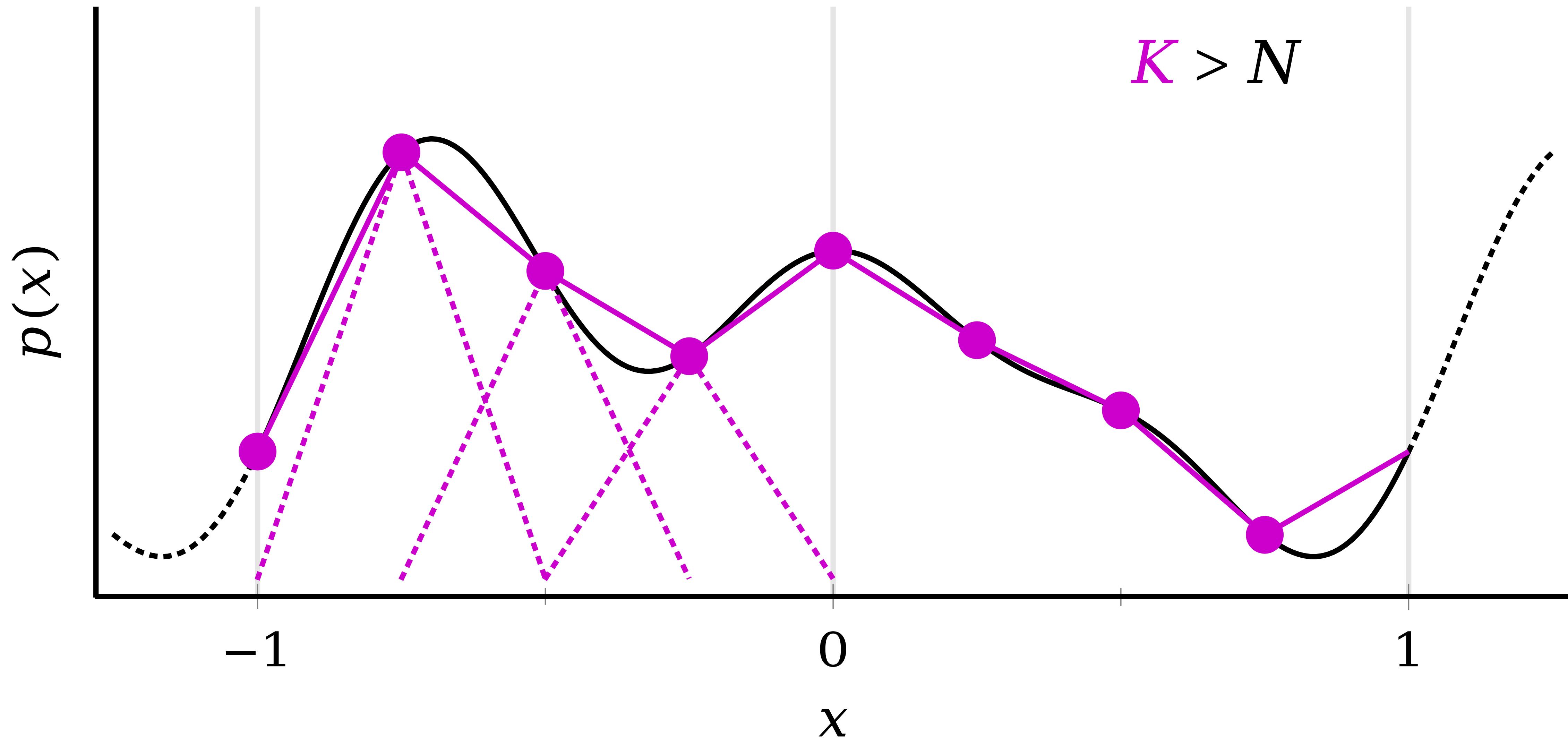






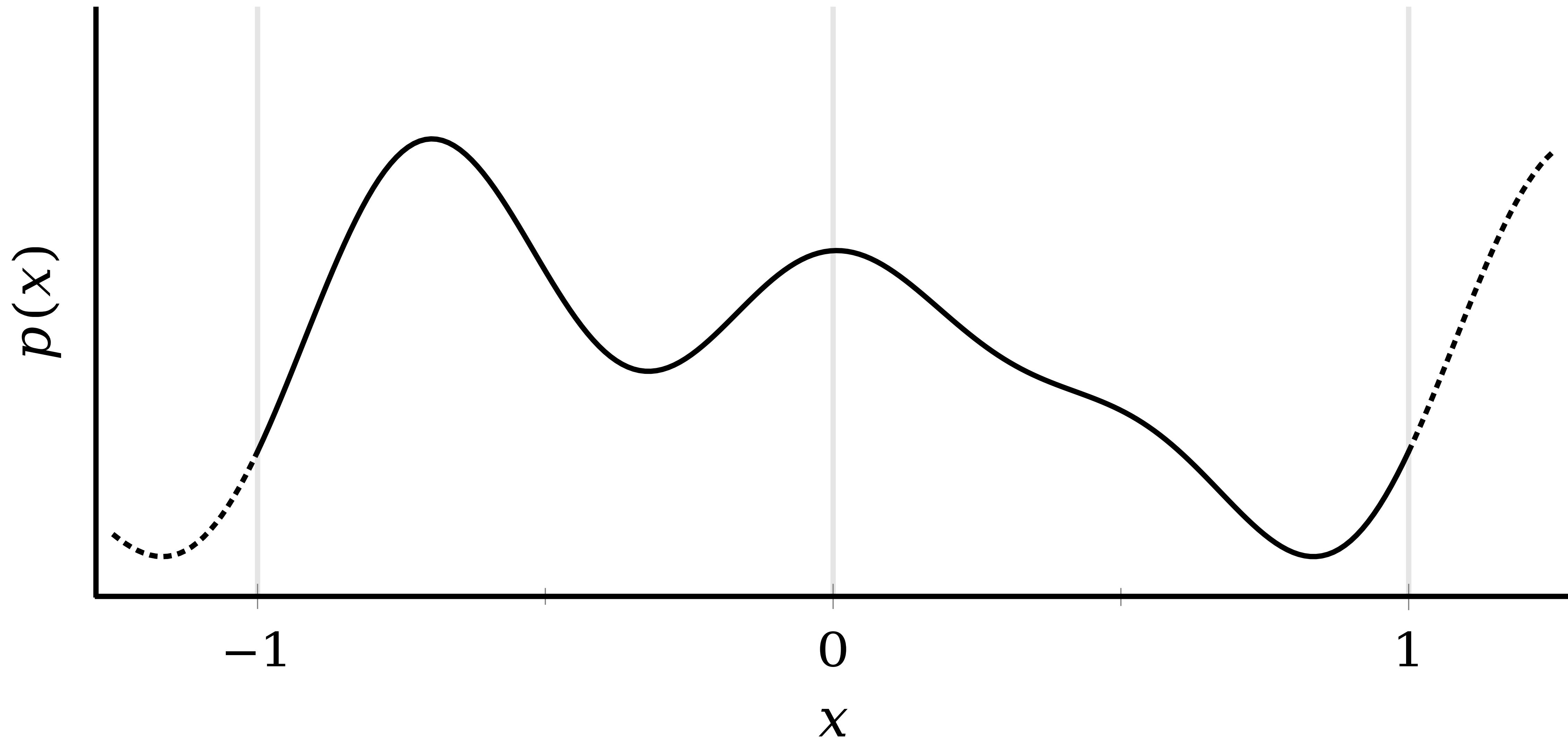


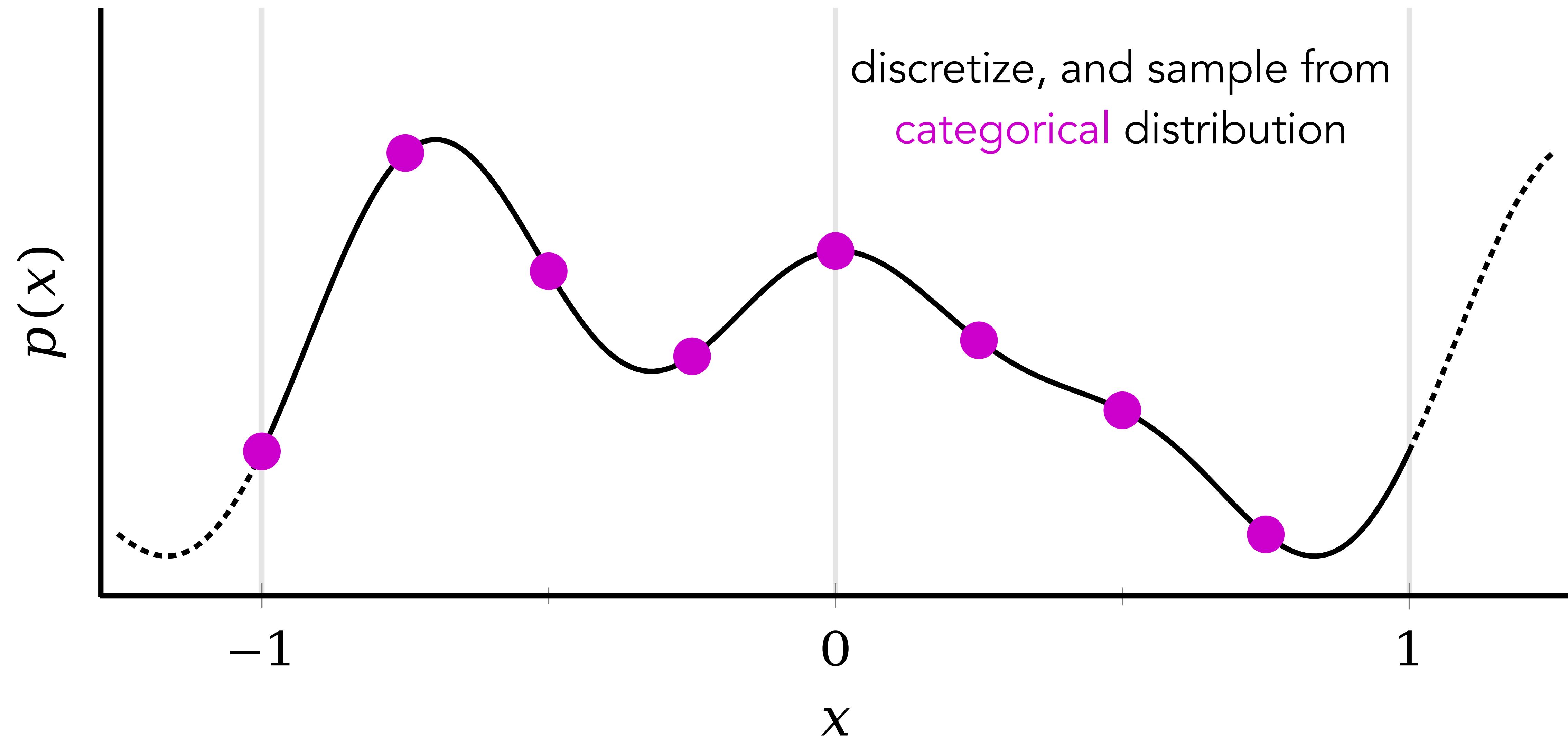


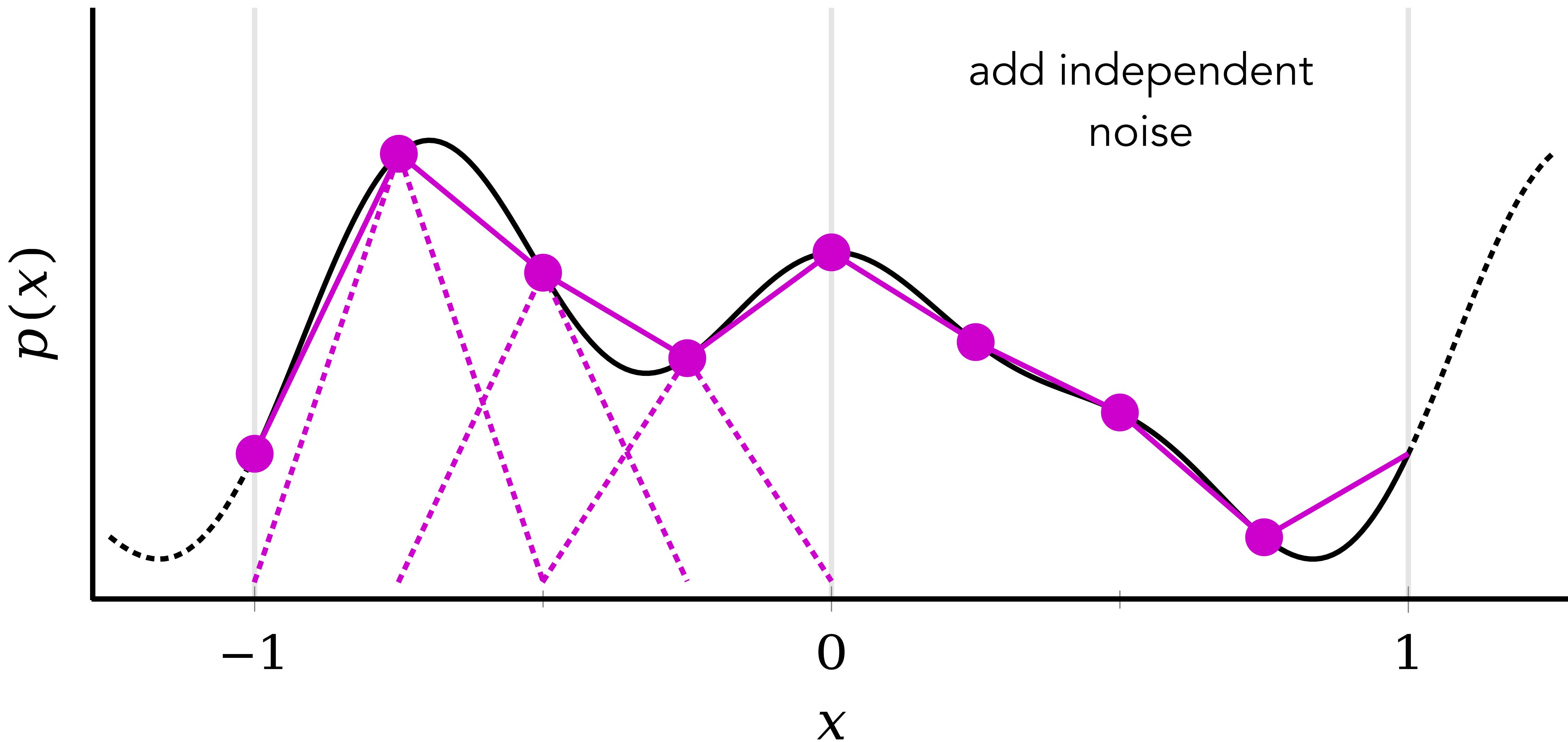


DSP vs. Discretized Approximate Ancestral Sampling

	DSP	DAAD
1. discretize	analog–digital converter	evaluate $p(x)$ (faster via DFT)
2. process/transmit discrete values	digital processing	sample from categorical
3. convolve with lowpass filter $w(x)$	digital–analog converter	add independent samples drawn from $w(x)$

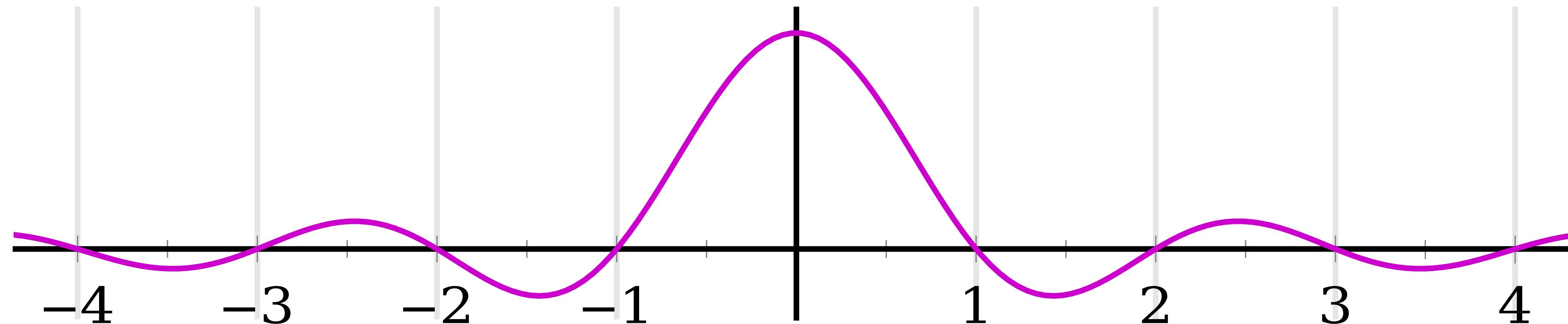






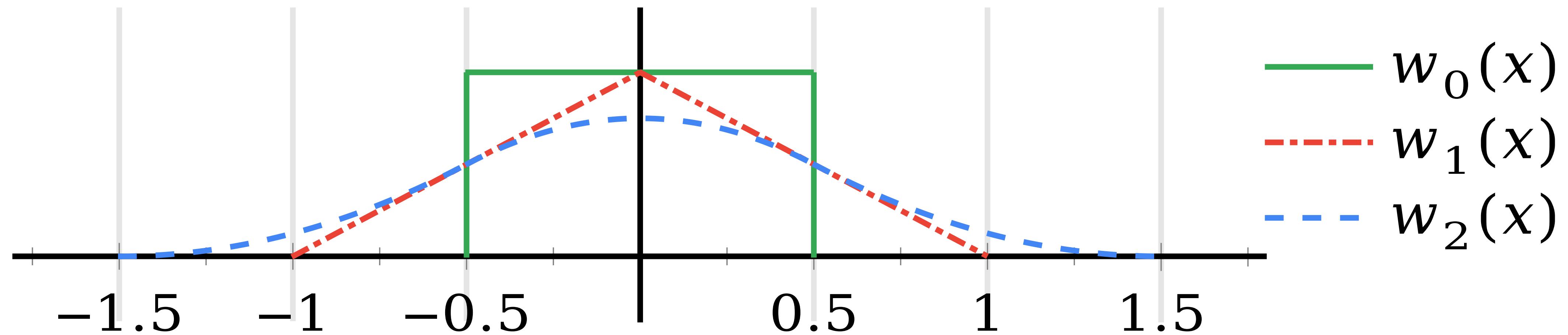
What reconstruction filter/density $w(x)$?

- Generally, we want $w(x)$ to be “close” to a *sinc* function.
- This is usually not possible → **filter design** problem.
- In DSP, we need to have finite support.
- In DAAD, we need non-negativity (because it must be a distribution).

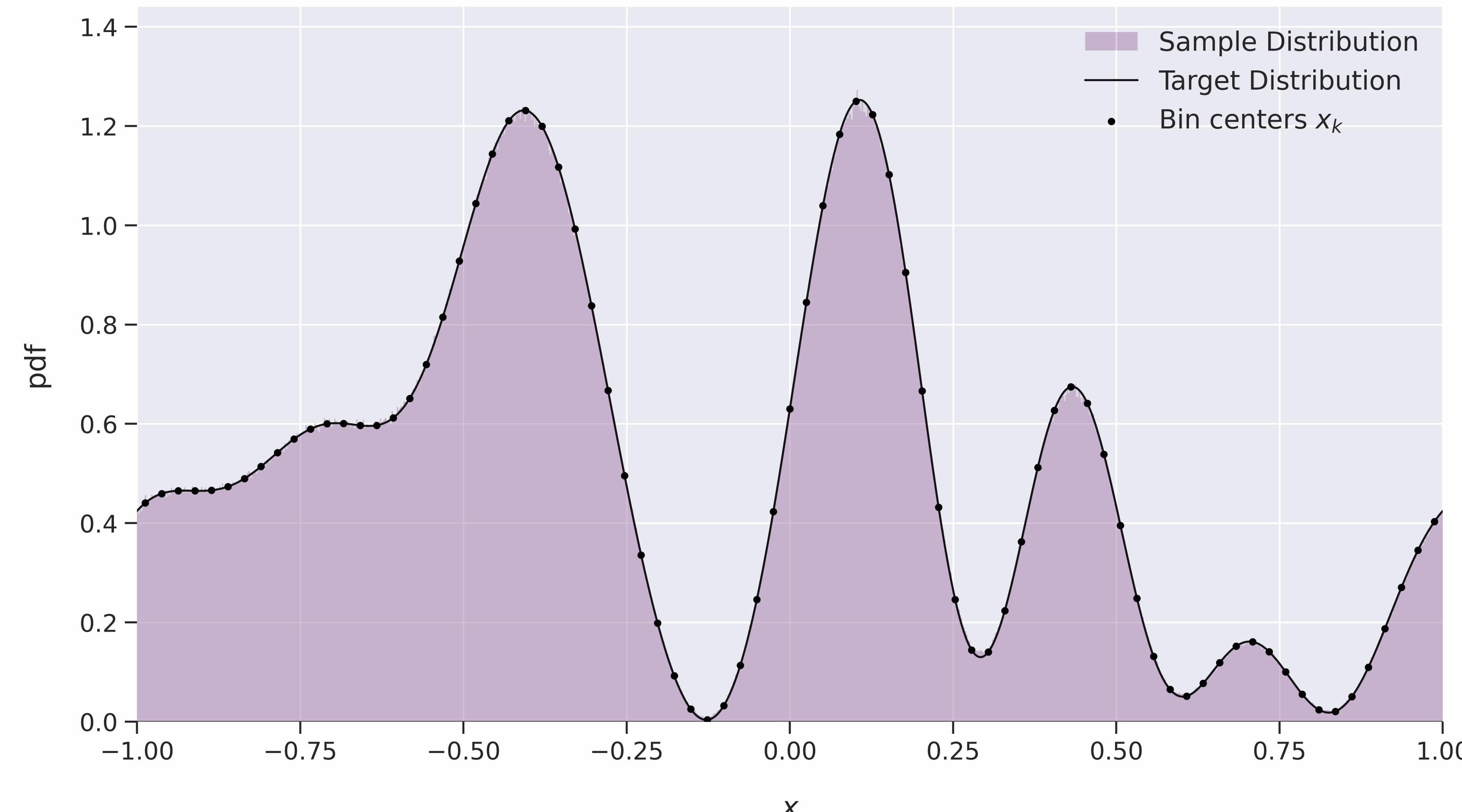


What reconstruction filter/density $w(x)$?

- B-splines are a simple and nice choice for DAAD.
- i -th order B-spline kernel: convolve a box function $i - 1$ times with itself.
- Simply take $i - 1$ i.i.d. samples from a uniform distribution, and add them.



Discretized Approximate Ancestral Sampling N=10 - K=80



Conclusion

- DAAD is a simple method for sampling from band-limited distributions.
- Strong analogy to Digital Signal Processing, but some of the practical requirements are inherently different (finite support vs. non-negativity).
- I've skipped many details here – in the paper, we also have:
 - Accuracy bounds for distribution of samples.
 - Empirical evaluation.
 - Refinement using MCMC with Langevin Dynamics.

Thank you!



[https://arxiv.org/abs/
2505.06098](https://arxiv.org/abs/2505.06098)