

Neural Distributed Compressor Does 'Binning'

Ezgi Ozyilkan

Neural Compression Workshop @ ICML 2023

Honolulu, HI | July 29, 2023

Neural Distributed Compressor Does 'Binning'

Ezgi Ozyilkan

Joint work with Johannes Ballé (Google Research)

Neural Compression Workshop @ ICML 2023

Honolulu, HI | July 29, 2023

Google Research



Neural Distributed Compressor Does 'Binning'

Ezgi Ozyilkan

Joint work with Johannes Ballé (Google Research) and Elza Erkip (NYU)

Neural Compression Workshop @ ICML 2023

Honolulu, HI | July 29, 2023

Google Research

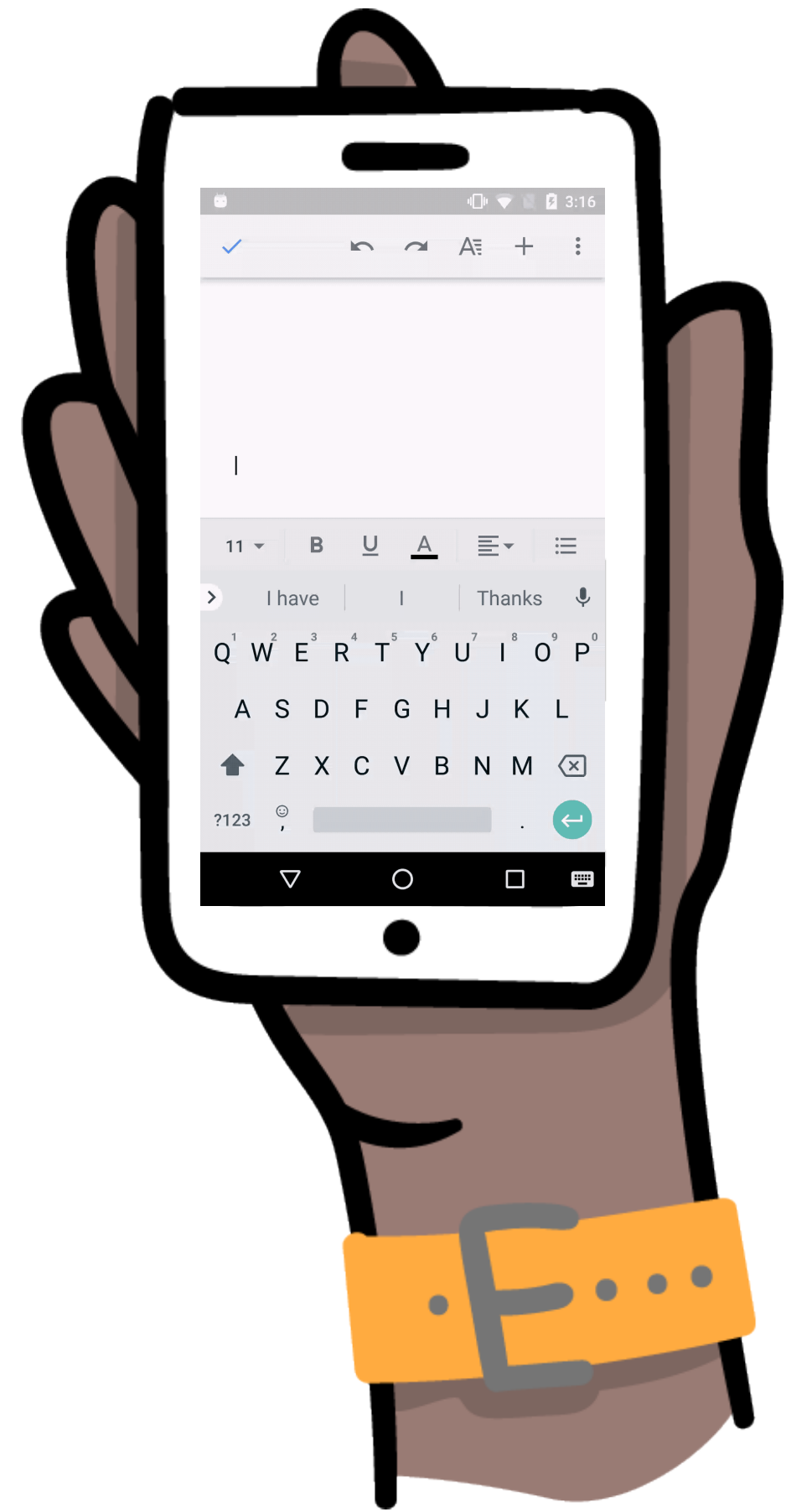


NYU

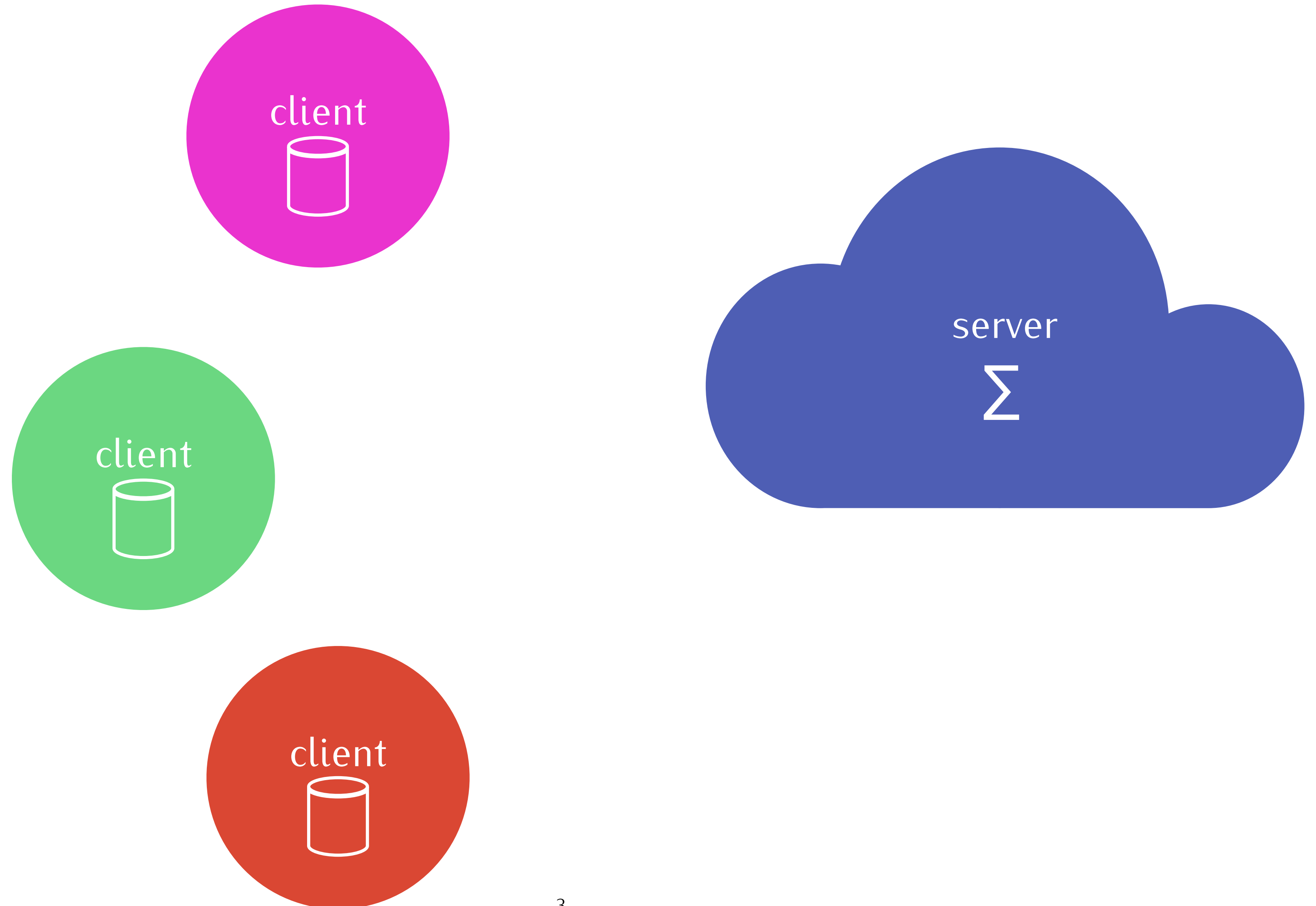
**TANDON SCHOOL
OF ENGINEERING**

Distributed Source Coding

Motivation: Distributed Source Coding



e.g., next-word prediction

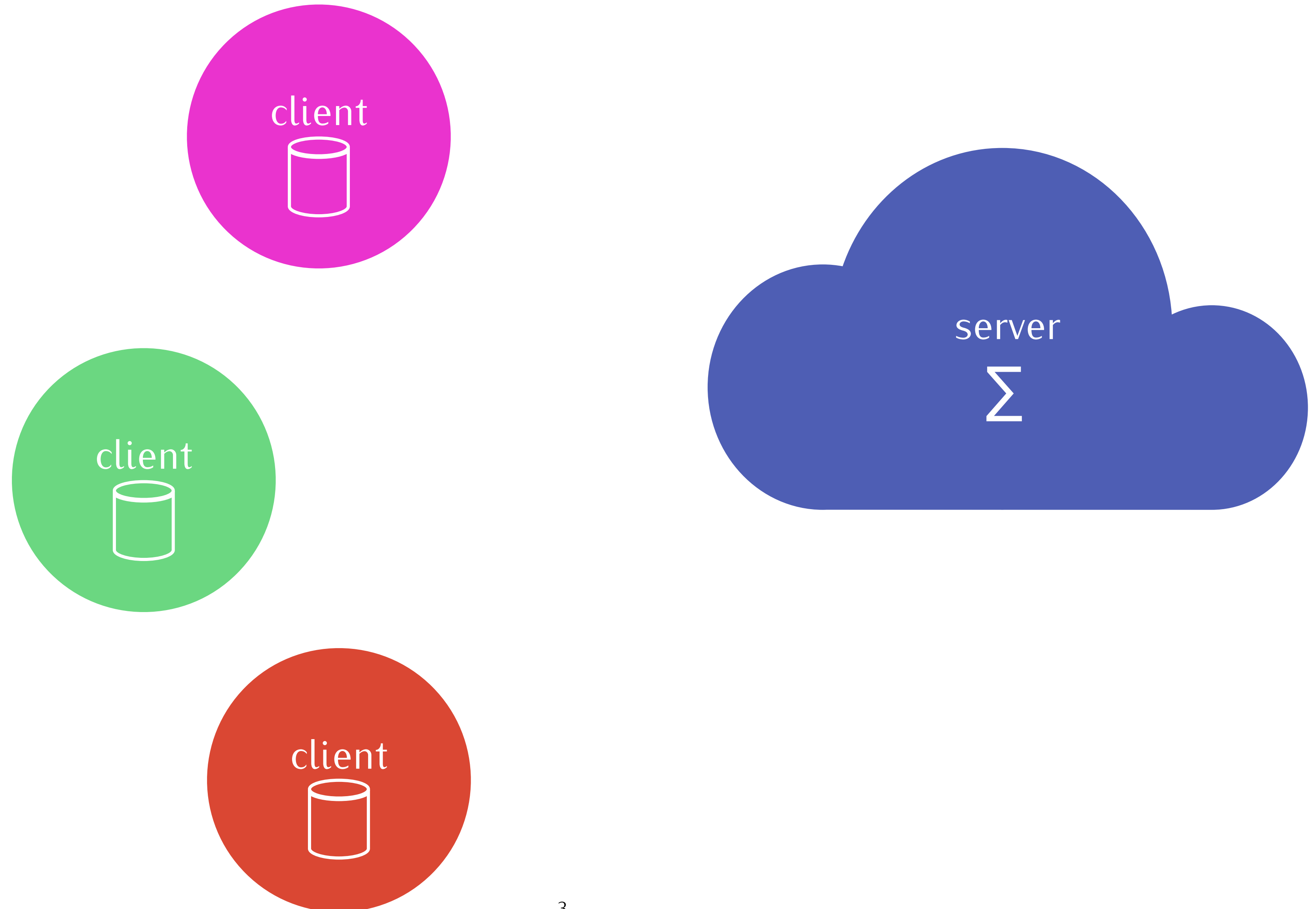


Motivation: Distributed Source Coding

Federated learning.



e.g., next-word prediction

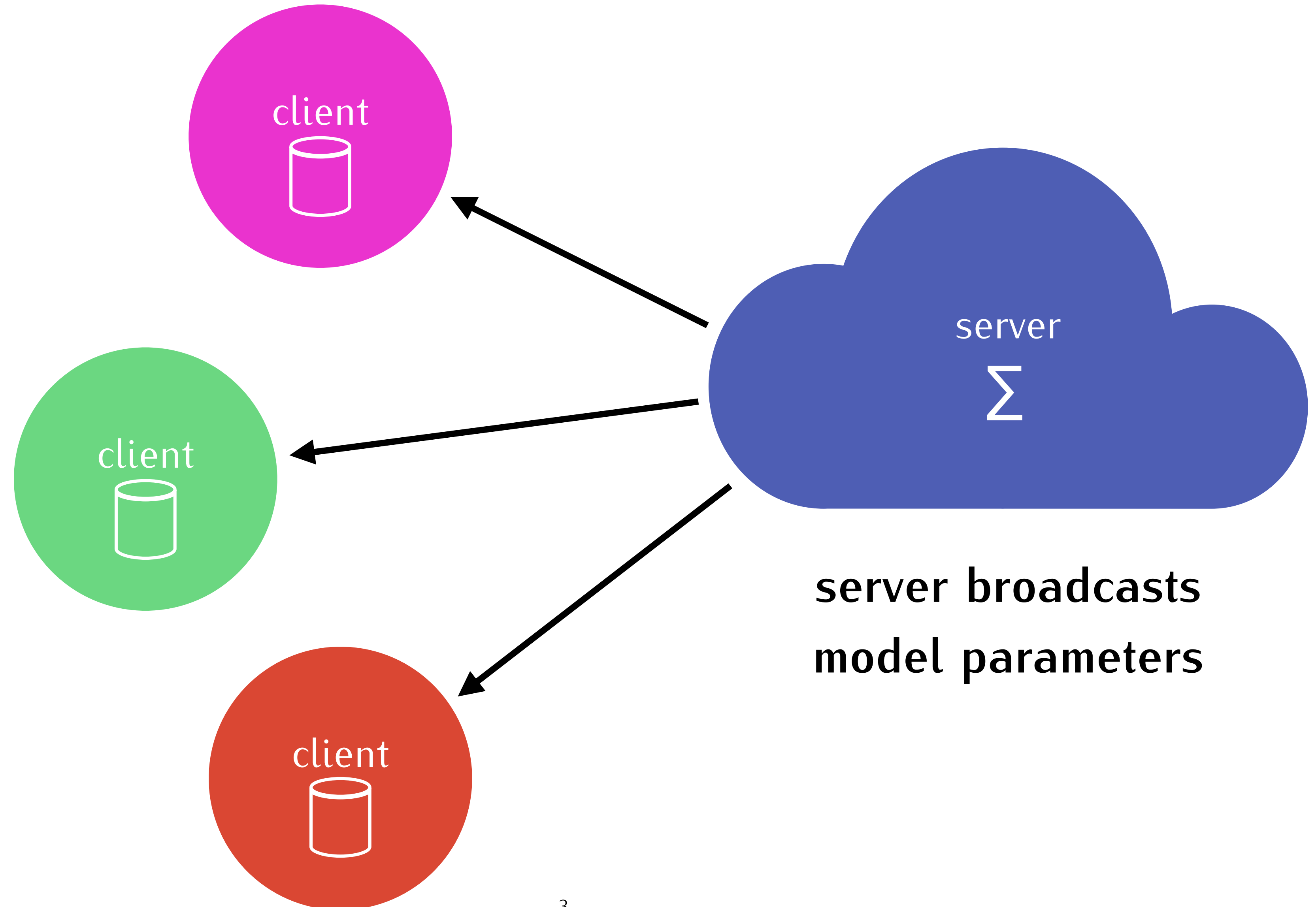


Motivation: Distributed Source Coding

Federated learning.



e.g., next-word prediction



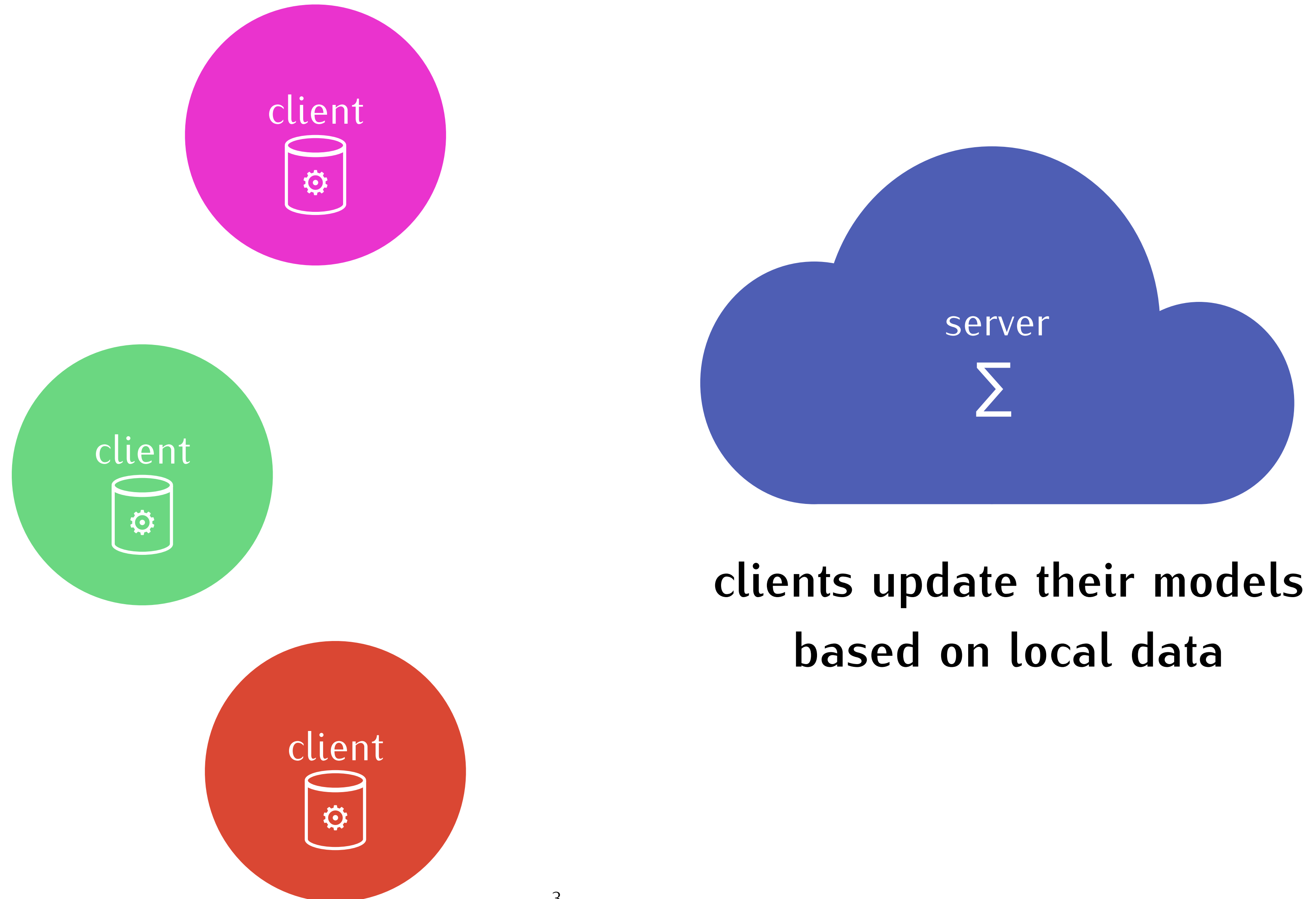
server broadcasts
model parameters

Motivation: Distributed Source Coding

Federated learning.



e.g., next-word prediction

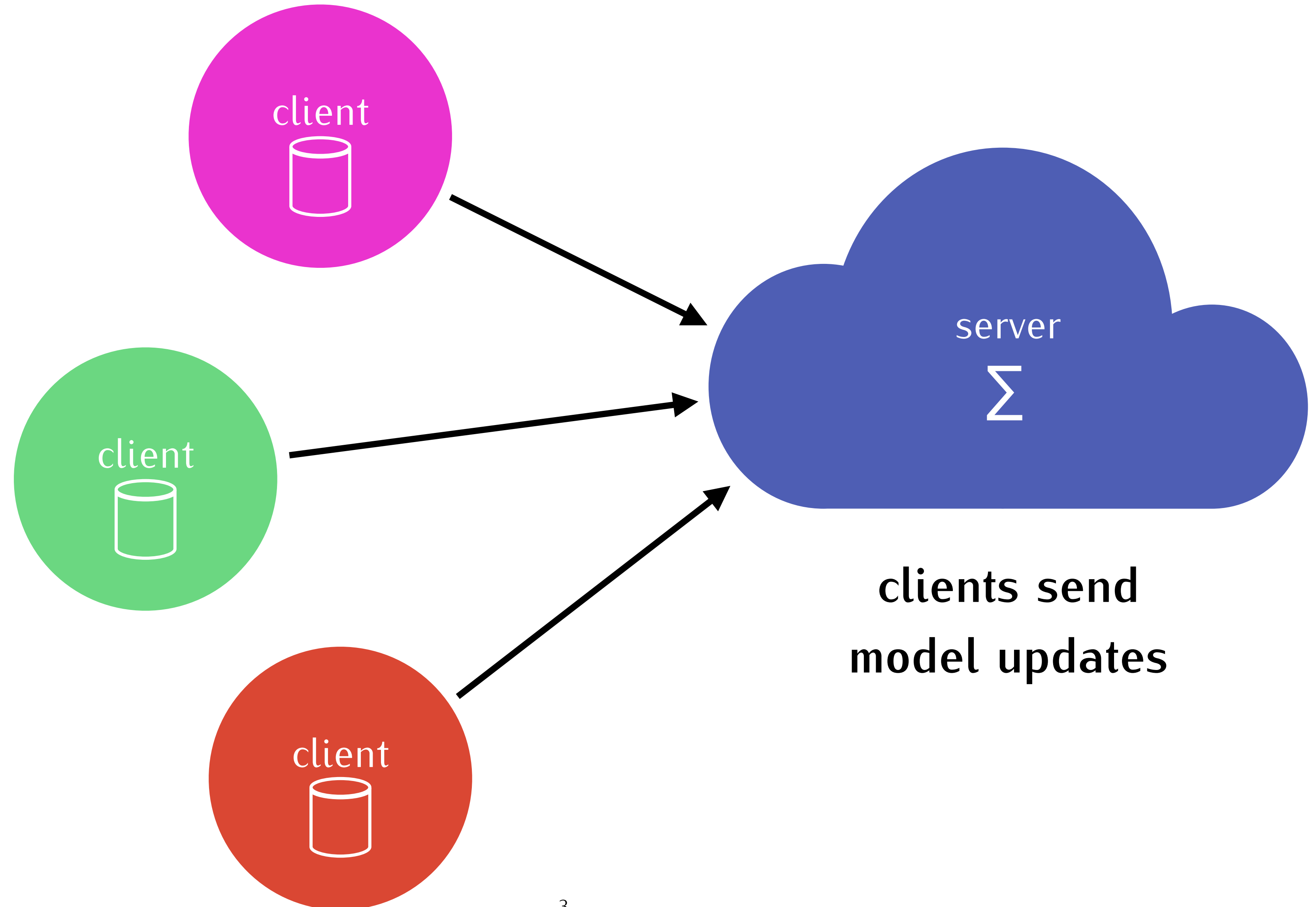


Motivation: Distributed Source Coding

Federated learning.

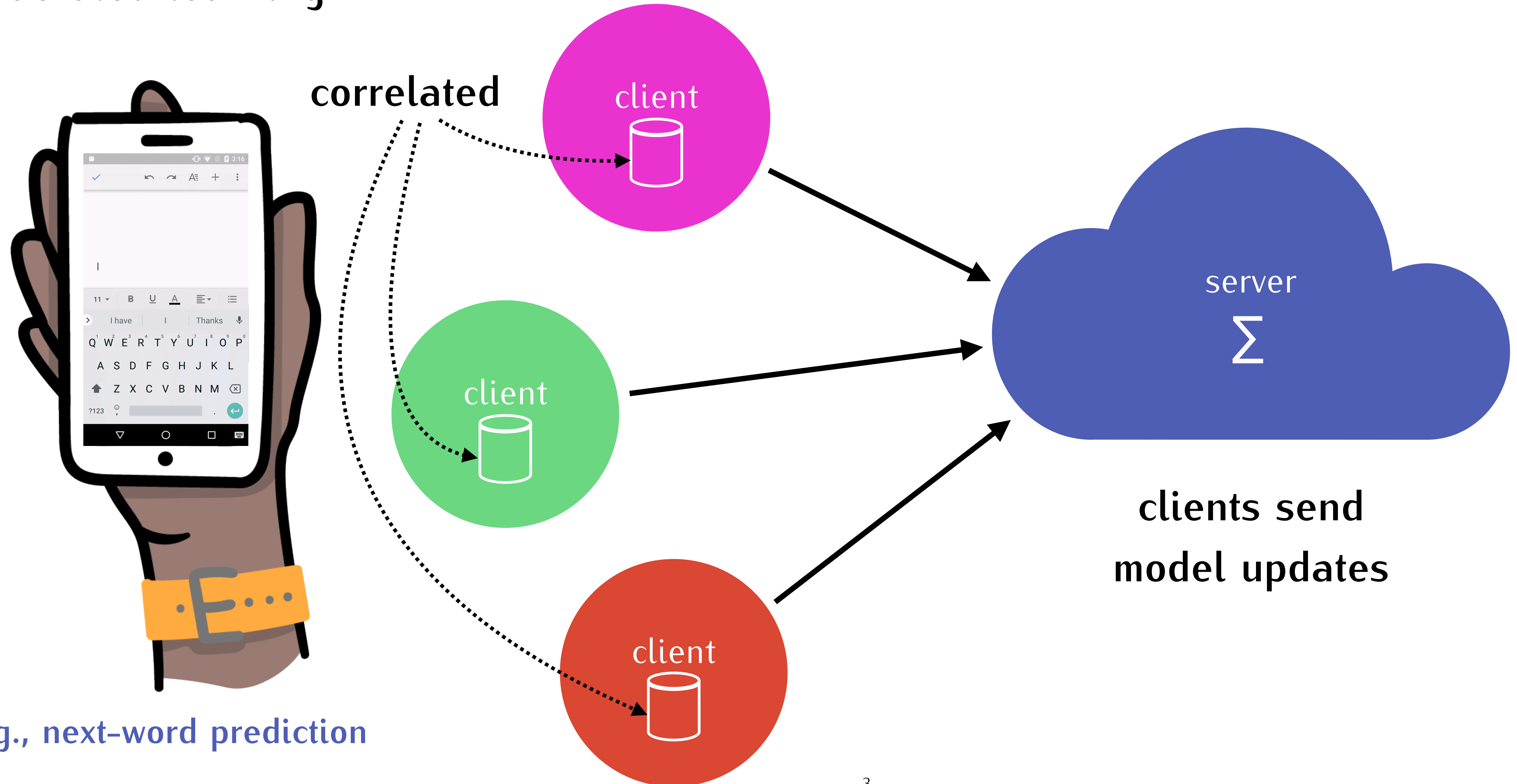


e.g., next-word prediction



Motivation: Distributed Source Coding

Federated learning.



50th year Commemorative Special Issue of the Transactions of Information Theory notes that

50th year Commemorative Special Issue of the Transactions of Information Theory notes that

“[...] despite the existence of potential applications, the conceptual importance of distributed source coding has not been mirrored in practical data compression.”

S. Verdú, “Fifty years of Shannon theory”, *IEEE Transactions on Information Theory*, 1998.

50th year Commemorative Special Issue of the Transactions of Information Theory notes that

“[...] despite the existence of potential applications, the conceptual importance of distributed source coding has not been mirrored in practical data compression.”

S. Verdú, “Fifty years of Shannon theory”, *IEEE Transactions on Information Theory*, 1998.

50th year Commemorative Special Issue of the Transactions of Information Theory notes that

“[...] despite the existence of potential applications, the conceptual importance of distributed source coding has not been mirrored in practical data compression.”

Still, the case after 25 years.

S. Verdú, “Fifty years of Shannon theory”, *IEEE Transactions on Information Theory*, 1998.

50th year Commemorative Special Issue of the Transactions of Information Theory notes that

“[...] despite the existence of potential applications, the conceptual importance of distributed source coding has not been mirrored in practical data compression.”

Still, the case after 25 years.

Particularly, for general sources.

S. Verdú, “Fifty years of Shannon theory”, *IEEE Transactions on Information Theory*, 1998.

50th year Commemorative Special Issue of the Transactions of Information Theory notes that

“[...] despite the existence of potential applications, the conceptual importance of distributed source coding has not been mirrored in practical data compression.”

Still, the case after 25 years.

Particularly, for *general sources*.

Learning-based compressors (e.g., Ballé et al., 2017) may help.

S. Verdú, “Fifty years of Shannon theory”, *IEEE Transactions on Information Theory*, 1998.

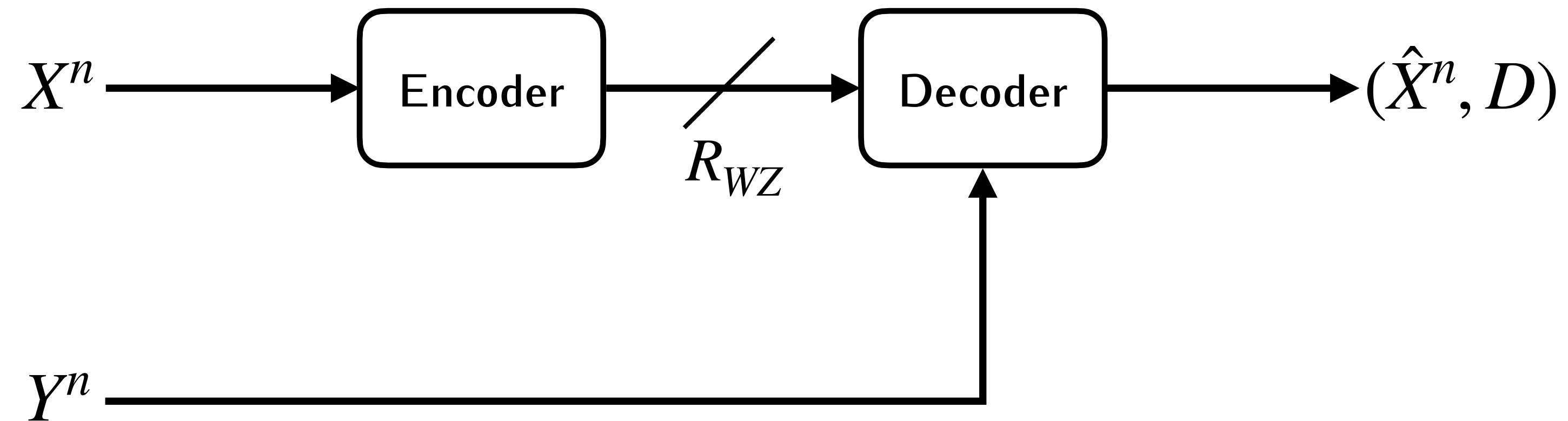
J. Ballé et al., “End-to-end Optimized Image Compression”, *International Conference on Learning Representations (ICLR)*, 2017.

Simpler special case: Rate-distortion (R-D) with side information

Simpler special case: Rate-distortion (R-D) with side information
Known as the Wyner-Ziv (WZ) problem.

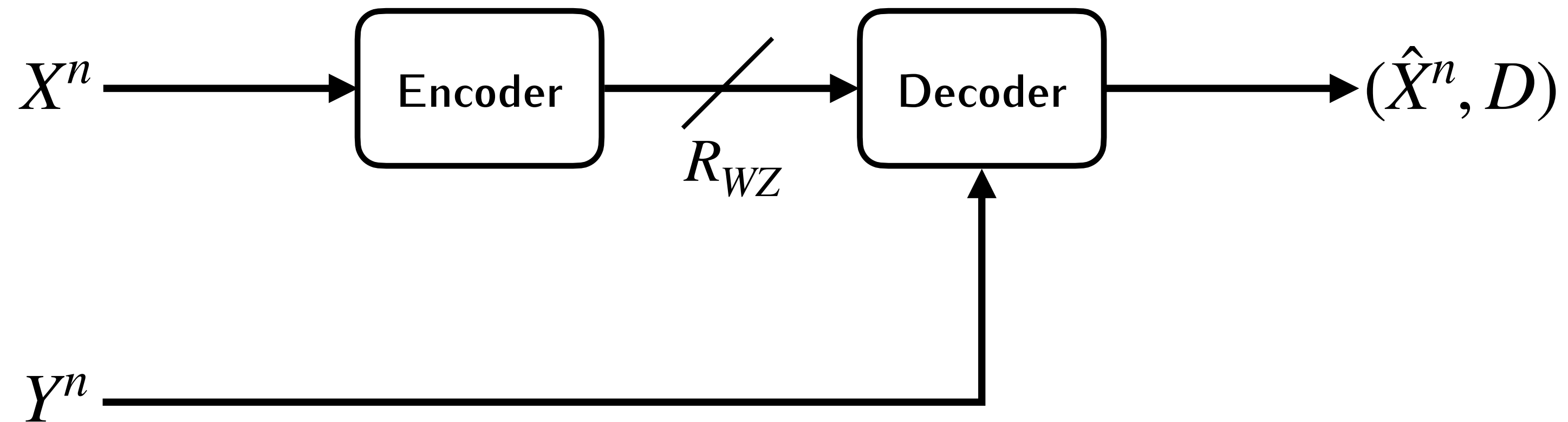
Simpler special case: Rate-distortion (R-D) with side information

Known as the Wyner-Ziv (WZ) problem.



Simpler special case: Rate-distortion (R-D) with side information

Known as the Wyner-Ziv (WZ) problem.



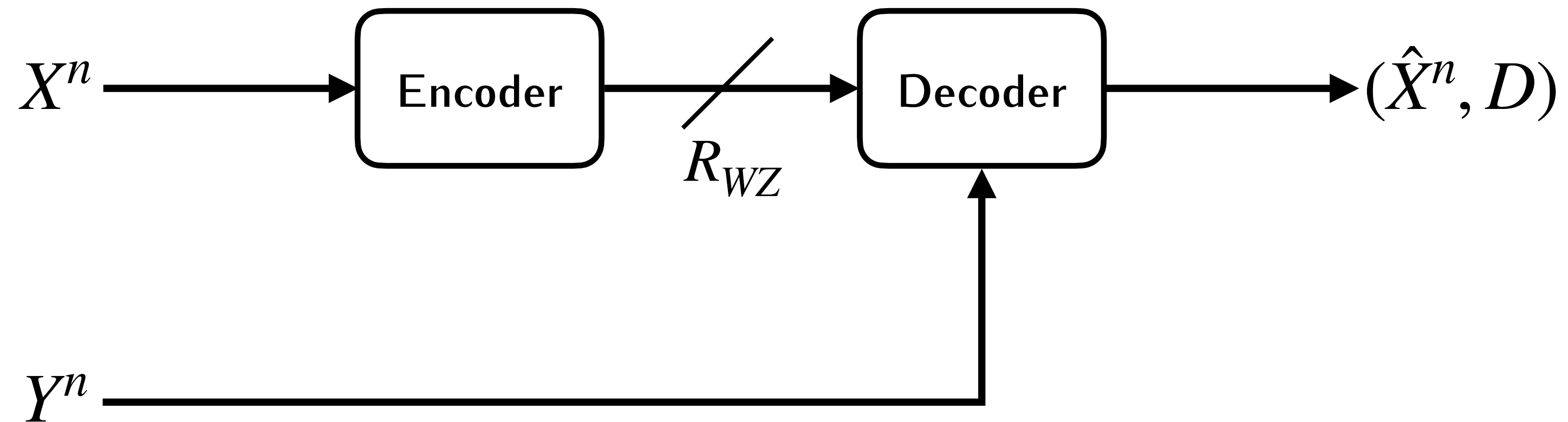
Theorem. Let (X, Y) be correlated i.i.d. $\sim p(x, y)$, and let $d(x, \hat{x})$ be a distortion measure. The R-D function for X when Y available at the decoder is:

$$R_{WZ}(D) = \min(I(X; U) - I(Y; U)),$$

where the minimization is over all $p(u|x)$ and all functions $g(u, y)$ satisfying $\mathbb{E}_{p(x,y)p(u|x)} d(x, g(u, y)) \leq D$.

Simpler special case: Rate-distortion (R-D) with side information

Known as the Wyner-Ziv (WZ) problem.



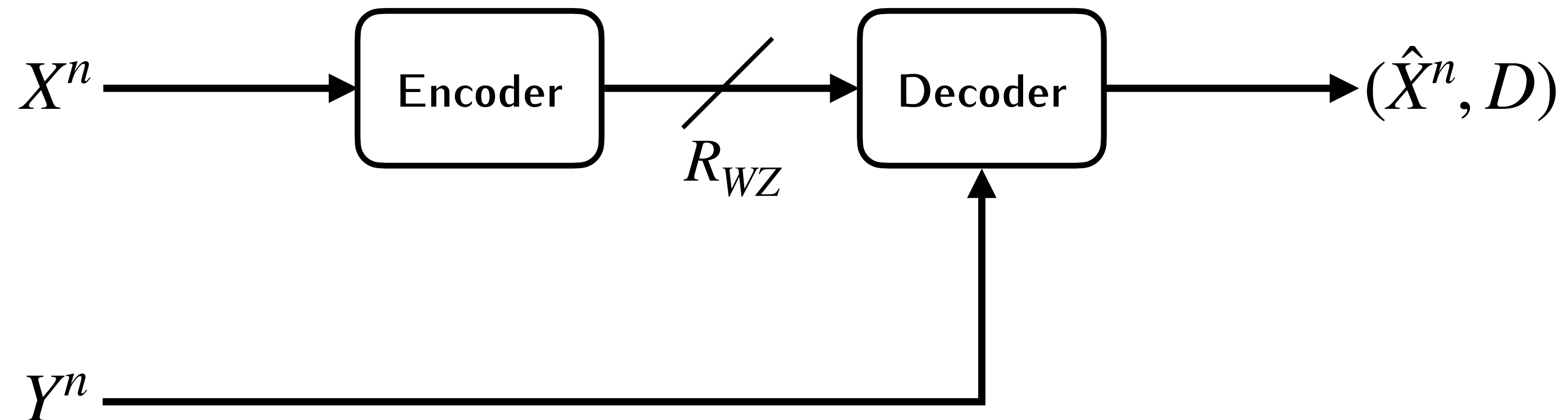
Theorem. Let (X, Y) be correlated i.i.d. $\sim p(x, y)$, and let $d(x, \hat{x})$ be a distortion measure. The R-D function for X when Y available at the decoder is:

$$R_{WZ}(D) = \min(I(X; U) - I(Y; U)),$$

where the minimization is over all $p(u|x)$ and all functions $g(u, y)$ satisfying $\mathbb{E}_{p(x,y)p(u|x)} d(x, g(u, y)) \leq D$.

Simpler special case: Rate-distortion (R-D) with side information

Known as the Wyner-Ziv (WZ) problem.



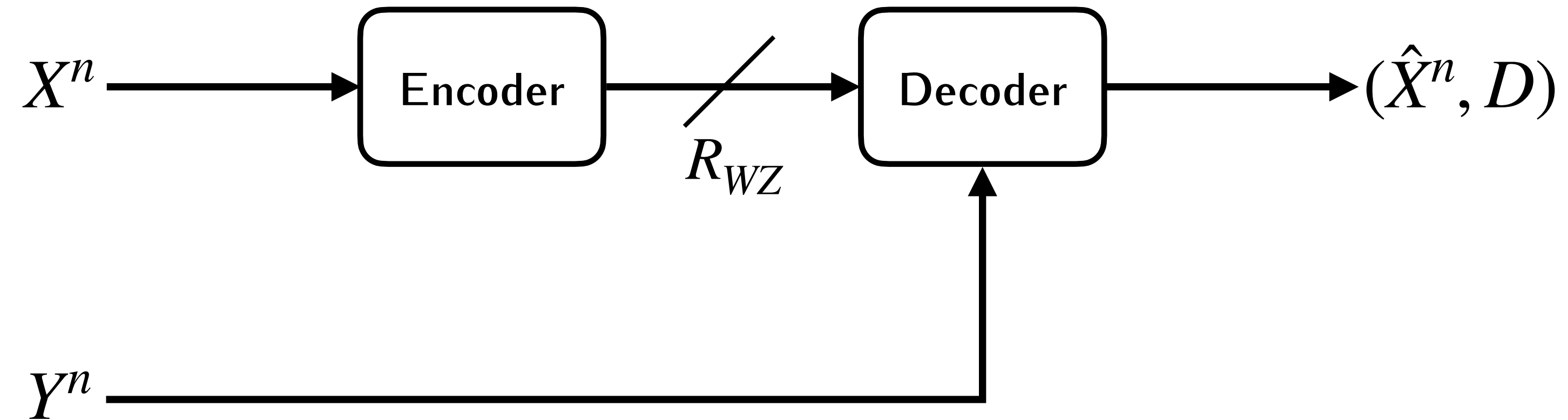
Theorem. Let (X, Y) be correlated i.i.d. $\sim p(x, y)$, and let $d(x, \hat{x})$ be a distortion measure. The R-D function for X when Y available at the decoder is:

$$R_{WZ}(D) = \min(I(X; U) - I(Y; U)),$$

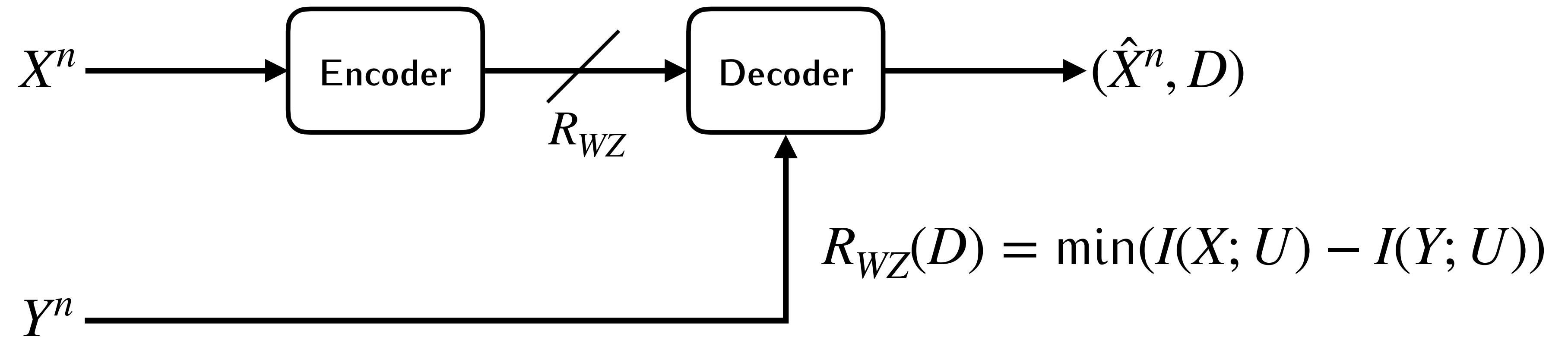
where the minimization is over all $p(u|x)$ and all functions $g(u, y)$ satisfying $\mathbb{E}_{p(x,y)p(u|x)} d(x, g(u, y)) \leq D$.

Wyner-Ziv achievability

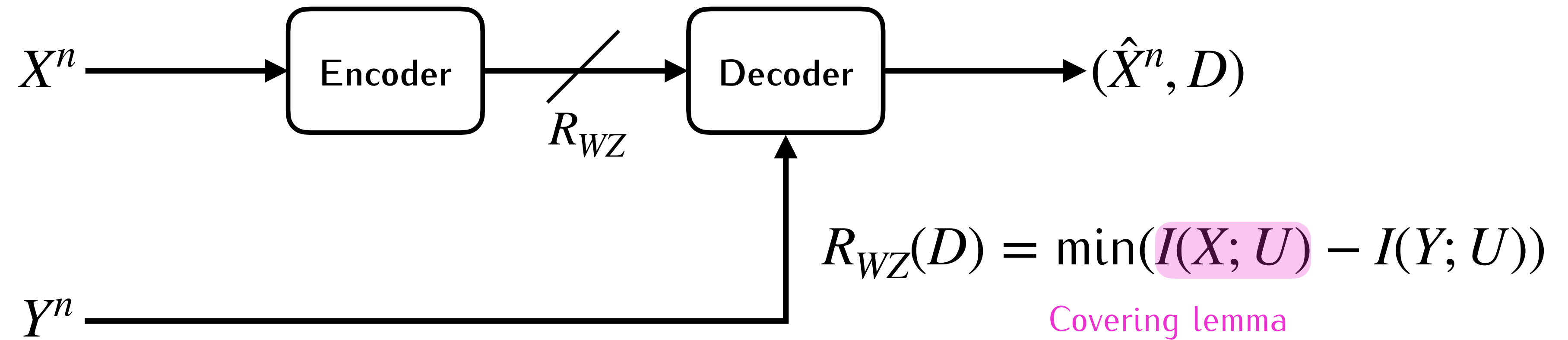
Wyner-Ziv achievability



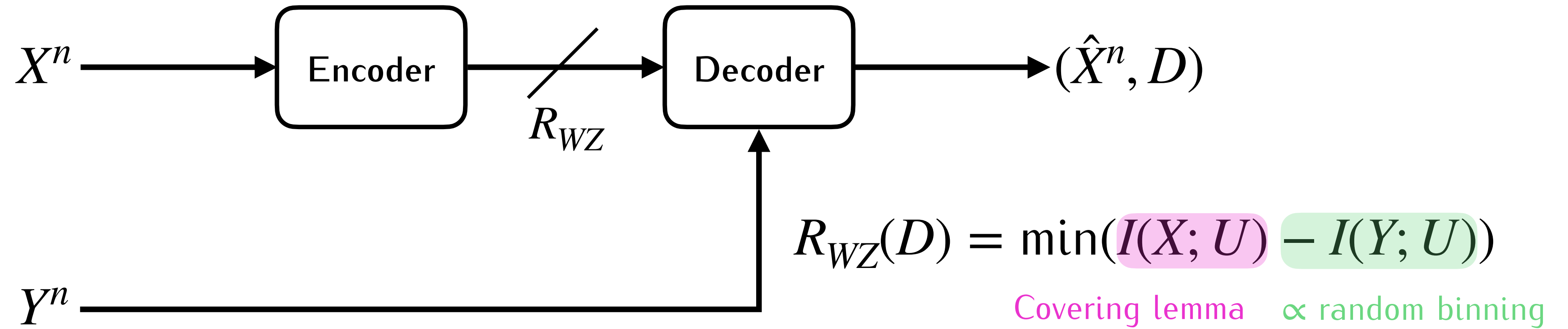
Wyner-Ziv achievability



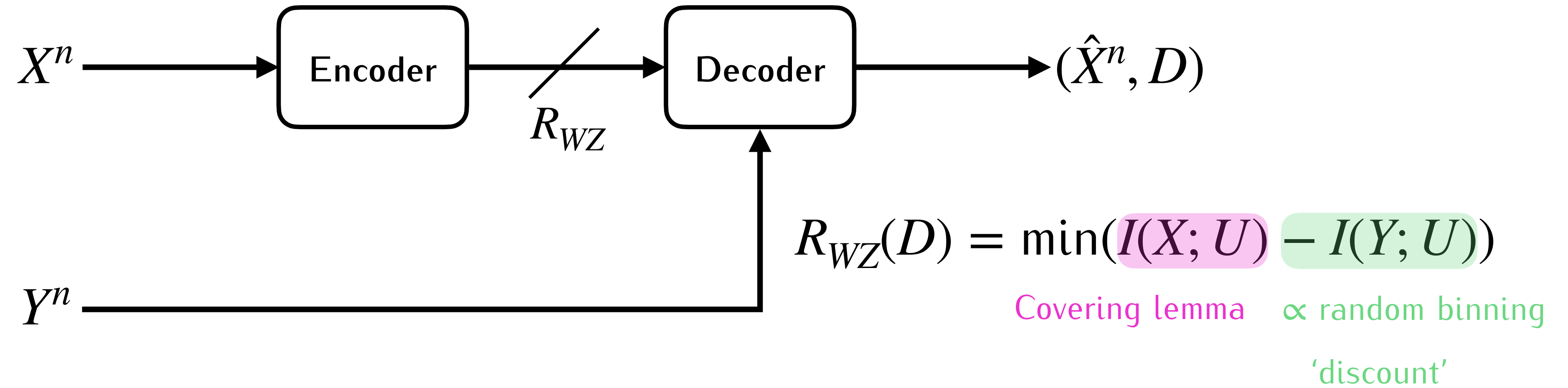
Wyner-Ziv achievability



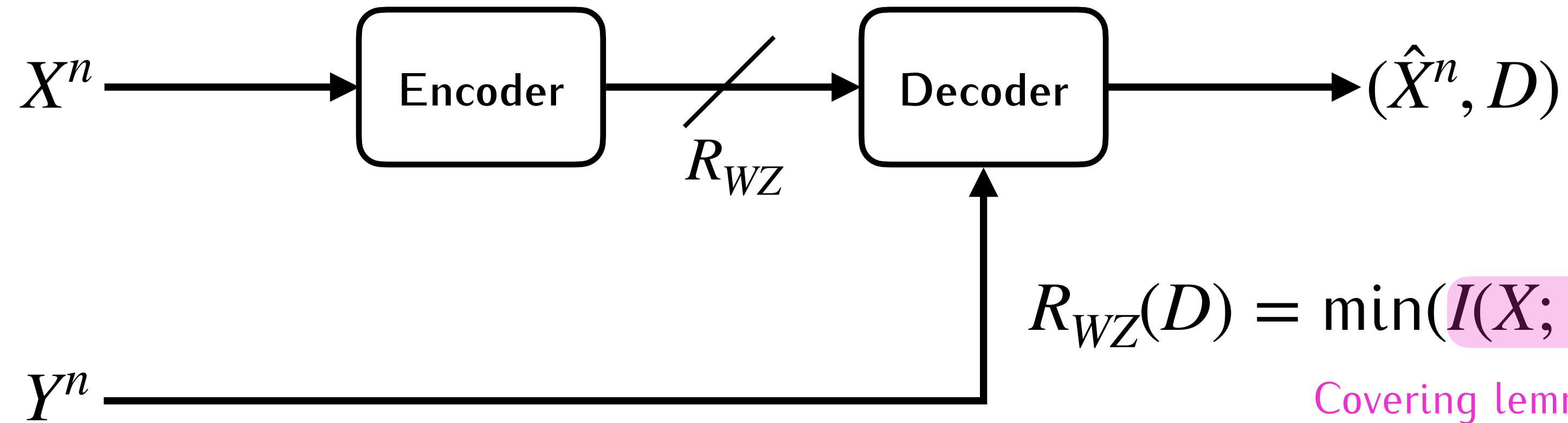
Wyner-Ziv achievability



Wyner-Ziv achievability

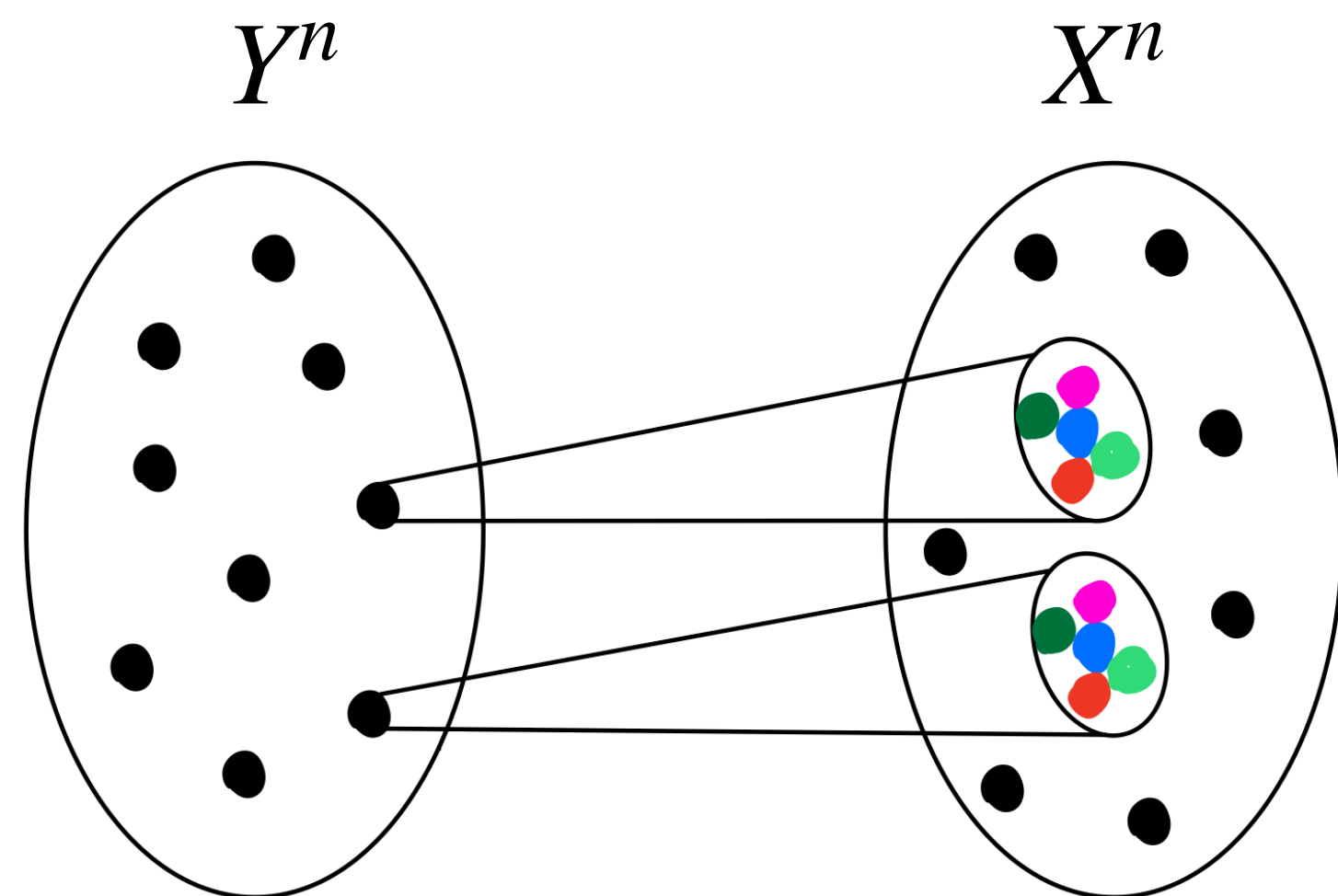


Wyner-Ziv achievability

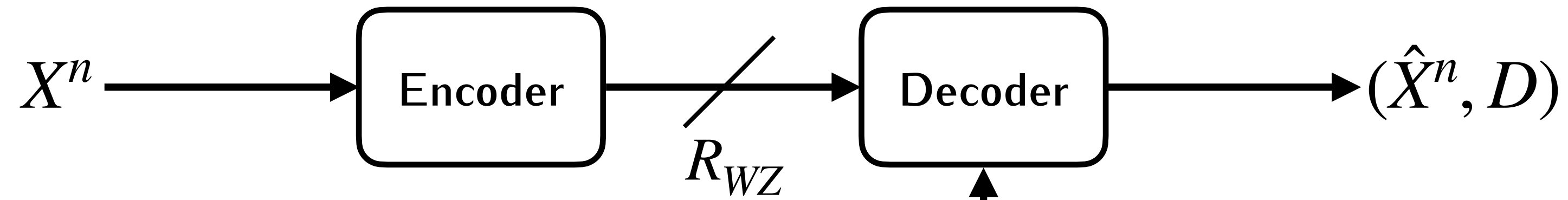


$$R_{WZ}(D) = \min(I(X; U) - I(Y; U))$$

Covering lemma \propto random binning
 'discount'



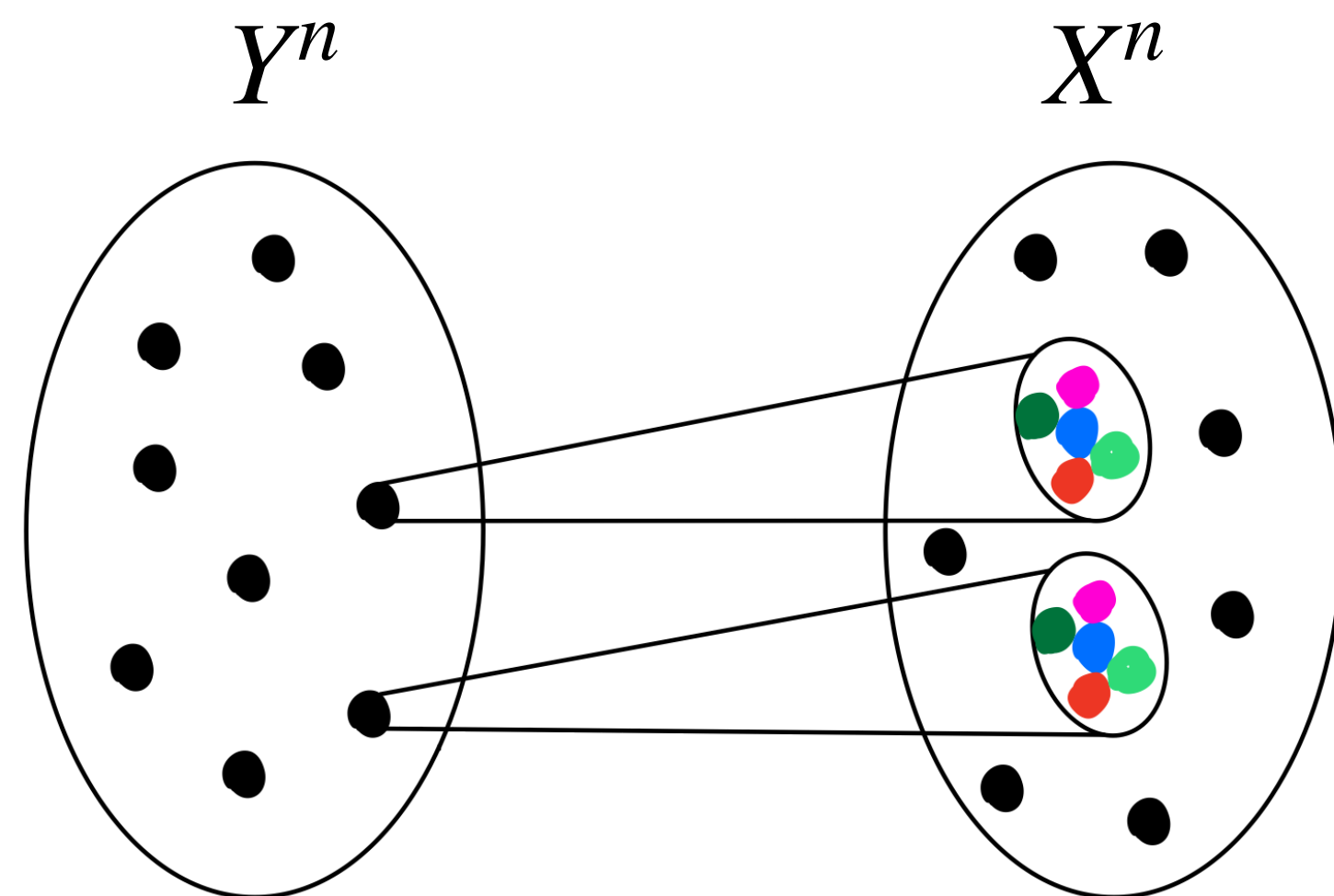
Wyner-Ziv achievability



$$R_{WZ}(D) = \min(I(X; U) - I(Y; U))$$

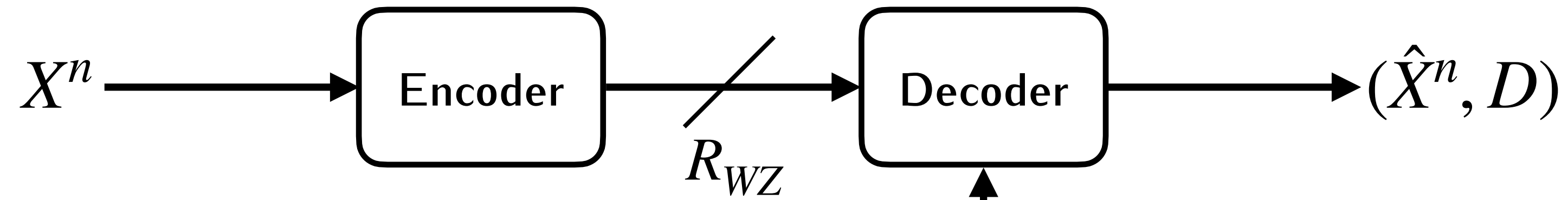
Covering lemma \propto random binning

'discount'



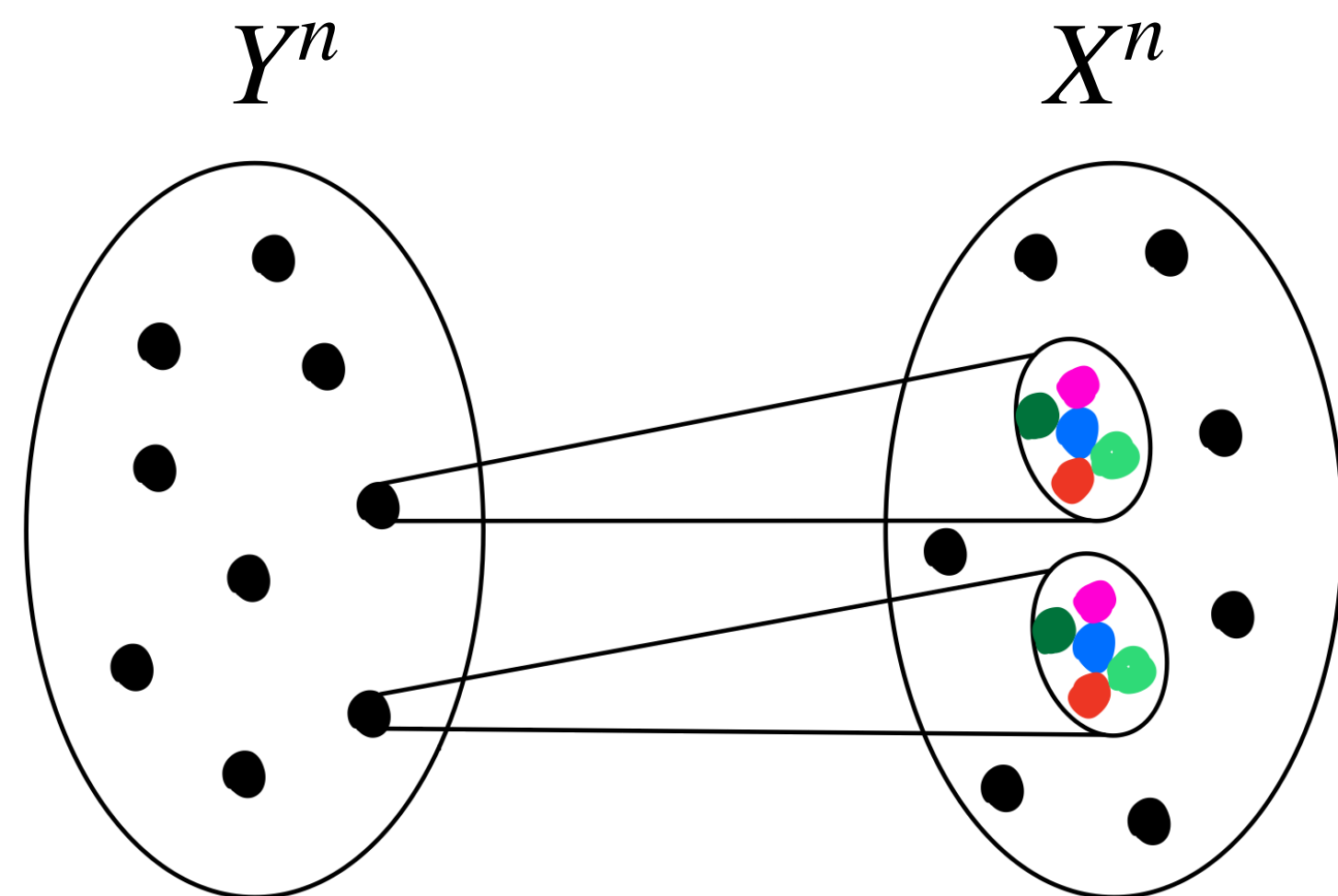
For X^n , send the color within the “fan”.

Wyner-Ziv achievability



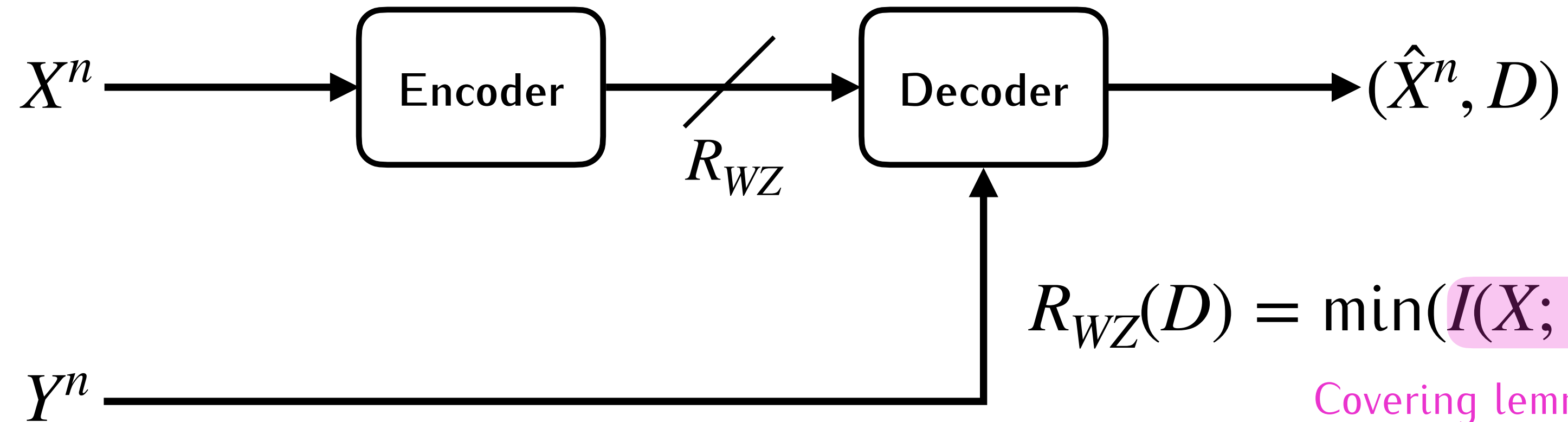
$$R_{WZ}(D) = \min(I(X; U) - I(Y; U))$$

Covering lemma \propto random binning
'discount'



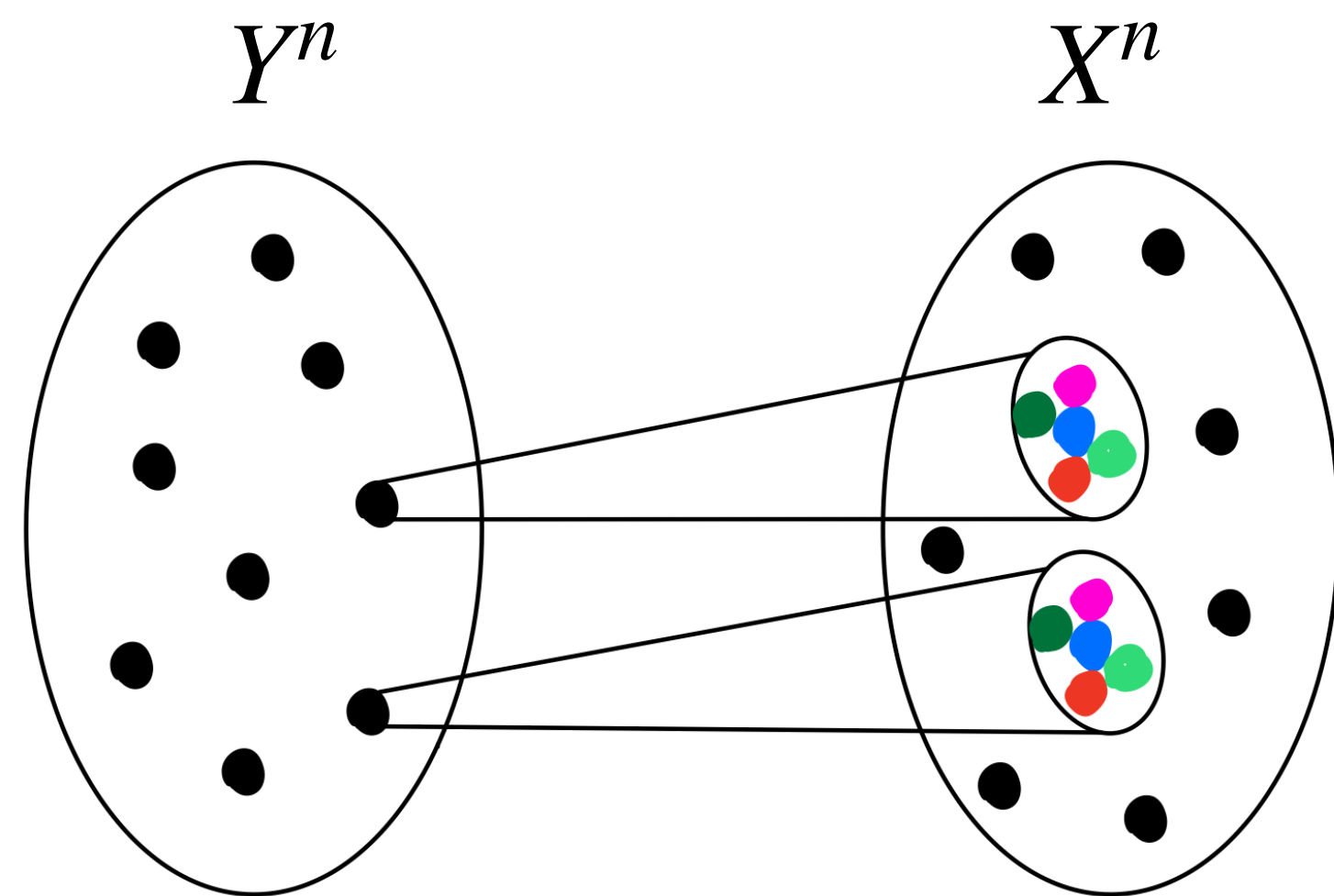
For X^n , send the color within the “fan”.
 \implies **binning**.

Wyner-Ziv achievability



$$R_{WZ}(D) = \min(I(X; U) - I(Y; U))$$

Covering lemma ∝ random binning
‘discount’



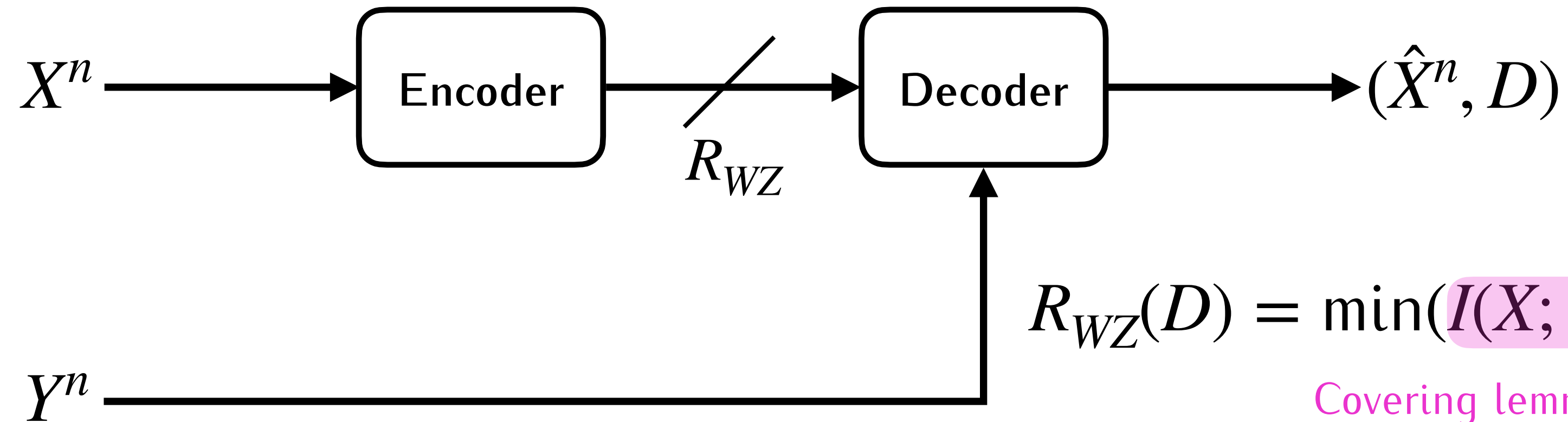
For X^n , send the color within the “fan”.

⇒ **binning**.

In quadratic-Gaussian setup,

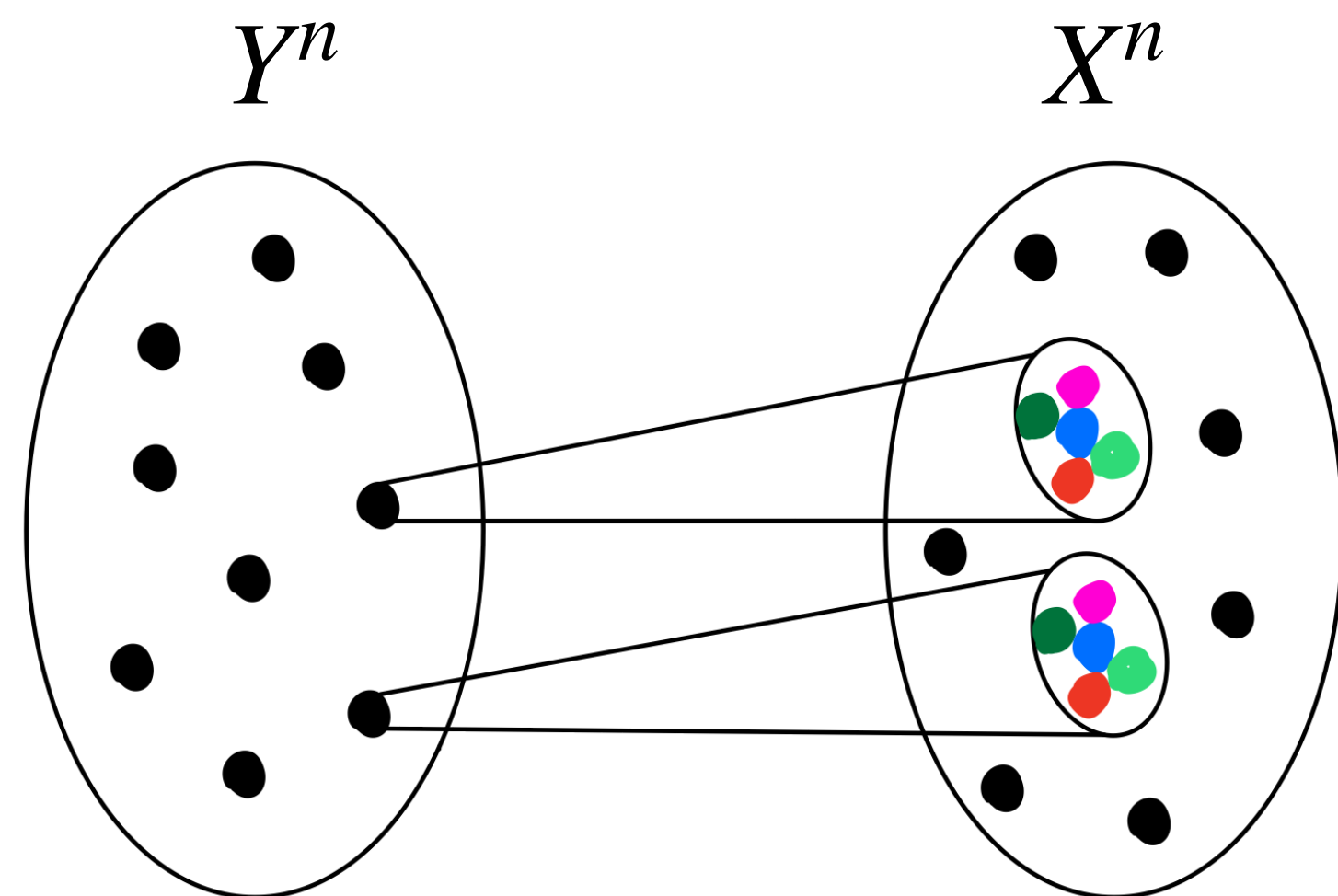
$$g(u, y) = a^*y + \beta[u - a^*y].$$

Wyner-Ziv achievability



$$R_{WZ}(D) = \min(I(X; U) - I(Y; U))$$

Covering lemma ∝ random binning
'discount'



For X^n , send the color within the “fan”.

⇒ **binning**.

In quadratic-Gaussian setup,
 $g(u, y) = a^*y + \beta[u - a^*y]$.

linear!

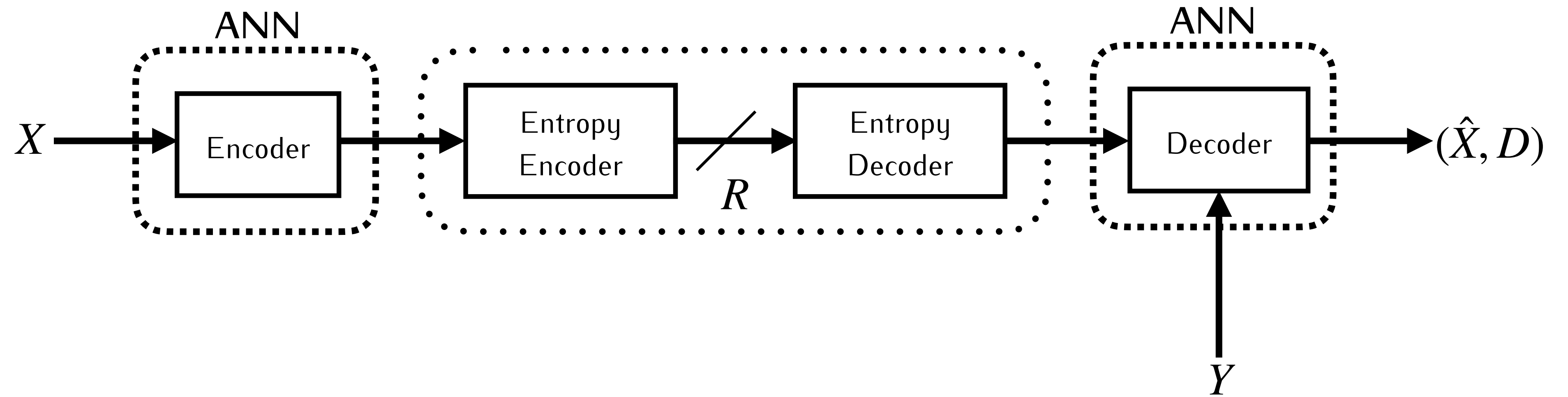
Operational schemes

Operational schemes

With Artificial Neural Networks (ANNs).

Operational schemes

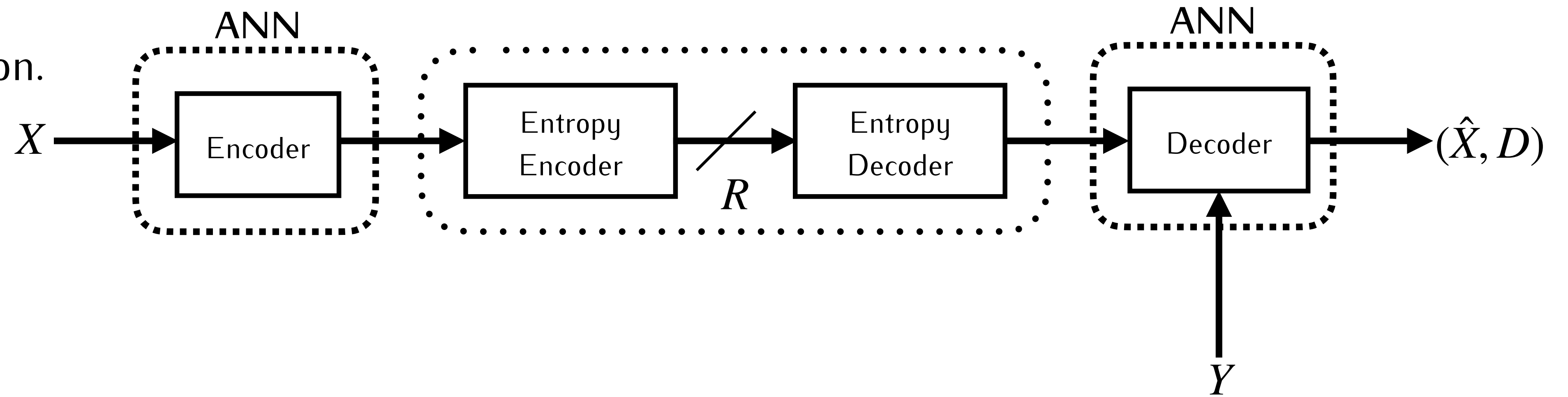
With Artificial Neural Networks (ANNs).



Operational schemes

With Artificial Neural Networks (ANNs).

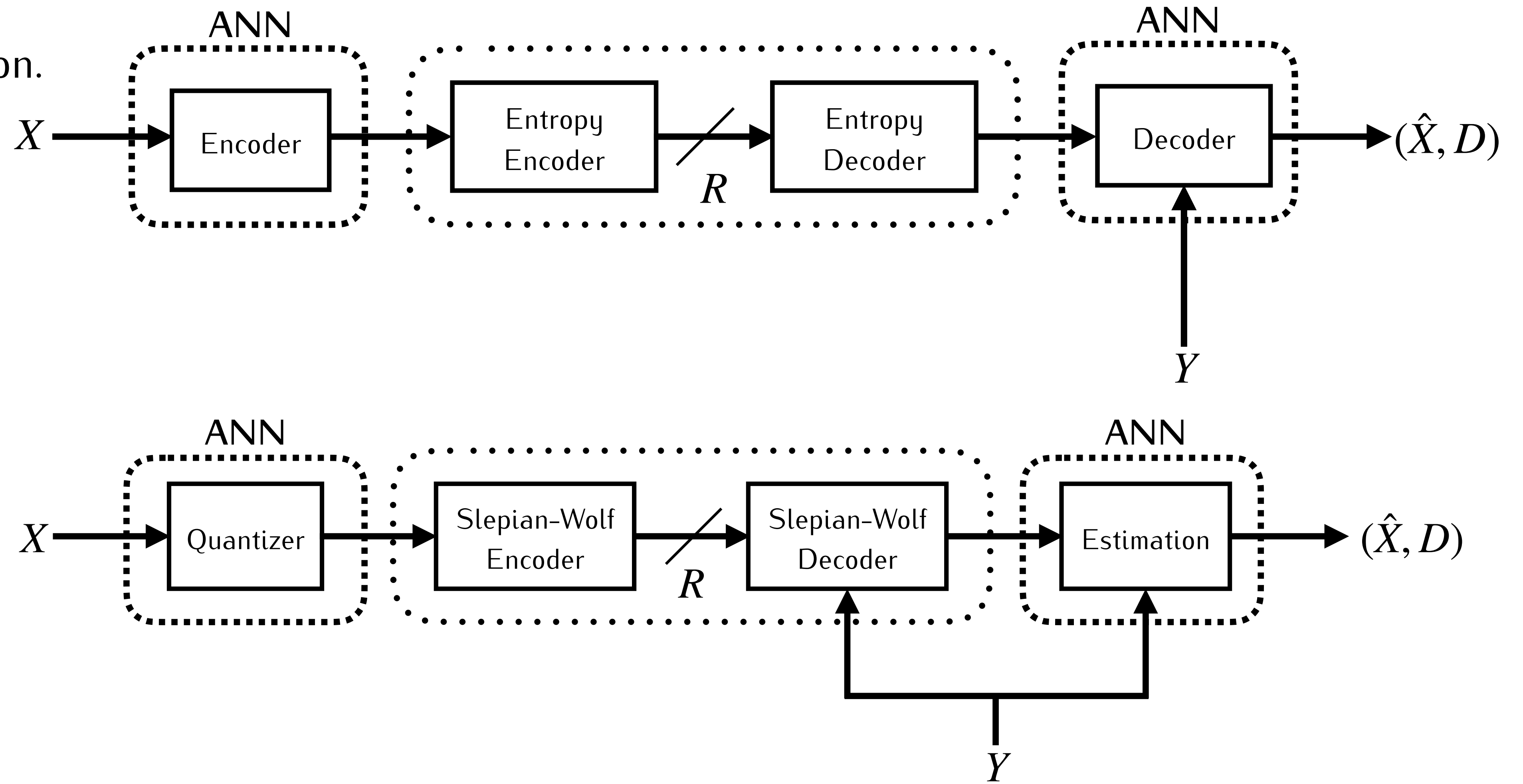
Marginal formulation.



Operational schemes

With Artificial Neural Networks (ANNs).

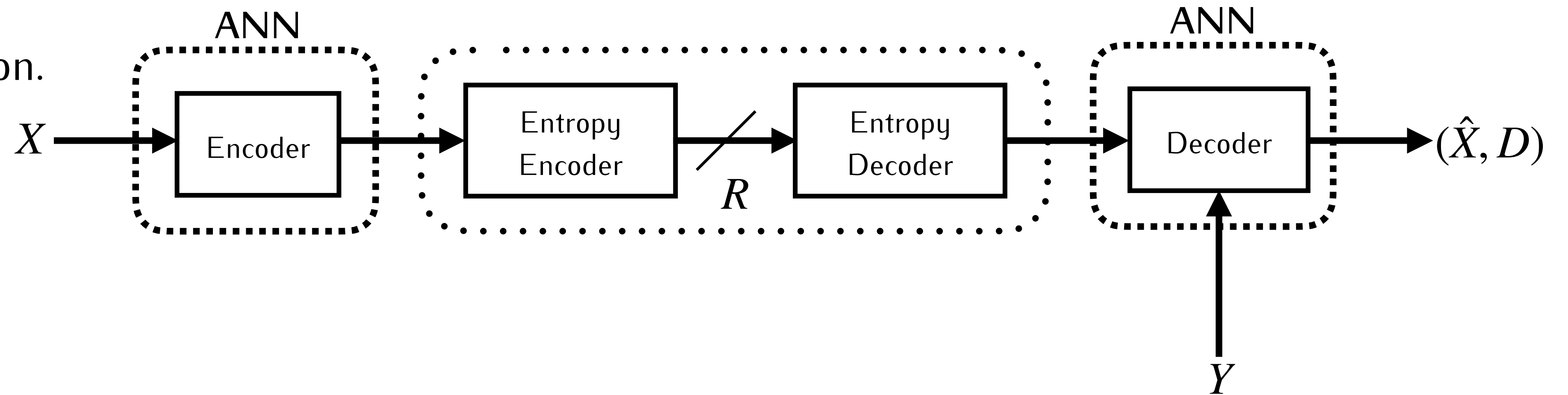
Marginal formulation.



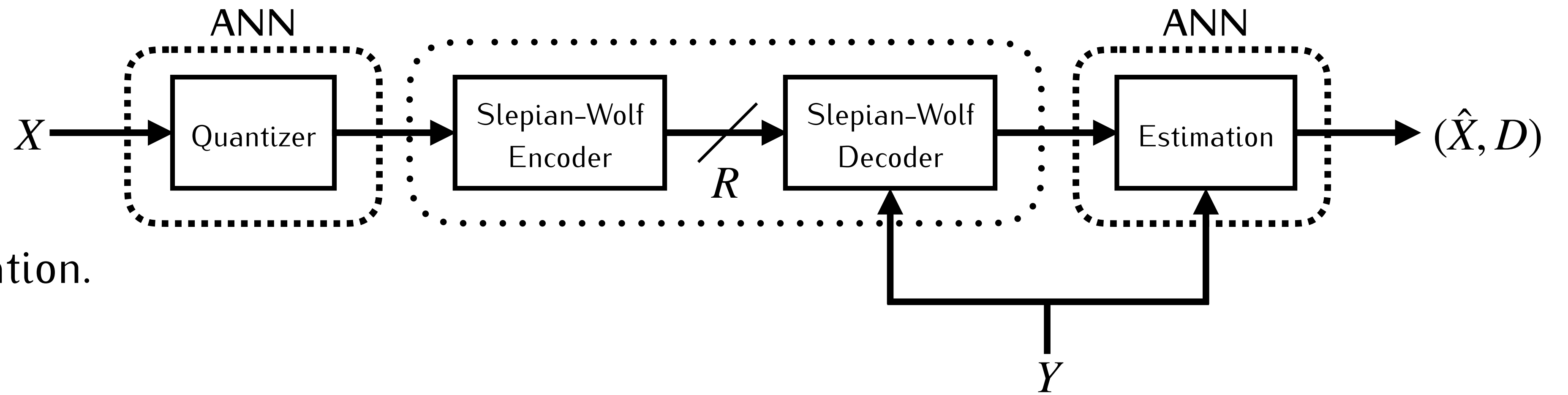
Operational schemes

With Artificial Neural Networks (ANNs).

Marginal formulation.



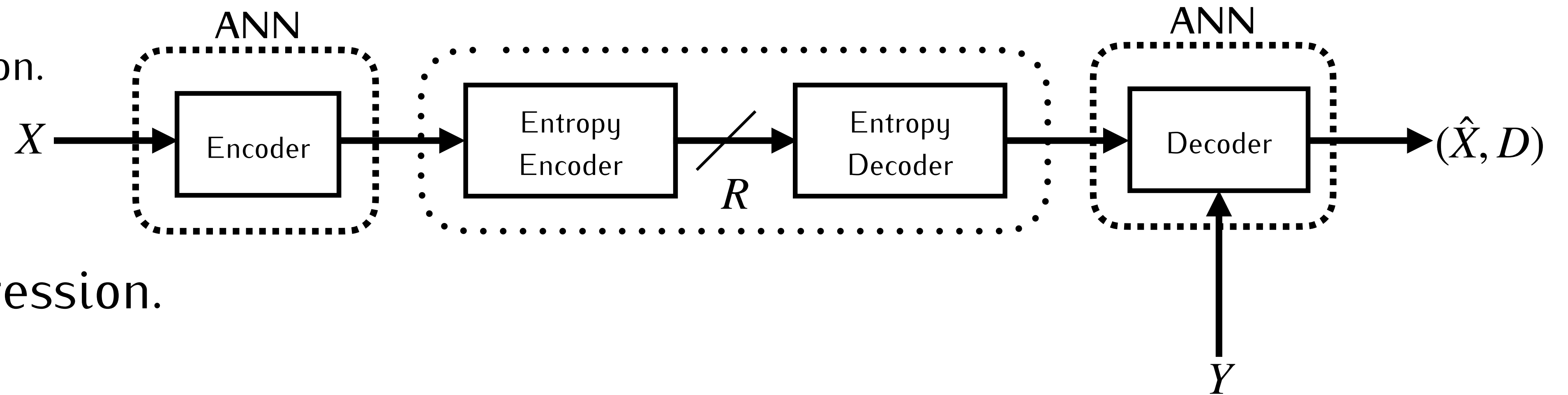
Conditional formulation.



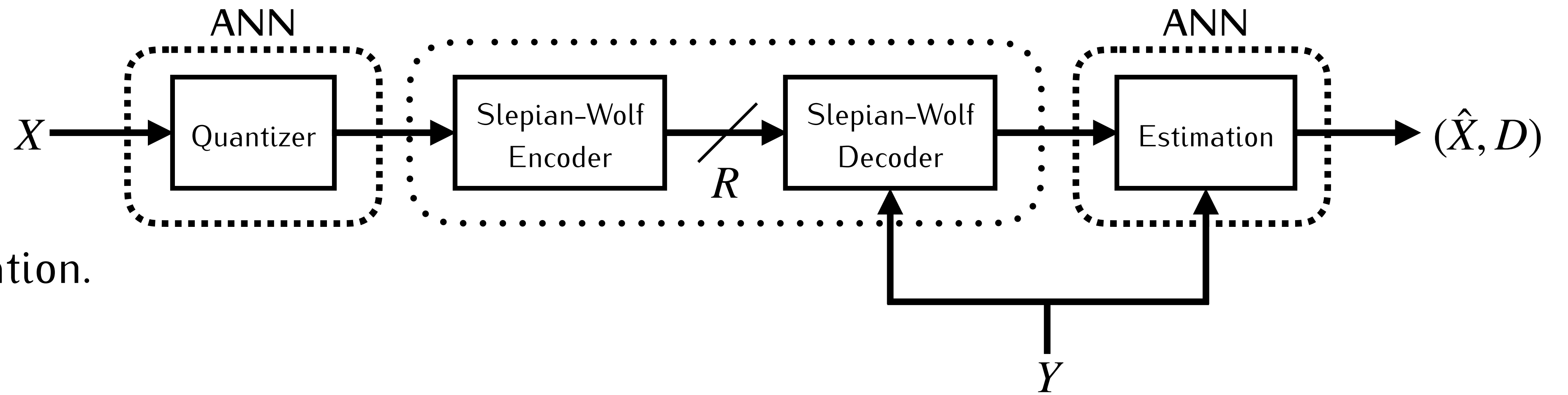
Operational schemes

With Artificial Neural Networks (ANNs).

Marginal formulation.



One-shot compression.

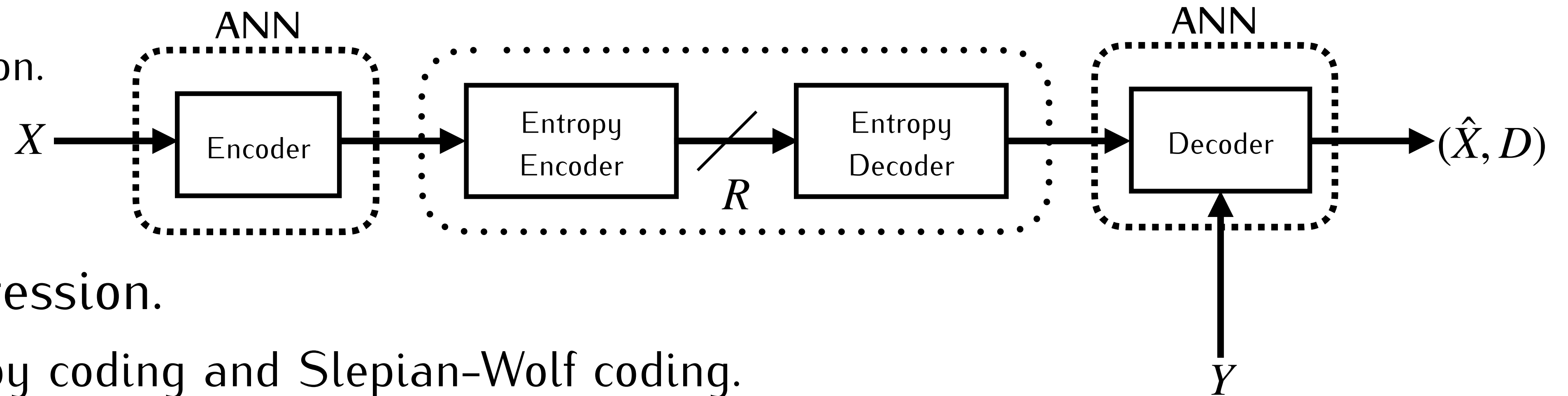


Conditional formulation.

Operational schemes

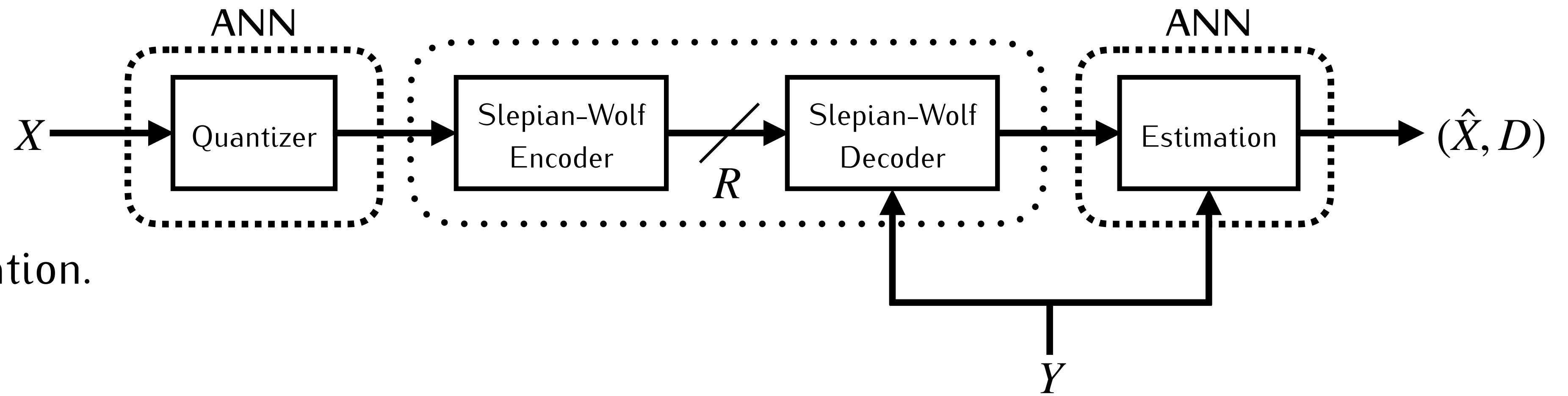
With Artificial Neural Networks (ANNs).

Marginal formulation.



One-shot compression.

High-order entropy coding and Slepian-Wolf coding.



Conditional formulation.

Neural parametrization for Wyner-Ziv

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

$$I(X; U) - I(Y; U) = I(X; U | Y) = \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{p(u|y)} \right].$$

$U - X - Y$

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

$$I(X; U) - I(Y; U) = I(X; U | Y) = \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{p(u|y)} \right].$$

$U - X - Y$

- For test time, set encoder output as $u = \operatorname{argmax}_v p_{\theta}(v|x)$, and have U as discrete.

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

$$I(X; U) - I(Y; U) = I(X; U | Y) = \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{p(u|y)} \right].$$

$U - X - Y$

- For test time, set encoder output as $u = \operatorname{argmax}_v p_{\theta}(v|x)$, and have U as discrete.
- Choose one of two variational upper bounds:

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

$$I(X; U) - I(Y; U) = I(X; U | Y) = \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{p(u|y)} \right].$$

$U - X - Y$

- For test time, set encoder output as $u = \operatorname{argmax}_v p_{\theta}(v|x)$, and have U as discrete.
- Choose one of two variational upper bounds:

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u)} \right],$$

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u|y)} \right].$$

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

$$I(X; U) - I(Y; U) = I(X; U | Y) = \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{p(u|y)} \right].$$

$U - X - Y$

- For test time, set encoder output as $u = \operatorname{argmax}_v p_{\theta}(v|x)$, and have U as discrete.
- Choose one of two variational upper bounds:

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u)} \right],$$

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u|y)} \right].$$

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

$$I(X; U) - I(Y; U) = I(X; U | Y) = \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{p(u|y)} \right].$$

$U - X - Y$

- For test time, set encoder output as $u = \operatorname{argmax}_v p_{\theta}(v|x)$, and have U as discrete.
- Choose one of two variational upper bounds:

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u)} \right], \quad \text{marginal}$$

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u|y)} \right].$$

Neural parametrization for Wyner-Ziv

- Assume that during training, the encoder in achievability is represented by $p_{\theta}(u|x)$,

$$I(X; U) - I(Y; U) = I(X; U | Y) = \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{p(u|y)} \right].$$

$U - X - Y$

- For test time, set encoder output as $u = \operatorname{argmax}_v p_{\theta}(v|x)$, and have U as discrete.
- Choose one of two variational upper bounds:

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u)} \right], \quad \text{marginal}$$

$$I(X; U | Y) \leq \mathbb{E} \left[\log \frac{p_{\theta}(u|x)}{q_{\xi}(u|y)} \right]. \quad \text{conditional}$$

Neural parametrization for Wyner-Ziv

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{p_\theta(u|x)}{q_\xi(u)} + \lambda \cdot d(x, g_\phi(u, y)) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{p_\theta(u|x)}{q_\zeta(u|y)} + \lambda \cdot d(x, g_\phi(u, y)) \right].$$

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{p_\theta(u|x)}{q_\xi(u)} + \lambda \cdot d(x, g_\phi(u, y)) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{p_\theta(u|x)}{q_\zeta(u|y)} + \lambda \cdot d(x, g_\phi(u, y)) \right].$$

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{p_\theta(u|x)}{q_\xi(u)} + \lambda \cdot d(x, g_\phi(u, y)) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{p_\theta(u|x)}{q_\zeta(u|y)} + \lambda \cdot d(x, g_\phi(u, y)) \right].$$

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{q_\xi(u)} + \lambda \cdot d(x, g_\phi(u, y)) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{\overset{\text{quantizer}}{p_\theta(u|x)}}{q_\zeta(u|y)} + \lambda \cdot d(x, g_\phi(u, y)) \right].$$

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{q_\xi(u)} + \lambda \cdot d(x, \overset{\text{decoder}}{g_\phi(u, y)}) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{\overset{\text{quantizer}}{p_\theta(u|x)}}{q_\zeta(u|y)} + \lambda \cdot d(x, \overset{\text{de-quantizer}}{g_\phi(u, y)}) \right].$$

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\xi(u)}} + \lambda \cdot d(x, \overset{\text{decoder}}{g_\phi(u, y)}) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{\overset{\text{quantizer}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\zeta(u|y)}} + \lambda \cdot d(x, \overset{\text{de-quantizer}}{g_\phi(u, y)}) \right].$$

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\xi(u)}} + \lambda \cdot d(x, \underset{\text{decoder}}{g_\phi(u, y)}) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{\overset{\text{quantizer}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\zeta(u|y)}} + \lambda \cdot d(x, \underset{\text{de-quantizer}}{g_\phi(u, y)}) \right].$$

- Define all models $p_\theta(u|x)$, $q_\xi(u)$ and $q_\zeta(u|y)$ as **discrete** distributions with probabilities:

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\xi(u)}} + \lambda \cdot d(x, \underset{\text{decoder}}{g_\phi(u, y)}) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\zeta(u|y)}} + \lambda \cdot d(x, \underset{\text{decoder}}{g_\phi(u, y)}) \right].$$

- Define all models $p_\theta(u|x)$, $q_\xi(u)$ and $q_\zeta(u|y)$ as **discrete** distributions with probabilities:

$$P_k = \frac{\exp \alpha_k}{\sum_{i=1}^K \exp \alpha_i}.$$

Neural parametrization for Wyner-Ziv

- Relax the constrained formulation of Wyner-Ziv theorem using Lagrange multipliers:

$$L_m(\theta, \phi, \xi) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\xi(u)}} + \lambda \cdot d(x, \underset{\text{decoder}}{g_\phi(u, y)}) \right],$$

$$L_c(\theta, \phi, \zeta) = \mathbb{E} \left[\log \frac{\overset{\text{encoder}}{p_\theta(u|x)}}{\underset{\text{entropy coder}}{q_\zeta(u|y)}} + \lambda \cdot d(x, \underset{\text{de-quantizer}}{g_\phi(u, y)}) \right].$$

- Define all models $p_\theta(u|x)$, $q_\xi(u)$ and $q_\zeta(u|y)$ as **discrete** distributions with probabilities:

$$P_k = \frac{\exp \alpha_k}{\sum_{i=1}^K \exp \alpha_i}.$$

- This keeps the parametric families as general as possible, and does not impose any structure.

Neural parametrization for Wyner-Ziv

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

For example, $\frac{\partial}{\partial \boldsymbol{\theta}} \mathbb{E}[l_{\boldsymbol{\theta}}(x, y)] \approx \frac{1}{|B|} \sum_{(x, y) \in B} \frac{\partial l_{\boldsymbol{\theta}}(x, y)}{\partial \boldsymbol{\theta}}$, where $l_{\boldsymbol{\theta}}$ is a sample loss with parameters $\boldsymbol{\theta}$.

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

For example, $\frac{\partial}{\partial \theta} \mathbb{E}[l_{\theta}(x, y)] \approx \frac{1}{|B|} \sum_{(x, y) \in B} \frac{\partial l_{\theta}(x, y)}{\partial \theta}$, where l_{θ} is a sample loss with parameters θ .

- To draw samples u from $p_{\theta}(u | x)$, use Gumbel-max 'trick' that is:

E. J. Gumbel, "Statistical theory of extreme values and some practical applications: a series of lectures", *US Department of Commerce*, 1954.

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

For example, $\frac{\partial}{\partial \theta} \mathbb{E}[l_{\theta}(x, y)] \approx \frac{1}{|B|} \sum_{(x, y) \in B} \frac{\partial l_{\theta}(x, y)}{\partial \theta}$, where l_{θ} is a sample loss with parameters θ .

- To draw samples u from $p_{\theta}(u | x)$, use Gumbel-max 'trick' that is:

$$\arg \max_{k \in 1, \dots, K} \{ \alpha_k + G_k \} .$$

E. J. Gumbel, "Statistical theory of extreme values and some practical applications: a series of lectures", *US Department of Commerce*, 1954.

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

For example, $\frac{\partial}{\partial \theta} \mathbb{E}[l_{\theta}(x, y)] \approx \frac{1}{|B|} \sum_{(x, y) \in B} \frac{\partial l_{\theta}(x, y)}{\partial \theta}$, where l_{θ} is a sample loss with parameters θ .

- To draw samples u from $p_{\theta}(u | x)$, use Gumbel-max 'trick' that is:

$$\arg \max_{k \in 1, \dots, K} \{\alpha_k + G_k\} .$$

- **Problem:** the derivative of $\arg \max$ is 0 almost everywhere.

E. J. Gumbel, "Statistical theory of extreme values and some practical applications: a series of lectures", *US Department of Commerce*, 1954.

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

For example, $\frac{\partial}{\partial \theta} \mathbb{E}[l_{\theta}(x, y)] \approx \frac{1}{|B|} \sum_{(x, y) \in B} \frac{\partial l_{\theta}(x, y)}{\partial \theta}$, where l_{θ} is a sample loss with parameters θ .

- To draw samples u from $p_{\theta}(u | x)$, use Gumbel-max 'trick' that is:

$$\arg \max_{k \in 1, \dots, K} \{ \alpha_k + G_k \} .$$

- **Problem:** the derivative of $\arg \max$ is 0 almost everywhere.
- Need continuous relaxation of $\arg \max$ during training.

E. J. Gumbel, "Statistical theory of extreme values and some practical applications: a series of lectures", *US Department of Commerce*, 1954.

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

For example, $\frac{\partial}{\partial \theta} \mathbb{E}[l_{\theta}(x, y)] \approx \frac{1}{|B|} \sum_{(x, y) \in B} \frac{\partial l_{\theta}(x, y)}{\partial \theta}$, where l_{θ} is a sample loss with parameters θ .

- To draw samples u from $p_{\theta}(u | x)$, use Gumbel-max ‘trick’ that is:

$$\arg \max_{k \in 1, \dots, K} \{\alpha_k + G_k\} .$$

- **Problem:** the derivative of $\arg \max$ is 0 almost everywhere.
- Need continuous relaxation of $\arg \max$ during training.
 - Opt for *softmax* (differentiable!).

E. J. Gumbel, “Statistical theory of extreme values and some practical applications: a series of lectures”, *US Department of Commerce*, 1954.

Neural parametrization for Wyner-Ziv

- Optimize learnable parameters with stochastic gradient descent (SGD).
- SGD replaces $\mathbb{E}(\cdot)$ by averages over batches of samples B .

For example, $\frac{\partial}{\partial \theta} \mathbb{E}[l_{\theta}(x, y)] \approx \frac{1}{|B|} \sum_{(x, y) \in B} \frac{\partial l_{\theta}(x, y)}{\partial \theta}$, where l_{θ} is a sample loss with parameters θ .

- To draw samples u from $p_{\theta}(u | x)$, use Gumbel-max ‘trick’ that is:

$$\arg \max_{k \in 1, \dots, K} \{ \alpha_k + G_k \} .$$

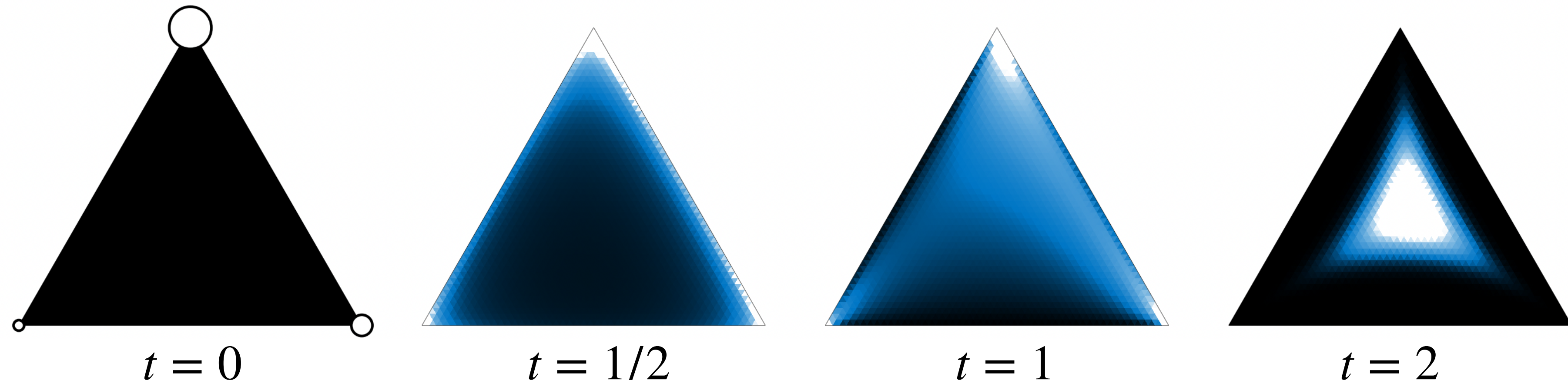
- **Problem:** the derivative of $\arg \max$ is 0 almost everywhere.
- Need continuous relaxation of $\arg \max$ during training.
 - Opt for *softmax* (differentiable!).
 - Use Gumbel-softmax ‘trick’ by Maddison et al.

E. J. Gumbel, “Statistical theory of extreme values and some practical applications: a series of lectures”, *US Department of Commerce*, 1954.

C. Maddison et al., “The concrete distribution: a continuous relaxation of discrete random variables”, *ICLR*, 2017.

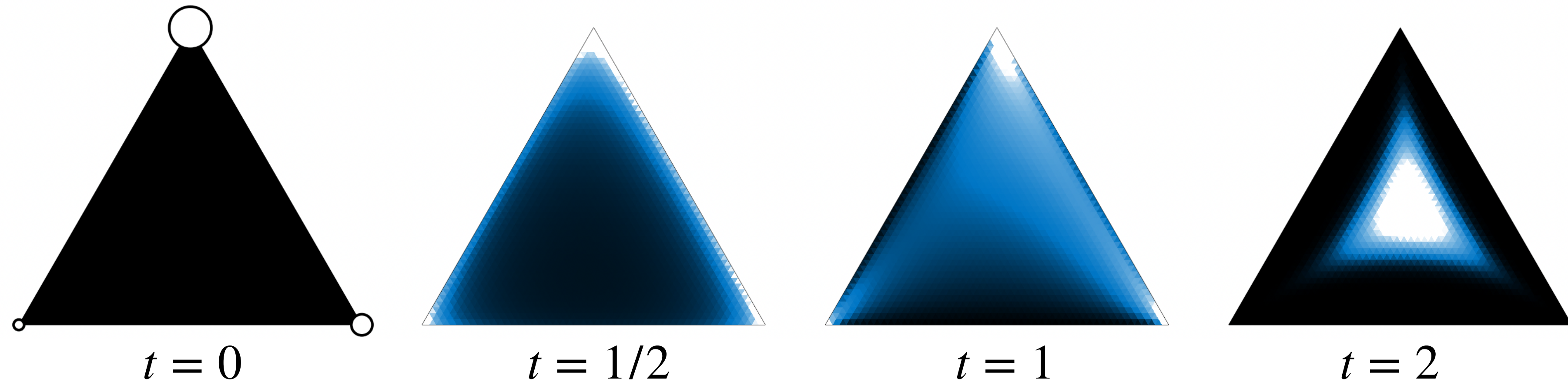
Gumbel-softmax 'trick'

Gumbel-softmax 'trick'



- Concrete distribution (with temperature t) relaxes sampling from a discrete distribution.

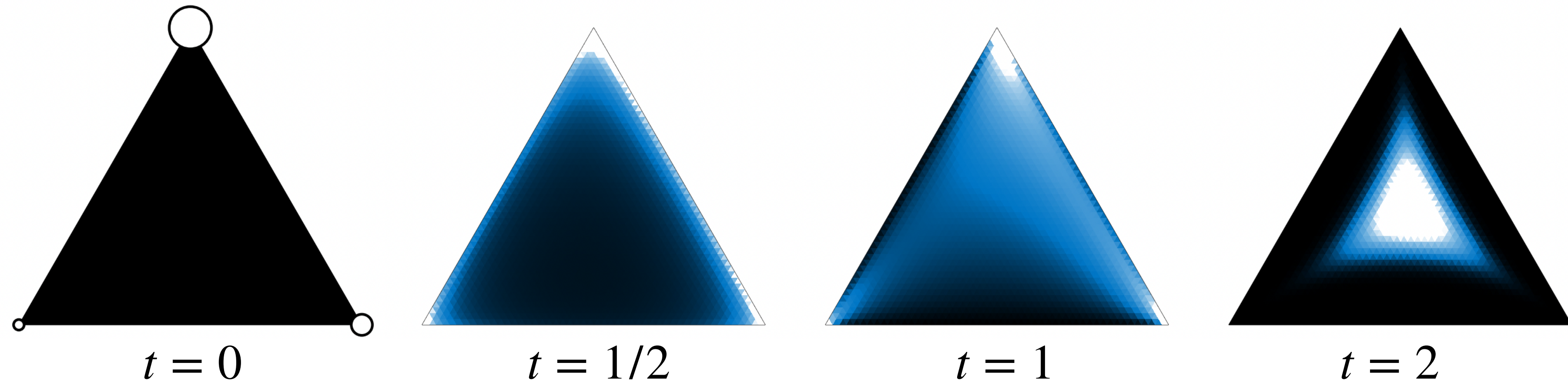
Gumbel-softmax 'trick'



- Concrete distribution (with temperature t) relaxes sampling from a discrete distribution.
- Rather than sampling an index U , sample a vector \mathbf{U} :

$$U_k = \frac{\exp((\alpha_k + G_k) / t)}{\sum_{i=1}^n \exp((\alpha_i + G_i) / t)} .$$

Gumbel-softmax 'trick'

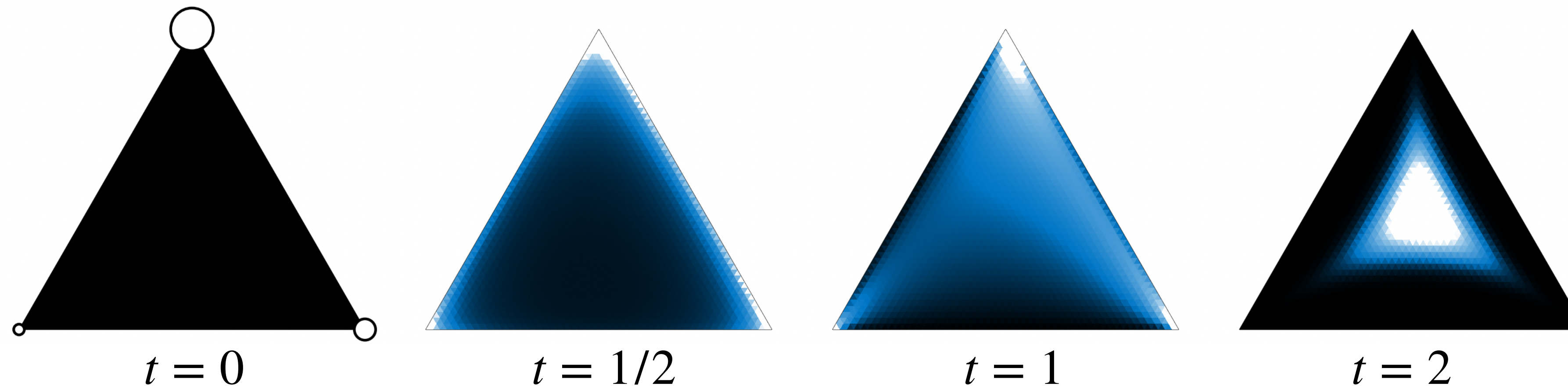


- Concrete distribution (with temperature t) relaxes sampling from a discrete distribution.
- Rather than sampling an index U , sample a vector \mathbf{U} :

$$U_k = \frac{\exp((\alpha_k + G_k) / t)}{\sum_{i=1}^n \exp((\alpha_i + G_i) / t)} .$$

\propto soft max

Gumbel-softmax 'trick'



- Concrete distribution (with temperature t) relaxes sampling from a discrete distribution.
- Rather than sampling an index U , sample a vector \mathbf{U} :

$$U_k = \frac{\exp((\alpha_k + G_k) / t)}{\sum_{i=1}^n \exp((\alpha_i + G_i) / t)} .$$

\propto soft max

- As $t \rightarrow 0^+$, soft max \rightarrow arg max.
 - Concrete distribution \rightarrow discrete distribution.

Figure taken from C. Maddison et al., "The concrete distribution: a continuous relaxation of discrete random variables", *ICLR*, 2017.

Evaluation

Evaluation

- Wyner-Ziv formula has a **closed-form solution in few special cases.**

Evaluation

- Wyner-Ziv formula has a **closed-form solution in few special cases.**
- To evaluate how close we can get to the R-D bound, we choose:

Evaluation

- Wyner-Ziv formula has a **closed-form solution in few special cases**.
- To evaluate how close we can get to the R-D bound, we choose:
 - Let X and Y be correlated, zero-mean and stationary Gaussian memoryless sources.
 - Let $d(\cdot)$ be mean-squared error.

Evaluation

- Wyner-Ziv formula has a **closed-form solution in few special cases**.
- To evaluate how close we can get to the R-D bound, we choose:
 - Let X and Y be correlated, zero-mean and stationary Gaussian memoryless sources.
 - Let $d(\cdot)$ be mean-squared error.
- Wyner-Ziv R-D function then is:

$$R_{WZ}(D) = \frac{1}{2} \log \left(\frac{\sigma_{x|y}^2}{D} \right), \quad 0 \leq D \leq \sigma_{x|y}^2.$$

Evaluation

- Wyner-Ziv formula has a **closed-form solution in few special cases**.
- To evaluate how close we can get to the R-D bound, we choose:
 - Let X and Y be correlated, zero-mean and stationary Gaussian memoryless sources.
 - Let $d(\cdot)$ be mean-squared error.
- Wyner-Ziv R-D function then is:

$$R_{WZ}(D) = \frac{1}{2} \log \left(\frac{\sigma_{x|y}^2}{D} \right), \quad 0 \leq D \leq \sigma_{x|y}^2.$$

- Consider correlation patterns of $X = Y + N$ and $Y = X + N$.

Evaluation

- Wyner-Ziv formula has a **closed-form solution in few special cases**.
- To evaluate how close we can get to the R-D bound, we choose:
 - Let X and Y be correlated, zero-mean and stationary Gaussian memoryless sources.
 - Let $d(\cdot)$ be mean-squared error.
- Wyner-Ziv R-D function then is:

$$R_{WZ}(D) = \frac{1}{2} \log \left(\frac{\sigma_{x|y}^2}{D} \right), \quad 0 \leq D \leq \sigma_{x|y}^2.$$

- Consider correlation patterns of $X = Y + N$ and $Y = X + N$.
- The neural compressor **does not make any assumptions** on the source distribution.

Evaluation

- Wyner-Ziv formula has a **closed-form solution in few special cases**.
- To evaluate how close we can get to the R-D bound, we choose:
 - Let X and Y be correlated, zero-mean and stationary Gaussian memoryless sources.
 - Let $d(\cdot)$ be mean-squared error.
- Wyner-Ziv R-D function then is:

$$R_{WZ}(D) = \frac{1}{2} \log \left(\frac{\sigma_{x|y}^2}{D} \right), \quad 0 \leq D \leq \sigma_{x|y}^2.$$

- Consider correlation patterns of $X = Y + N$ and $Y = X + N$.
- The neural compressor **does not make any assumptions** on the source distribution.
 - The model parameters $\{\theta, \phi, \xi, \zeta\}$ are learned in a data-driven way.

Results

Learned compressor recovers binning.

Results

Learned compressor recovers binning.

Learned encoder:

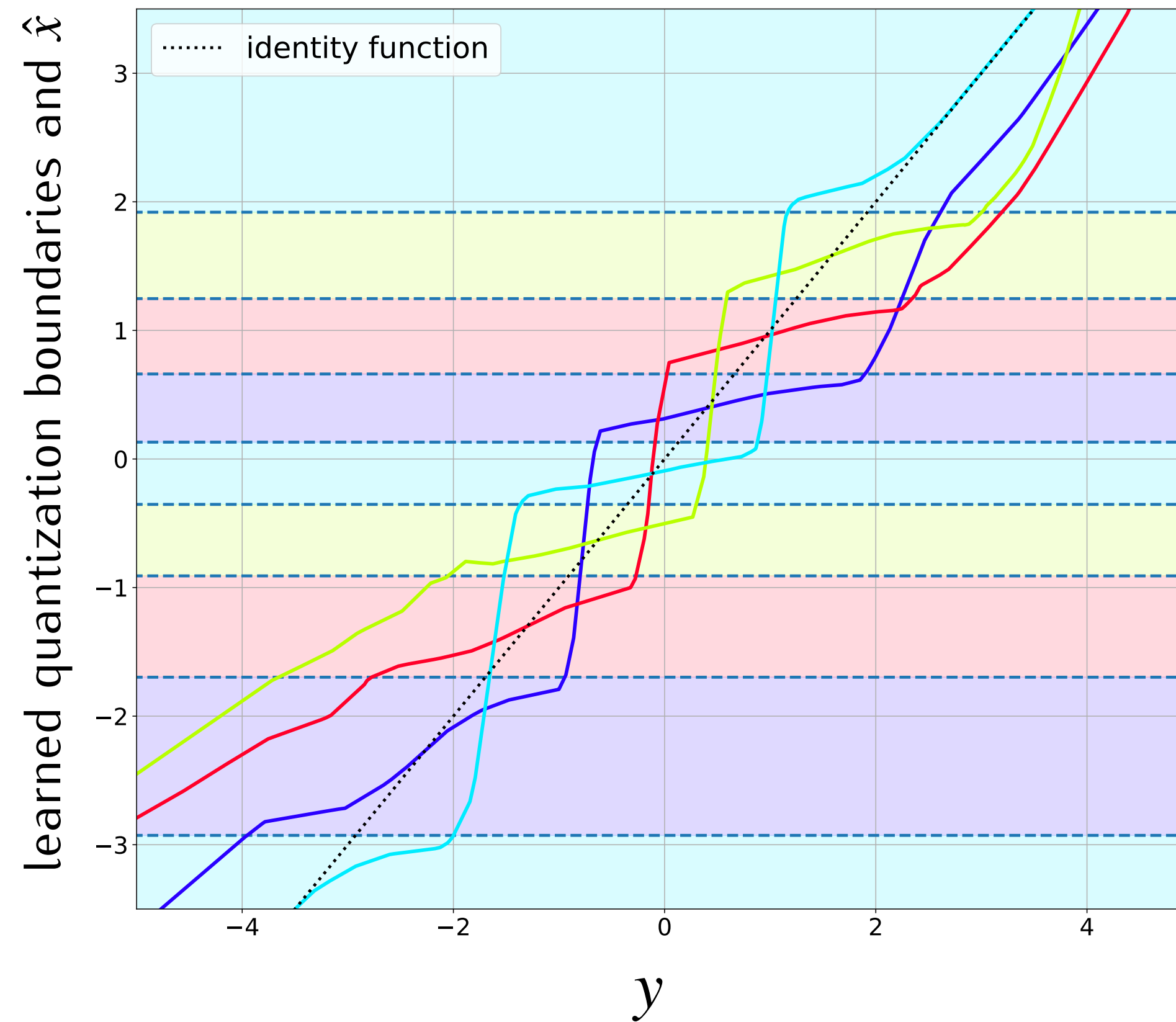
$$u = \arg \max_v p_{\theta}(v | x)$$

Results

Learned compressor recovers binning.

Learned encoder:

$$u = \arg \max_v p_{\theta}(v | x)$$



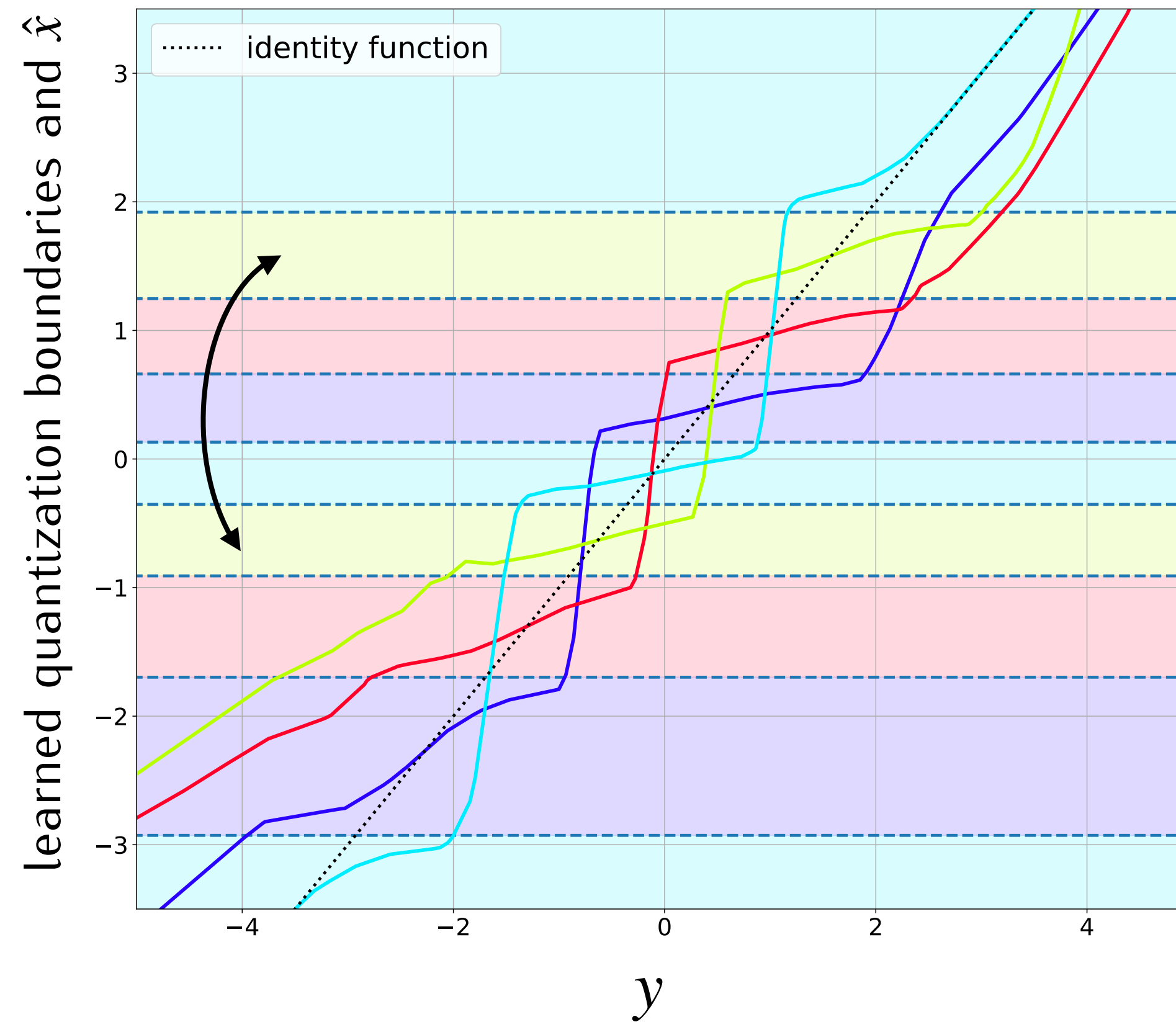
Marginal formulation.

$$X = Y + N \text{ with } Y \sim N(0,1) \text{ and } N \sim N(0,10^{-1}).$$

Results

Learned compressor recovers binning.

Learned encoder:
 $u = \arg \max_v p_\theta(v | x)$



Marginal formulation.

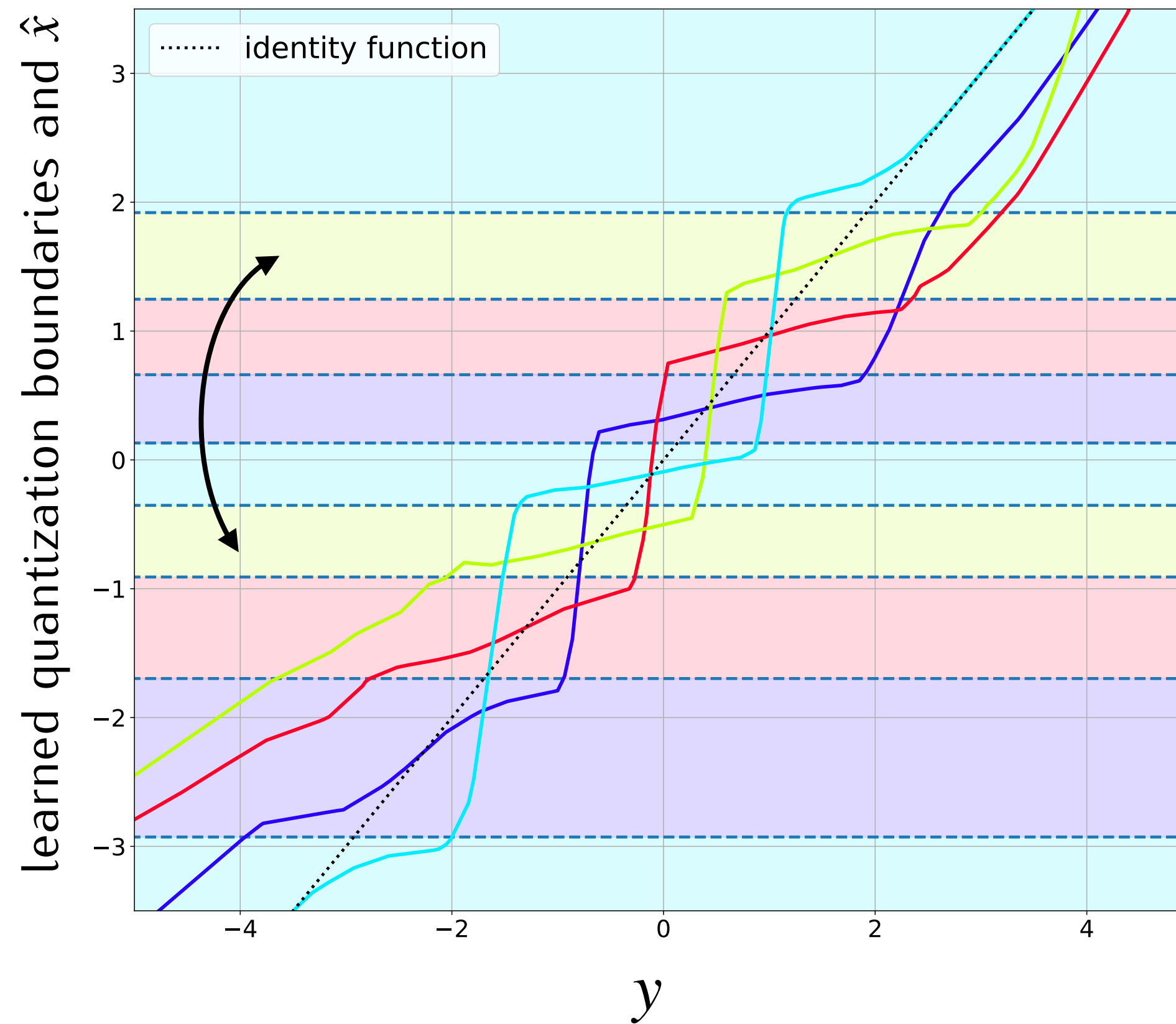
$X = Y + N$ with $Y \sim N(0,1)$ and $N \sim N(0,10^{-1})$.

Results

Learned compressor recovers binning.

Learned encoder:
 $u = \arg \max_v p_\theta(v | x)$

same index



Marginal formulation.

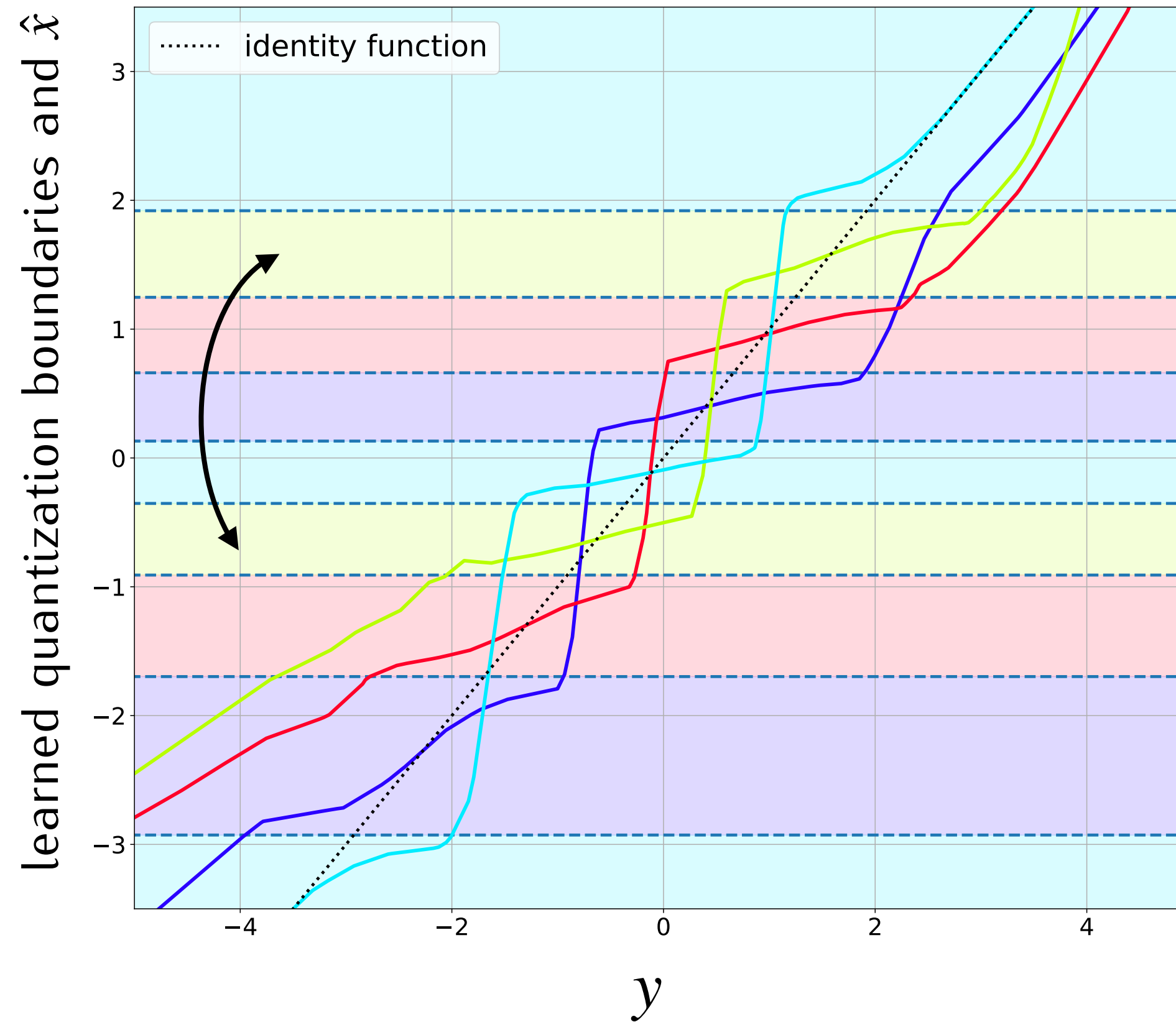
$X = Y + N$ with $Y \sim N(0,1)$ and $N \sim N(0,10^{-1})$.

Results

Learned compressor recovers binning.

Learned encoder:
 $u = \arg \max_v p_\theta(v | x)$

same index
 \implies binning.



Marginal formulation.

$X = Y + N$ with $Y \sim N(0,1)$ and $N \sim N(0,10^{-1})$.

Results

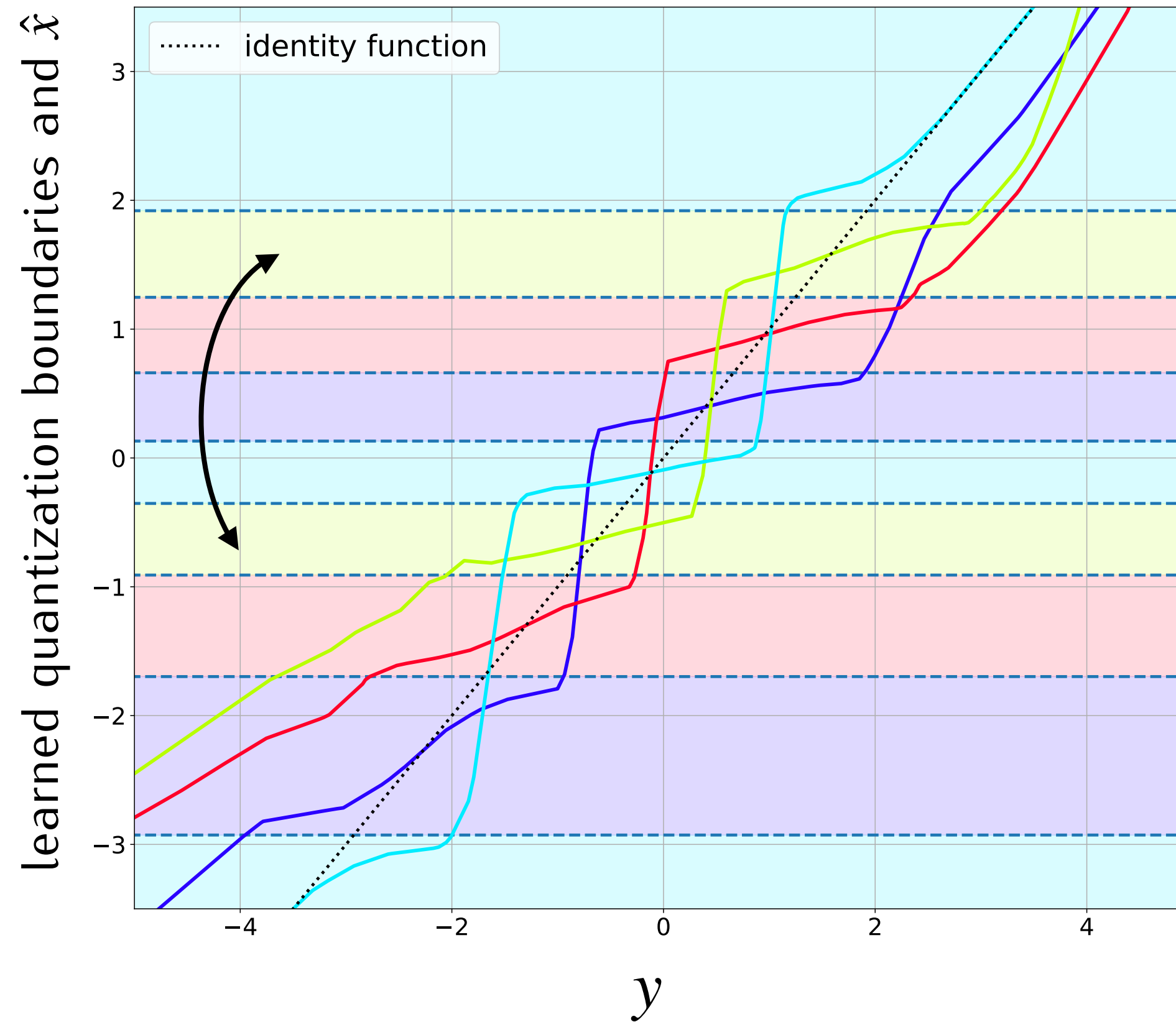
Learned compressor recovers binning.

Learned decoder:

$$\hat{x} = g_{\phi}(u, y)$$

Learned encoder:

$$u = \arg \max_v p_{\theta}(v | x)$$



same index

\implies binning.

Marginal formulation.

$$X = Y + N \text{ with } Y \sim N(0,1) \text{ and } N \sim N(0,10^{-1}).$$

Results

Learned compressor recovers binning.

Learned decoder:

$$\hat{x} = g_{\phi}(u, y)$$

In quadratic-Gaussian WZ setup, the optimal decoder does:

$$\hat{x} = (1 - \beta) \cdot y + \beta \cdot u,$$

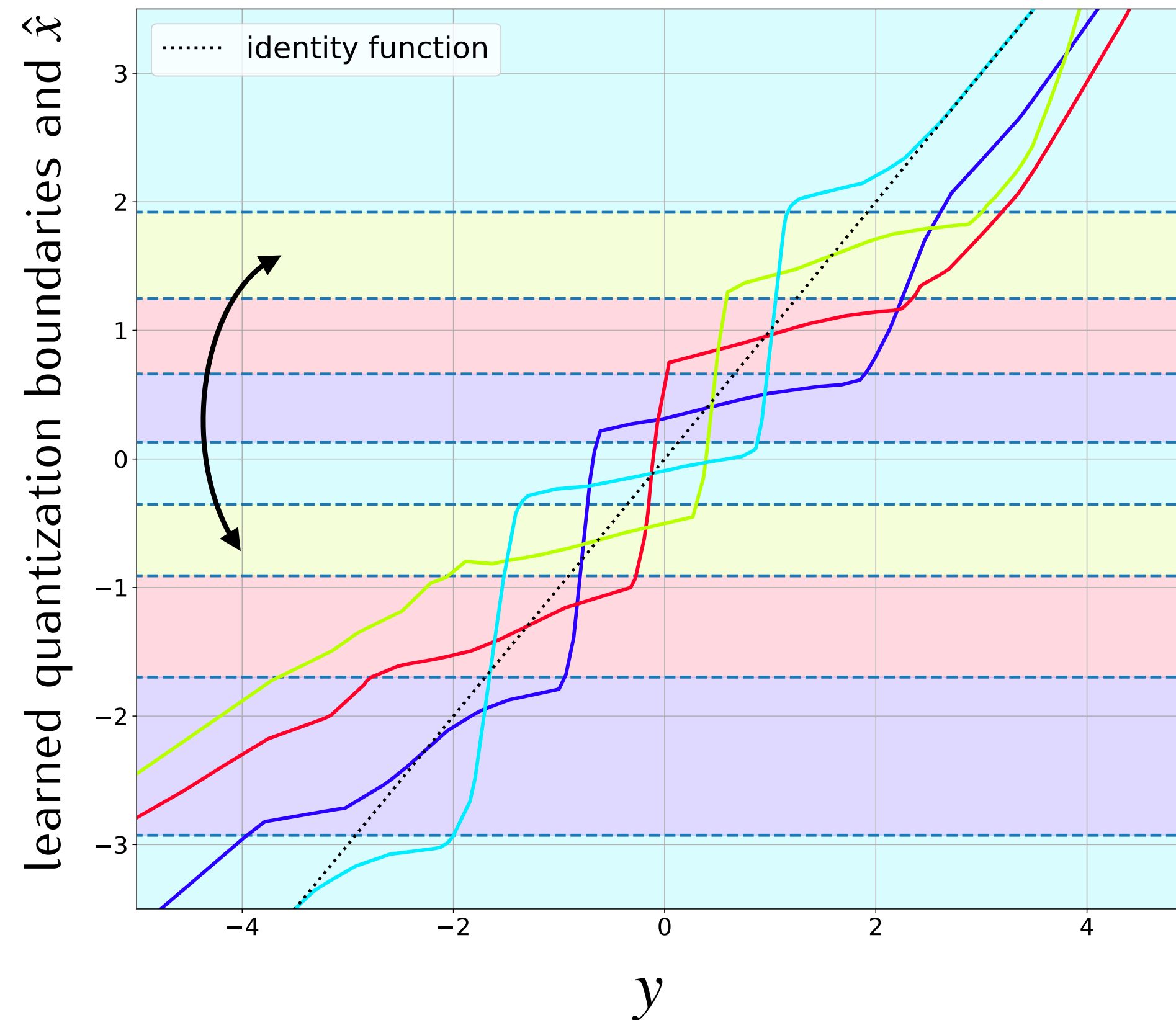
where $\beta \propto \sigma_n^2$.

Learned encoder:

$$u = \arg \max_v p_{\theta}(v | x)$$

same index

\implies binning.



Marginal formulation.

$$X = Y + N \text{ with } Y \sim N(0,1) \text{ and } N \sim N(0,10^{-1}).$$

Results

Learned compressor recovers binning.

Learned decoder:

$$\hat{x} = g_{\phi}(u, y)$$

In quadratic-Gaussian WZ setup, the optimal decoder does:

$$\hat{x} = (1 - \beta) \cdot y + \beta \cdot u,$$

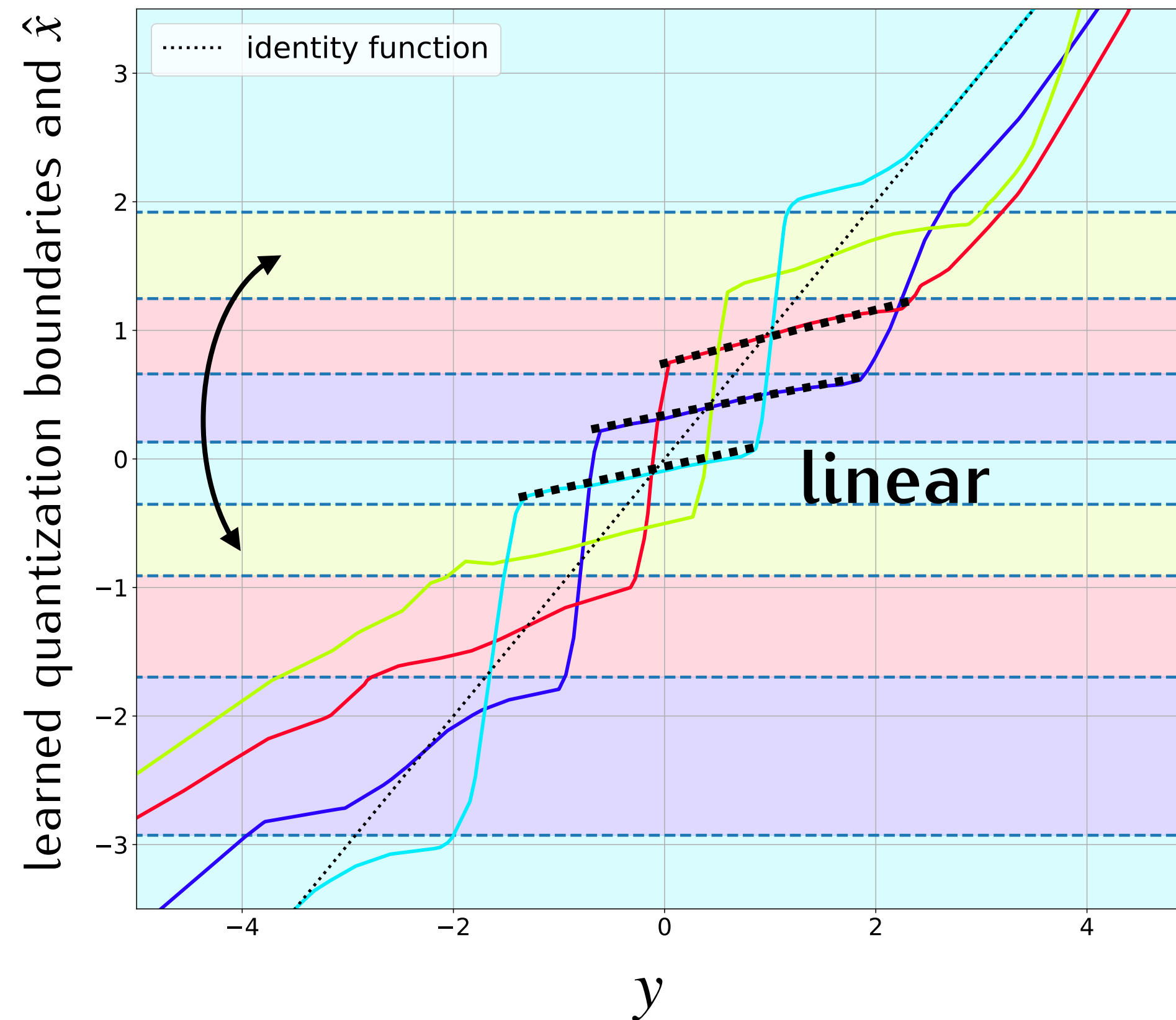
where $\beta \propto \sigma_n^2$.

Recovers optimal reconstruction function.

Learned encoder:

$$u = \arg \max_v p_{\theta}(v | x)$$

same index
 \implies binning.



Marginal formulation.

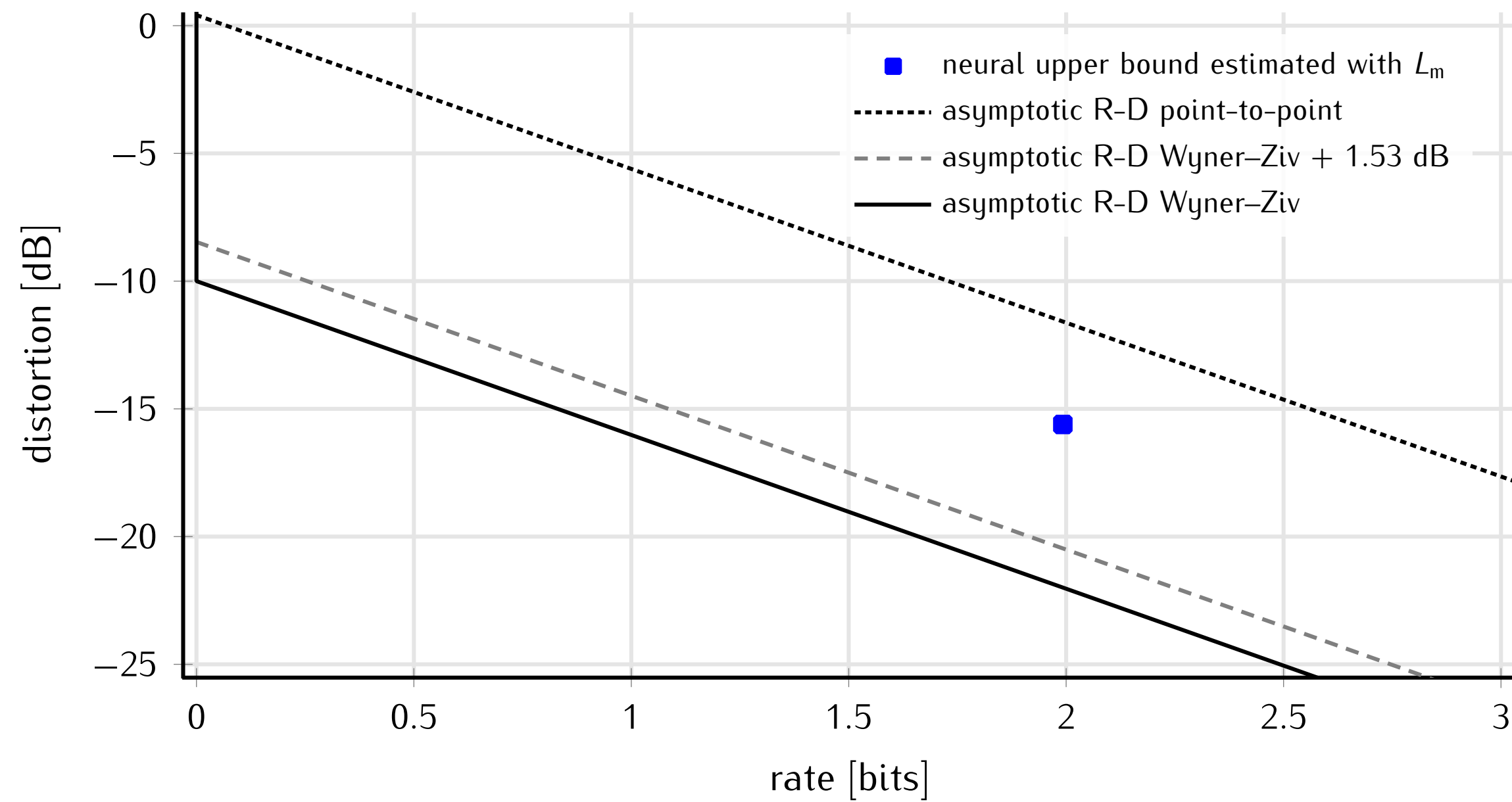
$$X = Y + N \text{ with } Y \sim N(0,1) \text{ and } N \sim N(0,10^{-1}).$$

Results

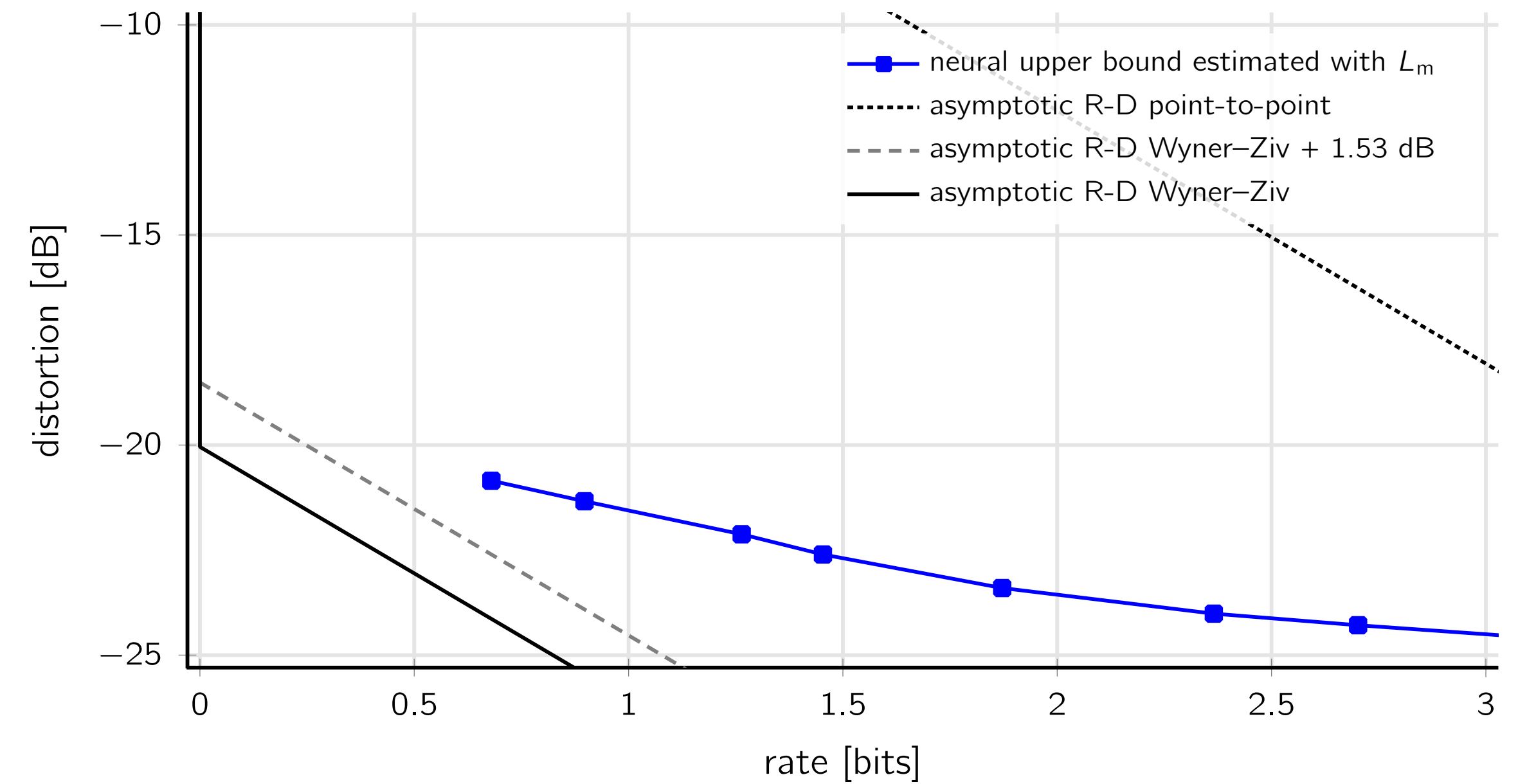
R-D performances.

Results

R-D performances.



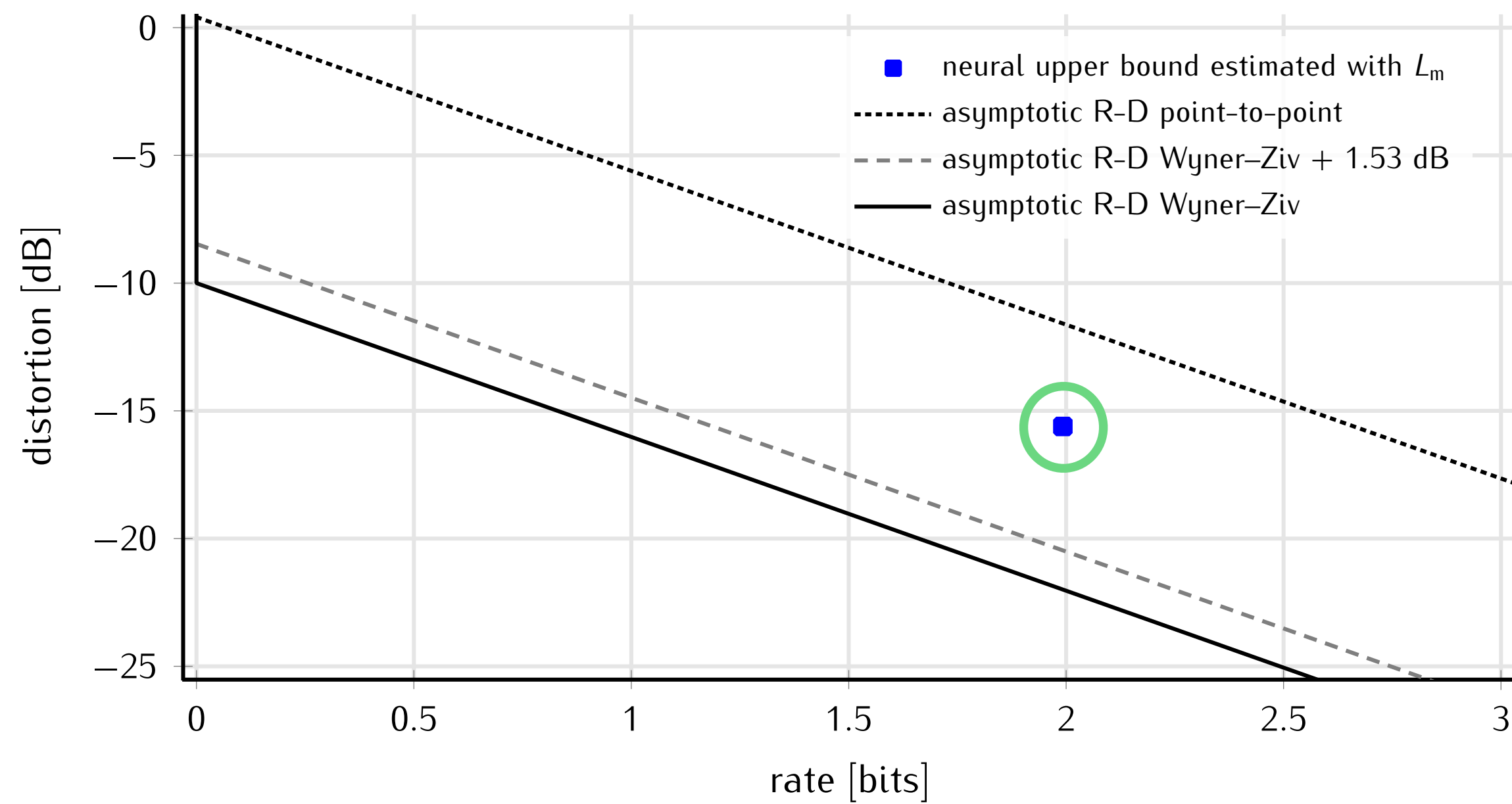
$X = Y + N$ with $Y \sim N(0,1)$ and $N \sim N(0,10^{-1})$.



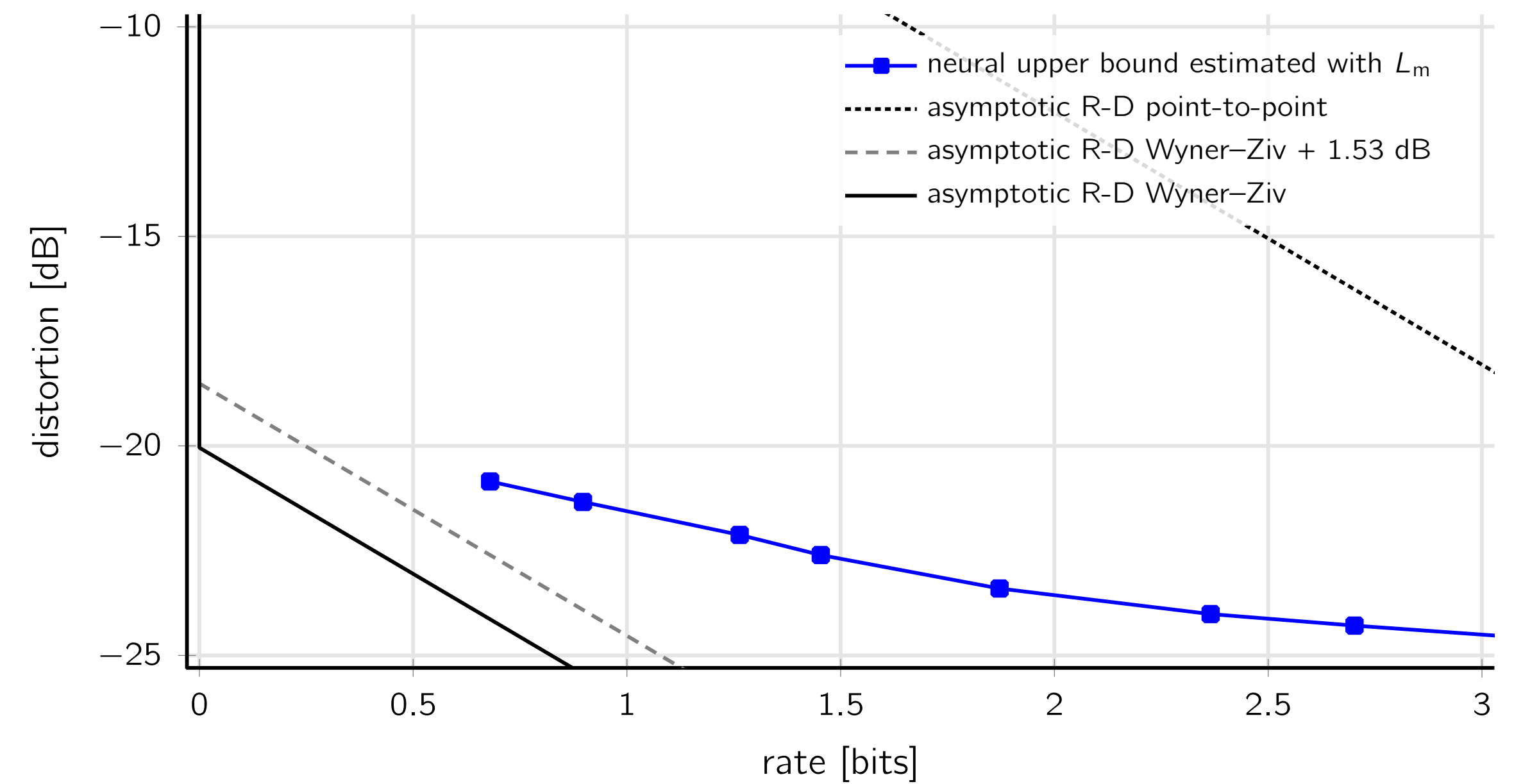
$Y = X + N$ with $X \sim N(0,1)$ and $N \sim N(0,10^{-2})$.

Results

R-D performances.



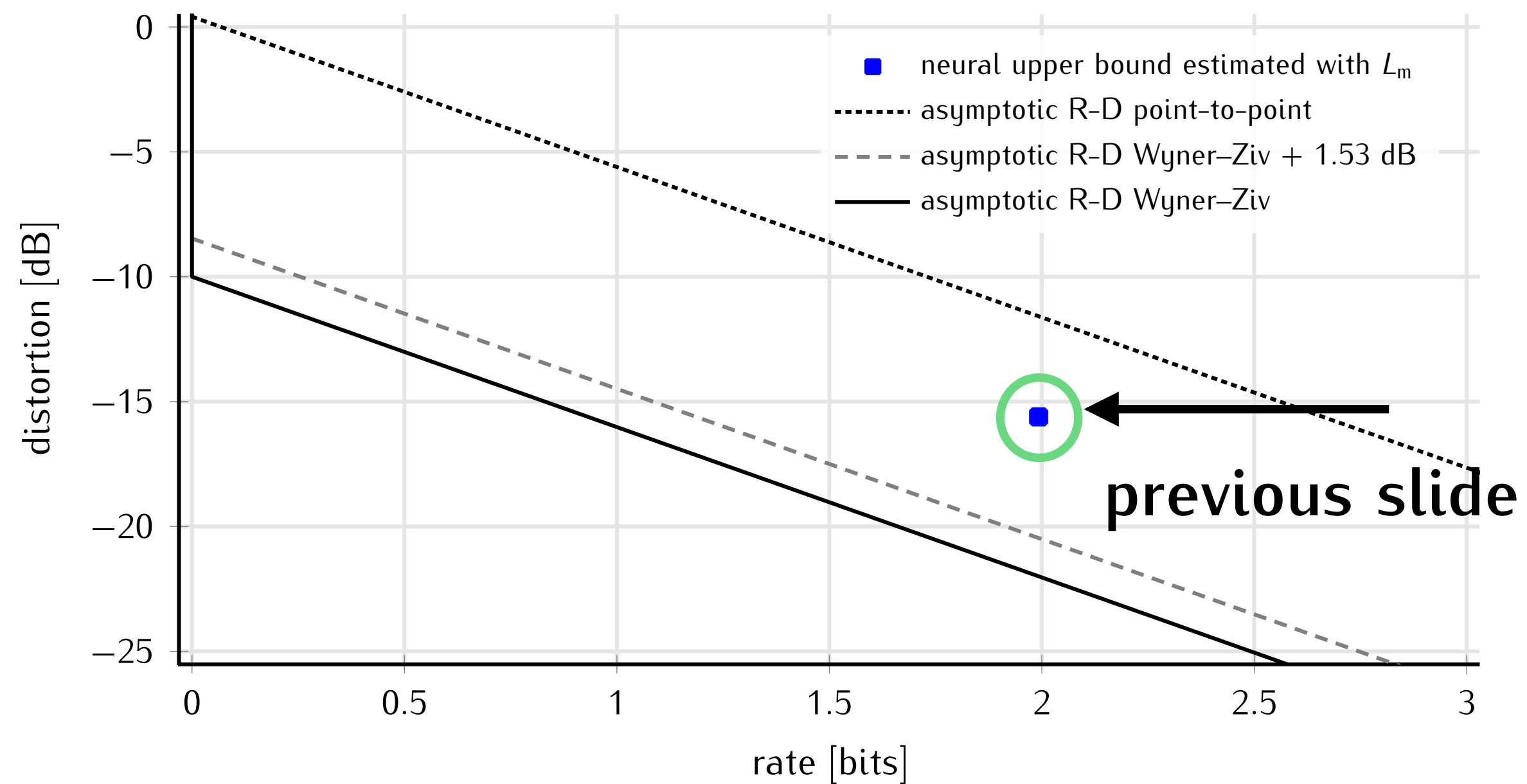
$X = Y + N$ with $Y \sim N(0,1)$ and $N \sim N(0,10^{-1})$.



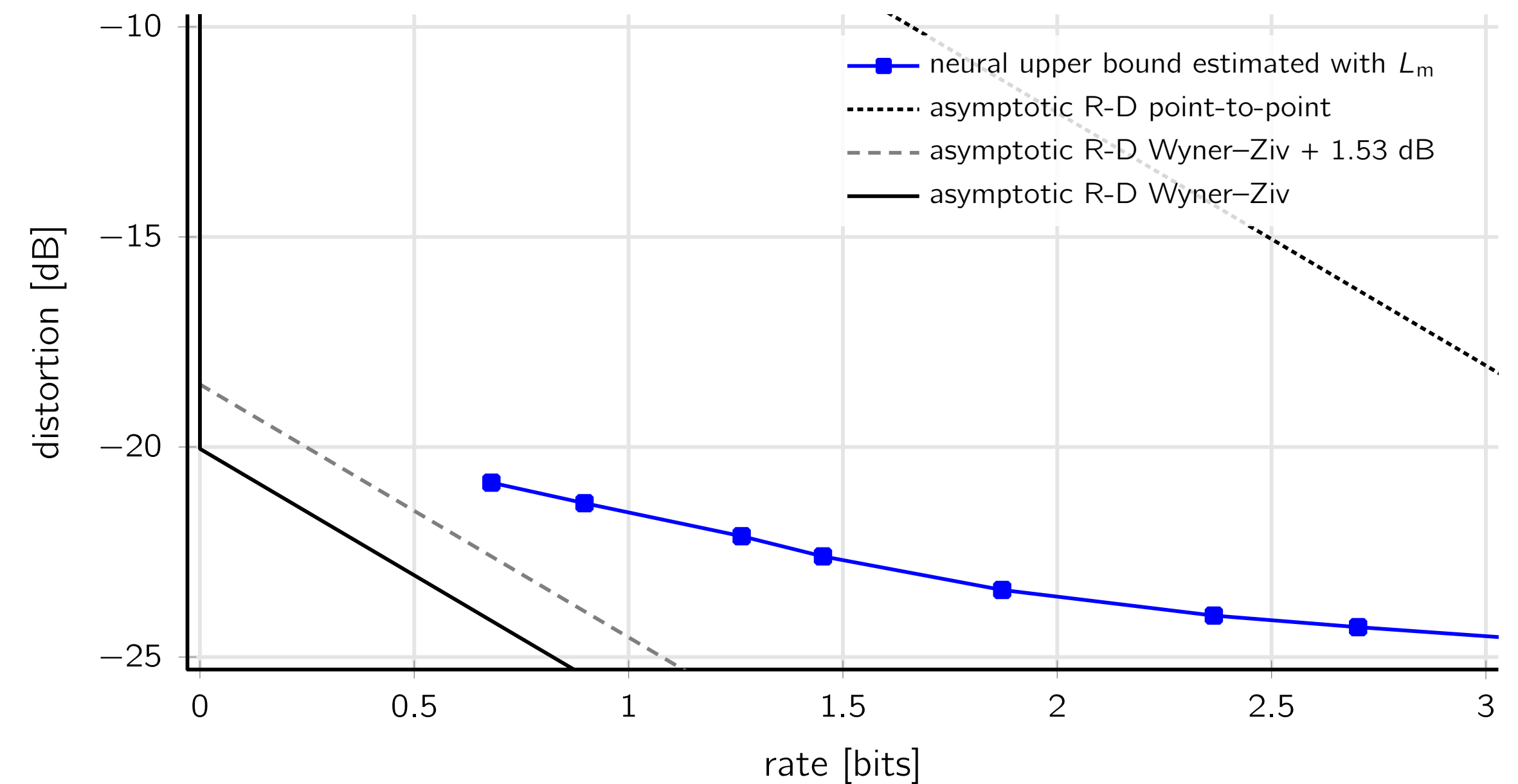
$Y = X + N$ with $X \sim N(0,1)$ and $N \sim N(0,10^{-2})$.

Results

R-D performances.



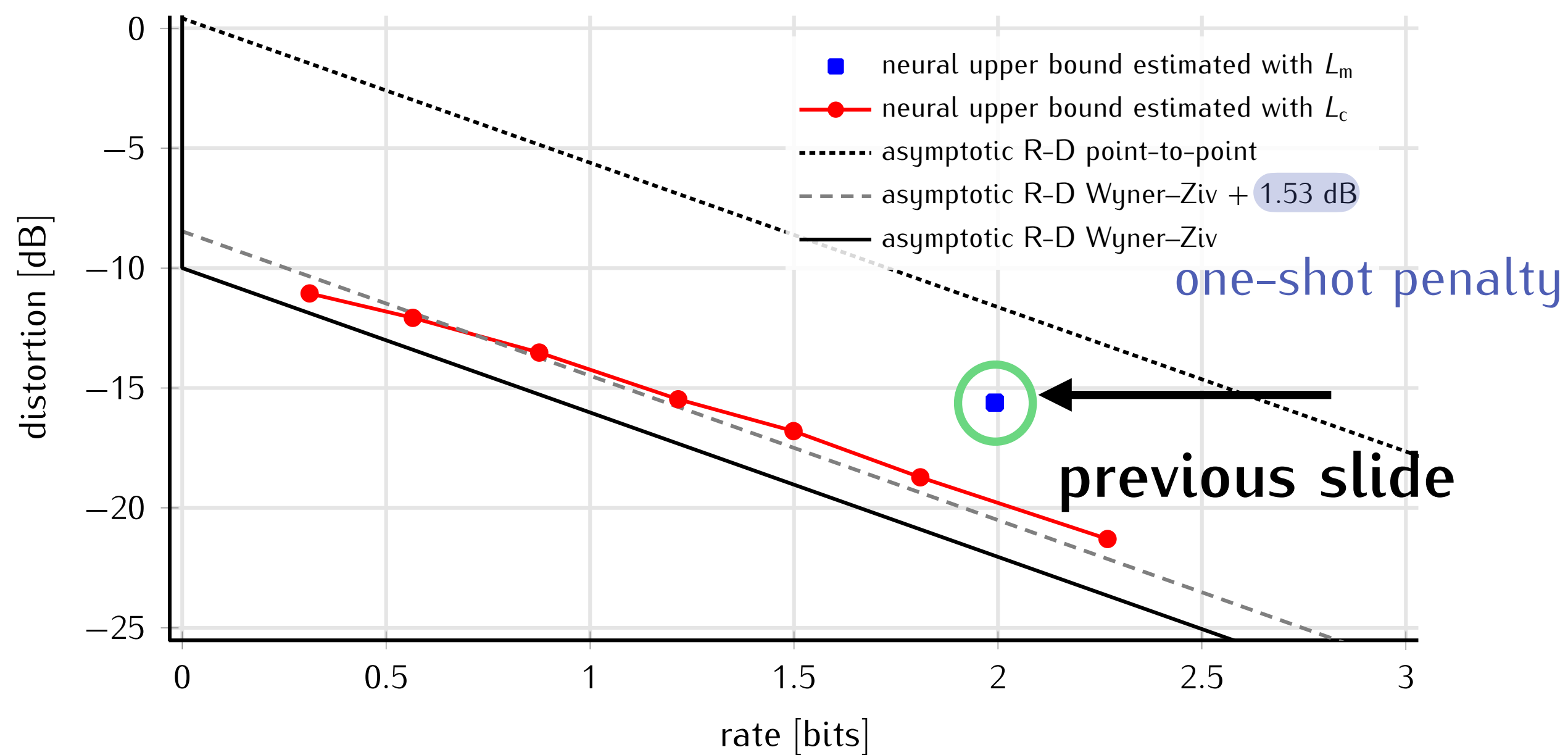
$$X = Y + N \text{ with } Y \sim N(0,1) \text{ and } N \sim N(0,10^{-1}) .$$



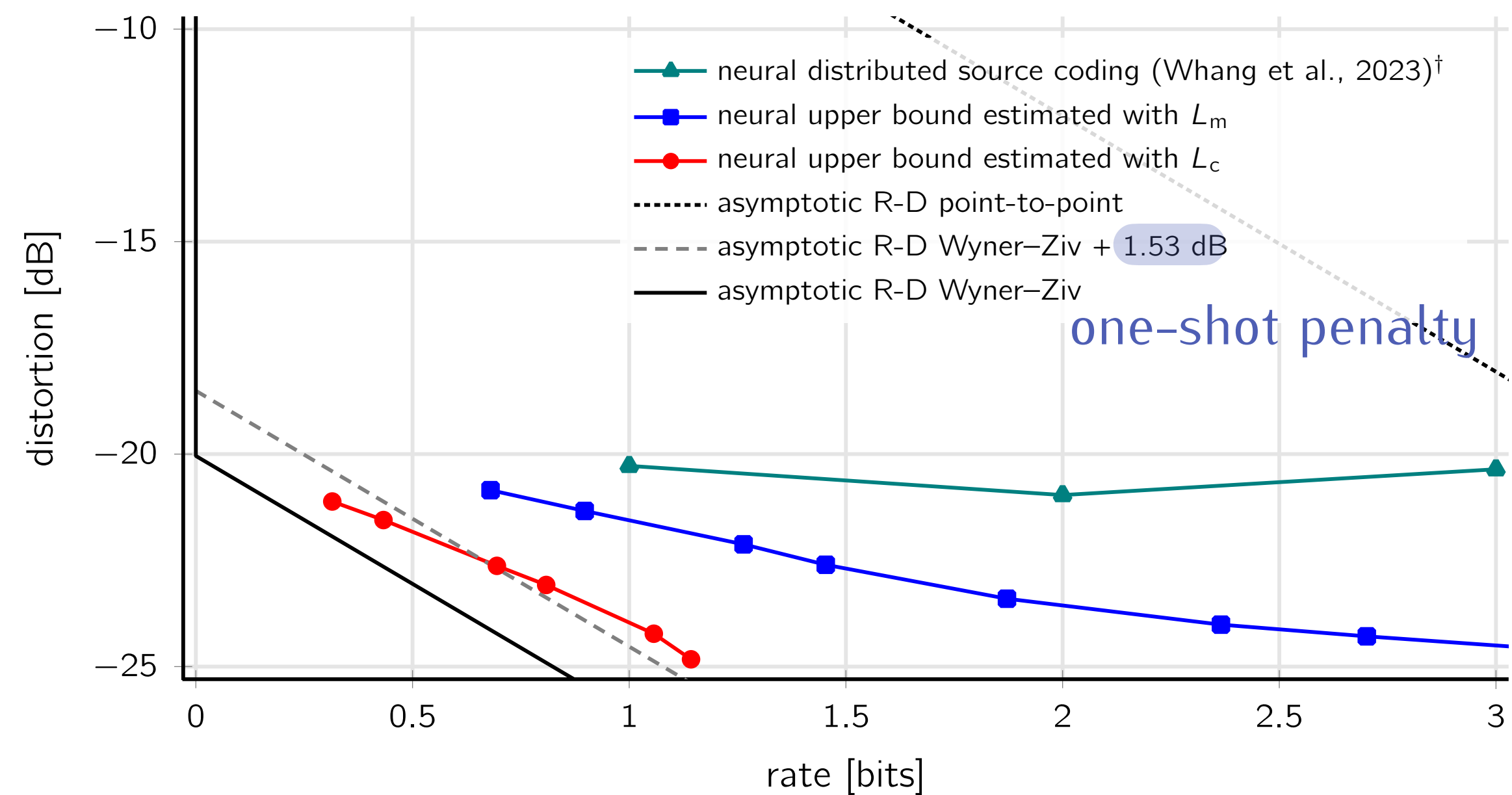
$$Y = X + N \text{ with } X \sim N(0,1) \text{ and } N \sim N(0,10^{-2}) .$$

Results

R-D performances.



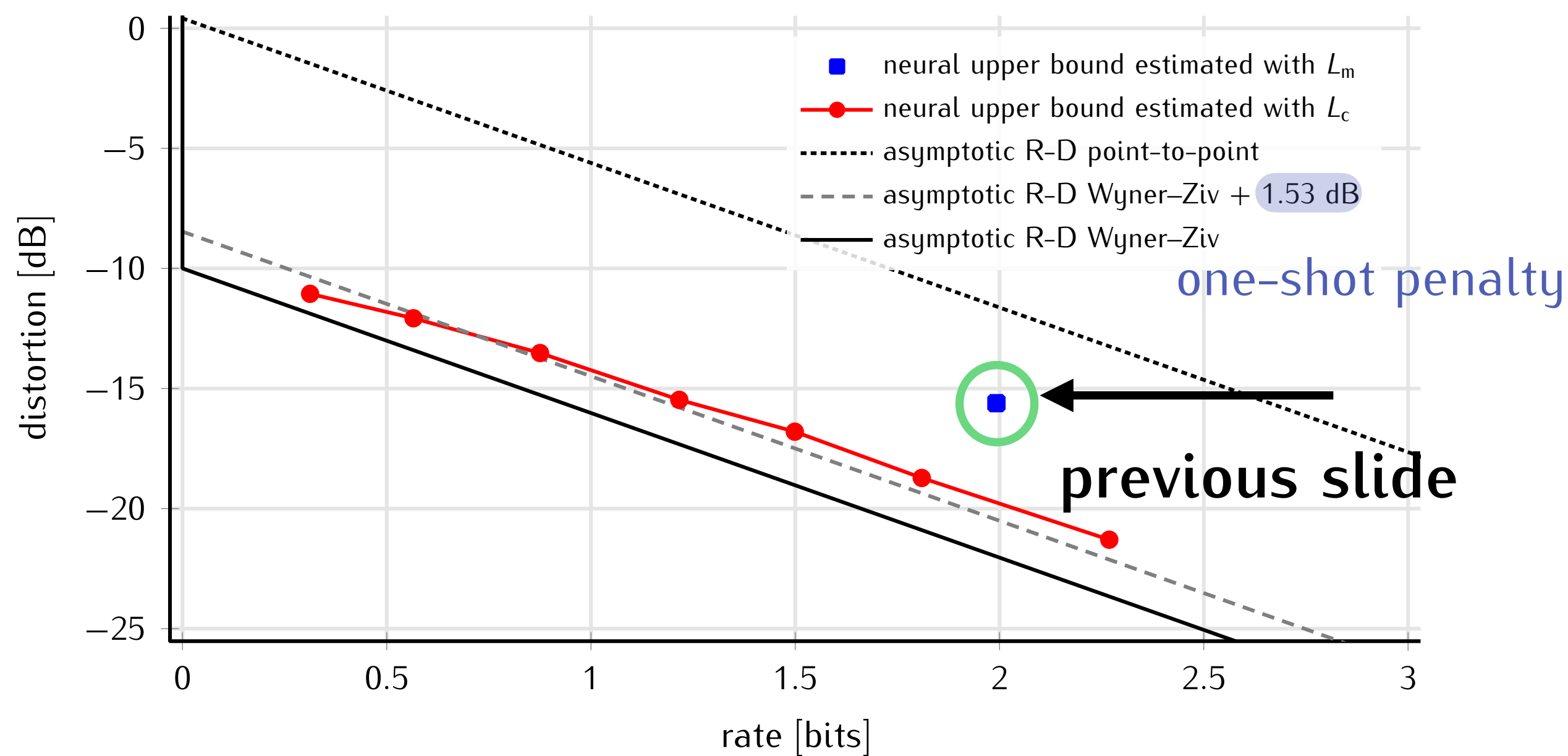
$$X = Y + N \text{ with } Y \sim N(0,1) \text{ and } N \sim N(0,10^{-1}) .$$



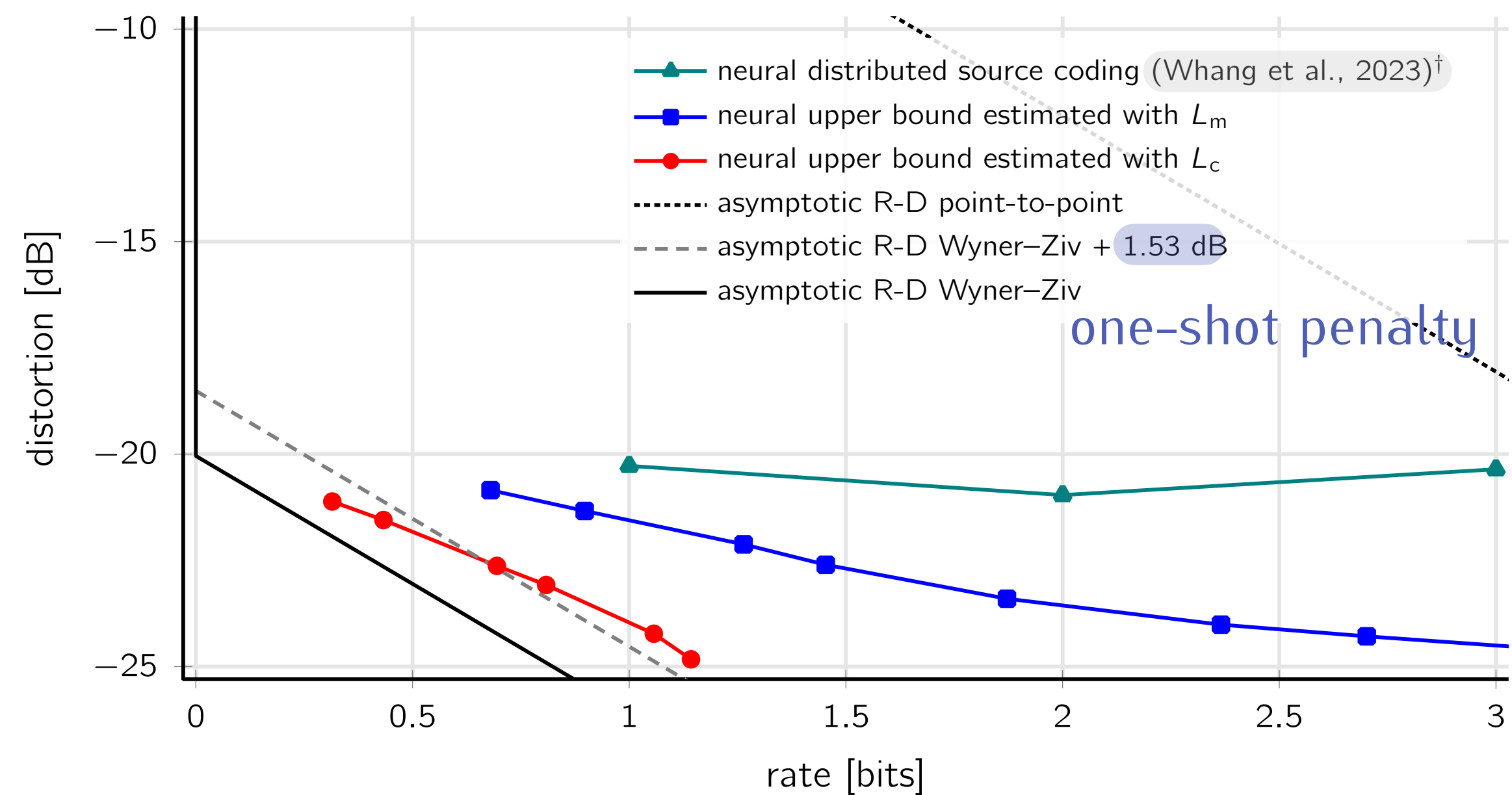
$$Y = X + N \text{ with } X \sim N(0,1) \text{ and } N \sim N(0,10^{-2}) .$$

Results

R-D performances.



$$X = Y + N \text{ with } Y \sim N(0,1) \text{ and } N \sim N(0,10^{-1}) .$$



$$Y = X + N \text{ with } X \sim N(0,1) \text{ and } N \sim N(0,10^{-2}) .$$

[†]J. Whang, A. Nagle, A. Acharya, H. Kim, and A. G. Dimakis, "Neural distributed source coding", <https://arxiv.org/abs/2106.02797>, 2023.

Take-home messages

Take-home messages

- To close the gap between theory and practice in distributed source coding, **learned compression is a promising approach.**

Take-home messages

- To close the gap between theory and practice in distributed source coding, **learned compression is a promising approach.**
- In quadratic-Gaussian case, learned compressors recover some elements of the optimal theoretical solution.

Take-home messages

- To close the gap between theory and practice in distributed source coding, **learned compression is a promising approach.**
- In quadratic-Gaussian case, learned compressors recover some elements of the optimal theoretical solution.
 - Binning in the source space and linear decoding functions.

Take-home messages

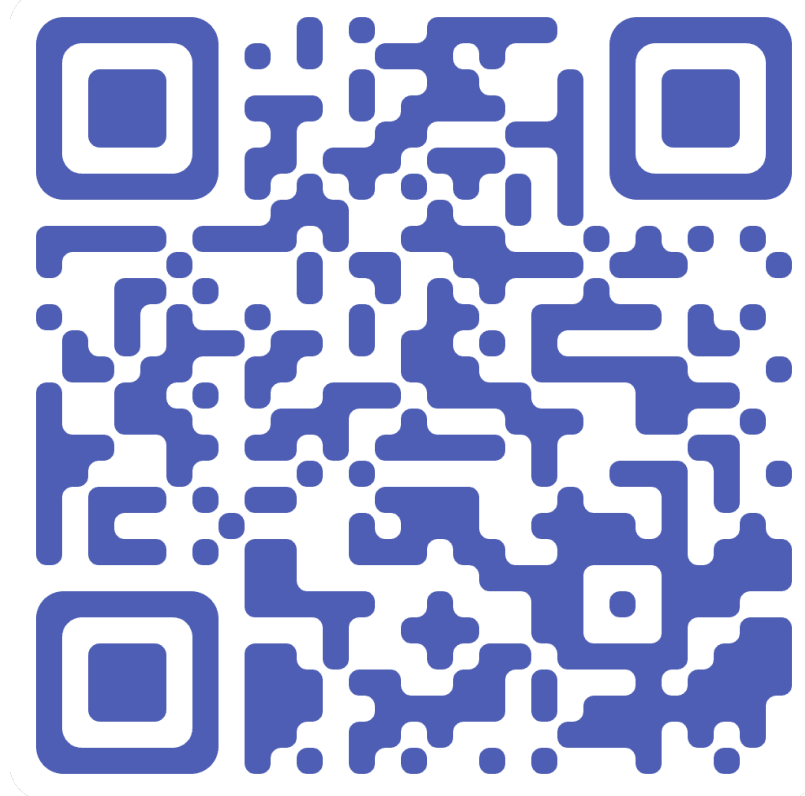
- To close the gap between theory and practice in distributed source coding, **learned compression is a promising approach.**
- In quadratic-Gaussian case, learned compressors recover some elements of the optimal theoretical solution.
 - Binning in the source space and linear decoding functions.
 - First-time **binning** emerges from learning.

Thank you. Questions?

Neural Distributed Compressor Does Binning

Ezgi Özyilkan*, Johannes Ballé†, Elza Erkip*

*NYU, †Google Research
ezgi.ozyilkan@nyu.edu



Neural Compression Workshop @ ICML 2023
Honolulu, HI | July 29, 2023

Thank you. Questions?

Neural Distributed Compressor Does Binning

Ezgi Özyilkan*, Johannes Ballé†, Elza Erkip*

*NYU, †Google Research
ezgi.ozyilkan@nyu.edu



Google Research



Neural Compression Workshop @ ICML 2023
Honolulu, HI | July 29, 2023

Presented at International Symposium of Information Theory (ISIT) 2023.

References

- S. Verdú, “Fifty years of Shannon theory”, *IEEE Transactions on Information Theory*, vol. 2, no. 5, p. 359–366, 1998.
- J. Ballé, V. Laparra, and E. P. Simoncelli, “End-to-end optimized image compression”, *International Conference on Learning Representations*, 2017.
- A. Wyner and J. Ziv, “The rate–distortion function for source coding with side information at the decoder”, *IEEE Transactions on Information Theory*, vol. 22, no. 1, pp. 1–10, 1976.
- D. Slepian and J. Wolf, “Noiseless coding of correlated information sources”, *IEEE Transactions on Information Theory*, vol. 19, no. 4, pp. 471– 480, 1973.
- E. J. Gumbel, “Statistical theory of extreme values and some practical applications: a series of lectures”, *US Department of Commerce*, 1954.
- C. J. Maddison, A. Mnih, and Y. W. Teh, “The concrete distribution: a continuous relaxation of discrete random variables”, *International Conference on Learning Representations*, 2017.
- J. Whang, A. Nagle, A. Acharya, H. Kim, and A. G. Dimakis, “Neural distributed source coding”, <https://arxiv.org/abs/2106.02797>, 2023.