

# Generalization and Optimization in Symmetry-Preserving ML: Sample Complexity and Implicit Bias

Wei Zhu

University of Massachusetts Amherst

Boston Symmetry Day  
MIT

November 3, 2023

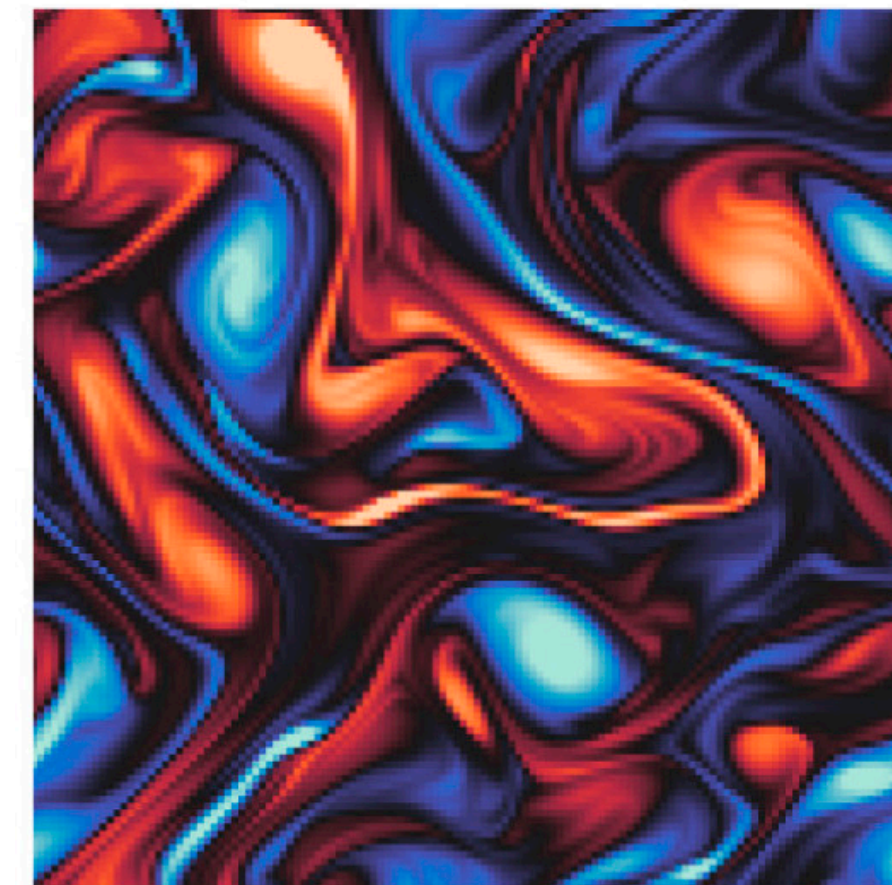
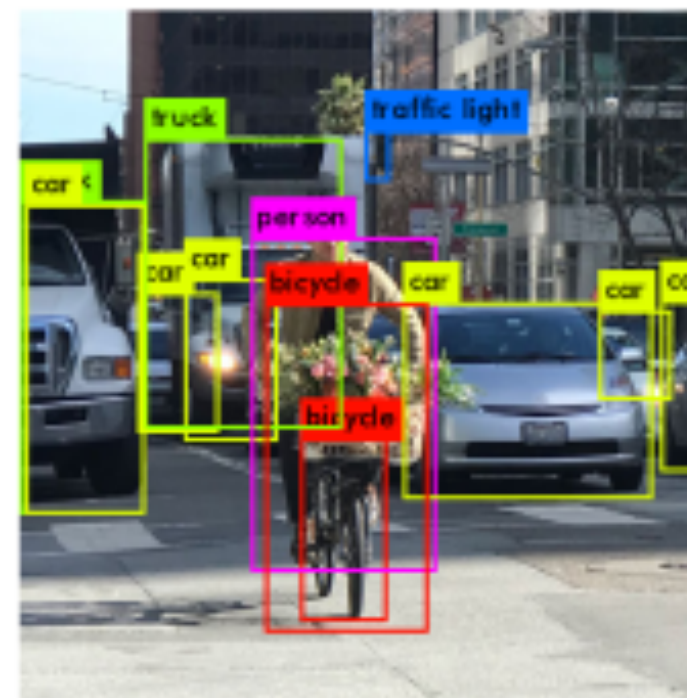
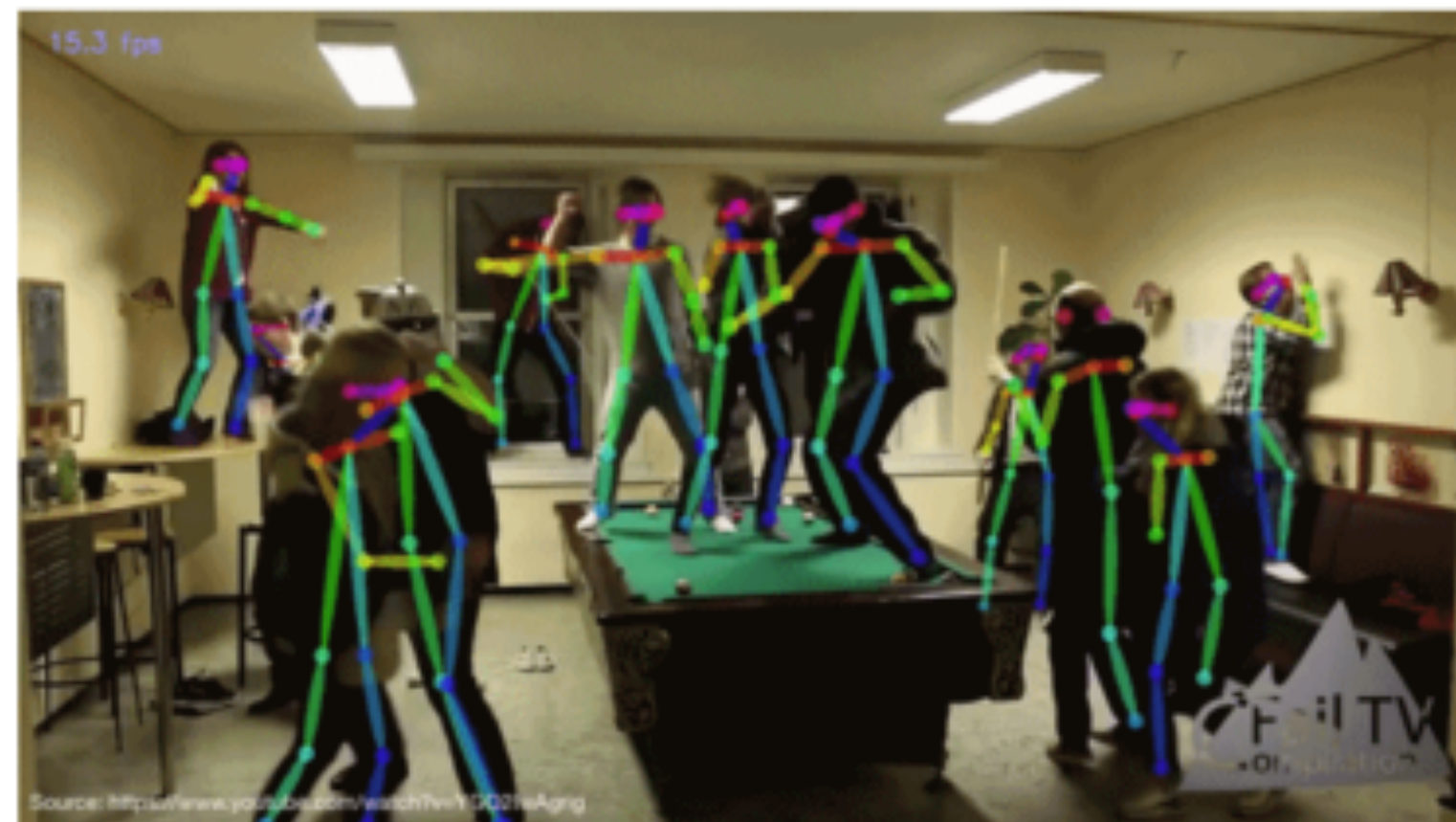
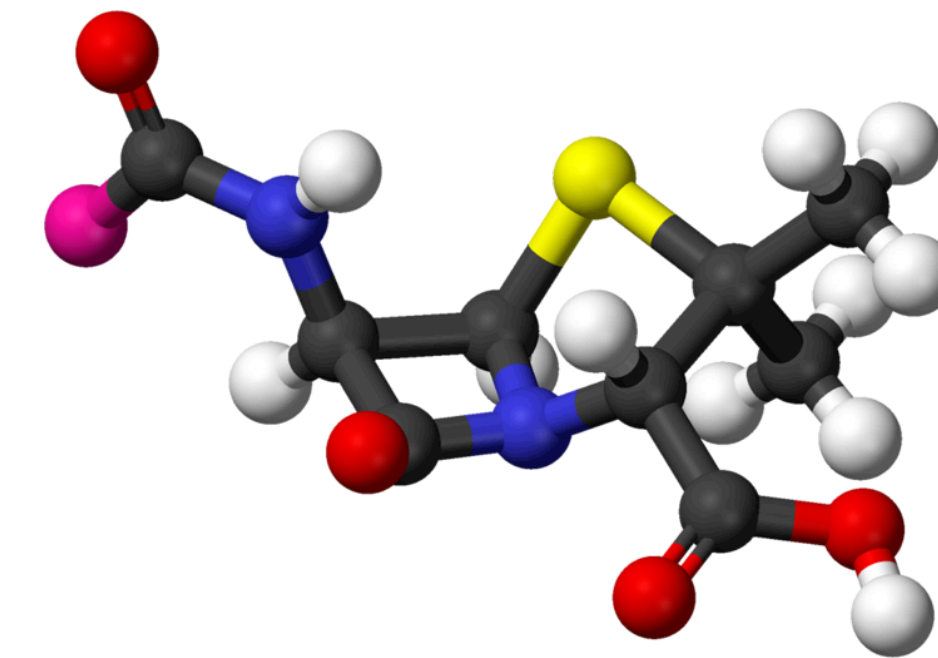


**Joint work** with many people, but mostly **Ziyu Chen** (UMass Amherst)





# Symmetry is everywhere



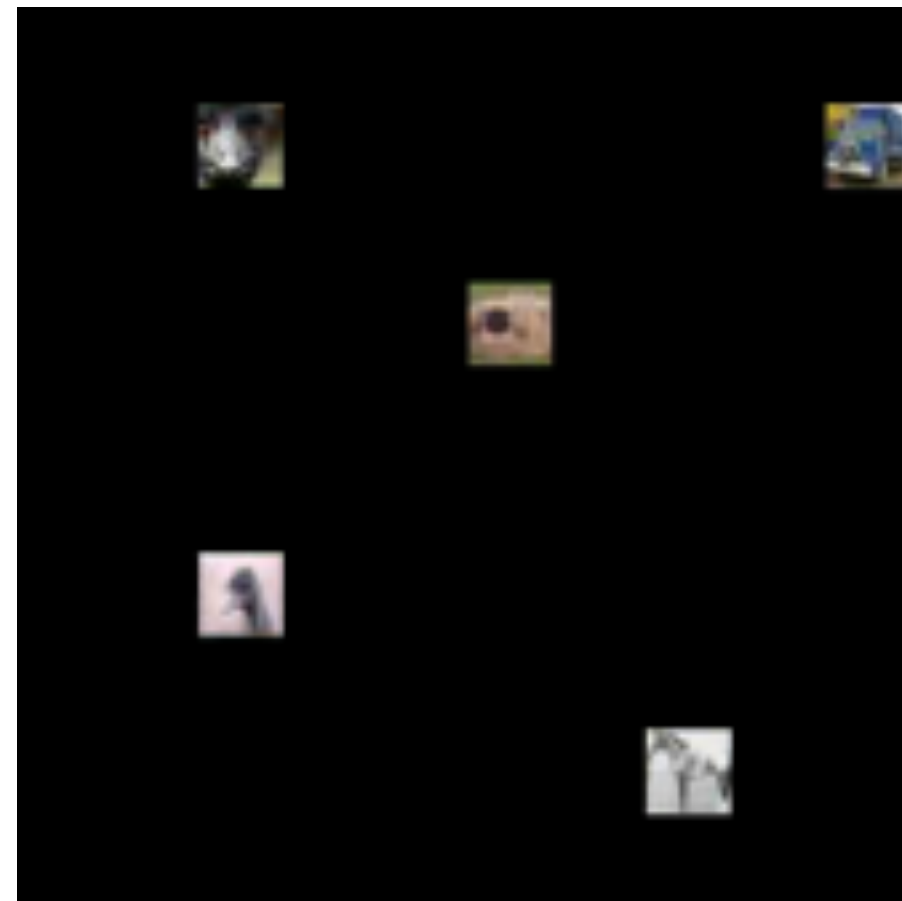
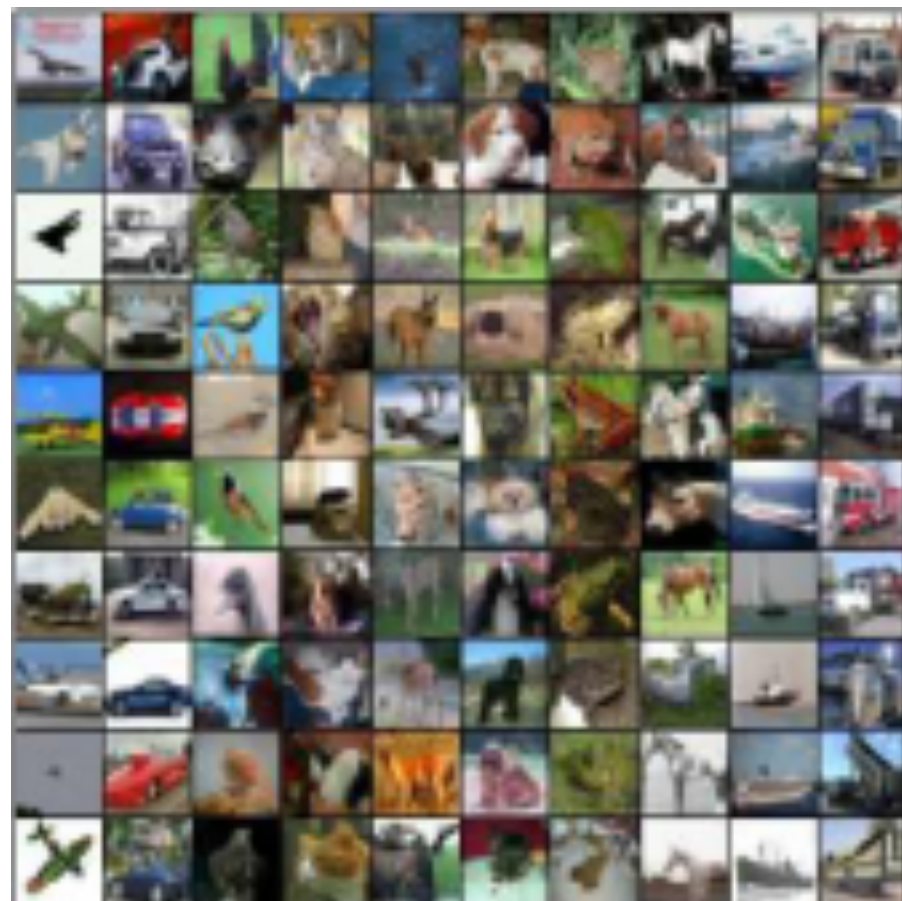


# Missing pieces



# Missing pieces

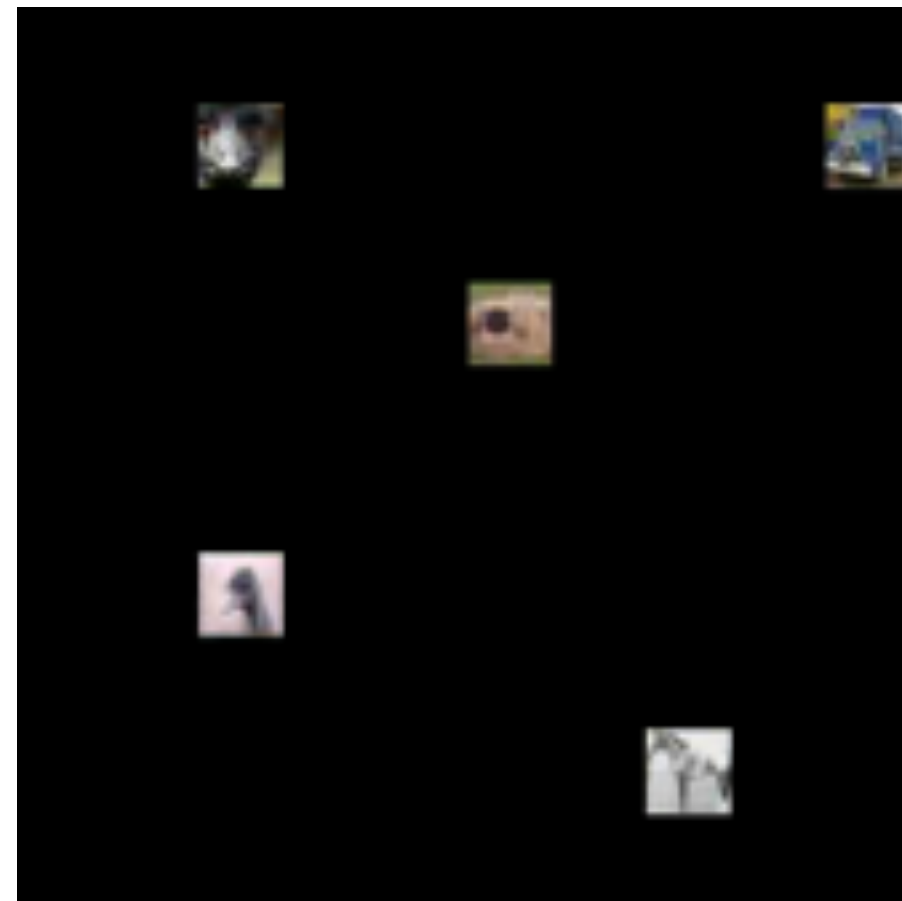
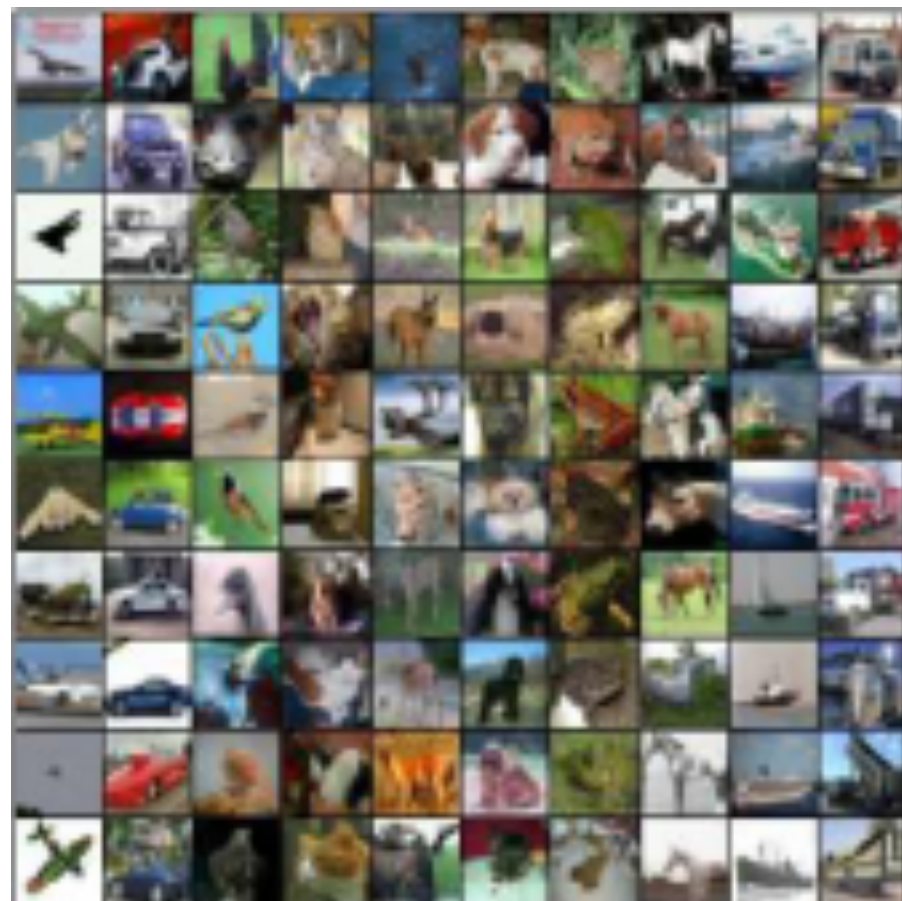
- **Exact quantification of the improvement**
  - **Sample complexity** and **error bound**.





# Missing pieces

- **Exact quantification of the improvement**
  - **Sample complexity** and **error bound**.
- **Does it converge? To what solution?**
  - **Training dynamics** of equivariant models

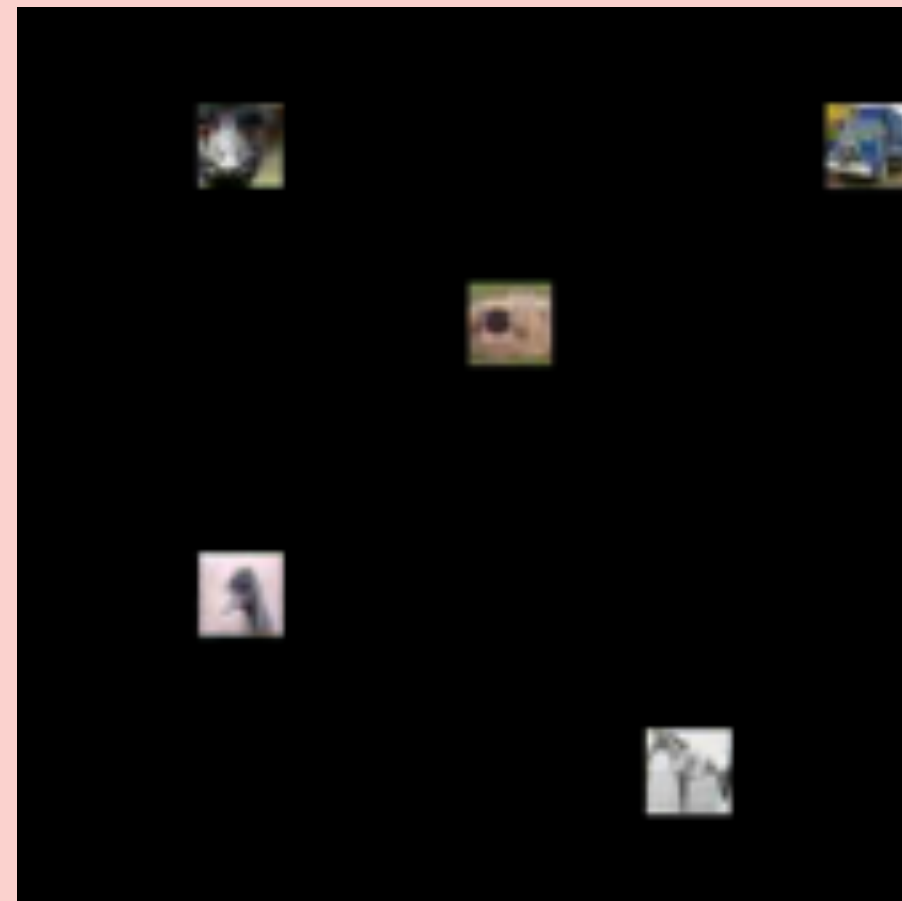


$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; S) = \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{W}), y_i)$$



# Missing pieces

- **Exact quantification of the improvement**
  - **Sample complexity** and **error bound**.



- **Does it converge? To what solution?**
  - **Training dynamics** of equivariant models

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; S) = \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{W}), y_i)$$



# Symmetry-preserving GANs and their improved sample complexity

- J. Birrell, M.A. Katsoulakis, L. Rey-Bellet, **W. Zhu**. “Structure-preserving GANs”. *ICML* (2022)
- Z. Chen, M.A. Katsoulakis, L. Rey-Bellet, **W. Zhu**. “Sample complexity of probability divergences under group symmetry”. *ICML* (2023)



# Generative adversarial networks (GANs)



StyleGAN2, Karras et al., CVPR 2020



StyleGAN3, Karras et al., NeurIPS 2021

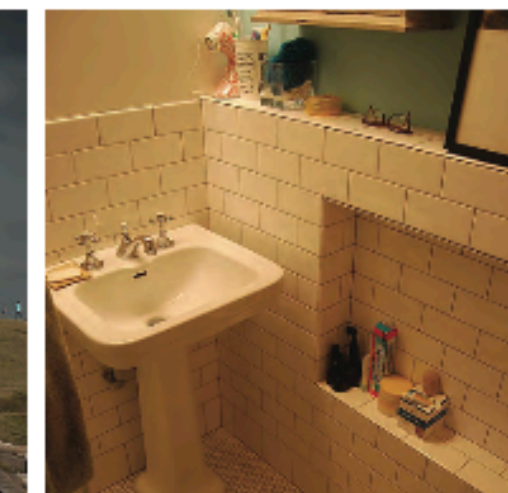
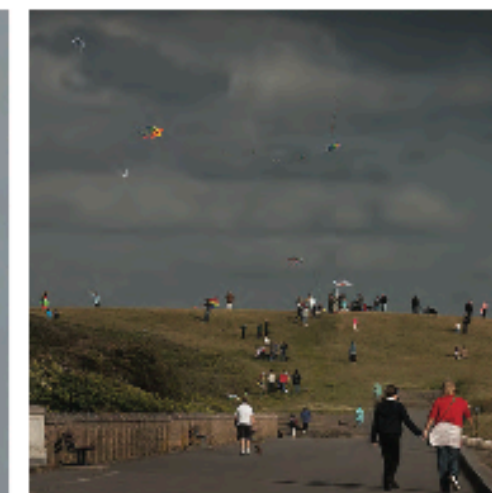
This small bird has a yellow crown and a white belly.

This bird has a blue crown with white throat and brown secondaries.

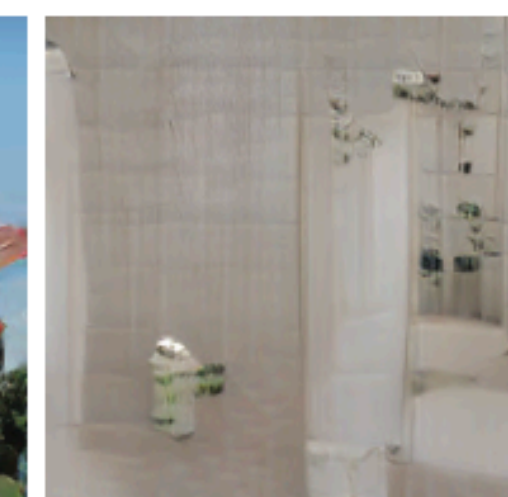
People at the park flying kites and walking.

The bathroom with the white tile has been cleaned.

Real images



Synthesized images



DM-GAN, Zhu et al., CVPR 2019



# Generative adversarial networks (GANs)

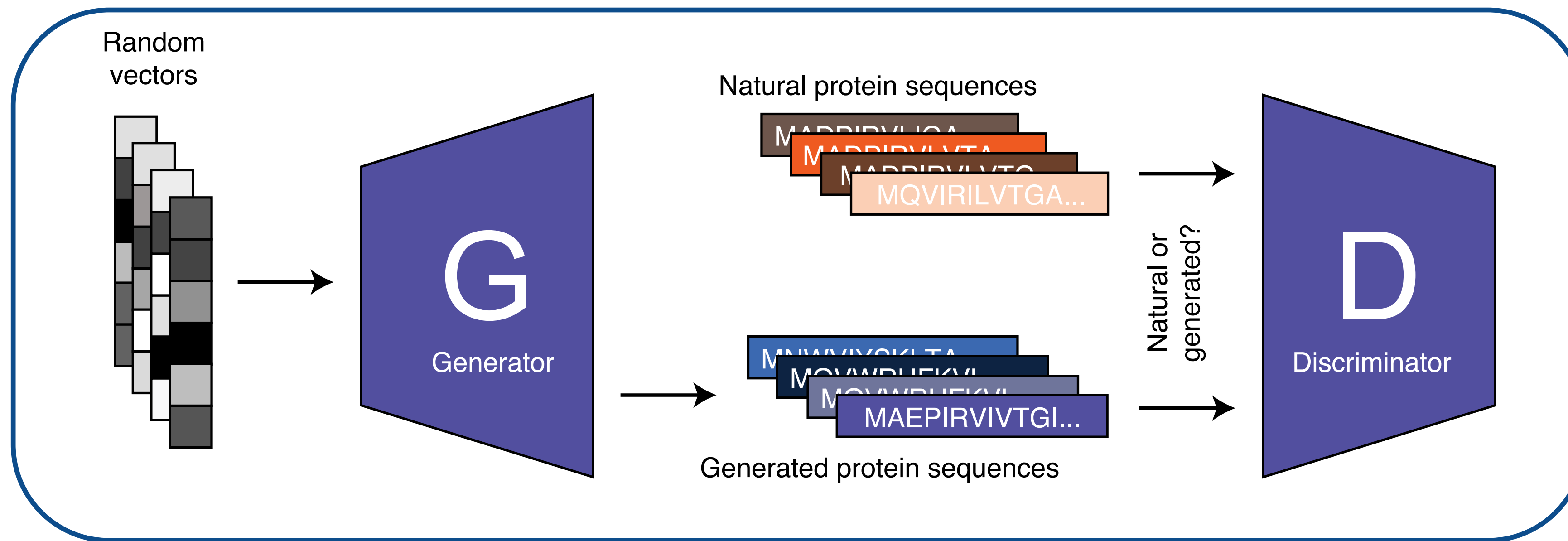


Figure: Repecka et al., *Nature Machine Intelligence* 2021



# Generative adversarial networks (GANs)

- GANs use a pair of networks to **learn (to sample from) an unknown probability distribution.**

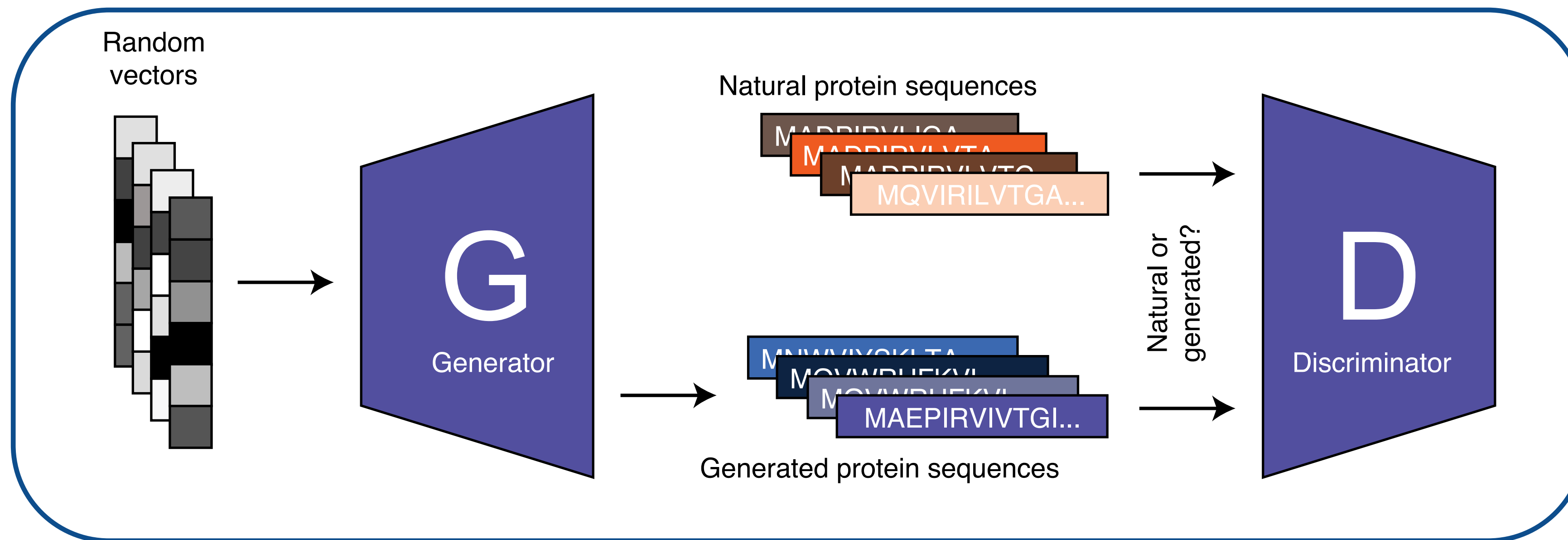


Figure: Repecka et al., *Nature Machine Intelligence* 2021

# Generative adversarial networks (GANs)

- GANs use a pair of networks to **learn (to sample from) an unknown probability distribution.**
- **Zero-sum game** between **discriminator** and **generator**—“the players”.

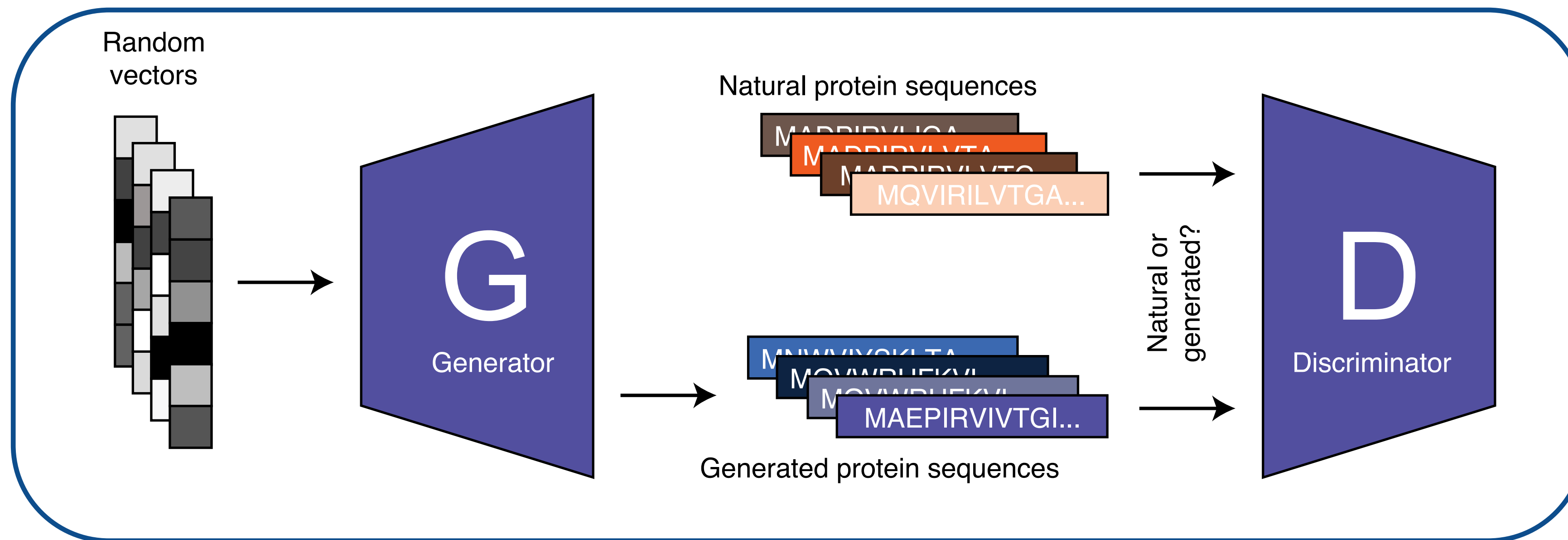


Figure: Repecka et al., *Nature Machine Intelligence* 2021



# Generative adversarial networks (GANs)

- GANs use a pair of networks to **learn (to sample from) an unknown probability distribution.**
- **Zero-sum game** between **discriminator** and **generator**—“the players”.
- **Game ends** when the players reach consensus: “fake data” looks like the “real” data.

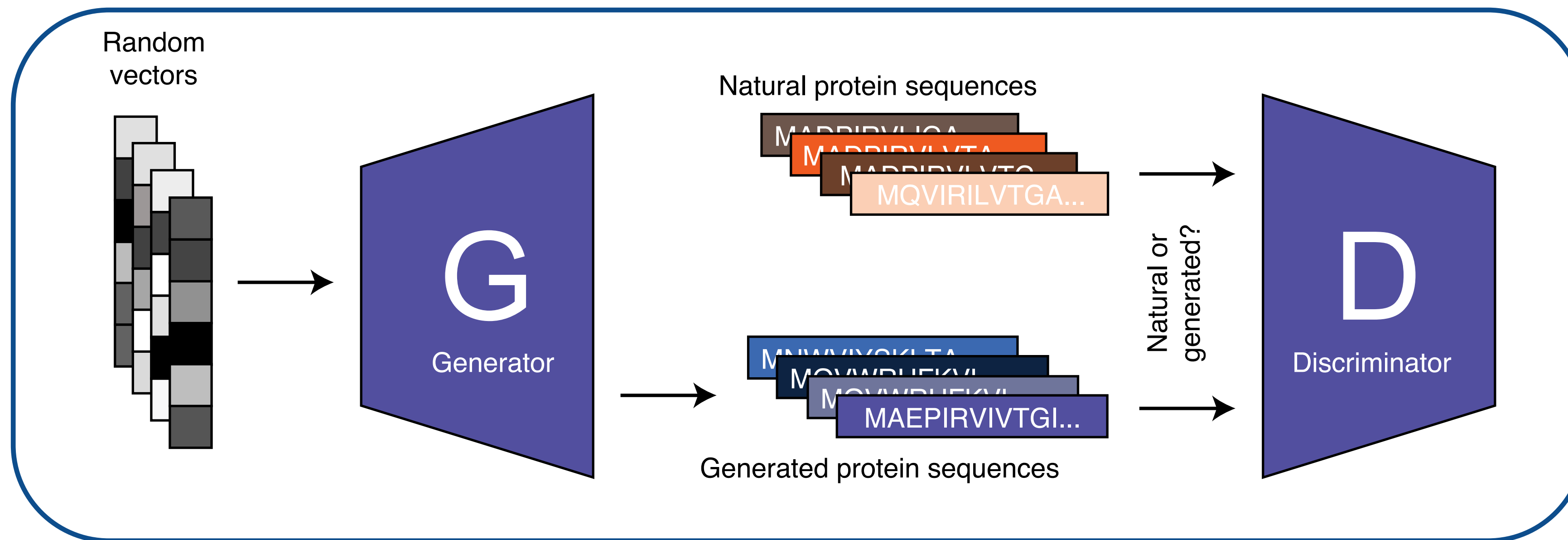
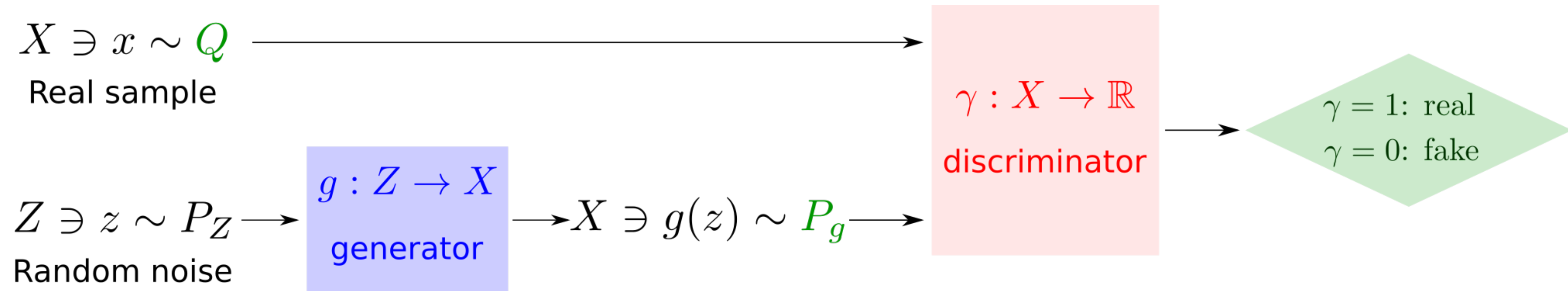


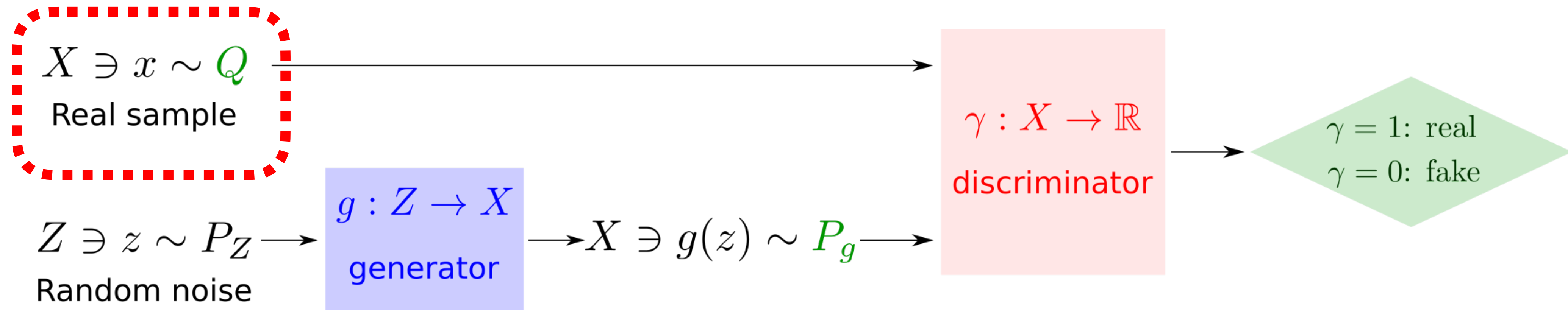
Figure: Repecka et al., *Nature Machine Intelligence* 2021

# Generative adversarial networks (GANs)

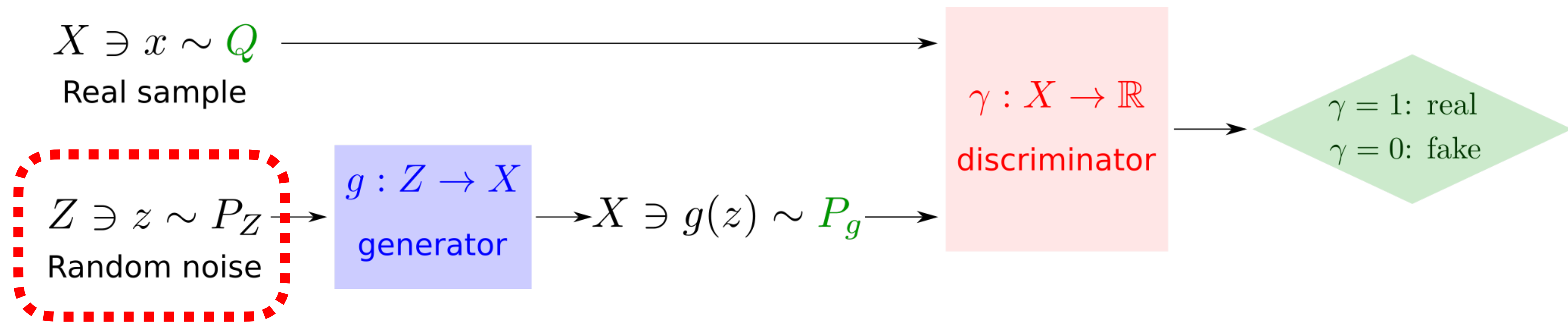




# Generative adversarial networks (GANs)

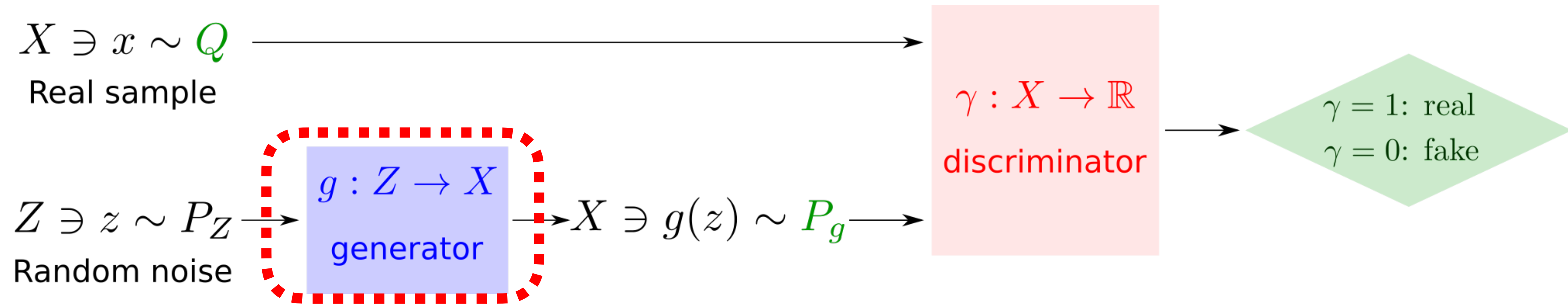


# Generative adversarial networks (GANs)

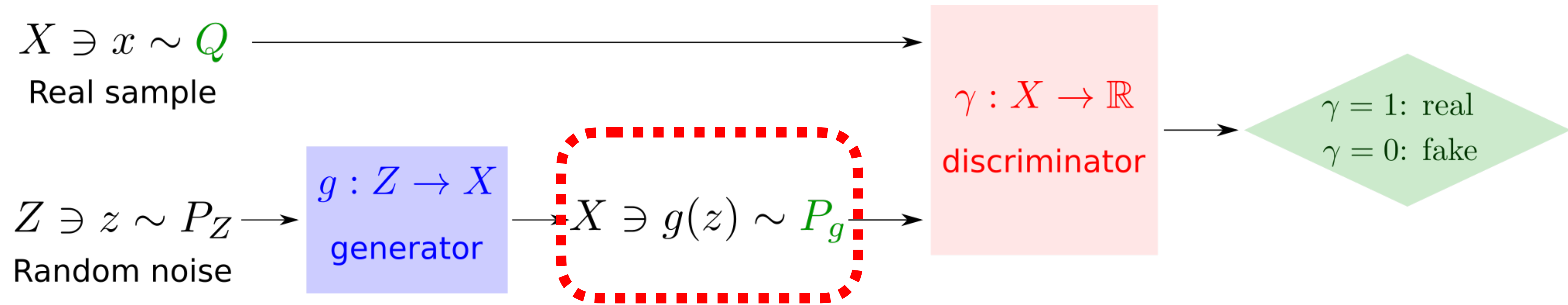




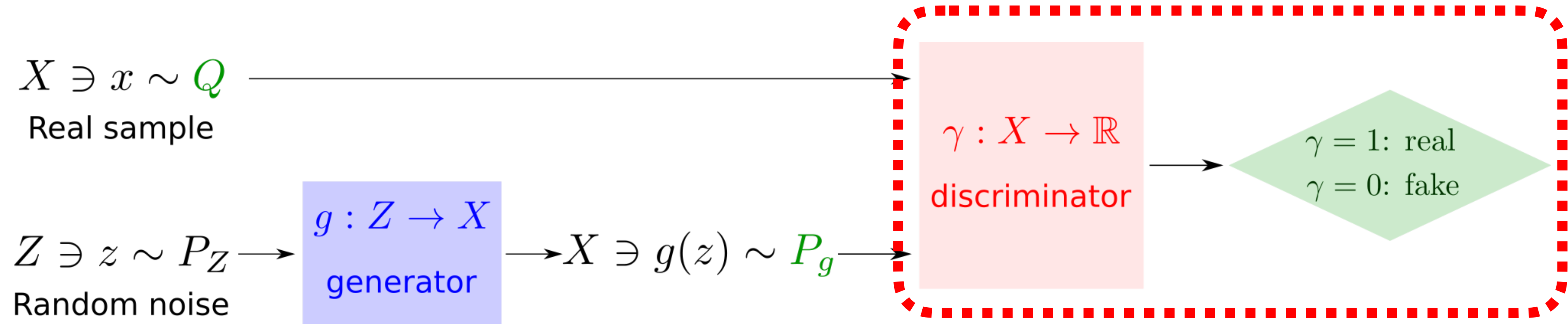
# Generative adversarial networks (GANs)



# Generative adversarial networks (GANs)

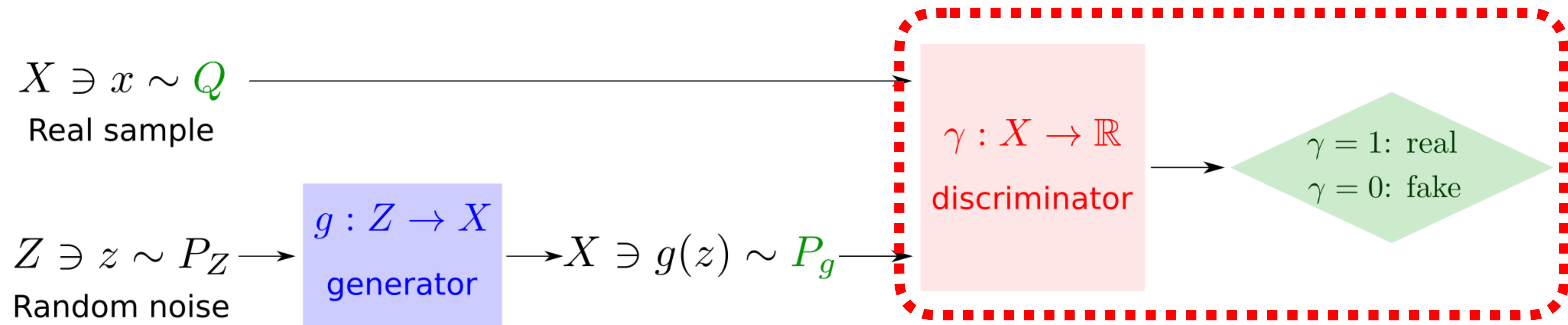


# Generative adversarial networks (GANs)





# Generative adversarial networks (GANs)



- Mathematically, GAN is minimizing some **divergence**,  $D_H^\Gamma(Q \| P_g)$ , between  $Q$  and  $P_g$ .
- $D_H^\Gamma(Q \| P_g) = \max_{\gamma \in \Gamma} H(\gamma; Q, P_g)$  is determined by  $H$  and **discriminators**  $\gamma \in \Gamma$ .

$$\min_{g \in G} D_H^\Gamma(Q \| P_g) = \min_{g \in G} \max_{\gamma \in \Gamma} H(\gamma; Q, P_g).$$

# GAN is “probability divergence” minimization

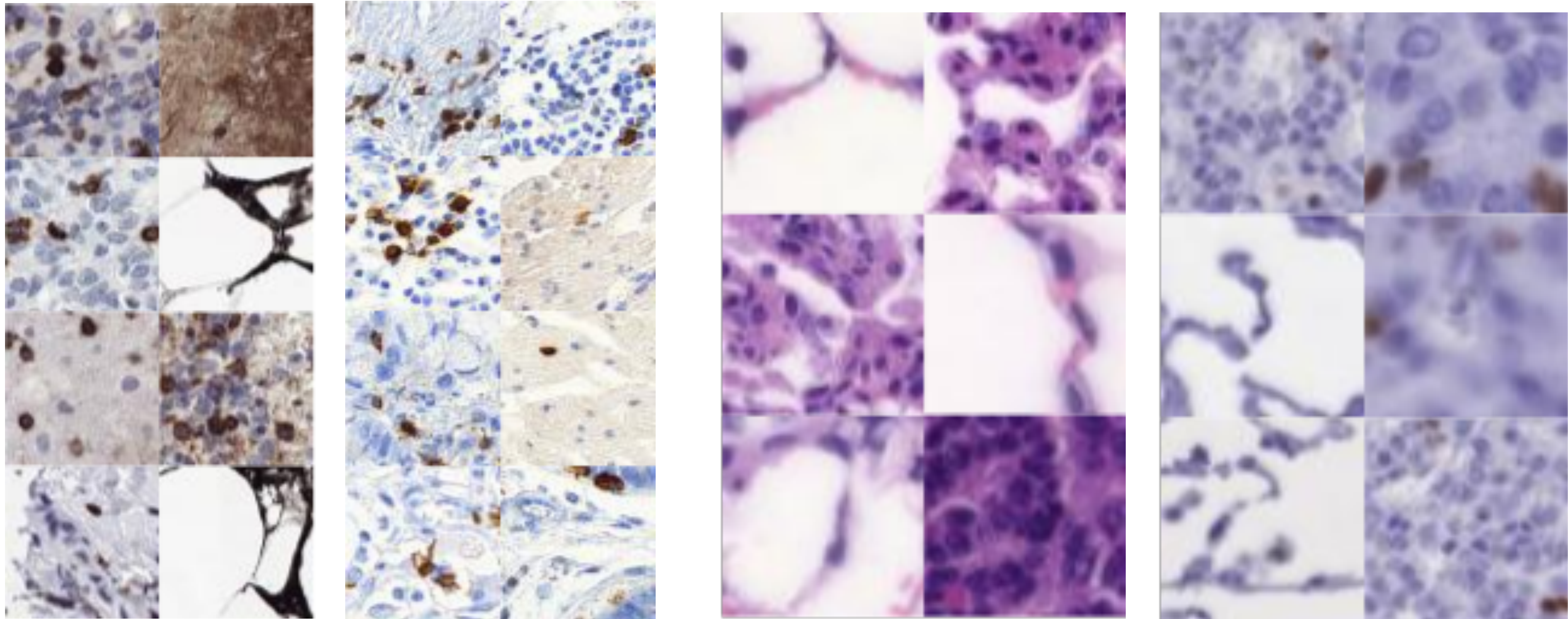
$$\min_{g \in G} D_H^\Gamma(Q \| P_g) = \min_{g \in G} \max_{\gamma \in \Gamma} H(\gamma; Q, P_g).$$

- The original GAN [Goodfellow et al., 2014]: Jensen–Shannon divergence (JSD).
- $f$ -divergences:  $D_f(Q \| P) = \sup_{\gamma \in \mathcal{M}_b(X)} \{ \mathbb{E}_Q[\gamma] - \mathbb{E}_P[f^*(\gamma)] \}$ . (KL, JSD, etc.)
- $\Gamma$ -IPM:  $W^\Gamma(Q \| P) = \sup_{\gamma \in \Gamma} \{ \mathbb{E}_Q[\gamma] - \mathbb{E}_P[\gamma] \}$ . (TV, Dudley metric, Wasserstein-1, MMD)
- Wasserstein metric and Sinkhorn divergence.



# Structured target data & distribution $Q$

$Q$



LYSTO<sup>1</sup>

ANHIR<sup>2</sup>

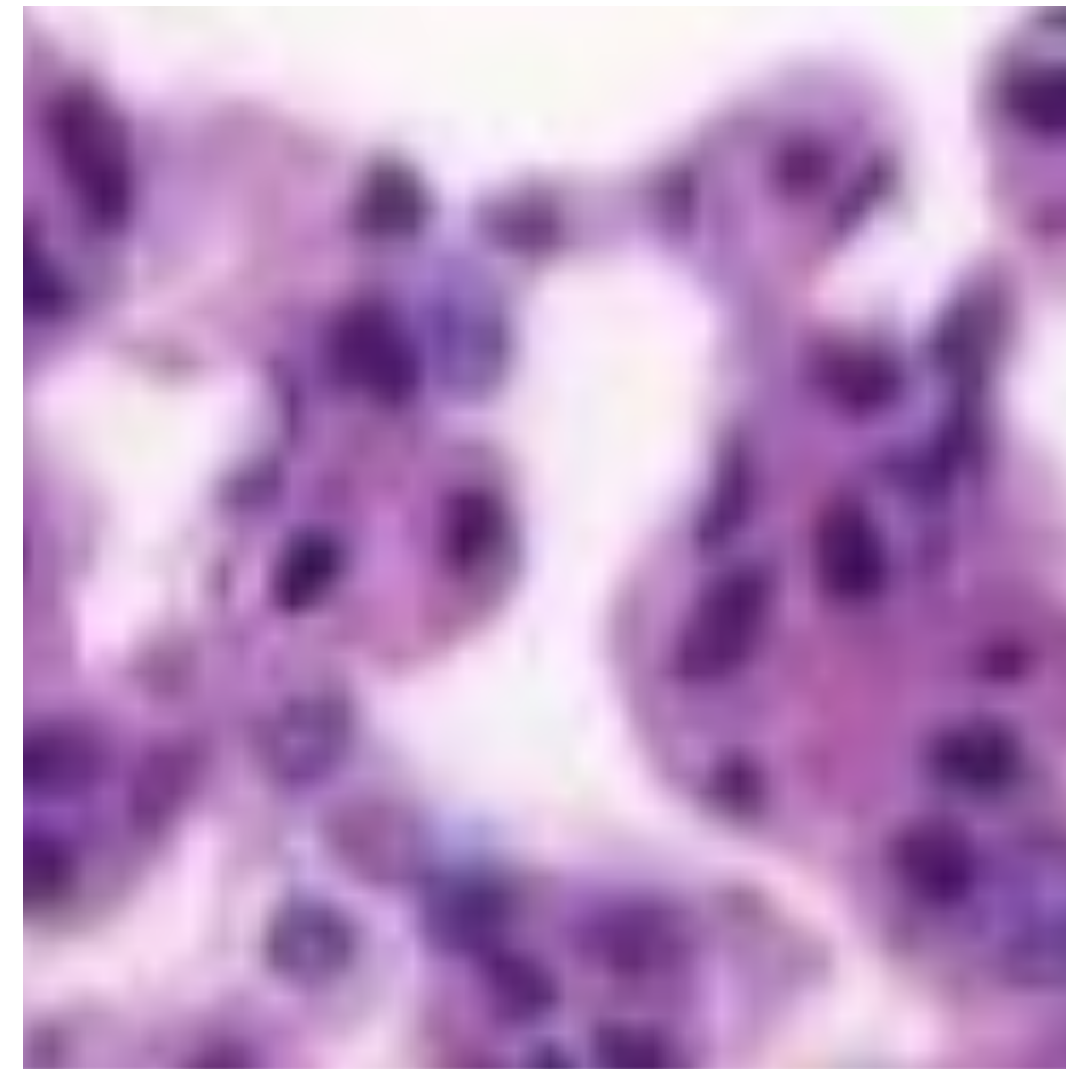
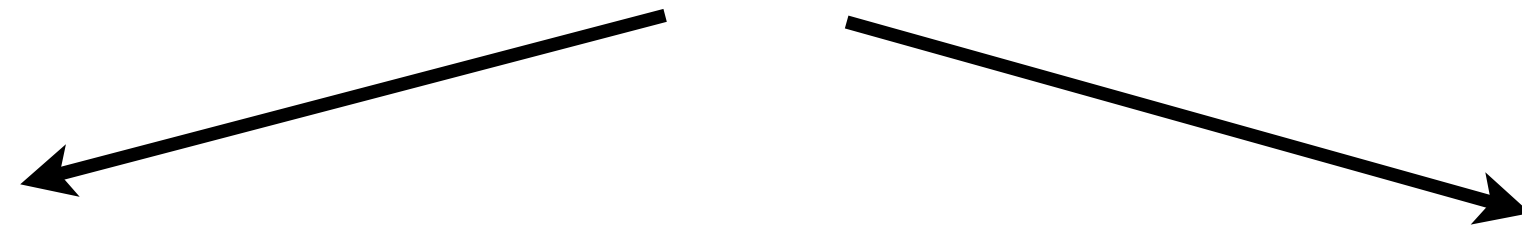
1. Ciompi et al., Zenodo 2019

2. Borovec et al., IEEE Transactions on Medical Imaging 2020



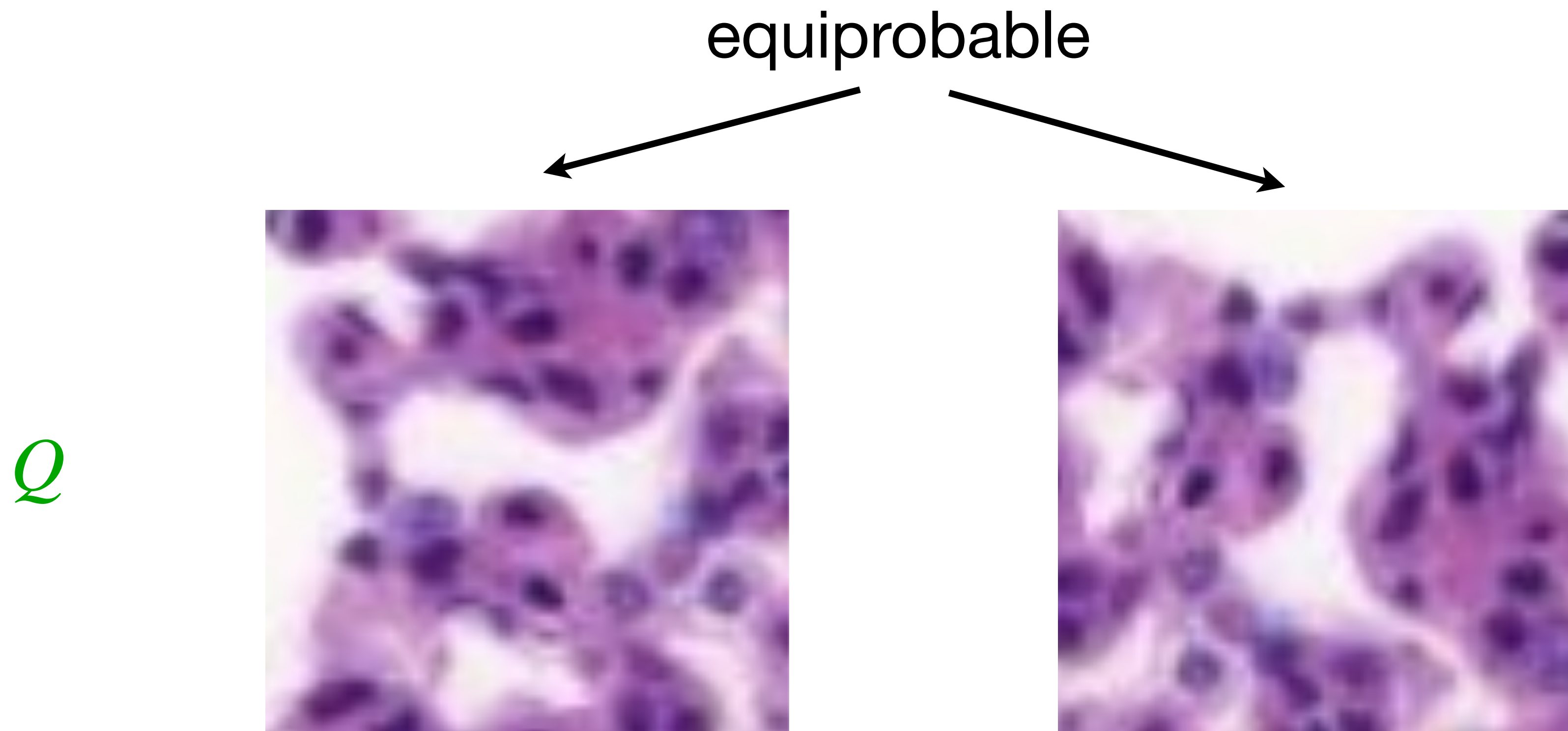
# Structured target data & distribution $Q$

equiprobable



$Q$

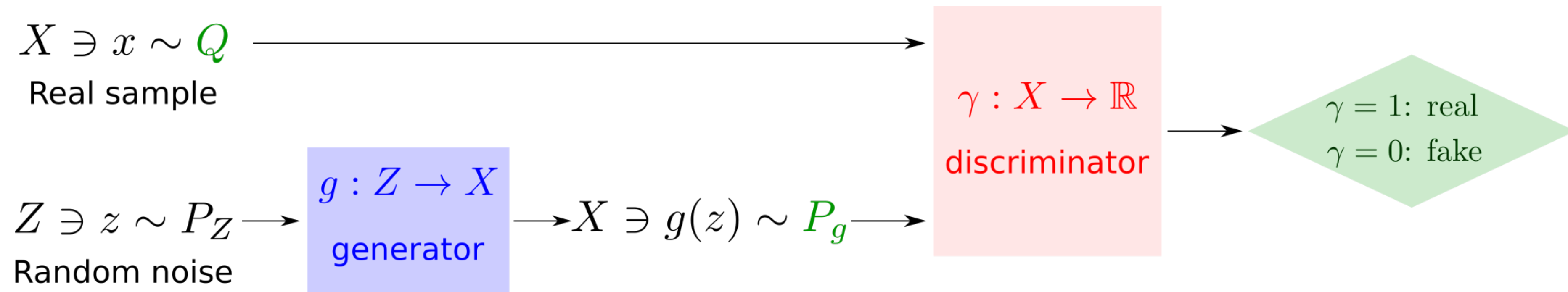
# Structured target data & distribution $Q$



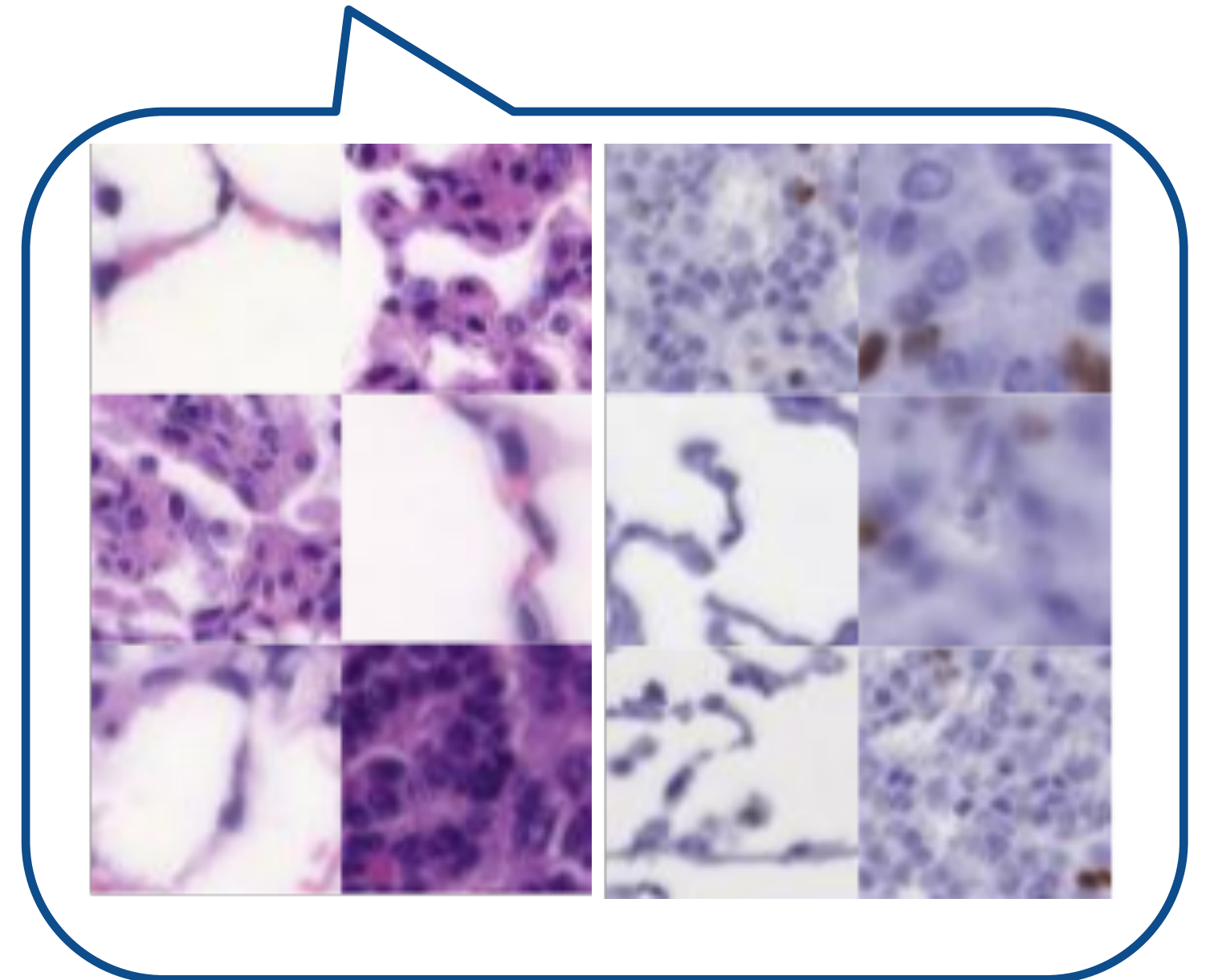
**Question:** how to build **embedded structure** into GAN players (generators and discriminators) for **data-efficient** distribution learning?



# GAN with embedded structure

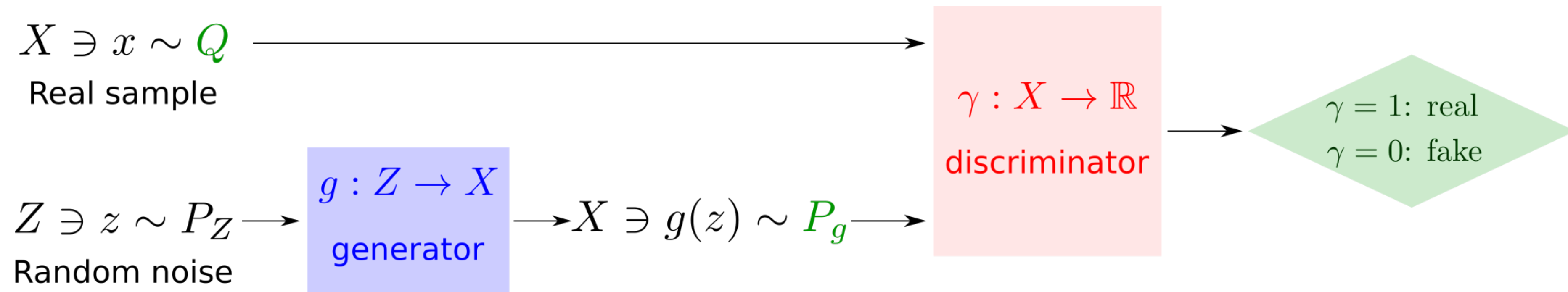


$$\min_{g \in G} D^\Gamma(Q \| P_g) = \min_{g \in G} \max_{\gamma \in \Gamma} H(\gamma; Q, P_g), \quad \underline{Q \text{ is } \Sigma\text{-invariant}}$$



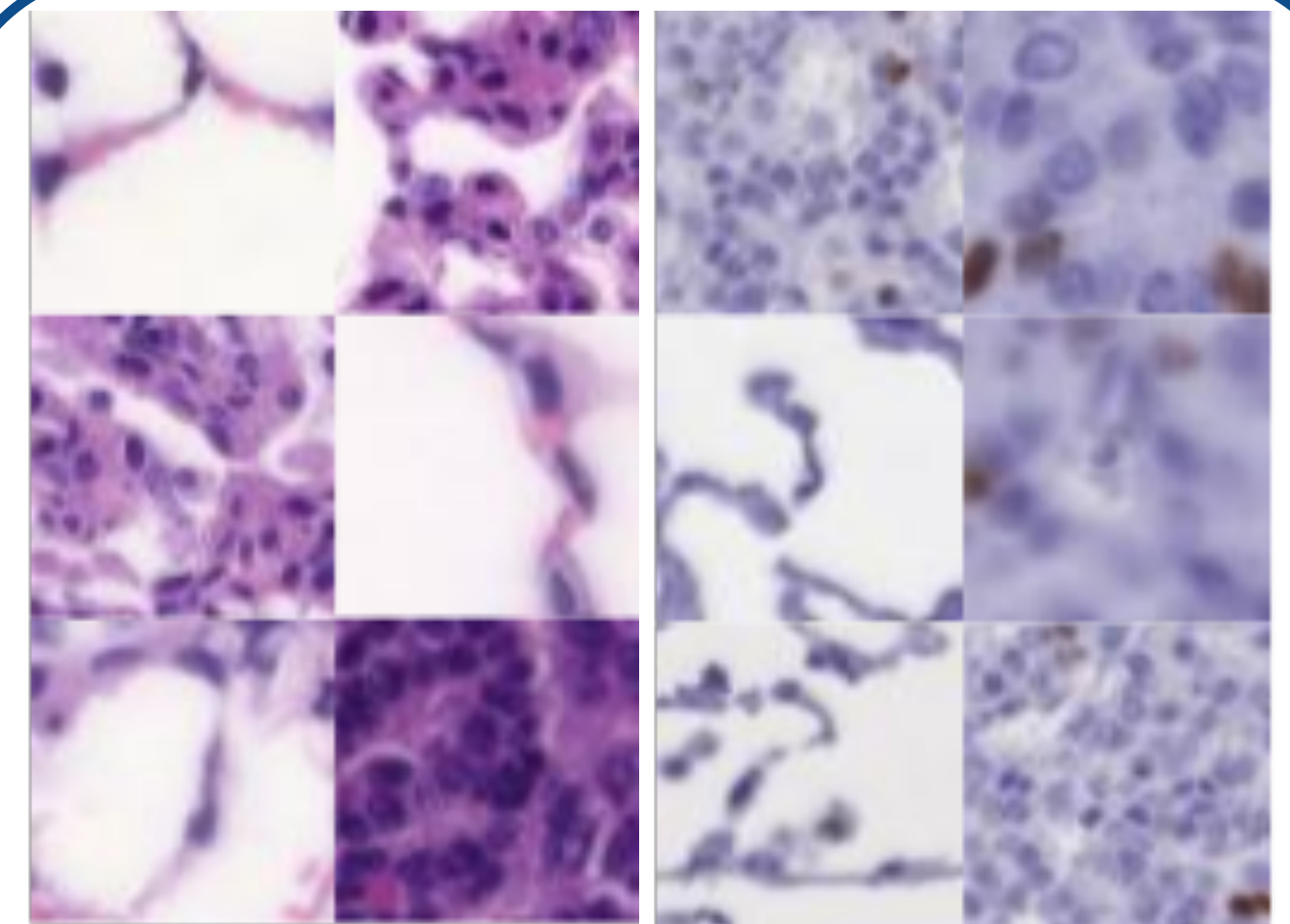


# GAN with embedded structure

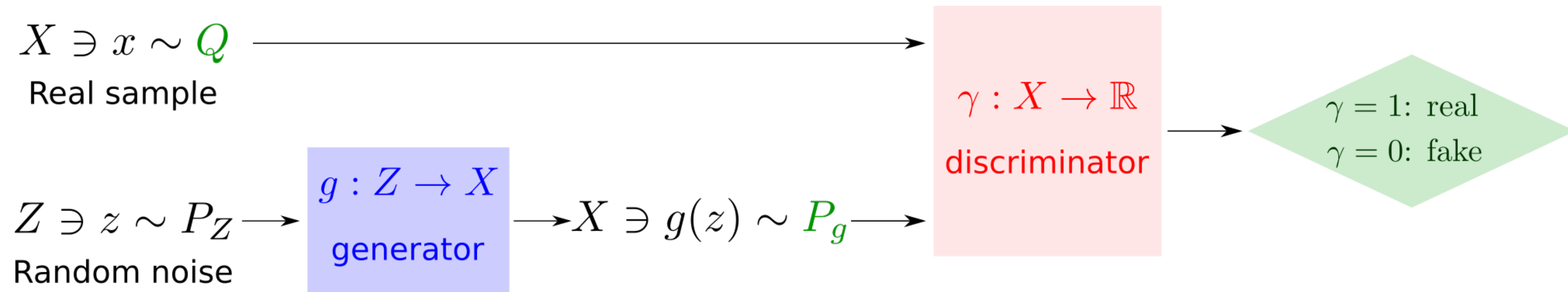


$$\min_{g \in G} D^\Gamma(Q \| P_g) = \min_{g \in G} \max_{\gamma \in \Gamma} H(\gamma; Q, P_g), \quad \underline{Q \text{ is } \Sigma\text{-invariant}}$$

- Target distribution  $Q$  is invariant under a group  $\Sigma$ .

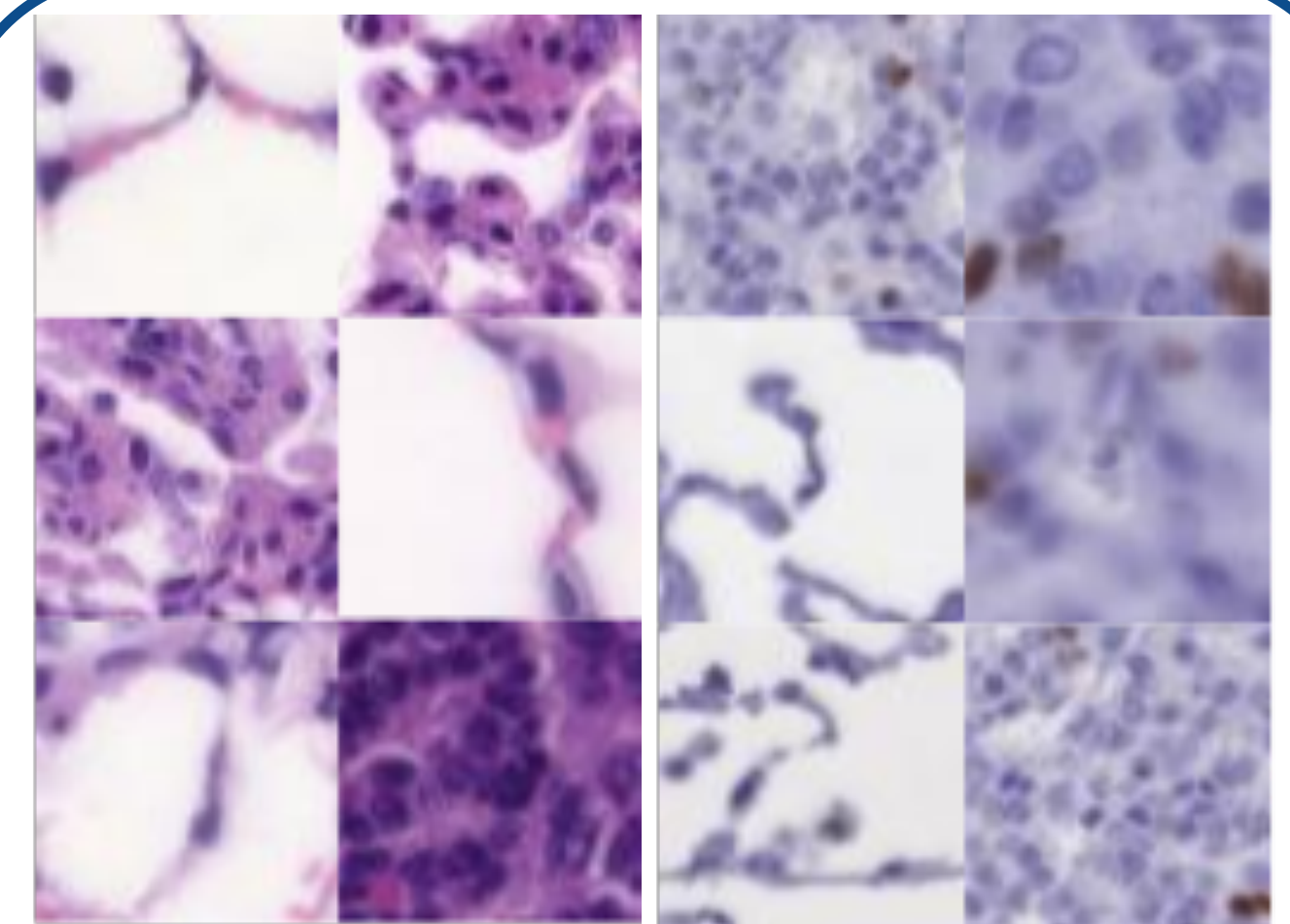


# GAN with embedded structure

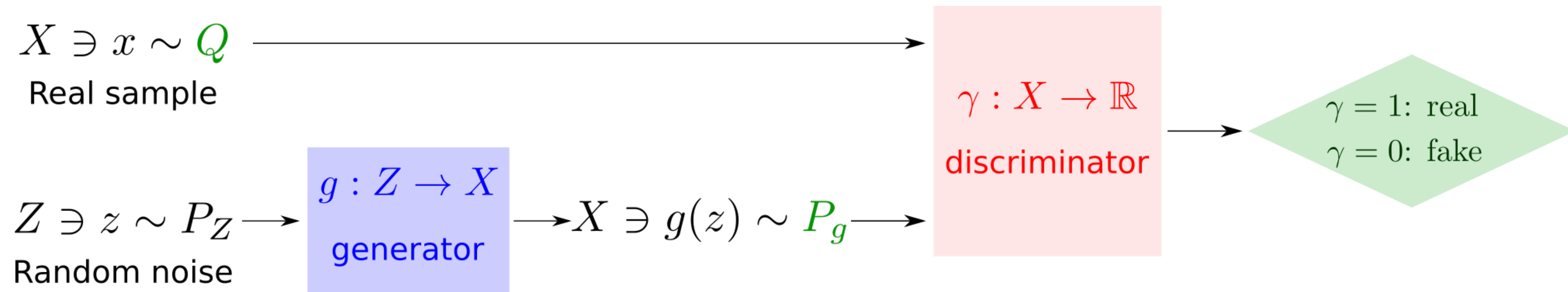


$$\min_{g \in G} D^\Gamma(Q \| P_g) = \min_{g \in G} \max_{\gamma \in \Gamma} H(\gamma; Q, P_g), \quad \underline{Q \text{ is } \Sigma\text{-invariant}}$$

- Target distribution  $Q$  is invariant under a group  $\Sigma$ .
- $\Sigma$ : rotation, reflection, permutation, etc.

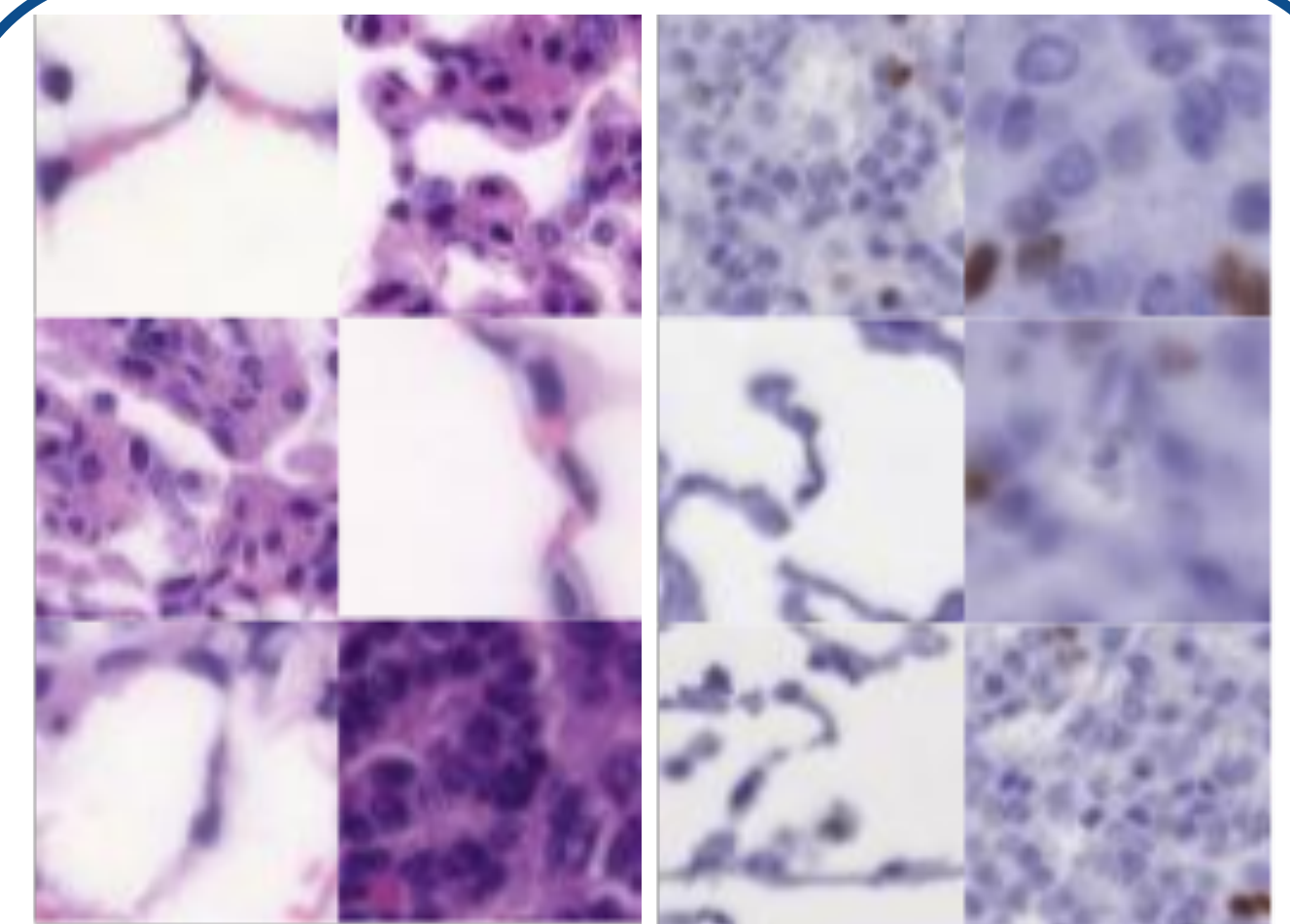


# GAN with embedded structure



$$\min_{g \in G} D^\Gamma(Q \| P_g) = \min_{g \in G} \max_{\gamma \in \Gamma} H(\gamma; Q, P_g), \quad \underline{Q \text{ is } \Sigma\text{-invariant}}$$

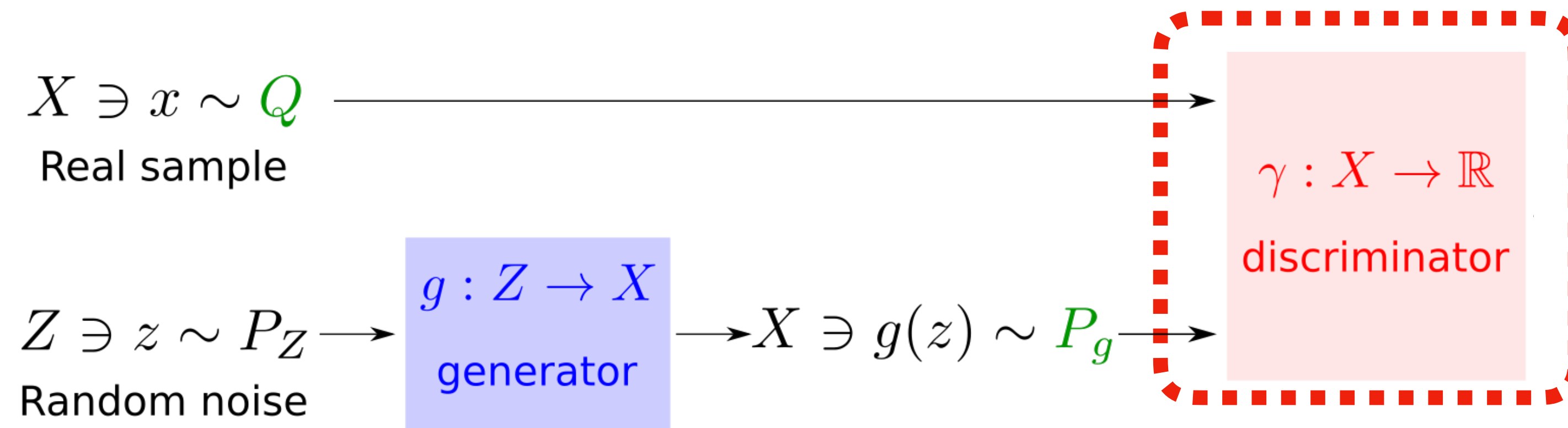
- Target distribution  $Q$  is invariant under a group  $\Sigma$ .
- $\Sigma$ : rotation, reflection, permutation, etc.
- **How to incorporate structure into  $g$  and  $\gamma$ ?**



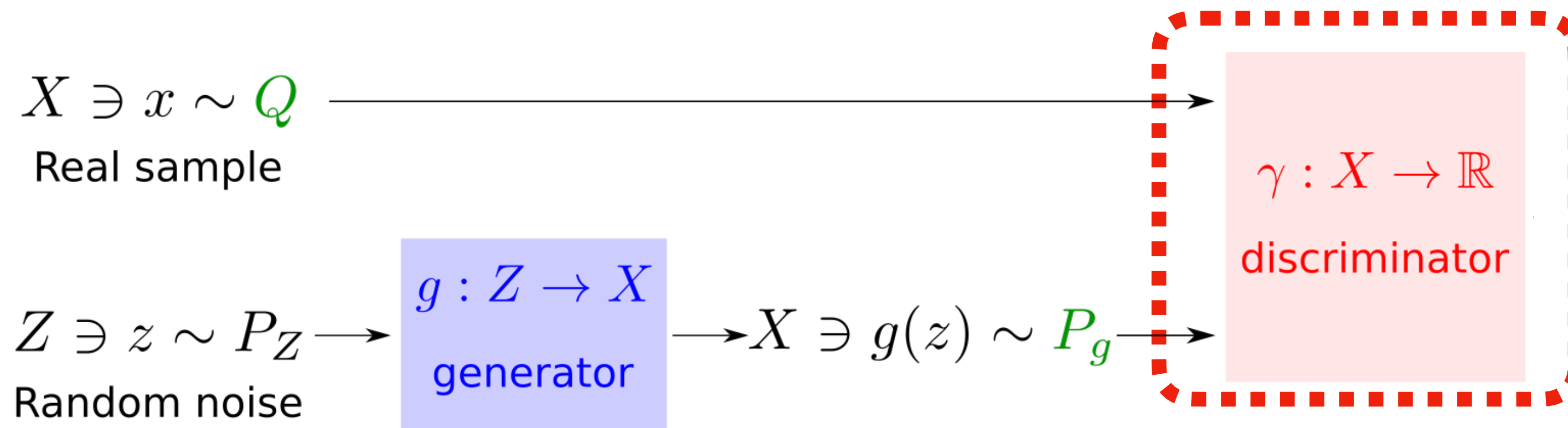


# **Theorem 1: “smarter” discriminator**

# Theorem 1: “smarter” discriminator



# Theorem 1: “smarter” discriminator



**Theorem** [Birrell, Katsoulakis, Rey-Bellet, **Z.**, *ICML* 2022]

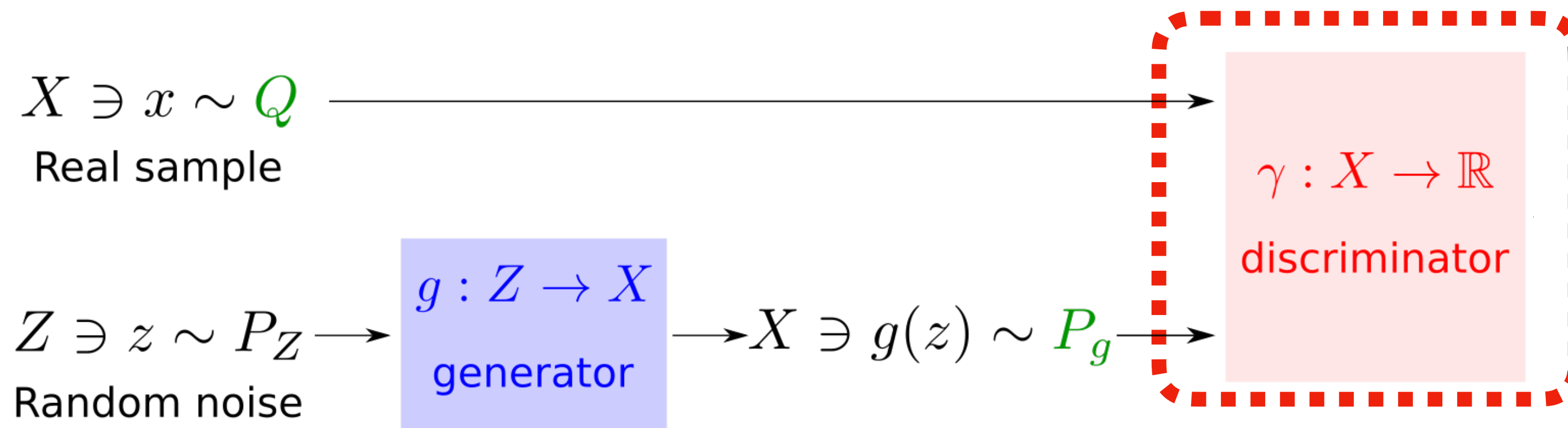
Under mild assumptions on  $\Sigma$  and  $\Gamma$ , if the distributions  $P, Q$  are  $\Sigma$ -invariant, then

$$D^\Gamma(Q||P) = D^{\Gamma_\Sigma^{\text{inv}}}(Q||P) = \sup_{\gamma \in \Gamma_\Sigma^{\text{inv}}} H(\gamma; Q, P),$$

- $\Gamma_\Sigma^{\text{inv}} \subset \Gamma$  is the subset of  $\Sigma$ -invariant “smarter” discriminators



# Theorem 1: “smarter” discriminator



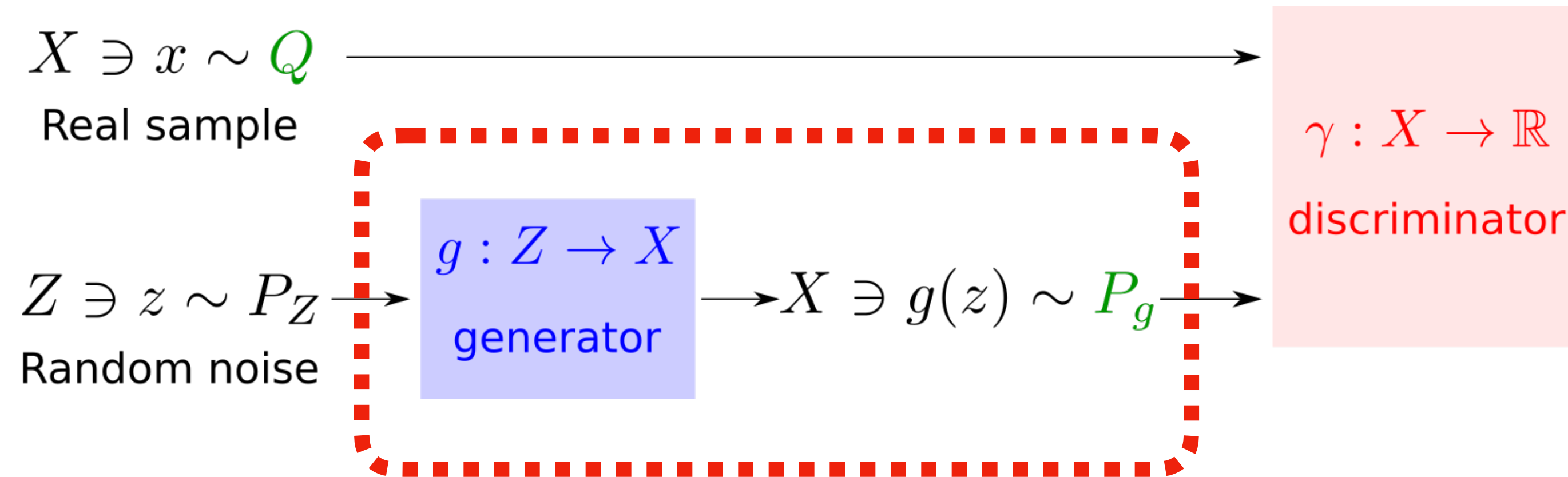
**Theorem** [Birrell, Katsoulakis, Rey-Bellet, **Z.**, *ICML* 2022]

Under mild assumptions on  $\Sigma$  and  $\Gamma$ , if the distributions  $P, Q$  are  $\Sigma$ -invariant, then

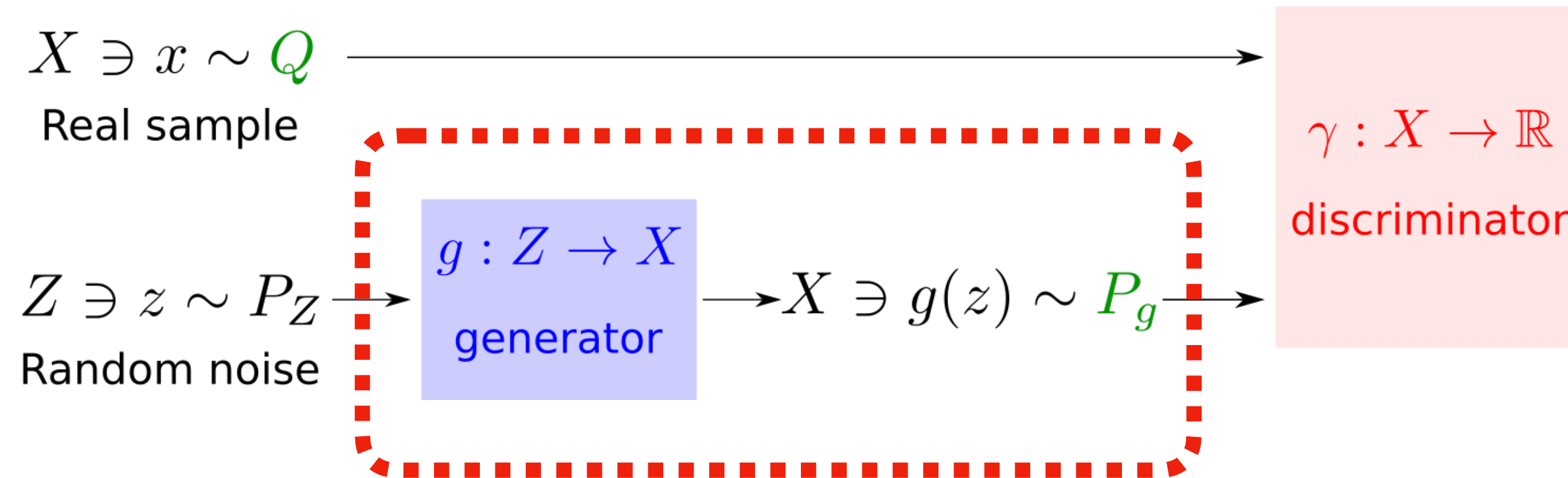
$$D^\Gamma(Q||P) = D^{\Gamma_\Sigma^{\text{inv}}}(Q||P) = \sup_{\gamma \in \Gamma_\Sigma^{\text{inv}}} H(\gamma; Q, P),$$

- $\Gamma_\Sigma^{\text{inv}} \subset \Gamma$  is the subset of  $\Sigma$ -invariant “smarter” discriminators
- $\Gamma_\Sigma^{\text{inv}}$  serves as an **unbiased regularization** to prevent **discriminator overfitting**.

# Theorem 2: “smarter” generator



# Theorem 2: “smarter” generator

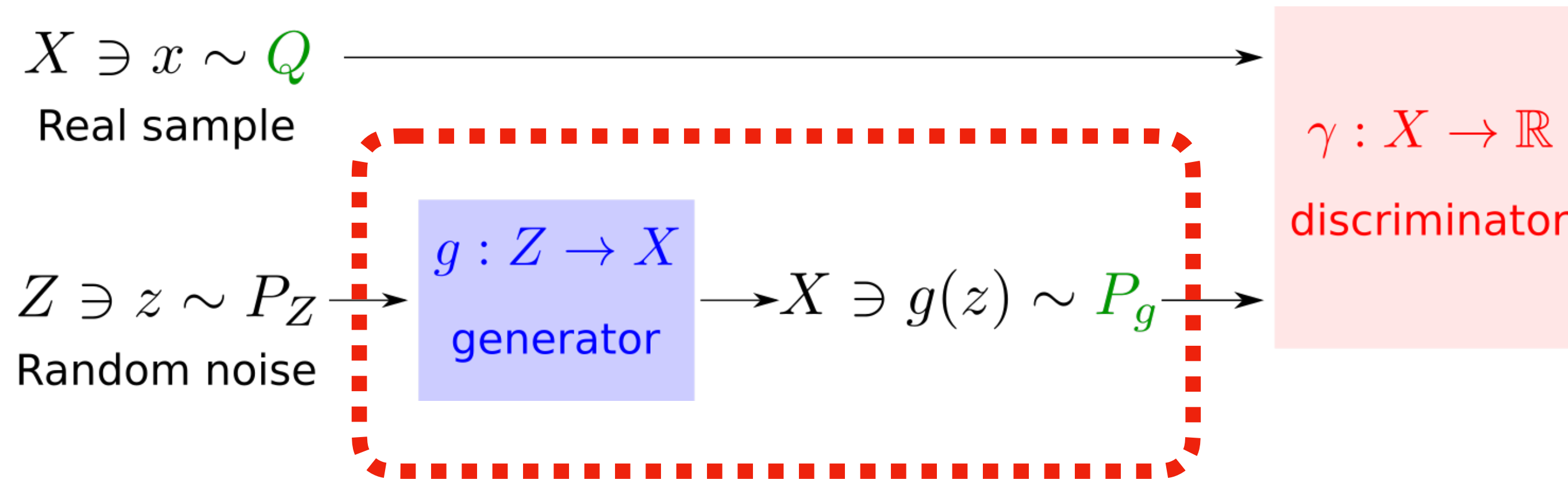


**Theorem** [Birrell, Katsoulakis, Rey-Bellet, **Z.**, *ICML* 2022]

If  $P_Z$  is  $\Sigma$ -invariant and  $g : Z \rightarrow X$  is  $\Sigma$ -equivariant, the generated measure  $P_g$  is  $\Sigma$ -invariant.



# Theorem 2: “smarter” generator



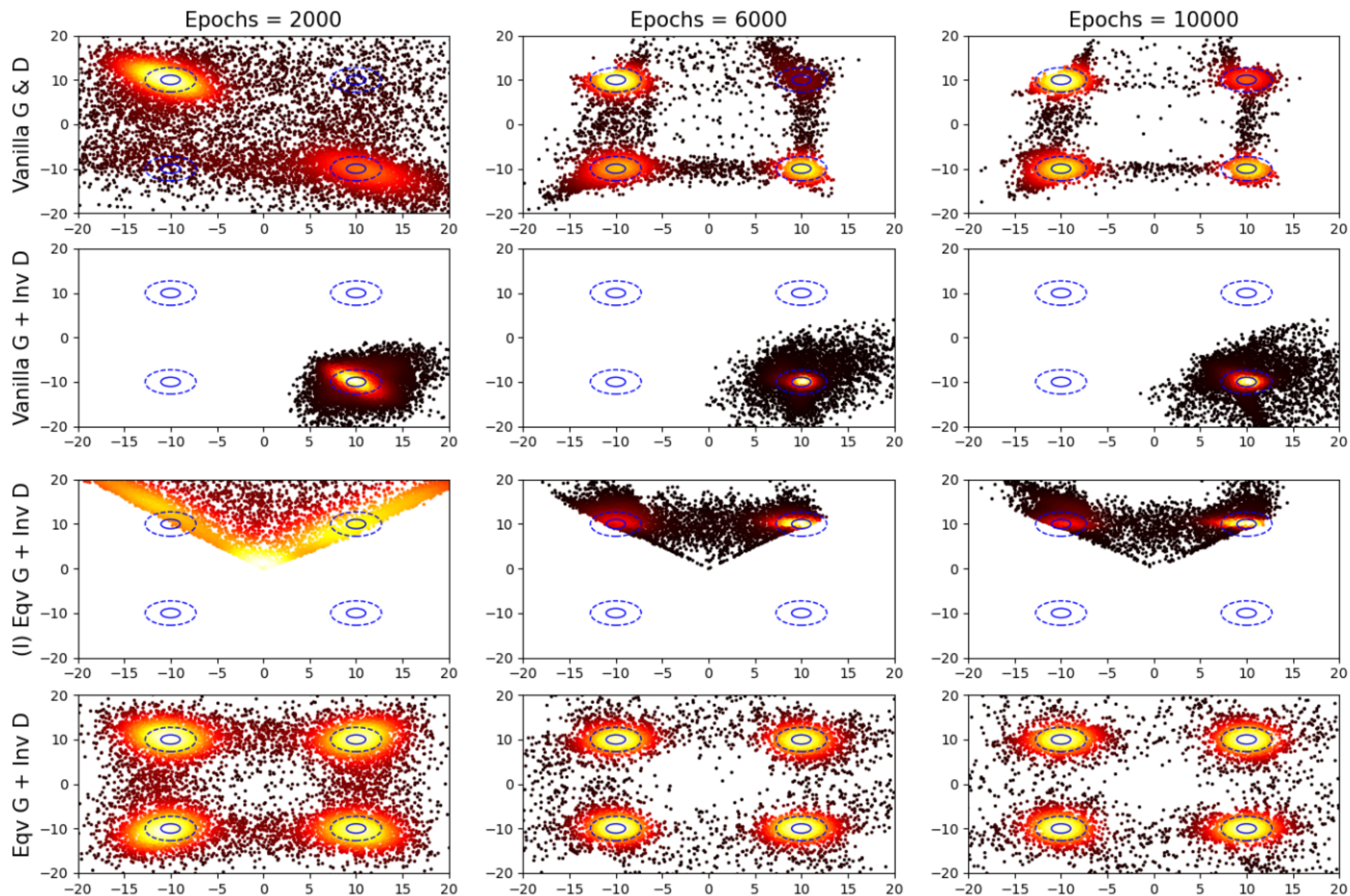
**Theorem** [Birrell, Katsoulakis, Rey-Bellet, **Z.**, *ICML* 2022]

If  $P_Z$  is  $\Sigma$ -invariant and  $g : Z \rightarrow X$  is  $\Sigma$ -equivariant, the generated measure  $P_g$  is  $\Sigma$ -invariant.

- Structure information embedded in the “smarter” generator **and** noise source.



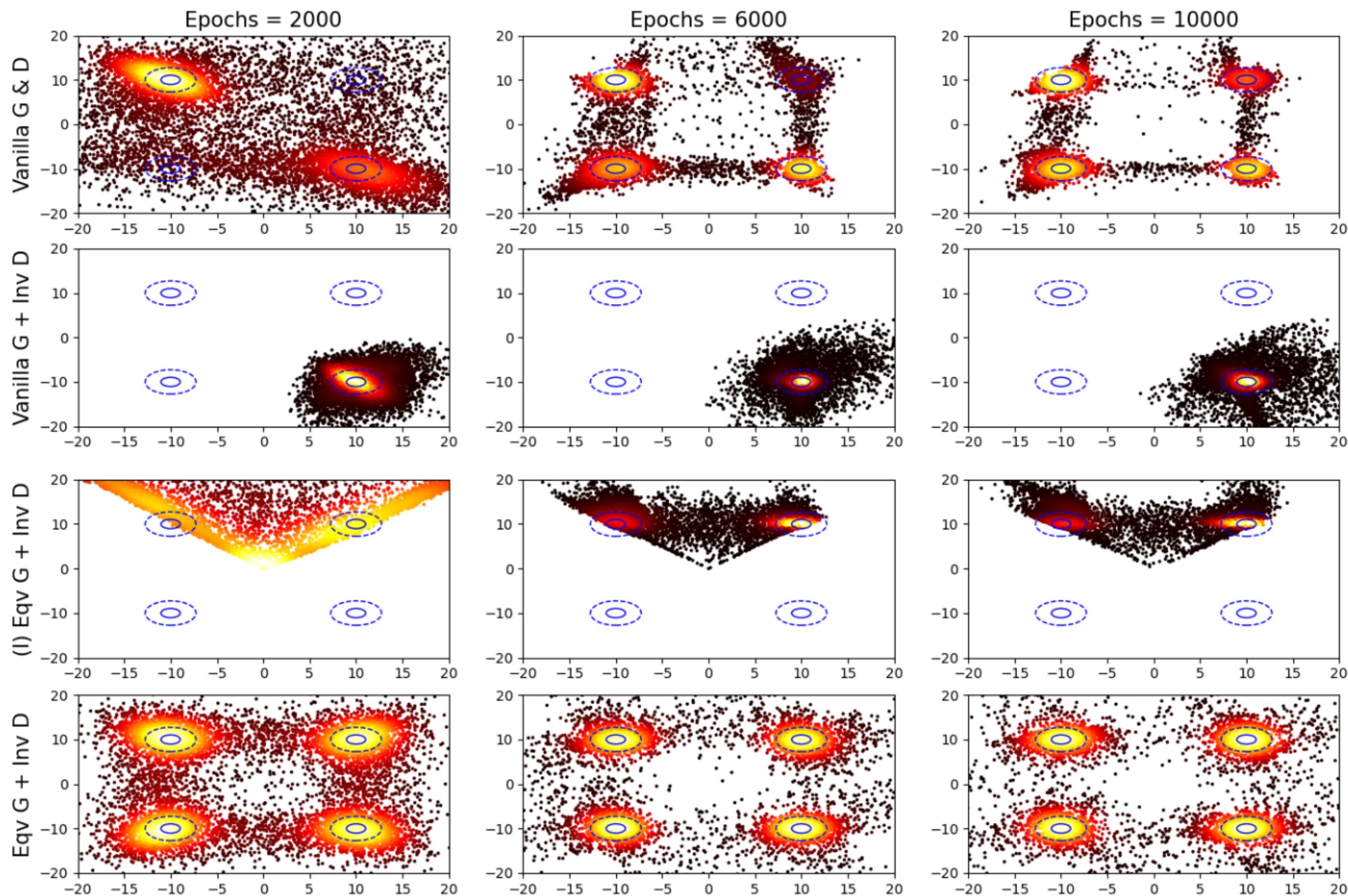
# Two “smart” players





# Two “smart” players

“Ignorant”  
players need lots of  
data, lots of time...  
(the usual GANs)

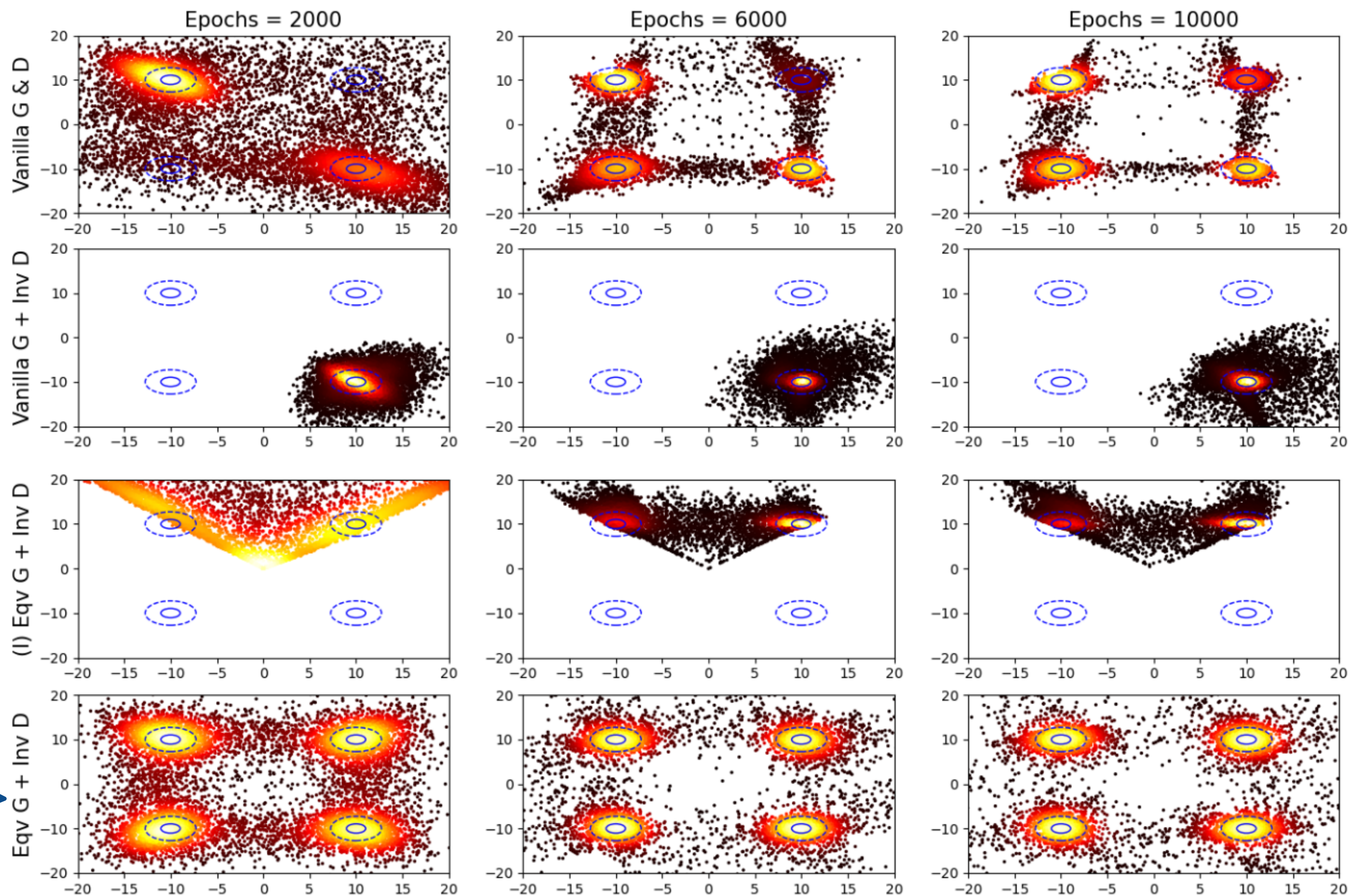




# Two “smart” players

“Ignorant”  
players need lots of  
data, lots of time...  
(the usual GANs)

“Smart” players  
learn faster and better  
(our GANs)



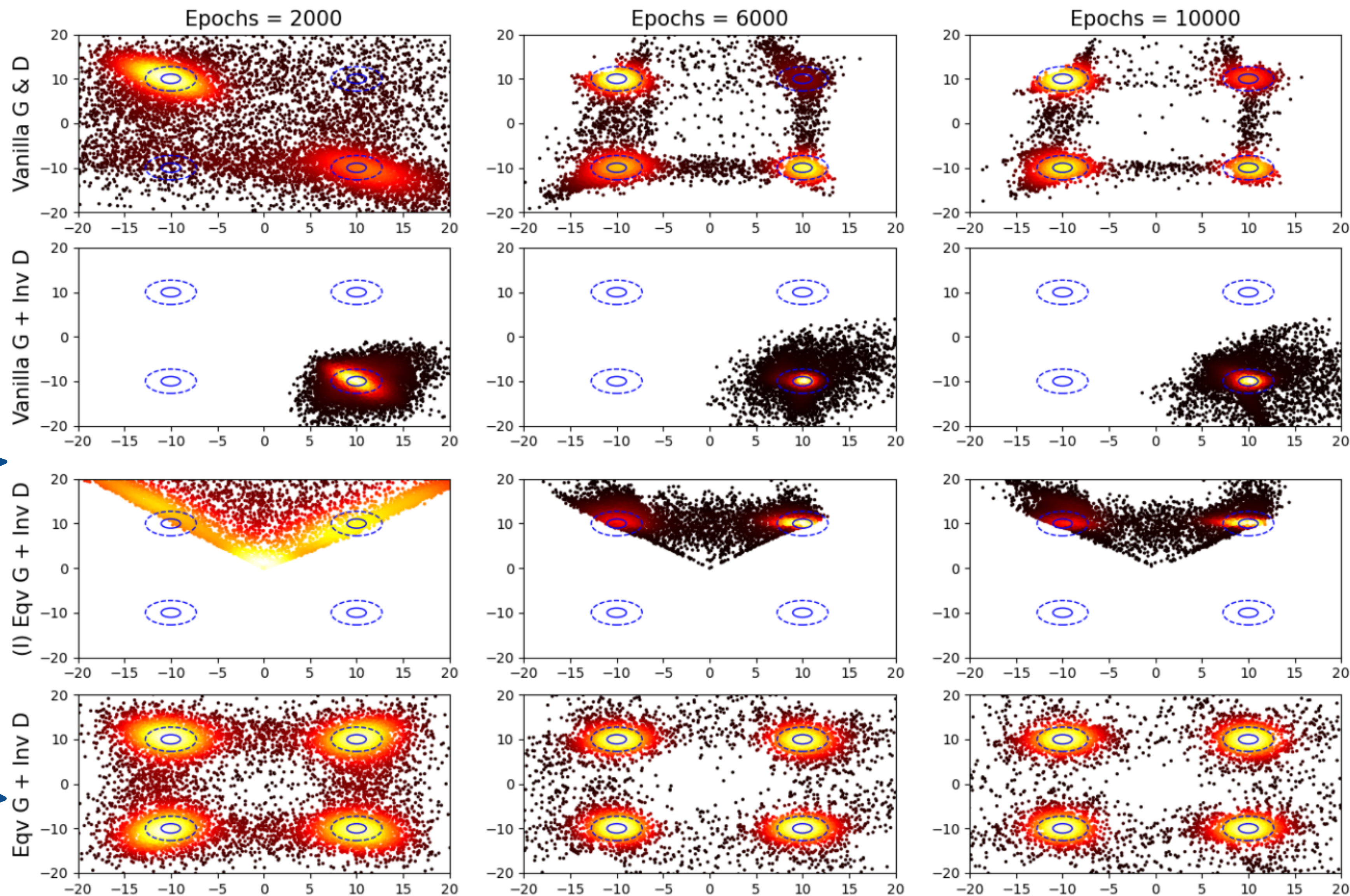


# Two “smart” players

“Ignorant”  
players need lots of  
data, lots of time...  
(the usual GANs)

Players need to be  
“equally smart”: no  
weak links!

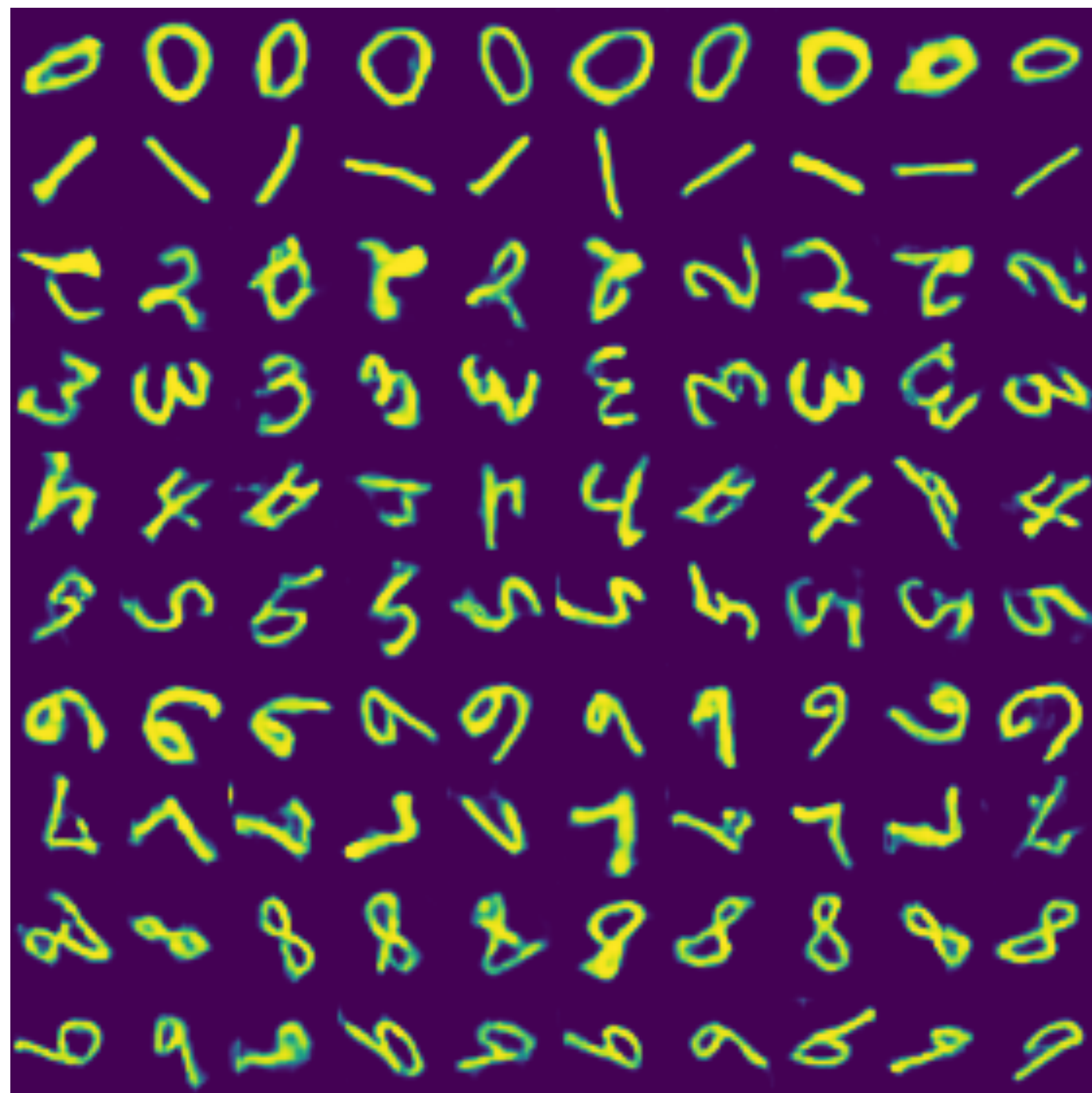
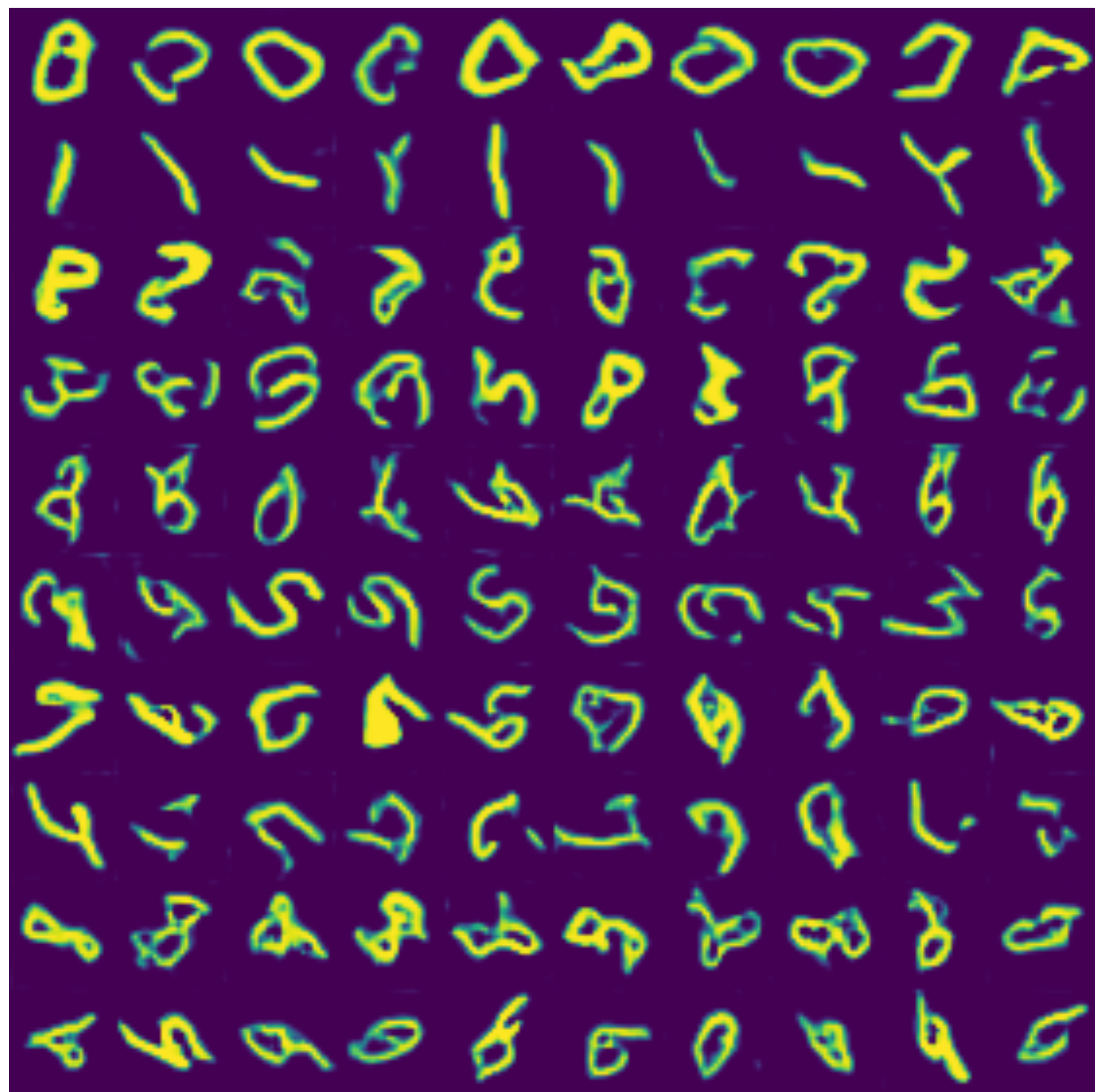
“Smart” players  
learn faster and better  
(our GANs)





# RotMNIST with 1% training samples

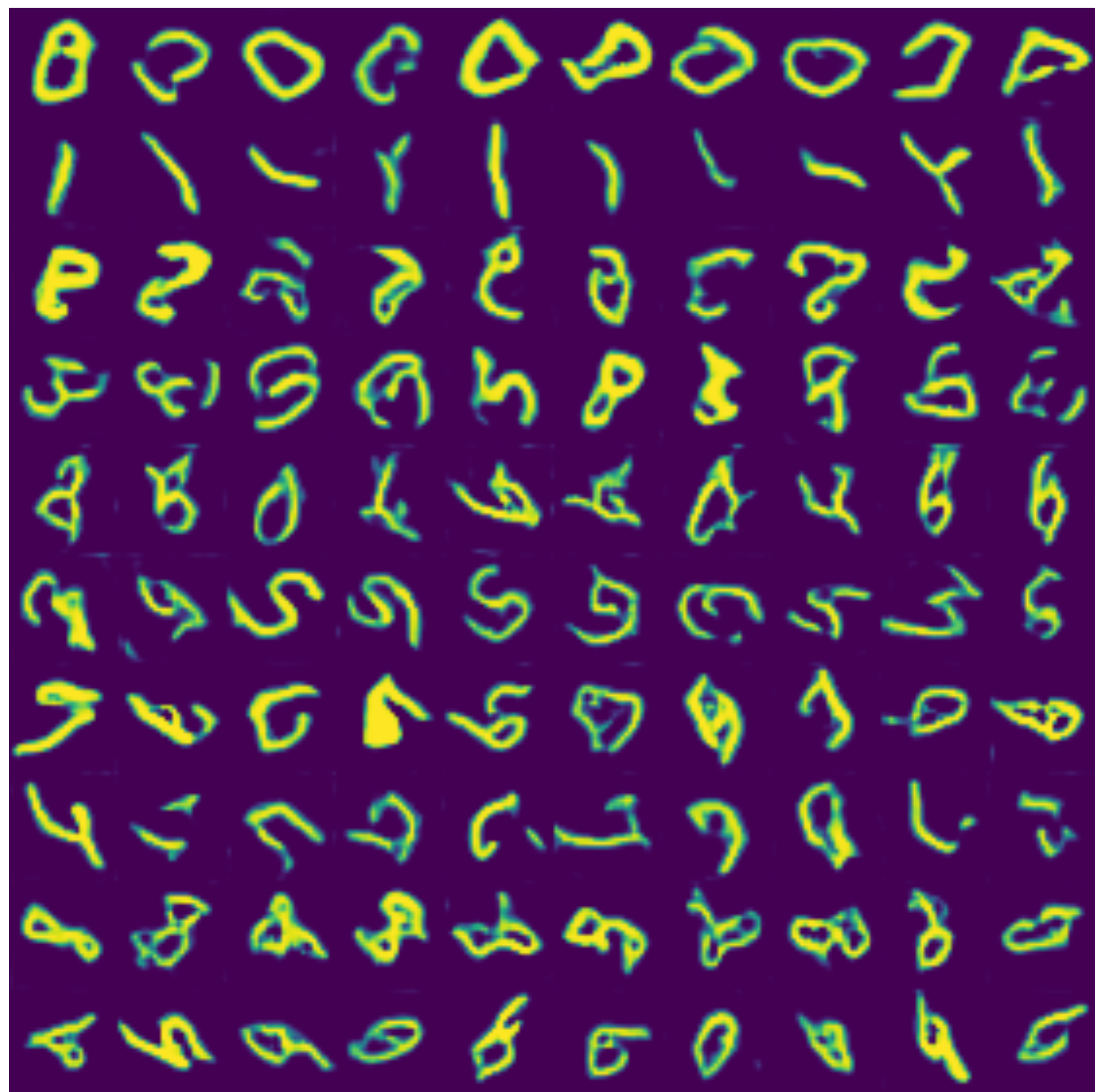
“Ignorant” players



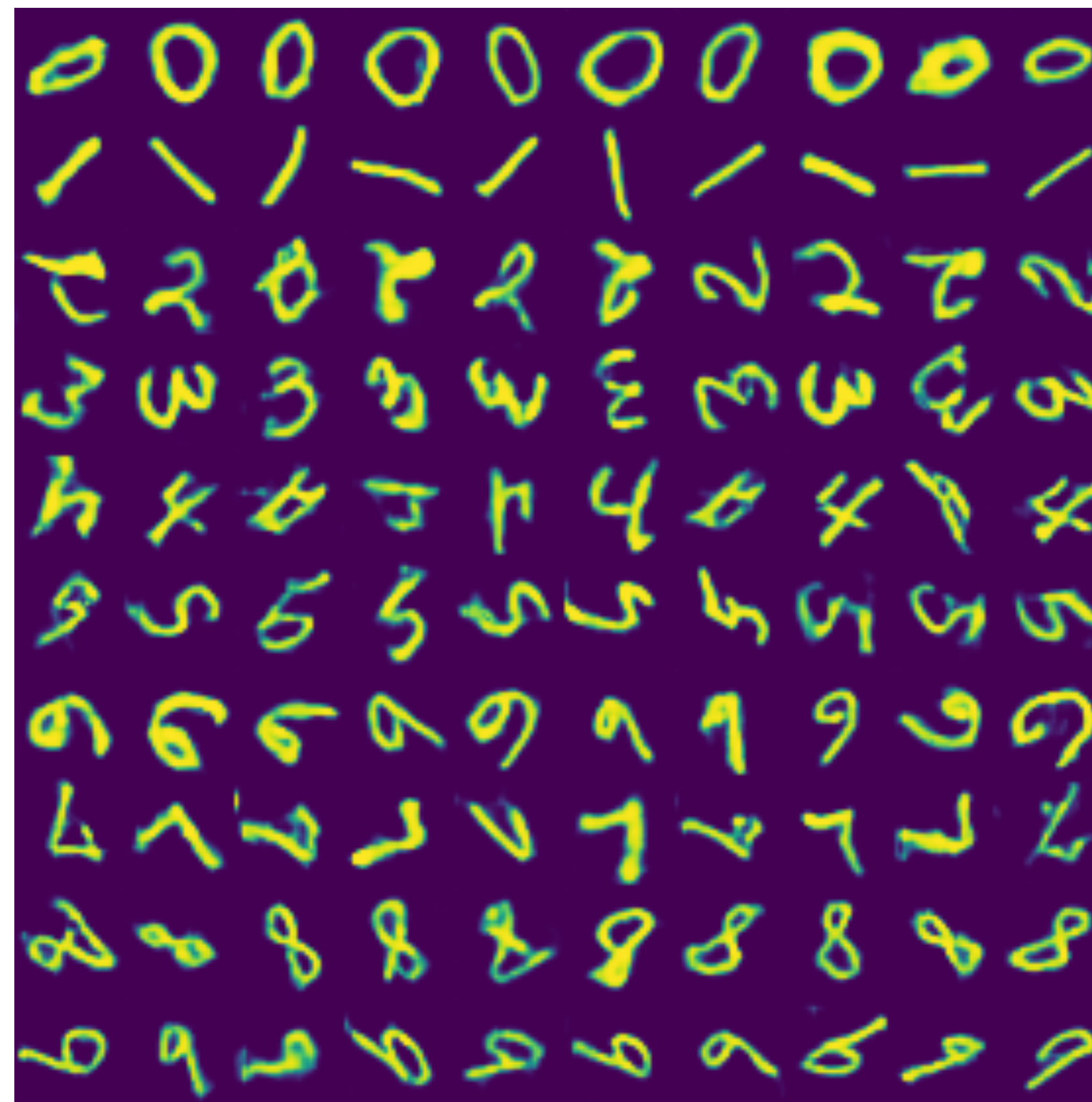


# RotMNIST with 1% training samples

“Ignorant” players



“Smart” players



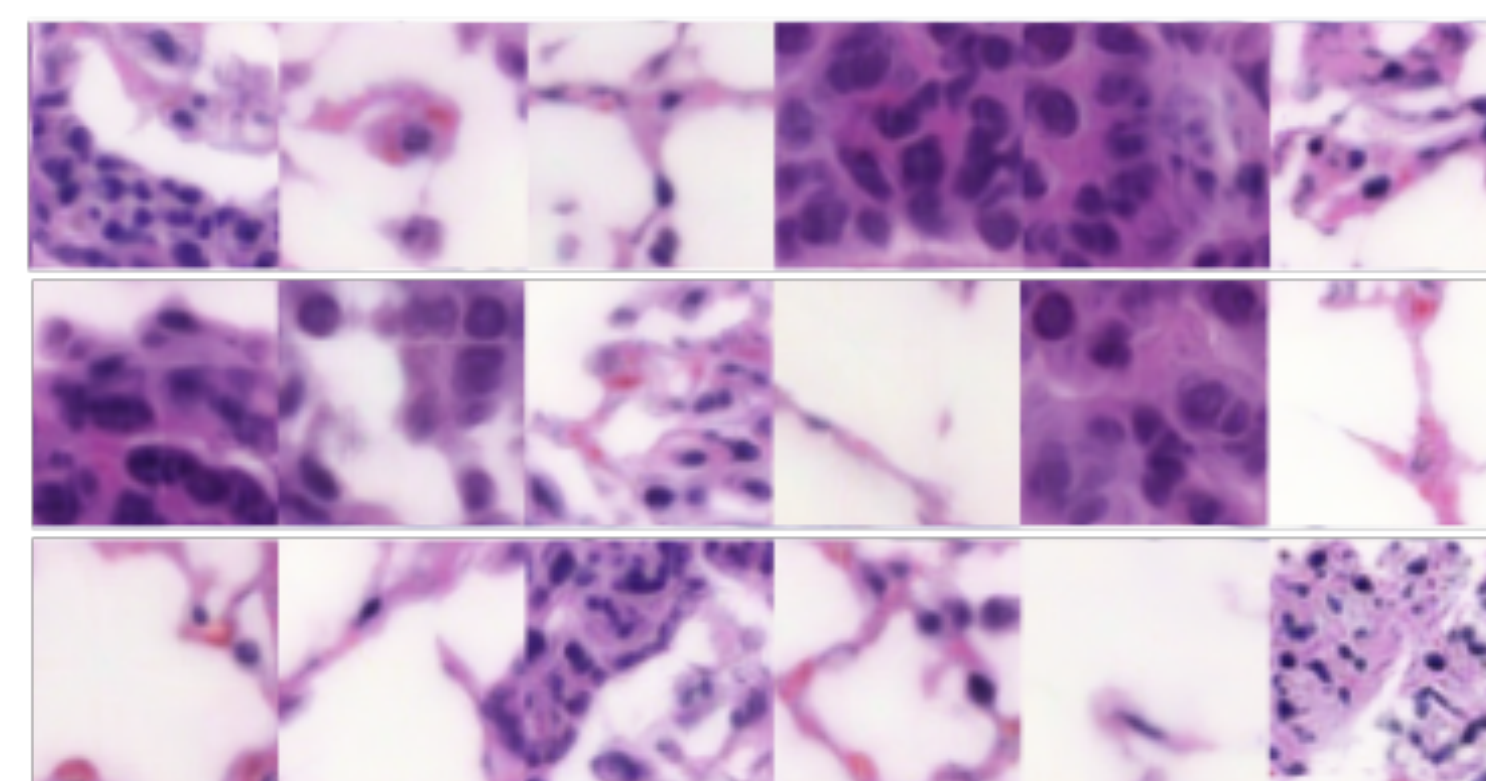
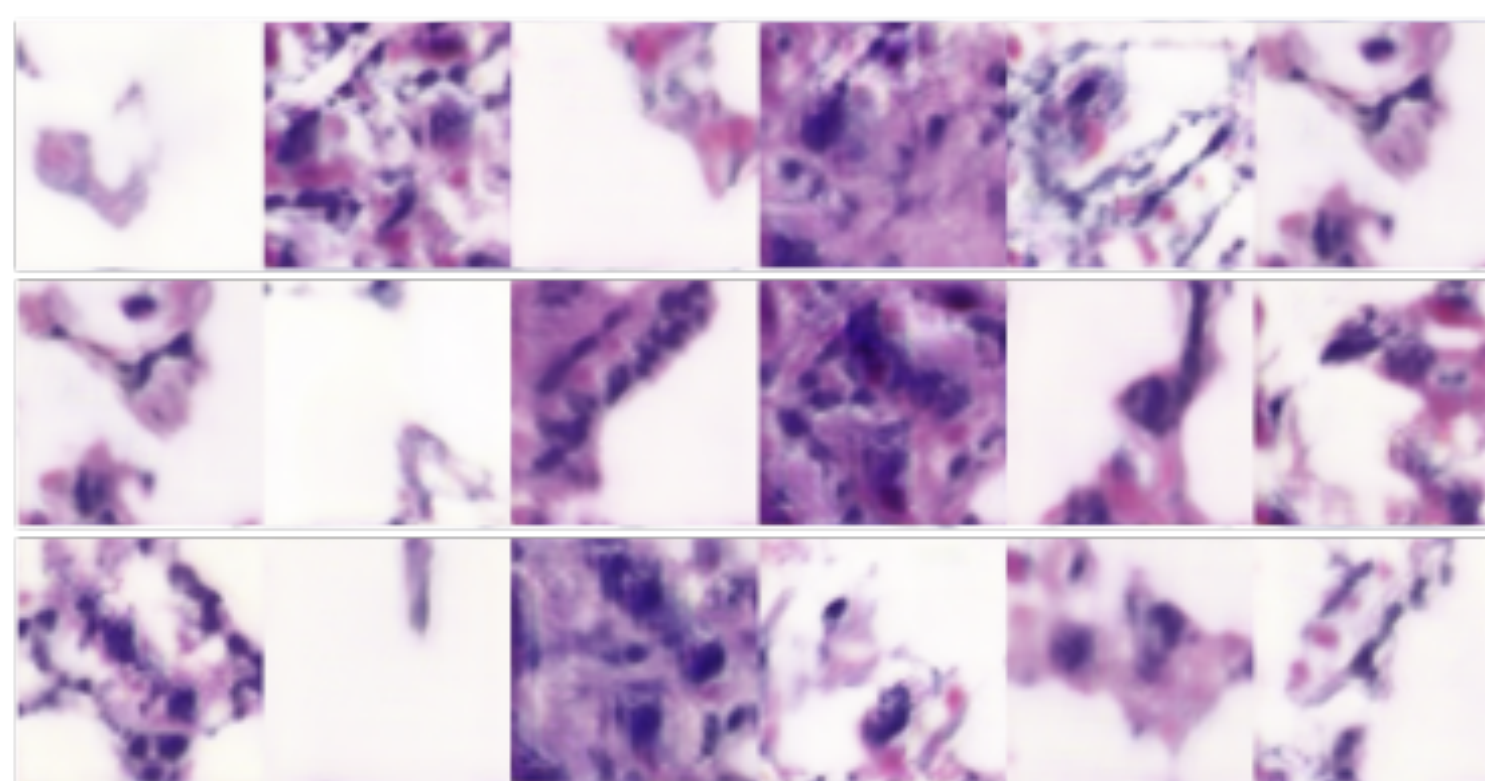
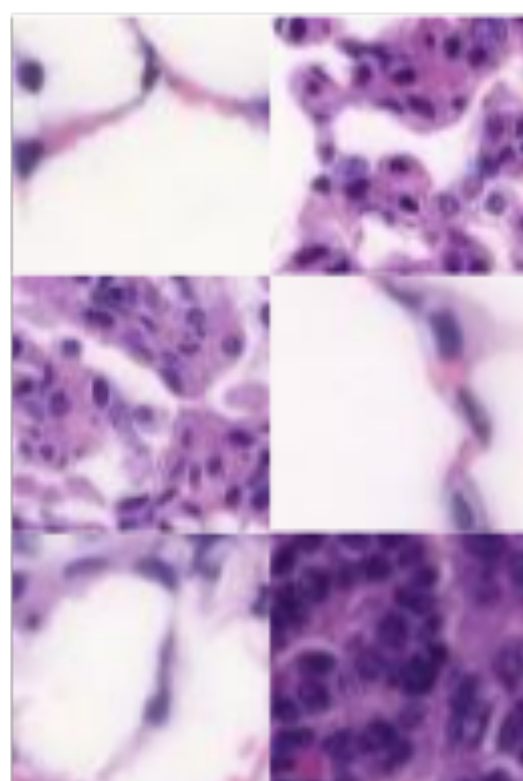
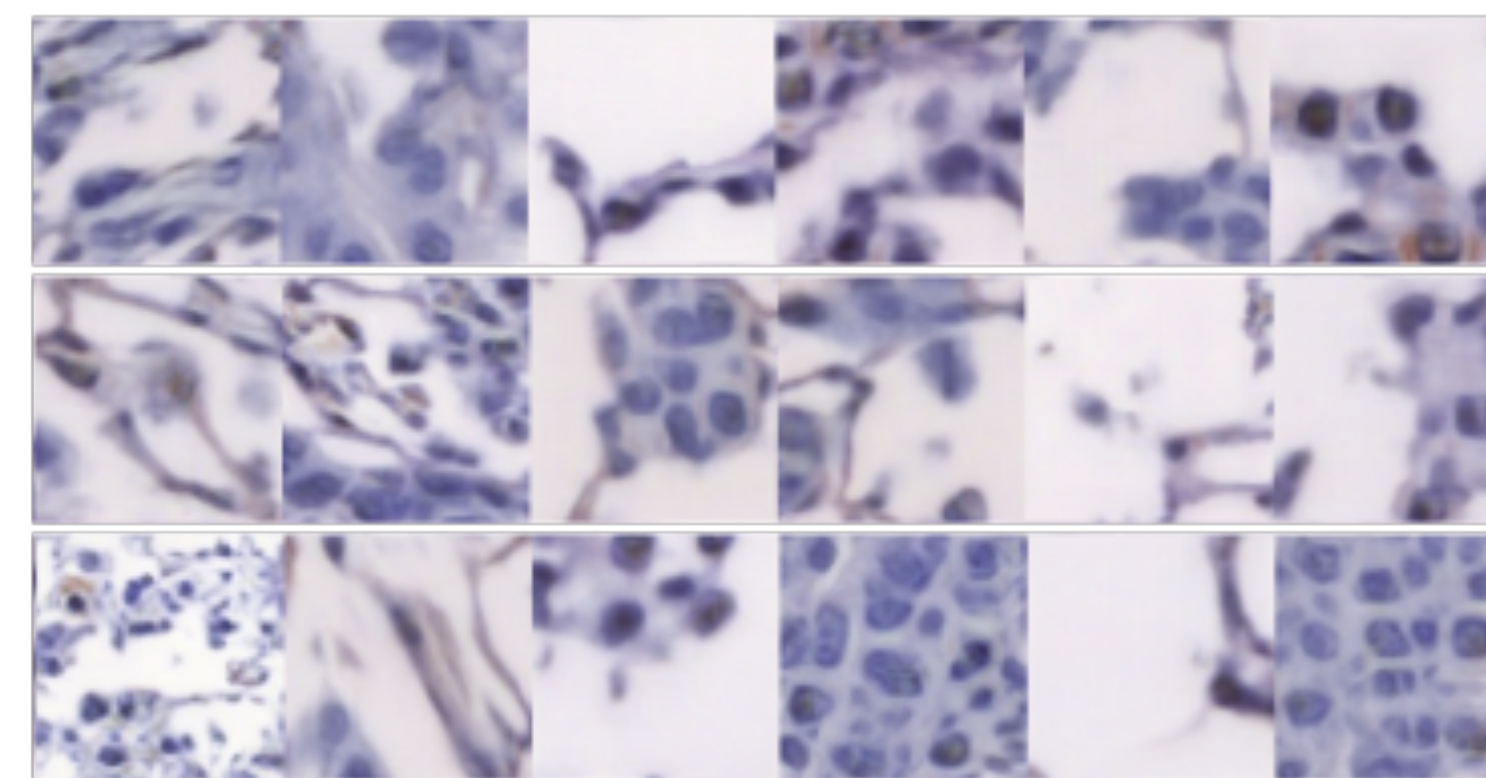
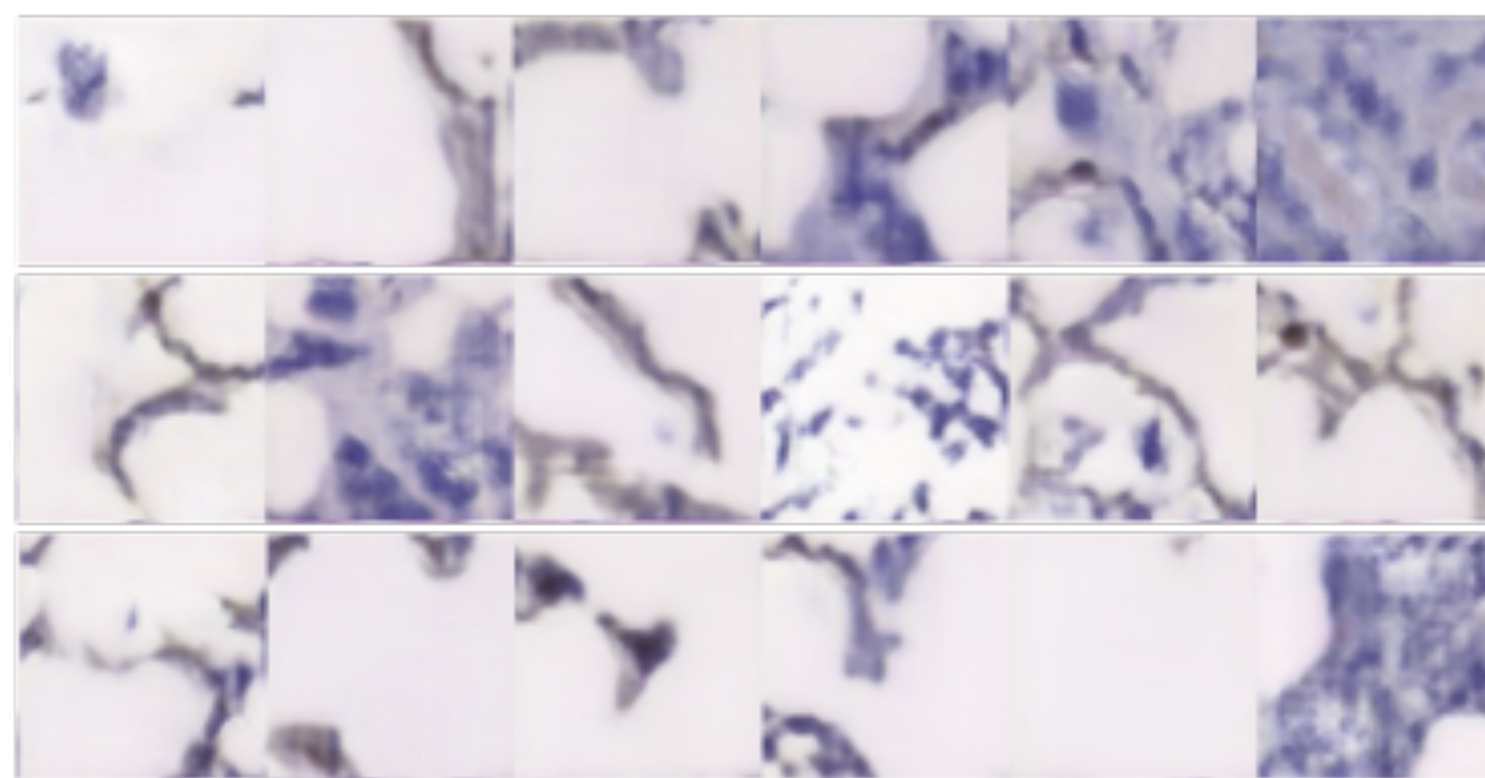
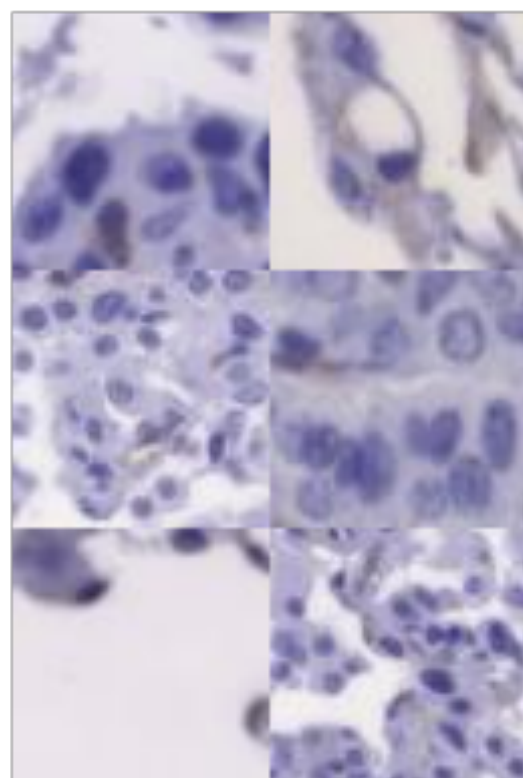


# Medical images (ANHIR)

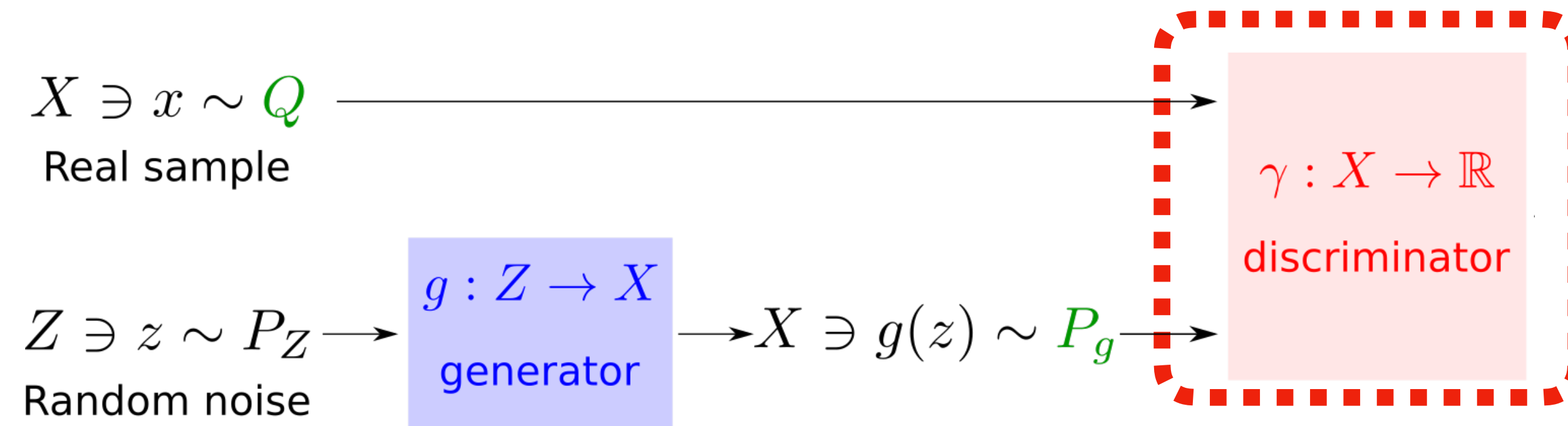
“Ignorant”  
players

“Smart” players

Real Samples

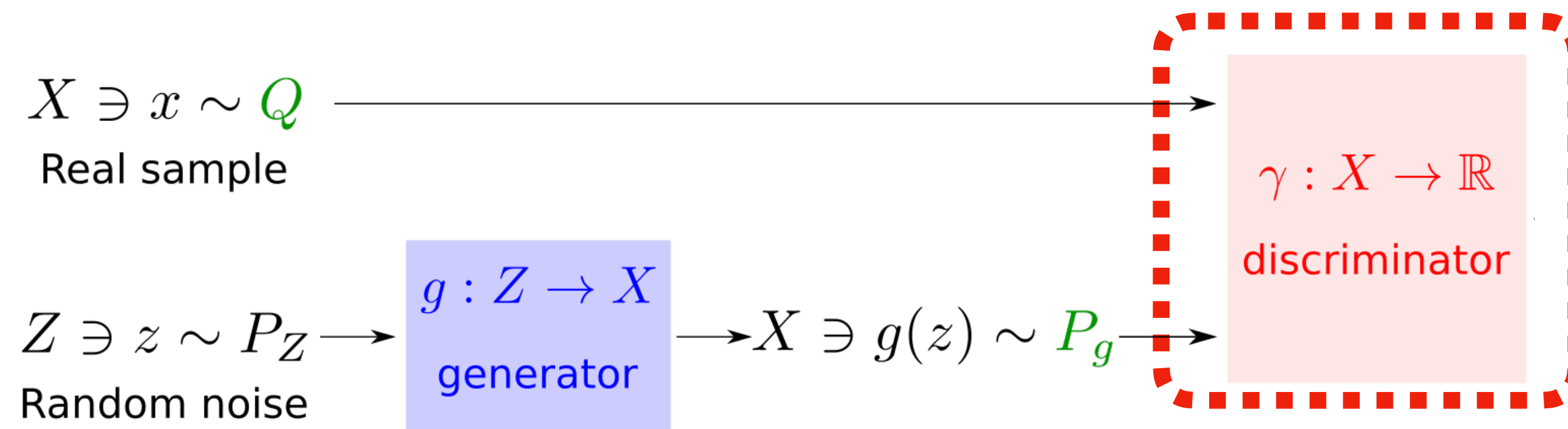


# What is the reason behind the improvement?





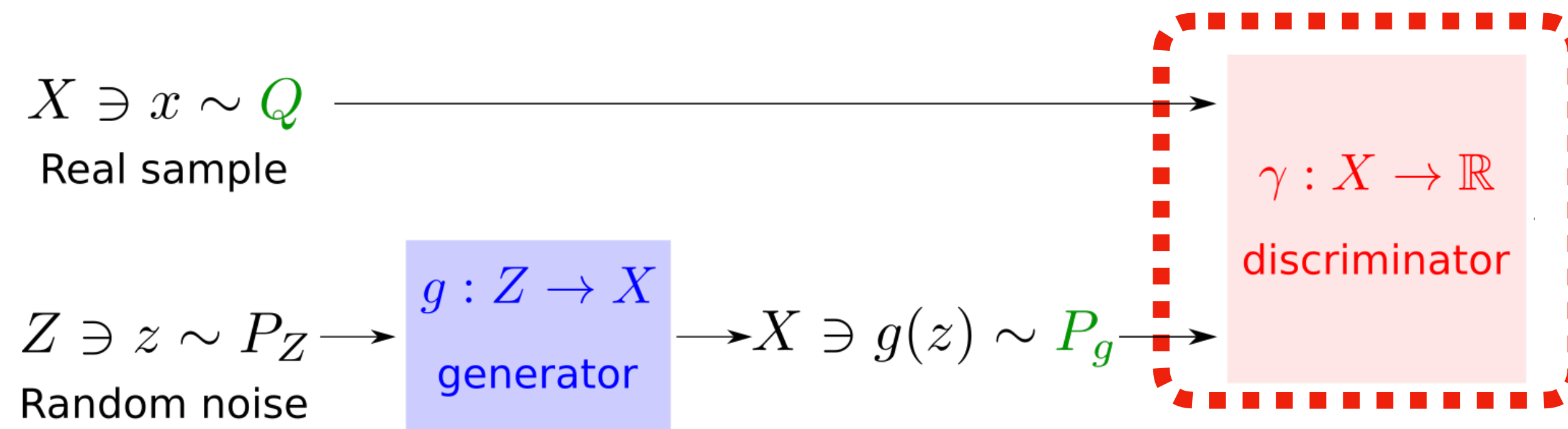
# What is the reason behind the improvement?



$$P, Q \text{ are } \Sigma\text{-invariant} \implies D^\Gamma(Q||P) = D^{\Gamma_\Sigma^{\text{inv}}}(Q||P),$$

- Reducing  $\Gamma$  to  $\Gamma_\Sigma^{\text{inv}}$  provides a better **empirical estimation** for  $D^\Gamma(Q||P)$ .

# What is the reason behind the improvement?



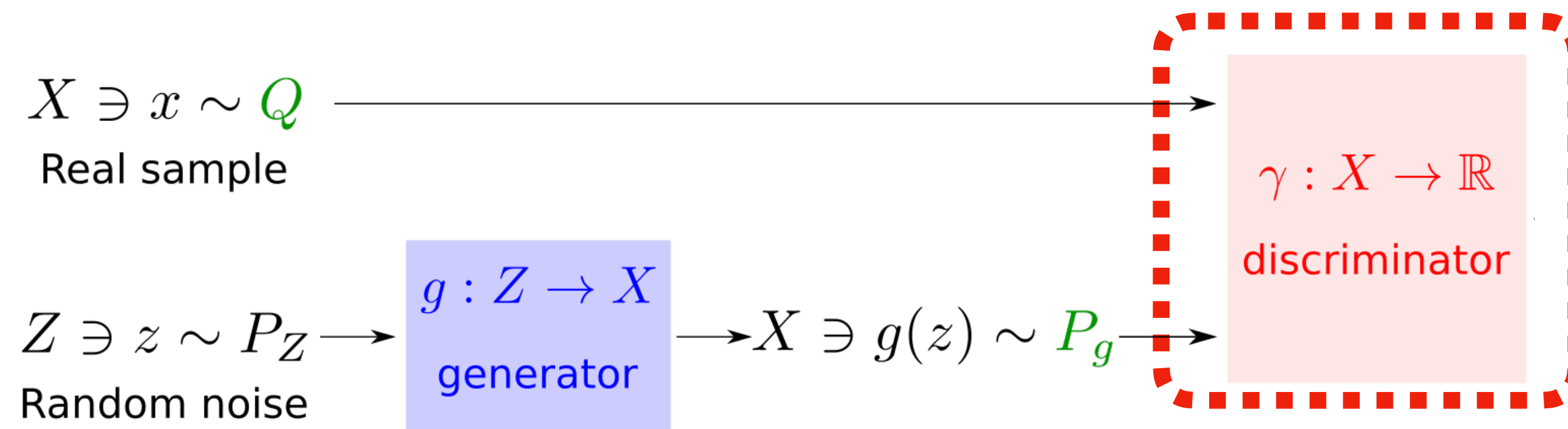
$$P, Q \text{ are } \Sigma\text{-invariant} \implies D^\Gamma(Q||P) = D^{\Gamma_\Sigma^{\text{inv}}}(Q||P),$$

- Reducing  $\Gamma$  to  $\Gamma_\Sigma^{\text{inv}}$  provides a better **empirical estimation** for  $D^\Gamma(Q||P)$ .

- $(x_1, \dots, x_m) \sim P, (y_1, \dots, y_n) \sim Q \implies$  **Empirical measures**  $P_m = \frac{1}{m} \sum_{i=1}^m \delta_{x_i}, Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$



# What is the reason behind the improvement?



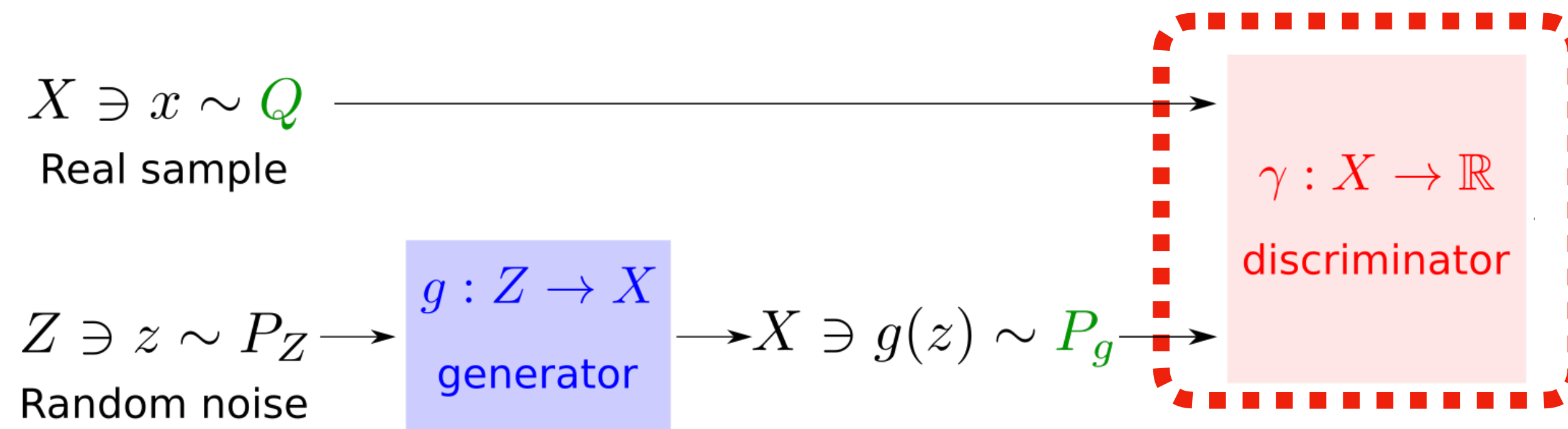
$$P, Q \text{ are } \Sigma\text{-invariant} \implies D^\Gamma(Q \| P) = D^{\Gamma_\Sigma^{\text{inv}}}(Q \| P),$$

- Reducing  $\Gamma$  to  $\Gamma_\Sigma^{\text{inv}}$  provides a better **empirical estimation** for  $D^\Gamma(Q \| P)$ .

- $(x_1, \dots, x_m) \sim P, (y_1, \dots, y_n) \sim Q \implies$  **Empirical measures**  $P_m = \frac{1}{m} \sum_{i=1}^m \delta_{x_i}, Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$

- $D^\Gamma(Q \| P) \approx \cancel{D^\Gamma(Q_n \| P_m)} = D^{\Gamma_\Sigma^{\text{inv}}}(Q_n \| P_m)$

# What is the reason behind the improvement?



$$P, Q \text{ are } \Sigma\text{-invariant} \implies D^\Gamma(Q \| P) = D^{\Gamma_\Sigma^{\text{inv}}}(Q \| P),$$

- Reducing  $\Gamma$  to  $\Gamma_\Sigma^{\text{inv}}$  provides a better **empirical estimation** for  $D^\Gamma(Q \| P)$ .

- $(x_1, \dots, x_m) \sim P, (y_1, \dots, y_n) \sim Q \implies$  **Empirical measures**  $P_m = \frac{1}{m} \sum_{i=1}^m \delta_{x_i}, Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$

- $D^\Gamma(Q \| P) \approx \cancel{D^\Gamma(Q_n \| P_m)} = D^{\Gamma_\Sigma^{\text{inv}}}(Q_n \| P_m)$

**Question:** How much **more accurate** is the new estimation?

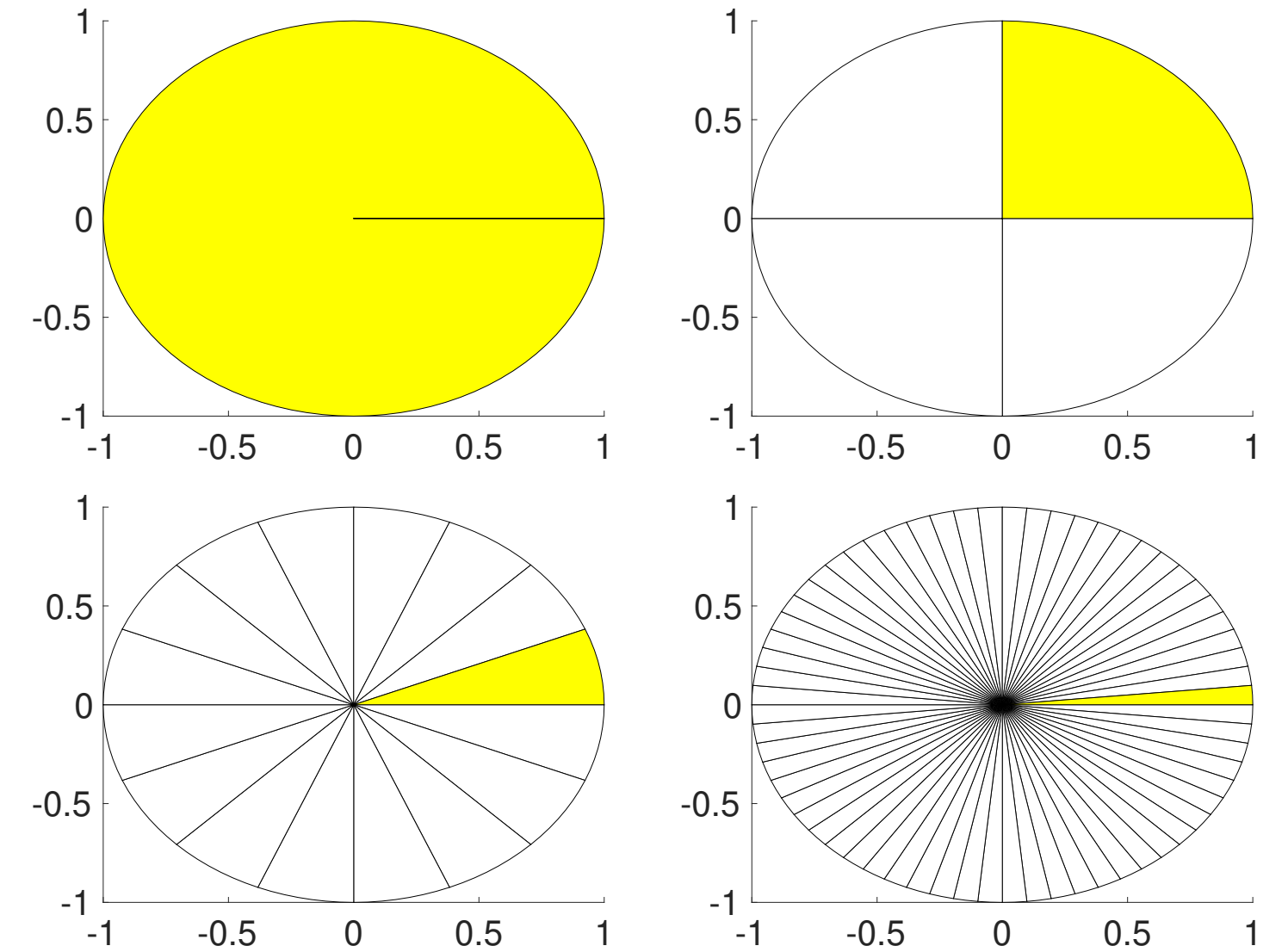


# Wasserstein-1 metric

- $W(Q, P) = \sup_{\gamma \in \Gamma} \{ \mathbb{E}_Q[\gamma] - \mathbb{E}_P[\gamma] \}$ .  $\Gamma = \text{Lip}_1(X)$
- **Estimator:**  $W^\Sigma(Q_n, P_m) = \sup_{\gamma \in \Gamma_\Sigma^{\text{inv}}} \{ \mathbb{E}_{Q_n}[\gamma] - \mathbb{E}_{P_m}[\gamma] \}$

# Wasserstein-1 metric

- $W(Q, P) = \sup_{\gamma \in \Gamma} \{ \mathbb{E}_Q[\gamma] - \mathbb{E}_P[\gamma] \}$ .  $\Gamma = \text{Lip}_1(X)$
- **Estimator:**  $W^\Sigma(Q_n, P_m) = \sup_{\gamma \in \Gamma_\Sigma^{\text{inv}}} \{ \mathbb{E}_{Q_n}[\gamma] - \mathbb{E}_{P_m}[\gamma] \}$



**Theorem** [Chen, Katsoulakis, Rey-Bellet, **Z.**, *ICML* 2023]

$X = \Sigma \times X_0$  bounded in  $\mathbb{R}^d$ , and  $P, Q \in \mathcal{P}_\Sigma(X)$  are  $\Sigma$ -invariant. With high probability,

- when  $d \geq 2$ :  $\forall s > 0$ ,  $\left| W(Q, P) - W^\Sigma(Q_n, P_m) \right| \leq C \left( \left( \frac{1}{|\Sigma| m} \right)^{\frac{1}{d+s}} + \left( \frac{1}{|\Sigma| n} \right)^{\frac{1}{d+s}} \right)$
- when  $d = 1$ :  $\left| W(Q, P) - W^\Sigma(Q_n, P_m) \right| \leq C \cdot \text{diam}(X_0) \left( \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right)$



# Maximum Mean Discrepancy (MMD)

- $\text{MMD}(Q, P) = \sup_{\gamma \in \Gamma} \{ \mathbb{E}_Q[\gamma] - \mathbb{E}_P[\gamma] \}$ .  $\Gamma$  is the unit ball in some **RKHS**  $\mathcal{H}$  with kernel  $k(x, y)$ .
- **Estimator:**  $\text{MMD}^\Sigma(Q_n, P_m) = \sup_{\gamma \in \Gamma_\Sigma^{\text{inv}}} \{ \mathbb{E}_{Q_n}[\gamma] - \mathbb{E}_{P_m}[\gamma] \}$

# Maximum Mean Discrepancy (MMD)

- $\text{MMD}(Q, P) = \sup_{\gamma \in \Gamma} \{ \mathbb{E}_Q[\gamma] - \mathbb{E}_P[\gamma] \}$ .  $\Gamma$  is the unit ball in some **RKHS**  $\mathcal{H}$  with kernel  $k(x, y)$ .
- **Estimator:**  $\text{MMD}^\Sigma(Q_n, P_m) = \sup_{\gamma \in \Gamma_\Sigma^{\text{inv}}} \{ \mathbb{E}_{Q_n}[\gamma] - \mathbb{E}_{P_m}[\gamma] \}$

**Theorem** [Chen, Katsoulakis, Rey-Bellet, **Z.**, *ICML* 2023]

$X = \Sigma \times X_0$  bounded in  $\mathbb{R}^d$ , and  $P, Q \in \mathcal{P}_\Sigma(X)$  are  $\Sigma$ -invariant. With high probability,

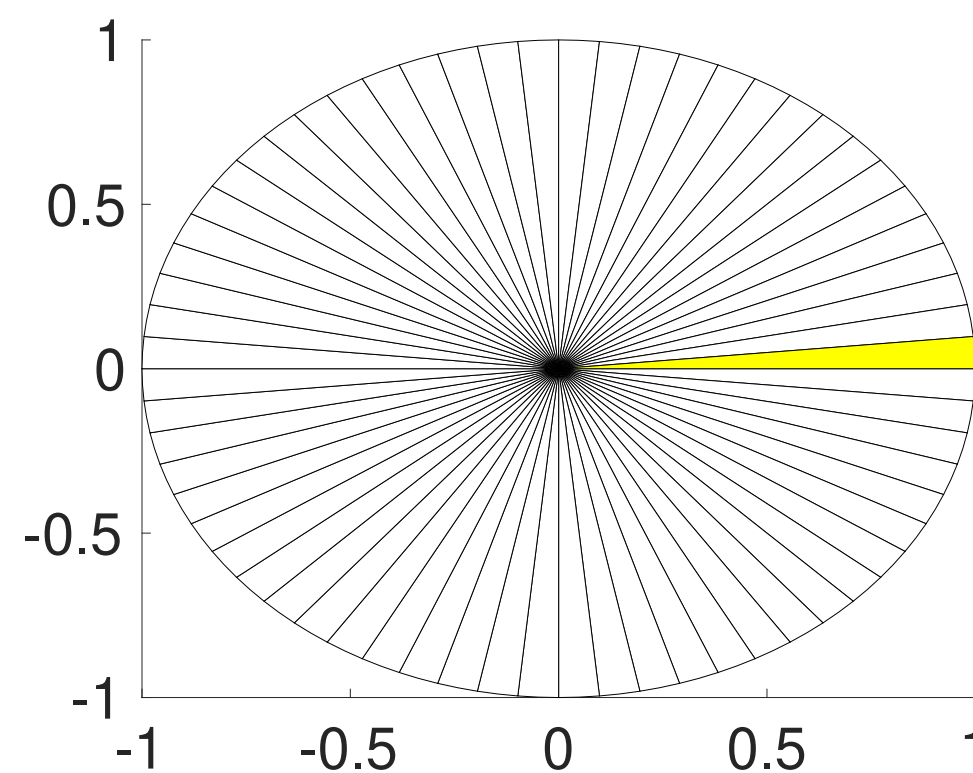
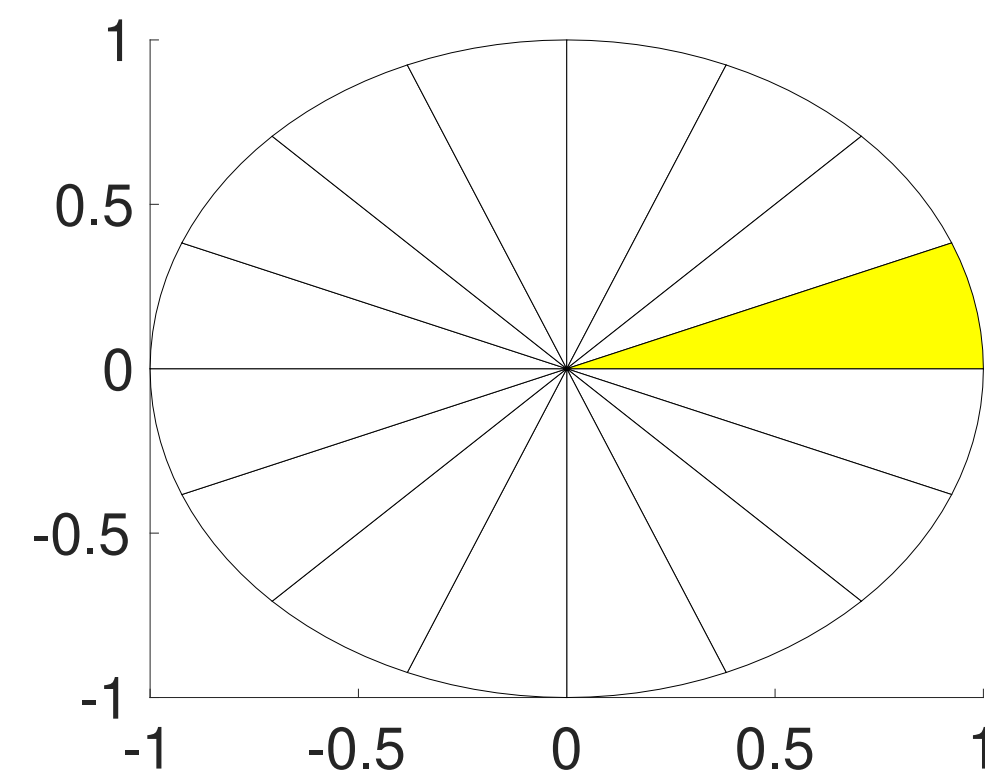
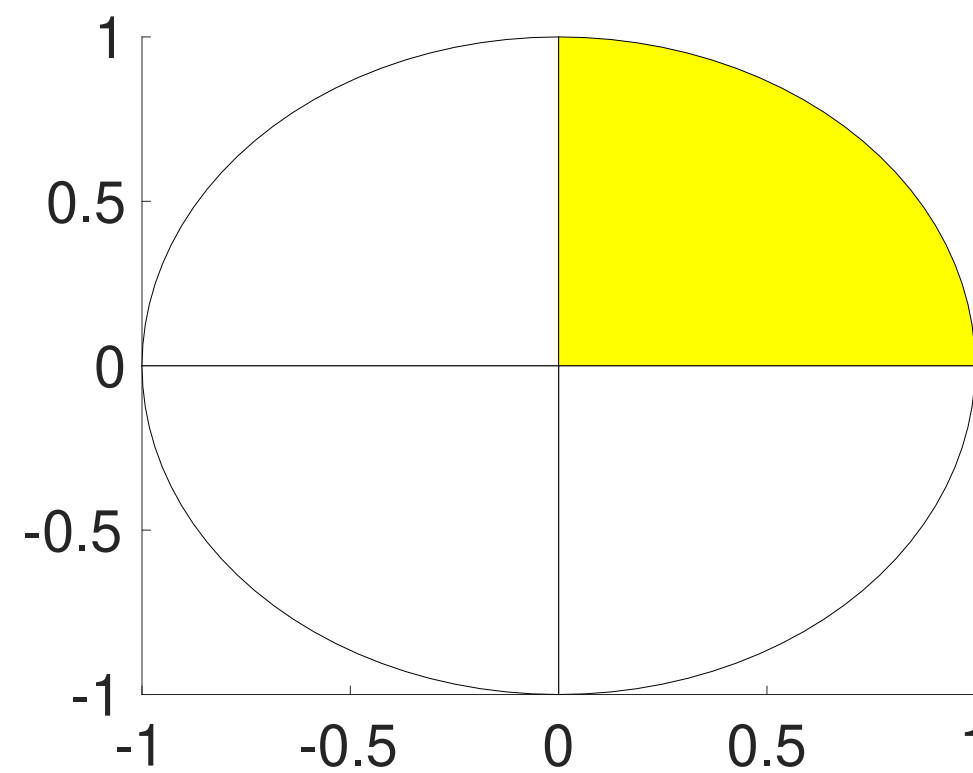
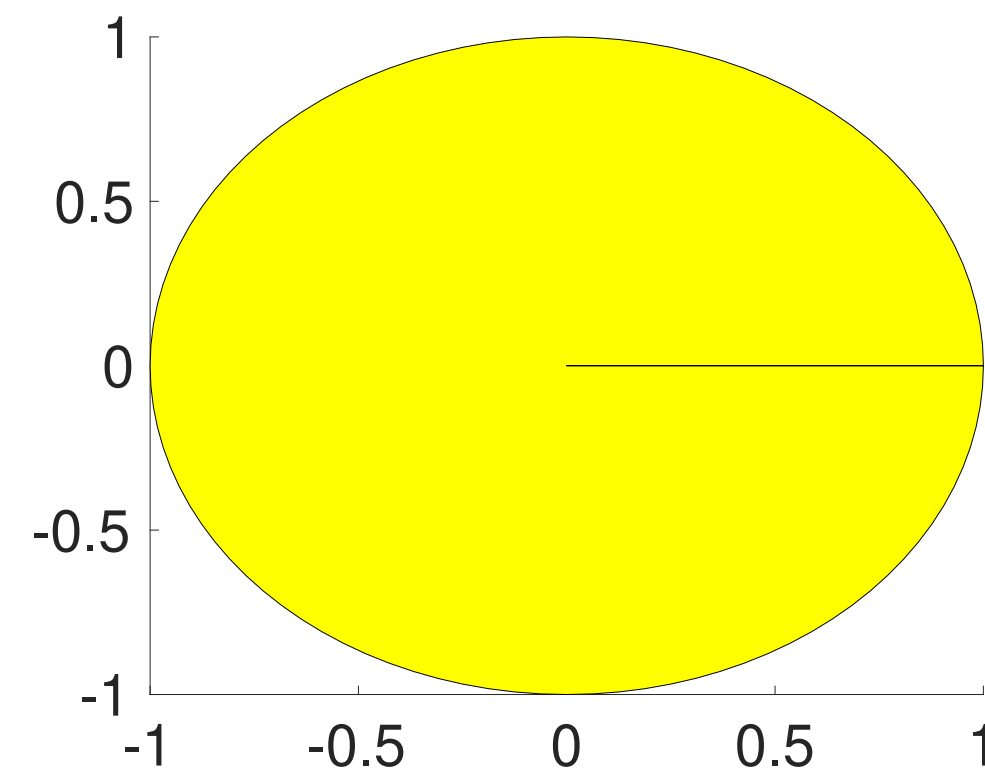
$$\left| \text{MMD}(Q, P) - \text{MMD}^\Sigma(Q_n, P_m) \right| = o \left( C_{\Sigma, k} \left( \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right) \right),$$

where  $C_{\Sigma, k} = \sqrt{a_{\Sigma, k} + \frac{1 - a_{\Sigma, k}}{|\Sigma|}}$ , and  $a_{\Sigma, k} \in (0, 1)$  depends on  $\Sigma$  and the **kernel**  $k(x, y)$ .



# Maximum Mean Discrepancy (MMD)

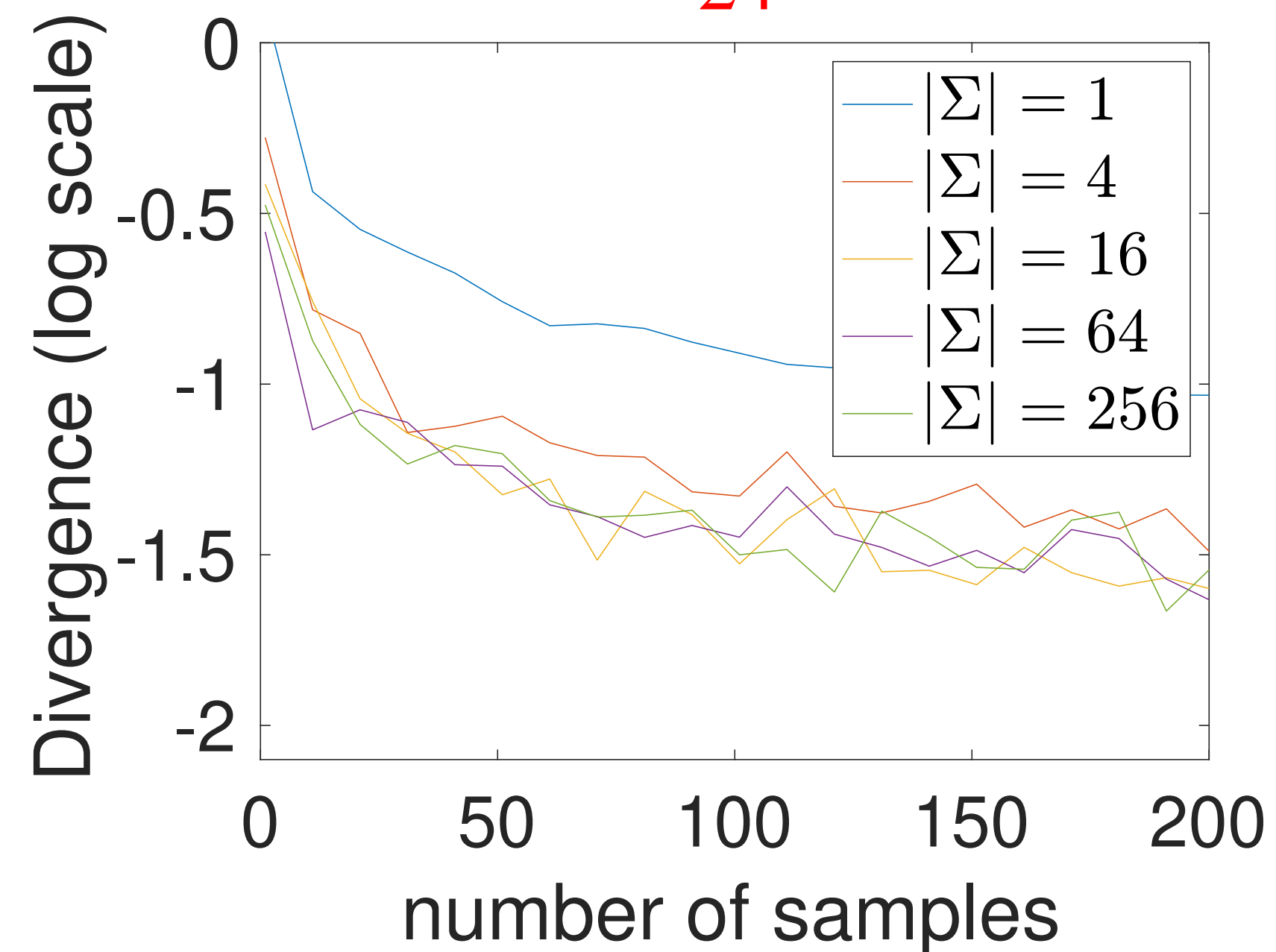
$$\left| \text{MMD}(Q, P) - \text{MMD}^{\Sigma}(Q_n, P_m) \right| = O \left( C_{\Sigma, k} \left( \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right) \right), \quad C_{\Sigma, k} = \sqrt{a_{\Sigma, k} + \frac{1 - a_{\Sigma, k}}{|\Sigma|}}$$



# Maximum Mean Discrepancy (MMD)

$$\left| \text{MMD}(Q, P) - \text{MMD}^{\Sigma}(Q_n, P_m) \right| = O \left( C_{\Sigma, k} \left( \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right) \right), \quad C_{\Sigma, k} = \sqrt{a_{\Sigma, k} + \frac{1 - a_{\Sigma, k}}{|\Sigma|}}$$

$$s = \frac{2\pi}{24}$$

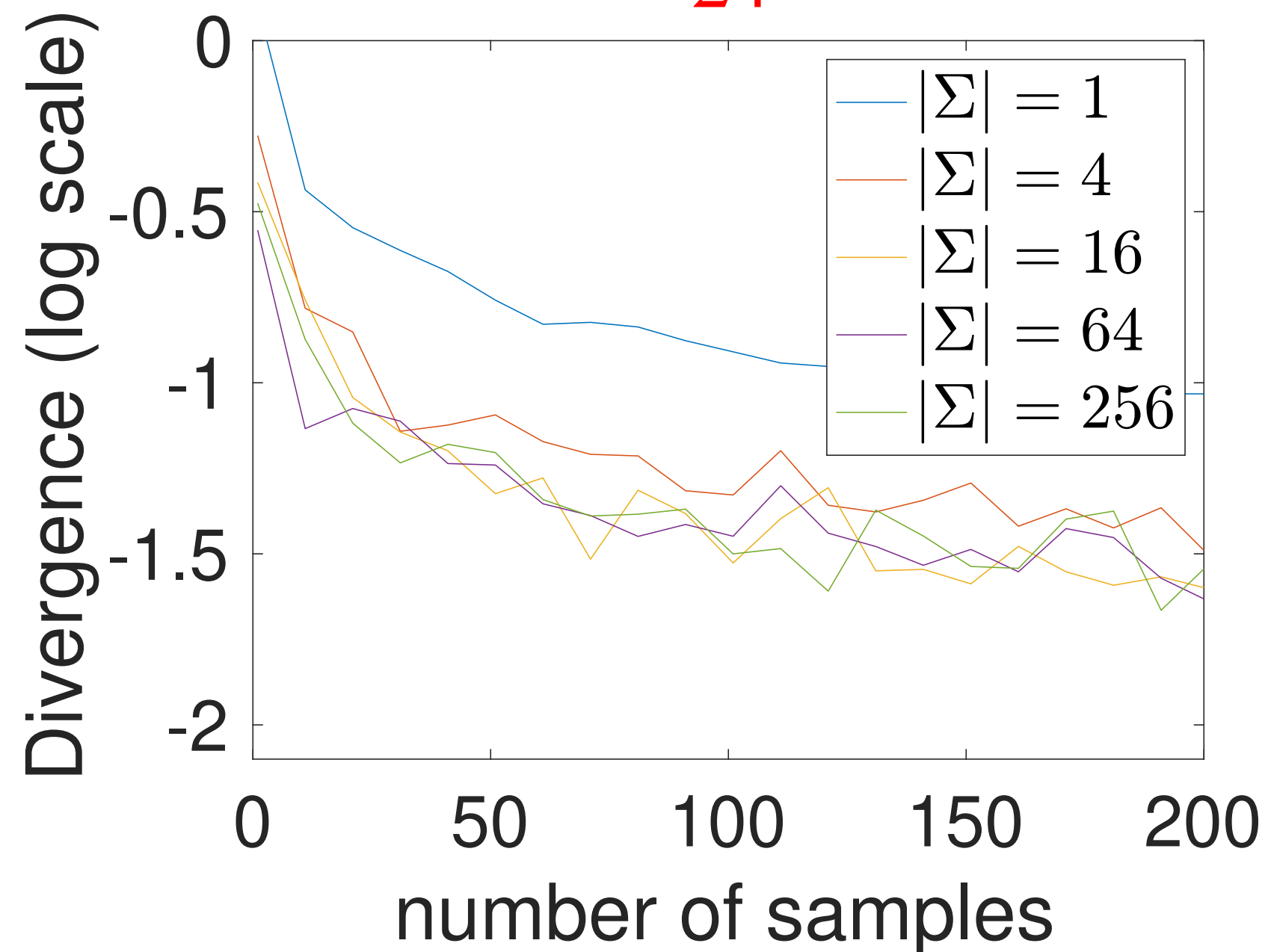




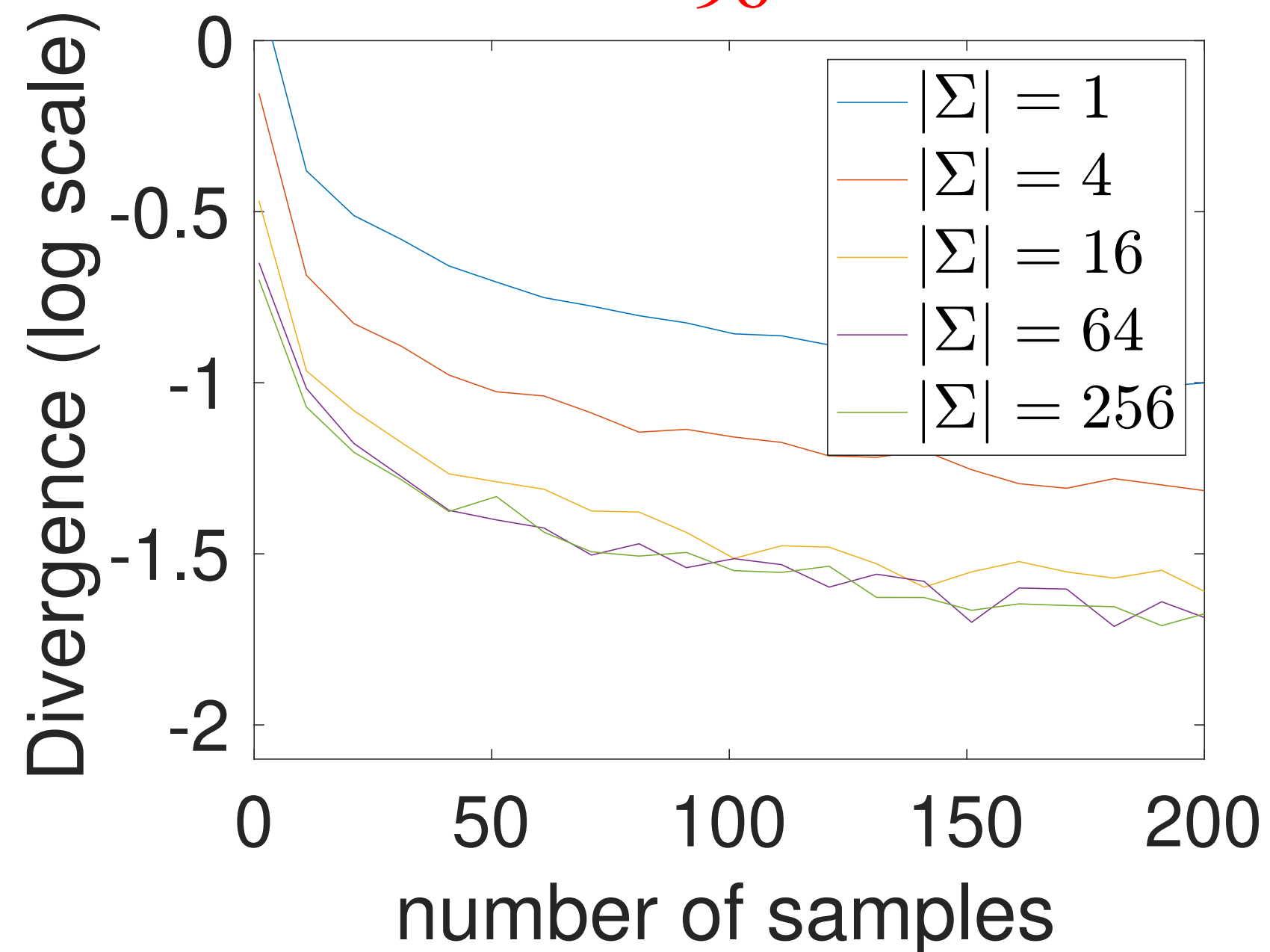
# Maximum Mean Discrepancy (MMD)

$$\left| \text{MMD}(Q, P) - \text{MMD}^{\Sigma}(Q_n, P_m) \right| = O \left( C_{\Sigma, k} \left( \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right) \right), \quad C_{\Sigma, k} = \sqrt{a_{\Sigma, k} + \frac{1 - a_{\Sigma, k}}{|\Sigma|}}$$

$$s = \frac{2\pi}{24}$$



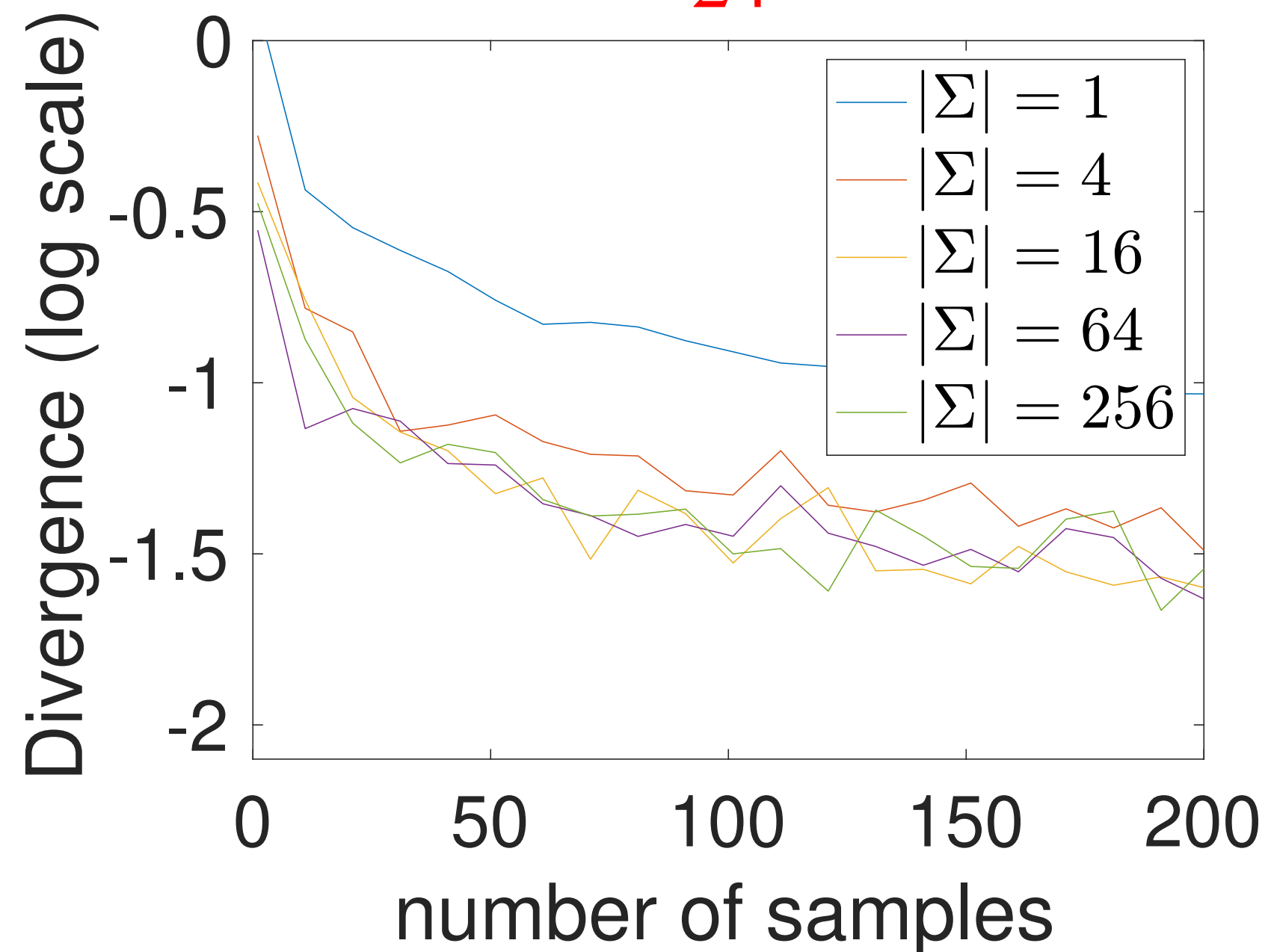
$$s = \frac{2\pi}{96}$$



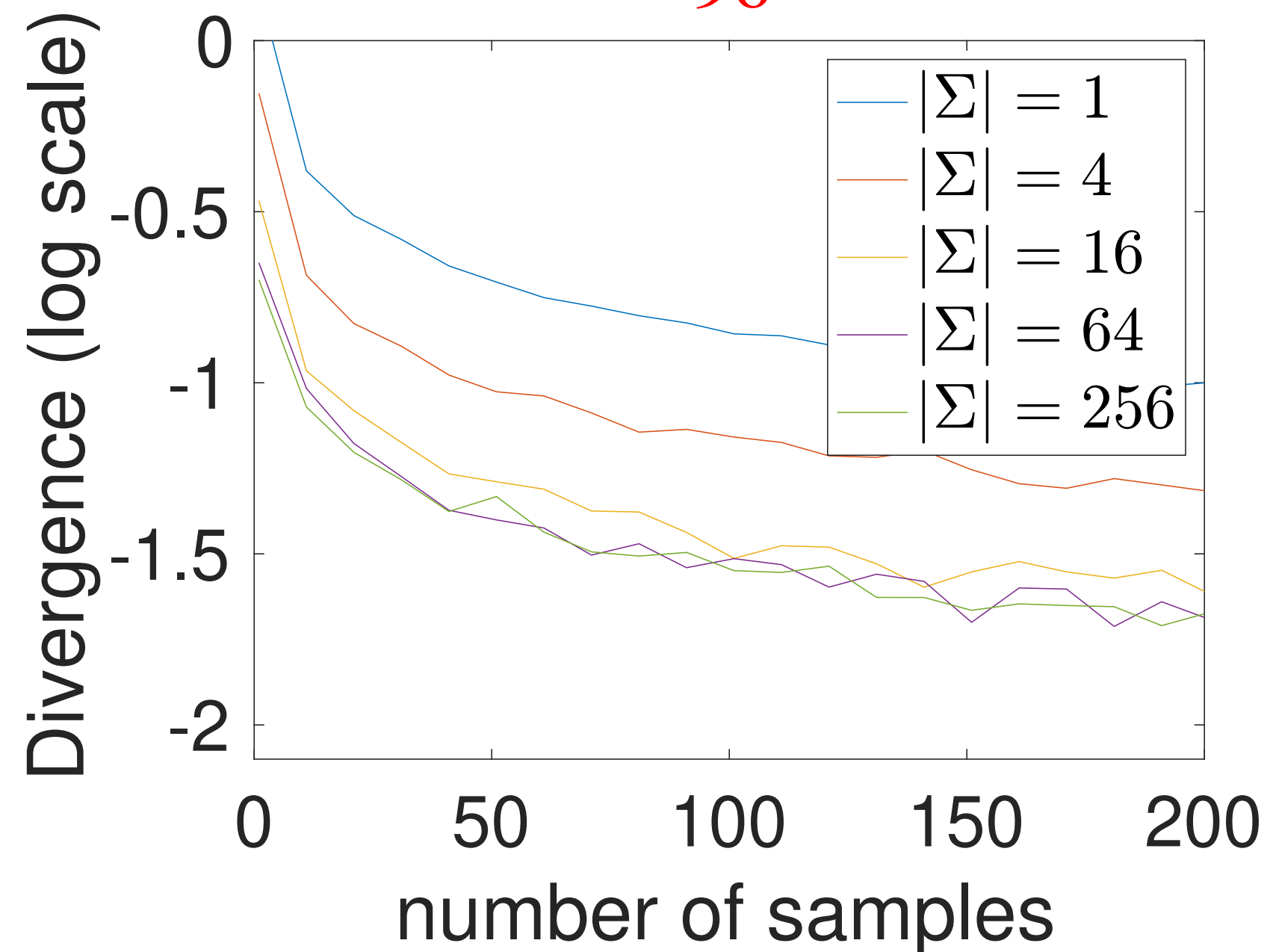
# Maximum Mean Discrepancy (MMD)

$$\left| \text{MMD}(Q, P) - \text{MMD}^{\Sigma}(Q_n, P_m) \right| = O \left( C_{\Sigma, k} \left( \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right) \right), \quad C_{\Sigma, k} = \sqrt{a_{\Sigma, k} + \frac{1 - a_{\Sigma, k}}{|\Sigma|}}$$

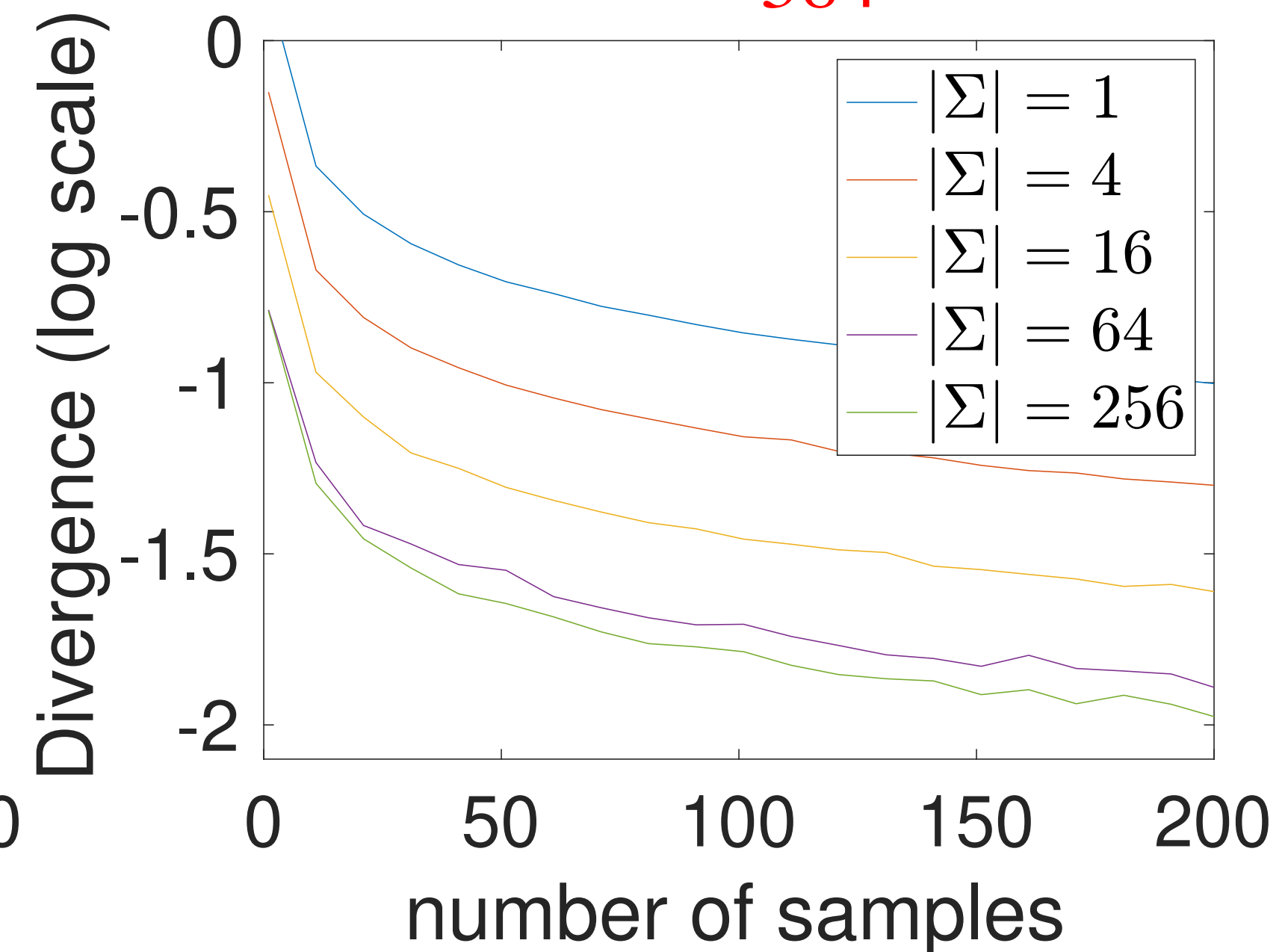
$$s = \frac{2\pi}{24}$$



$$s = \frac{2\pi}{96}$$

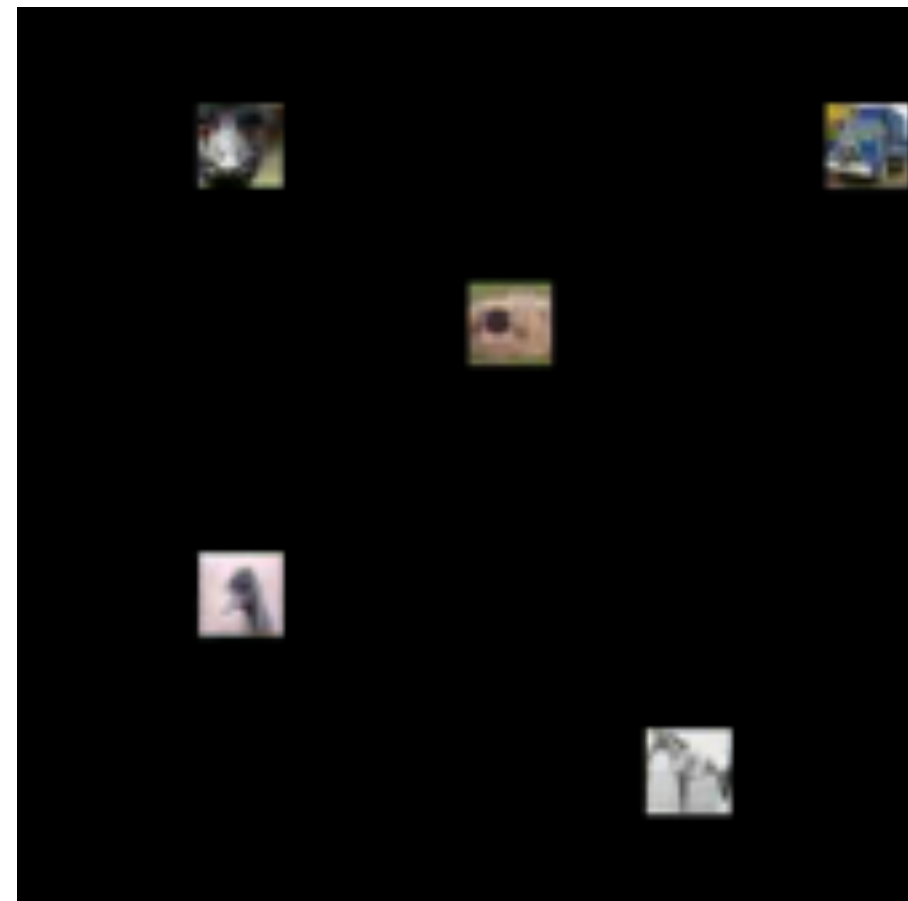
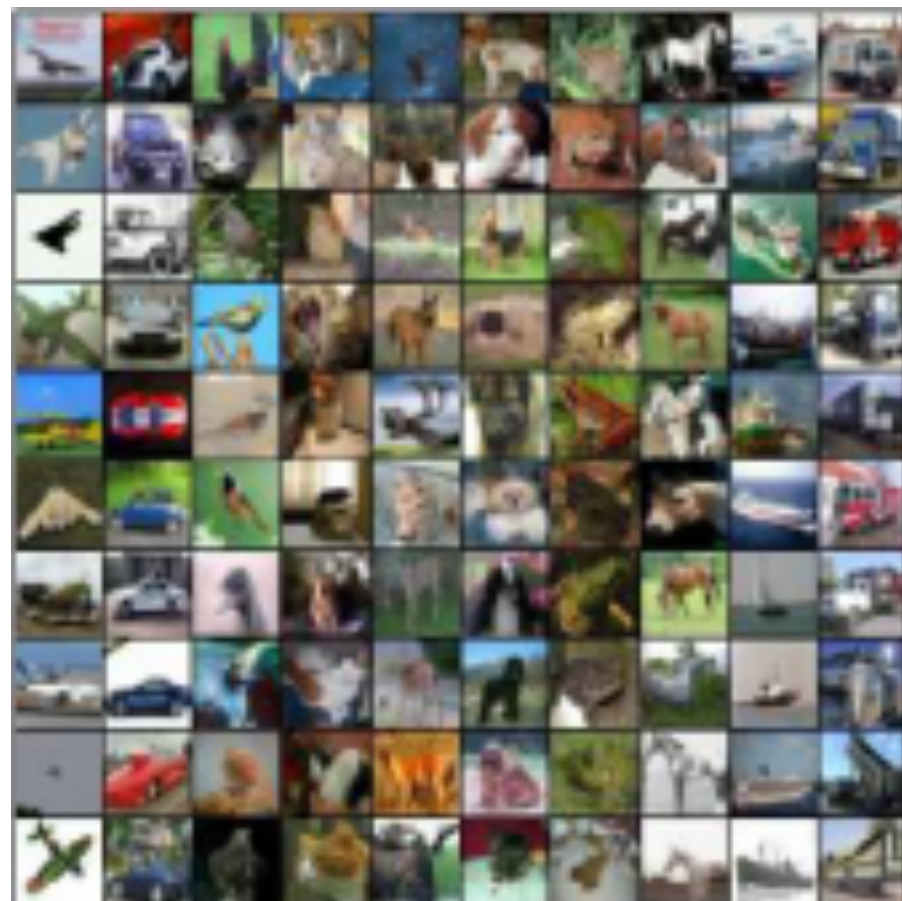


$$s = \frac{2\pi}{384}$$



# Missing pieces

- **Exact quantification of the improvement**
  - **Sample complexity** and **error bound**.



- **Does it converge? To what solution?**
  - **Training dynamics** of equivariant models

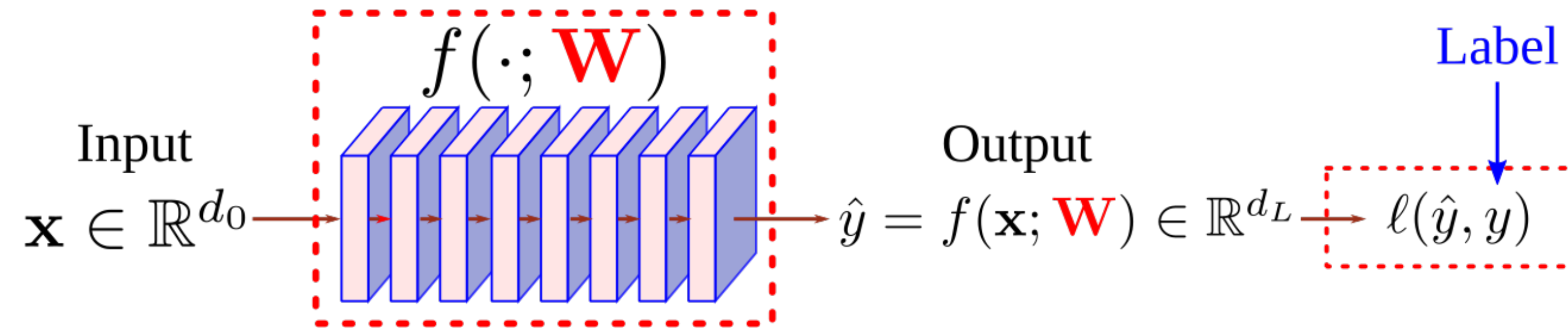
$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; S) = \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{W}), y_i)$$



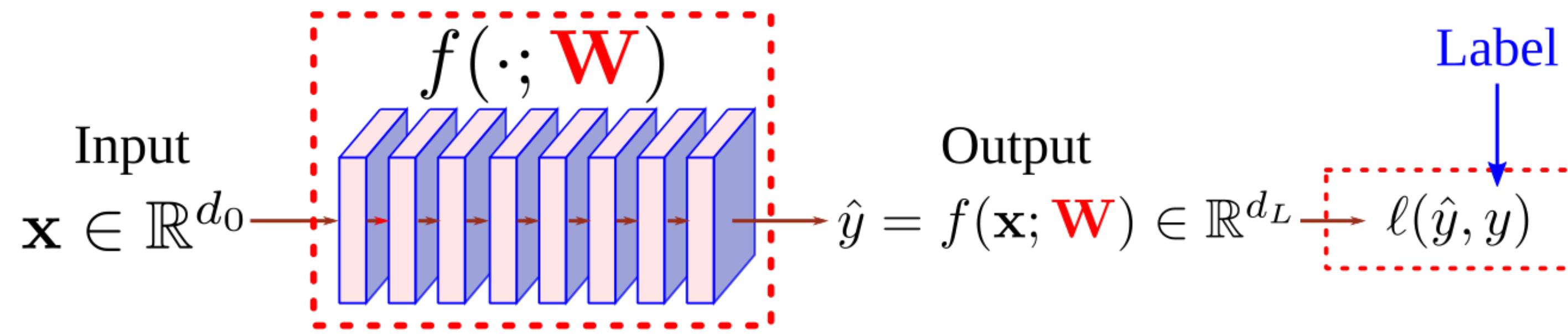
# Implicit bias of linear equivariant networks

- Z. Chen and **W. Zhu**. “On the implicit bias of linear equivariant steerable networks”. *NeurIPS* (2023)
- Inspired by [Lawrence et al., *ICML* 2022]

# Optimization (training) of G-CNN



# Optimization (training) of G-CNN

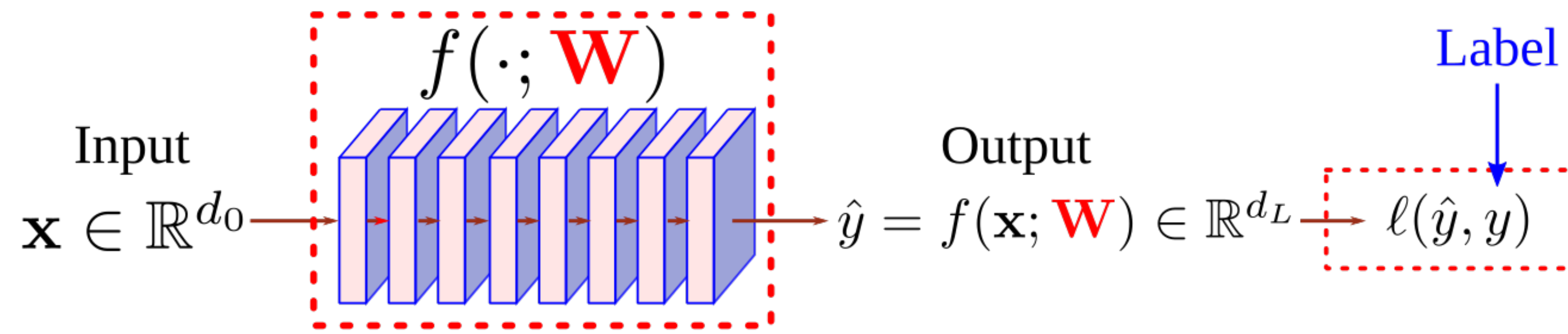


- Training a DNN on the a (labeled) data set  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ :

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; S) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{W}), y_i)$$



# Optimization (training) of G-CNN

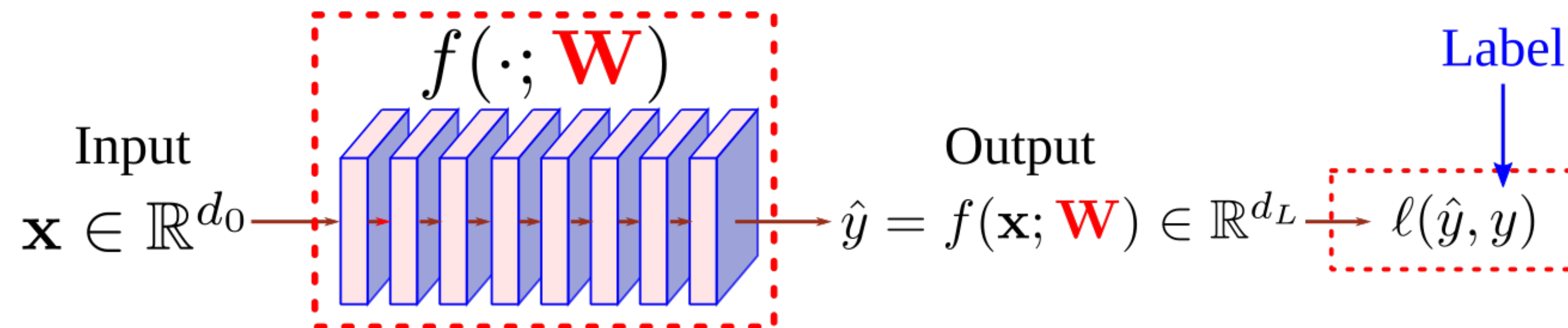


- Training a DNN on the a (labeled) data set  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ :

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; S) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{W}), y_i)$$

- Design  $f(\cdot; \mathbf{W})$  to respect group symmetry — **explicit regularization**.

# Optimization (training) of G-CNN



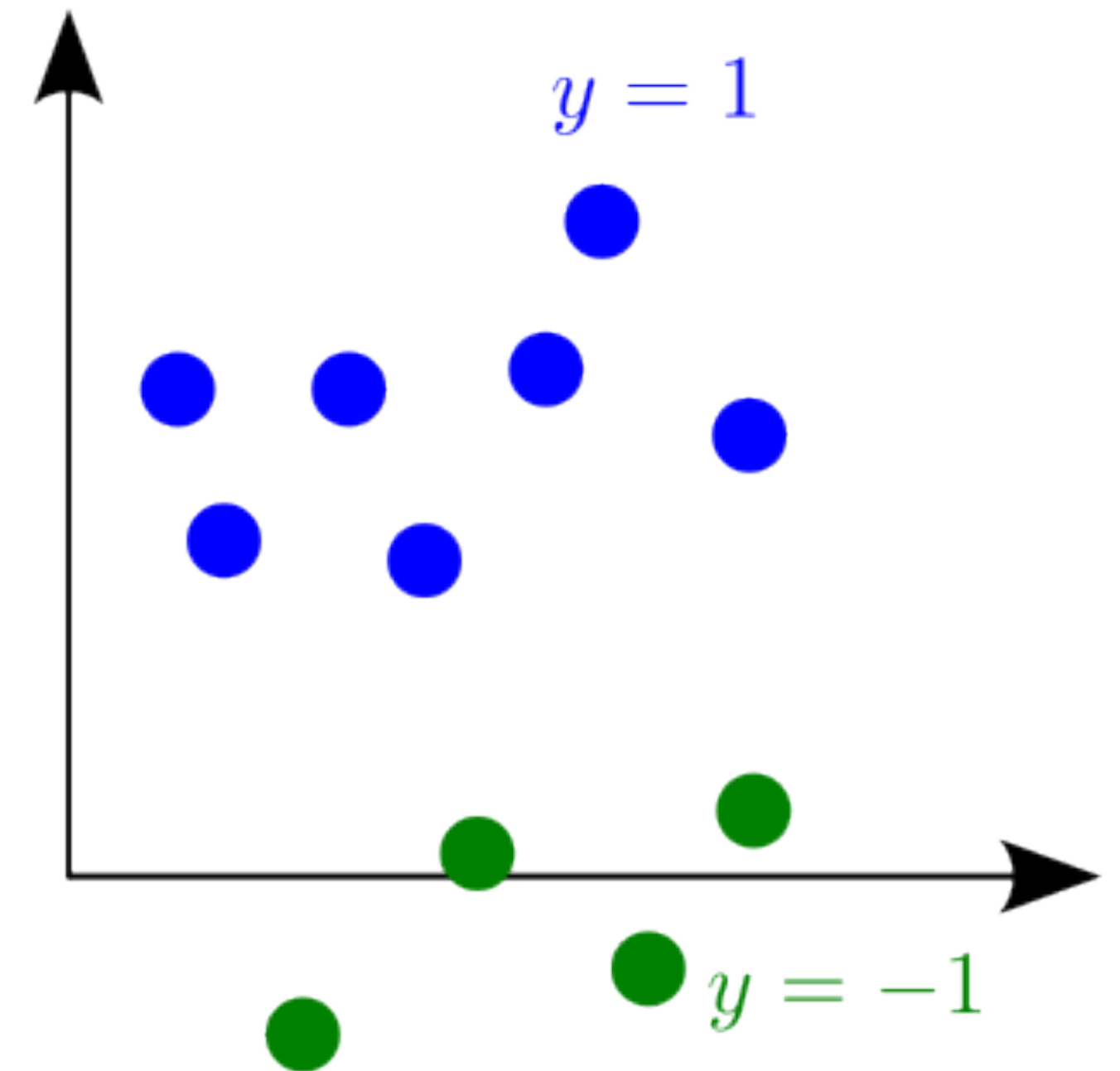
- Training a DNN on the a (labeled) data set  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ :

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; S) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{W}), y_i)$$

- Design  $f(\cdot; \mathbf{W})$  to respect group symmetry — **explicit regularization**.
- **Question:** when trained with gradient-based methods,
  - which solution does it converge to?
  - is it really better than non-equivariant models?

# Implicit regularization of training algorithms

Linear binary classification

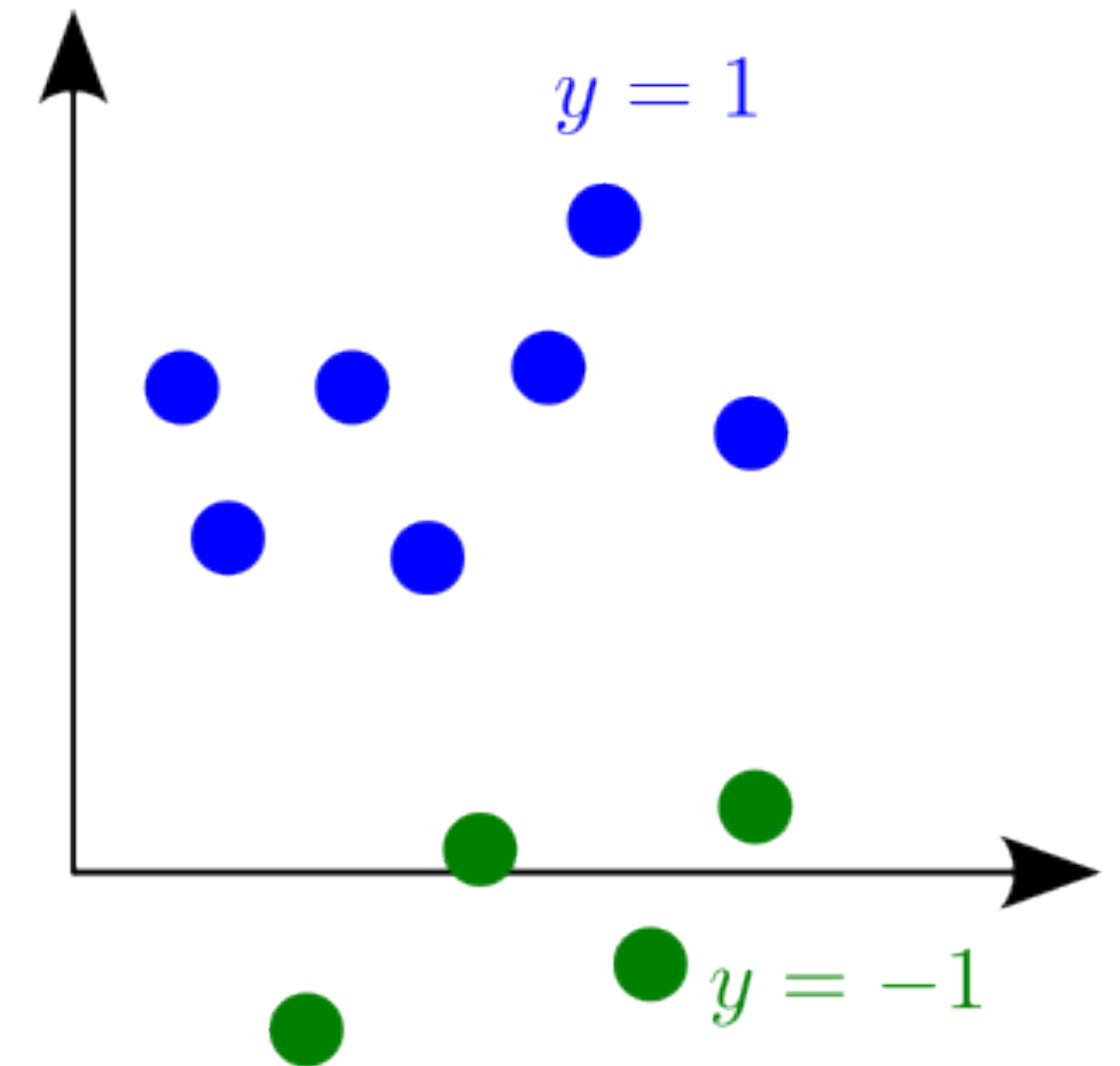




# Implicit regularization of training algorithms

- $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$ ,  $\mathbf{x}_i \in \mathbb{R}^{d_0}$  and  $y_i \in \{\pm 1\}$ .

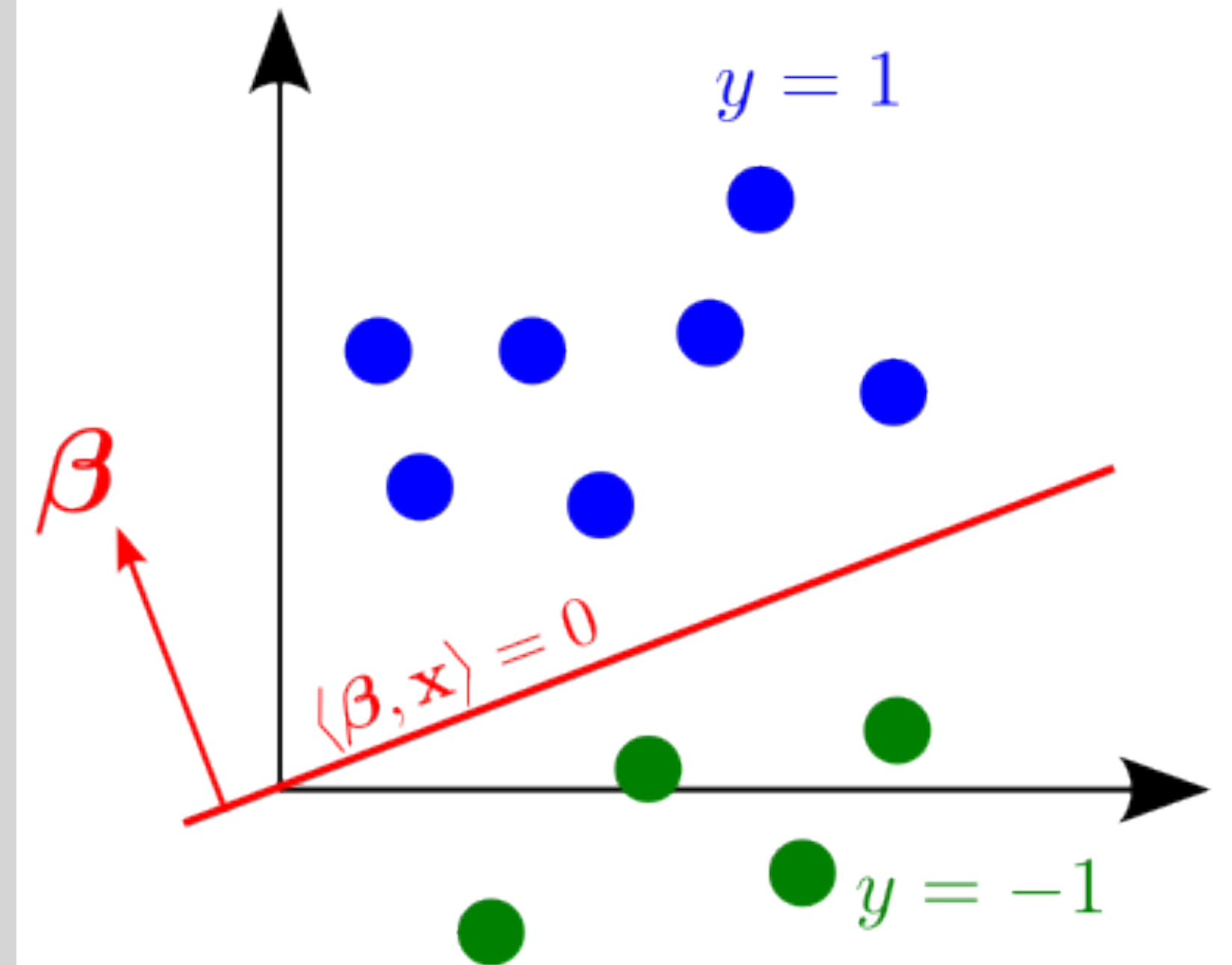
Linear binary classification



# Implicit regularization of training algorithms

- $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$ ,  $\mathbf{x}_i \in \mathbb{R}^{d_0}$  and  $y_i \in \{\pm 1\}$ .
- **Linearly separable:**  $\exists \beta^* \in \mathbb{R}^{d_0}$ , s.t  $y_i \langle \mathbf{x}_i, \beta^* \rangle \geq 1, \forall i \in [n]$ .

## Linear binary classification



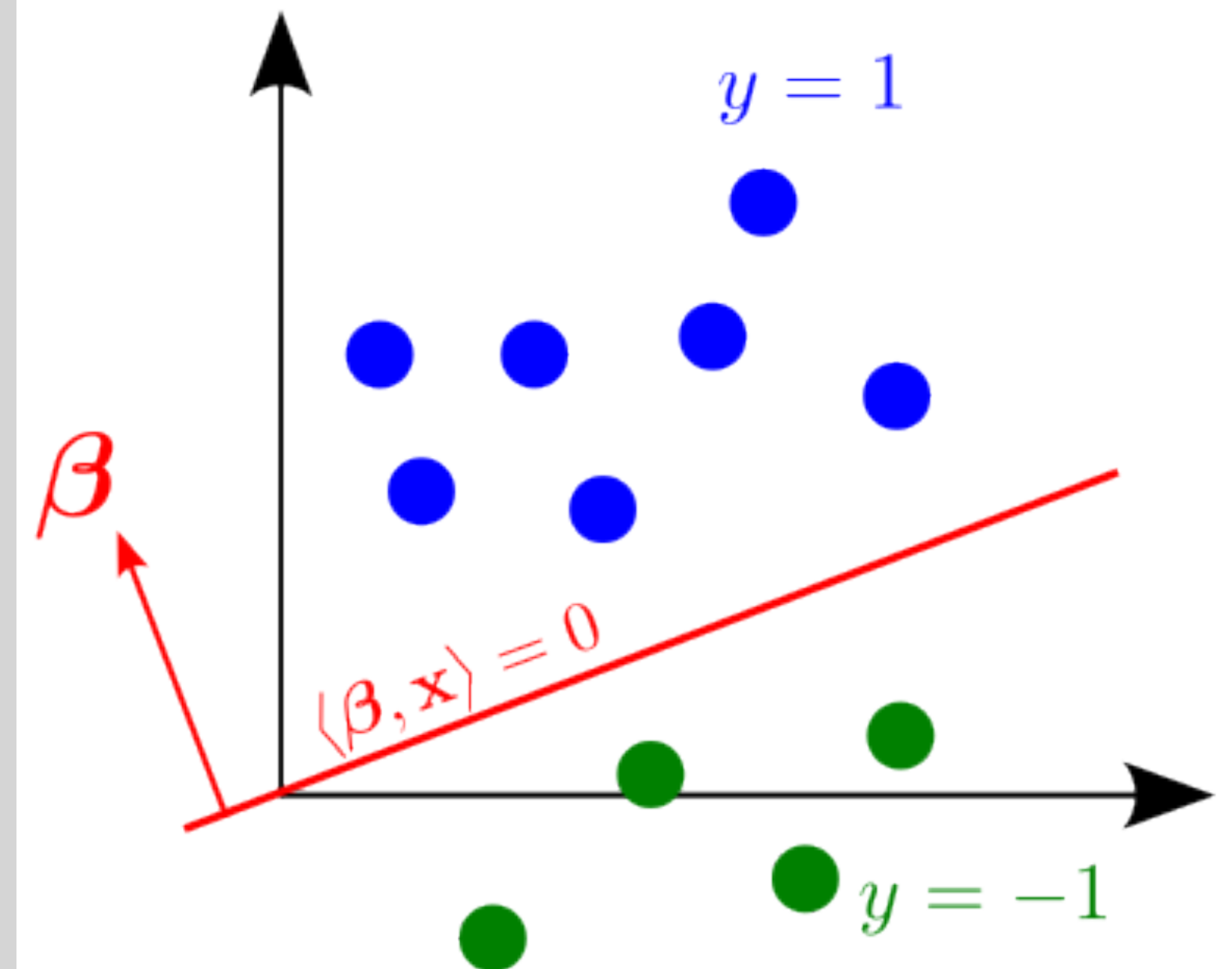
# Implicit regularization of training algorithms

- $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$ ,  $\mathbf{x}_i \in \mathbb{R}^{d_0}$  and  $y_i \in \{\pm 1\}$ .
- **Linearly separable**:  $\exists \beta^* \in \mathbb{R}^{d_0}$ , s.t  $y_i \langle \mathbf{x}_i, \beta^* \rangle \geq 1, \forall i \in [n]$ .
- Use **linear fully-connected (fc)** network to parameterize  $\langle \mathbf{x}, \beta^* \rangle$

$$f_{\text{fc}}(\mathbf{x}; \mathbf{W}) = \mathbf{w}_L^\top \mathbf{w}_{L-1}^\top \cdots \mathbf{w}_1^\top \mathbf{x} = \langle \mathbf{x}, \mathcal{P}_{\text{fc}}(\mathbf{W}) \rangle \stackrel{?}{\approx} \langle \mathbf{x}, \beta^* \rangle$$

$$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \cdots, \mathbf{w}_L], \quad \mathcal{P}_{\text{fc}}(\mathbf{W}) = \mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_L$$

## Linear binary classification





# Implicit regularization of training algorithms

- $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$ ,  $\mathbf{x}_i \in \mathbb{R}^{d_0}$  and  $y_i \in \{\pm 1\}$ .
- **Linearly separable**:  $\exists \beta^* \in \mathbb{R}^{d_0}$ , s.t  $y_i \langle \mathbf{x}_i, \beta^* \rangle \geq 1, \forall i \in [n]$ .
- Use **linear fully-connected (fc)** network to parameterize  $\langle \mathbf{x}, \beta^* \rangle$

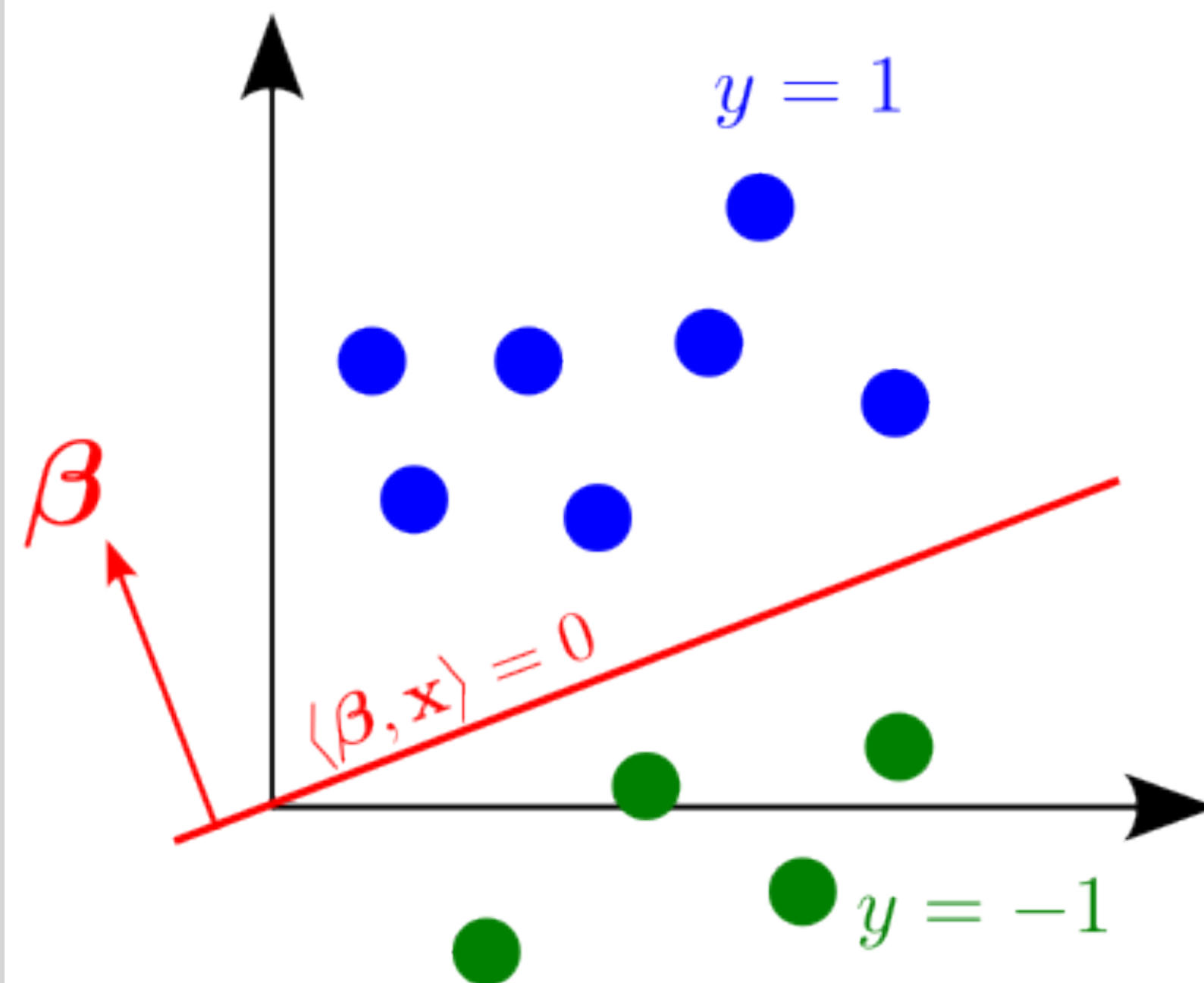
$$f_{\text{fc}}(\mathbf{x}; \mathbf{W}) = \mathbf{w}_L^\top \mathbf{w}_{L-1}^\top \cdots \mathbf{w}_1^\top \mathbf{x} = \langle \mathbf{x}, \mathcal{P}_{\text{fc}}(\mathbf{W}) \rangle \stackrel{?}{\approx} \langle \mathbf{x}, \beta^* \rangle$$

$$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \cdots, \mathbf{w}_L], \quad \mathcal{P}_{\text{fc}}(\mathbf{W}) = \mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_L$$

- Regression based on  $\ell_{\text{exp}}(\hat{y}, y) = \exp(-\hat{y}y)$

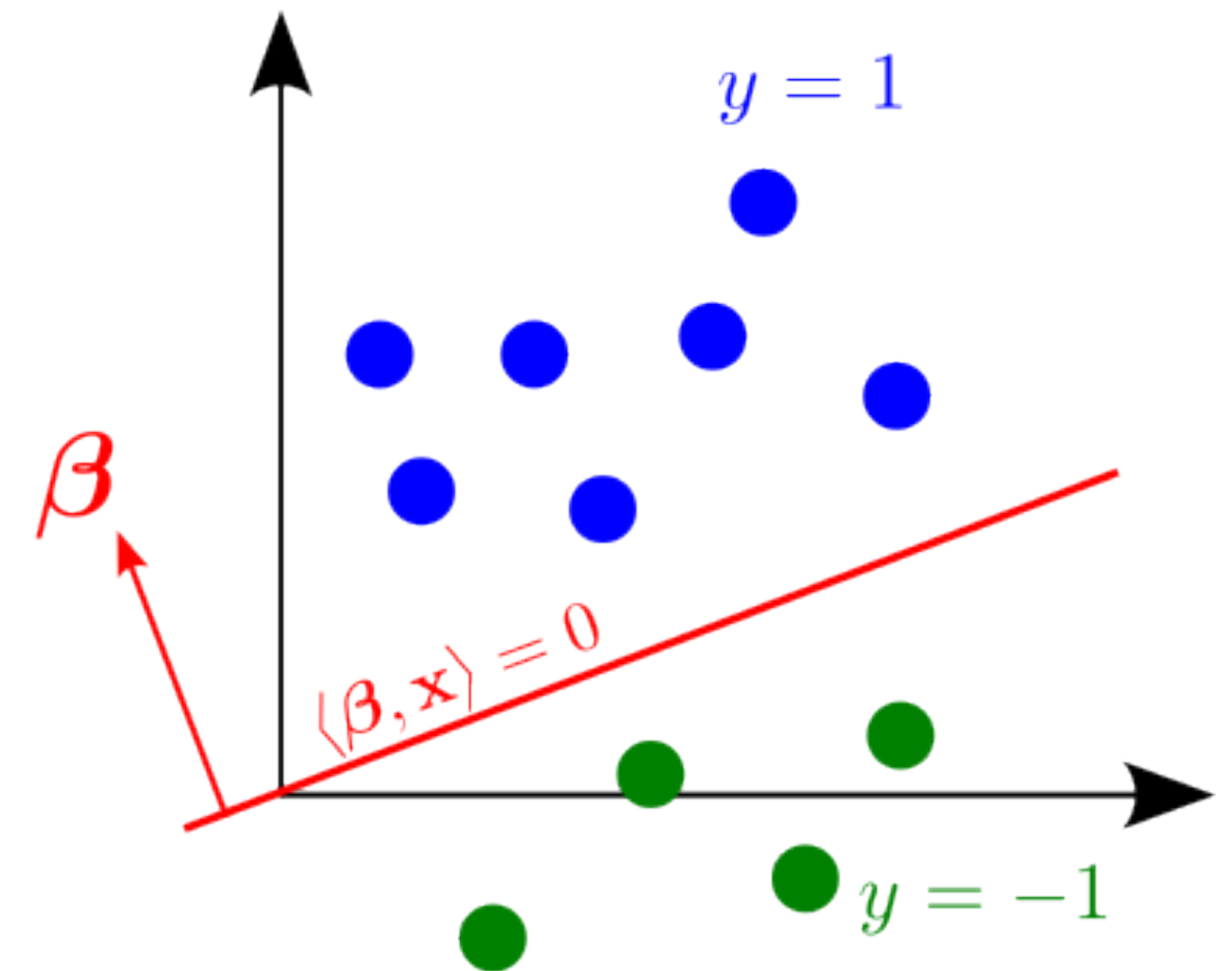
$$\min_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S) = \sum_{i=1}^n \ell_{\text{exp}} \left( \langle \mathbf{x}_i, \mathcal{P}_{\text{fc}}(\mathbf{W}) \rangle, y_i \right)$$

## Linear binary classification



# Implicit regularization of training algorithms

Linear binary classification



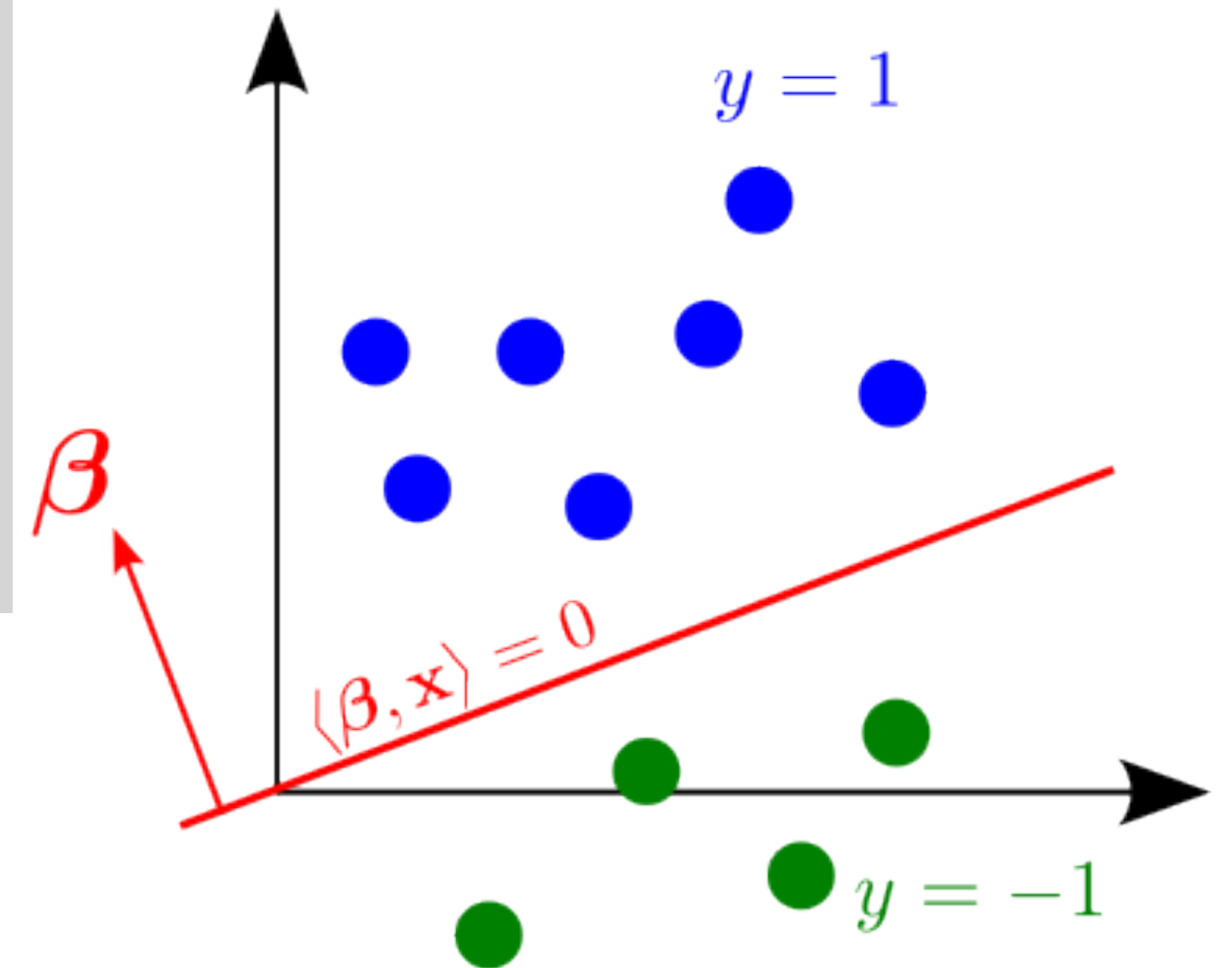
# Implicit regularization of training algorithms

$$\min_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S) = \sum_{i=1}^n \ell_{\text{exp}} \left( \langle \mathbf{x}_i, \mathcal{P}_{\text{fc}}(\mathbf{W}) \rangle, y_i \right)$$

- Trained under **gradient flow (GF)**:

$$\frac{d\mathbf{W}}{dt} = - \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S)$$

Linear binary classification





# Implicit regularization of training algorithms

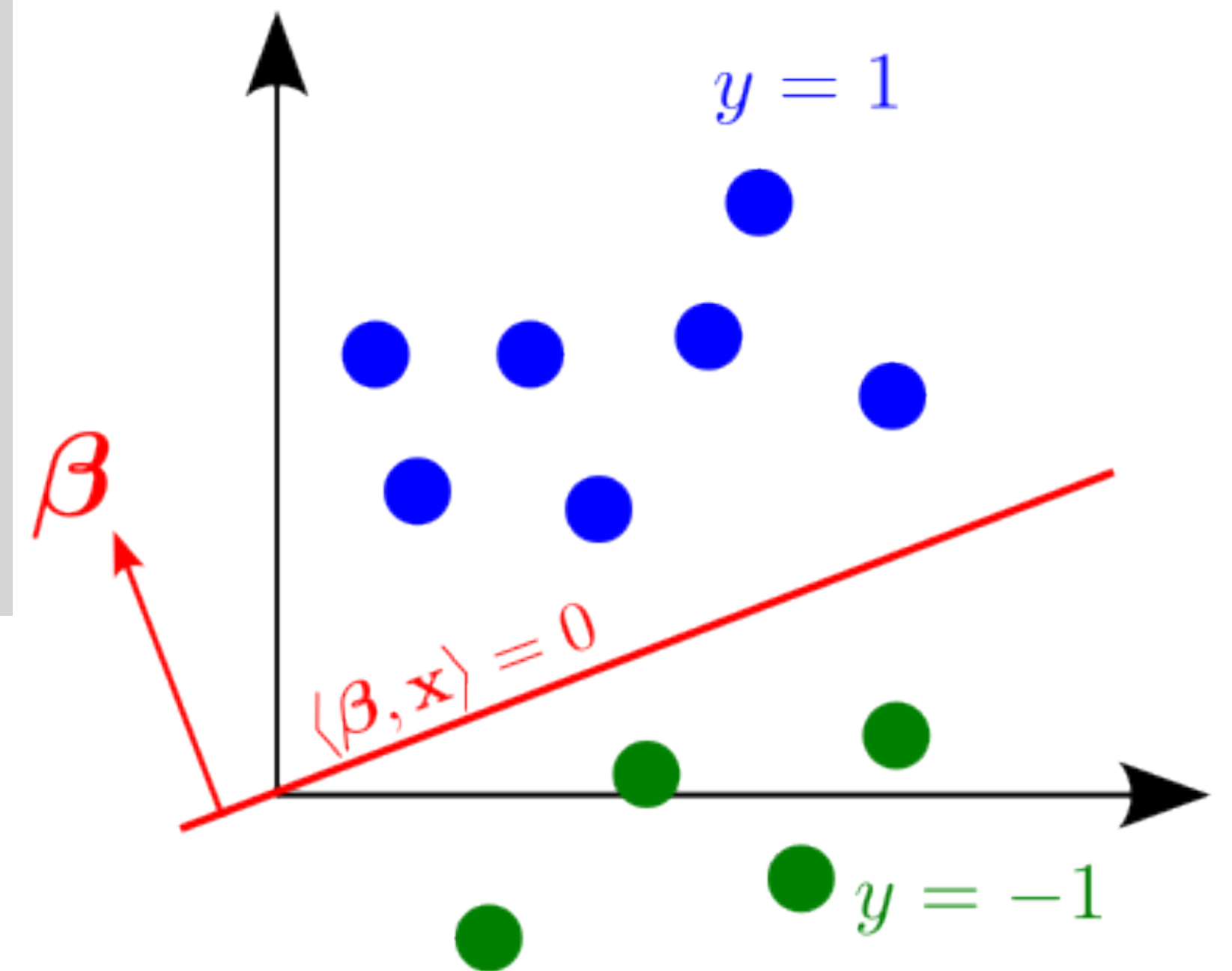
$$\min_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S) = \sum_{i=1}^n \ell_{\text{exp}} \left( \langle \mathbf{x}_i, \mathcal{P}_{\text{fc}}(\mathbf{W}) \rangle, y_i \right)$$

- Trained under **gradient flow (GF)**:

$$\frac{d\mathbf{W}}{dt} = - \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S)$$

**Question:** to what does  $\beta_{\text{fc}}(t) = \mathcal{P}_{\text{fc}}(\mathbf{W}(t))$  converge?

Linear binary classification



# Implicit regularization of training algorithms

$$\min_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S) = \sum_{i=1}^n \ell_{\text{exp}} \left( \langle \mathbf{x}_i, \mathcal{P}_{\text{fc}}(\mathbf{W}) \rangle, y_i \right)$$

- Trained under **gradient flow (GF)**:

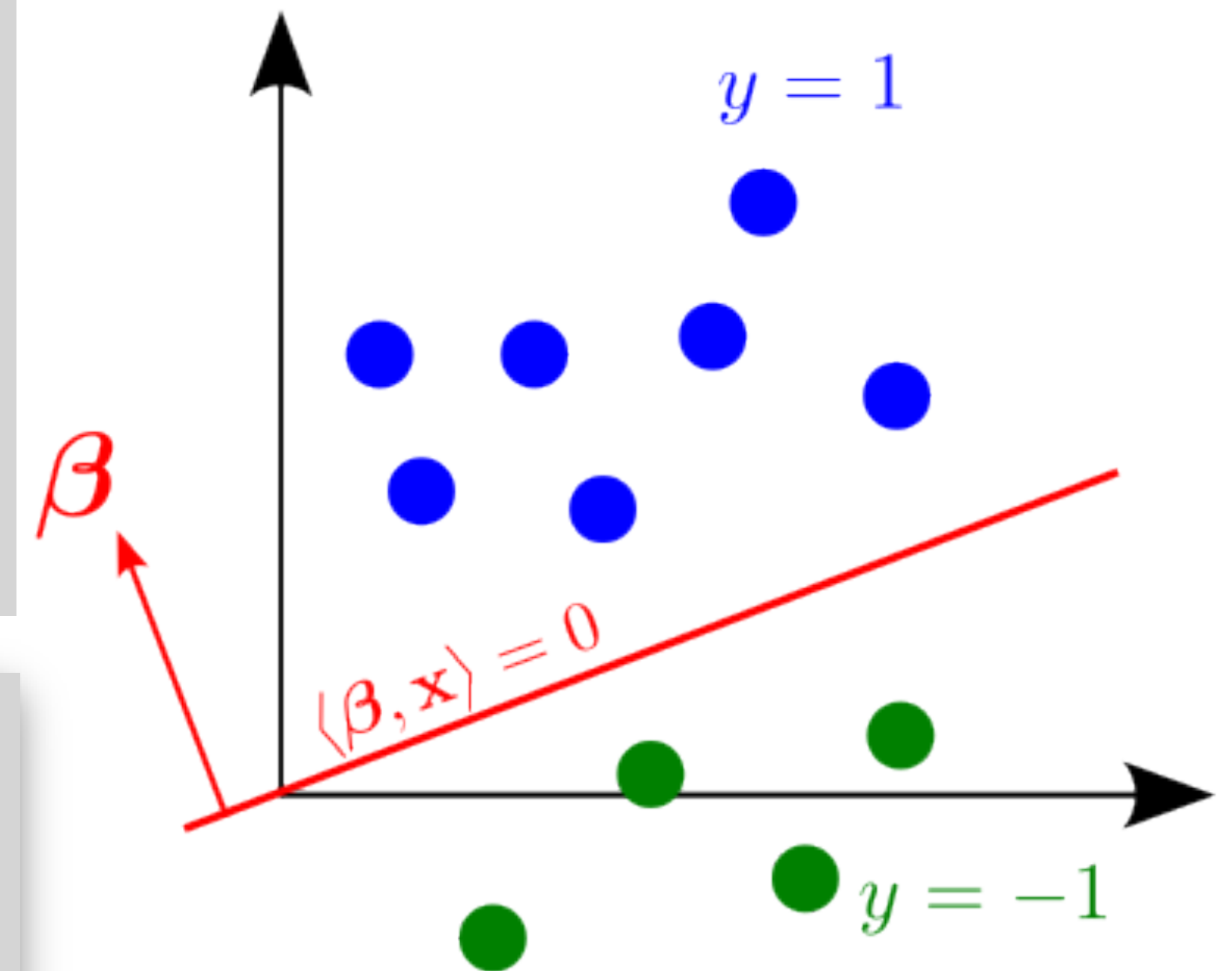
$$\frac{d\mathbf{W}}{dt} = - \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S)$$

**Question:** to what does  $\beta_{\text{fc}}(t) = \mathcal{P}_{\text{fc}}(\mathbf{W}(t))$  converge?

**Fact** [Ji and Telgarsky, *ICLR* 2018], [Yun et al., *ICLR* 2021]

- $\beta_{\text{fc}}^{\infty} = \lim_{t \rightarrow \infty} \beta_{\text{fc}}(t) / \|\beta_{\text{fc}}(t)\|$  exists.

Linear binary classification



# Implicit regularization of training algorithms

$$\min_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S) = \sum_{i=1}^n \ell_{\text{exp}} \left( \langle \mathbf{x}_i, \mathcal{P}_{\text{fc}}(\mathbf{W}) \rangle, y_i \right)$$

- Trained under **gradient flow (GF)**:

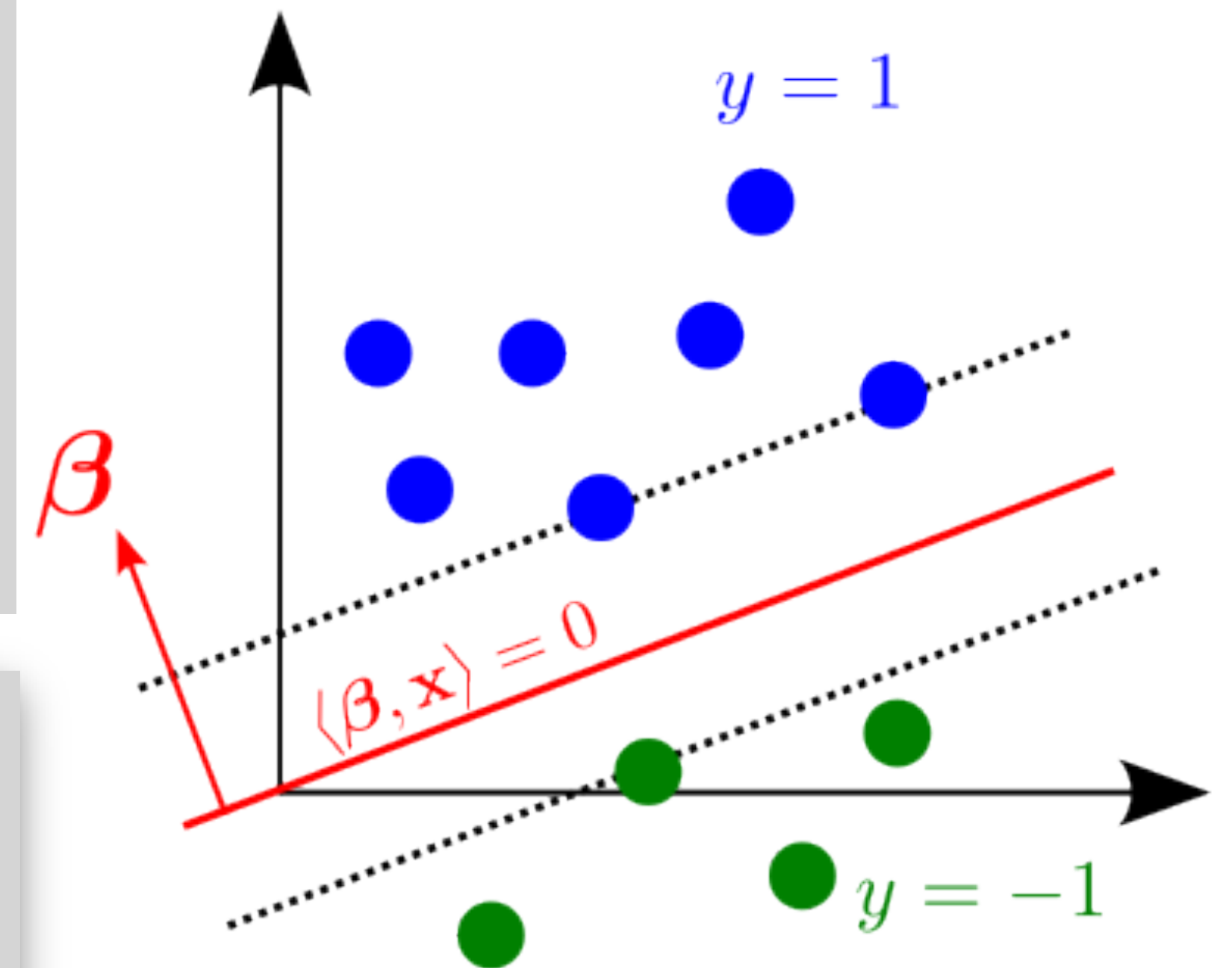
$$\frac{d\mathbf{W}}{dt} = - \nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{fc}}}(\mathbf{W}; S)$$

**Question:** to what does  $\beta_{\text{fc}}(t) = \mathcal{P}_{\text{fc}}(\mathbf{W}(t))$  converge?

**Fact** [Ji and Telgarsky, *ICLR* 2018], [Yun et al., *ICLR* 2021]

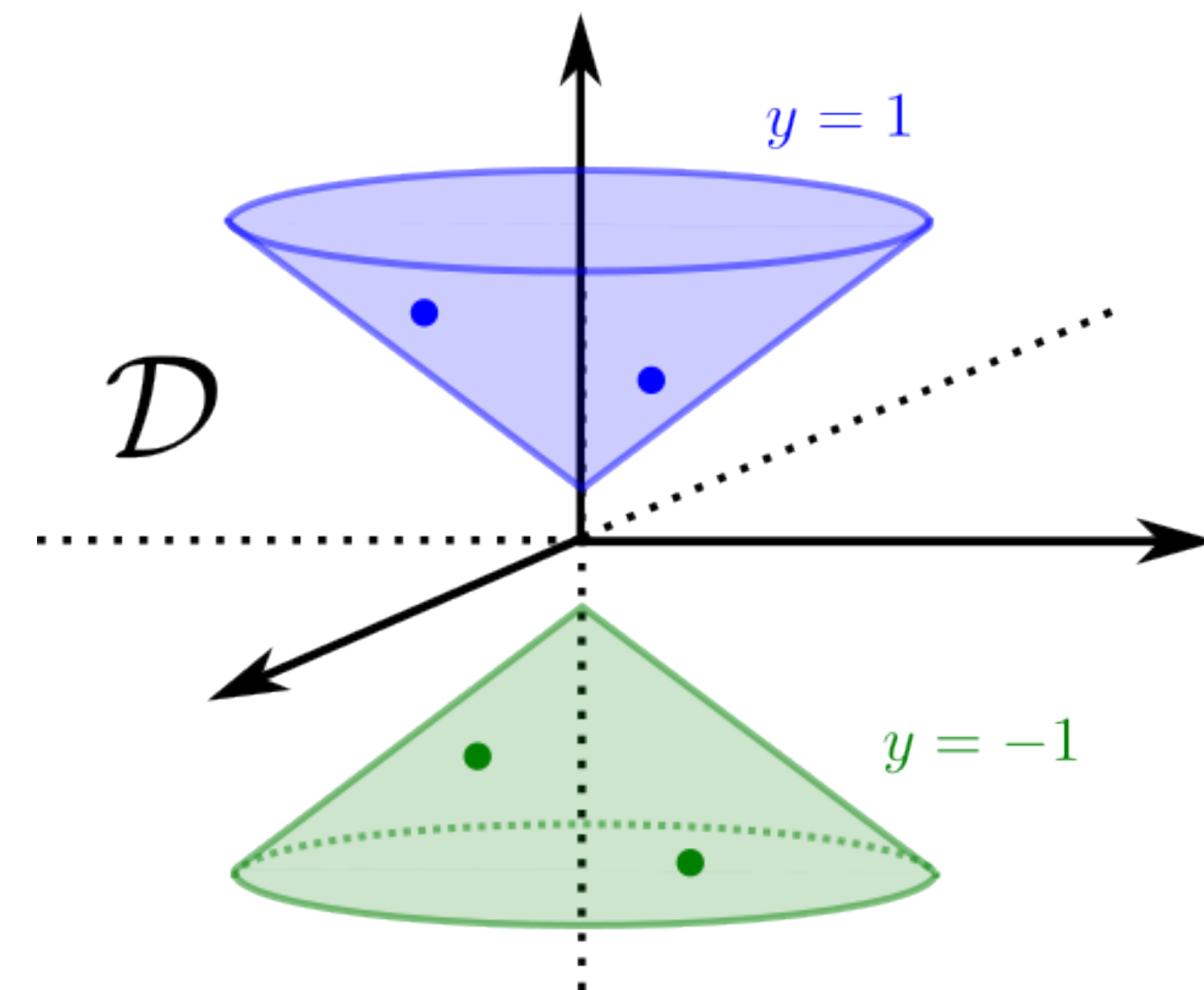
- $\beta_{\text{fc}}^{\infty} = \lim_{t \rightarrow \infty} \beta_{\text{fc}}(t) / \|\beta_{\text{fc}}(t)\|$  exists.
- $\beta_{\text{fc}}^{\infty}$  is the the **max- $L^2$ -margin** support vector machine (SVM).

Linear binary classification



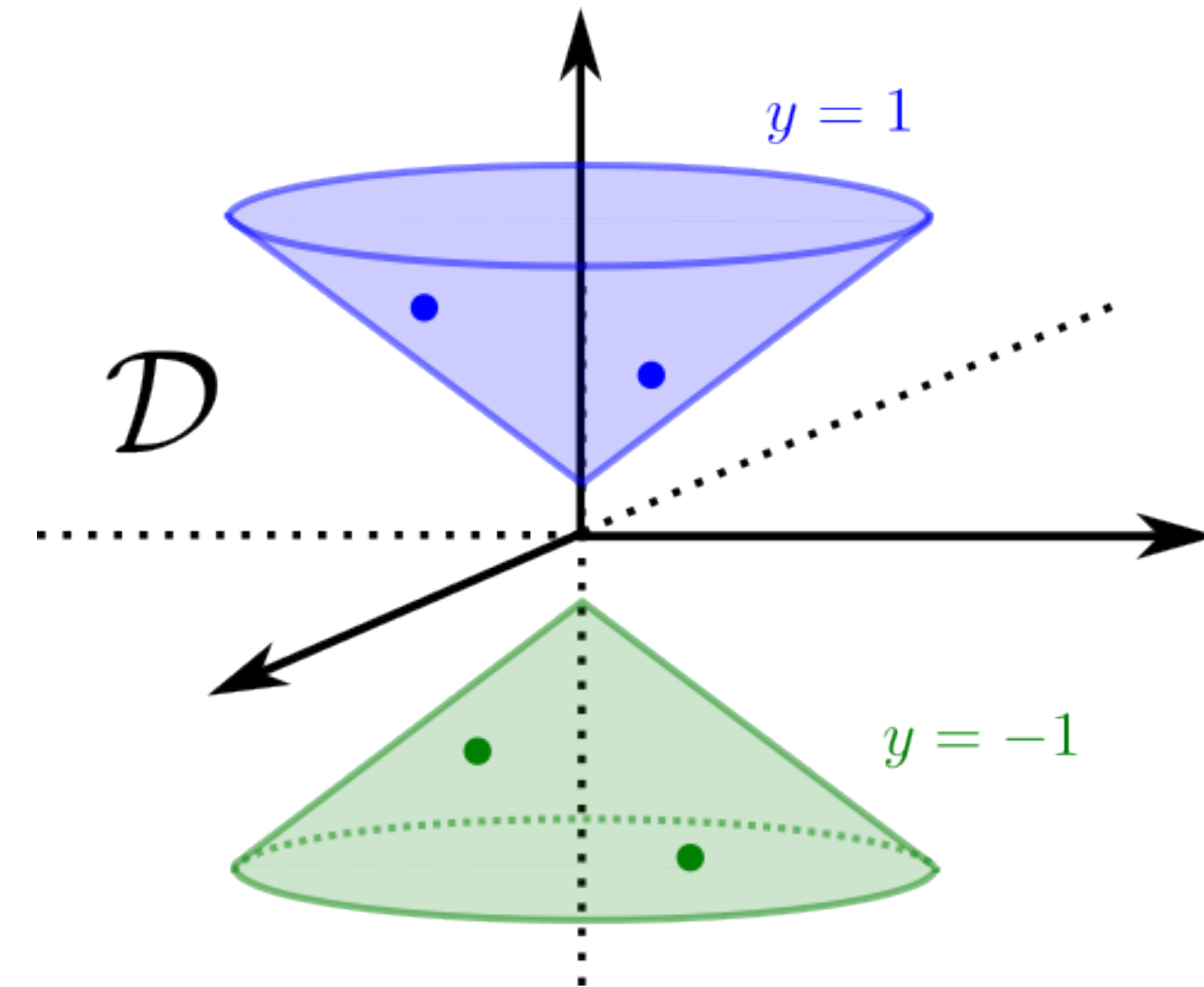


# Group-invariant binary classification



# Group-invariant binary classification

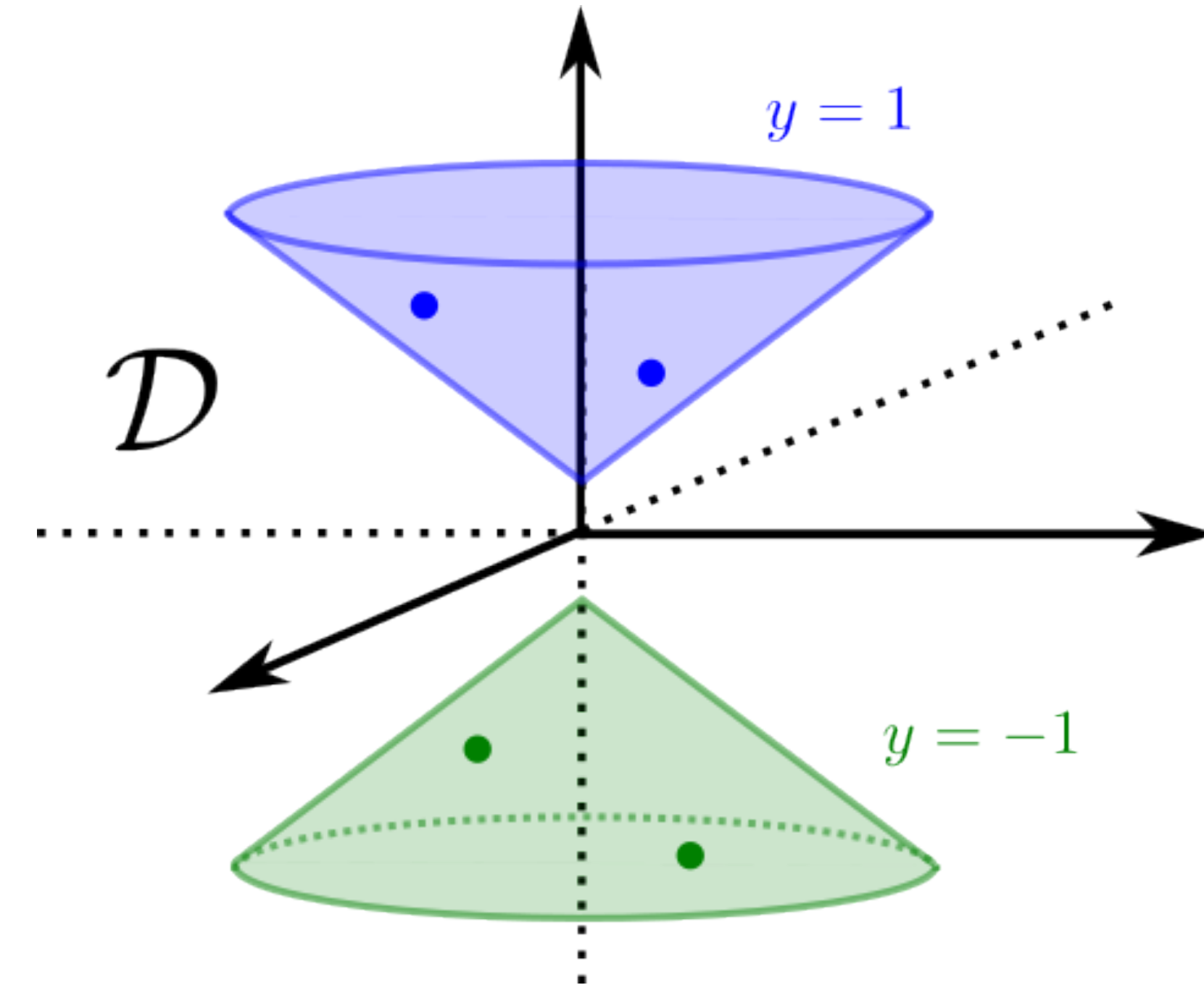
- Assume  $S \sim \mathcal{D}$ , and  $\mathcal{D}$  is invariant to a linear  $G$ -action.



# Group-invariant binary classification

- Assume  $S \sim \mathcal{D}$ , and  $\mathcal{D}$  is invariant to a linear  $G$ -action.
- Parameterize the invariant linear predictor  $\beta$  using a **G-CNN**,

$$f_{\text{inv}}(\mathbf{x}; \mathbf{W}) = \langle \mathbf{x}, \mathcal{P}_{\text{inv}}(\mathbf{W}) \rangle$$





# Group-invariant binary classification

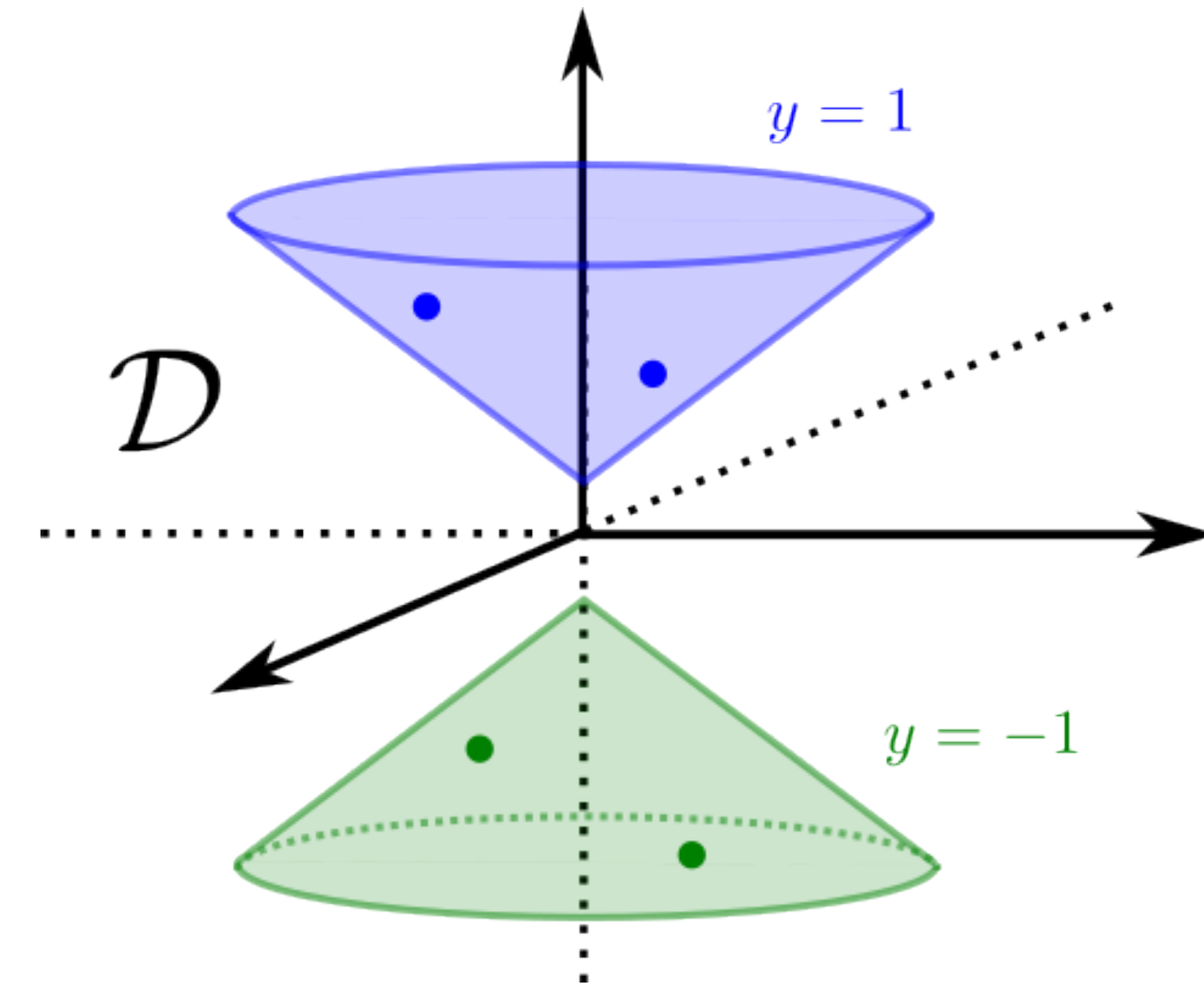
- Assume  $S \sim \mathcal{D}$ , and  $\mathcal{D}$  is invariant to a linear  $G$ -action.

- Parameterize the invariant linear predictor  $\beta$  using a **G-CNN**,

$$f_{\text{inv}}(\mathbf{x}; \mathbf{W}) = \langle \mathbf{x}, \mathcal{P}_{\text{inv}}(\mathbf{W}) \rangle$$

- Regression:  $\min_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{inv}}}(\mathbf{W}; S) = \sum_{i=1}^n \ell_{\text{exp}} \left( \langle \mathbf{x}_i, \mathcal{P}_{\text{inv}}(\mathbf{W}) \rangle, y_i \right)$

- Gradient flow:  $\frac{d\mathbf{W}}{dt} = -\nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{inv}}}(\mathbf{W}; S)$



# Group-invariant binary classification

- Assume  $S \sim \mathcal{D}$ , and  $\mathcal{D}$  is invariant to a linear  $G$ -action.

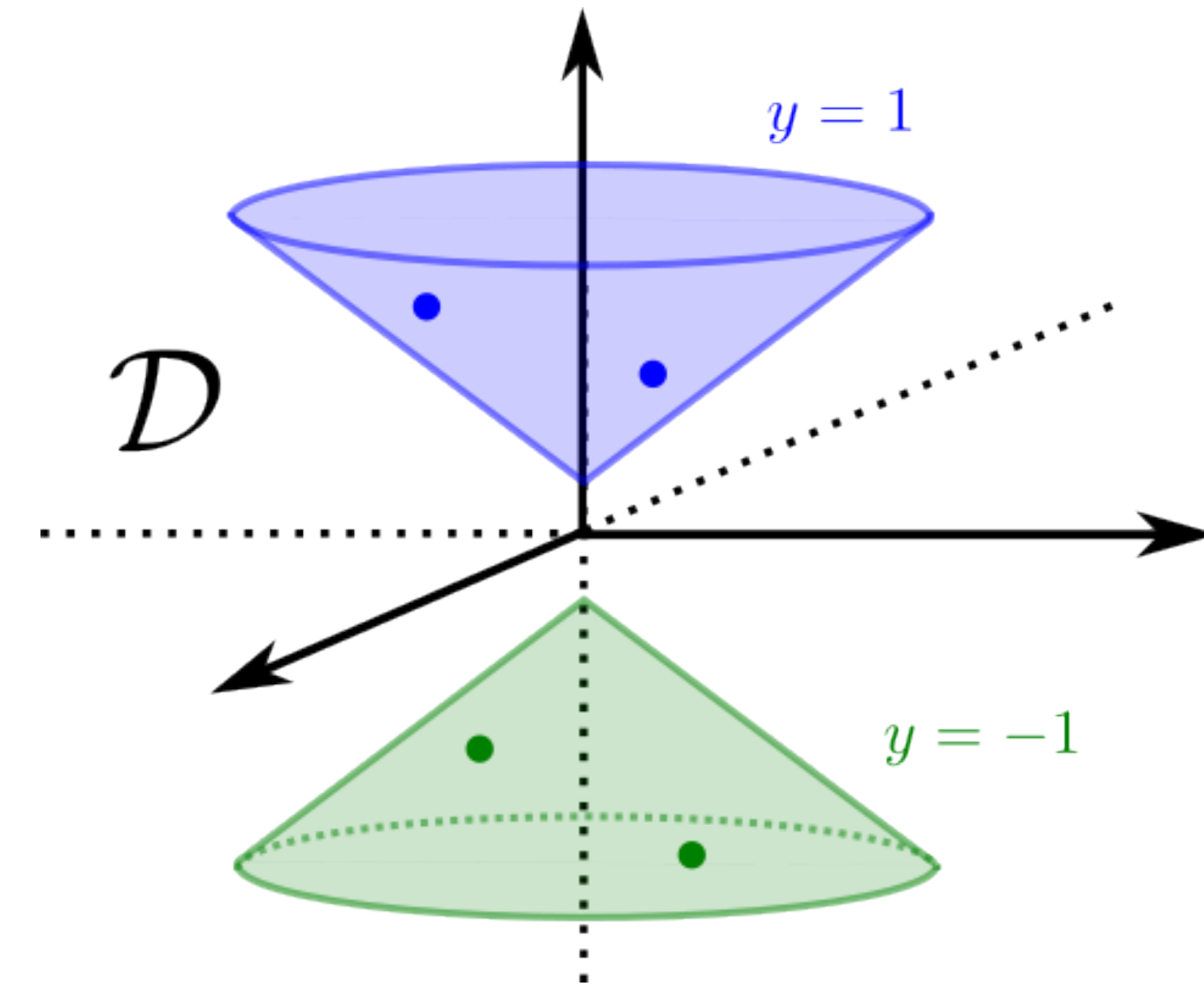
- Parameterize the invariant linear predictor  $\beta$  using a **G-CNN**,

$$f_{\text{inv}}(\mathbf{x}; \mathbf{W}) = \langle \mathbf{x}, \mathcal{P}_{\text{inv}}(\mathbf{W}) \rangle$$

- Regression:  $\min_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{inv}}}(\mathbf{W}; S) = \sum_{i=1}^n \ell_{\text{exp}} \left( \langle \mathbf{x}_i, \mathcal{P}_{\text{inv}}(\mathbf{W}) \rangle, y_i \right)$

- Gradient flow:  $\frac{d\mathbf{W}}{dt} = -\nabla_{\mathbf{W}} \mathcal{L}_{\mathcal{P}_{\text{inv}}}(\mathbf{W}; S)$

**Question:** to what does  $\beta_{\text{inv}}(t) = \mathcal{P}_{\text{inv}}(\mathbf{W}(t))$  converge?



# Group-invariant binary classification



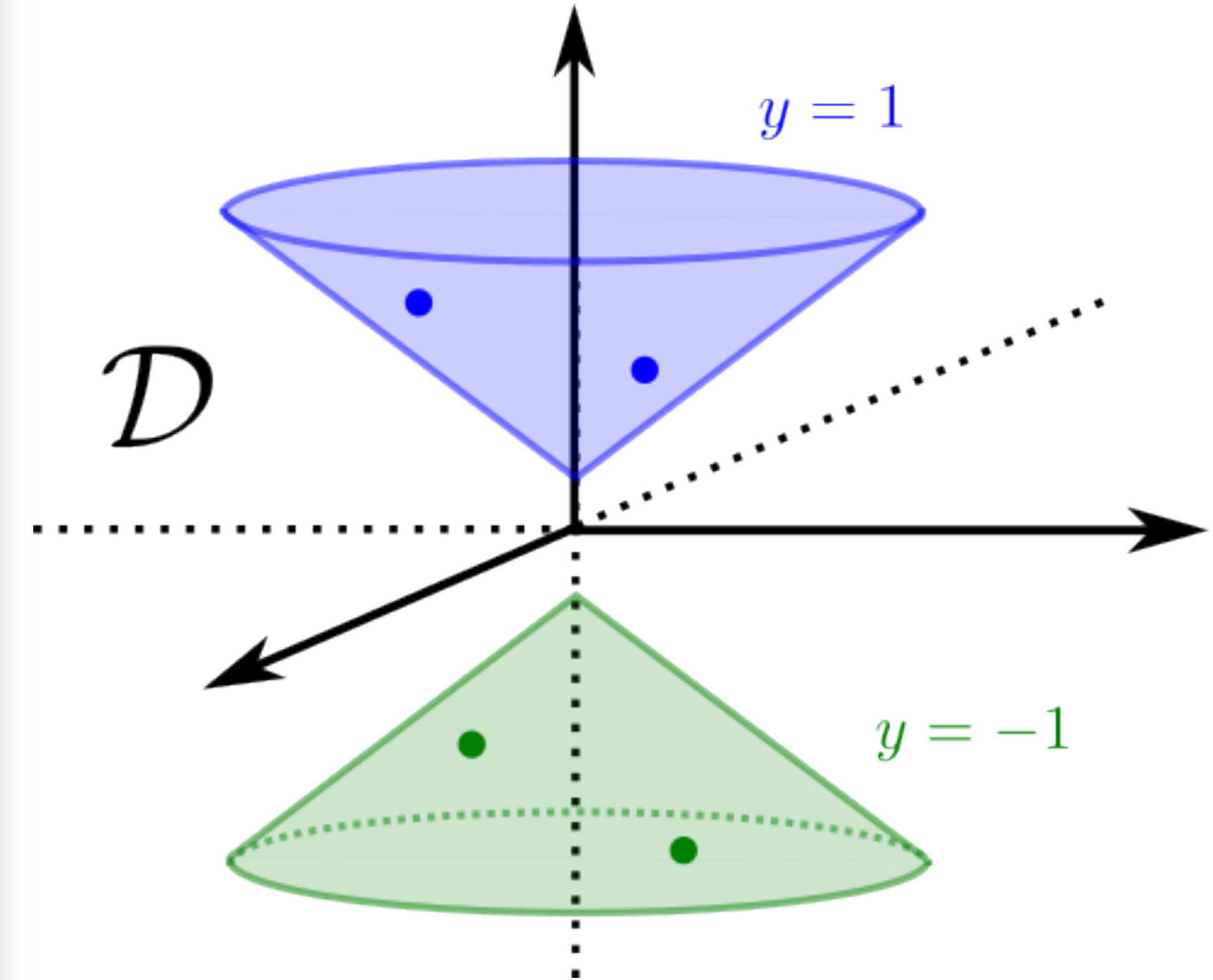
# Group-invariant binary classification

**Question:** to what does  $\beta_{\text{inv}}(t) = \mathcal{P}_{\text{inv}}(\mathbf{W}(t))$  converge?

**Theorem** [Chen and Z., *NeurIPS* 2023]

If the input linear  $G$ -action is unitary, then

- $\beta_{\text{inv}}^{\infty} = \lim_{t \rightarrow \infty} \beta_{\text{inv}}(t) / \|\beta_{\text{inv}}(t)\|$  exists.



# Group-invariant binary classification

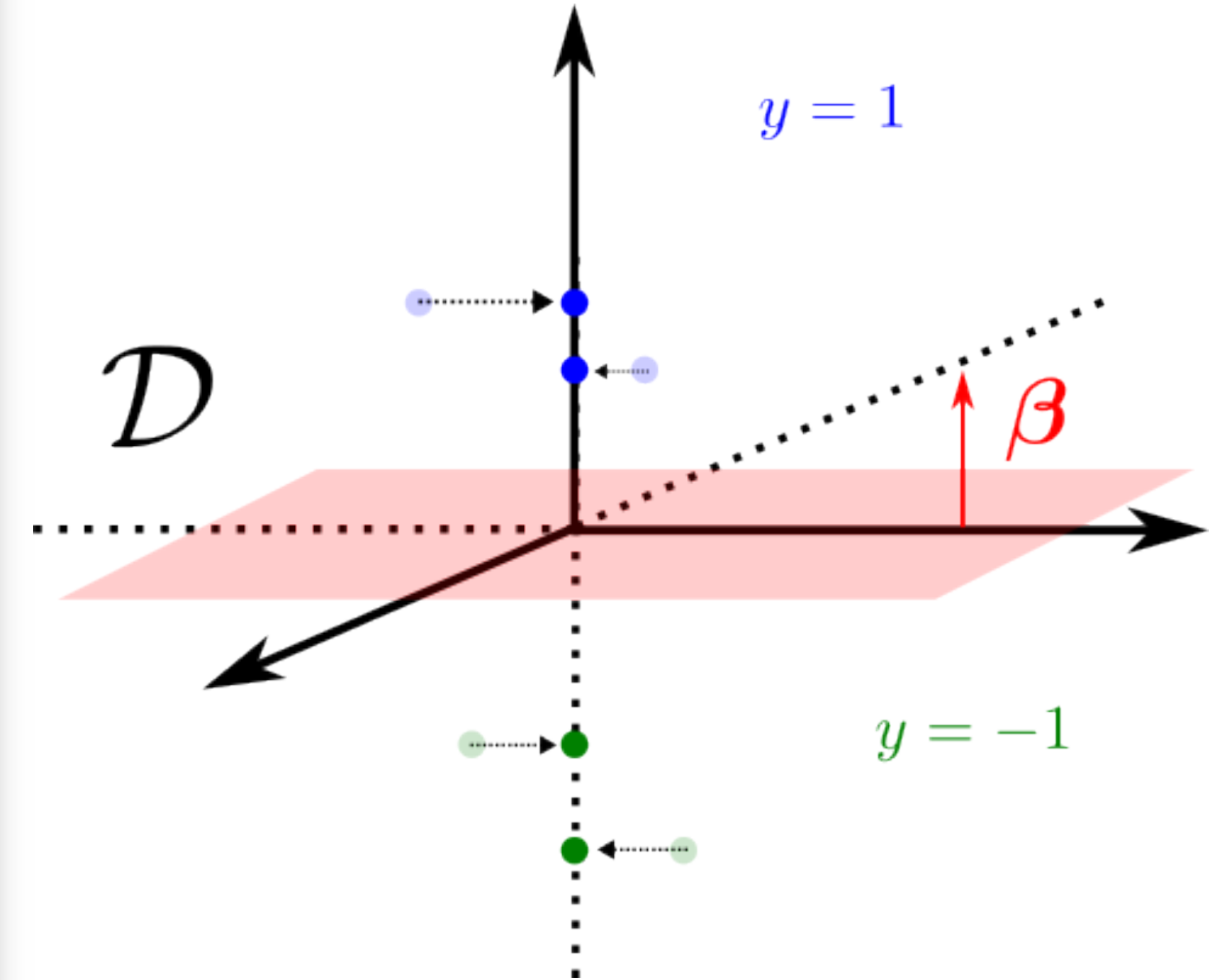
**Question:** to what does  $\beta_{\text{inv}}(t) = \mathcal{P}_{\text{inv}}(\mathbf{W}(t))$  converge?

**Theorem** [Chen and Z., *NeurIPS* 2023]

If the input linear  $G$ -action is unitary, then

- $\beta_{\text{inv}}^{\infty} = \lim_{t \rightarrow \infty} \beta_{\text{inv}}(t) / \|\beta_{\text{inv}}(t)\|$  exists.
- $\beta_{\text{inv}}^{\infty}$  is the the max-margin SVM on the **transformed dataset**

$$\bar{S} = \{(\bar{\mathbf{x}}_i, y_i) : i \in [n]\}, \text{ where } \bar{\mathbf{x}} = \frac{1}{|G|} \sum_{g \in G} g\mathbf{x}$$



# Group-invariant binary classification

**Question:** to what does  $\beta_{\text{inv}}(t) = \mathcal{P}_{\text{inv}}(\mathbf{W}(t))$  converge?

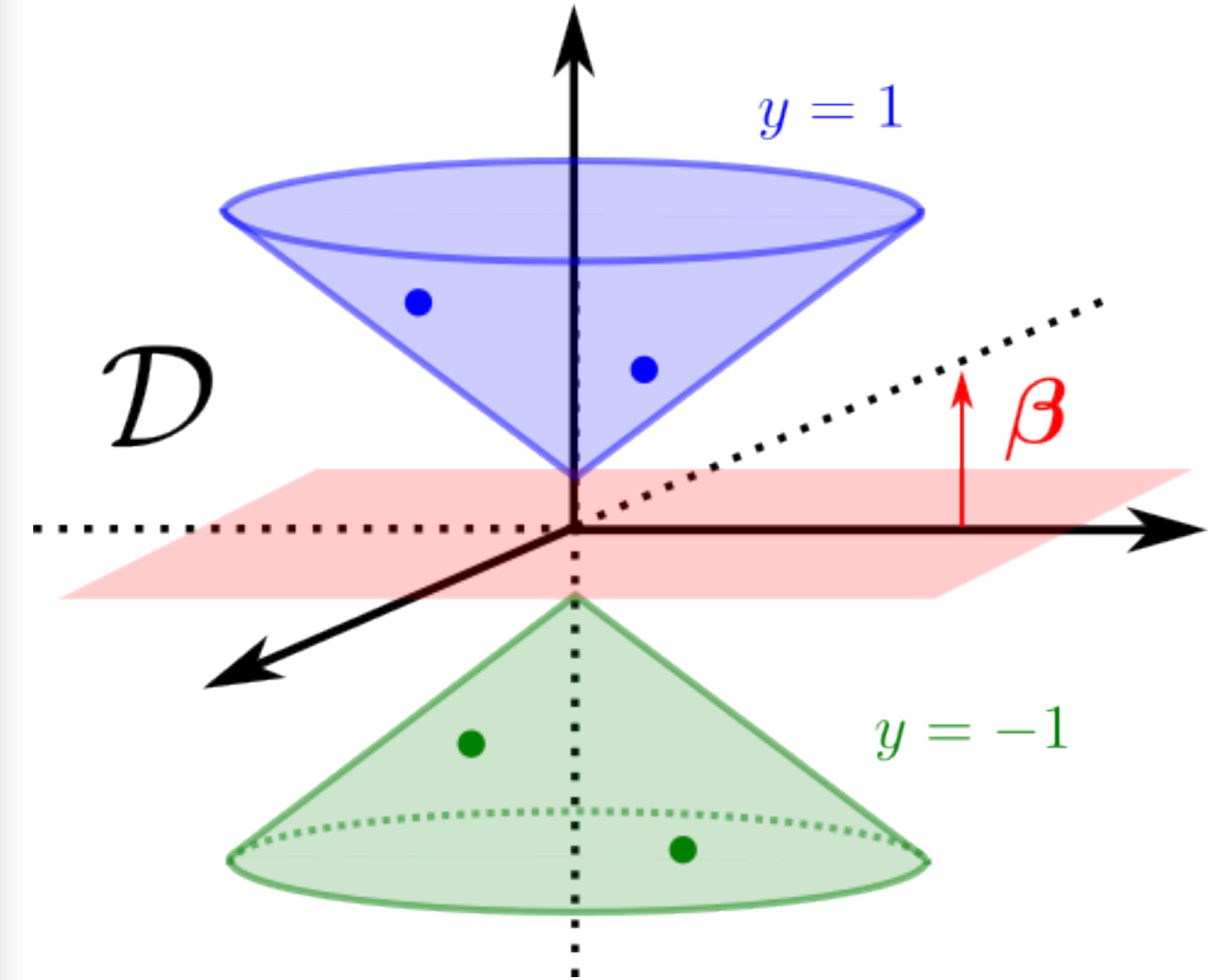
**Theorem** [Chen and Z., *NeurIPS* 2023]

If the input linear  $G$ -action is unitary, then

- $\beta_{\text{inv}}^{\infty} = \lim_{t \rightarrow \infty} \beta_{\text{inv}}(t) / \|\beta_{\text{inv}}(t)\|$  exists.
- $\beta_{\text{inv}}^{\infty}$  is the the max-margin SVM on the **transformed dataset**

$$\bar{S} = \{(\bar{\mathbf{x}}_i, y_i) : i \in [n]\}, \text{ where } \bar{\mathbf{x}} = \frac{1}{|G|} \sum_{g \in G} g\mathbf{x}$$

- $\beta_{\text{inv}}^{\infty}$  is the unique max-margin **invariant** SVM on the **original dataset**  $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$





# Implications

# Implications

## Corollary (G-CNN vs data augmentation)

- $\beta_{\text{inv}}^{\infty}$ : linear **G-CNN** trained on  $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$ .
- $\beta_{\text{fc}}^{\infty}$ : linear fully-connected network trained on  $S_{\text{aug}} = \{(g\mathbf{x}_i, y_i) : i \in [n], g \in G\}$ .

$$\beta_{\text{steer}}^{\infty} = \beta_{\text{fc}}^{\infty}$$

# Implications

## Corollary (G-CNN vs data augmentation)

- $\beta_{\text{inv}}^{\infty}$ : linear **G-CNN** trained on  $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$ .
- $\beta_{\text{fc}}^{\infty}$ : linear fully-connected network trained on  $S_{\text{aug}} = \{(g\mathbf{x}_i, y_i) : i \in [n], g \in G\}$ .

$$\beta_{\text{steer}}^{\infty} = \beta_{\text{fc}}^{\infty}$$

- **Equivariant neural networks** is equivalent to data augmentation.



# Implications

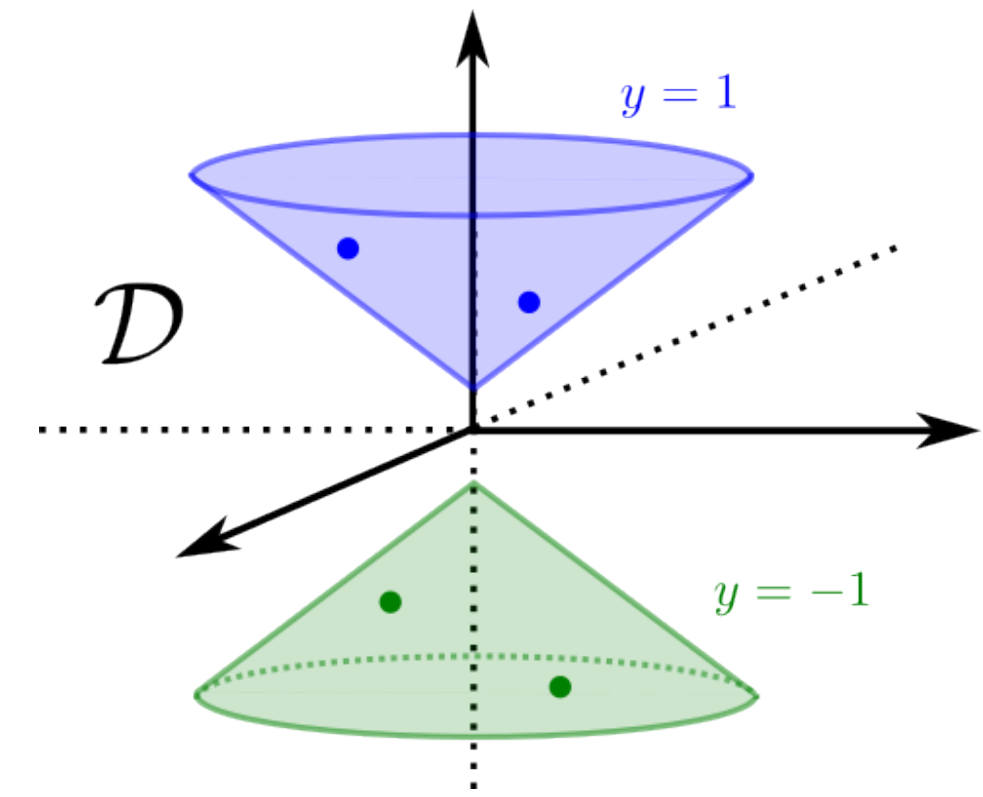
## Corollary (G-CNN vs data augmentation)

- $\beta_{\text{inv}}^{\infty}$ : linear **G-CNN** trained on  $S = \{(\mathbf{x}_i, y_i) : i \in [n]\}$ .
- $\beta_{\text{fc}}^{\infty}$ : linear fully-connected network trained on  $S_{\text{aug}} = \{(g\mathbf{x}_i, y_i) : i \in [n], g \in G\}$ .

$$\beta_{\text{steer}}^{\infty} = \beta_{\text{fc}}^{\infty}$$

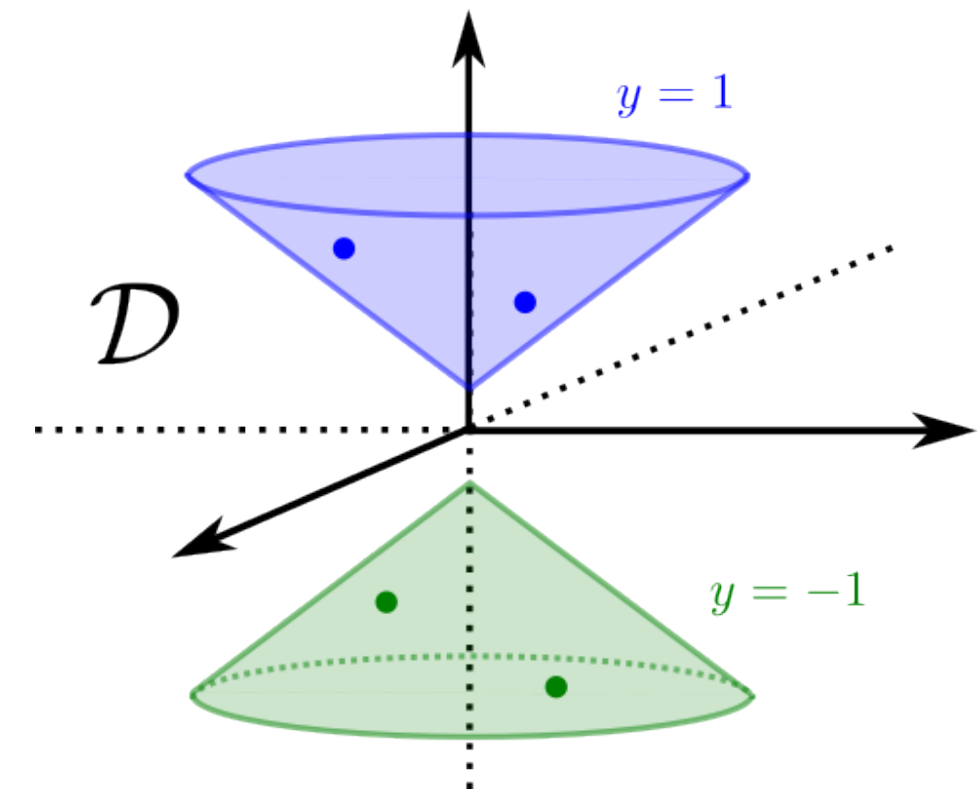
- **Equivariant neural networks** is equivalent to data augmentation.
- Caveat:
  - Full data augmentation on the entire group  $G$ .
  - Unitary input action.
  - Only linear models.

# Improved generalization



# Improved generalization

- **$G$ -invariant** distribution  $\mathcal{D}$  over  $\mathbb{R}^{d_0} \times \{\pm 1\}$ .
- $\mathcal{D}$  can be separated by an invariant classifier  $\beta_0$ .





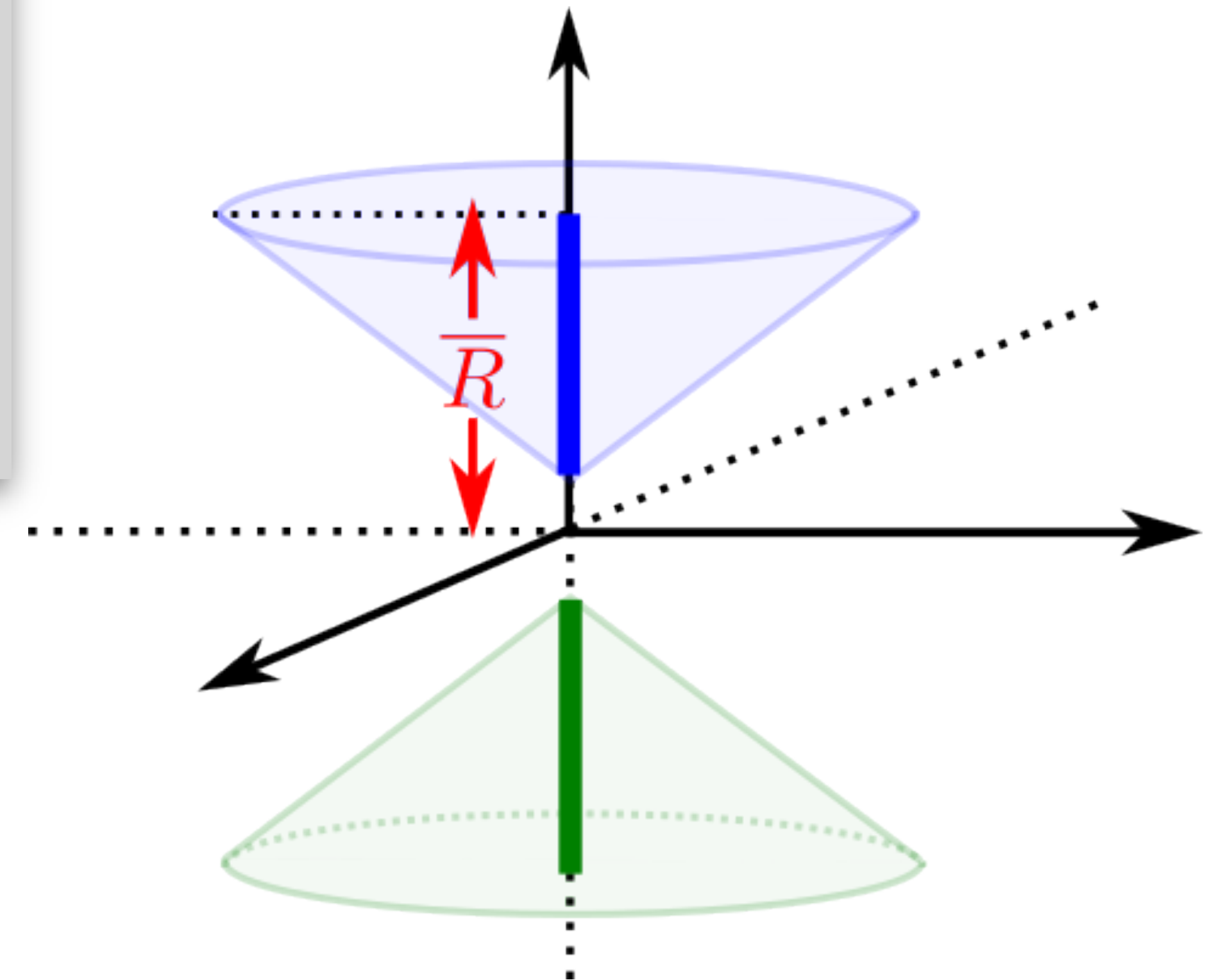
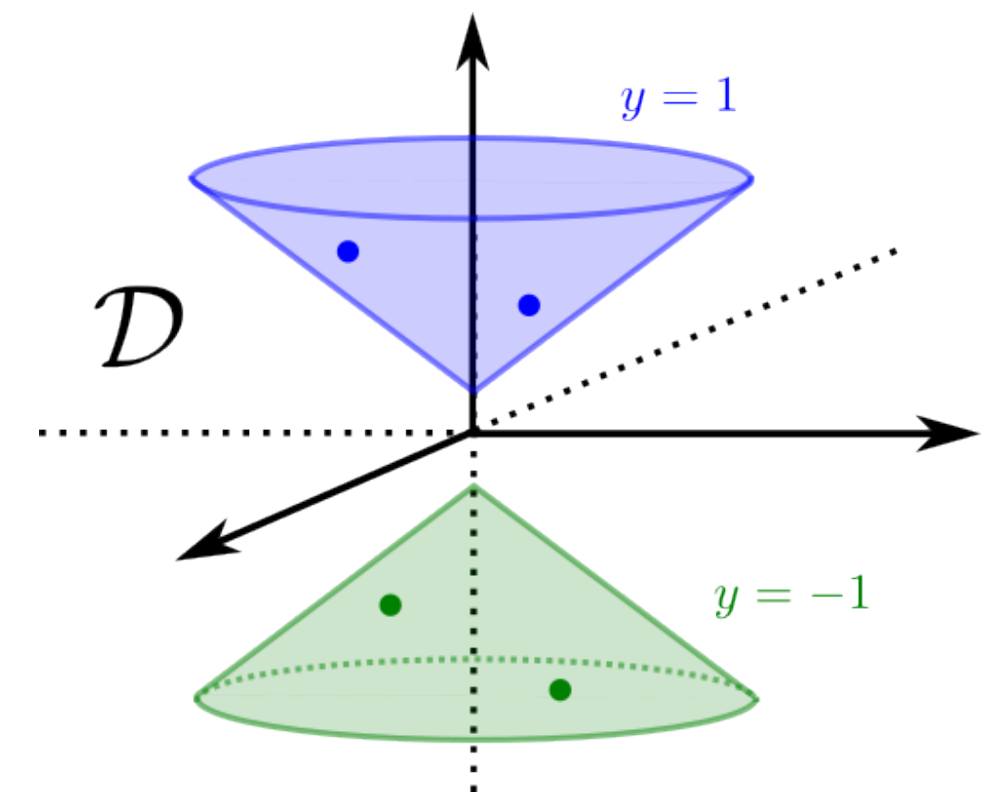
# Improved generalization

- $G$ -invariant distribution  $\mathcal{D}$  over  $\mathbb{R}^{d_0} \times \{\pm 1\}$ .
- $\mathcal{D}$  can be separated by an invariant classifier  $\beta_0$ .

**Theorem** [Chen and Z., *NeurIPS* 2023]

Let  $\bar{R} = \inf \{ r > 0 : \|\bar{\mathbf{x}}\| \leq r \}$ . For any  $\delta > 0$ , w.p. at least  $1 - \delta$ ,

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ y \neq \text{sign} \left( \left\langle \mathbf{x}, \beta_{\text{inv}}^\infty \right\rangle \right) \right] \leq \frac{2\bar{R}\|\beta_0\|}{\sqrt{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$



# Improved generalization

- **$G$ -invariant** distribution  $\mathcal{D}$  over  $\mathbb{R}^{d_0} \times \{\pm 1\}$ .
- $\mathcal{D}$  can be separated by an invariant classifier  $\beta_0$ .

**Theorem** [Chen and Z., *NeurIPS* 2023]

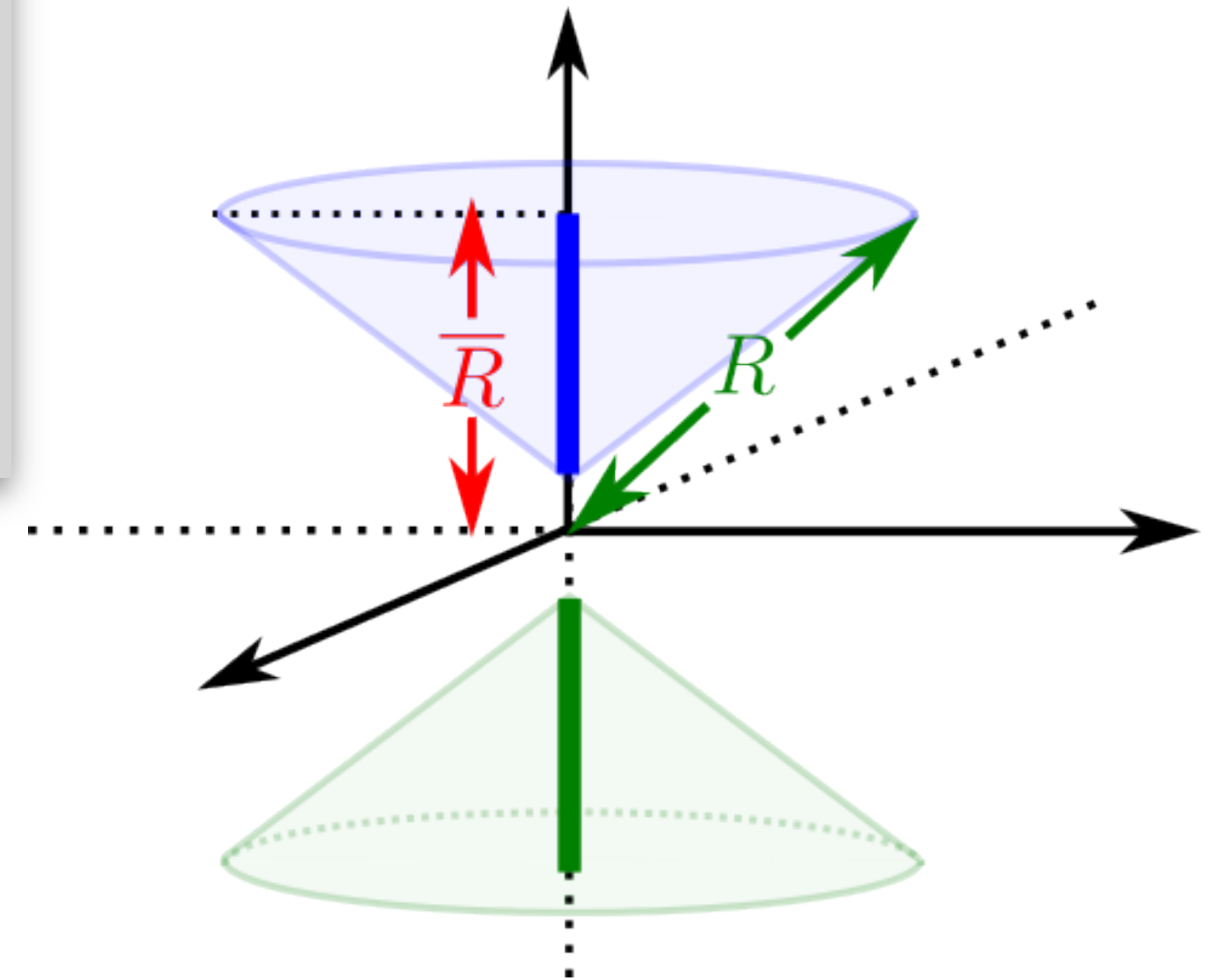
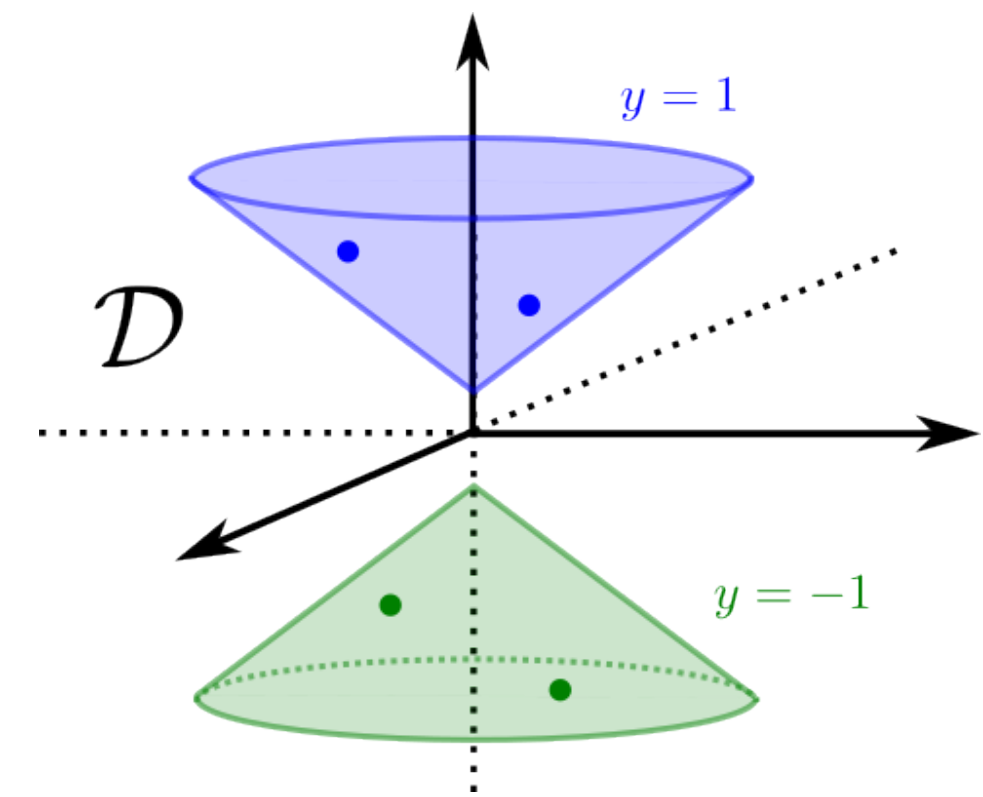
Let  $\bar{R} = \inf \{ r > 0 : \|\bar{\mathbf{x}}\| \leq r \}$ . For any  $\delta > 0$ , w.p. at least  $1 - \delta$ ,

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ y \neq \text{sign} \left( \langle \mathbf{x}, \beta_{\text{inv}}^\infty \rangle \right) \right] \leq \frac{2\bar{R}\|\beta_0\|}{\sqrt{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

**Remark:** In comparison, for fully-connected networks, we have

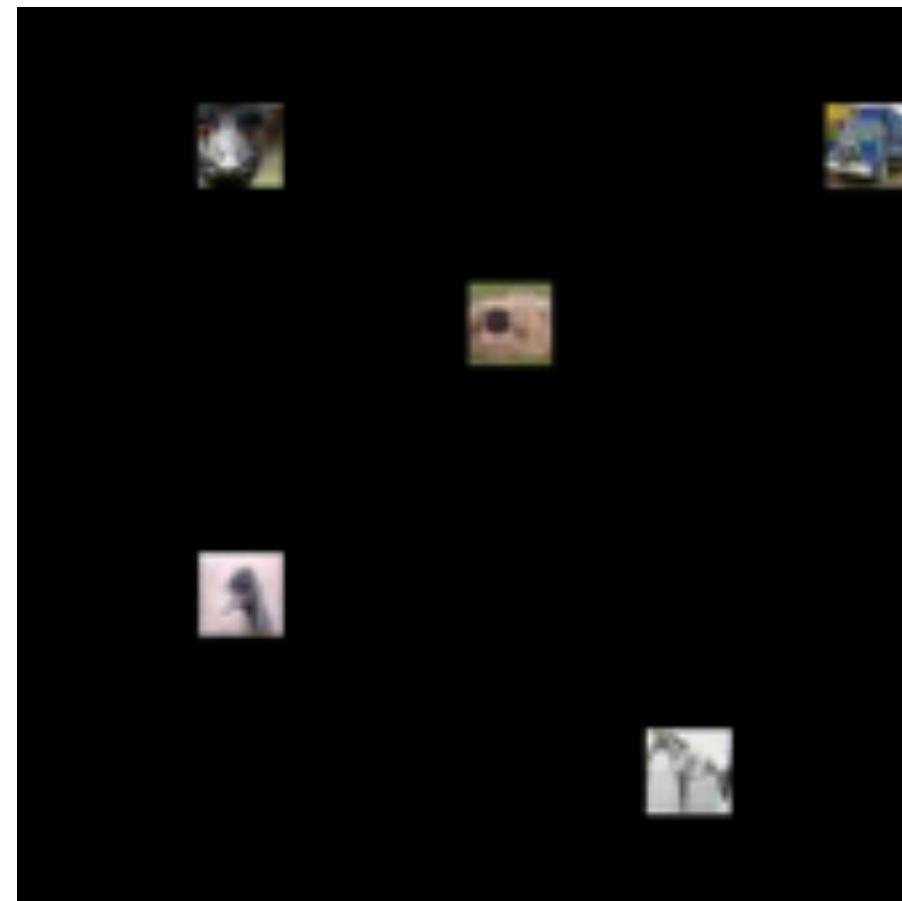
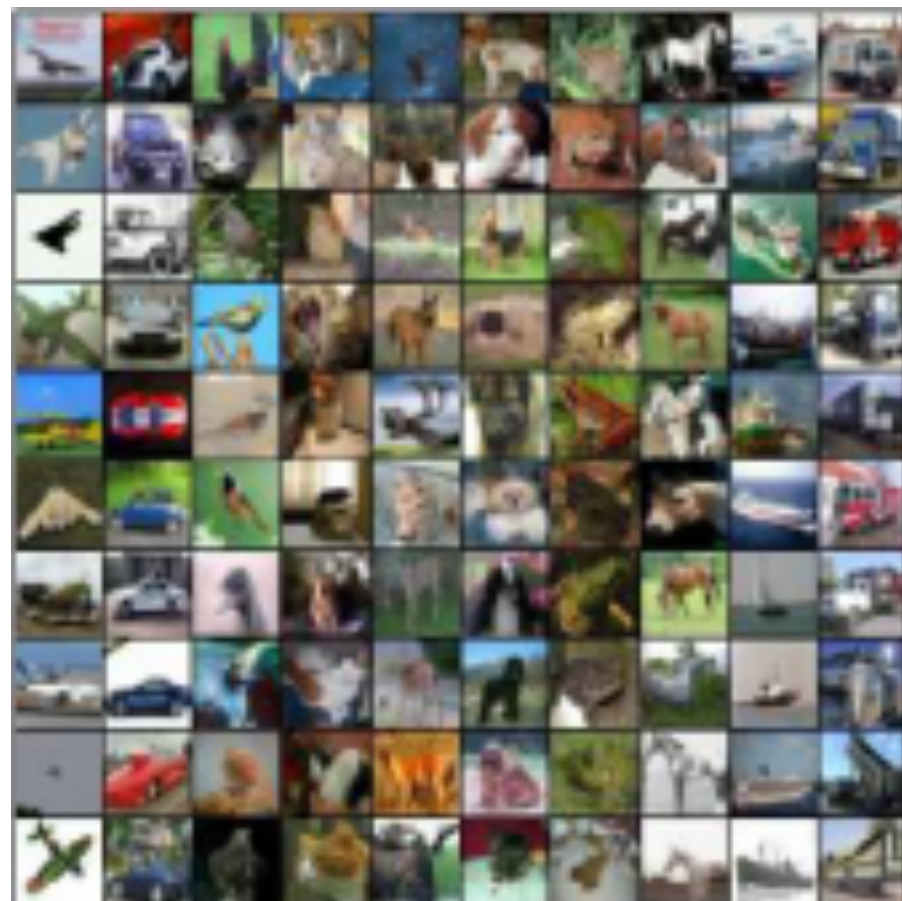
$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ y \neq \text{sign} \left( \langle \mathbf{x}, \beta_{\text{fc}}^\infty \rangle \right) \right] \leq \frac{2R\|\beta_0\|}{\sqrt{n}} + \sqrt{\frac{\log(1/\delta)}{2n}},$$

where  $R = \inf \{ r > 0 : \|\mathbf{x}\| \leq r \text{ with probability } 1 \} \geq \bar{R}$



# Conclusion

- **Exact quantification of the improvement**
  - **Sample complexity** and **error bound**.
- **Does it converge? To what solution?**
  - **Training dynamics** of equivariant models



$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; S) = \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{W}), y_i)$$

## Related papers

- J. Birrell, M.A. Katsoulakis, L. Rey-Bellet, **W. Zhu**. “Structure-preserving GANs”. *ICML* (2022)
- Z. Chen, M.A. Katsoulakis, L. Rey-Bellet, **W. Zhu**. “Sample complexity of probability divergences under group symmetry”. *ICML* (2023)
- Z. Chen and **W. Zhu**. “On the implicit bias of linear equivariant steerable networks: margin, generalization, and their equivalence to data augmentation”. *NeurIPS* (2023)

## Acknowledgement

NSF DMS-2052525, DMS-2140982, and DMS-2244976.



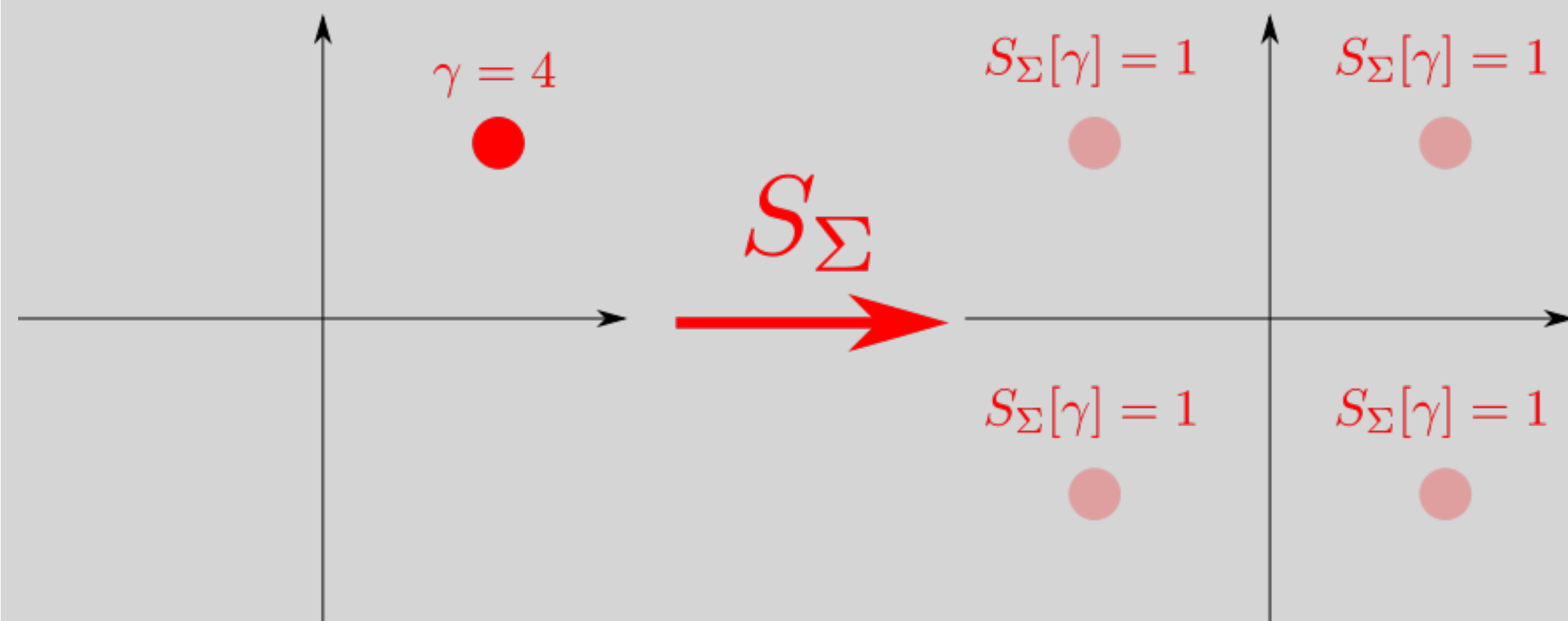


**Symmetrization operators  $S_\Sigma$  and  $S^\Sigma$**

# Symmetrization operators $S_\Sigma$ and $S^\Sigma$

- Symmetrization of functions:  $S_\Sigma : \mathcal{M}_b(X) \rightarrow \mathcal{M}_b(X)$ ,

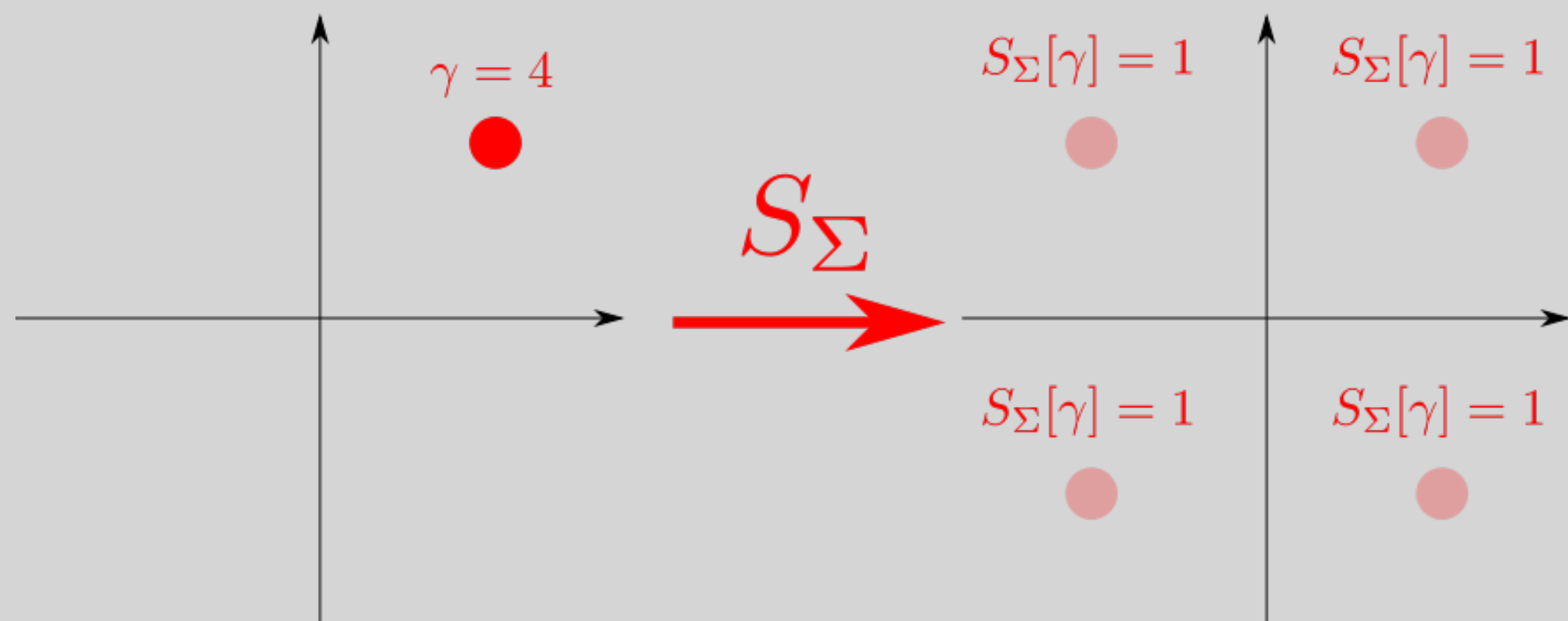
$$S_\Sigma[\gamma](x) = \int_\Sigma \gamma(T_{\sigma'}(x)) \mu_\Sigma(d\sigma') = E_{\mu_\Sigma}[\gamma \circ T_\sigma(x)].$$



# Symmetrization operators $S_\Sigma$ and $S^\Sigma$

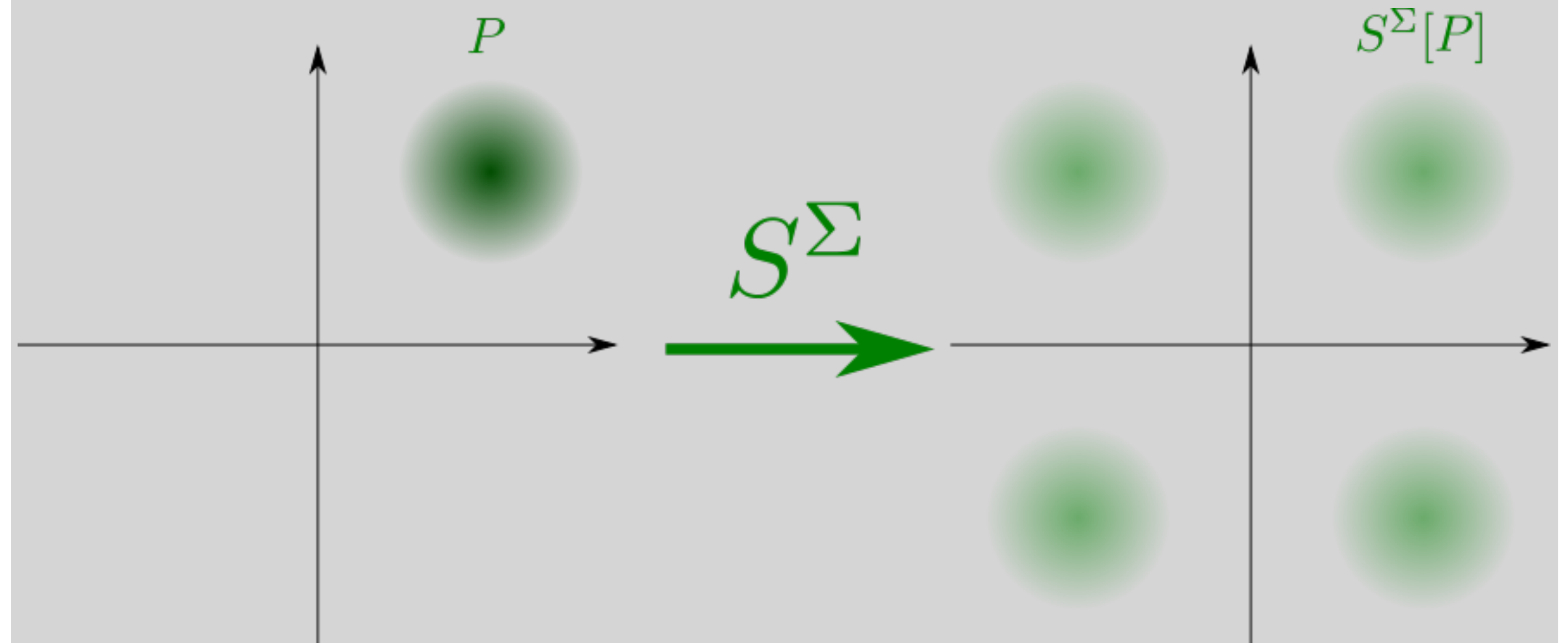
- Symmetrization of functions:  $S_\Sigma : \mathcal{M}_b(X) \rightarrow \mathcal{M}_b(X)$ ,

$$S_\Sigma[\gamma](x) = \int_\Sigma \gamma(T_{\sigma'}(x)) \mu_\Sigma(d\sigma') = E_{\mu_\Sigma}[\gamma \circ T_{\sigma'}(x)].$$



- Symmetrization of measures:  $S^\Sigma : \mathcal{P}(X) \rightarrow \mathcal{P}(X)$ ,

$$E_{S^\Sigma[P]}\gamma = \int_X S_\Sigma[\gamma](x) dP(x) = E_P S_\Sigma[\gamma], \forall \gamma \in \mathcal{M}_b(X).$$

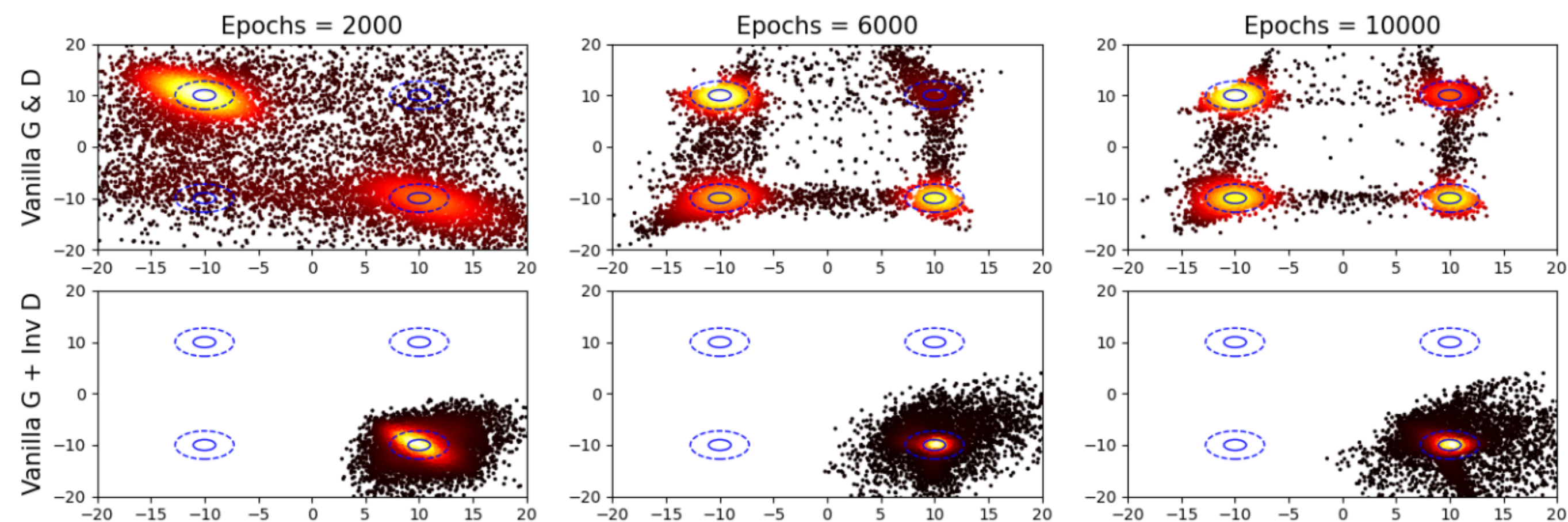


# Mode collapse — a warning

Theorem [Birrell, Katsoulakis, Rey-Bellet, Z., ICML 2022]

If  $S_{\Sigma}[\Gamma] \subset \Gamma$  and  $P, Q \in \mathcal{P}(X)$ , i.e., not necessarily  $\Sigma$ -invariant, then

$$D_{\Gamma_{\Sigma}^{\text{inv}}}(Q||P) = D^{\Gamma}(S^{\Sigma}[Q]||S^{\Sigma}[P]).$$





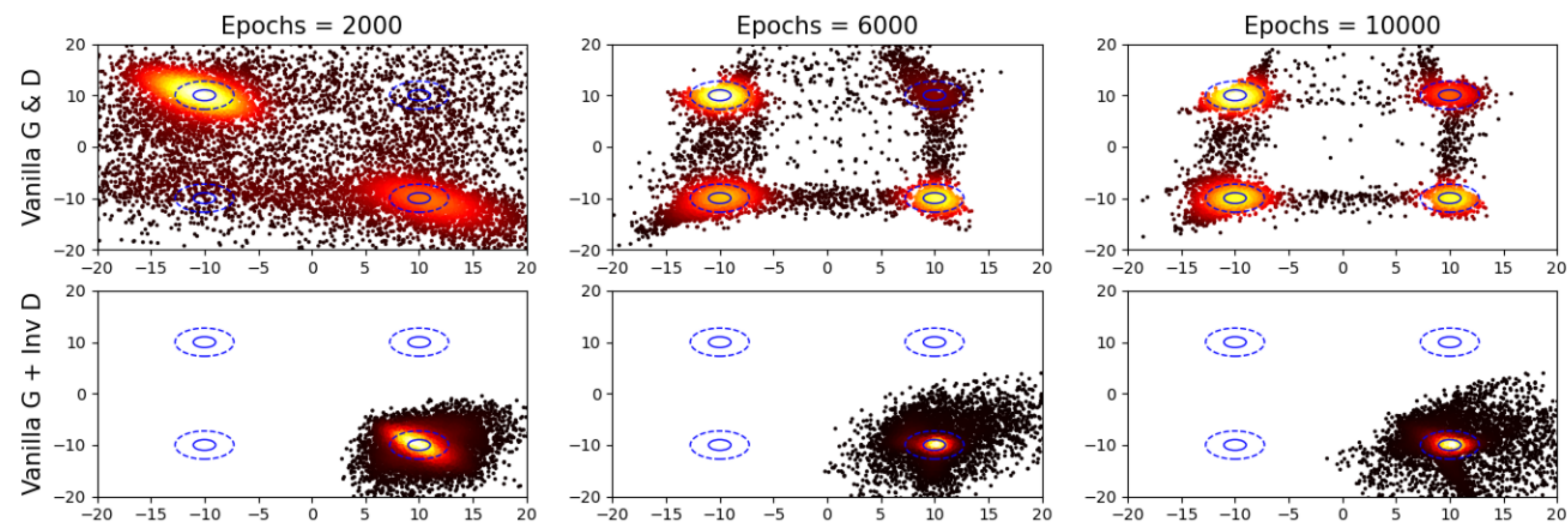
# Mode collapse — a warning

Theorem [Birrell, Katsoulakis, Rey-Bellet, Z., ICML 2022]

If  $S_\Sigma[\Gamma] \subset \Gamma$  and  $P, Q \in \mathcal{P}(X)$ , i.e., not necessarily  $\Sigma$ -invariant, then

$$D^{\Gamma_\Sigma^{\text{inv}}}(Q||P) = D^\Gamma(S^\Sigma[Q]||S^\Sigma[P]).$$

- Reducing  $\Gamma$  to  $\Gamma_\Sigma^{\text{inv}}$  might result in “mode collapse” if  $P_g$  is **NOT**  $\Sigma$ -invariant



# Mode collapse — a warning

Theorem [Birrell, Katsoulakis, Rey-Bellet, Z., ICML 2022]

If  $S_\Sigma[\Gamma] \subset \Gamma$  and  $P, Q \in \mathcal{P}(X)$ , i.e., not necessarily  $\Sigma$ -invariant, then

$$D_{\Gamma_\Sigma^{\text{inv}}}(Q||P) = D^\Gamma(S^\Sigma[Q]||S^\Sigma[P]).$$

- Reducing  $\Gamma$  to  $\Gamma_\Sigma^{\text{inv}}$  might result in “mode collapse” if  $P_g$  is **NOT**  $\Sigma$ -invariant
- The reason is as  $P_g$  only needs to equal  $Q$  after  $\Sigma$ -symmetrization.

