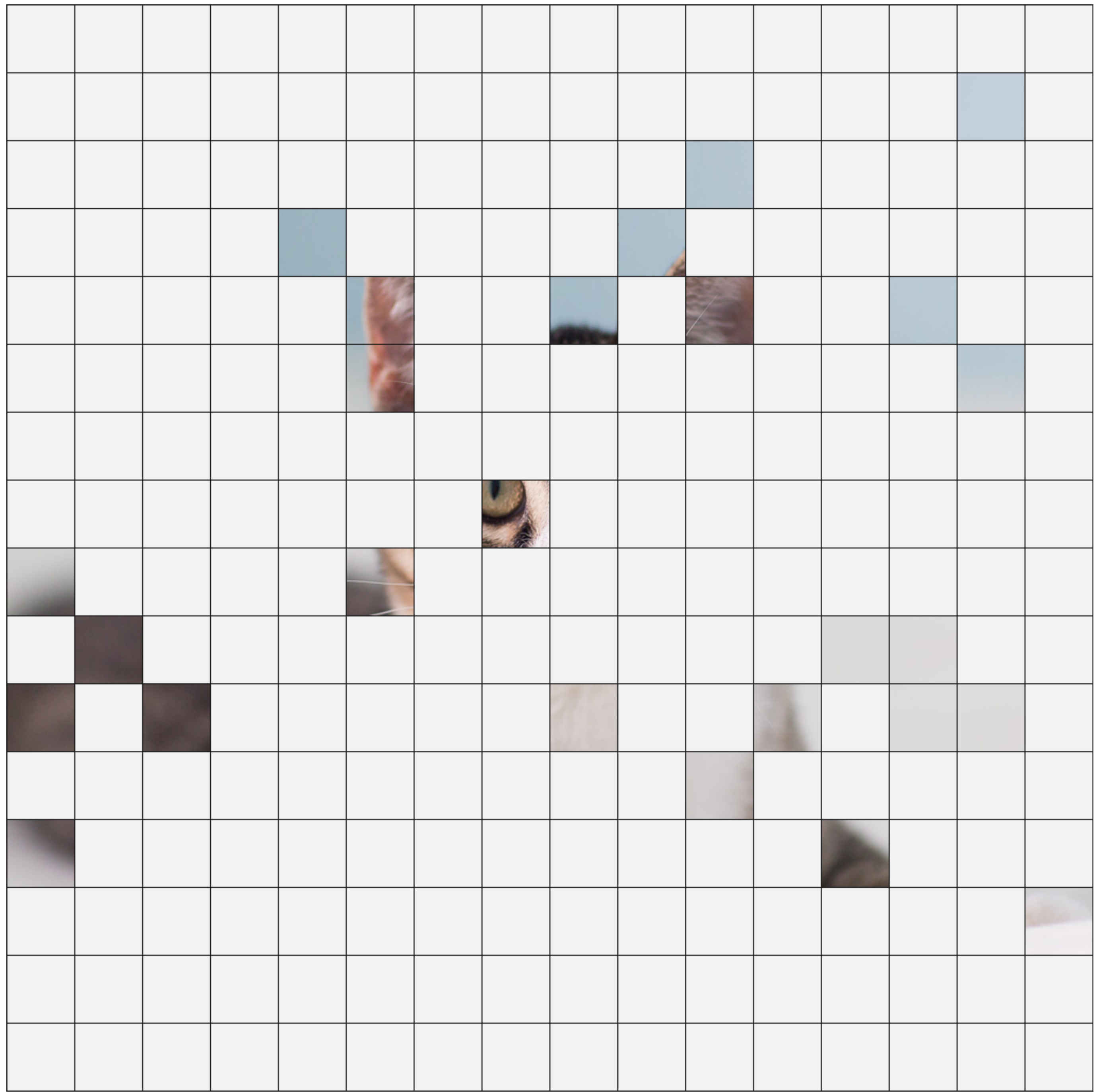


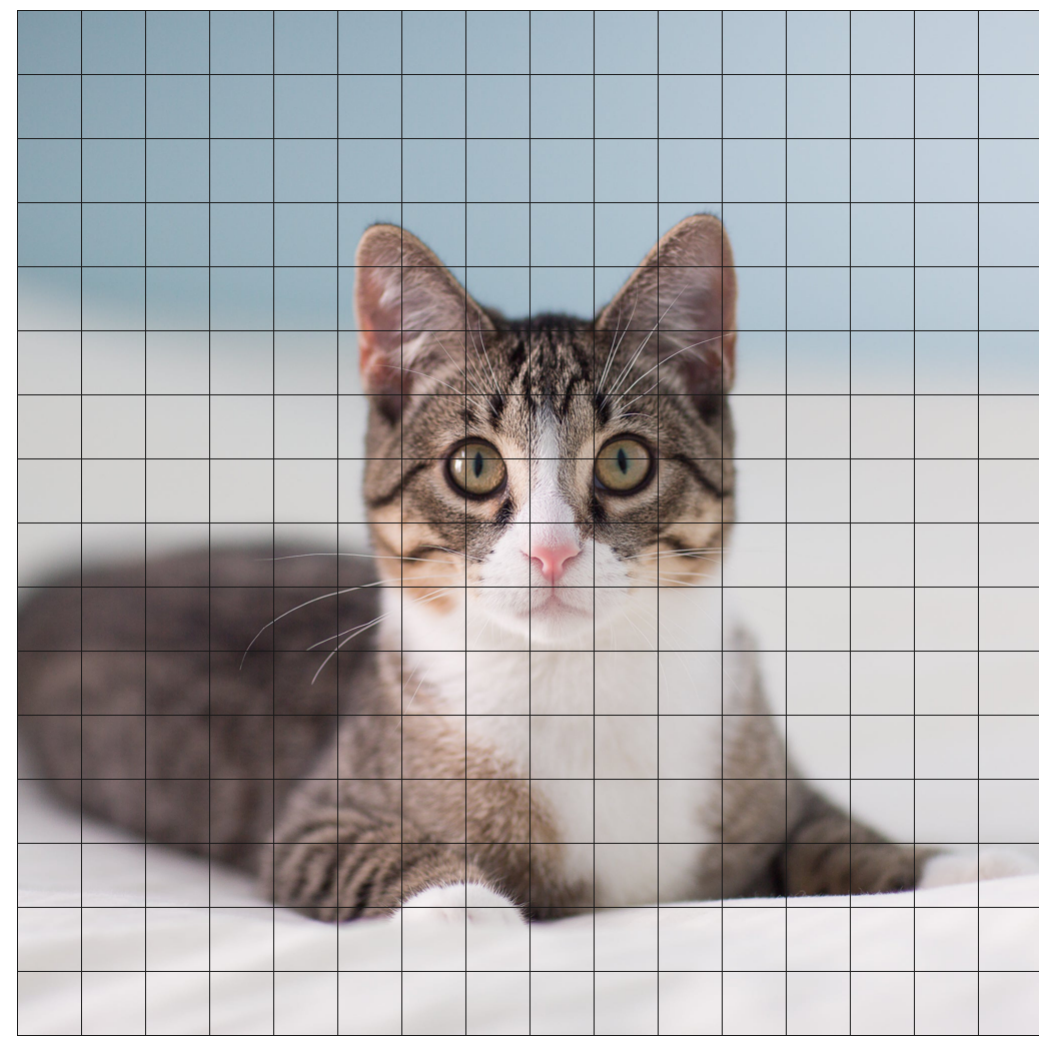
PixelTransformer: Sample Conditioned Signal Generation

Shubham Tulsiani (*Facebook AI Research*)

Abhinav Gupta (*Facebook AI Research, CMU*)

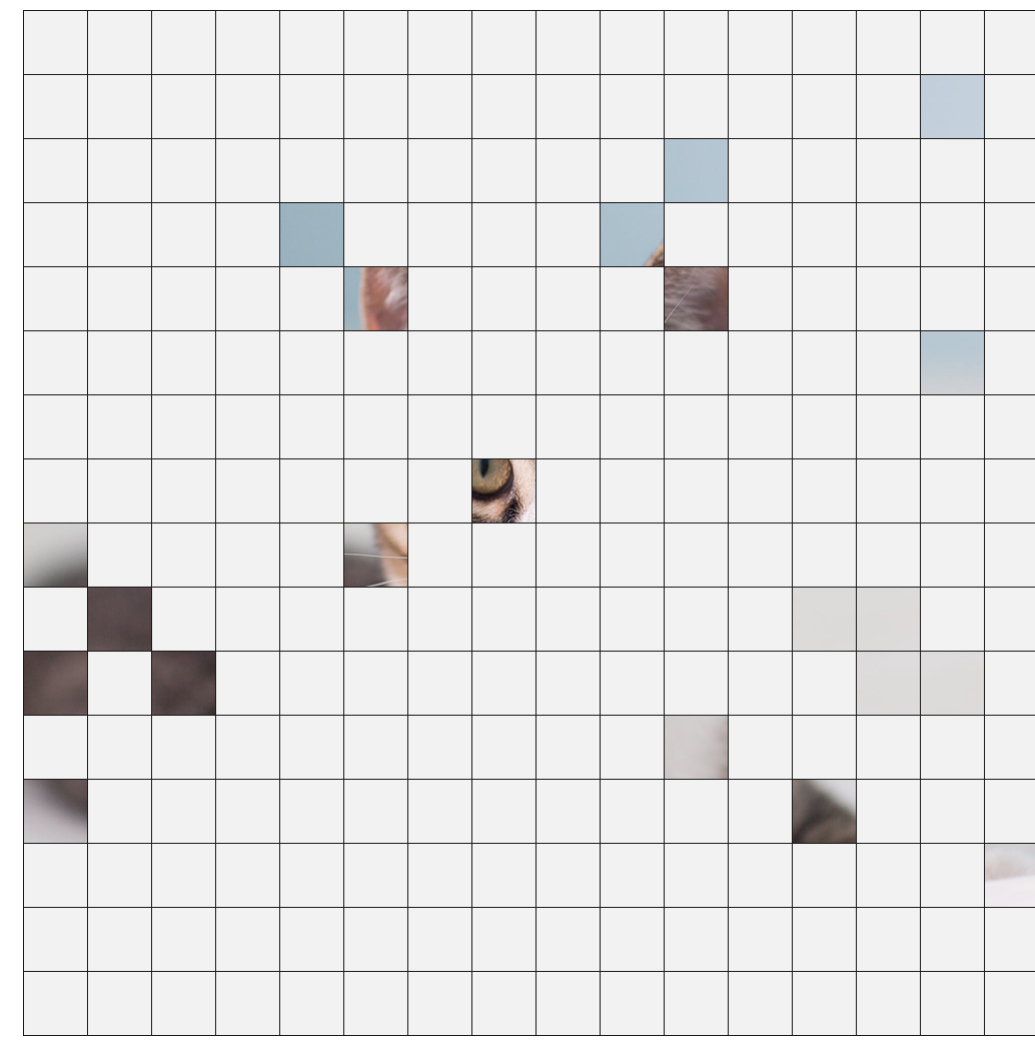


$p($



signal distribution

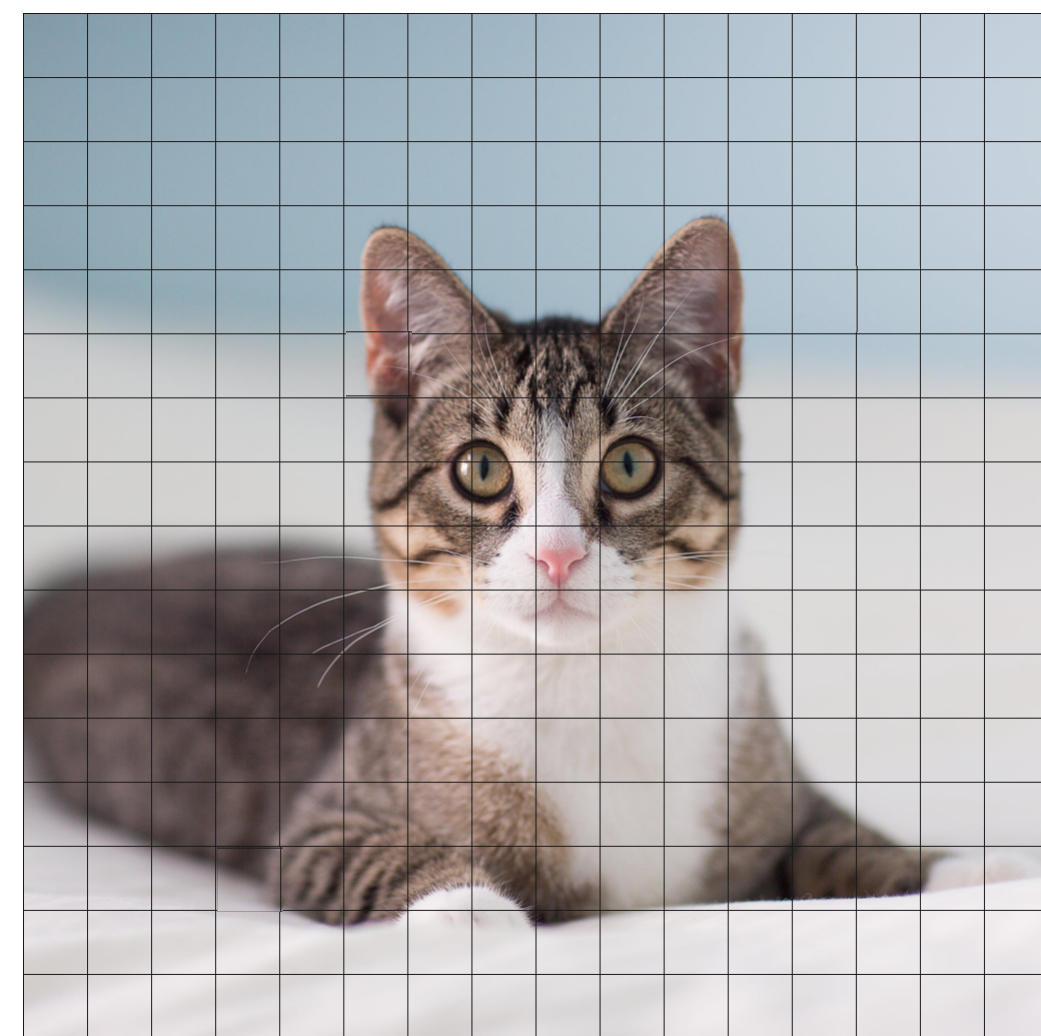
|



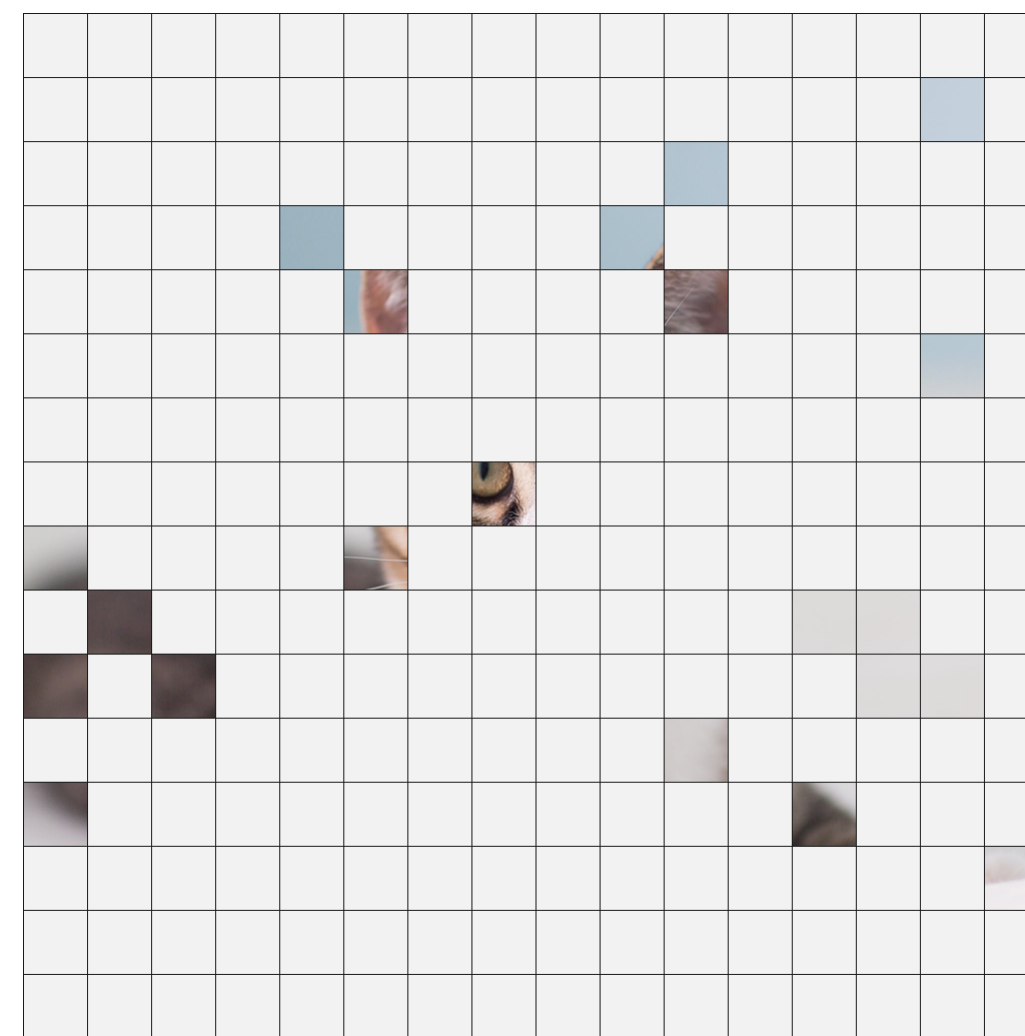
conditioned on **samples**

)

$p($



|

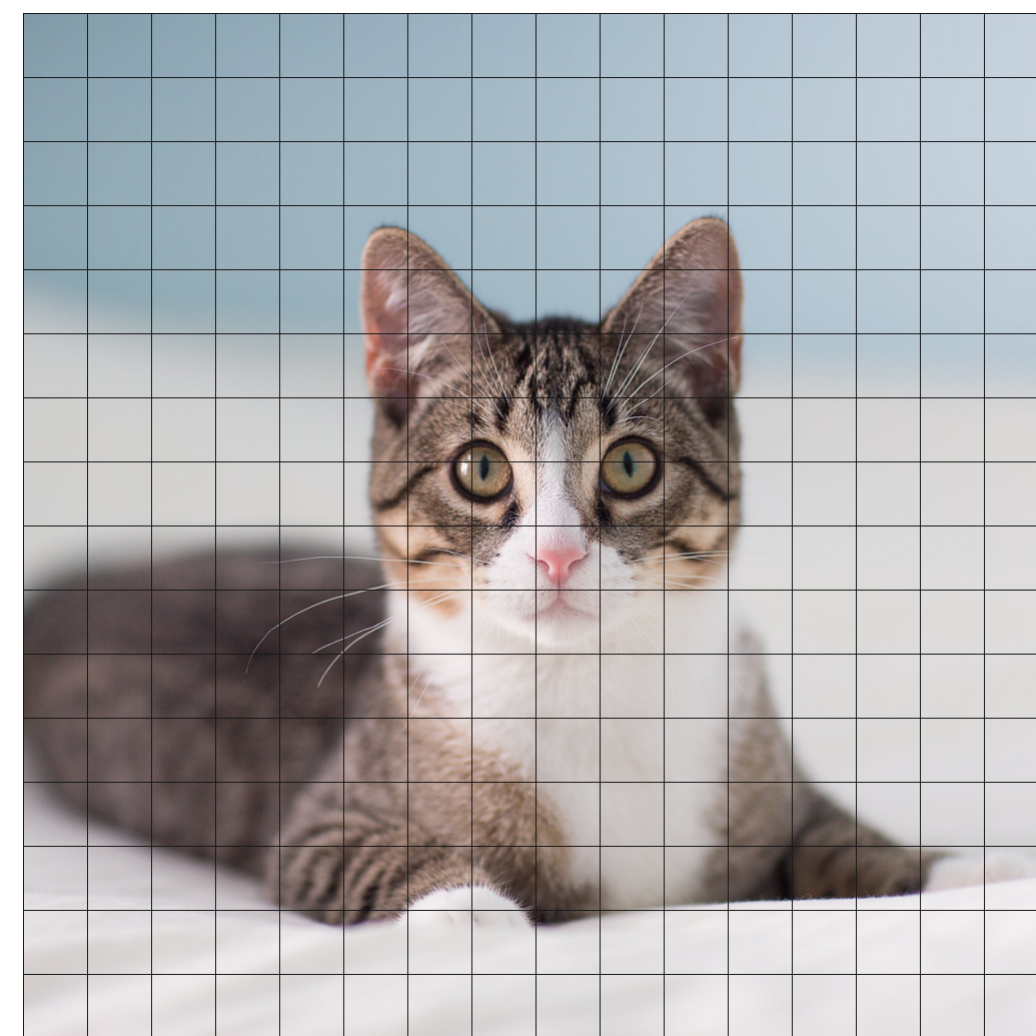


)

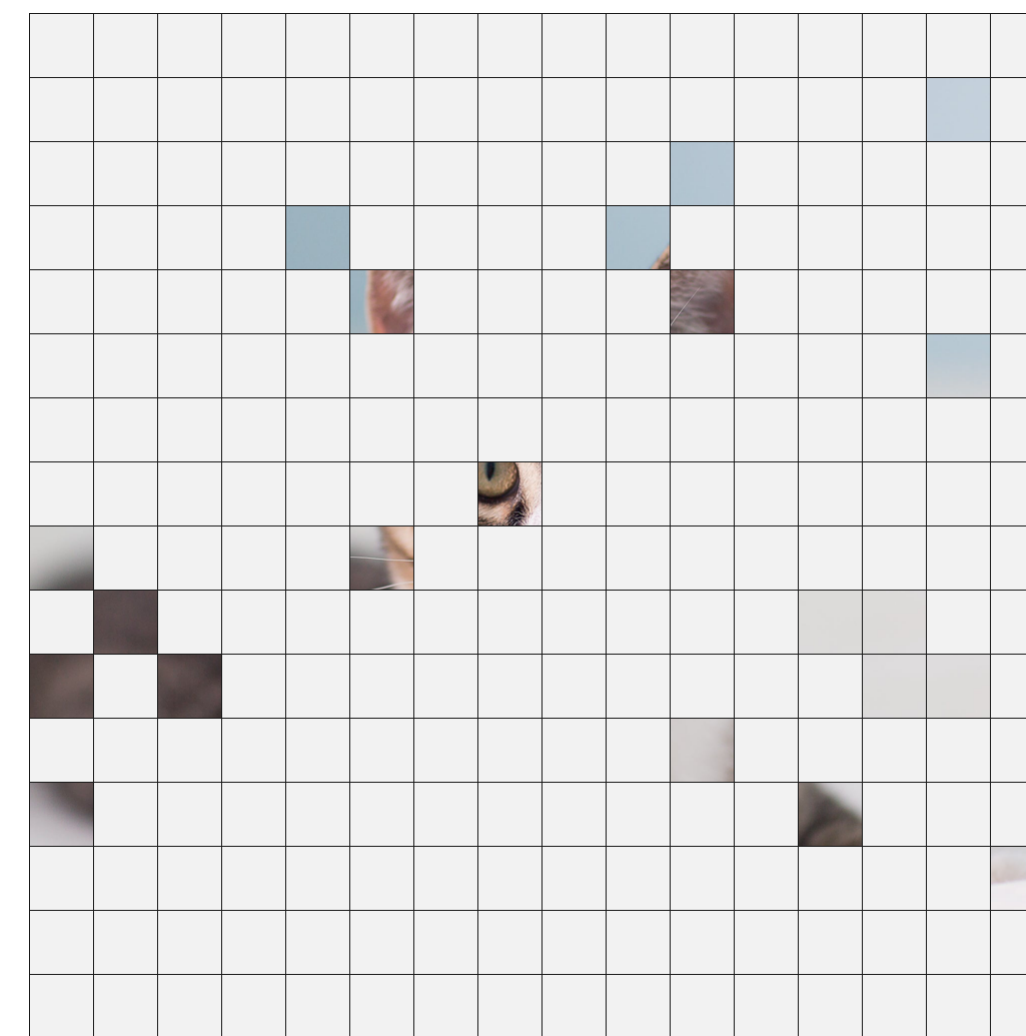
infer:

$p(I|S)$

$p($



$|$

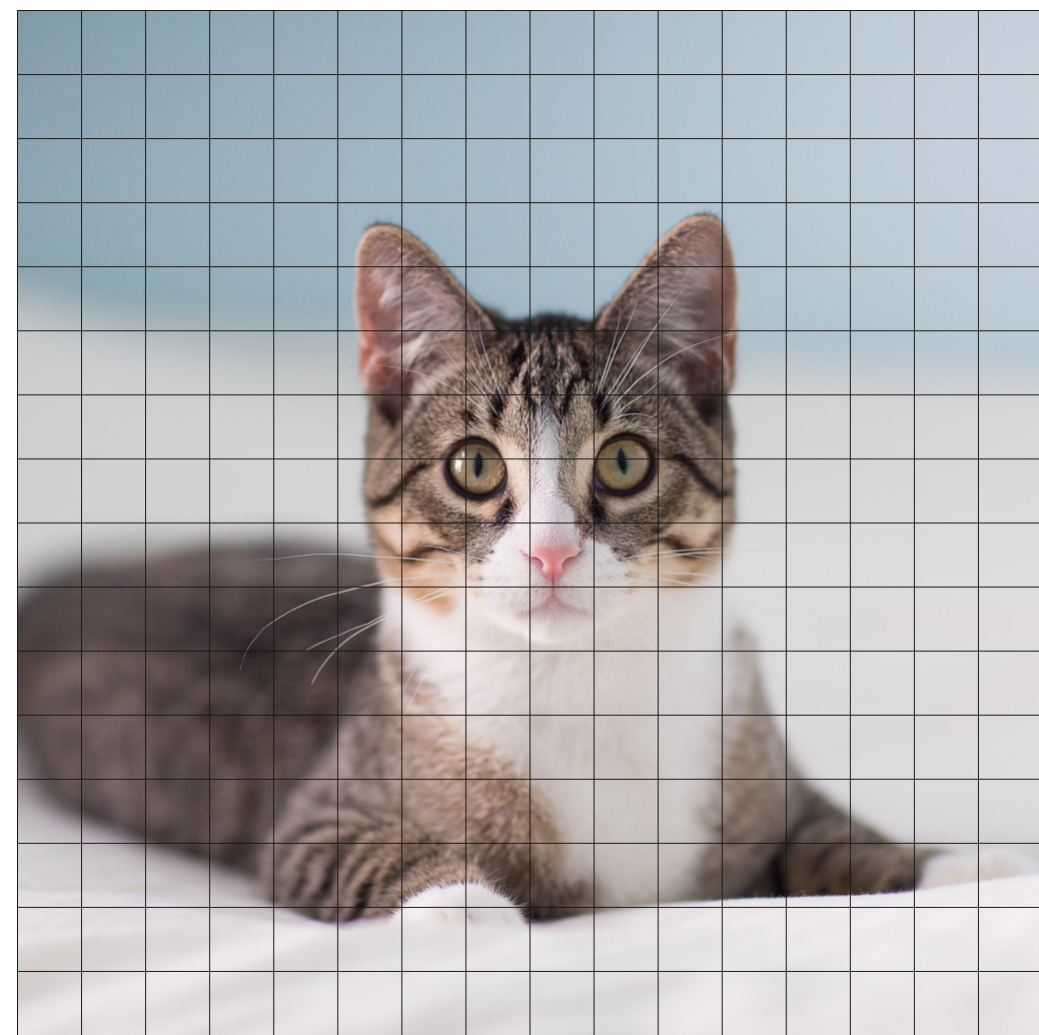


$)$

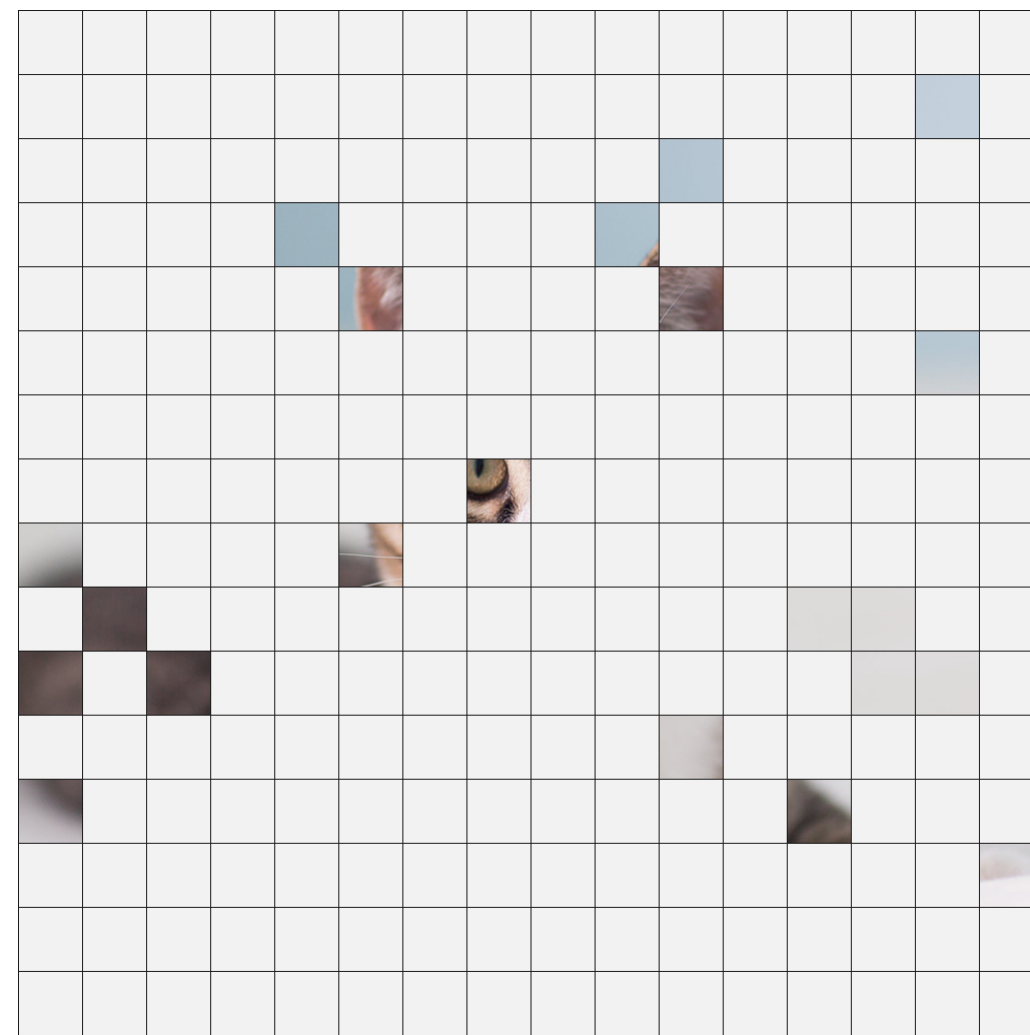
infer:

$$p(I|S) \equiv p(\text{[kitten face]} \text{ [blurred kitten face]} \dots \text{[blue square]} | S)$$

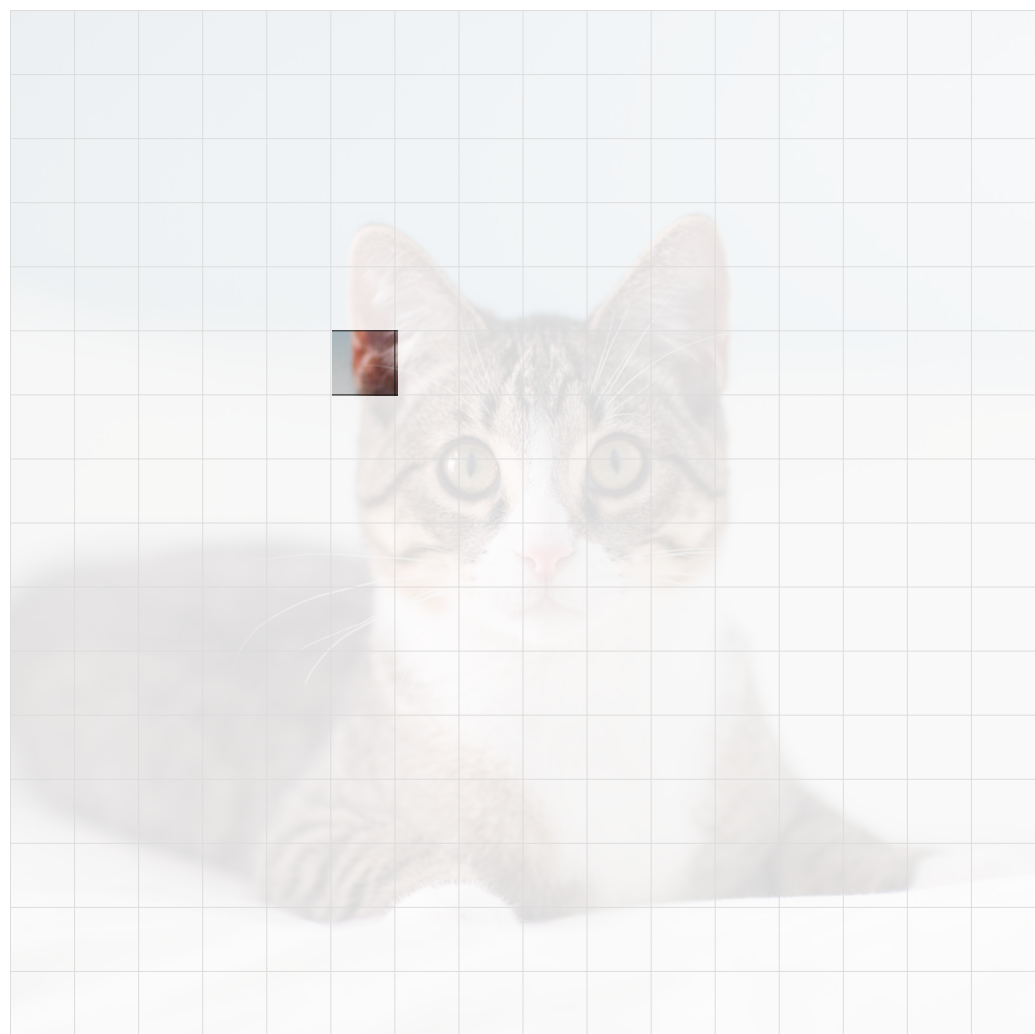
joint distribution over all pixel values

$p($ 

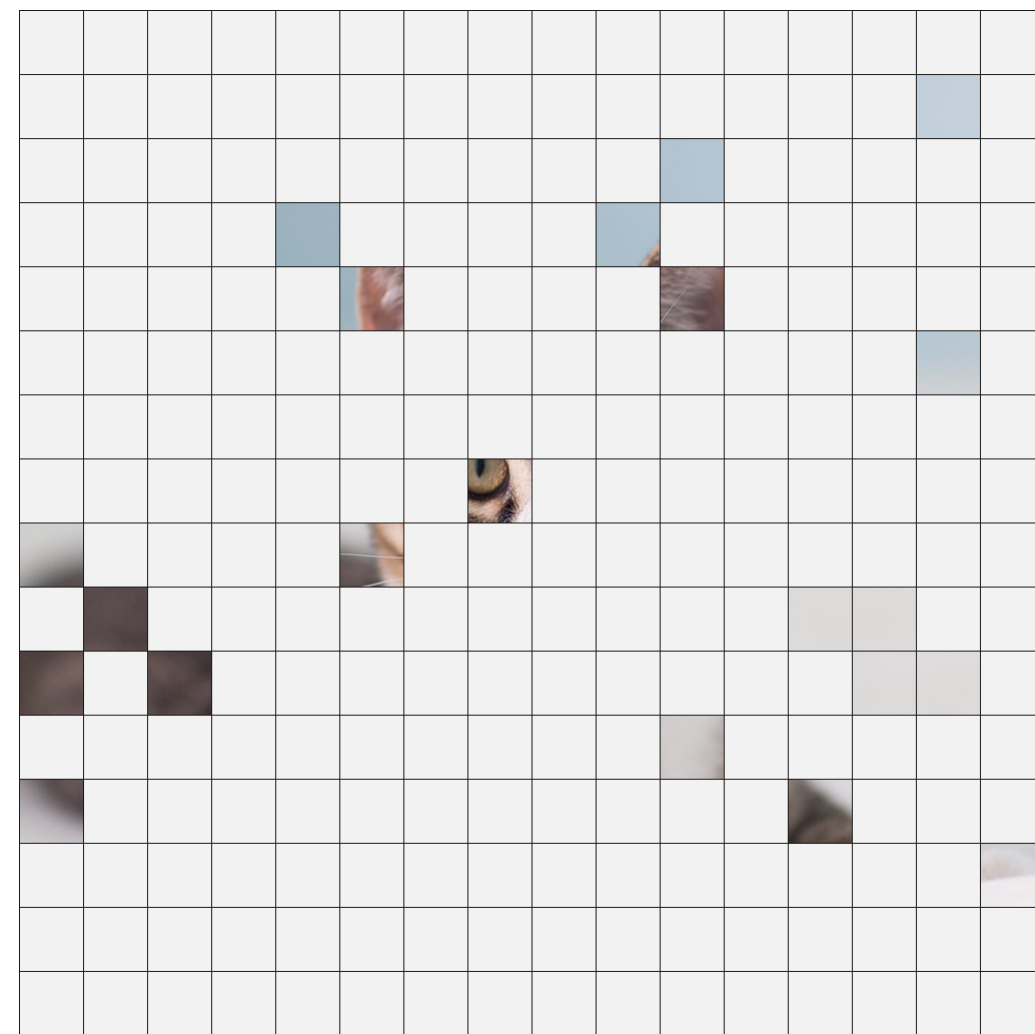
|



)

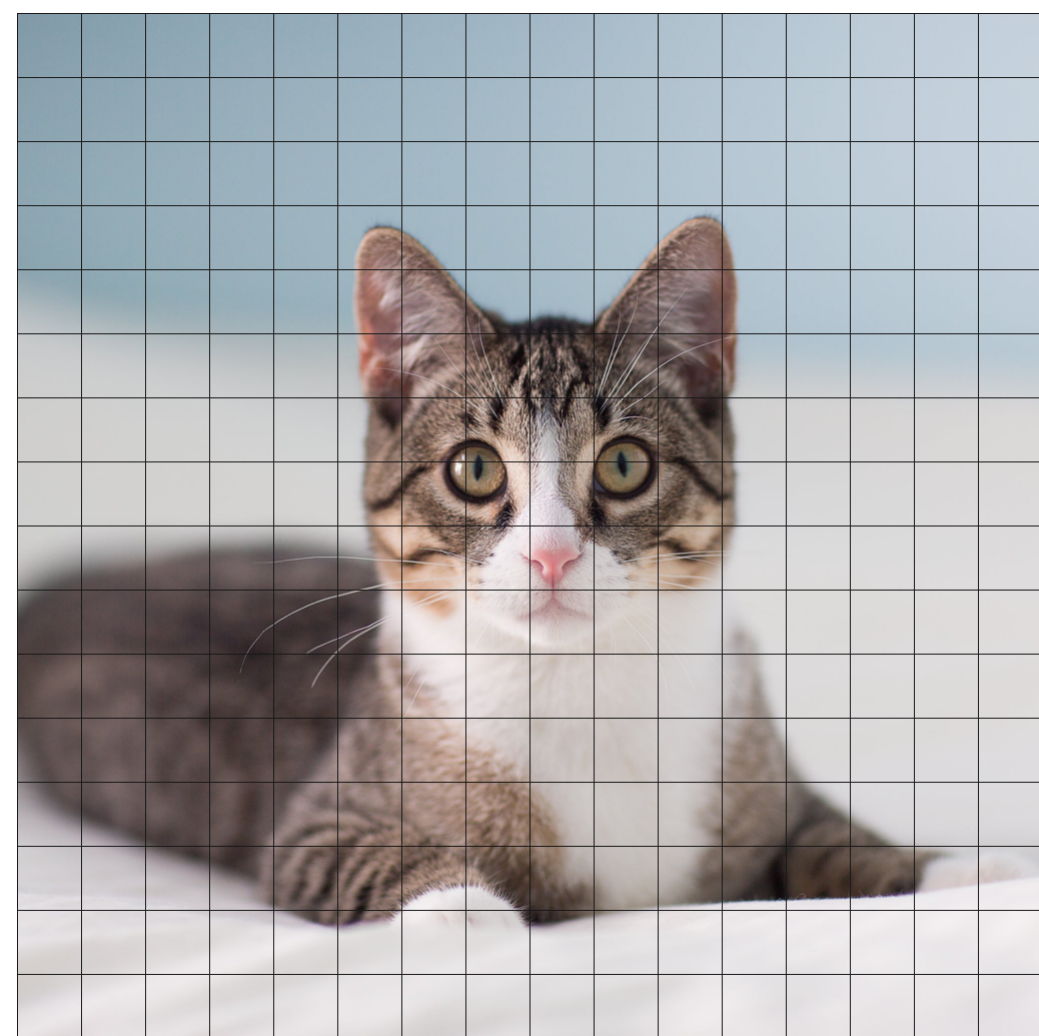
 \equiv $p($ 

|

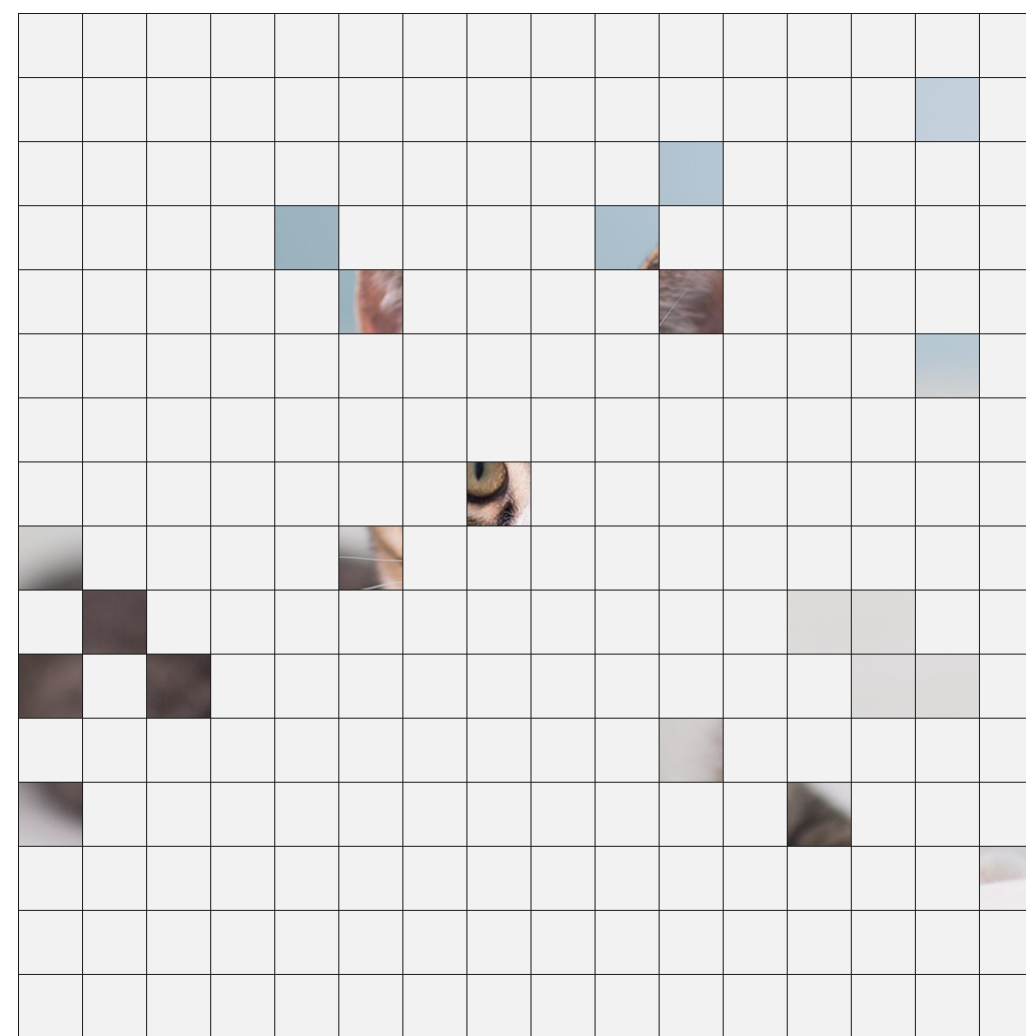


)

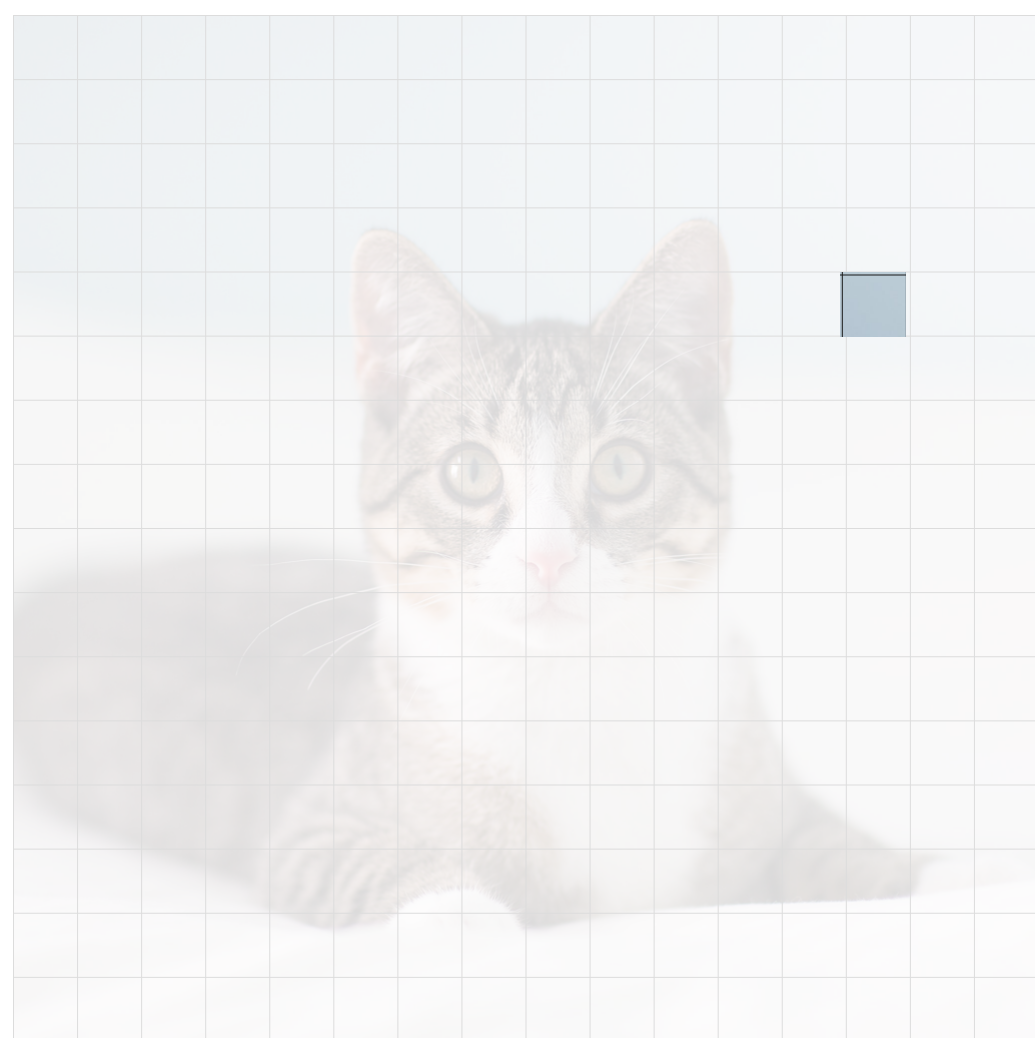
 $p(\mathbf{v}_{g_1} | S)$

$p($ 

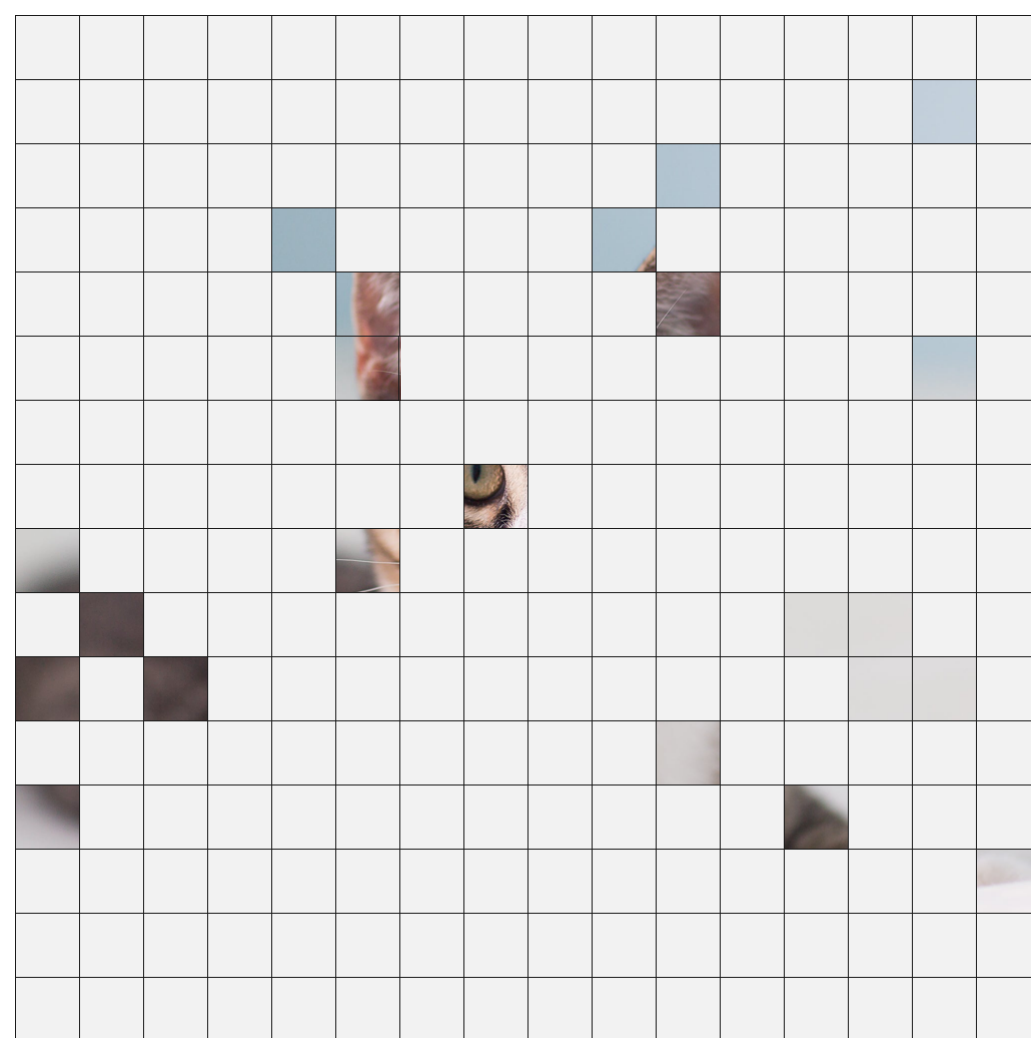
|



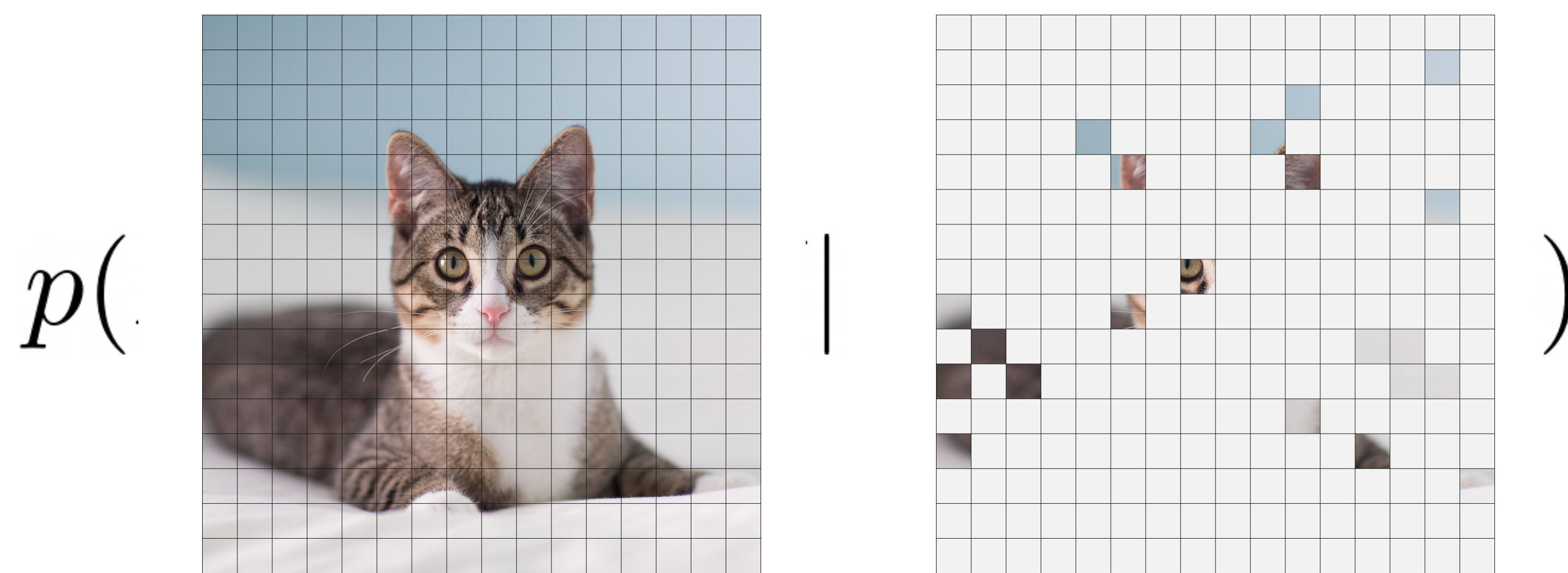
)

 \equiv $p(\mathbf{v}_{g_1} | S)$ $p($ 

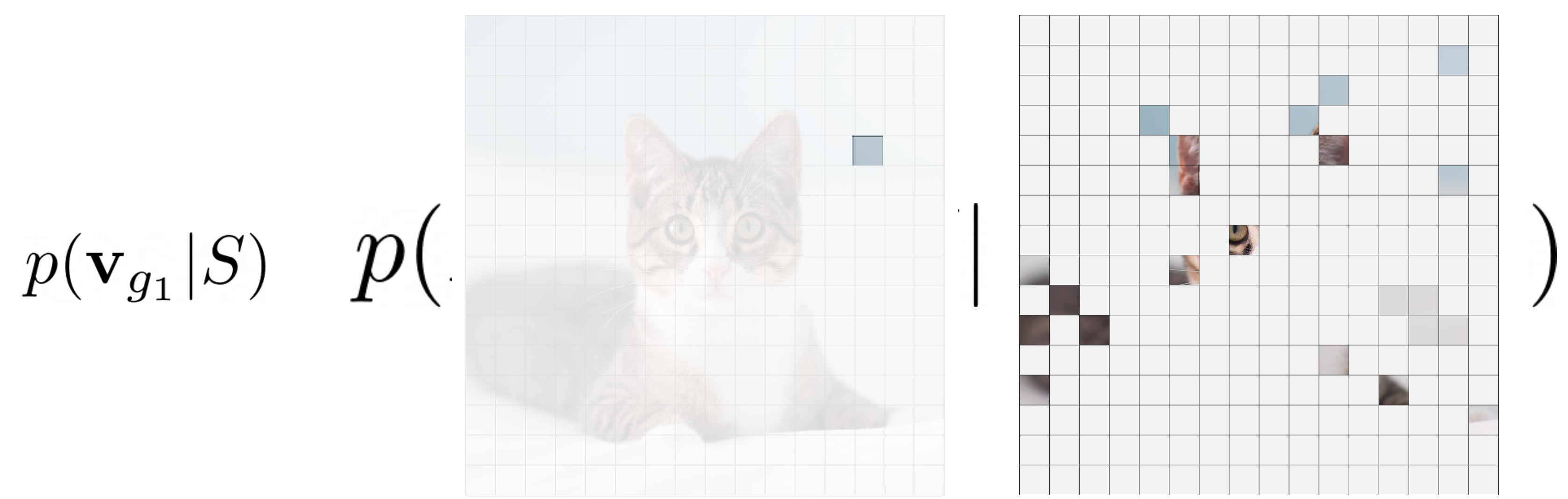
|



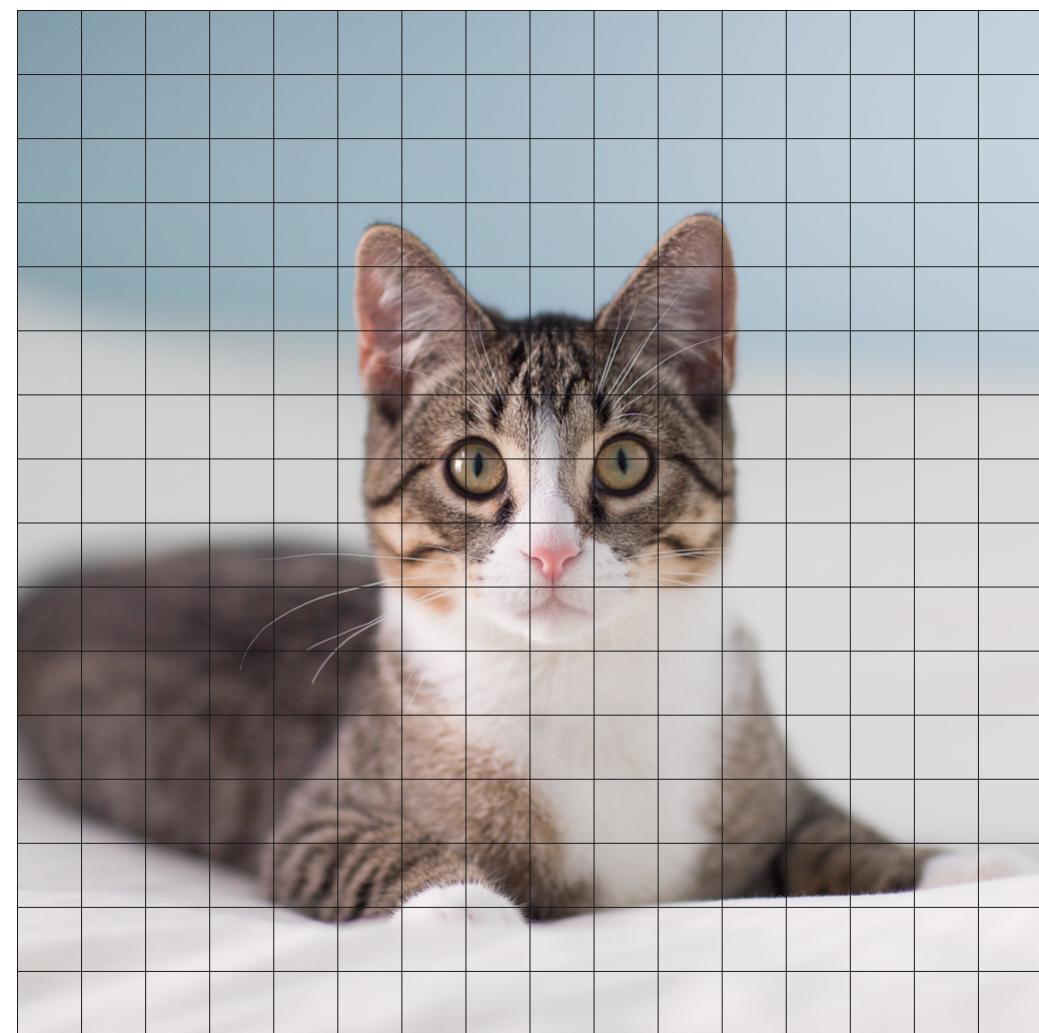
)



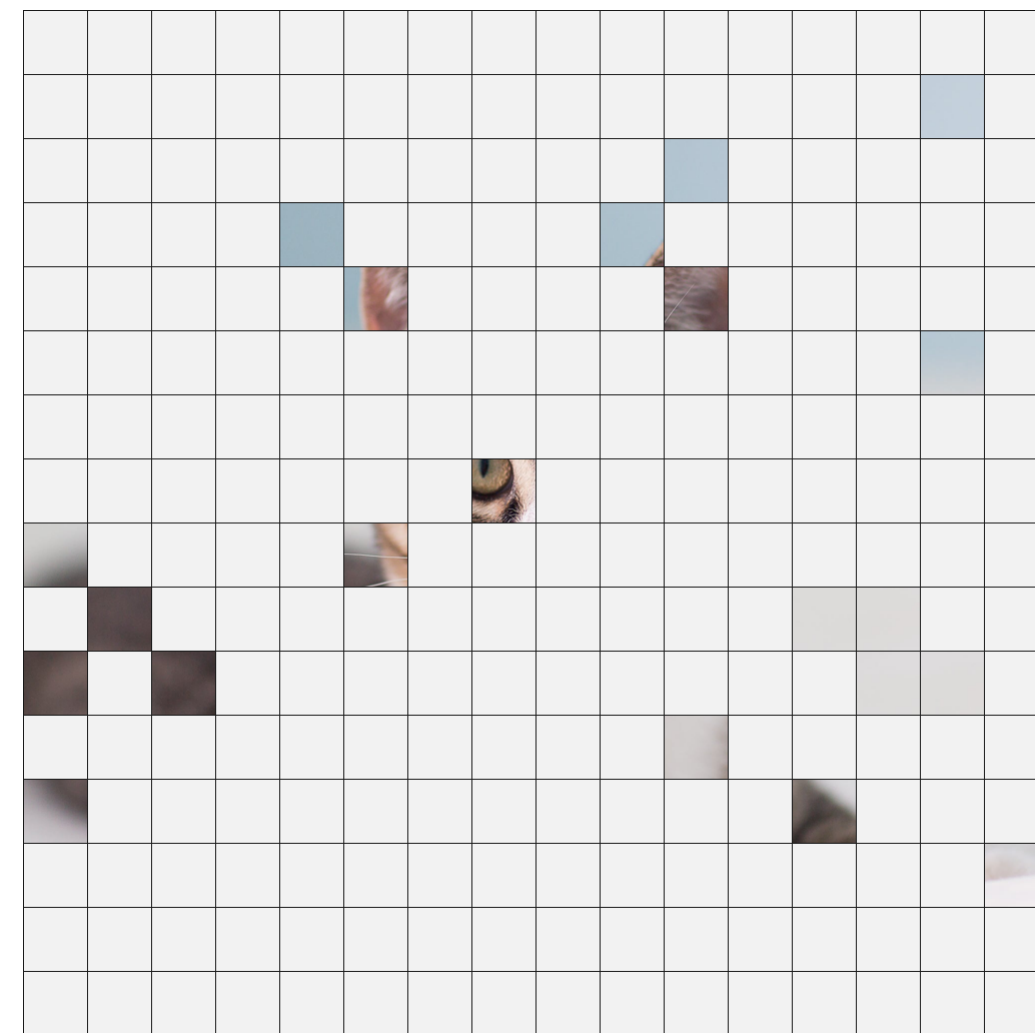
\equiv



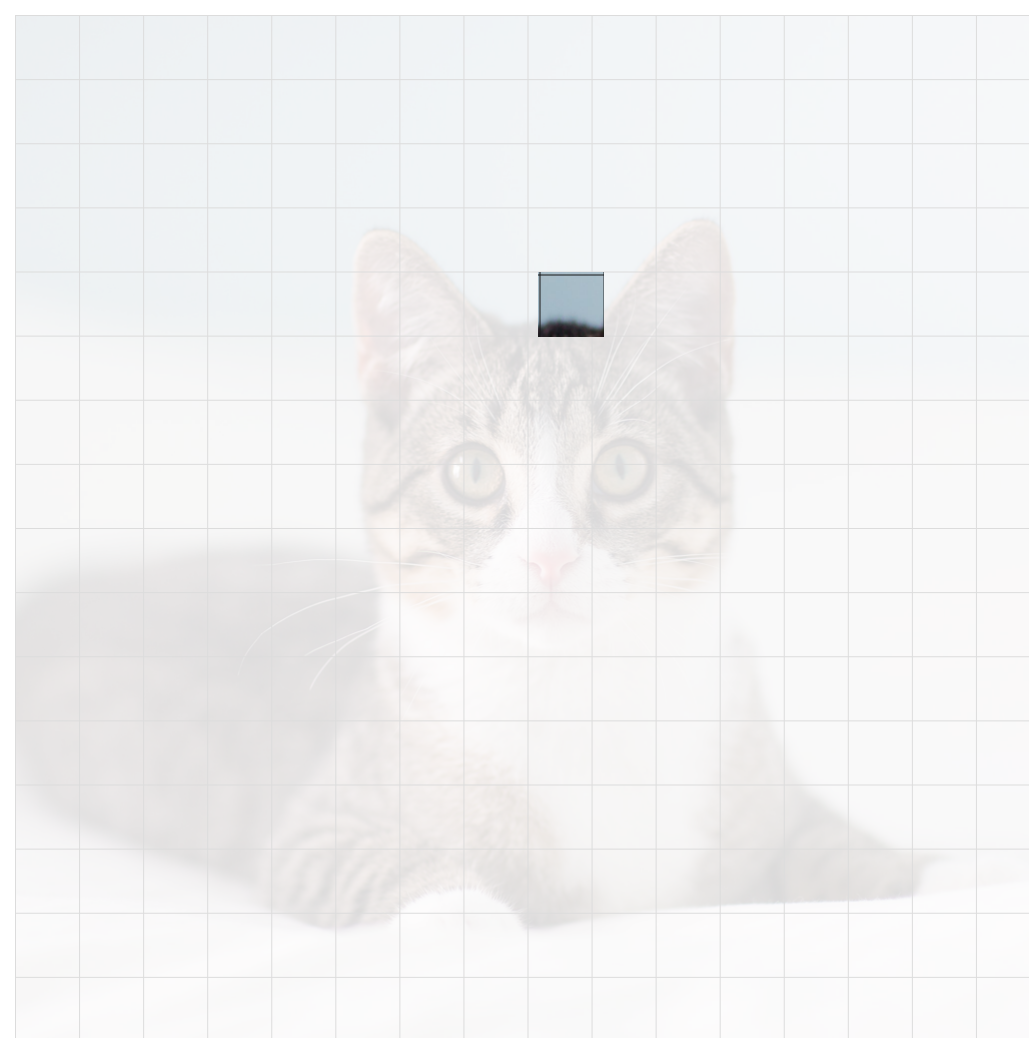
$p(\mathbf{v}_{g_2} \mid S_1)$

$p($ 

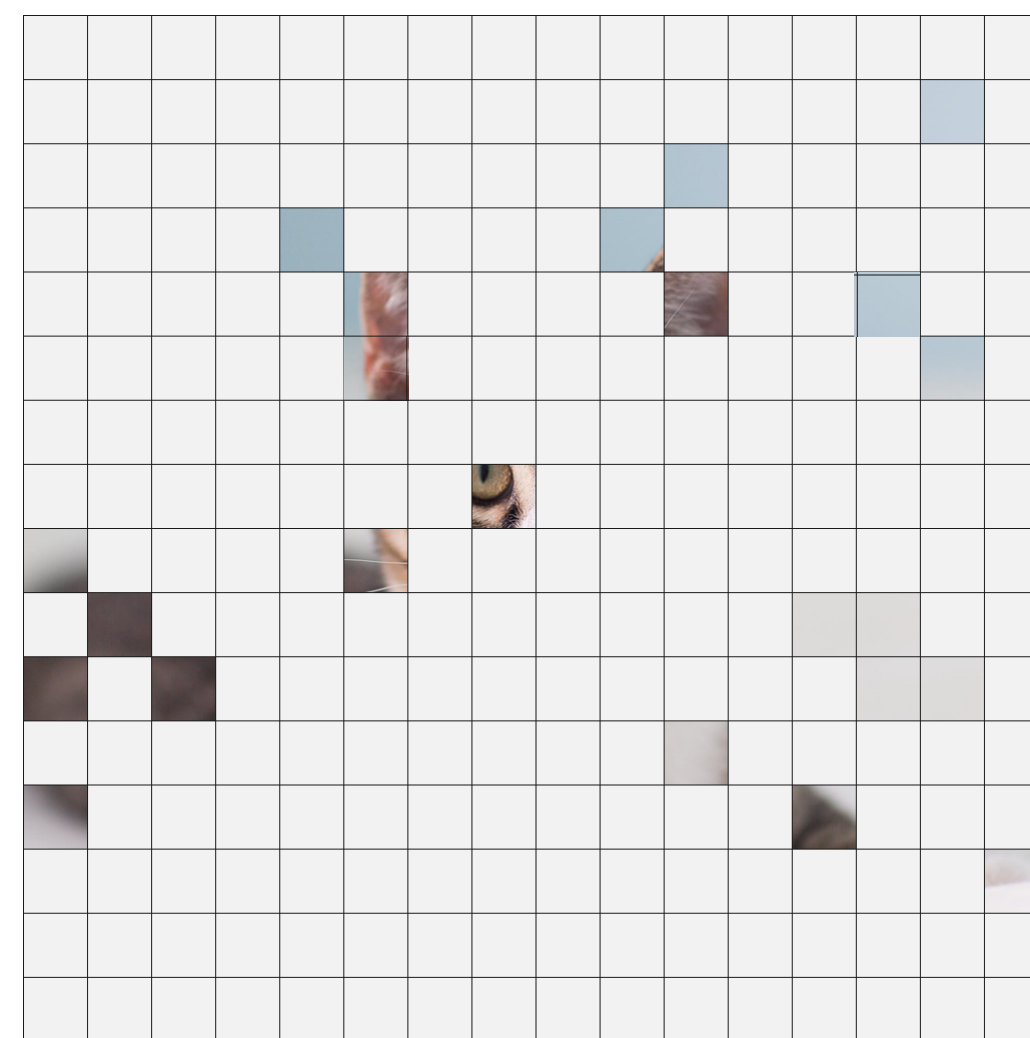
|



)

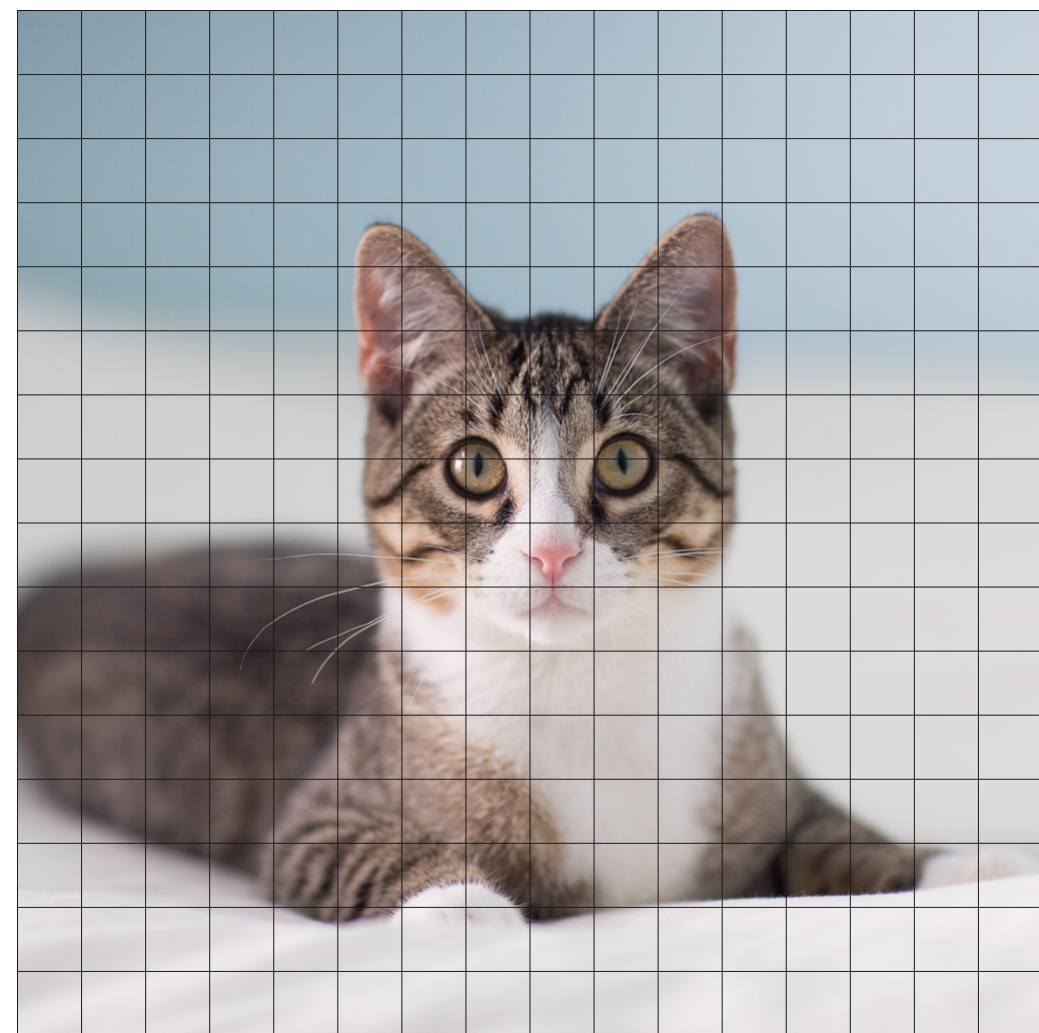
 \equiv $p(\mathbf{v}_{g_1} | S)$ $p(\mathbf{v}_{g_2} | S_1)$ $p($ 

|

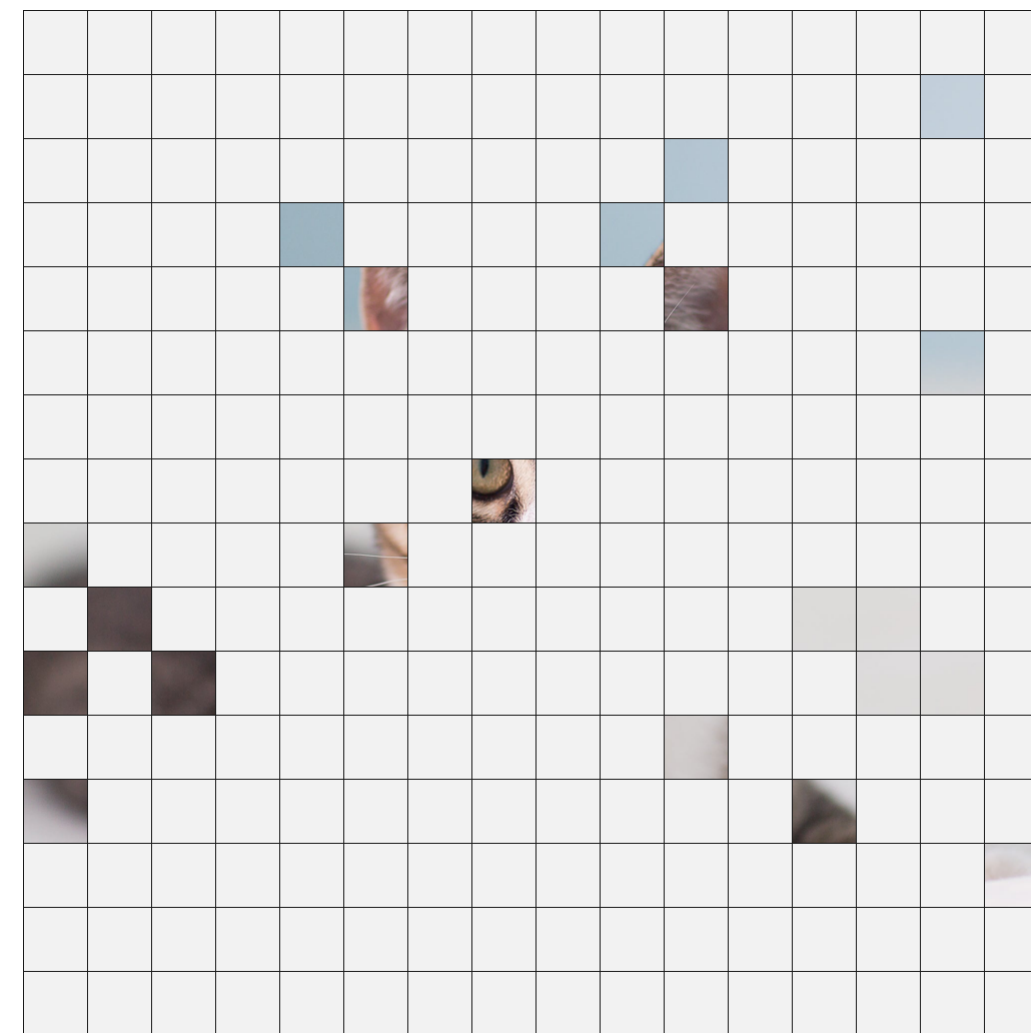


)

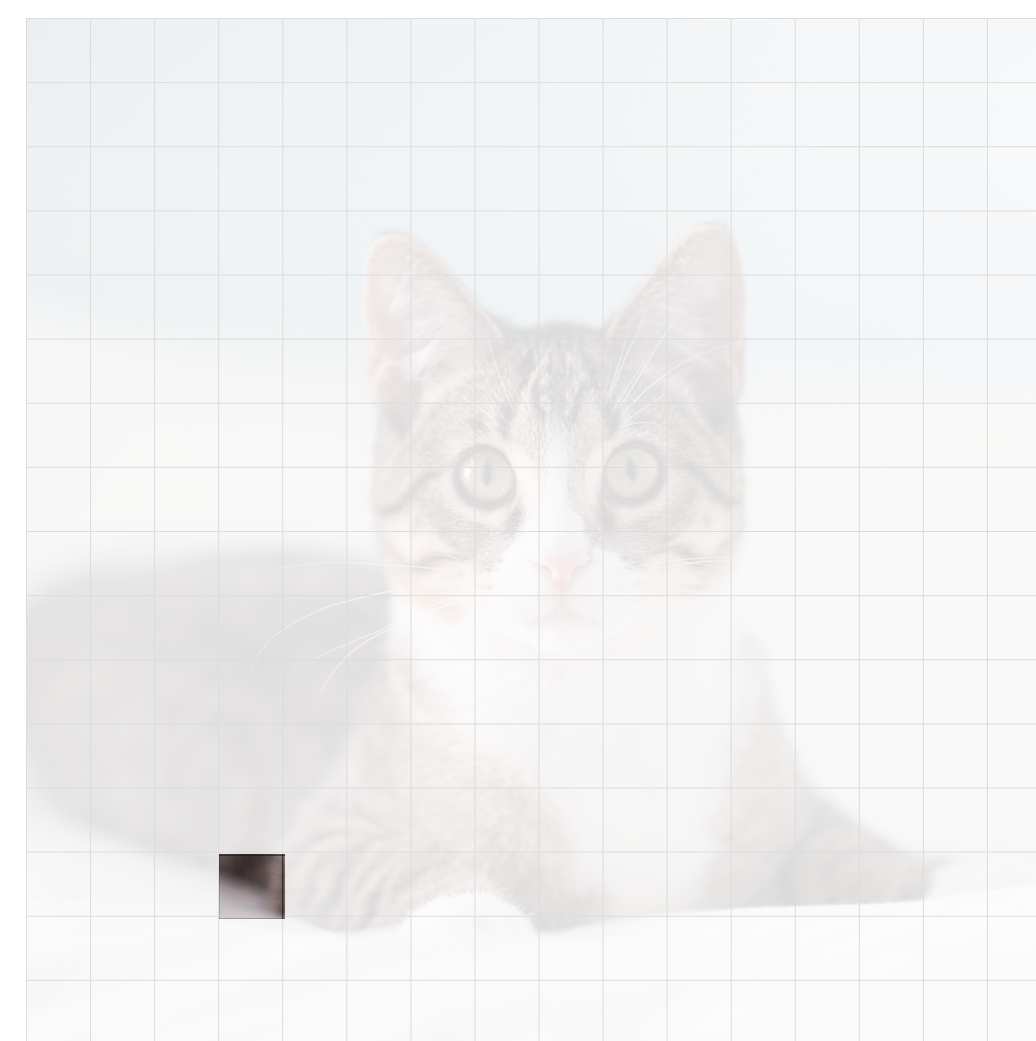
 $p(\mathbf{v}_{g_3} | S_2)$

$p($ 

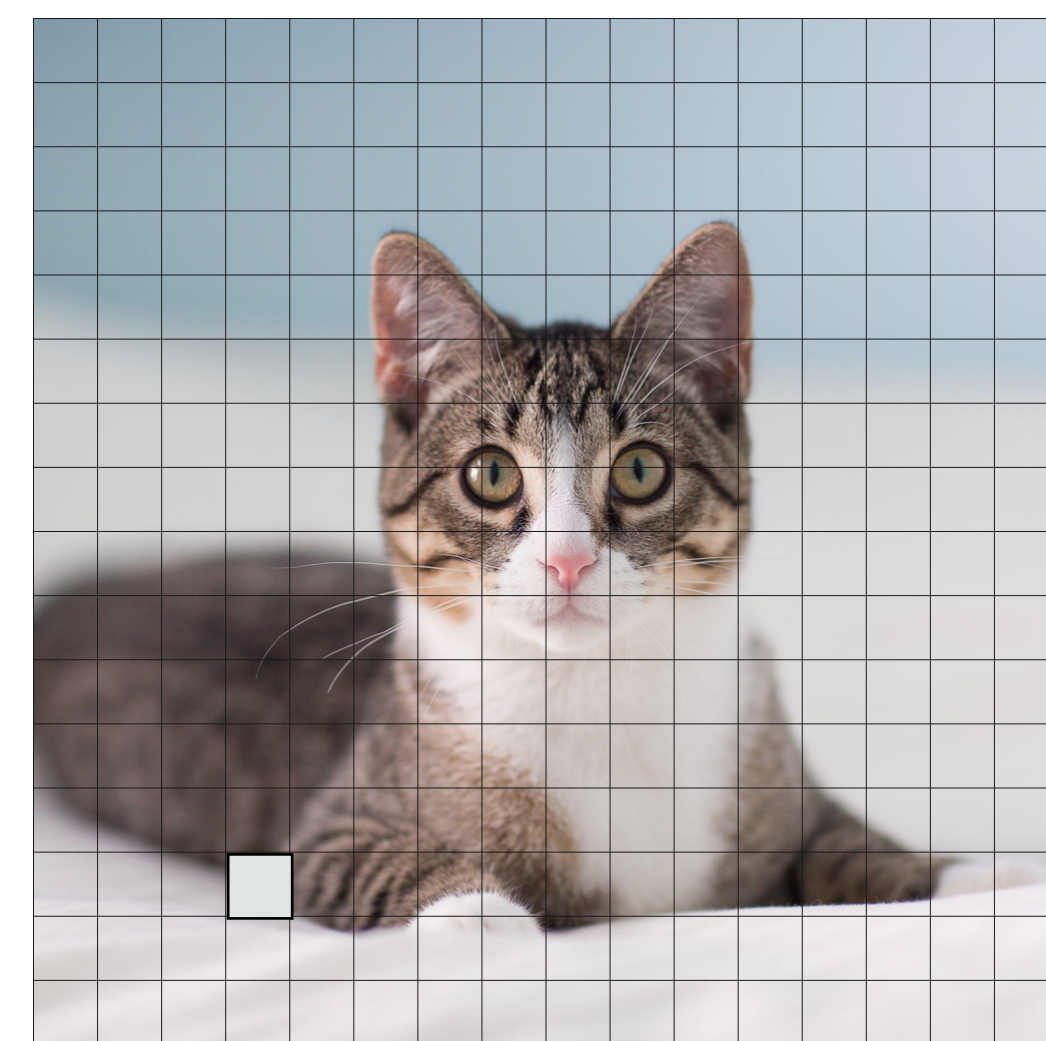
|



)

 \equiv $p(\mathbf{v}_{g_1} | S)$ $p(\mathbf{v}_{g_2} | S_1)$ $p(\mathbf{v}_{g_3} | S_2)$ $\dots p($ 

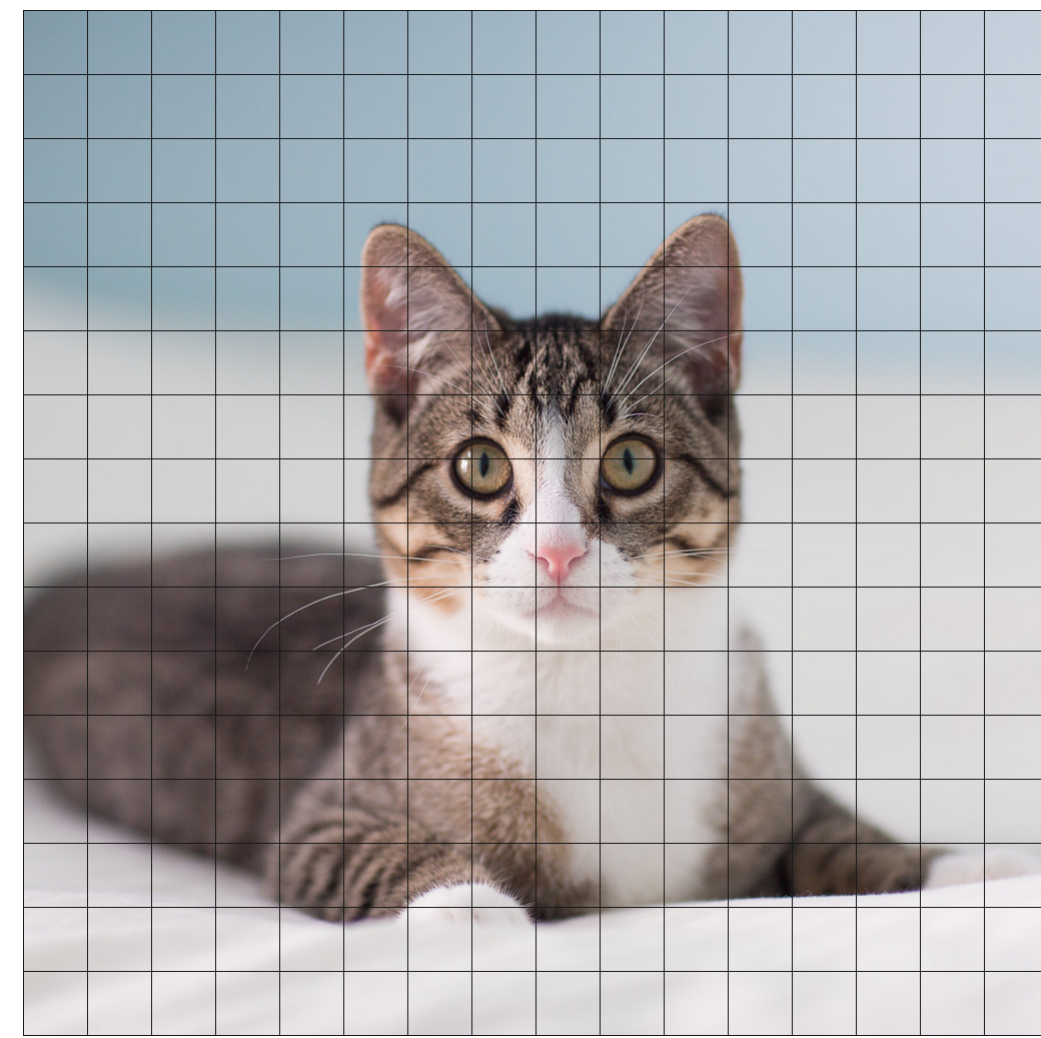
|



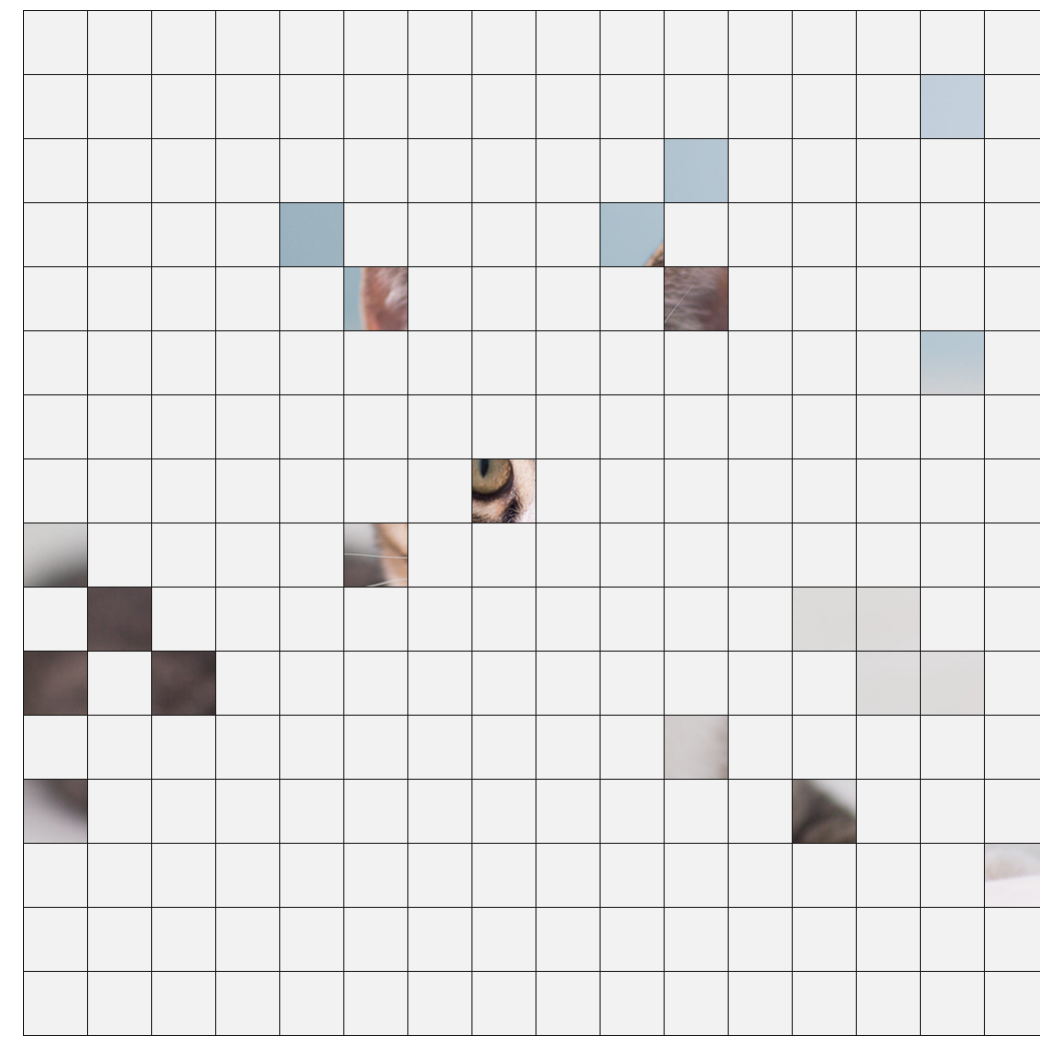
)

 $p(\mathbf{v}_{g_N} | S_{N-1})$

$p($



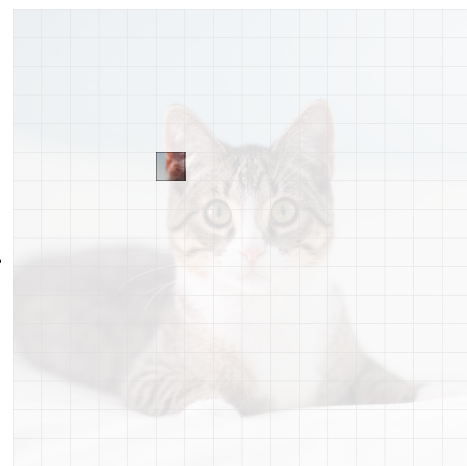
|



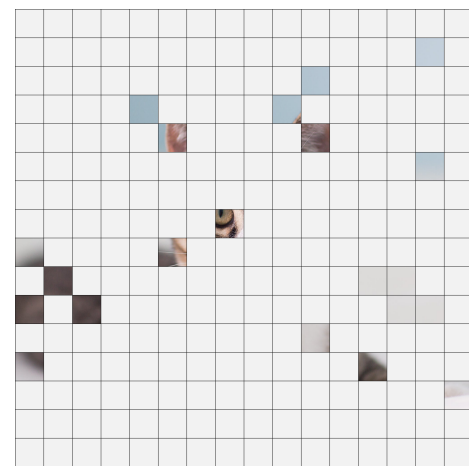
)

\equiv

$p($

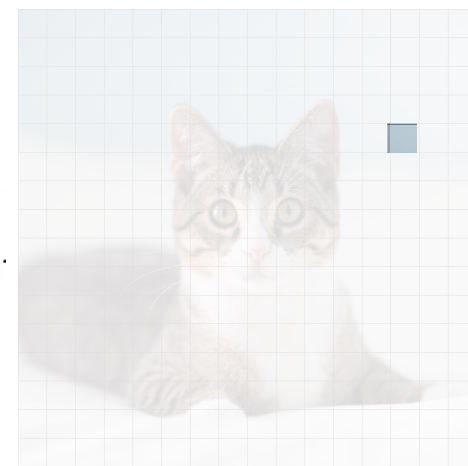


|

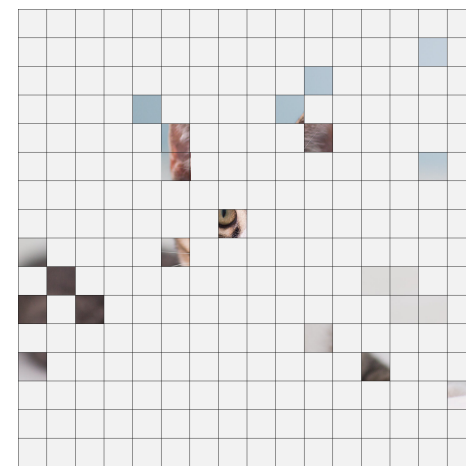


)

$p($

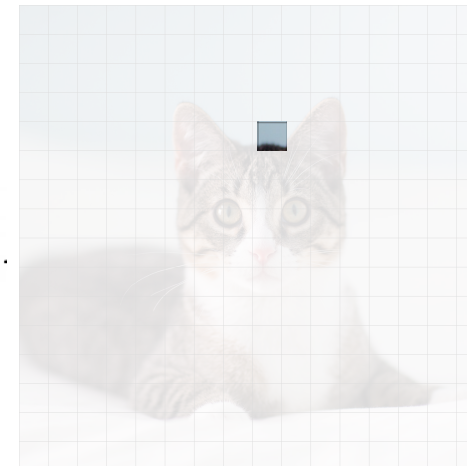


|

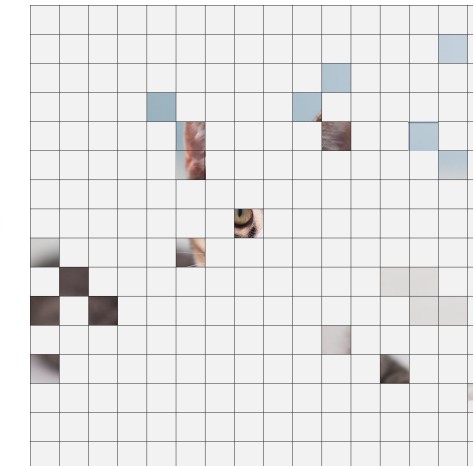


)

$p($

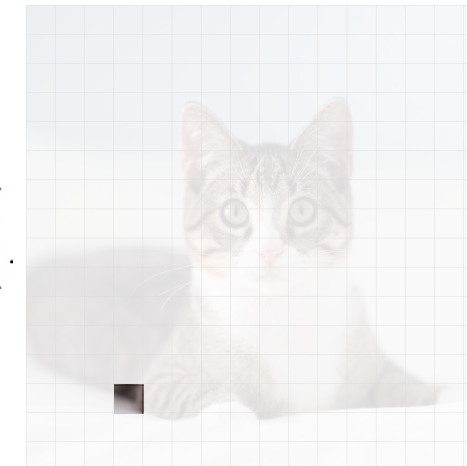


|

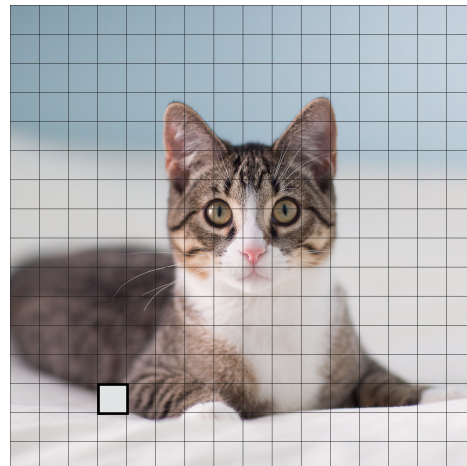


)

$\dots p($

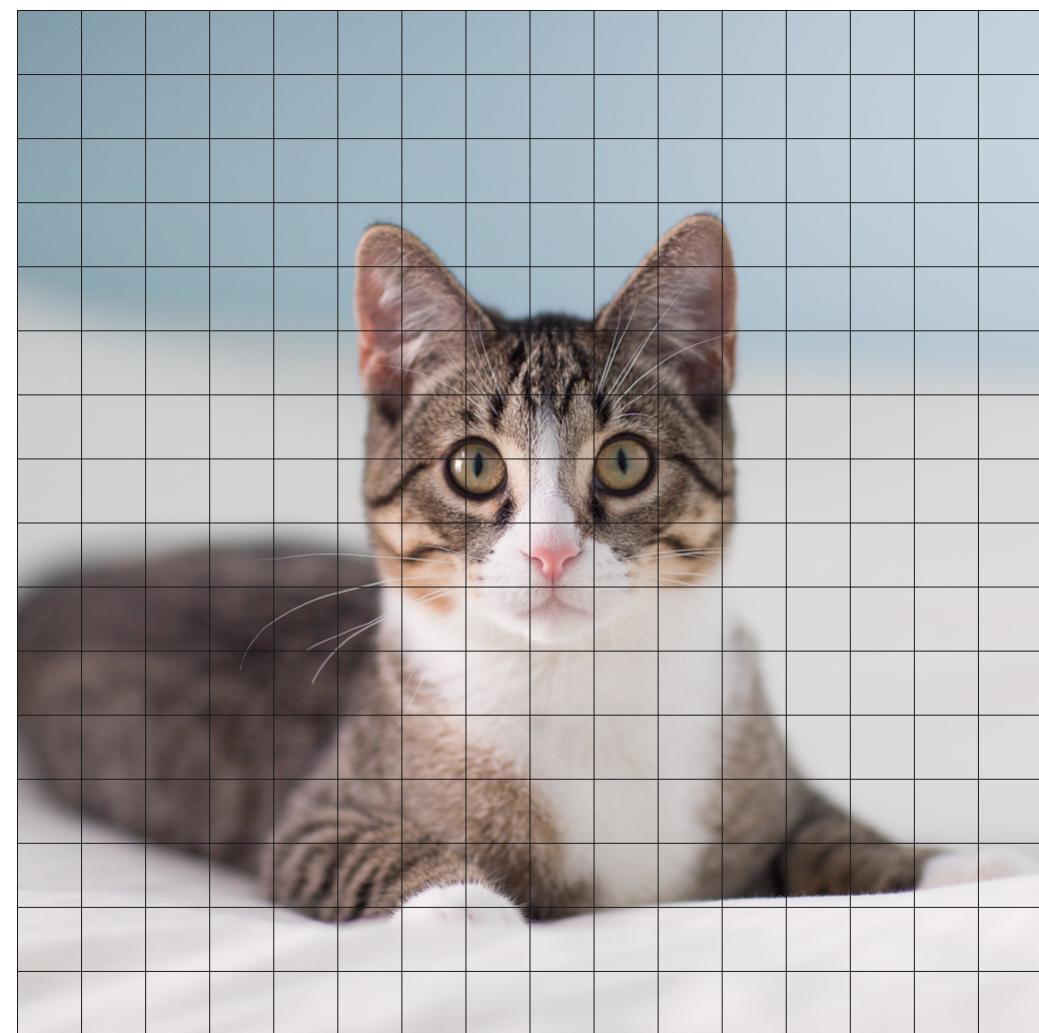


|

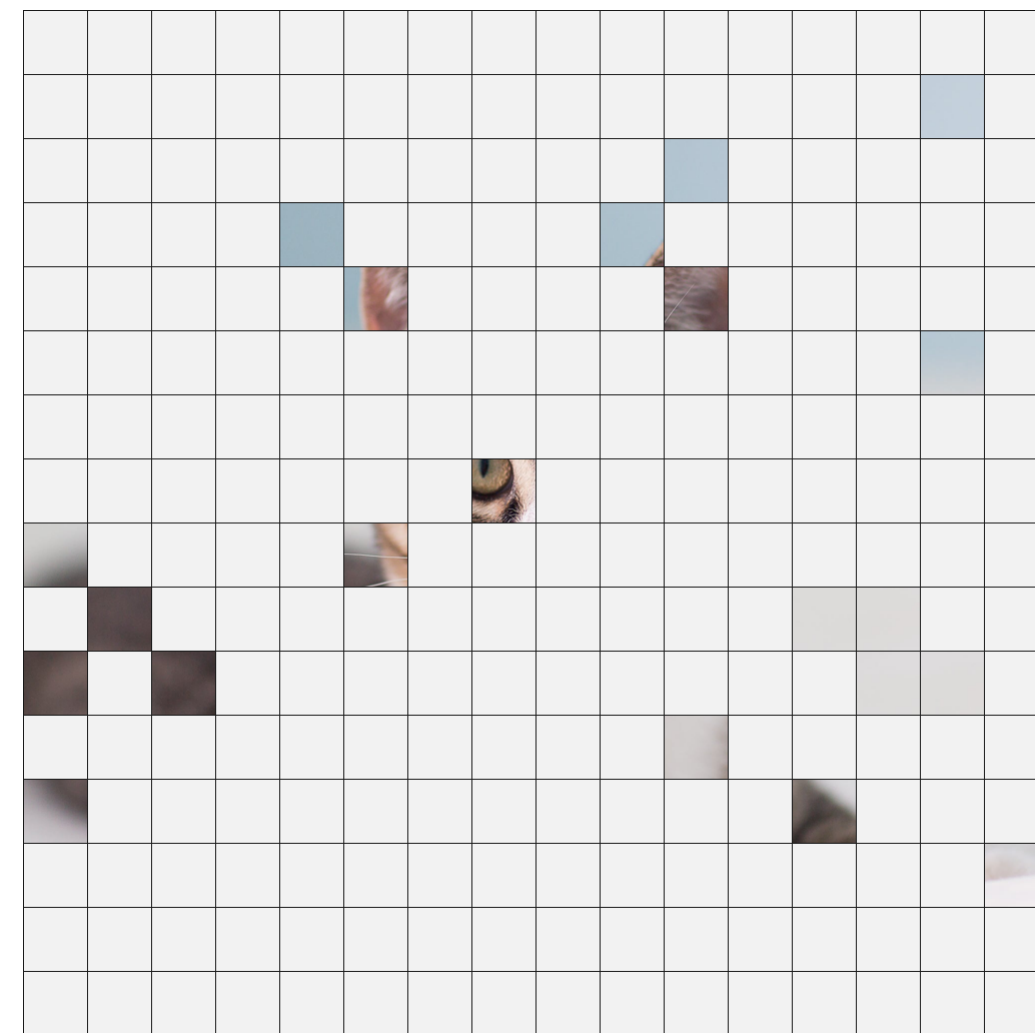


)

$p($



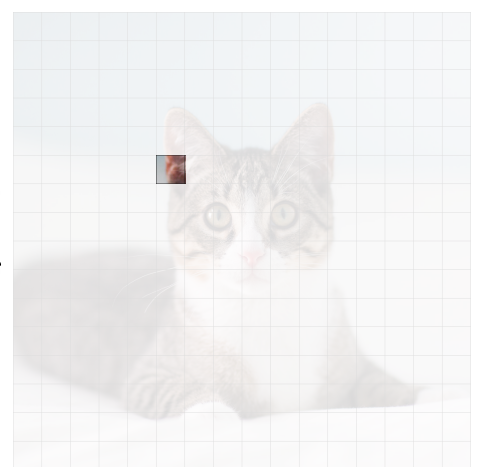
$|$



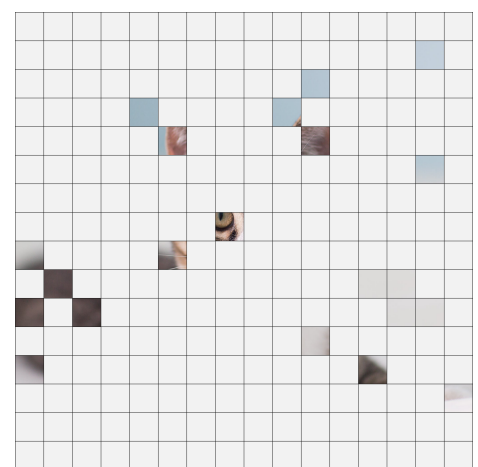
$)$

\equiv

$p($

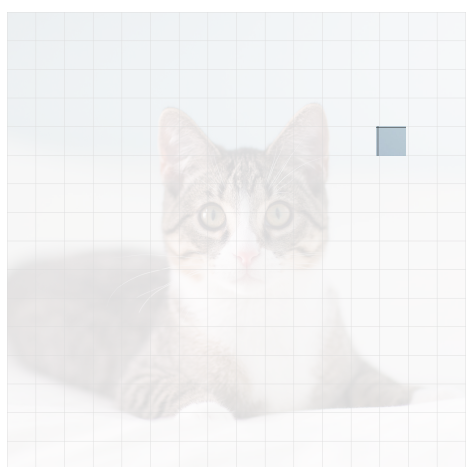


$|$

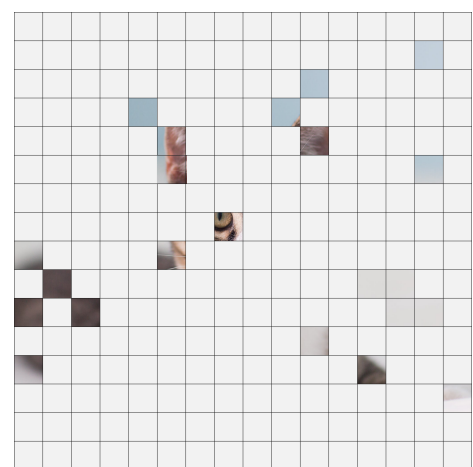


$)$

$p($

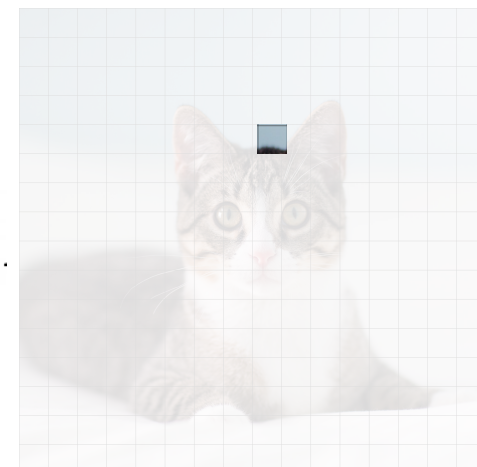


$|$

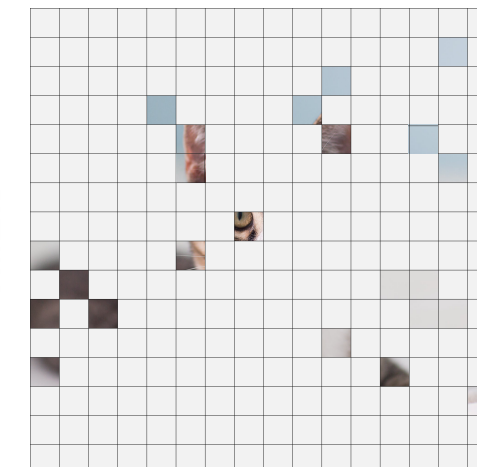


$)$

$p($



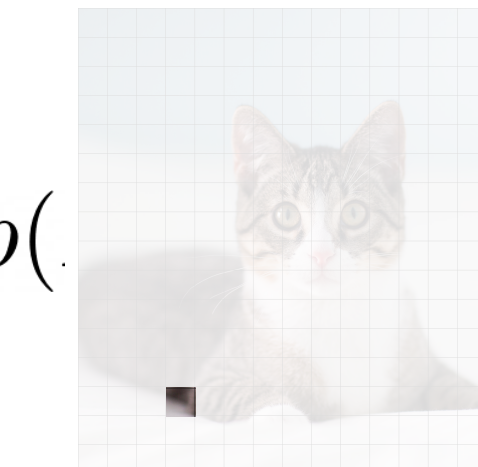
$|$



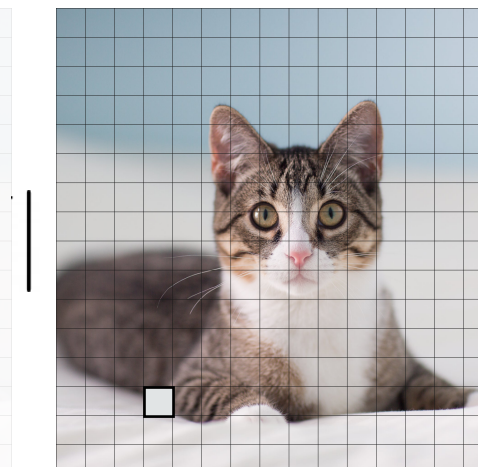
$)$

\dots

$p($



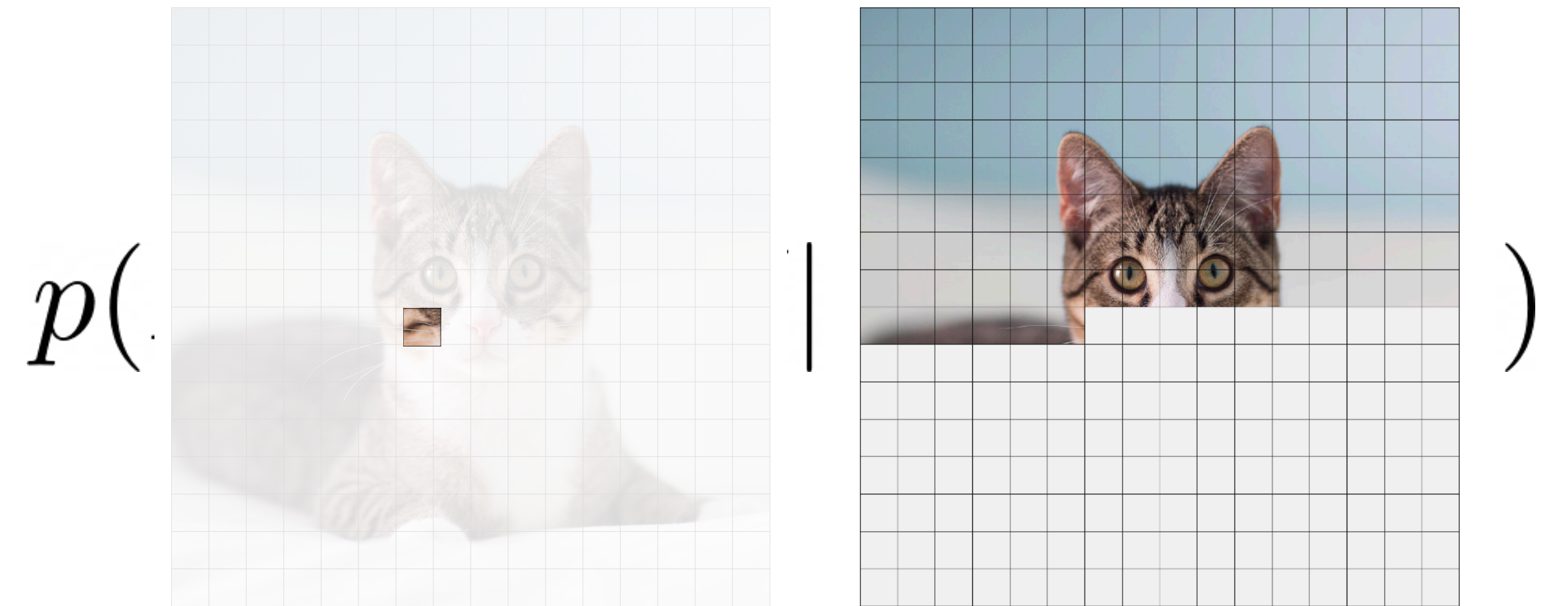
$|$



$)$

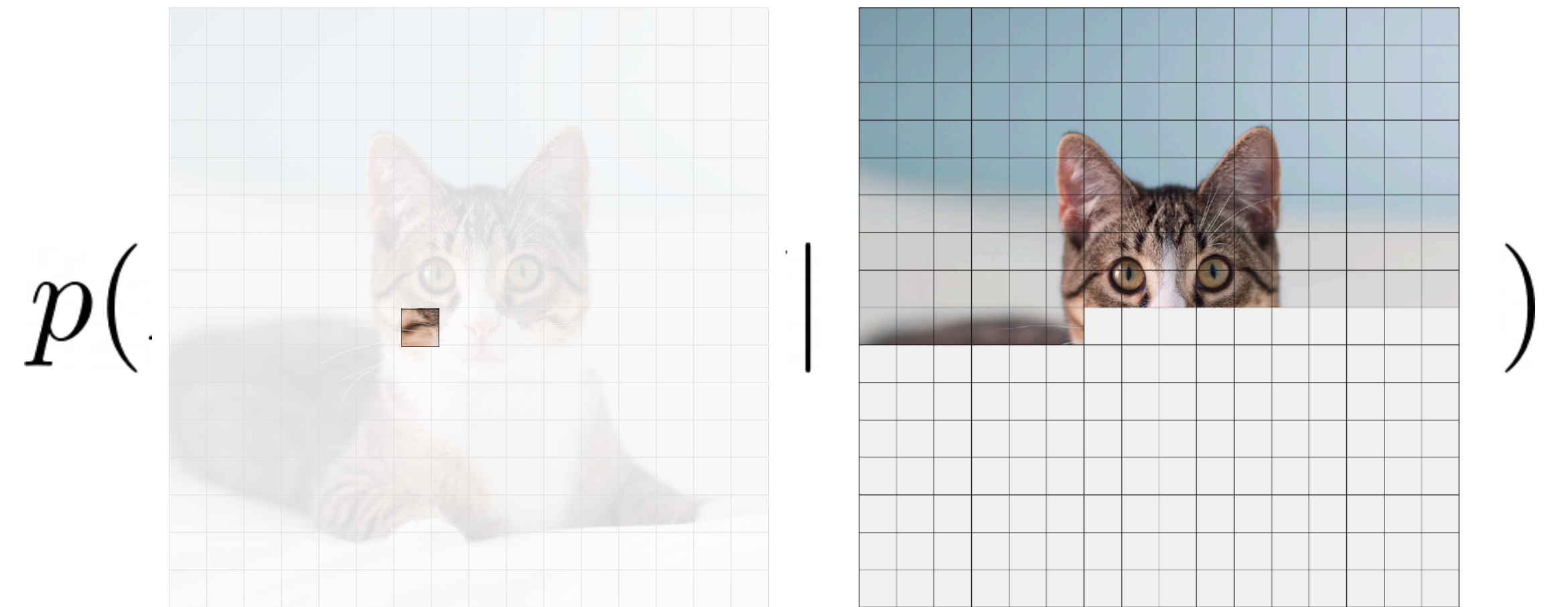
$p(\mathbf{v}_x | S)$ is all you need!

$p(\mathbf{v}_x | S)$ is all you need!



(Sequential) Autoregressive Models
e.g. PixelCNN, PixelRNN, PixelCNN++, Image-GPT

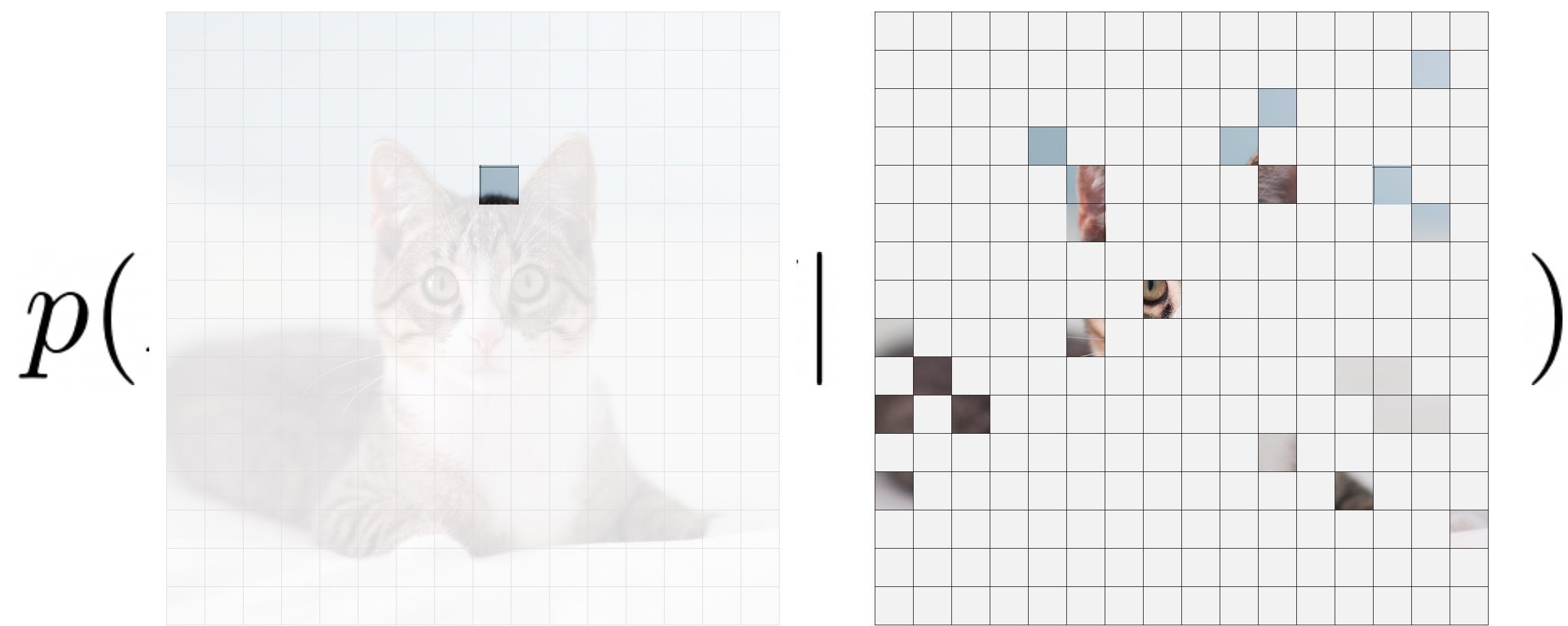
$p(\mathbf{v}_x | S)$ is all you need!



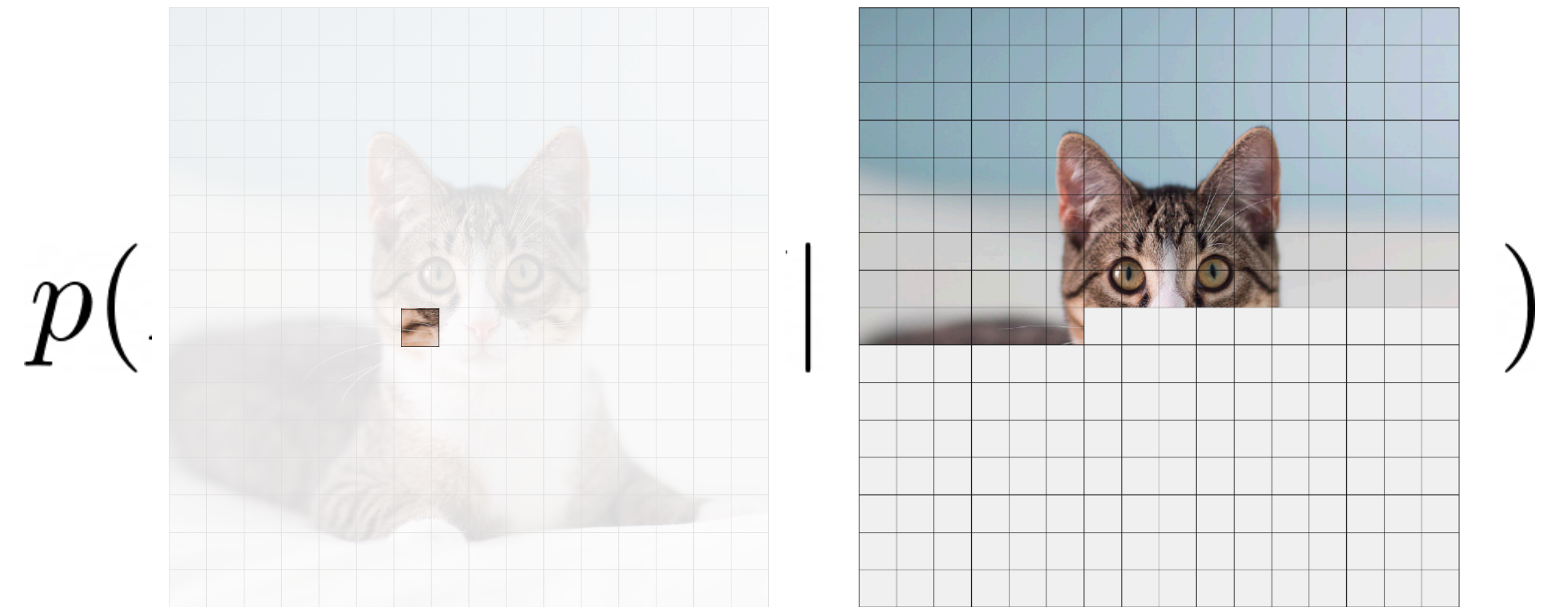
(Sequential) Autoregressive Models
e.g. PixelCNN, PixelRNN, PixelCNN++, Image-GPT

Ordered S; $x = \text{'next' pixel}$

$p(\mathbf{v}_x | S)$ is all you need!



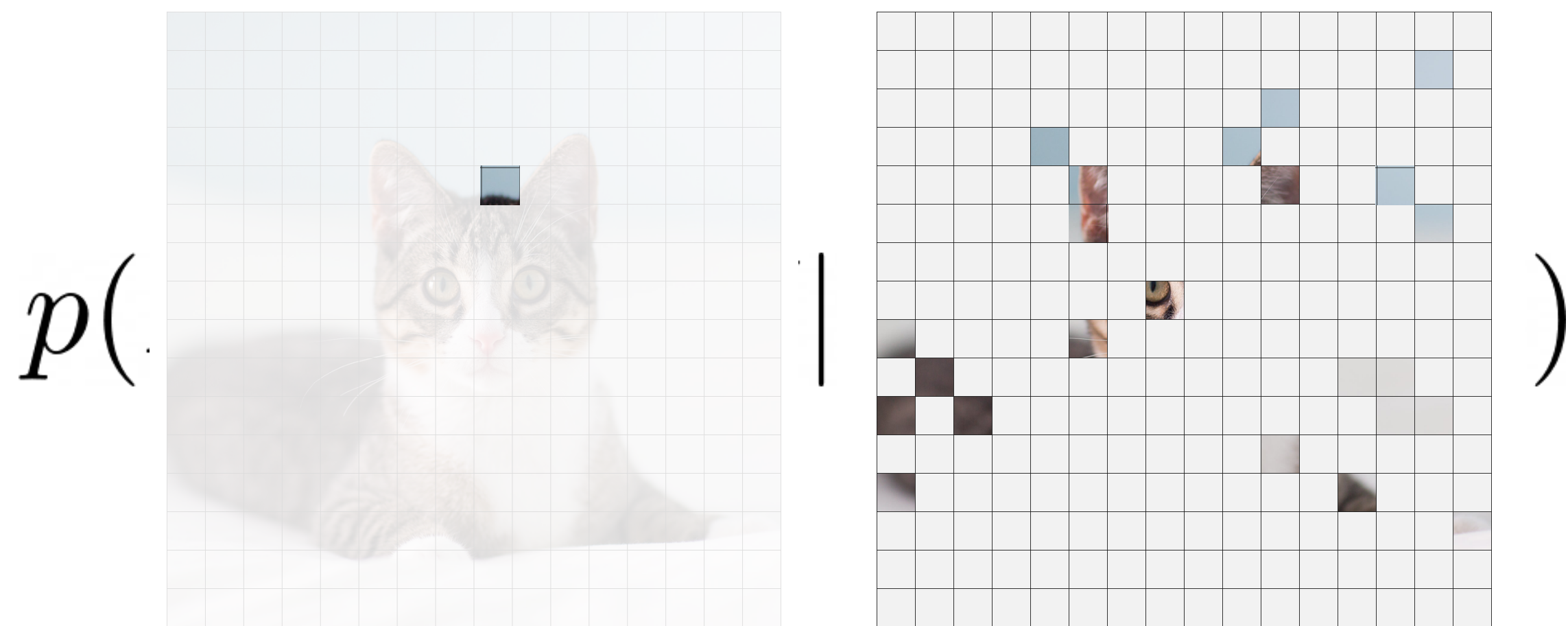
Ours



(Sequential) Autoregressive Models
e.g. PixelCNN, PixelRNN, PixelCNN++, Image-GPT

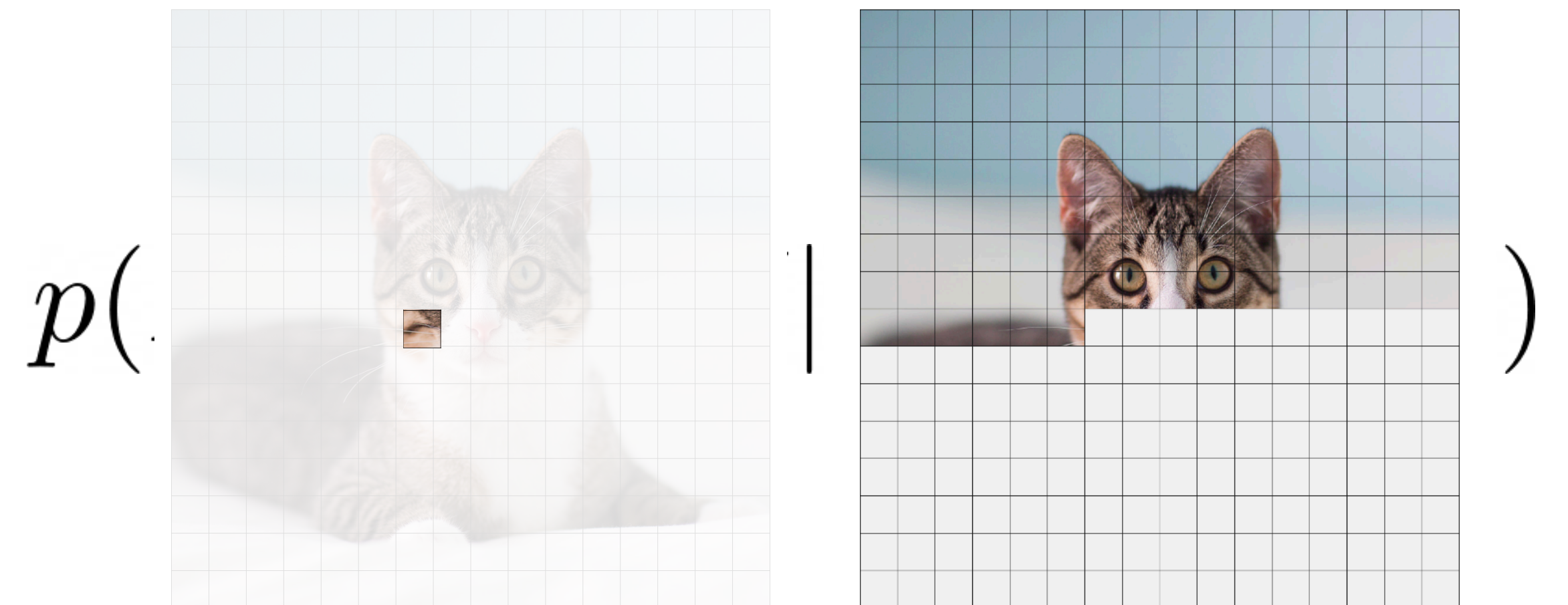
Ordered S; $x = \text{'next' pixel}$

$p(\mathbf{v}_x | S)$ is all you need!



Ours

Arbitrary S and x

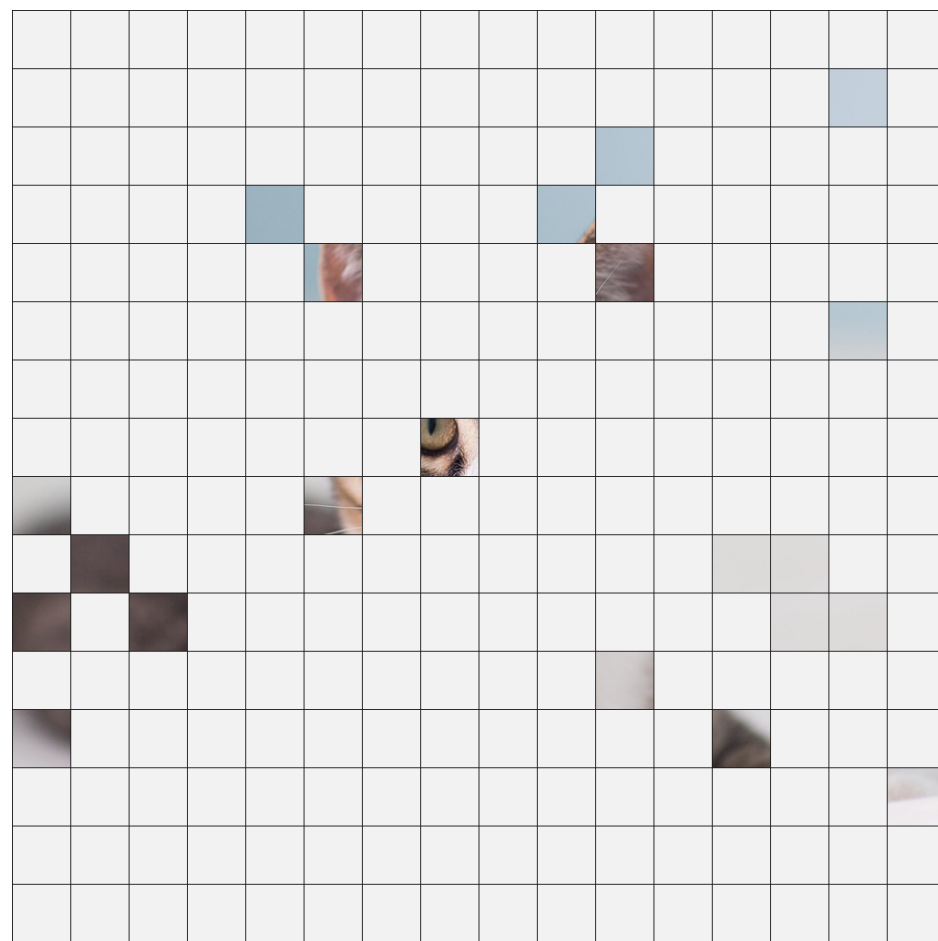


(Sequential) Autoregressive Models

e.g. PixelCNN, PixelRNN, PixelCNN++, Image-GPT

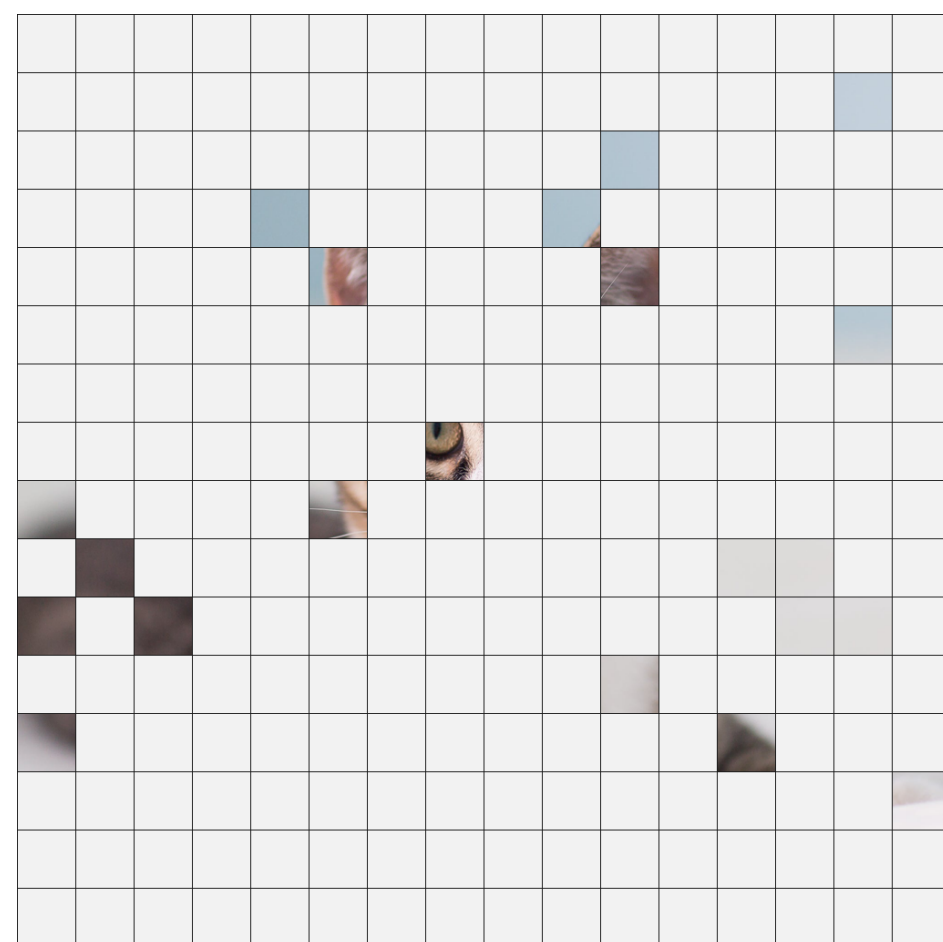
Ordered S; x = 'next' pixel

$p(\mathbf{v}_x | S)$ is all you need!

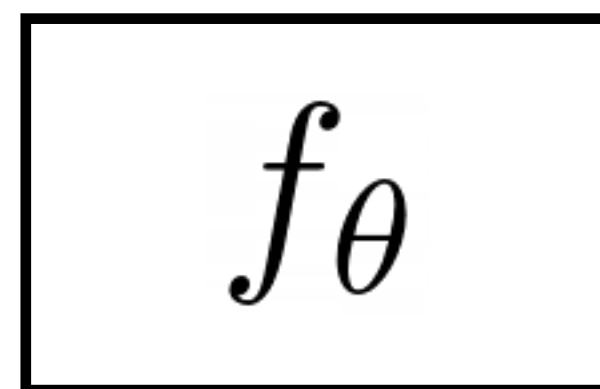


$$f_{\theta}$$

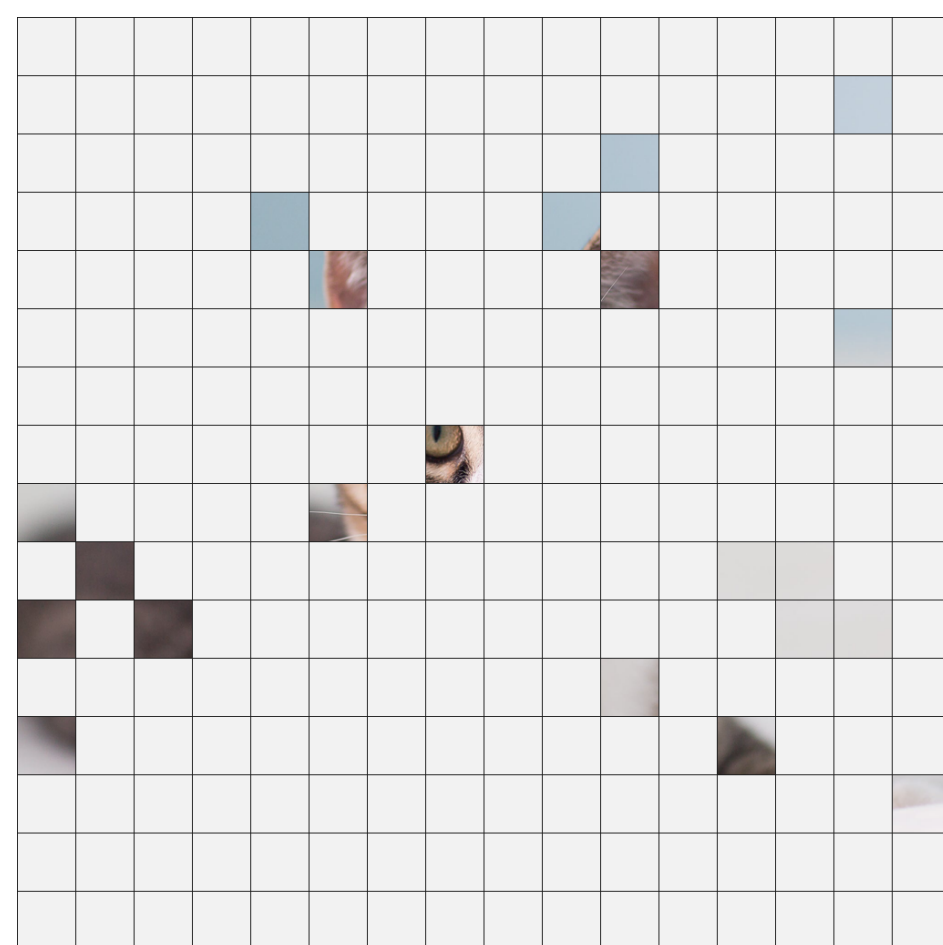
$p(\mathbf{v}_{\mathbf{x}} | \mathcal{S})$ is all you need!



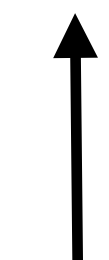
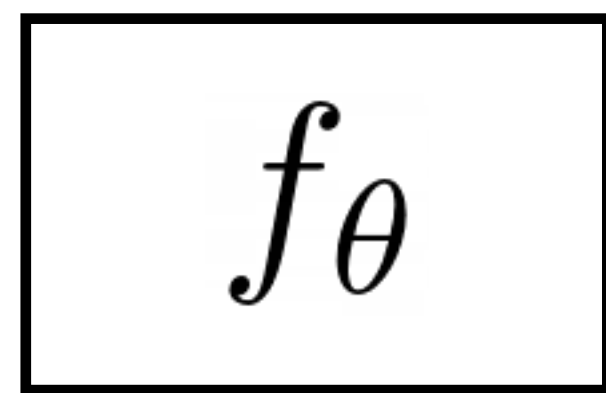
$\{(\mathbf{x}_k, \mathbf{v}_k)\}$



$p(\mathbf{v}_{\mathbf{x}} | \mathcal{S})$ is all you need!

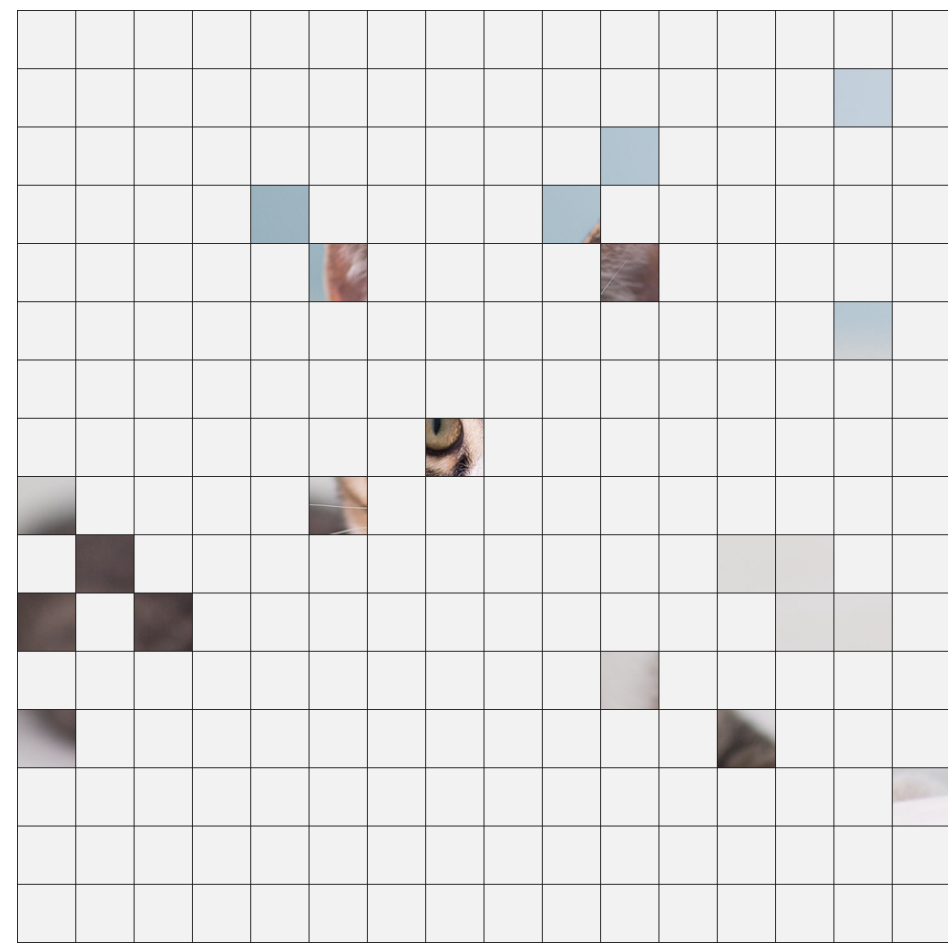


$\{(\mathbf{x}_k, \mathbf{v}_k)\}$



\mathbf{x}

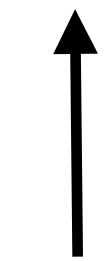
$p(\mathbf{v}_x | S)$ is all you need!



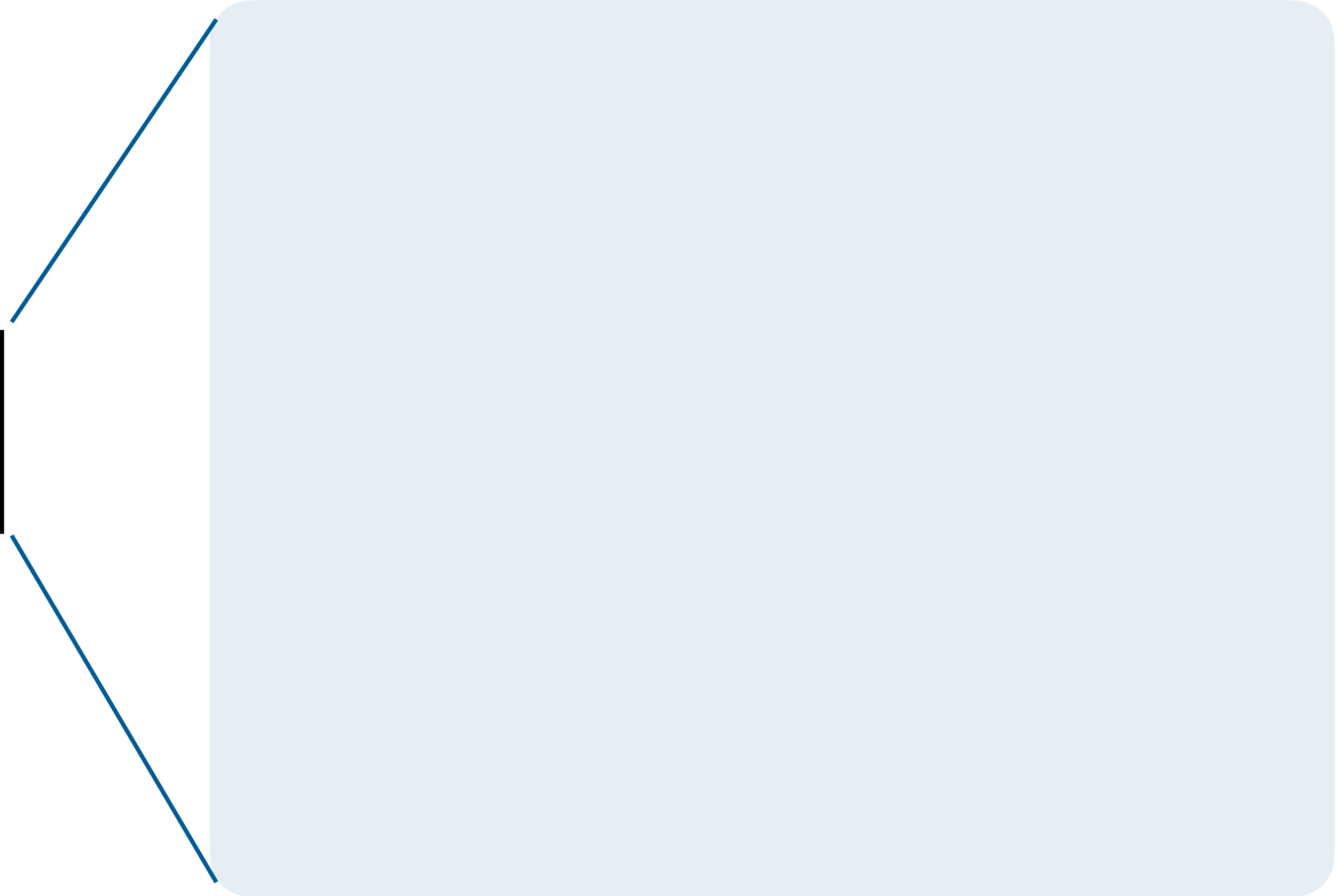
$\{(\mathbf{x}_k, \mathbf{v}_k)\}$



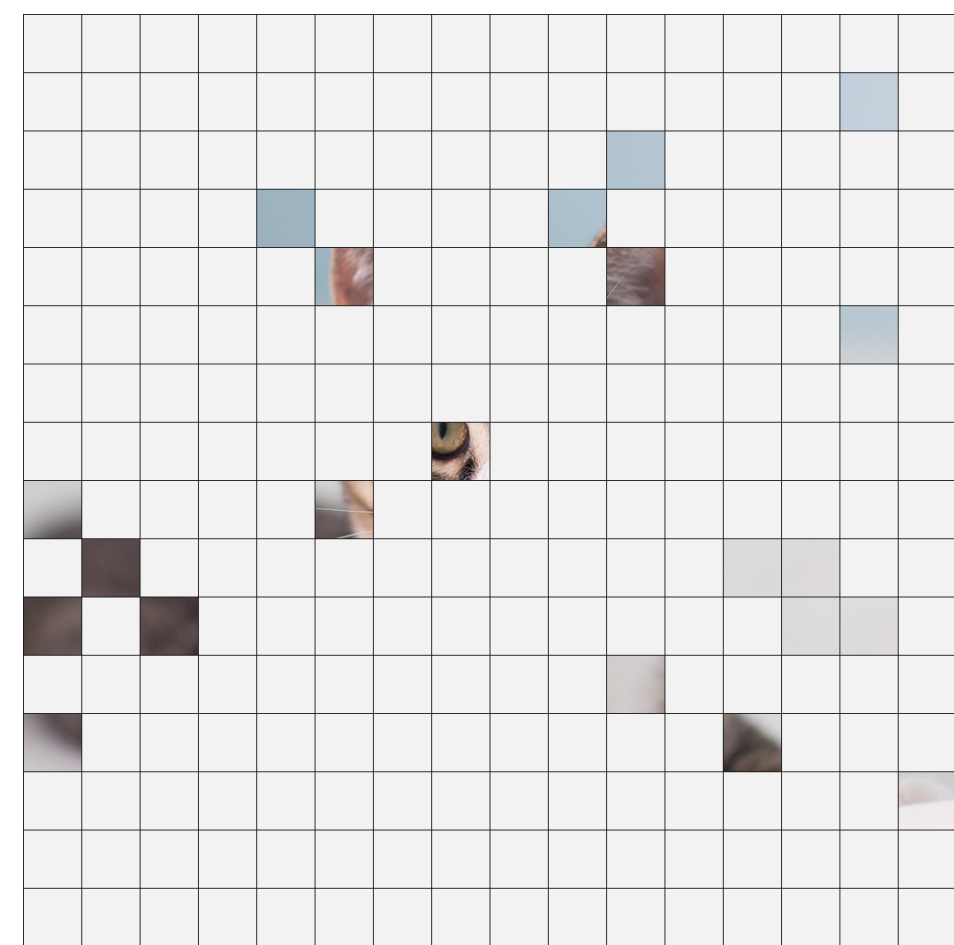
f_θ



\mathbf{x}



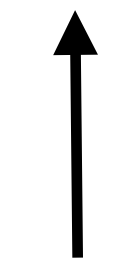
$p(\mathbf{v}_x | S)$ is all you need!



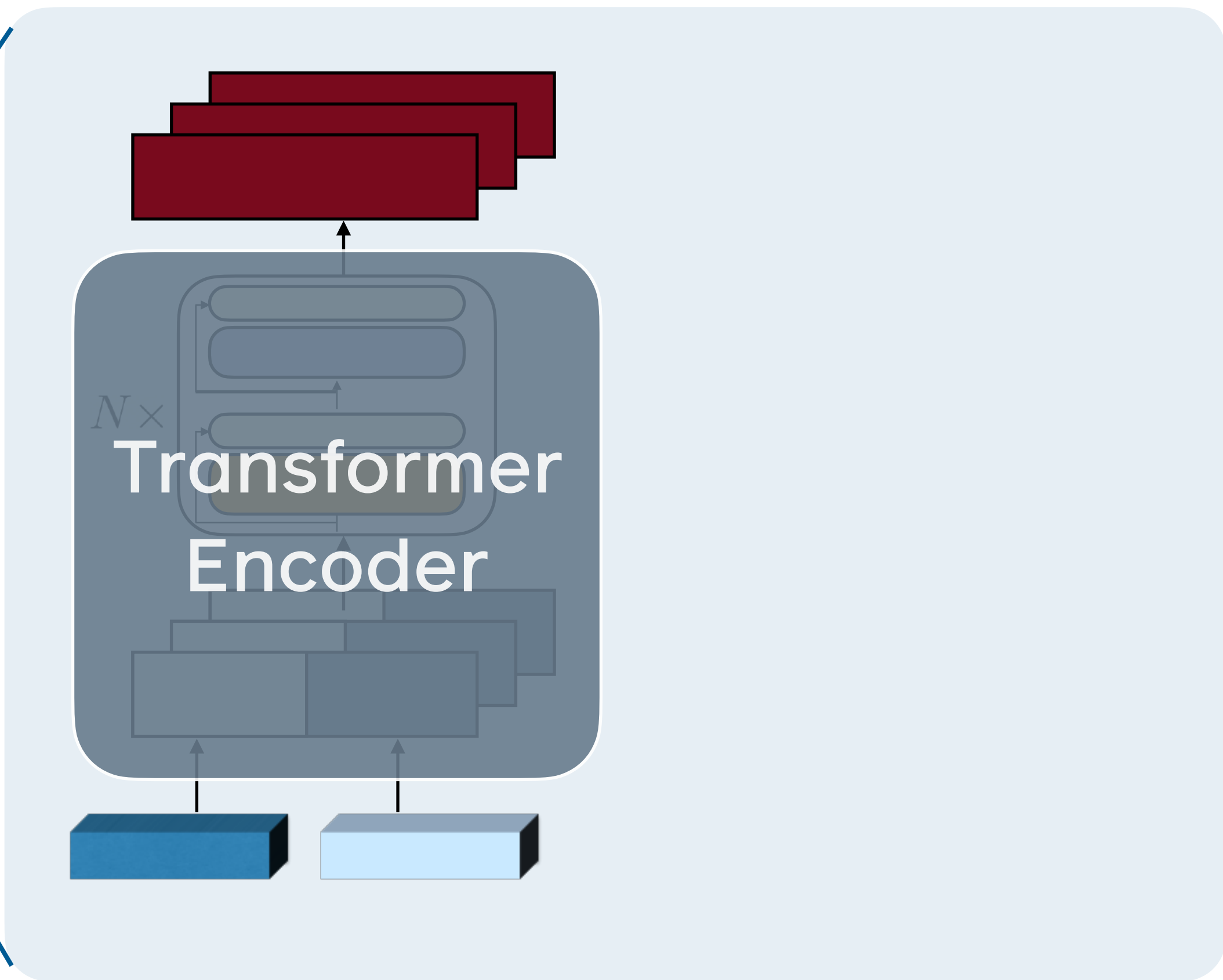
$\{(\mathbf{x}_k, \mathbf{v}_k)\}$



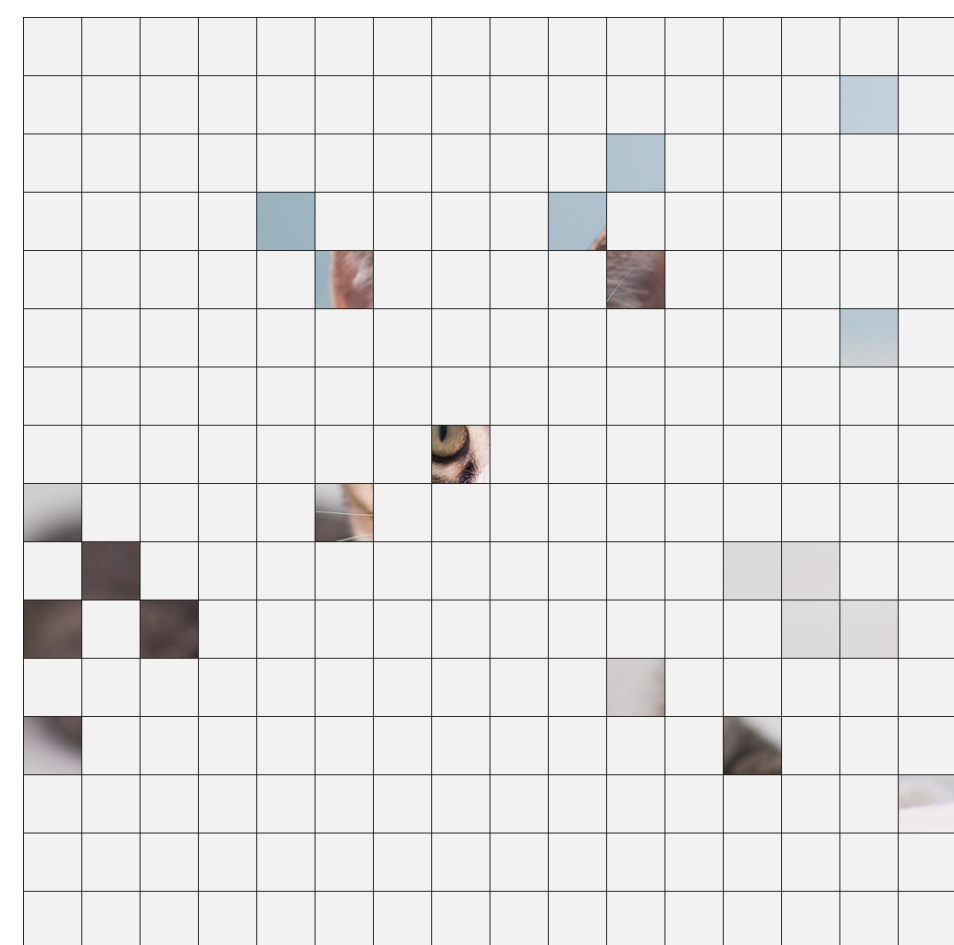
f_θ



\mathbf{x}



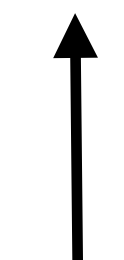
$p(\mathbf{v}_x | \mathcal{S})$ is all you need!



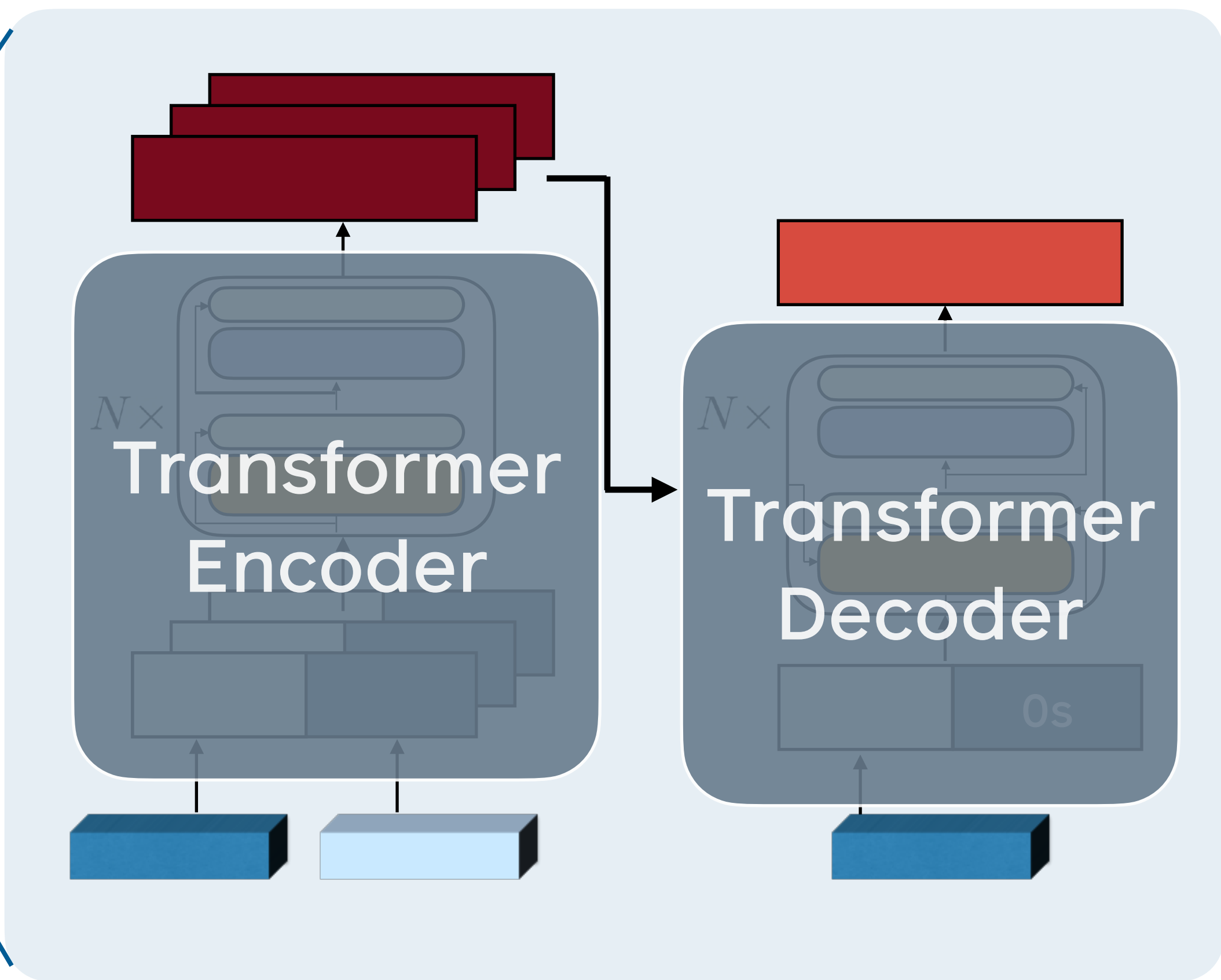
$\{(\mathbf{x}_k, \mathbf{v}_k)\}$



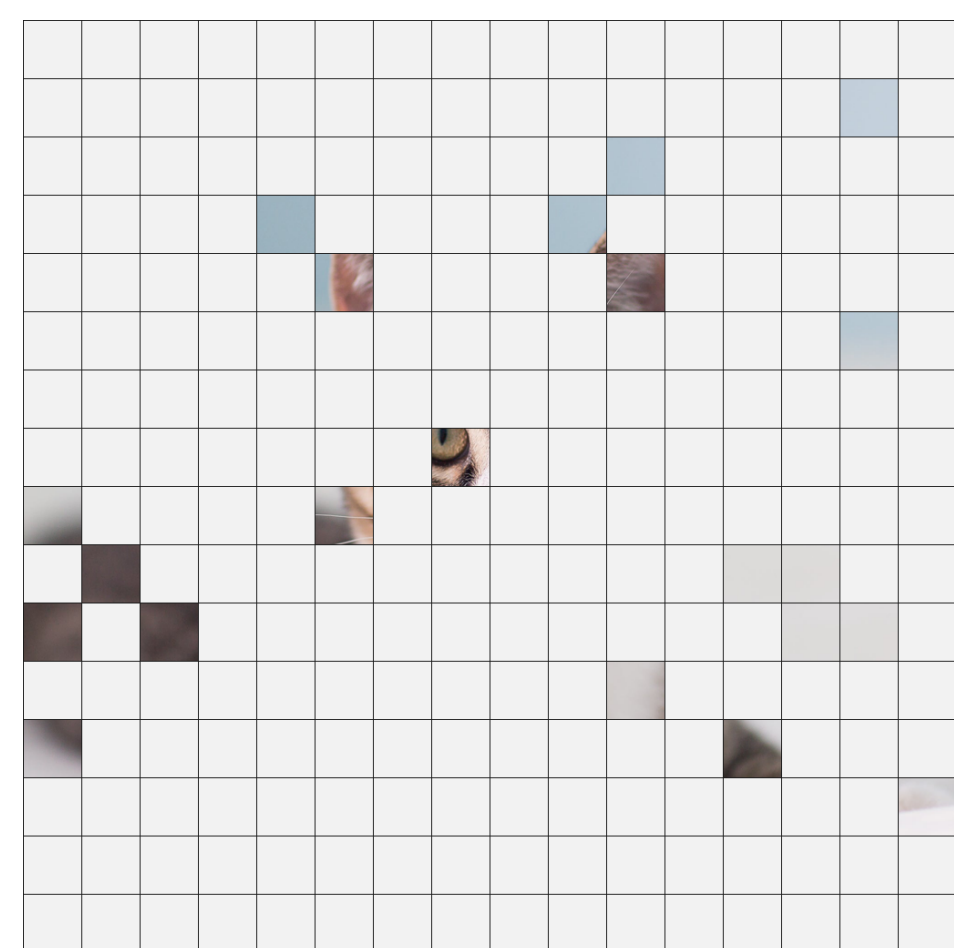
f_θ



\mathbf{x}



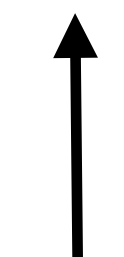
$p(\mathbf{v}_x | S)$ is all you need!



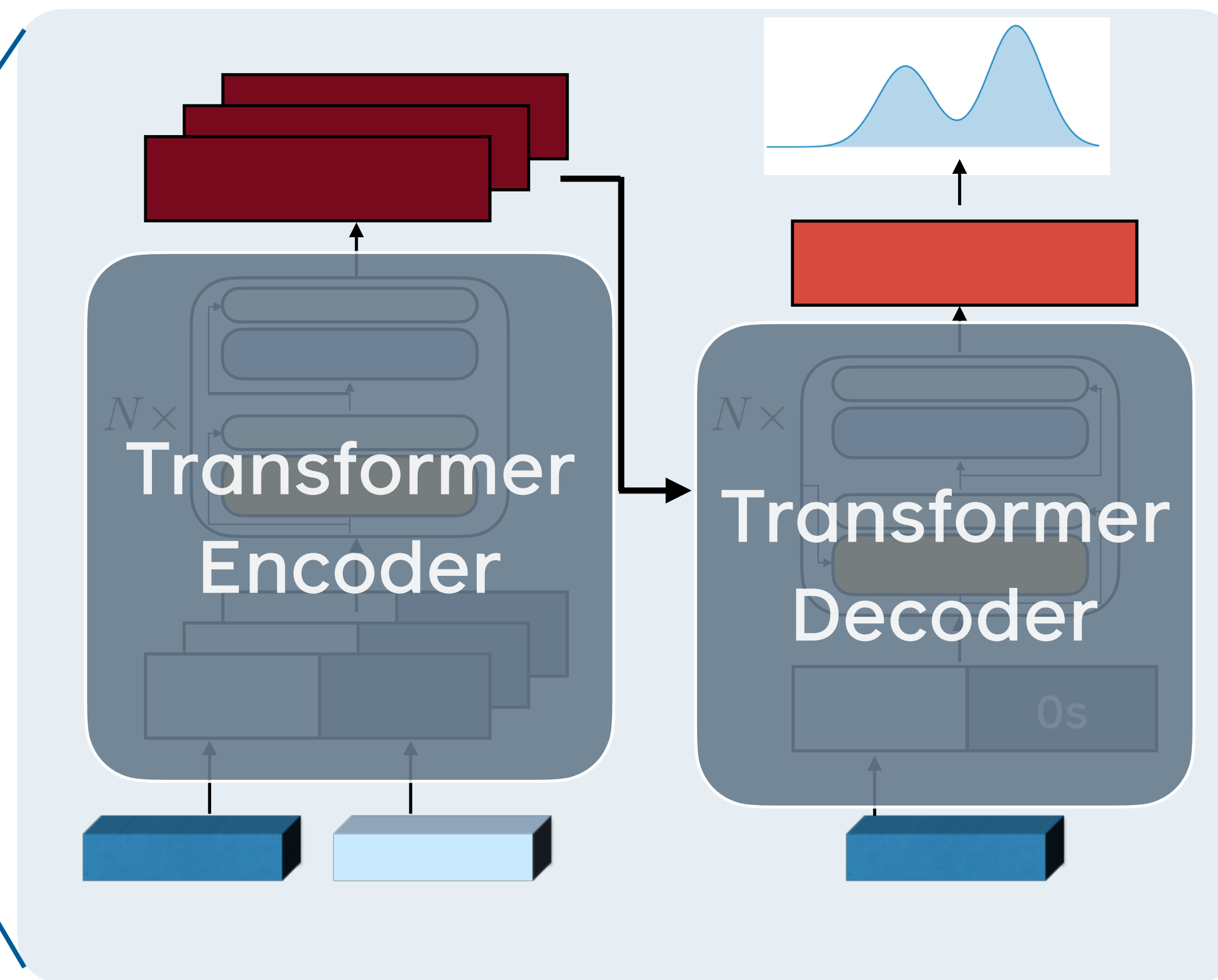
$\{(\mathbf{x}_k, \mathbf{v}_k)\}$



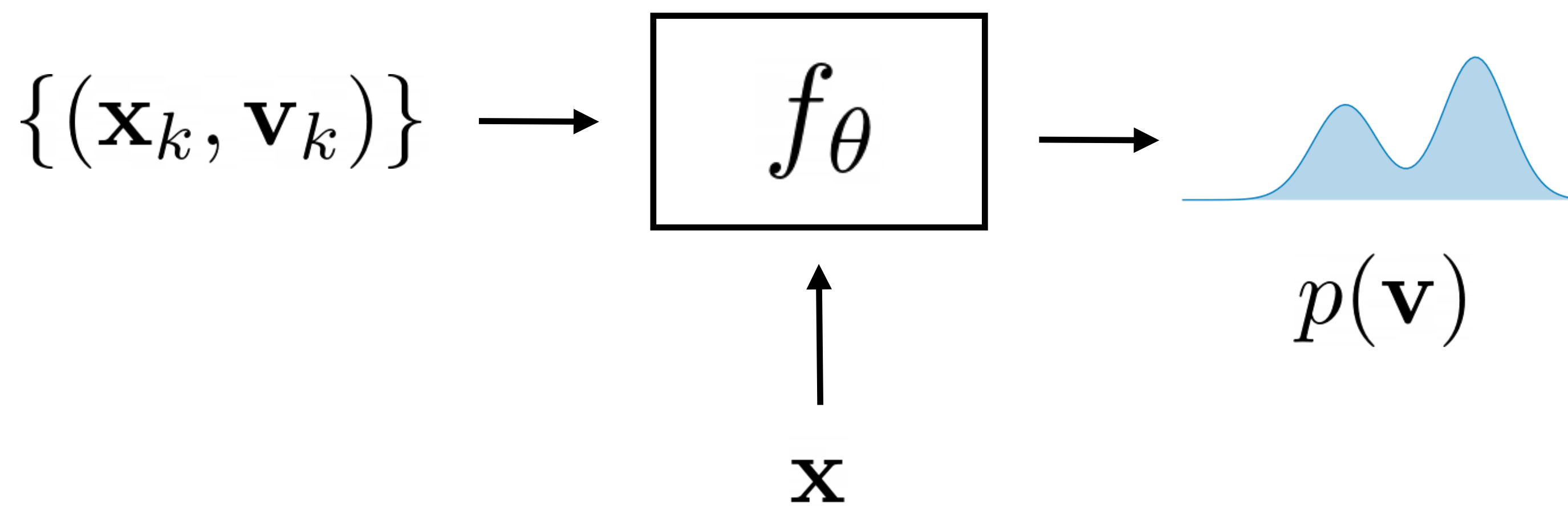
f_θ



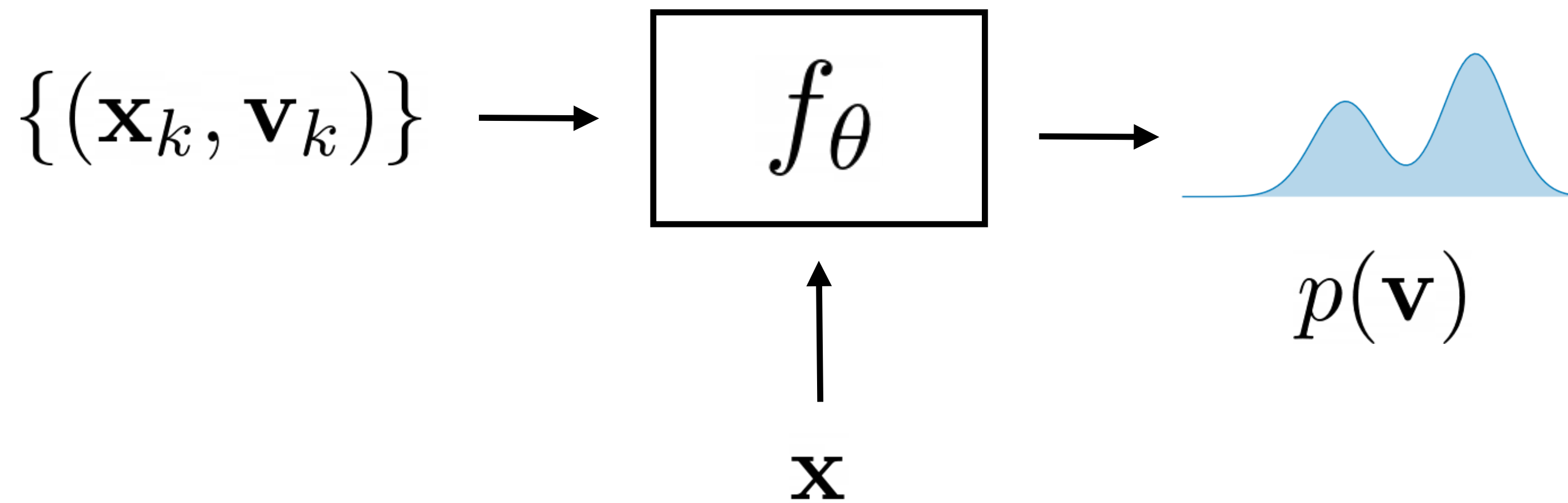
\mathbf{x}



$p(\mathbf{v}_{\mathbf{x}} | \mathcal{S})$ is all you need!

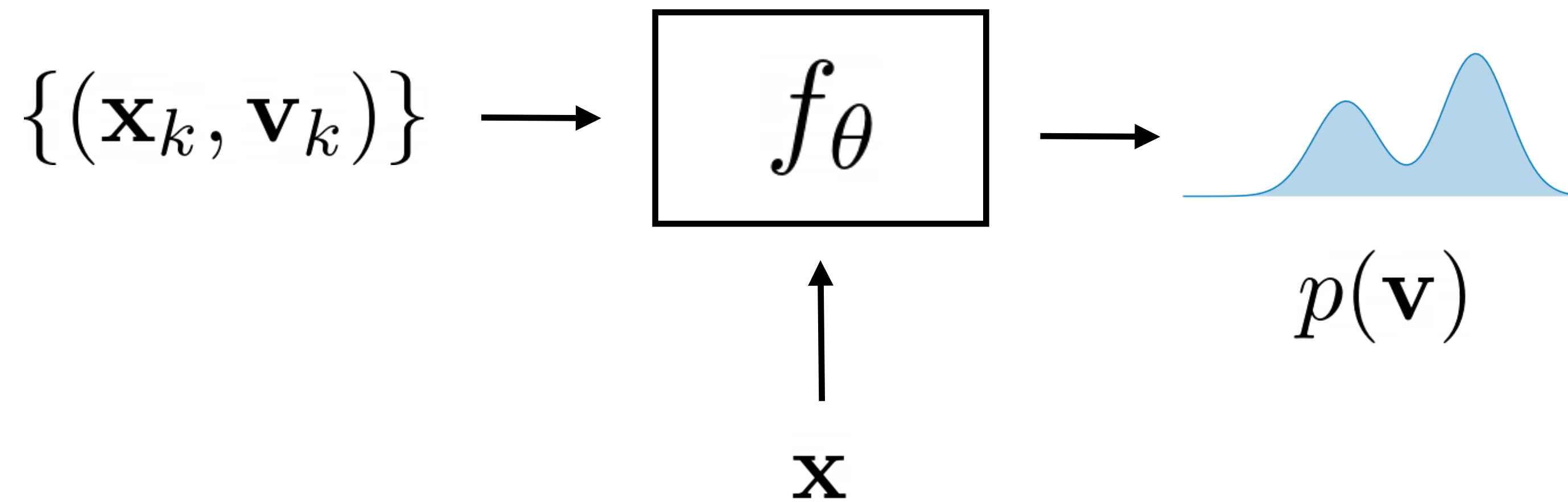


$p(\mathbf{v}_{\mathbf{x}} | S)$ is all you need!



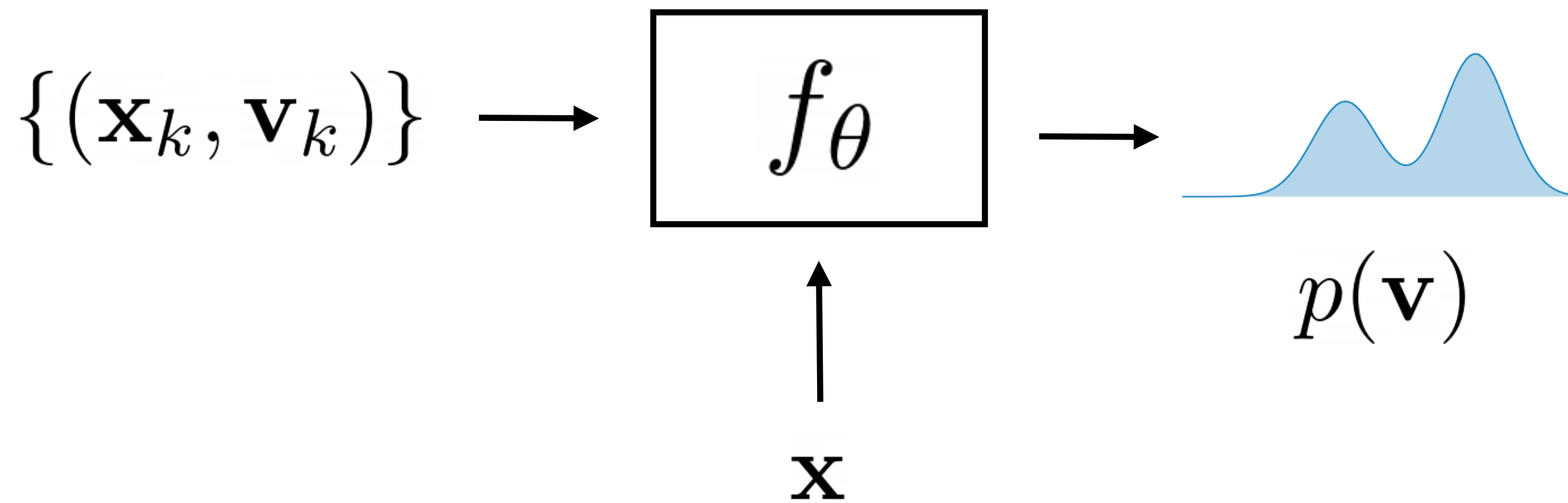
How do we enable learning?

$p(\mathbf{v}_{\mathbf{x}} | S)$ is all you need!



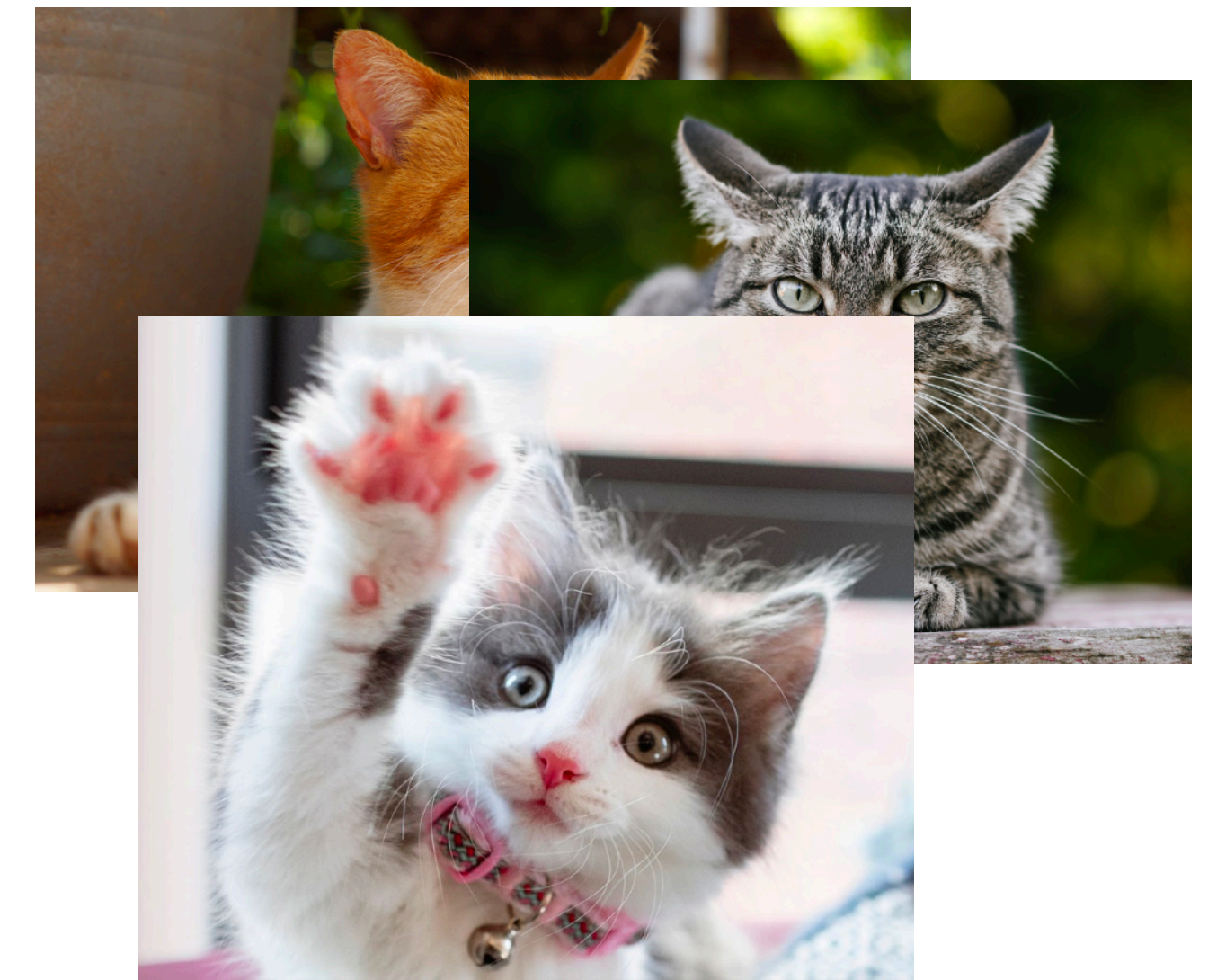
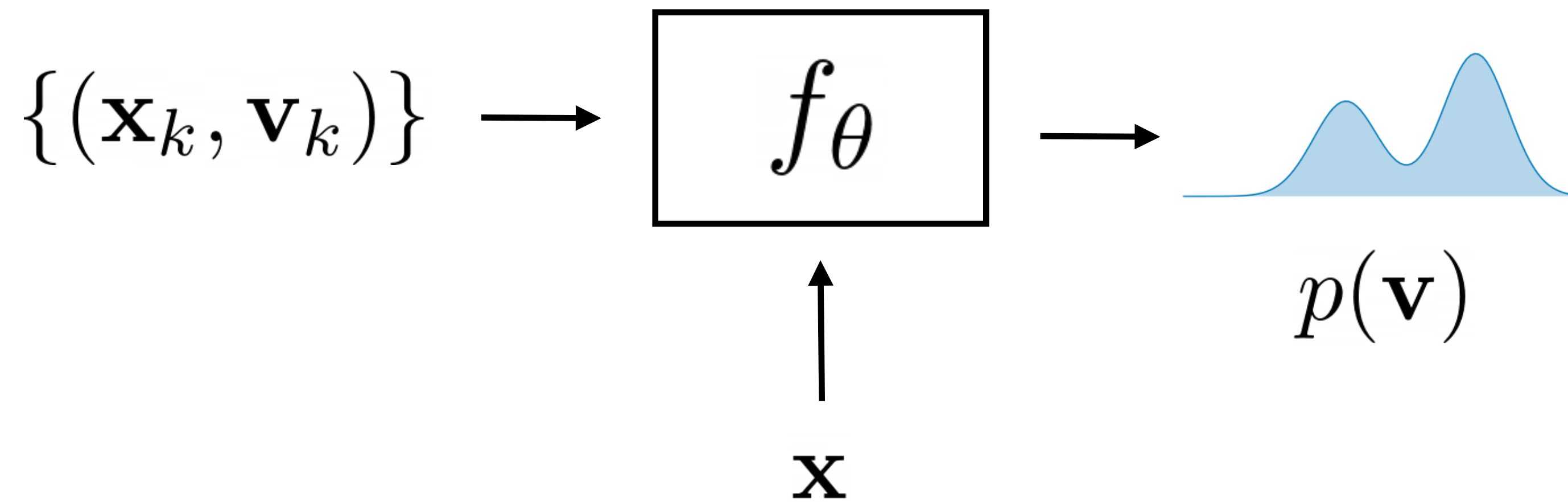
How do we enable learning?

$p(\mathbf{v}_{\mathbf{x}} | S)$ is all you need!



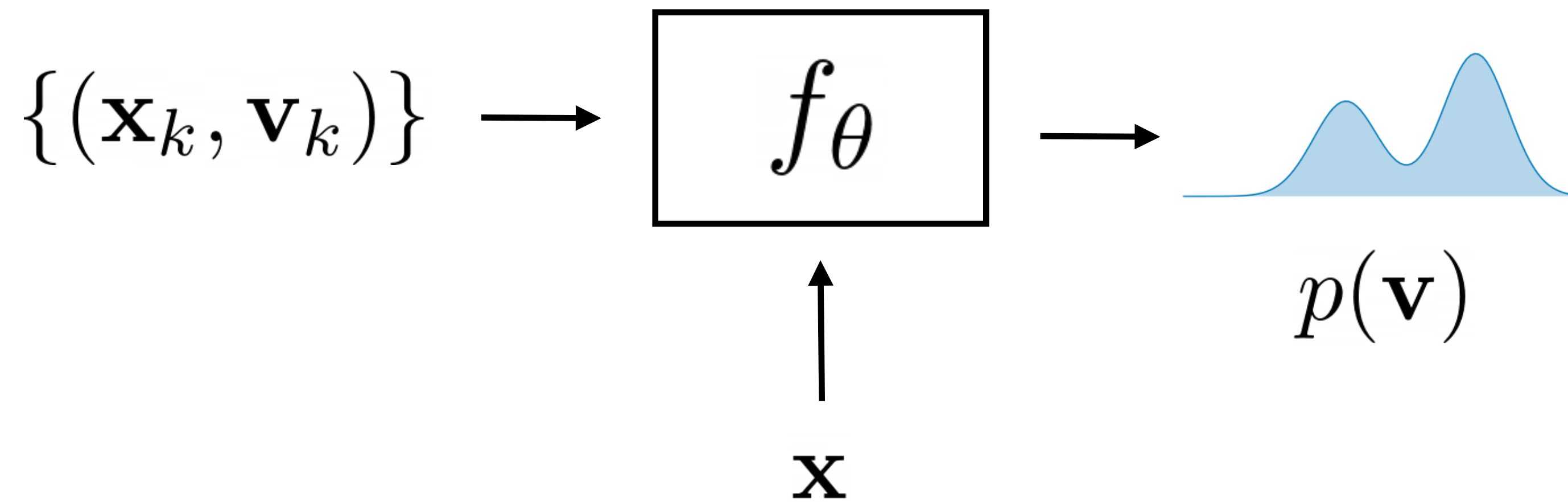
How do we enable learning?

$p(\mathbf{v}_{\mathbf{x}} | S)$ is all you need!



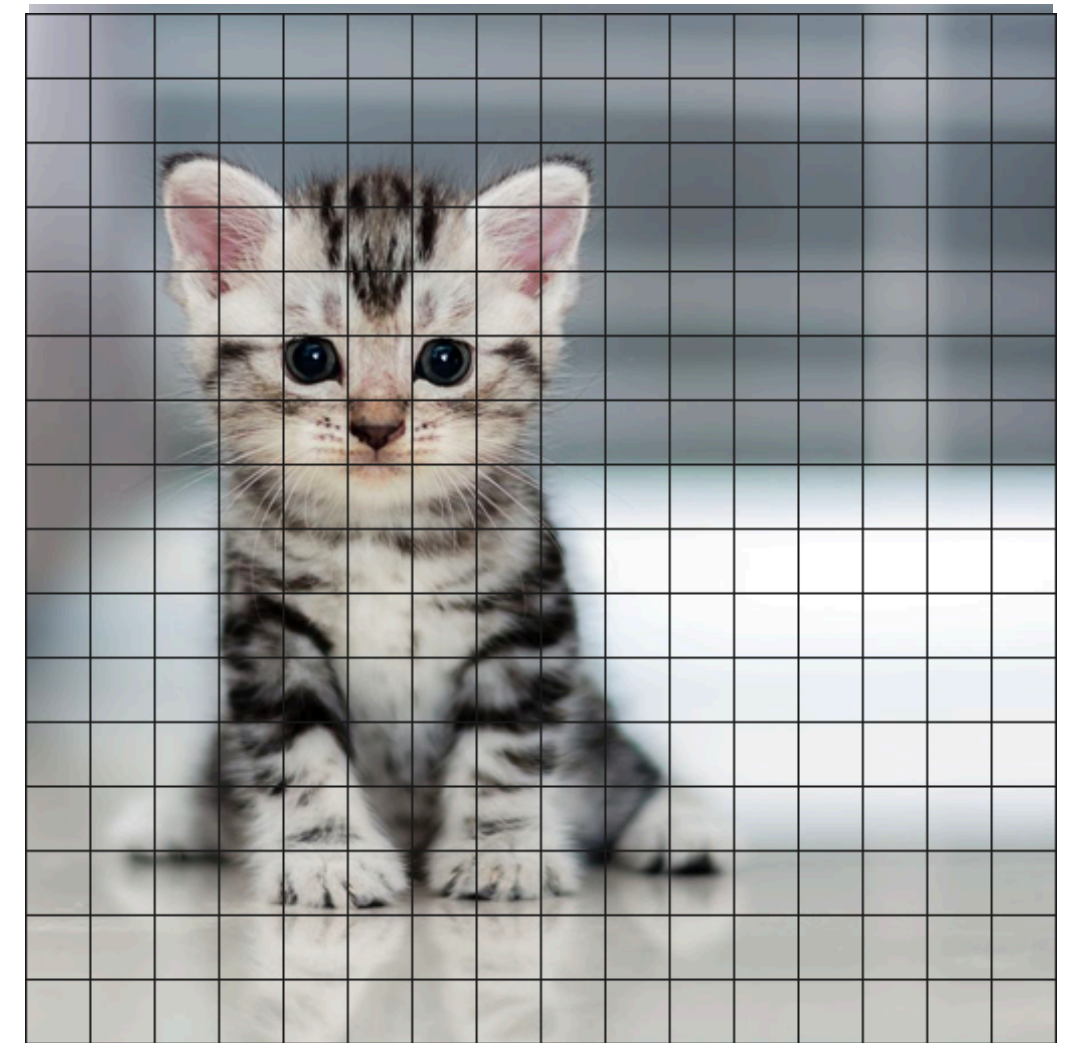
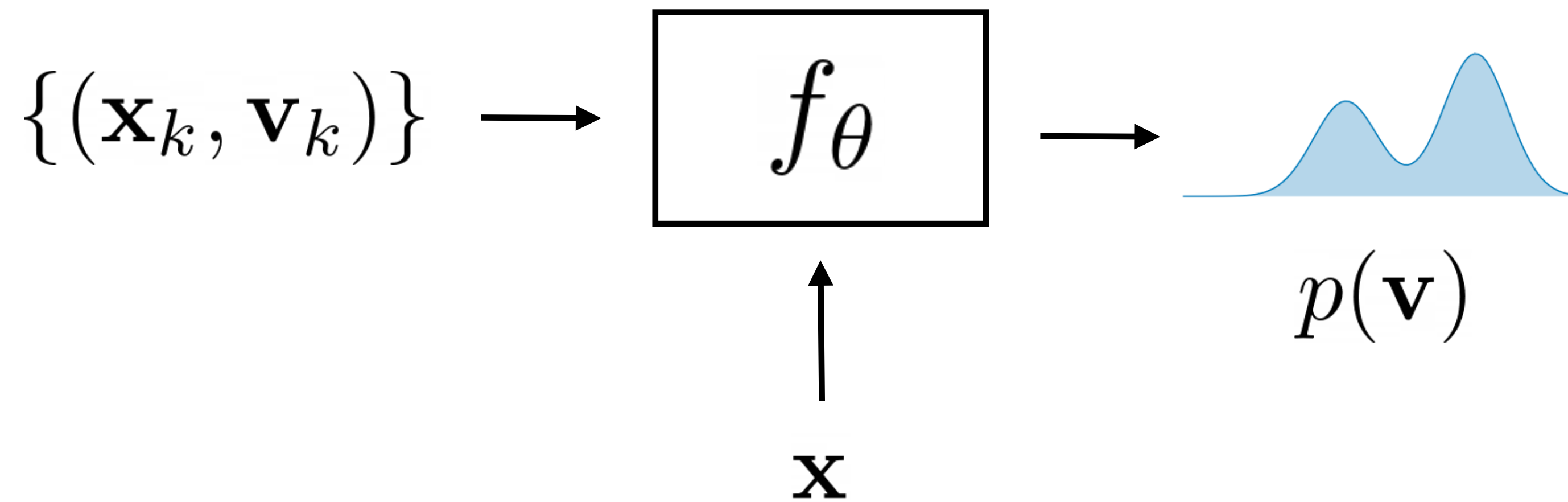
How do we enable learning?

$p(\mathbf{v}_{\mathbf{x}} | \mathcal{S})$ is all you need!

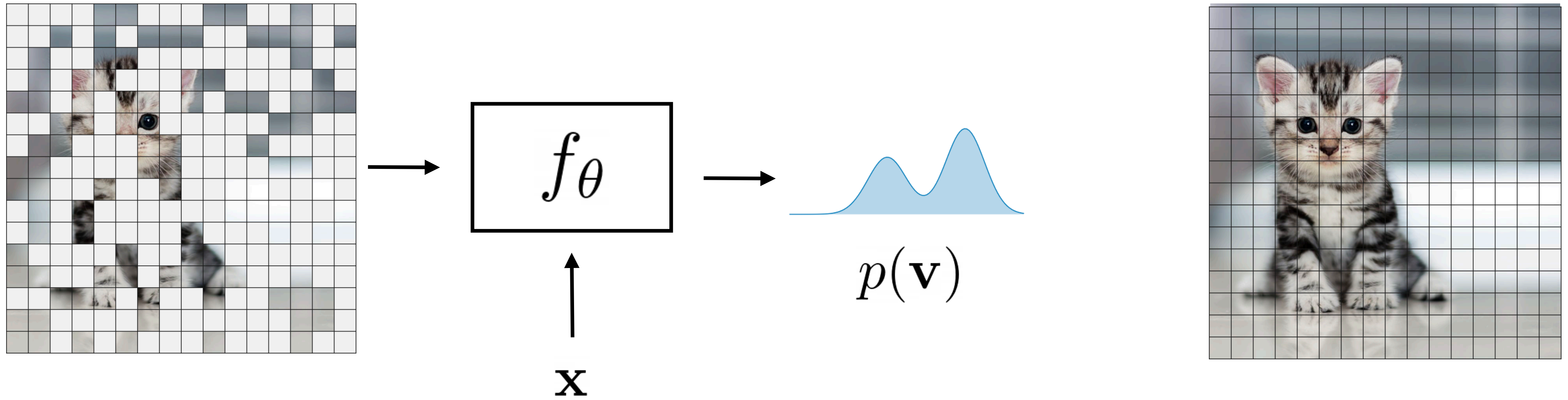


How do we enable learning?

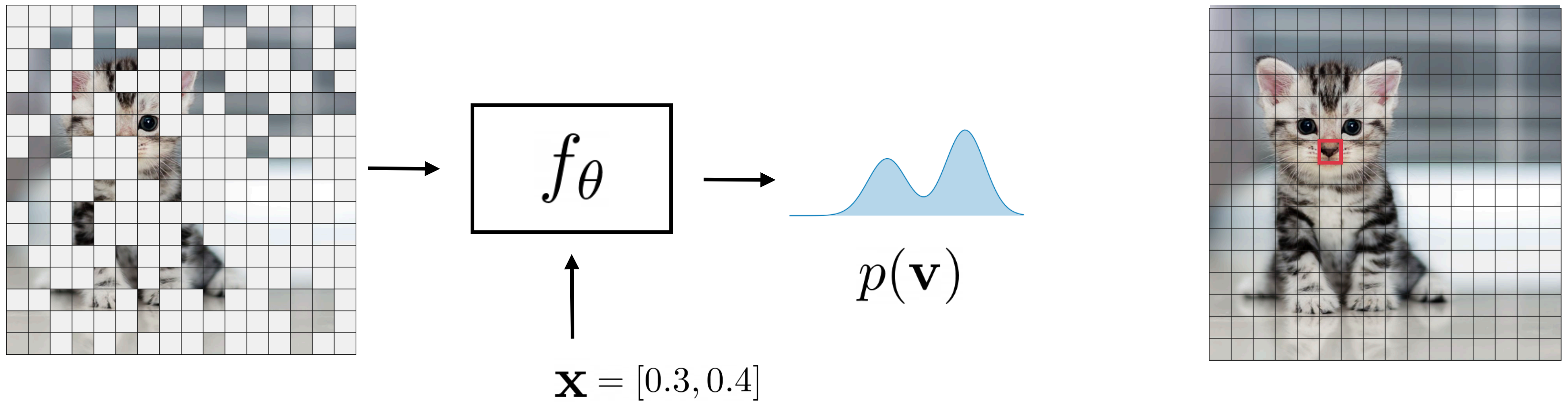
Learning from Self-supervision



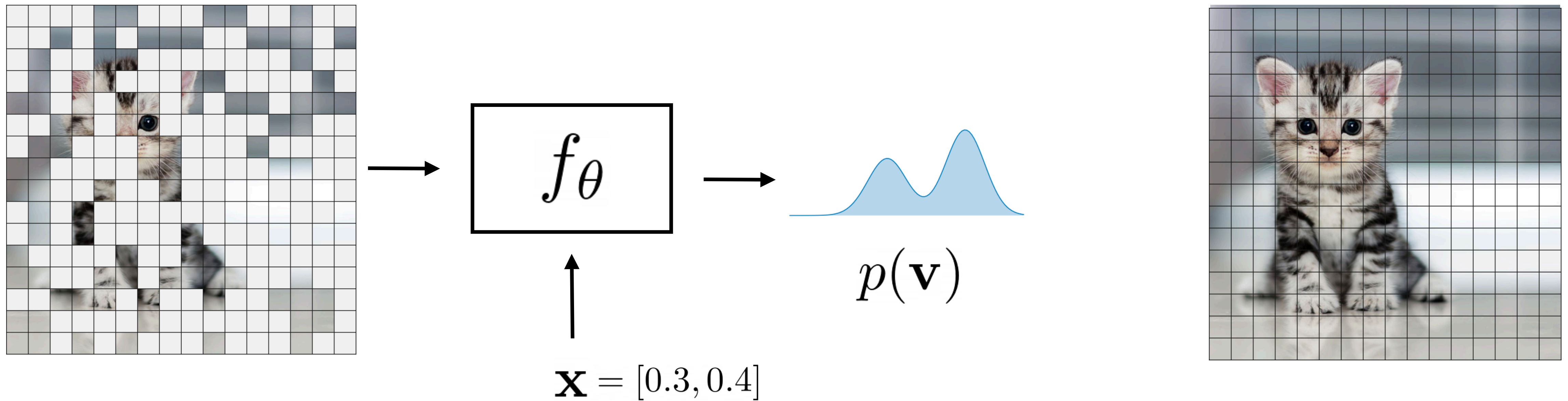
Learning from Self-supervision



Learning from Self-supervision

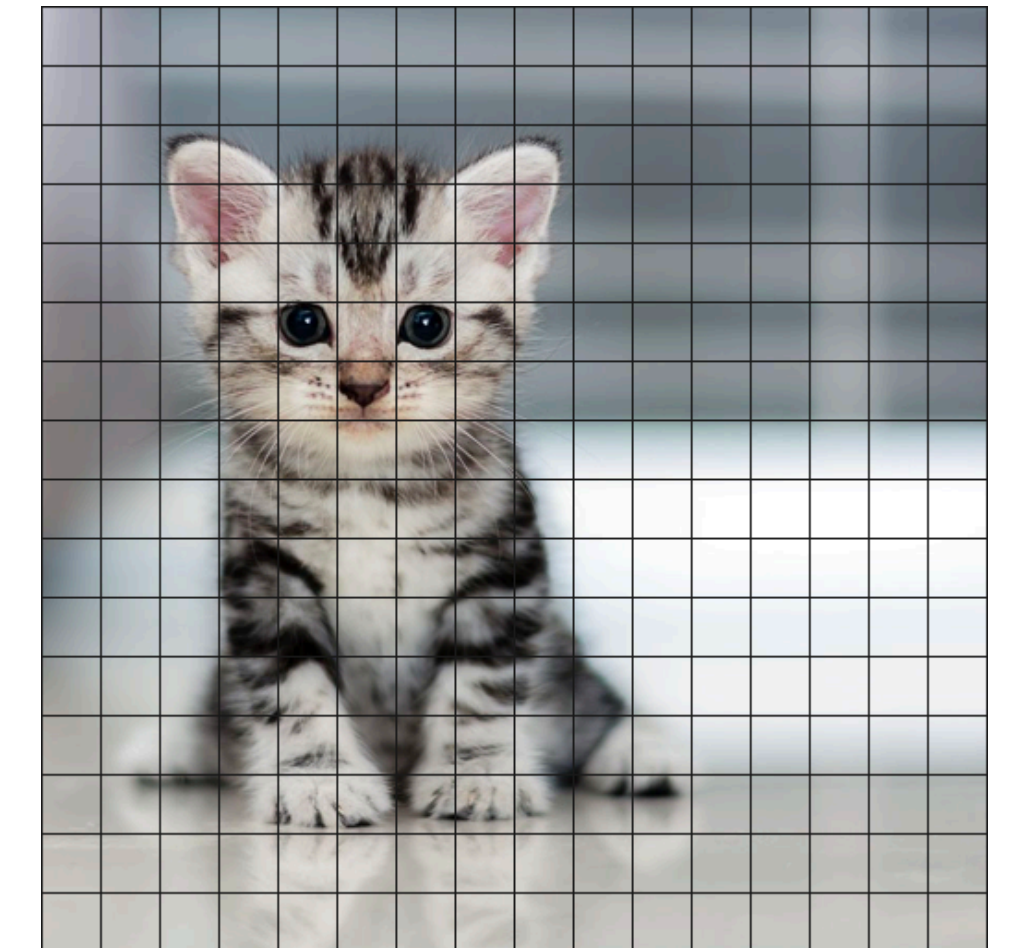
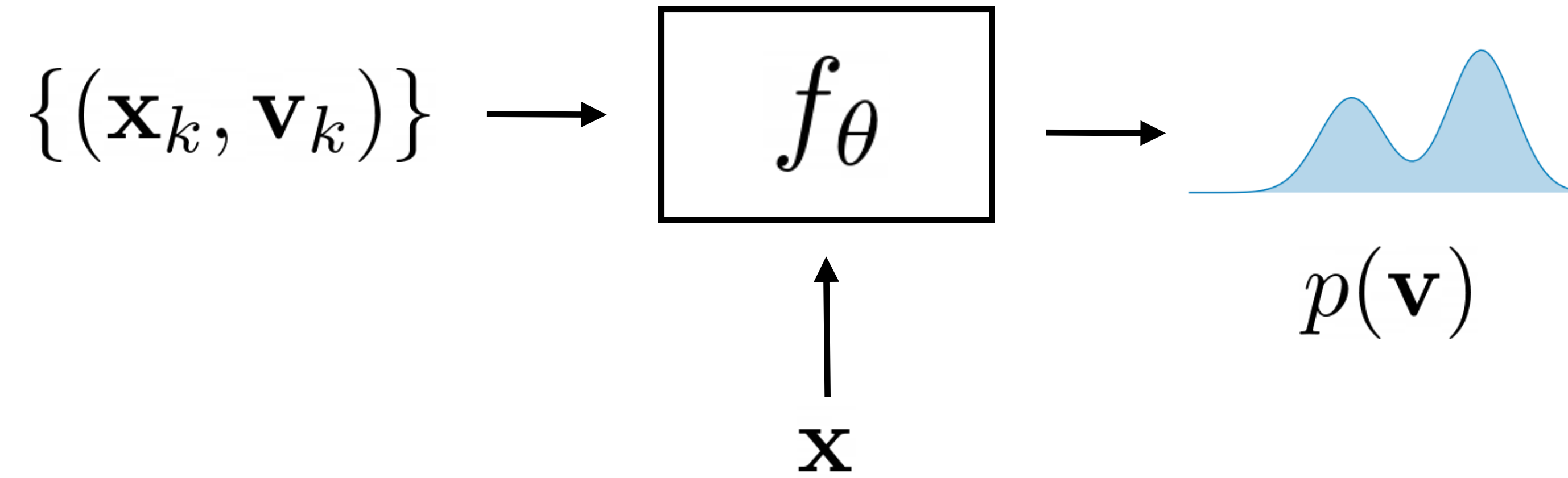


Learning from Self-supervision

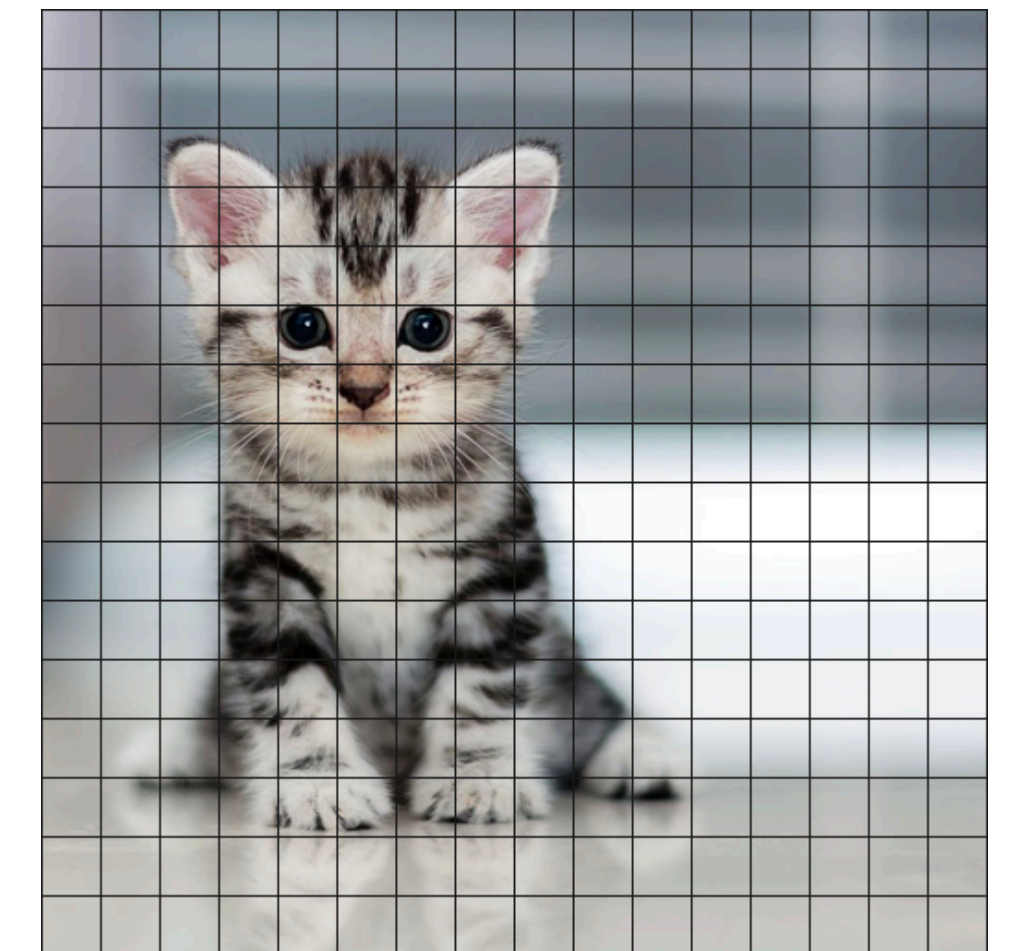
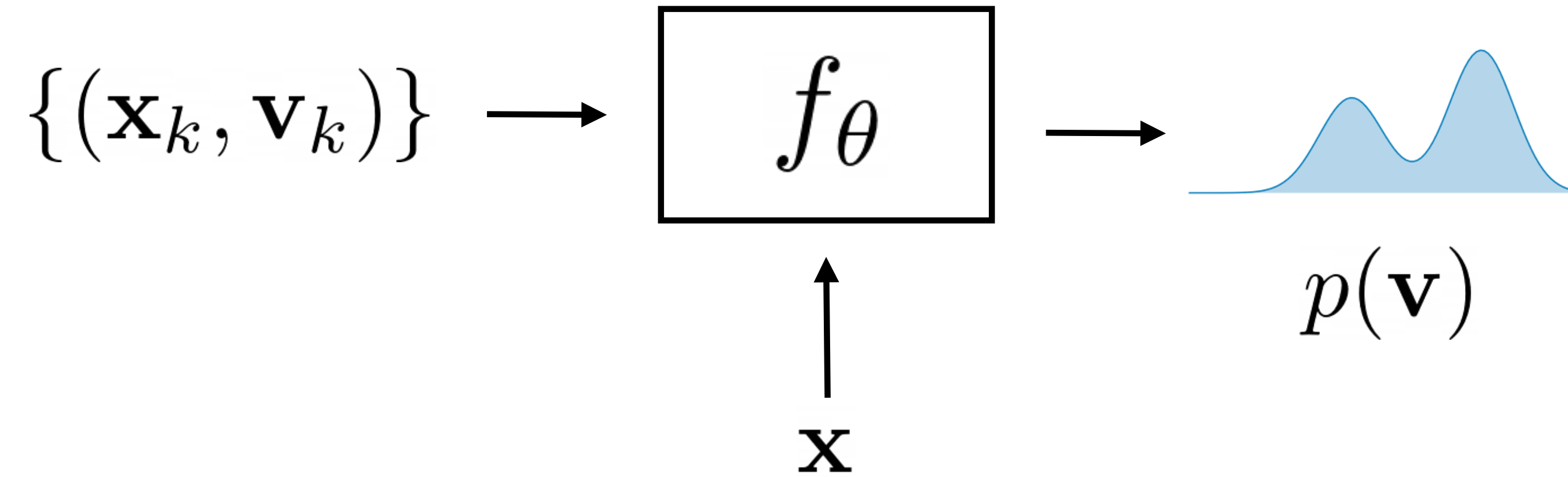


Training Objective: $-\log p(\text{img}; \text{distribution})$

Modeling Generic Spatial Signals



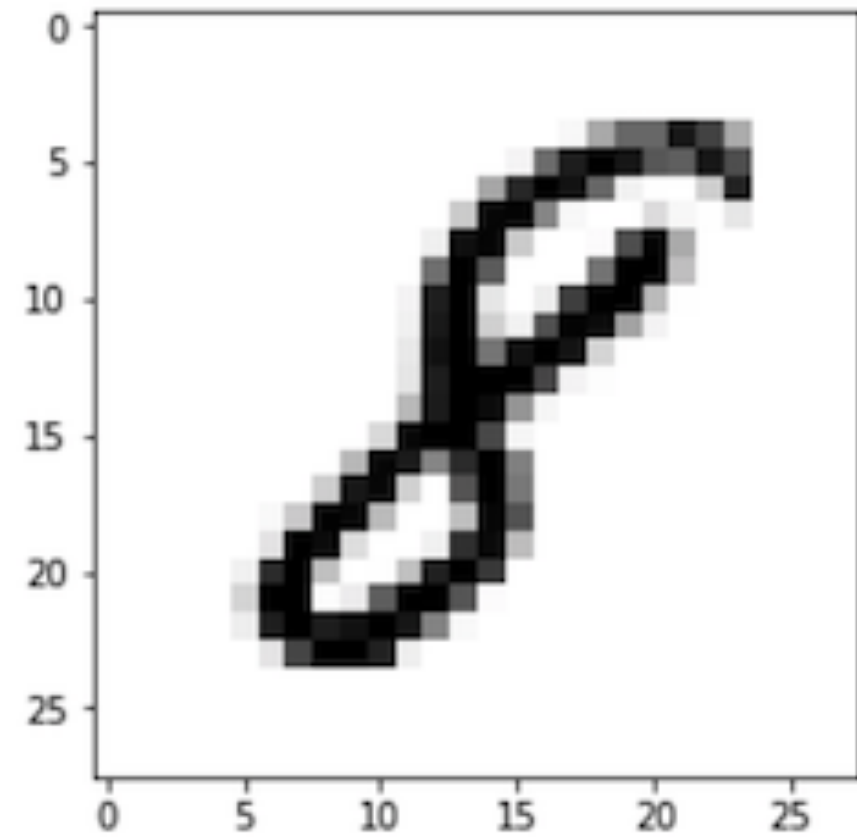
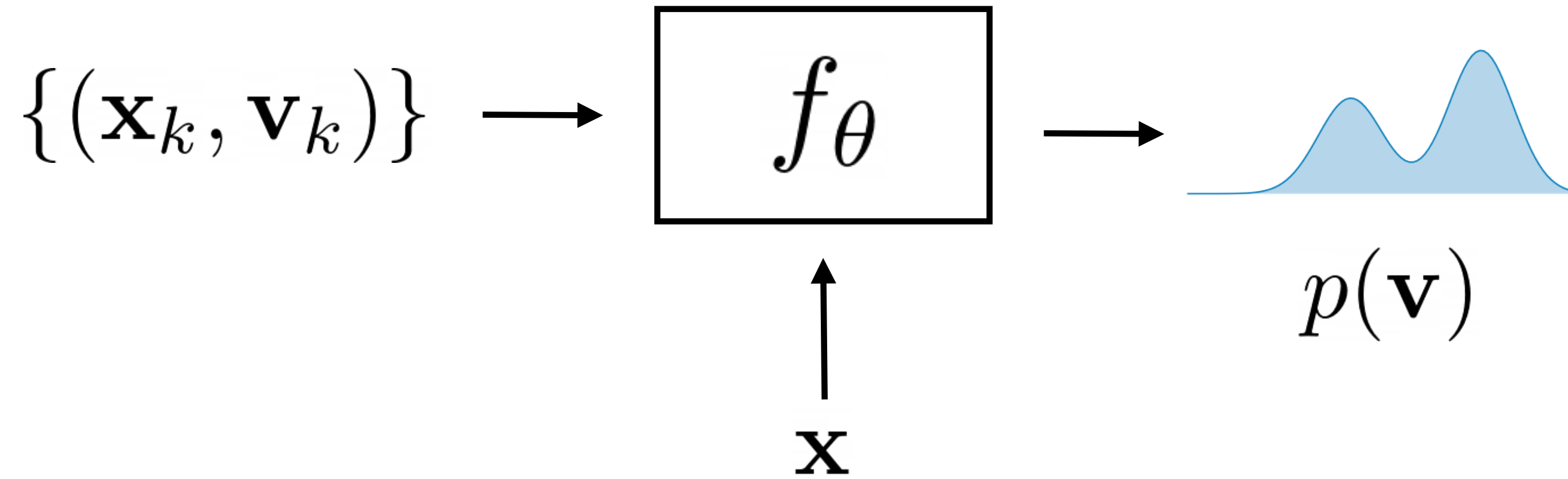
Modeling Generic Spatial Signals



$$\mathbf{v} \in \mathbb{R}^3$$

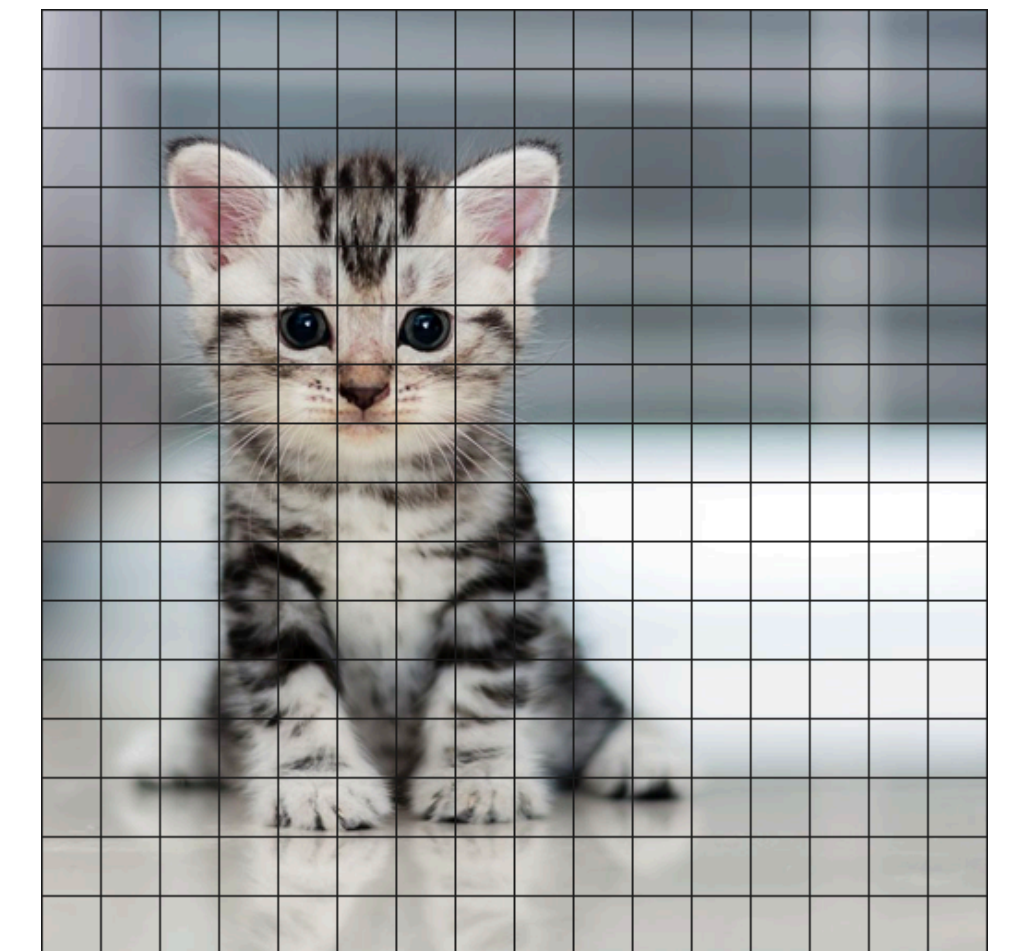
$$\mathbf{x} \in \mathbb{R}^2$$

Modeling Generic Spatial Signals



$$\mathbf{v} \in \mathbb{R}^1$$

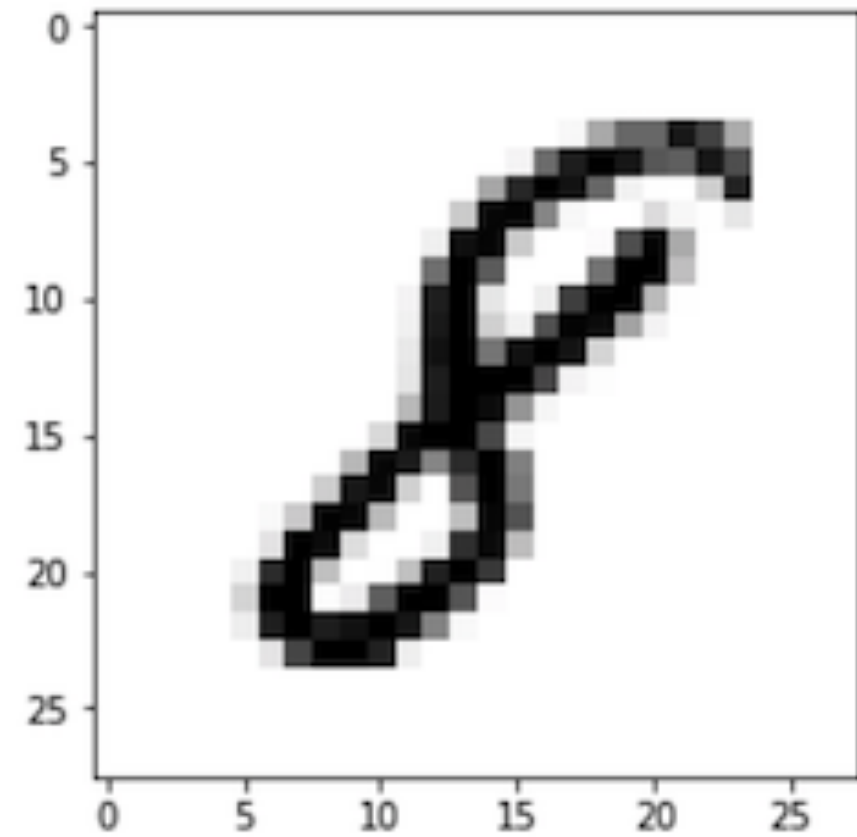
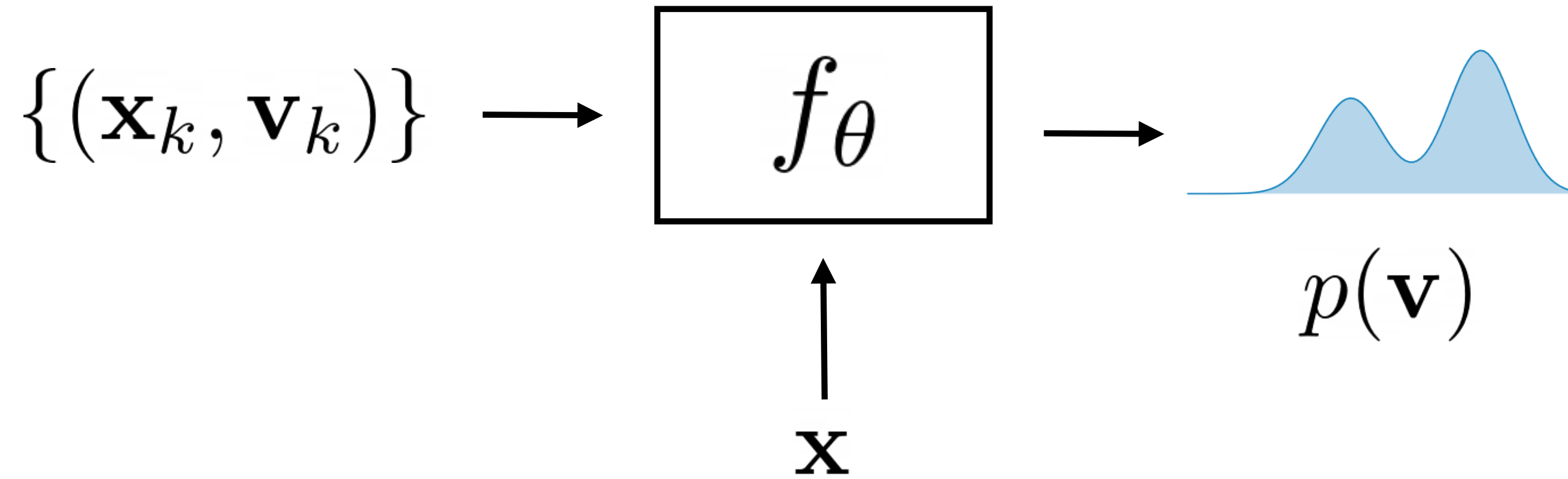
$$\mathbf{x} \in \mathbb{R}^2$$



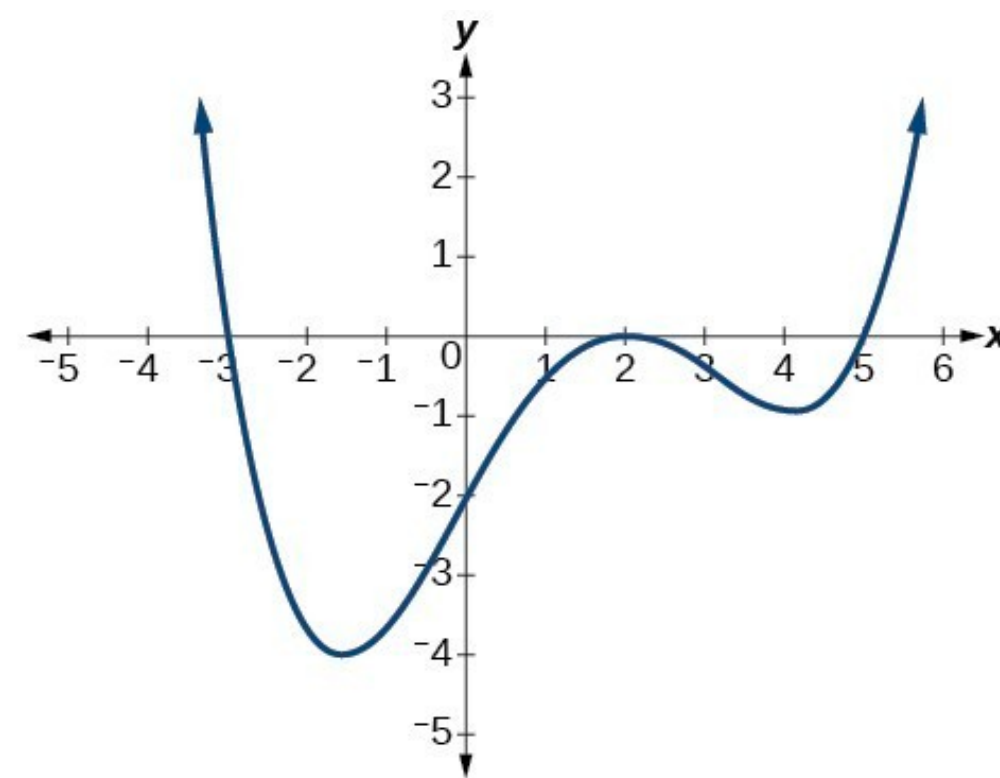
$$\mathbf{v} \in \mathbb{R}^3$$

$$\mathbf{x} \in \mathbb{R}^2$$

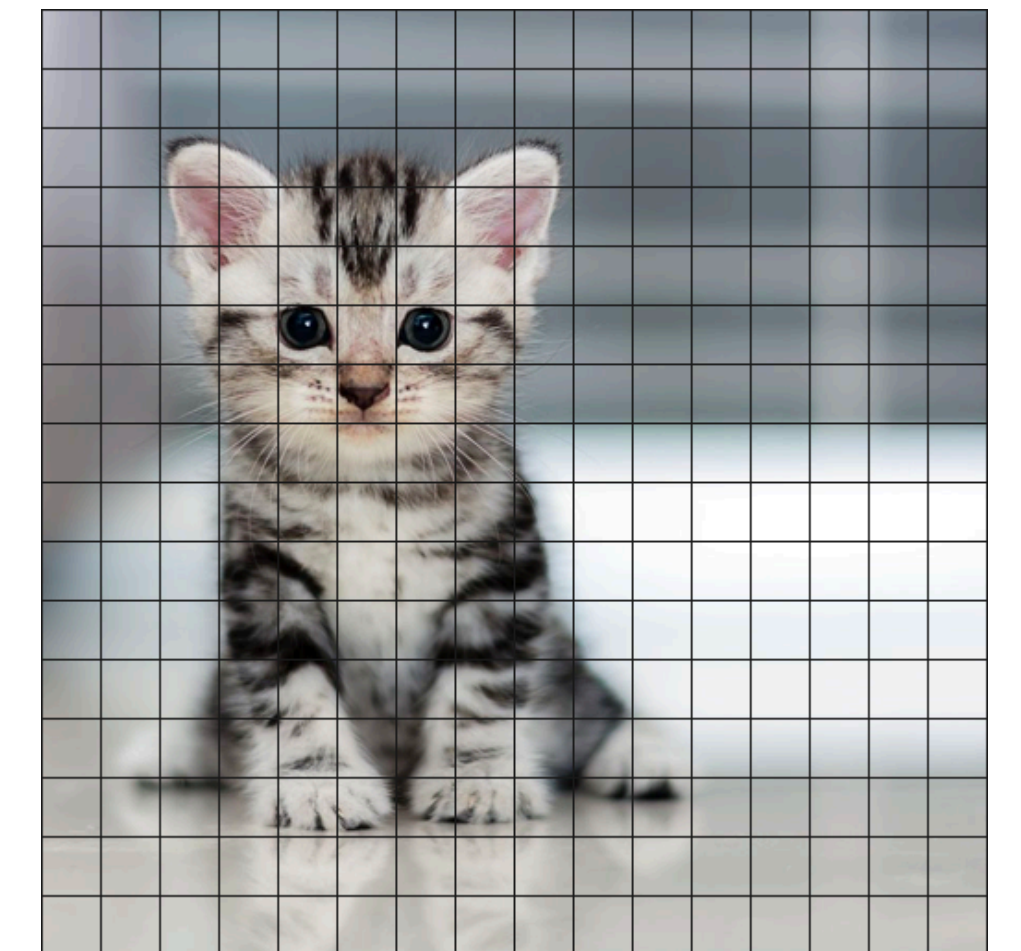
Modeling Generic Spatial Signals



$$\mathbf{v} \in \mathbb{R}^1$$
$$\mathbf{x} \in \mathbb{R}^2$$

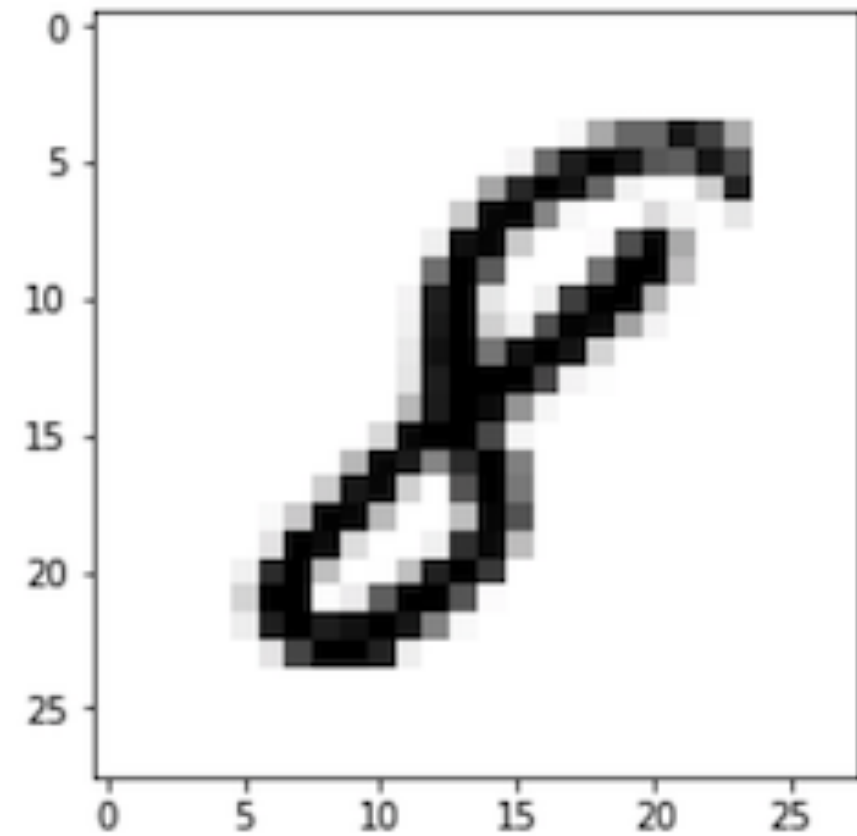
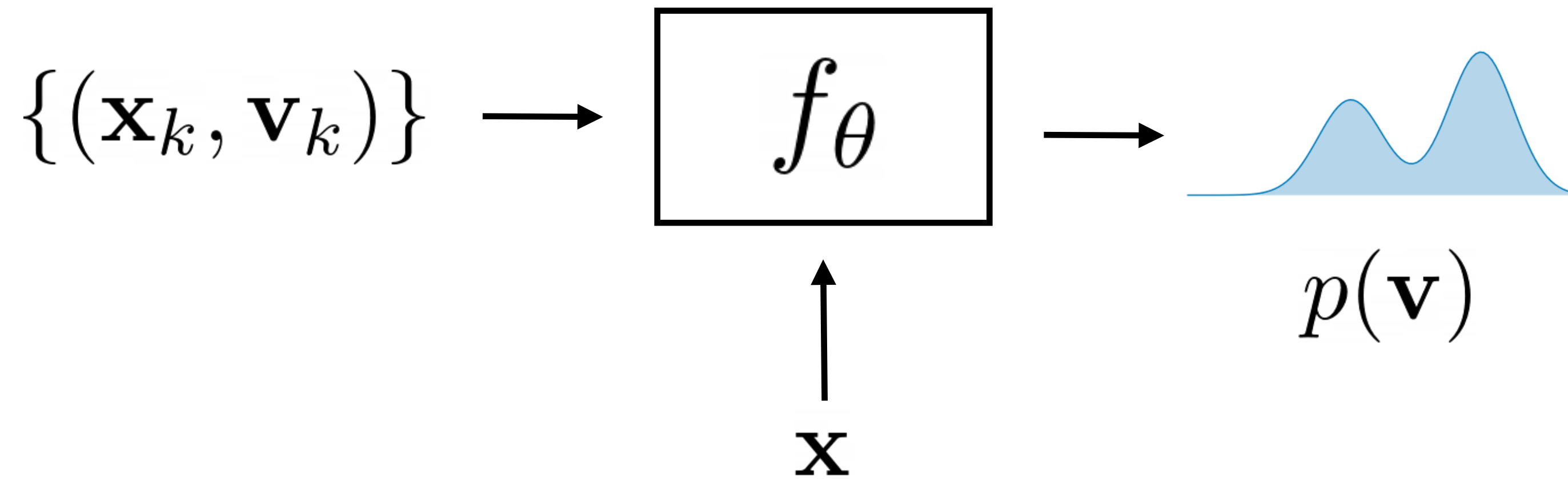


$$\mathbf{v} \in \mathbb{R}^1$$
$$\mathbf{x} \in \mathbb{R}^1$$

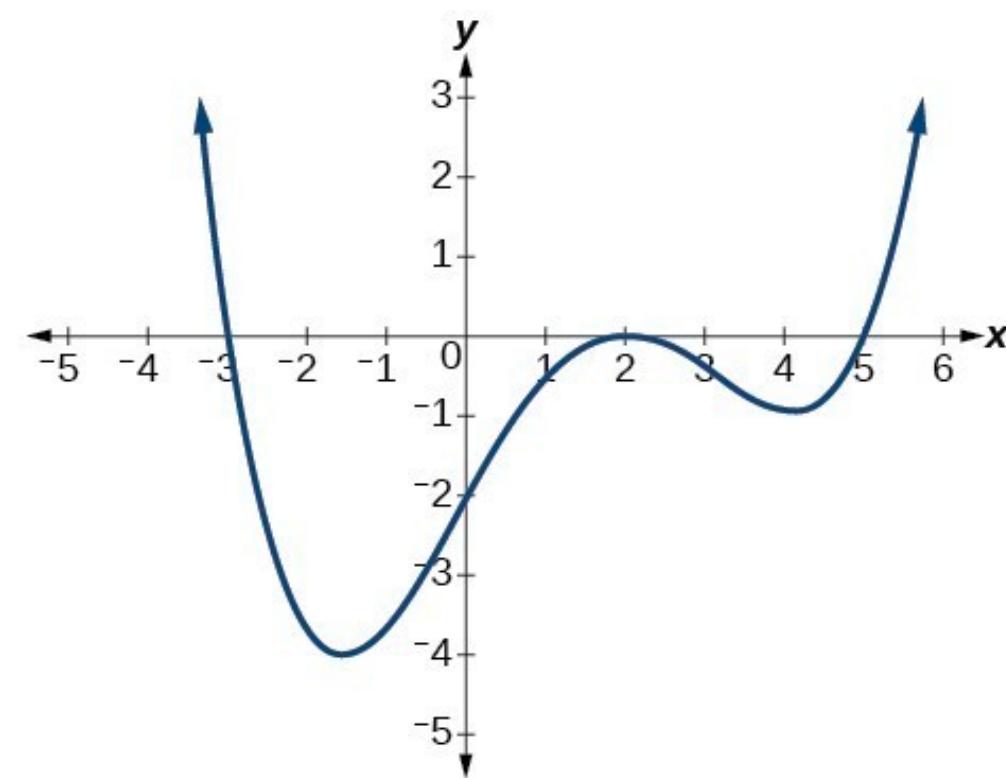


$$\mathbf{v} \in \mathbb{R}^3$$
$$\mathbf{x} \in \mathbb{R}^2$$

Modeling Generic Spatial Signals



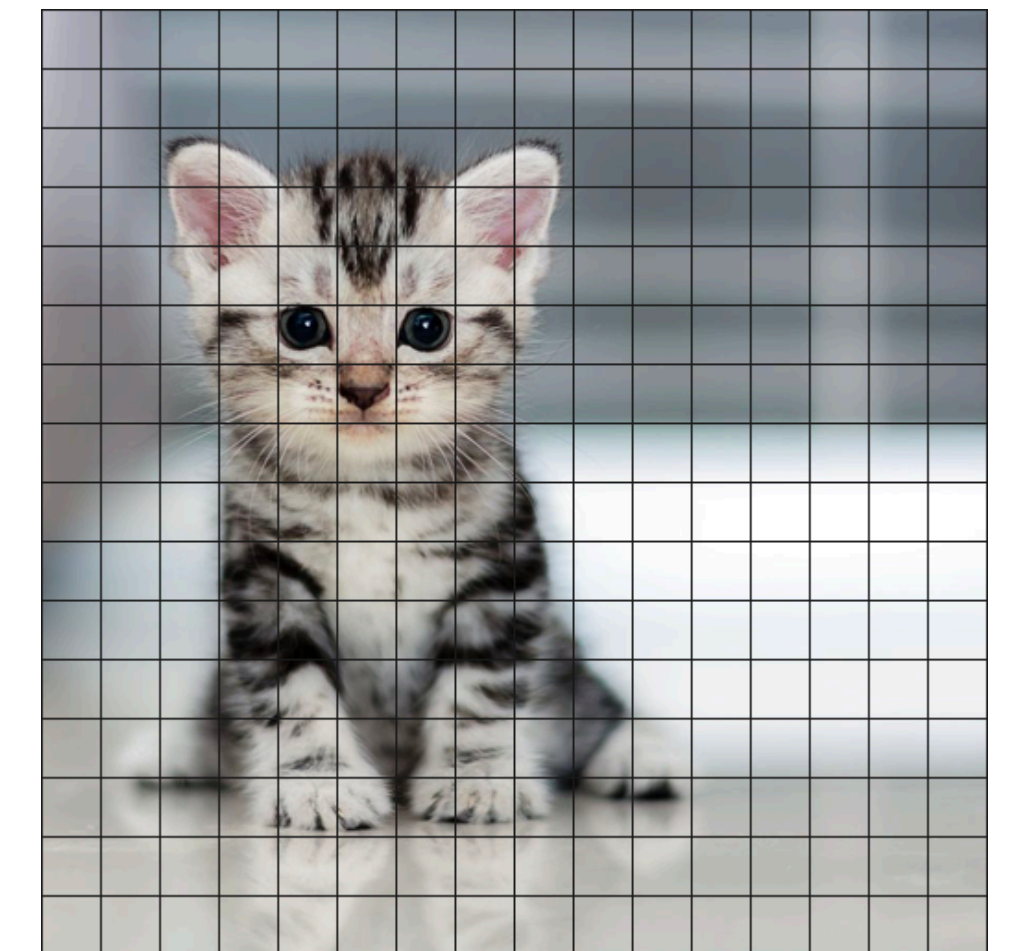
$$\mathbf{v} \in \mathbb{R}^1$$
$$\mathbf{x} \in \mathbb{R}^2$$



$$\mathbf{v} \in \mathbb{R}^1$$
$$\mathbf{x} \in \mathbb{R}^1$$



$$\mathbf{v} \in \mathbb{R}^1$$
$$\mathbf{x} \in \mathbb{R}^3$$



$$\mathbf{v} \in \mathbb{R}^3$$
$$\mathbf{x} \in \mathbb{R}^2$$

Image Synthesis



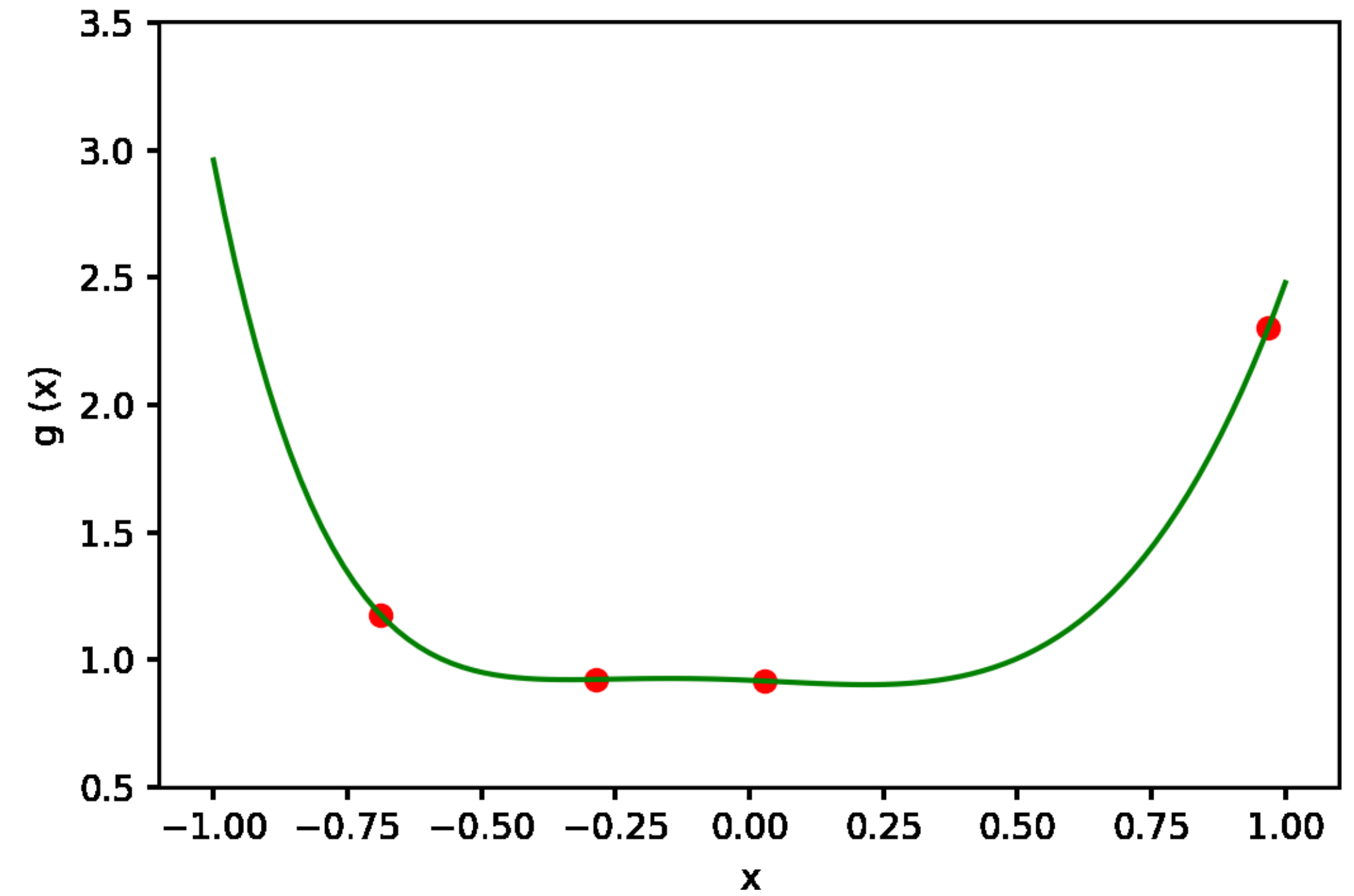
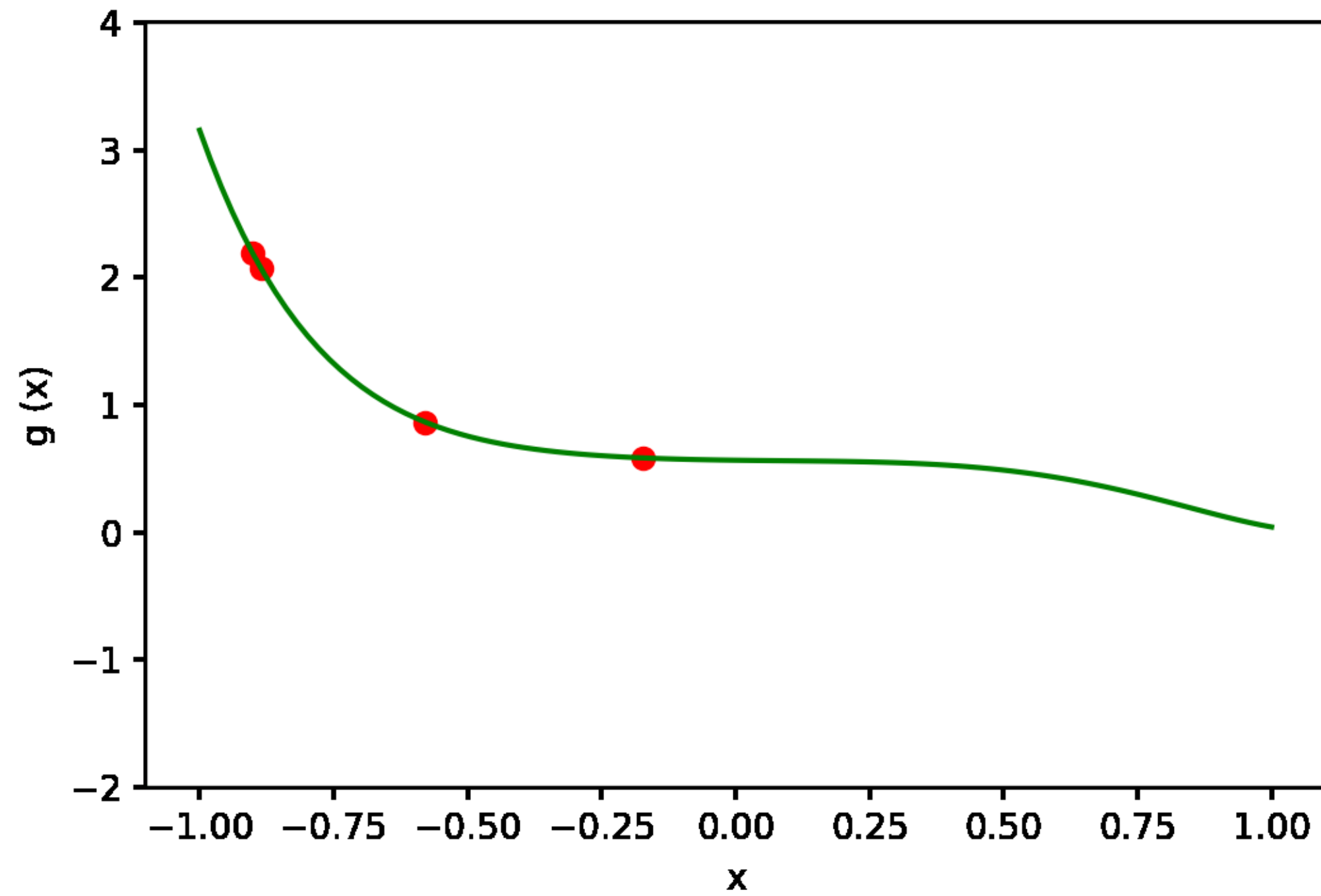
Image Synthesis



Image Synthesis

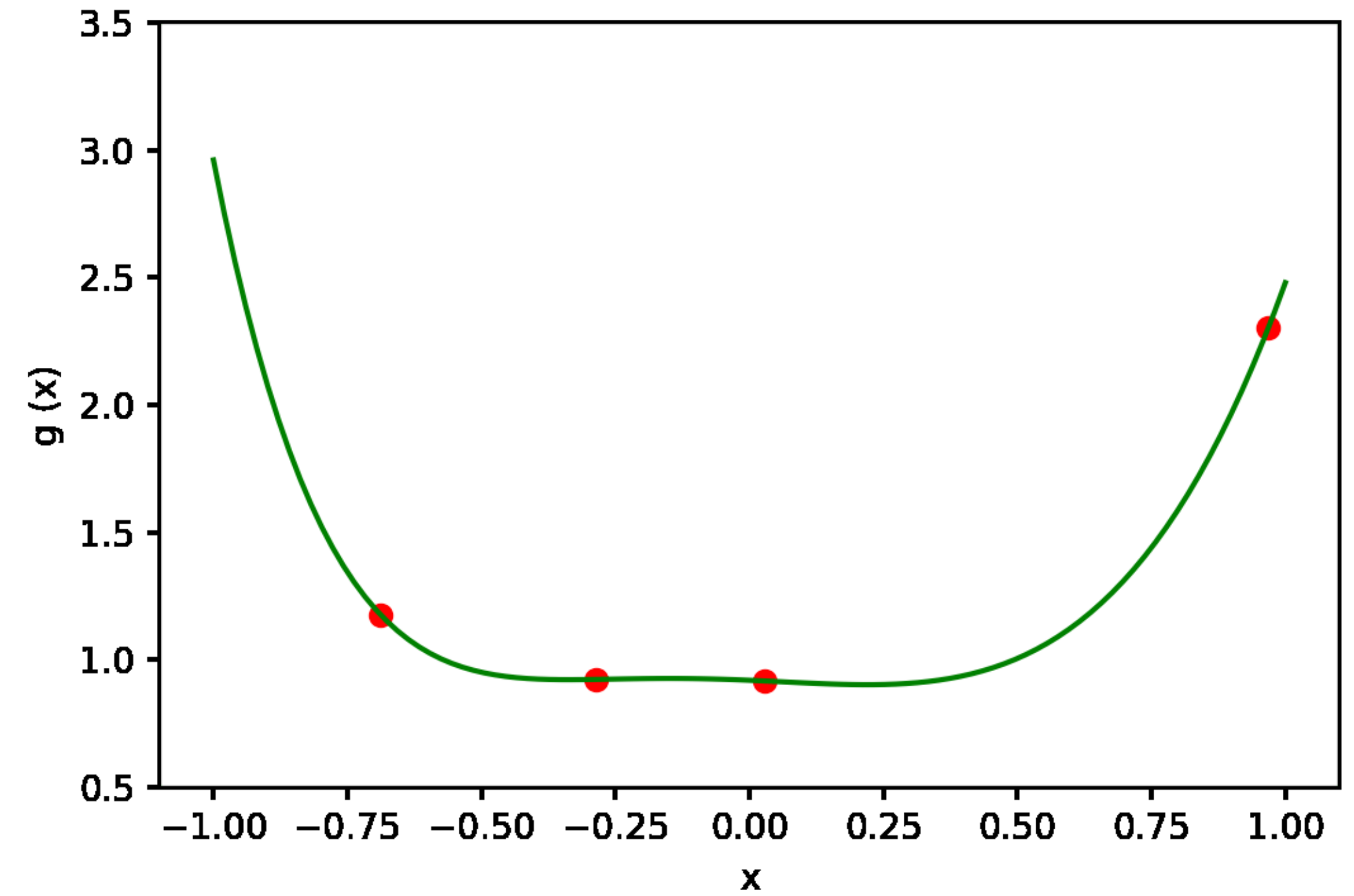
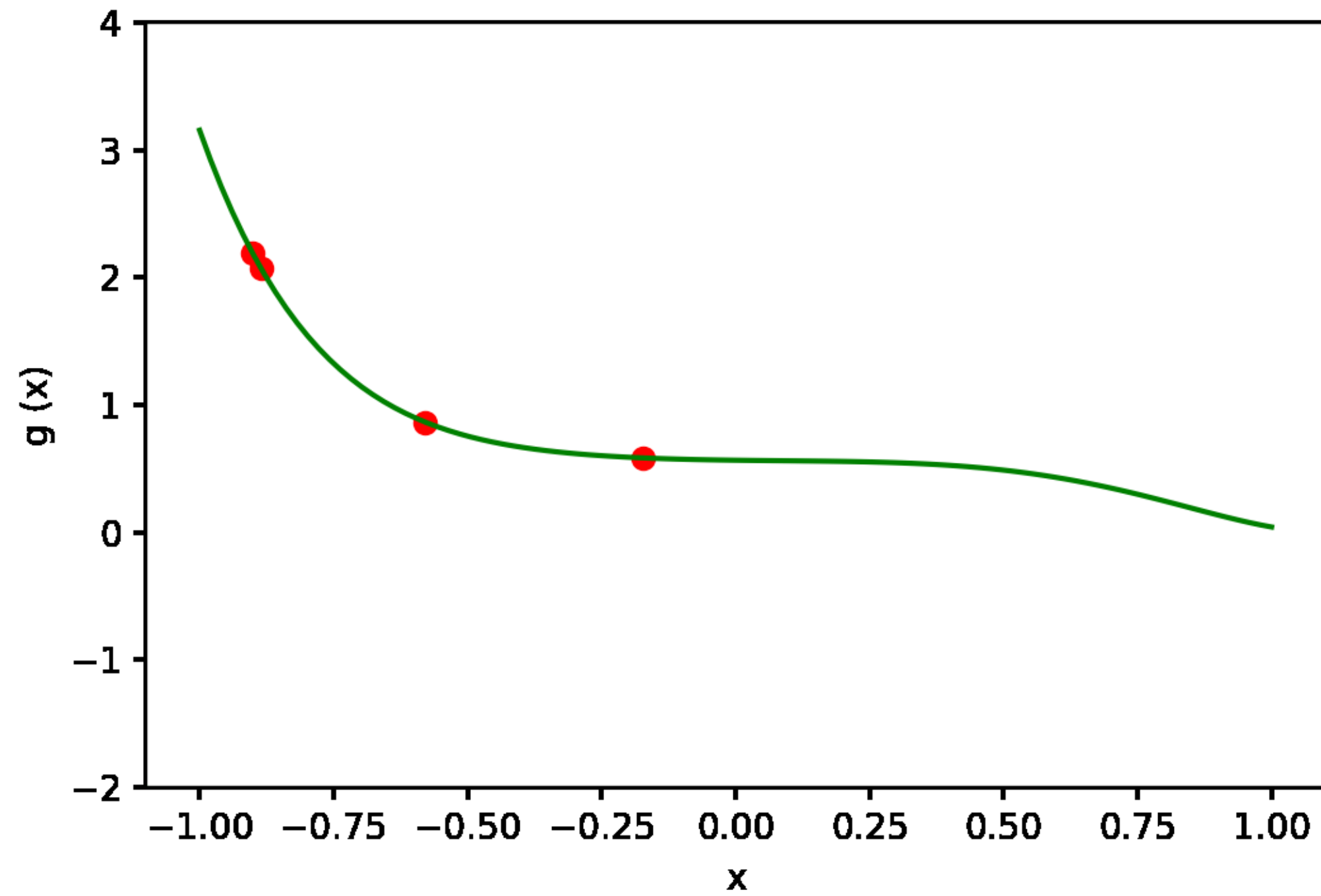


Polynomial Prediction



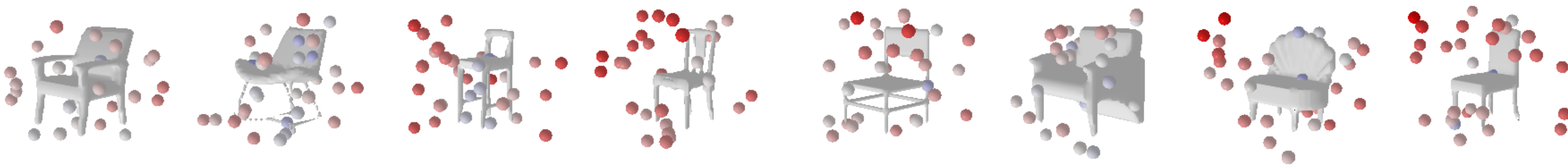
Training data: random degree-6 polynomials

Polynomial Prediction

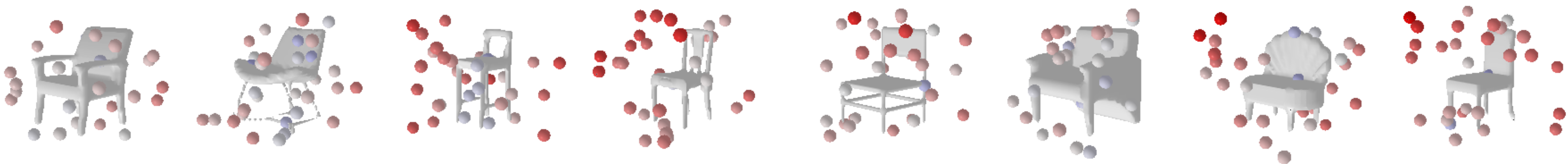


Training data: random degree-6 polynomials

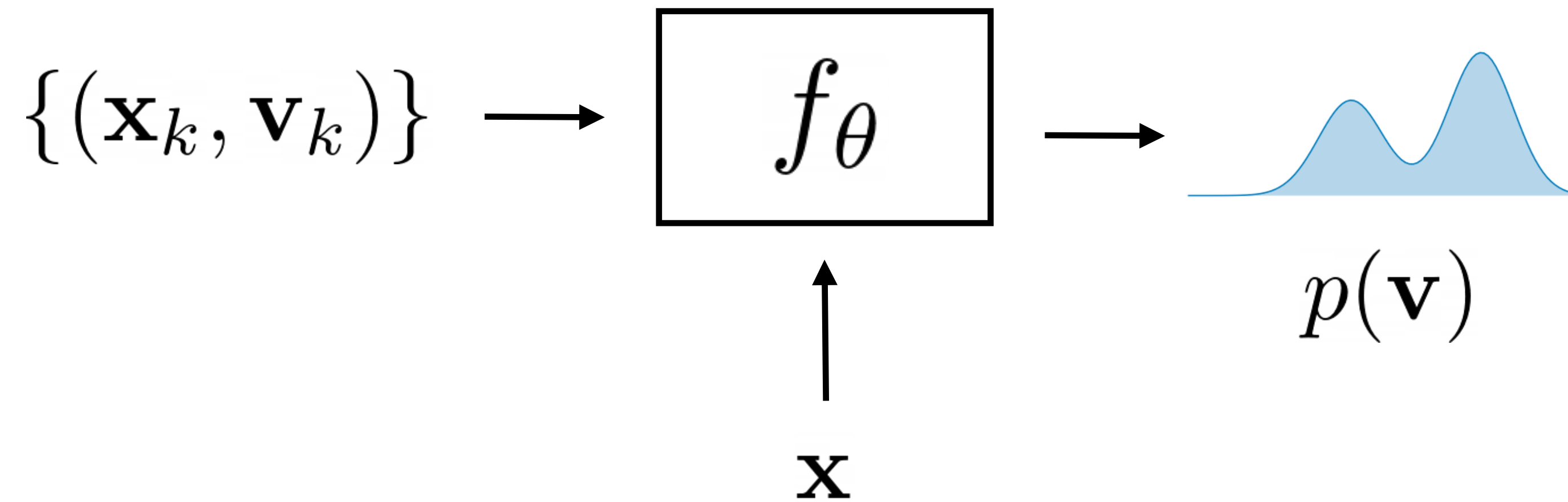
Shape Completion



Shape Completion



Thank you!



Project page: <https://shubhtuls.github.io/PixelTransformer/>