

# Neural Importance Sampling

by Müller et al., 2018

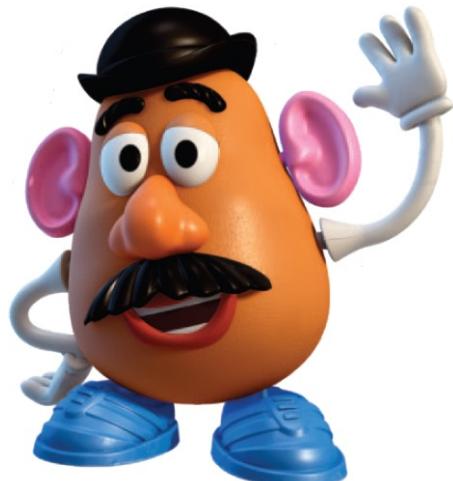
Jesus Solano and Federico D'Agostino

Machine Learning Seminar – Paper Presentation

University College London

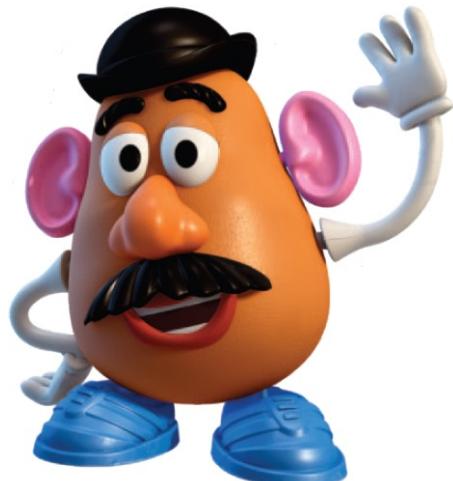
March 2022

# What do these two have in common?



$$I = \int_{\mathcal{P}} L_e(\mathbf{x}_0, \mathbf{x}_1) T(\bar{x}) W(\mathbf{x}_{k-1}, \mathbf{x}_k) d\bar{x}$$

# What do these two have in common?



$$I = \int_{\mathcal{P}} L_e(\mathbf{x}_0, \mathbf{x}_1) T(\bar{x}) W(\mathbf{x}_{k-1}, \mathbf{x}_k) d\bar{x}$$

And why does **Disney** need a ML research division?  
**... Let's first back up a bit to answer this**

# Contents

- What's in an animated motion picture movie?
- The problem of integration
- Integrating with neural networks: NICE (***Nonlinear Independent Components Estimation***, Dinh et al. 2014)
- Extending NICE: Neural Importance Sampling
- Experiments and applications to rendering
- Evaluation and outlook

**Nonlinear Independent Components Estimation**  
Dinh, Krueger and Bengio, 2014

# How are cartoons made?



Pre-production → Modelling → Texturing → ... → **Rendering** → Post-production

# Light-Transport simulation



- One of the main challenges in computer graphics

# Path Sampling and Path Guiding

Reduce estimation variance by "guiding" the construction of light paths



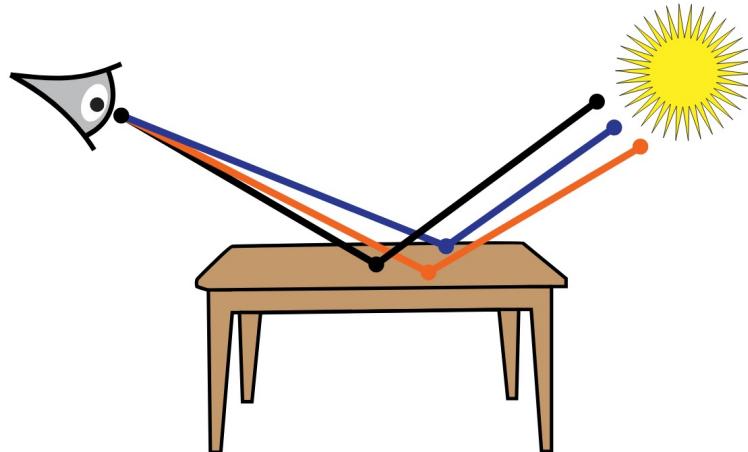
## Path Sampling

## Path Guiding

# Primary-Sample-Space Path Sampling

**GOAL**

Estimate amount of light reaching the camera after taking **ANY** of the possible paths



**Approach**

Integrating over all paths in the scene

Path Throughput

$$I = \int_{\mathcal{P}} L_e(\mathbf{x}_0, \mathbf{x}_1) T(\bar{x}) W(\mathbf{x}_{k-1}, \mathbf{x}_k) d\bar{x}$$

Emitted Radiance

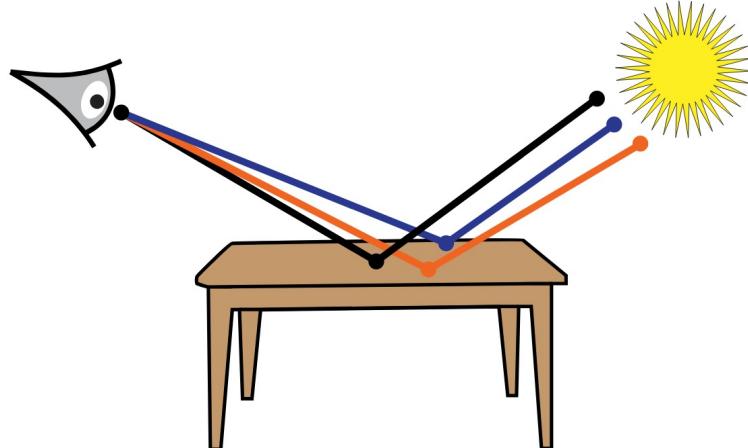
Camera Sensor response

Chain of positions of a **single light path** with  $k$  vertices

# Primary-Sample-Space Path Sampling

**GOAL**

Estimate amount of light reaching the camera after taking **ANY** of the possible paths



**Approach**

Integrating over all paths in the scene

$$I = \int_{\mathcal{P}} L_e(\mathbf{x}_0, \mathbf{x}_1) T(\bar{x}) W(\mathbf{x}_{k-1}, \mathbf{x}_k) d\bar{x}$$

**How**

Use Monte Carlo Integration \*

$$\langle I \rangle = \frac{1}{N} \sum_{j=1}^N \frac{L_e(\mathbf{x}_{j0}, \mathbf{x}_{j1}) T(\bar{x}_j) W(\mathbf{x}_{jk-1}, \mathbf{x}_{jk})}{q(\bar{x}_j)}$$

Joint probability density of generating all k vertices of a path

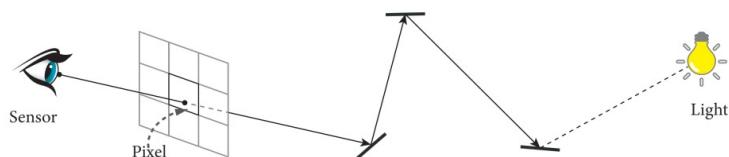
# Path Guiding

**GOAL**

Estimate radiative equilibrium of every surface point  $x$  in the direction  $\omega_0$

**Approach**

Integrating over all possible reflecting angles



Emitted Radiance

$$L_o(x, \omega_0) = L_e(x, \omega_0) + \int_{\Omega} L(x, \omega) f_s(x, \omega_0, \omega) |\cos \gamma| d\omega$$

Reflected Radiance

Incident Radiance

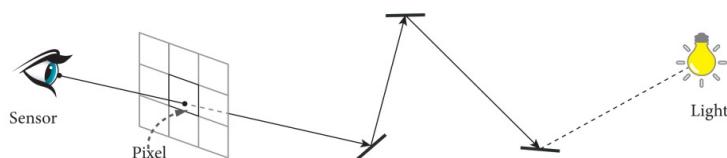
Camera Sensor response

Angle between  $w$  and surface normal

# Path Guiding

**GOAL**

Estimate radiative equilibrium of every surface point  $\mathbf{x}$  in the direction  $\omega_o$



**Approach**

Integrating over all possible reflecting angles

$$L_o(\mathbf{x}, \omega_o) = L_e(\mathbf{x}, \omega_o) + \int_{\Omega} L(\mathbf{x}, \omega) f_s(\mathbf{x}, \omega_o, \omega) |\cos \gamma| d\omega$$

**How?**

Use Monte Carlo Integration \*

$$\langle L_r(\mathbf{x}, \omega_o) \rangle = \frac{1}{N} \sum_{j=1}^N \frac{L(\mathbf{x}, \omega_j) f_s(\mathbf{x}, \omega_o, \omega_j) |\cos \gamma_j|}{q(\omega_j | \mathbf{x}, \omega_o)}$$

Density conditioned on position and angle (>5D)

# In summary...

$$L_o(\mathbf{x}, \omega_o, \lambda, t) = L_e(\mathbf{x}, \omega_o, \lambda, t) + \int_{\Omega} f_r(\mathbf{x}, \omega_i, \omega_o, \lambda, t) L_i(\mathbf{x}, \omega_i, \lambda, t) (\omega_i \cdot \mathbf{n}) d\omega_i$$

$$I = \int_{\mathcal{P}} L_e(\mathbf{x}_0, \bar{x}) T(\bar{x}) W(\mathbf{x}_{k-1}, \mathbf{x}_k) d\bar{x}$$

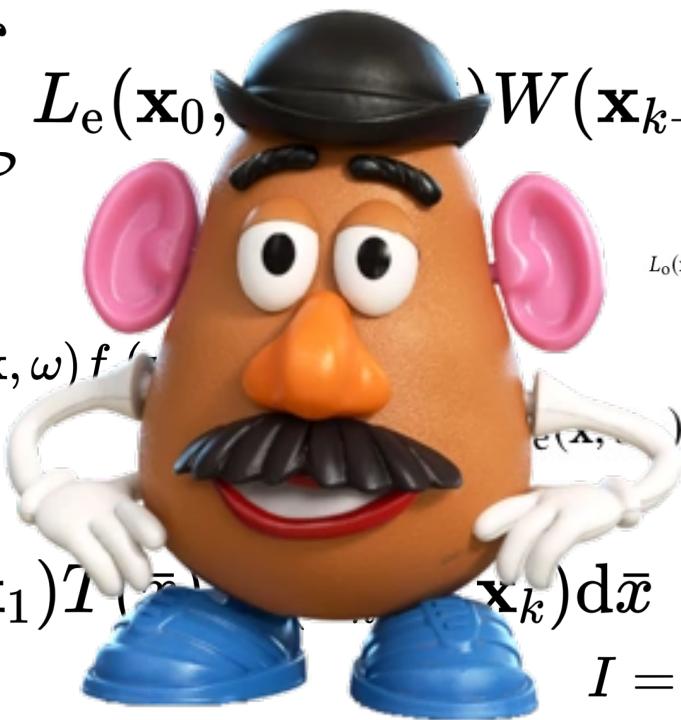
$$L_o(\mathbf{x}, \omega_o) = L_e(\mathbf{x}, \omega_o) + \int_{\Omega} L(\mathbf{x}, \omega) f_s(\mathbf{x}, \omega_o, \omega) | \cos \gamma | d\omega$$

$$L_o(\mathbf{x}, \omega_o) = L_e(\mathbf{x}, \omega_o) + \int_{\Omega} L(\mathbf{x}, \omega) f_s(\mathbf{x}, \omega_o, \omega) | \cos \gamma | d\omega$$

$$I = \int_{\mathcal{P}} L_e(\mathbf{x}_0, \mathbf{x}_1) T(\bar{x}) W(\mathbf{x}_{k-1}, \mathbf{x}_k) d\bar{x}$$

$$I = \int_{\mathcal{P}} L_e(\mathbf{x}_0, \mathbf{x}_1) T(\bar{x}) W(\mathbf{x}_{k-1}, \mathbf{x}_k) d\bar{x}$$

**Integrals everywhere!**



# The problem of integration

- Recall from previous lectures: integration is everywhere! (unfortunately)
  - Not only in ML!
- Analytical solutions may be intractable or non-existent
- Dimensionality only adds to the problem
  - In any strategy we previously encountered as the dimensionality of the integrand increases the approximation deteriorates
- Some approximation strategies are situation-specific
  - We can work our way around that...

# Background: MC with importance sampling

- The usual **importance sampling scenario**:

$$\mathbb{E}[f(x)] = \int f(x)p(x)dx = \int f(x)\frac{p(x)}{q(x)}q(x)dx \approx \frac{1}{N} \sum_{i=1}^N f(X_i)\frac{p(X_i)}{q(X_i)}$$

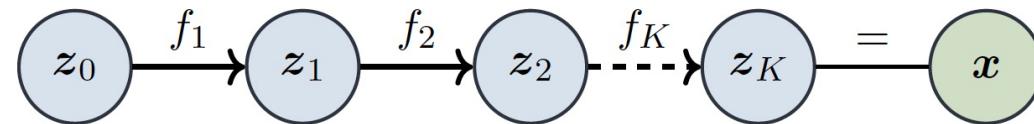
Where  $f(x)$  follows the distribution  $p(x)$ , and  $q(x)$  is a proposal distribution.

- **Our scenario:**

$$F = \int_{\mathcal{D}} f(x)dx = \int_{\mathcal{D}} \frac{f(x)}{q(x)}q(x)dx = \mathbb{E}\left[\frac{f(X)}{q(X)}\right] \approx \langle F \rangle_N = \frac{1}{N} \sum_{i=1}^N \frac{f(X_i)}{q(X_i)}$$

# Background: Normalizing Flows (NFs)

- Normalizing flows:



with  $z_0 \sim p_0$ , a “simple” base distribution and  $z_k = f_k(z_{k-1})$

Through the **change of variables trick**, we can express this in terms of probability distributions (*where we swapped  $f$  with  $h$  to be consistent with our previous notation*):

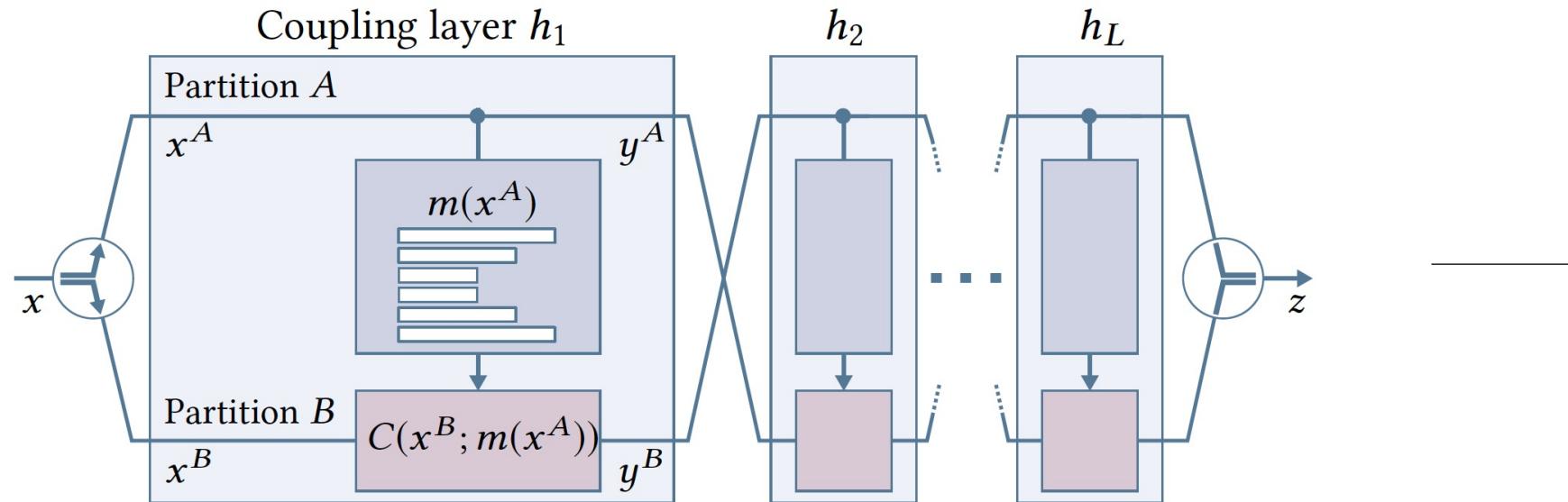
$$q(x; \theta) = \int \delta(x - h(z; \theta)) q(z) dz = q(z) \left| \det \left( \frac{\partial h(z; \theta)}{\partial z^T} \right) \right|^{-1}$$

# What do we want from an ideal integration algorithm?

- A parametric proposal distribution
  - Additional MC over  $q$  introduces bias
- Applicability to high dimensionality
  - i.e. an expressive proposal
- Efficient sampling and evaluation of the proposal
- If using NFs, a mapping that is easily invertible
  - i.e. whose Jacobian determinant can be computed efficiently

# MC Integration with NICE: Neural Importance Sampling

$$\langle F \rangle_N = \frac{1}{N} \sum_{i=1}^N \frac{f(X_i)}{q(X_i; \theta)} = \frac{1}{N} \sum_{i=1}^N \frac{p(X_i)F}{\boxed{q(X_i; \theta)}}$$



Let's break this down component by component

# Coupling Layer

## The trick to add Neural Networks

The **cost** of computing the determinant grows **super-linearly**

What can be done to reduce the computational cost?

Use transformations whose Jacobian is diagonal, lower **triangular** or upper **triangular**.

How to use this observation in Normalizing Flows and MC integration?

1

Neural network with triangular weight matrices and bijective activation functions

**BUT** this will constrain too much the possible architectures!

2

**Coupling Layers**

Family of functions with **triangular Jacobian**

# Coupling Layer

## The trick to add Neural Networks

What is a  
**coupling layer?**

Specific family of transformations whose determinant reduce to the product of diagonal terms

1

Split the input vector into two  
**disjoint partitions** ( $x^A, x^B$ ) of its dimensions



2

Apply the coupling transform  $C$  to ONLY one of the partitions



$$\begin{aligned}y^A &= x^A \\y^B &= C(x^B; m(x^A))\end{aligned}$$

3

Build output of **coupling layer** by concatenation



$$h(x) = y = (y^A, y^B)$$

# Coupling Layer

The trick to add Neural Networks

**Coupling  
Layer**

$$\begin{aligned}y^A &= x^A \\y^B &= C(x^B; m(x^A))\end{aligned}$$

$$\begin{aligned}y &= (y^A, y^B) \\&= h(x)\end{aligned}$$

Trivial inversion of  $h$ :

$$\begin{aligned}x^A &= y^A \\x^B &= C^{-1}(y^B; m(x^A)) = C^{-1}(y^B; m(y^A))\end{aligned}$$

What are the **basic** Coupling transforms?

1

**Additive  
Coupling Transform**

Translation of signal

$$C(x^B; t) = x^B + t$$

2

**Multiply-Add  
Coupling Transform**

Element-wise product plus translation

$$C(x^B; s, t) = x^B \odot e^s + t$$

# NICE ideas, summarized

**Nonlinear Independent Components Estimation**  
Dinh, Krueger and Bengio, 2014

- Composable transformations: 
$$h = h_i \circ \dots \circ h_2 \circ h_1$$
- Disjoint partitions of :

$$h_i(x) = y = (y^A, y^B) \longrightarrow y^A = x^A \quad y^B = C(x^B; m(x^A))$$

With:  $C : \mathbb{R}^{|B|} \times m(\mathbb{R}^{|A|}) \rightarrow \mathbb{R}^{|B|}$

- Tractable inverses thanks to smart use of disjoint partitions:

$$x^A = y^A$$

$$x^B = C^{-1}(y^B; m(x^A)) = C^{-1}(y^B; m(y^A))$$

# Piecewise-Polynomial Coupling Layers

Coupling transforms ( $C$ ) could be **more complex!**

$$C(x^B; m(x^A)) = (C_1(x_1^B; m(x^A)), \dots, C_{|B|}(x_{|B|}^B; m(x^A)))^T$$

Transform each dimension independently



We can interpretate this warping function as a *cumulative distribution function* (CDF)

How to produce each  $C_i$  ?

Neural Network to output the corresponding unnormalized probability density ( $q_i$ )

Construct  $C_i$  by integration

1

Piecewise-Linear Coupling Transform

$$C_i(x_i^B; Q) = \int_0^{x_i^B} q_i(t) dt = \alpha Q_{ib} + \sum_{k=1}^{b-1} Q_{ik}$$

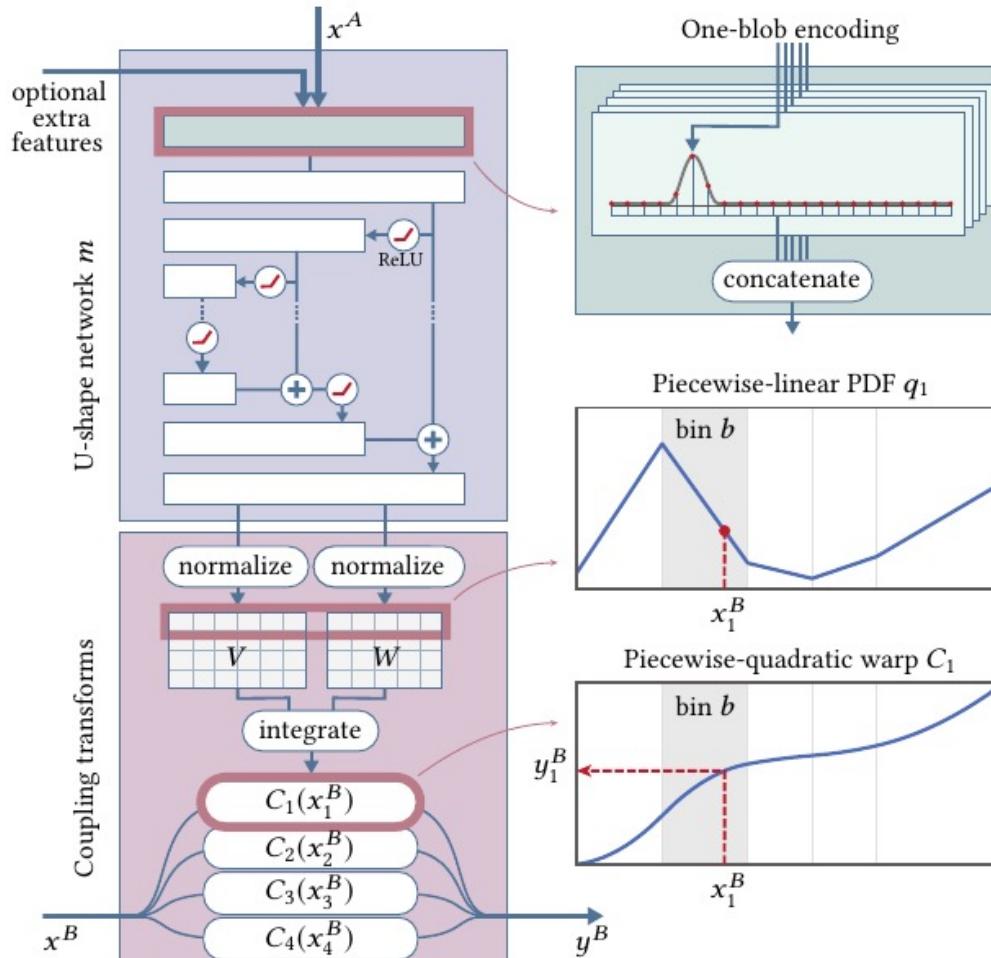
2

Piecewise-Quadratic Coupling Transform

$$C_i(x_i^B; W, V) = \frac{\alpha^2}{2}(V_{ib+1} - V_{ib})W_{ib} + \alpha V_{ib}W_{ib} + \sum_{k=1}^{b-1} \frac{V_{ik} + V_{ik+1}}{2}W_{ik}$$

# Piecewise-Polynomial Coupling Layers

How does the whole **coupling layer** look like?

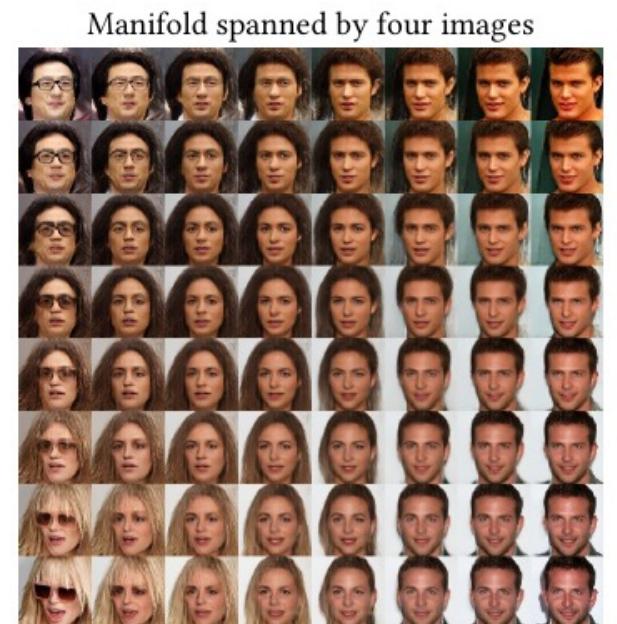
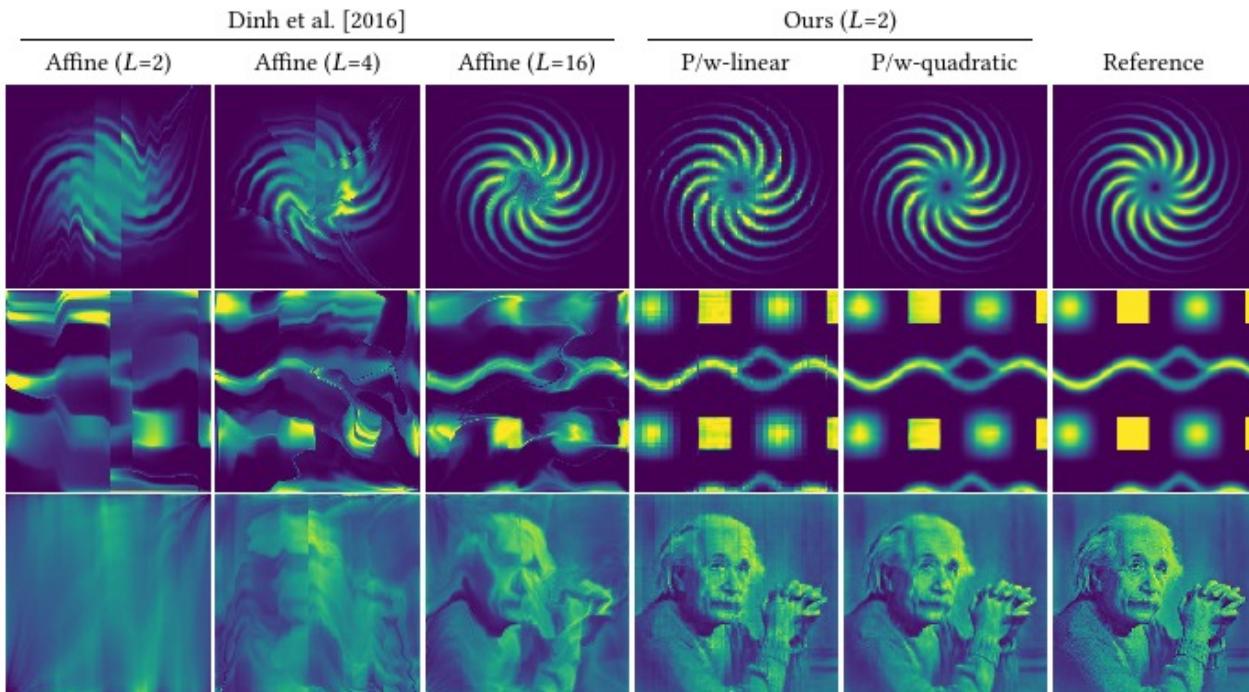


Generalization of one-hot encoding

A kernel is used to activate multiple adjacent matrix instead of a single one

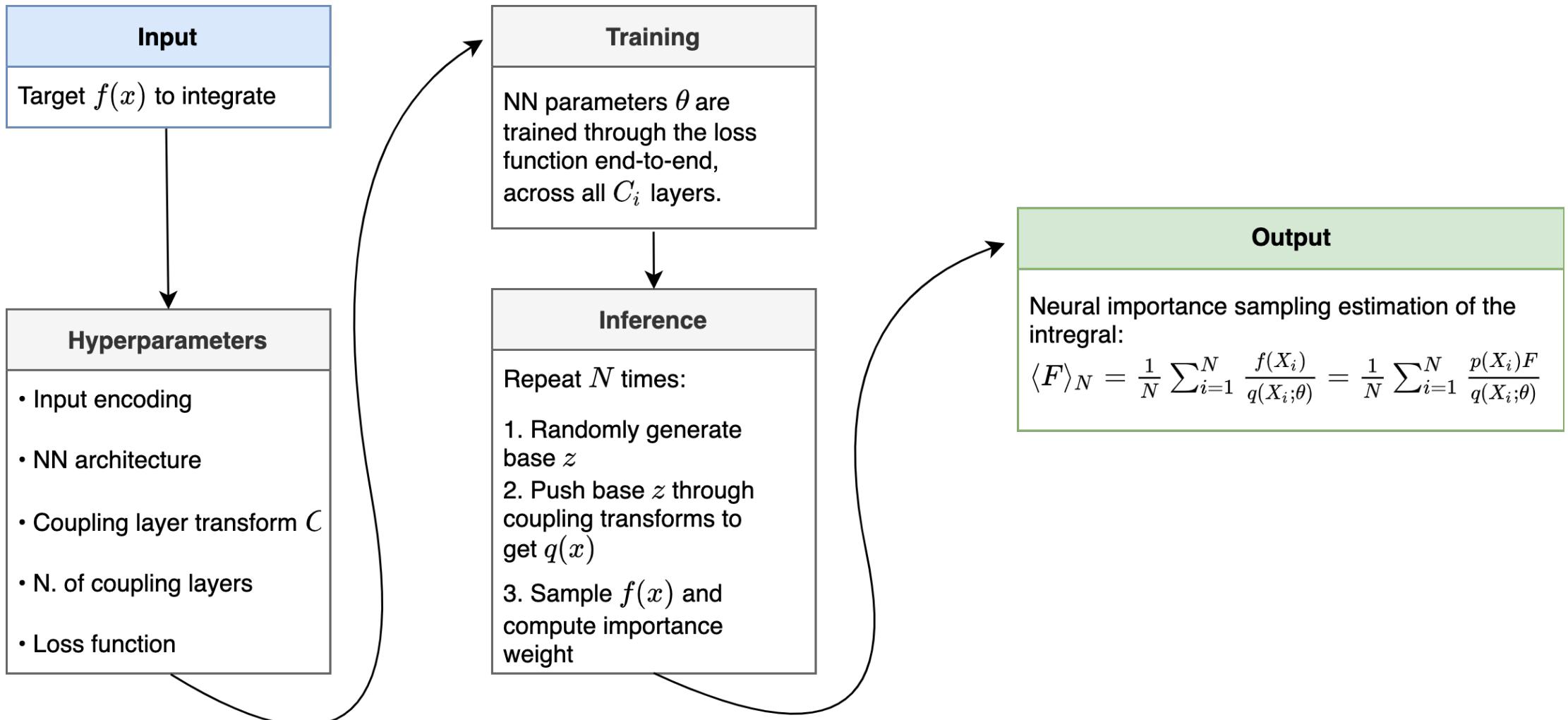
# Performance of Coupling Layers

Piecewise coupling layers achieve superior performance compared with affine transforms on **low-dimensional** regression problems

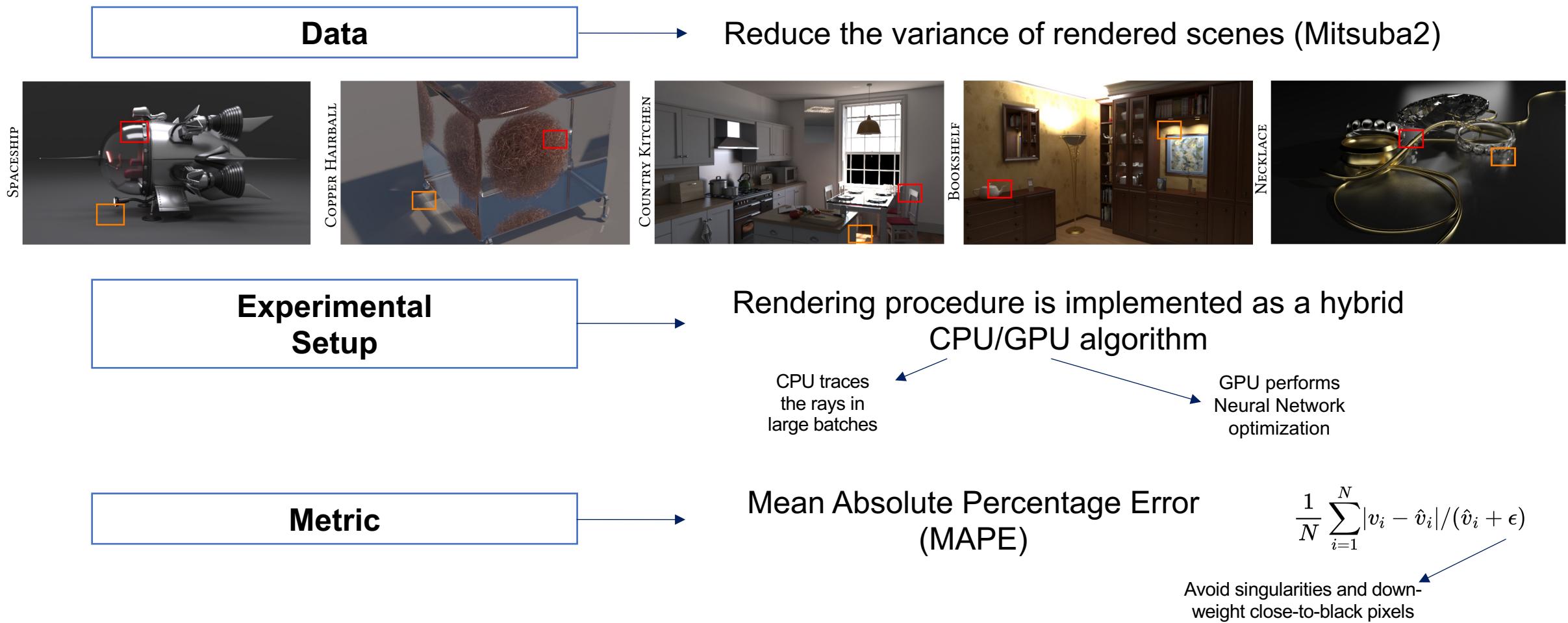


Learning CelebFaces attributes with  
Piecewise-quadratic coupling layers:  
a **high-dimensional** density-estimation problem

# The full algorithm

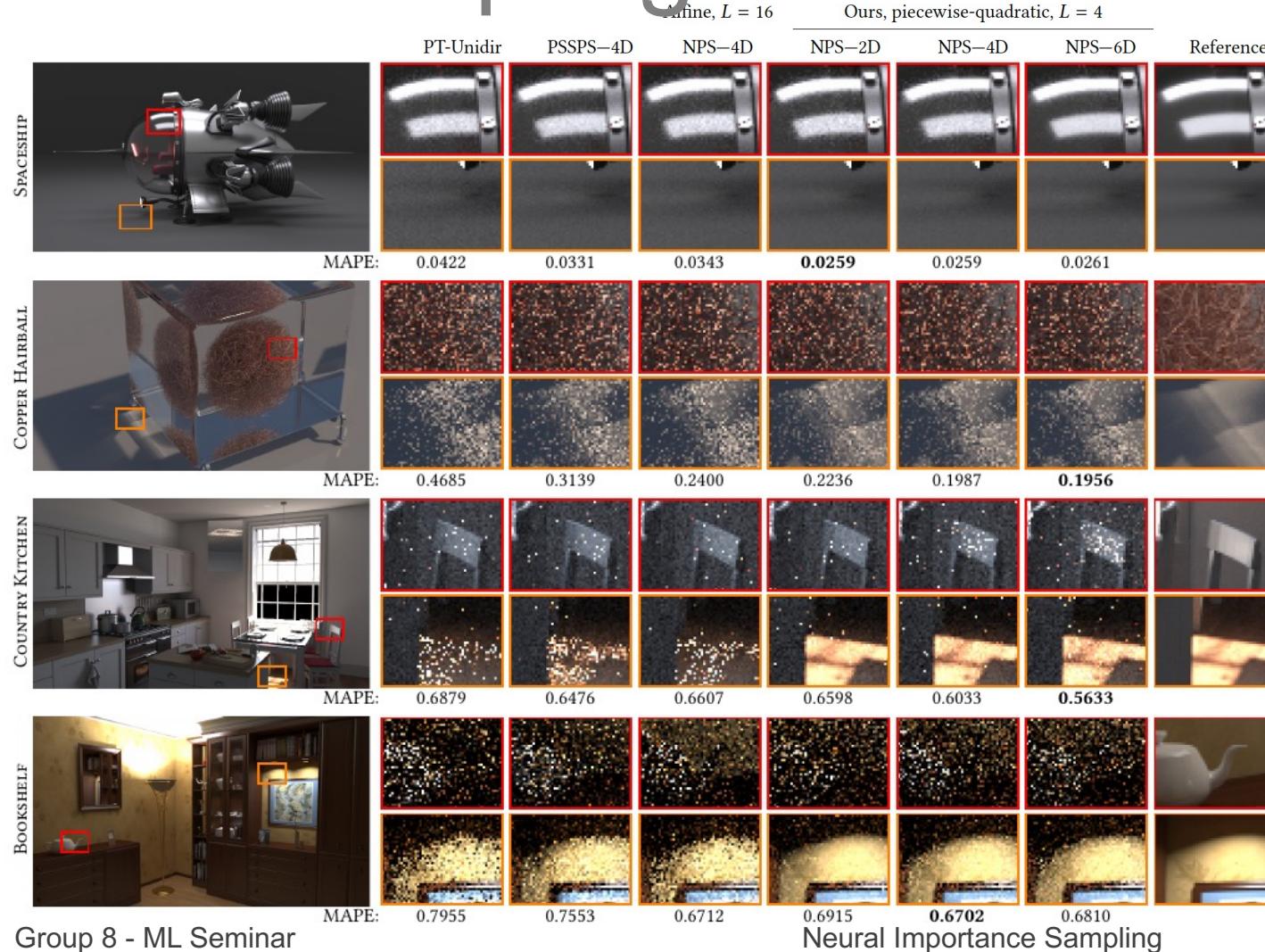


# Evaluation



# Results

## Path Sampling



1 **Neural path sampling outperforms** previous approaches in terms of MAPE, **variance is reduced** for challenging scenes

2 **Neural path sampling only offers improvement beyond the 4D case if paths stay Coherent**

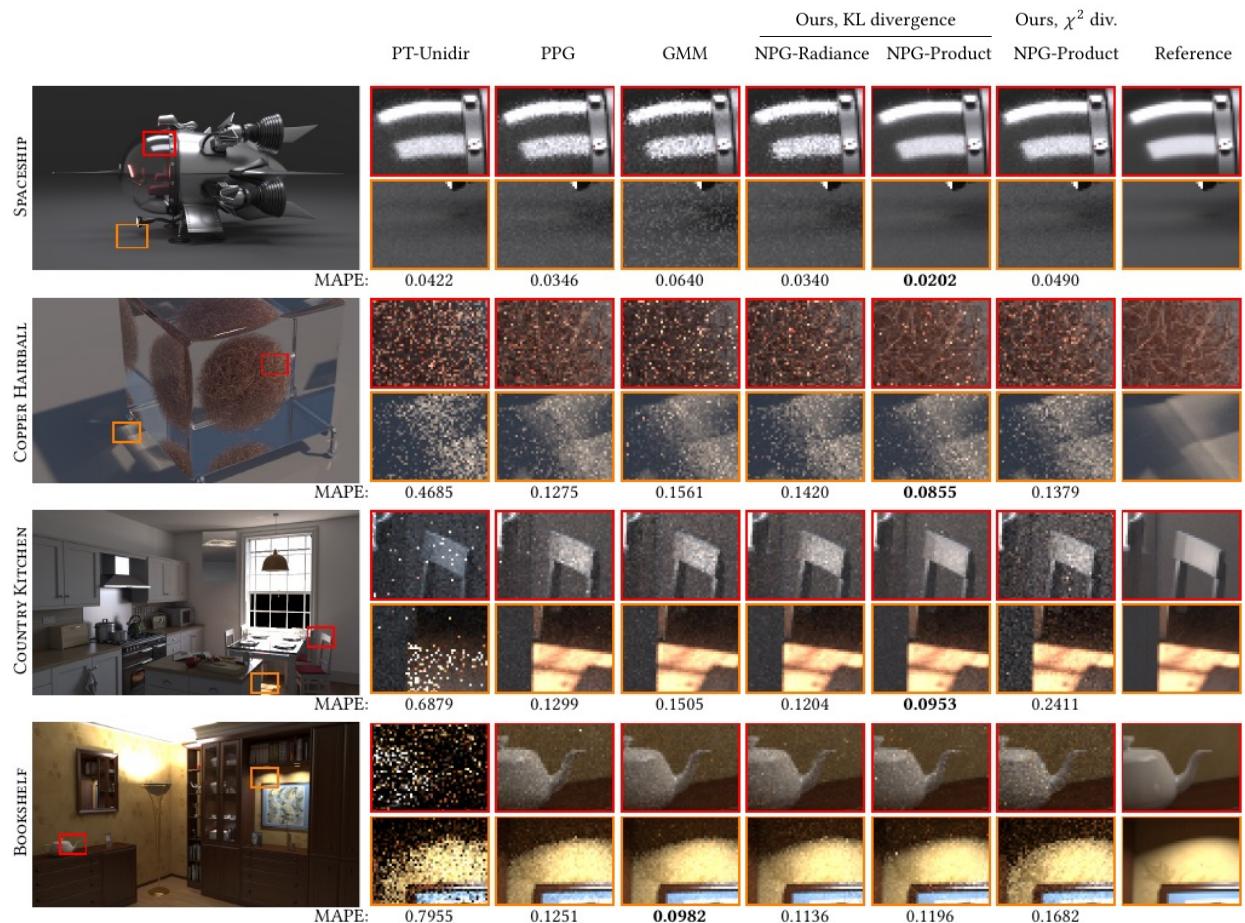
# Results

## Path Guiding

**Neural path guiding** outperforms previous models even without considering texture properties of materials (*NPG-radiance*), but only overall radiance.

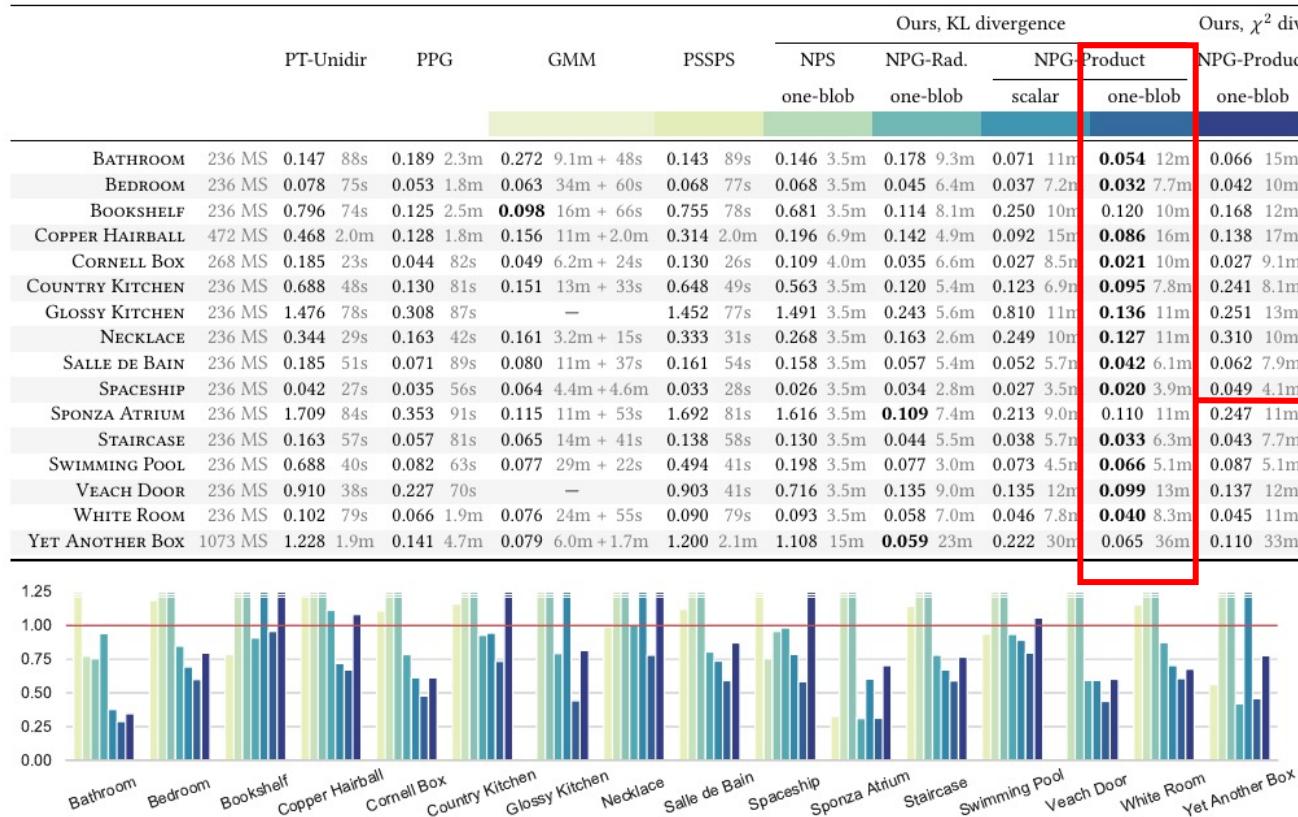
1

**Neural path guiding** performs even better while accounting for texture properties. This is called *NPG-Product* because it applies multiple importance sampling to combine different densities capturing properties of the materials



# Advantages of Neural Importance Sampling

- A collection of **Neural Networks** can parametrize **complex sampling densities**.
- Piecewise-polynomial transforms further **improve overall expressiveness**
- The **performance gain** in challenging scenarios could be **big**.



Neural path sampling performs better in almost every scene

# Drawbacks

- Expensive compared to other rendering models
  - A different neural network per coupling layer
  - More expensive coupling transforms than previous research
- Application scenarios are limited
  - Need to satisfy cost/benefits trade-off
  - Prohibitive for simple tasks
- Does not extend to multiple integrals
  - Cannot neglect the normalising factor
- More difficult to examine failure cases of the model
  - Relevant in the rendering case

$$\langle F \rangle_N = \frac{1}{N} \sum_{i=1}^N \frac{f(X_i)}{q(X_i; \theta)} = \frac{1}{N} \sum_{i=1}^N \frac{p(X_i)F}{q(X_i; \theta)}$$

# Let's do a fair comparison...

Neural Importance Sampling **performs better....**

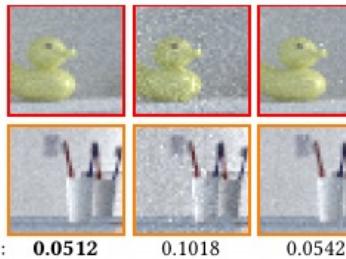
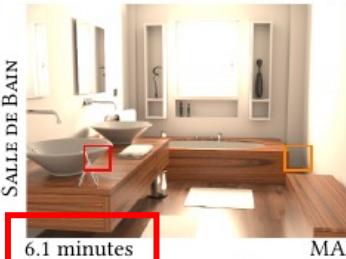
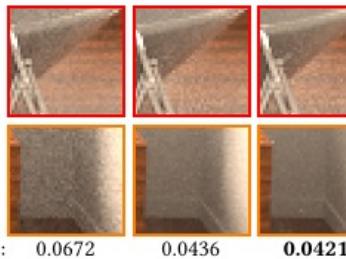
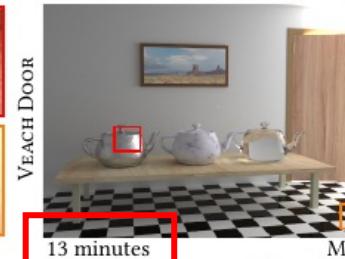
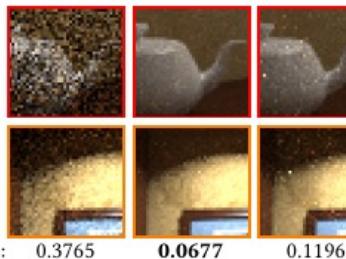
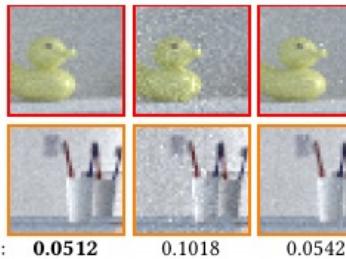
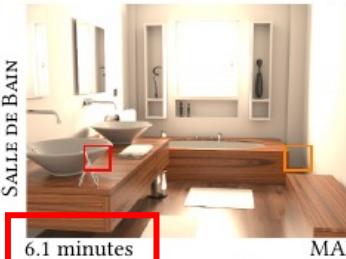
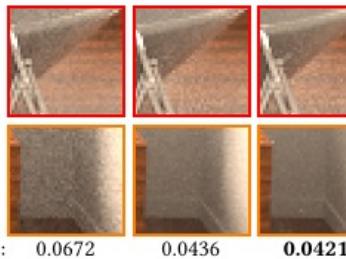
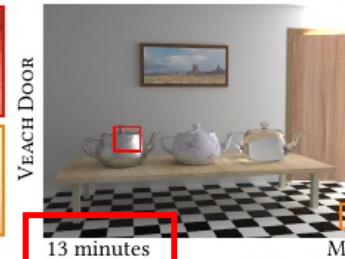
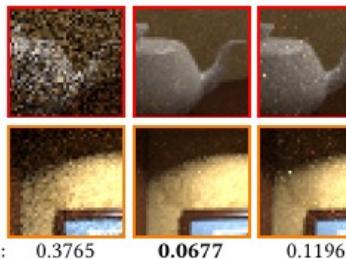
**BUT**

What is the computational **overhead** of Neural Importance Sampling?



Let's run all the approaches the **same amount of time**, which one performs **better**?



	PT-Unidir	PPG	NPG-Product		
BATHROOM					
	12 minutes	MAPE: Mega samples:	0.0512 2164	0.1018 860	0.0542 236
SALLE DE BAIN					
	6.1 minutes	MAPE: Mega samples:	0.0672 1875	0.0436 661	0.0421 236
BOOKSHELF					
	10 minutes	MAPE: Mega samples:	0.3765 2157	0.0677 664	0.1196 236
WHITE ROOM					
	8.3 minutes	MAPE: Mega samples:	0.0396 1604	0.0396 668	0.0400 236
WEECH DOOR					
	13 minutes	MAPE: Mega samples:	0.2317 5413	0.0857 1540	0.0995 236
GLOSSY KITCHEN					
	11 minutes	MAPE: Mega samples:	1.3179 2222	0.1343 1081	0.1363 236

# Further Extensions

1

Extend the piecewise-polynomial coupling layer to higher degree



Piecewise-**cubic** warp!

**Cubic-spline Flows**  
Durkan et al., 2019

2

Improve density estimation by proposing a coupling layer in the rational domain



Warping made of the **quotient of two quadratic polynomials**

**Neural Spline Flows**  
Durkan et al., 2019

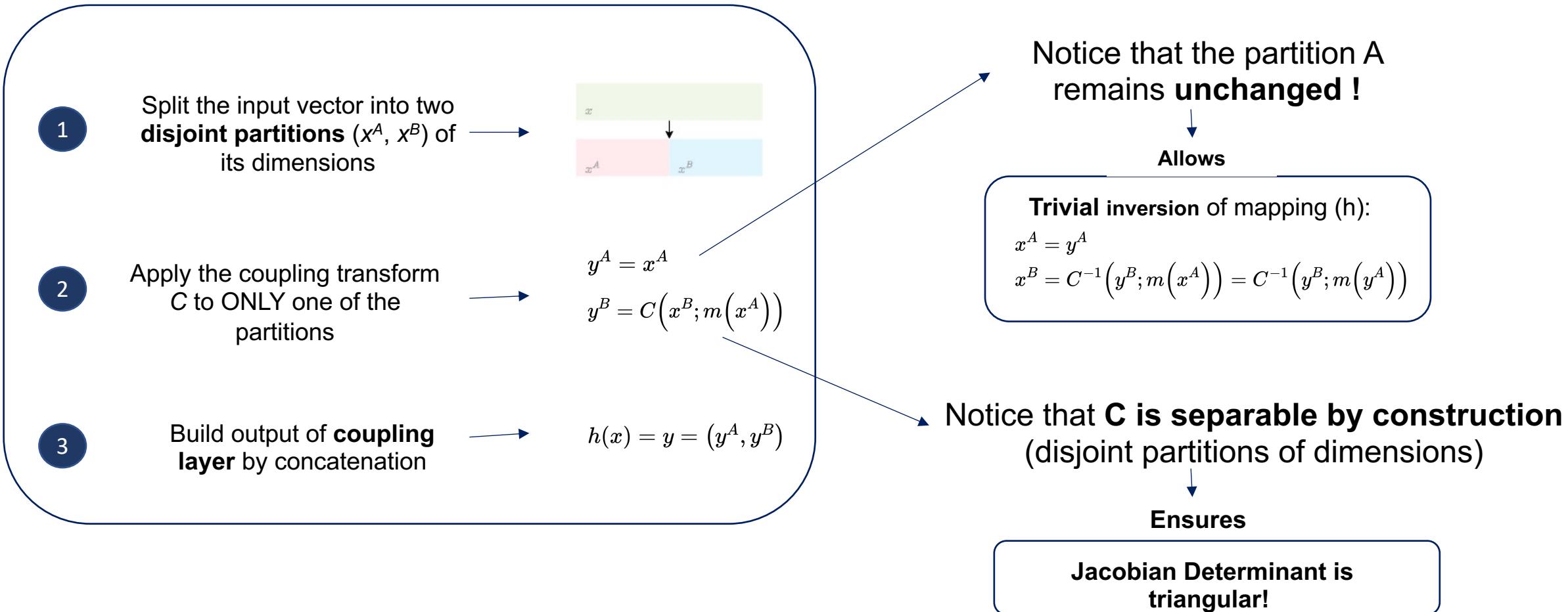
# Takeaways

- **Integration** is a problem in a surprising variety of fields
  - **ML can help** solve these problems!
- Monte Carlo Integration with **NICE reduces the estimated variance** of the integrand.
- **NFs** can be made **more expressive** without overloading the Jacobian determinant term
  - The idea of coupling layers can be extended to **Neural Networks!**
  - Coupling layers can be of **various types**, adding expressiveness
- While Neural Importance Sampling is impractical, it outperforms the state-of-the-art approaches!

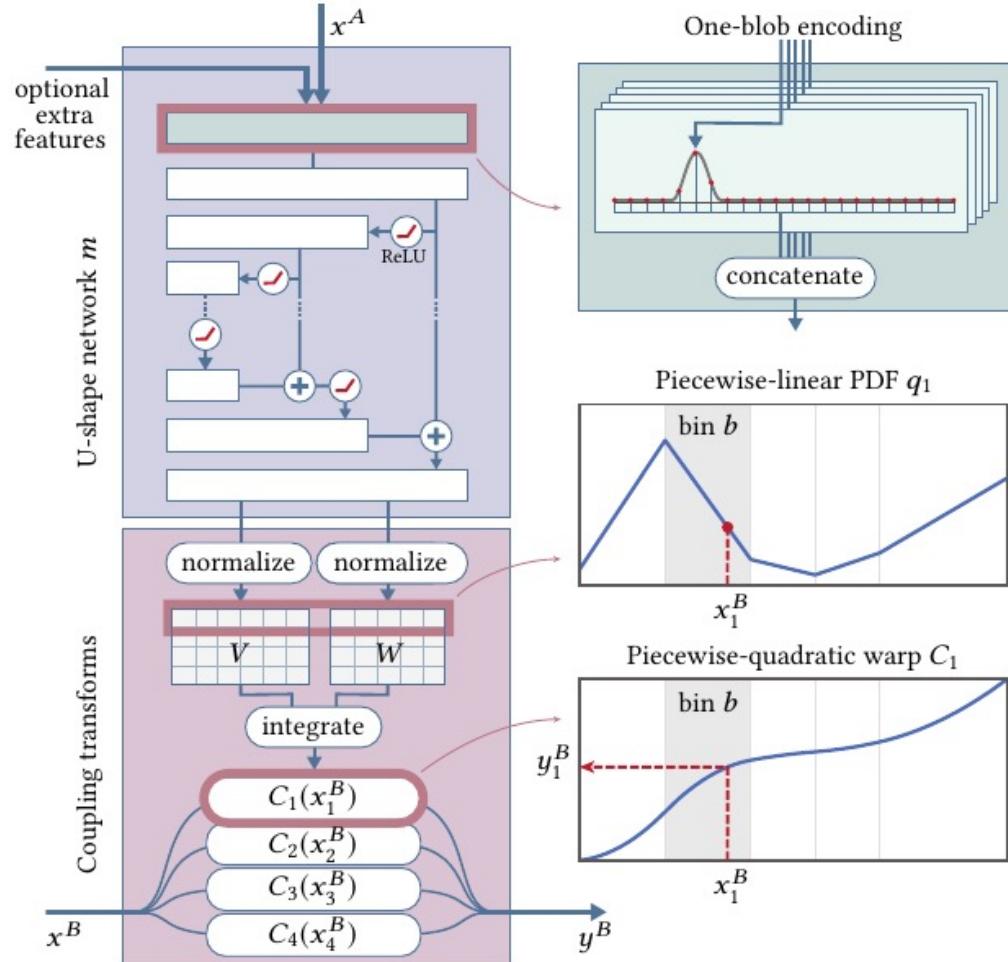
# References

- Dinh, L., Krueger, D., & Bengio, Y. (2015). NICE: Non-linear Independent Components Estimation. *ArXiv:1410.8516 [Cs]*.  
<http://arxiv.org/abs/1410.8516>
- Durkan, C., Bekasov, A., Murray, I., & Papamakarios, G. (2019a). Cubic-Spline Flows. *ArXiv:1906.02145 [Cs, Stat]*.  
<http://arxiv.org/abs/1906.02145>
- Durkan, C., Bekasov, A., Murray, I., & Papamakarios, G. (2019b). Neural Spline Flows. *ArXiv:1906.04032 [Cs, Stat]*.  
<http://arxiv.org/abs/1906.04032>
- Müller, T., McWilliams, B., Rousselle, F., Gross, M., & Novák, J. (2019). Neural Importance Sampling. *ArXiv:1808.03856 [Cs, Stat]*.  
<http://arxiv.org/abs/1808.03856>

# Extra: Why should we transform only one partition when building coupling layers?



# Extra: More details about One-blob Encoding



Generalization  
of one-hot  
encoding

**Gaussian Kernel**  
 $\sigma = 1/k$ , where  $k = \#$  of bins

A **kernel** is used to activate multiple adjacent matrix instead of a single one

Fancy properties...

- It shuts down certain parts of the linear path of the network, allowing it to **specialize the model on various sub-domains** of the input.
- In contrast to one-hot encoding, the **one-blob encoding is lossless**; it captures the exact position of the continue variable.

# Extra: the effect of different loss functions on the model

## (1) Kullback-Leibler-Divergence

- Notice the **importance weight** before the gradient of the log likelihood
- **Minimizing the KL divergence via** gradient descent is equivalent to **minimizing the negative log likelihood weighted by MC estimates of F.**
- This is perhaps the most “natural” loss to use, but **it does not explicitly control the generation of outliers** which will increase the variance of the estimator.

$$\begin{aligned} D_{\text{KL}}(p\|q; \theta) &= \int_{\Omega} p(x) \log \frac{p(x)}{q(x; \theta)} dx \\ &= \int_{\Omega} p(x) \log p(x) dx - \underbrace{\int_{\Omega} p(x) \log q(x; \theta) dx}_{\text{Cross entropy}}. \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} D_{\text{KL}}(p\|q; \theta) &= -\nabla_{\theta} \int_{\Omega} p(x) \log q(x; \theta) dx \\ &= \mathbb{E} \left[ -\frac{p(X)}{q(X; \theta)} \nabla_{\theta} \log q(X; \theta) \right] \end{aligned}$$

# Extra: the effect of different loss functions on the model

## (2) $\chi^2$ -Divergence

- Notice how the gradient is weighted by the **square** of the importance weight
- Minimising the variance** of the estimator is **equivalent**, upon close inspection, to **minimising the Pearson chi-squared divergence** between the proposal and ground truth
- This estimate is **more conservative than DKL**, i.e. it penalises more strongly low densities of  $q$  in regions where  $p$  is high density.

$$\begin{aligned}\mathbb{V}\left[\frac{p(X)}{q(X; \theta)}\right] &= \mathbb{E}\left[\frac{p(X)^2}{q(X; \theta)^2}\right] - \mathbb{E}\left[\frac{p(X)}{q(X; \theta)}\right]^2 \\ &= \int_{\Omega} \frac{p(x)^2}{q(x; \theta)} dx - \underbrace{\left(\int_{\Omega} p(x) dx\right)^2}_1. \\ \nabla_{\theta} \mathbb{V}\left[\frac{p(X)}{q(X; \theta)}\right] &= \nabla_{\theta} \int_{\Omega} \frac{p(x)^2}{q(x; \theta)} dx \\ &= \int_{\Omega} p(x)^2 \nabla_{\theta} \frac{1}{q(x; \theta)} dx \\ &= \int_{\Omega} -\frac{p(x)^2}{q(x; \theta)} \nabla_{\theta} \log q(x; \theta) dx \\ &= \mathbb{E}\left[-\left(\frac{p(X)}{q(X; \theta)}\right)^2 \nabla_{\theta} \log q(X; \theta)\right].\end{aligned}$$