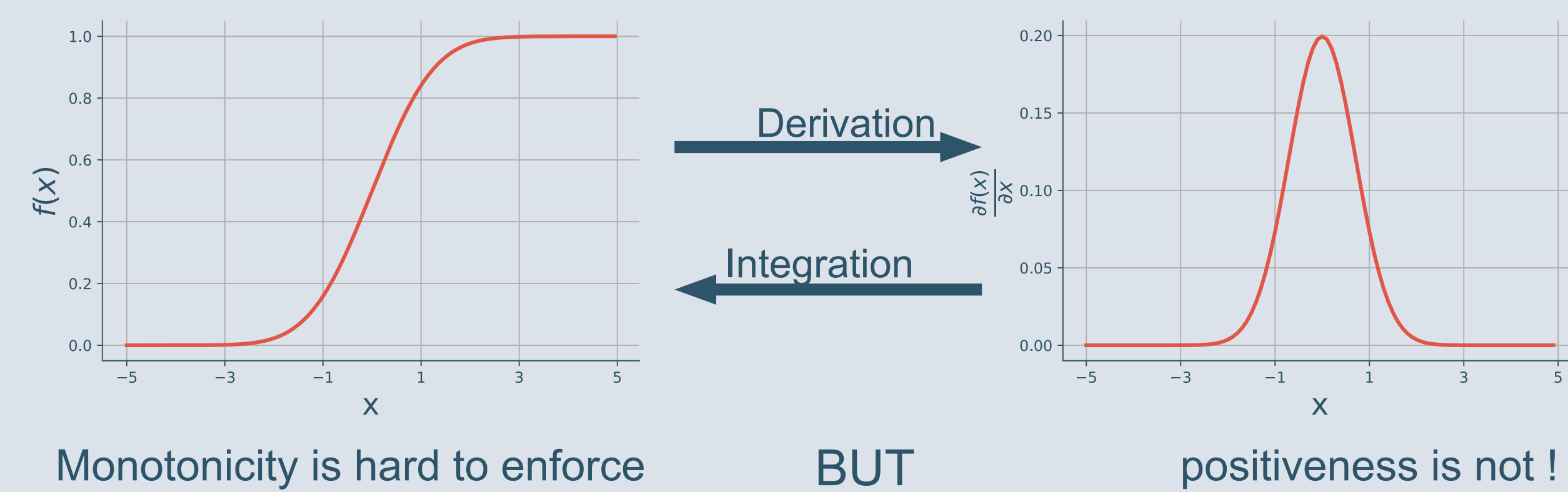


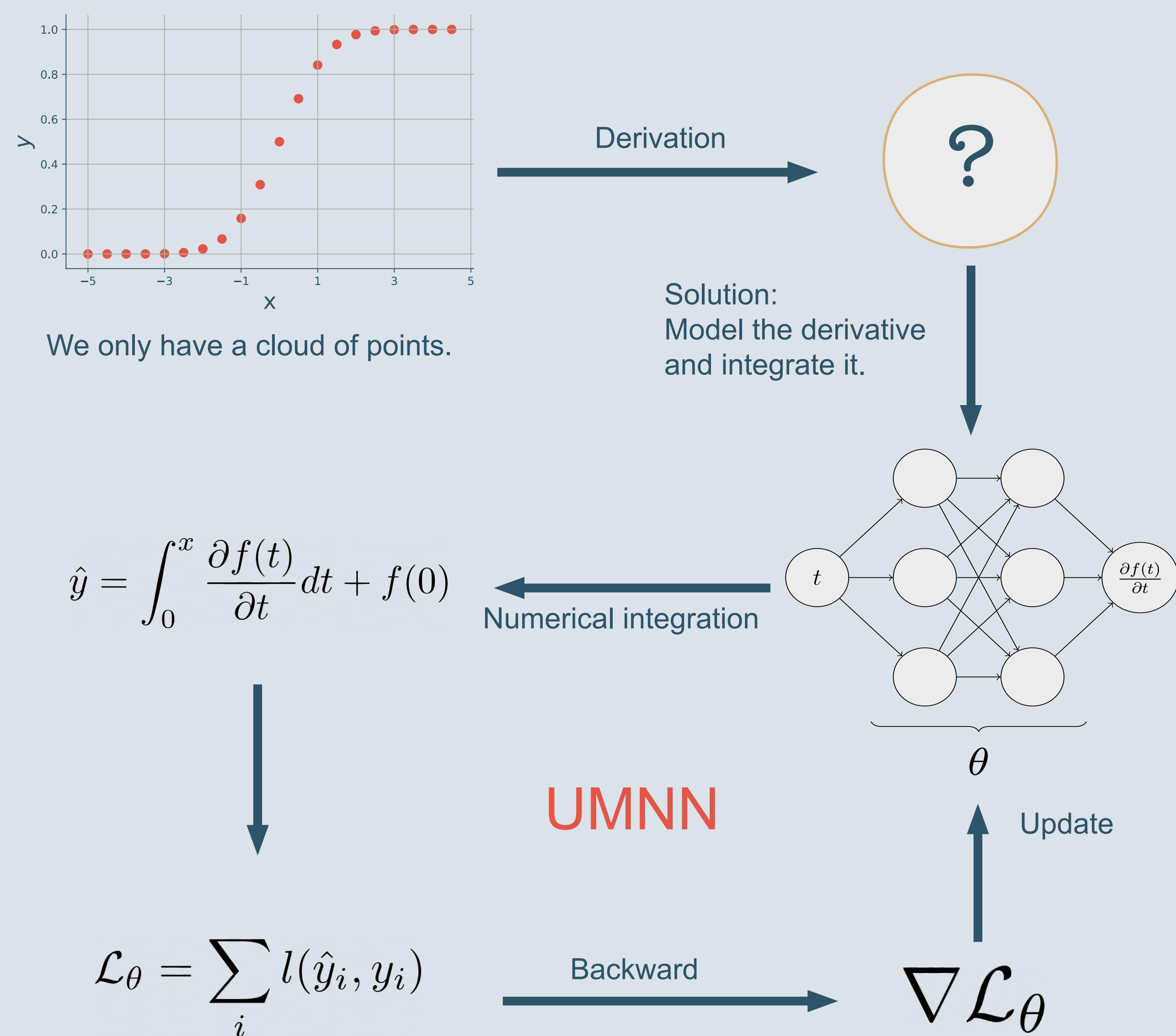
What ? UMNN is a new architecture to model monotonic functions.
How ? The strictly positive scalar output of a neural network is numerically integrated.
Applications ? We combine UMNNs with autoregressive flows to perform density estimation.

Monotonicity

In theory



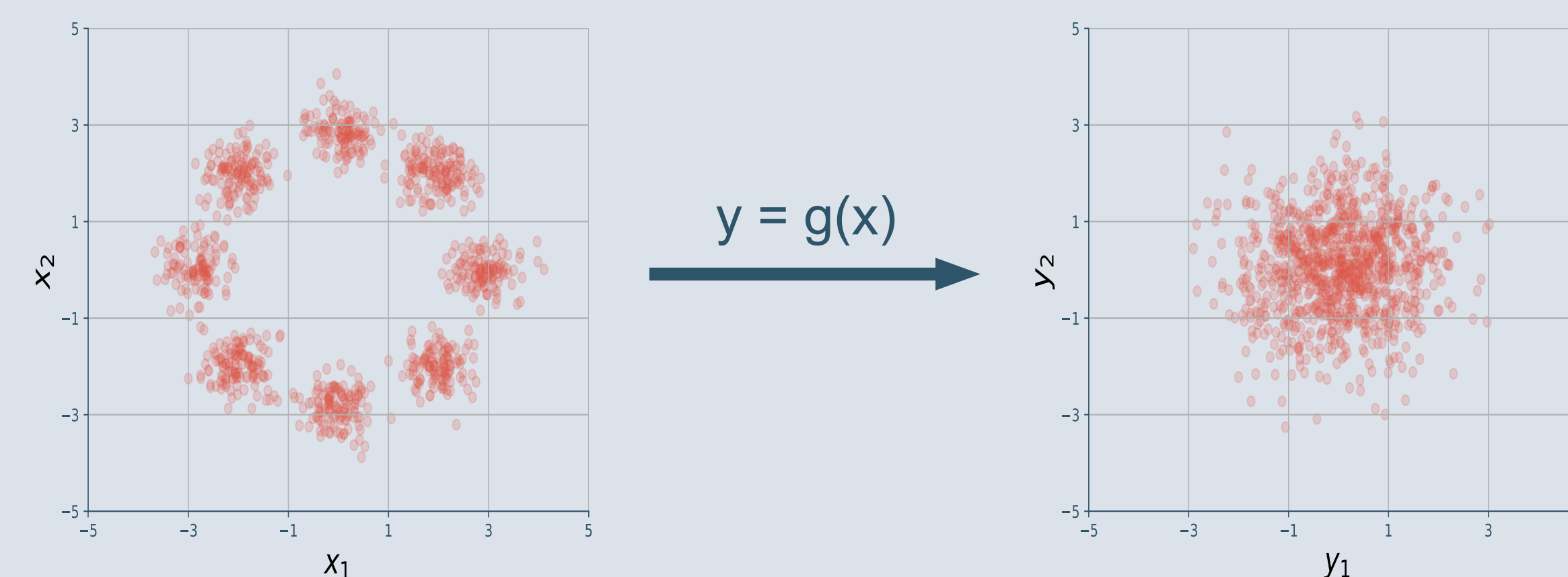
In practice



Change of variables

Let g be a bijective function, x a random variable and let y defined as $g(x)$. The change of variables theorem states that:

$$f_Y(y) = f_X(g^{-1}(y)) \left| \det(J_{g^{-1}}) \right|$$



A bijective transformation can be built by the combination of an autoregressive architecture with a UMNN.

Autoregressivity

Autoregressive transformations are commonly used to build bijective transformations.

$$g(\mathbf{x}; \theta) = [g^1(x_1; \theta) \quad \dots \quad g^i(\mathbf{x}_{1:i}; \theta) \quad \dots \quad g^d(\mathbf{x}_{1:d}; \theta)]$$

The induced multivariate density can be expressed by the chain rule:

$$p(\mathbf{x}; \theta) = p(x_1; \theta) \prod_{i=1}^{d-1} p(x_{i+1} | \mathbf{x}_{1:i}; \theta).$$

We combine UMNN with autoregressive transformations as:

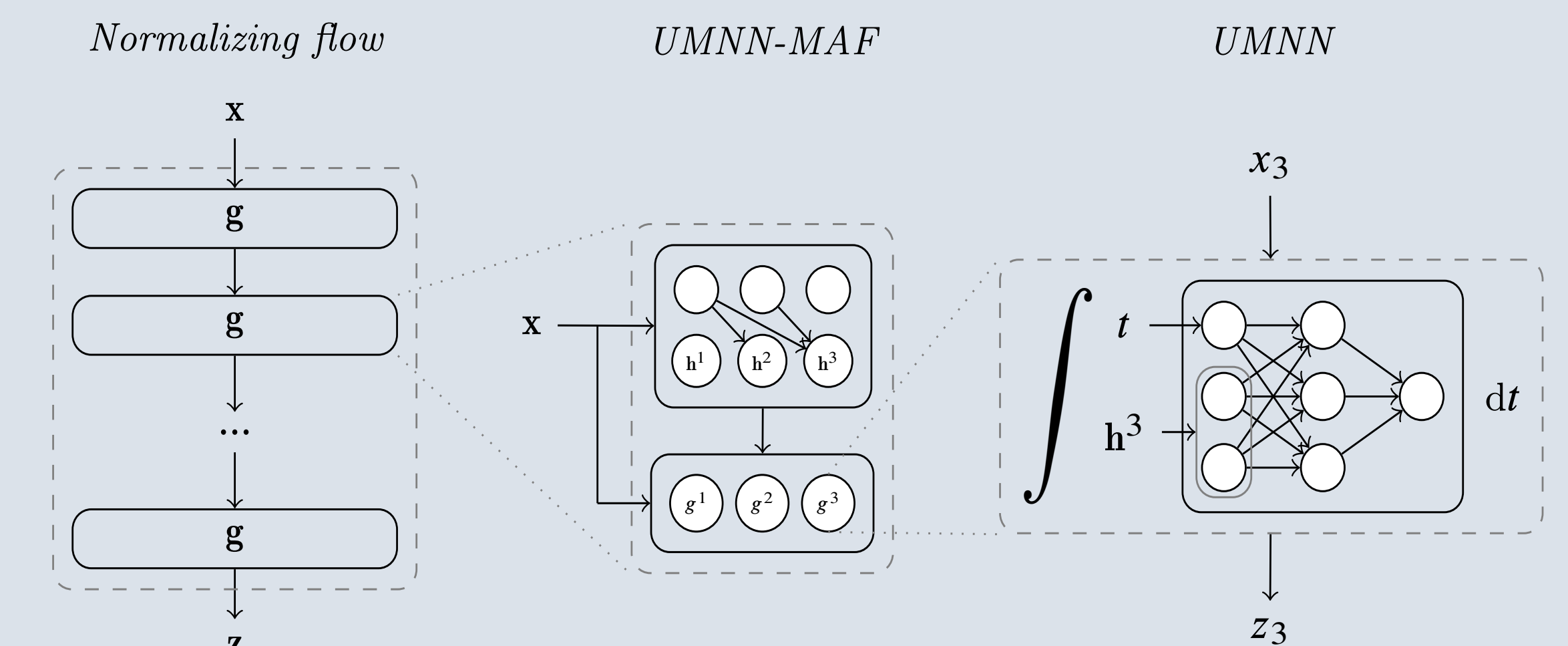
$$g^i(\mathbf{x}_{1:i}; \theta) = F^i(x_i, \mathbf{h}^i(\mathbf{x}_{1:i-1}; \phi^i); \psi^i) = \int_0^{x_i} f^i(t, \mathbf{h}^i(\mathbf{x}_{1:i-1}; \phi^i); \psi^i) + \beta^i(\mathbf{h}^i(\mathbf{x}_{1:i-1}; \phi^i))$$

UMNN-MAF leads to a simple expression of the Jacobian:

$$\log p(\mathbf{x}; \theta) = \log p_Z(g(\mathbf{x}; \theta)) + \sum_{i=1}^d \log f^i(x_i, \mathbf{h}^i(\mathbf{x}_{1:i-1}))$$

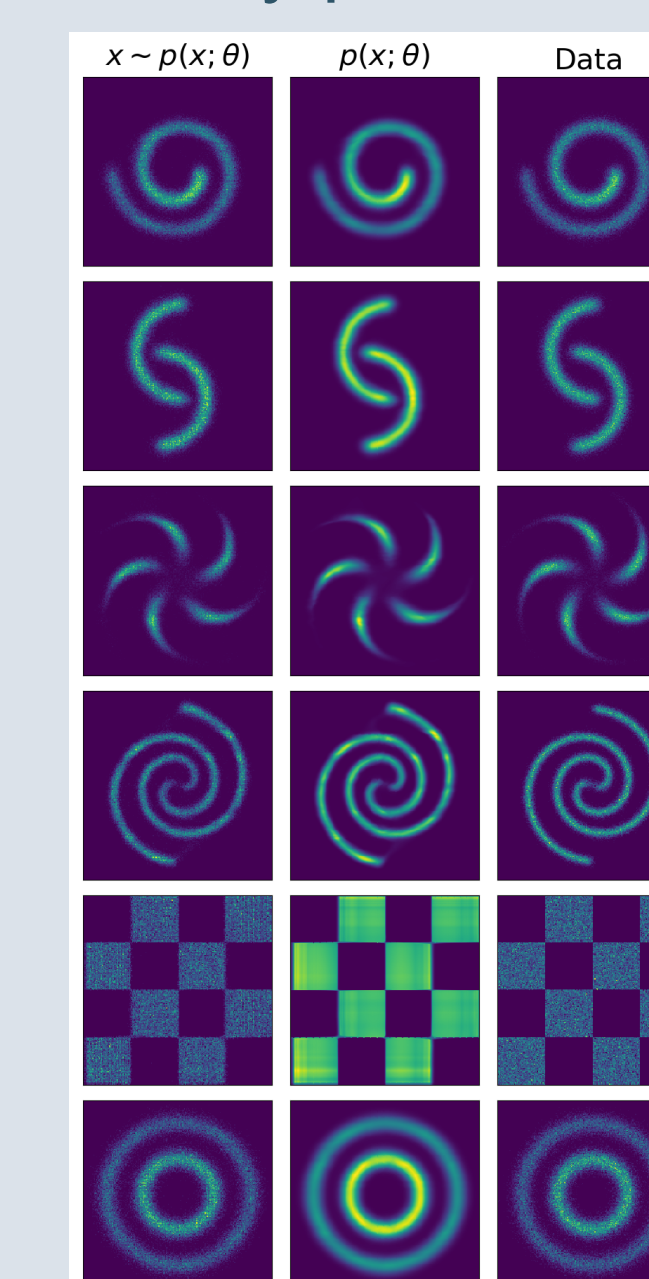
Architecture

We combine the UMNN architecture with an autoregressive network to represent multi-dimensional bijective transformations.

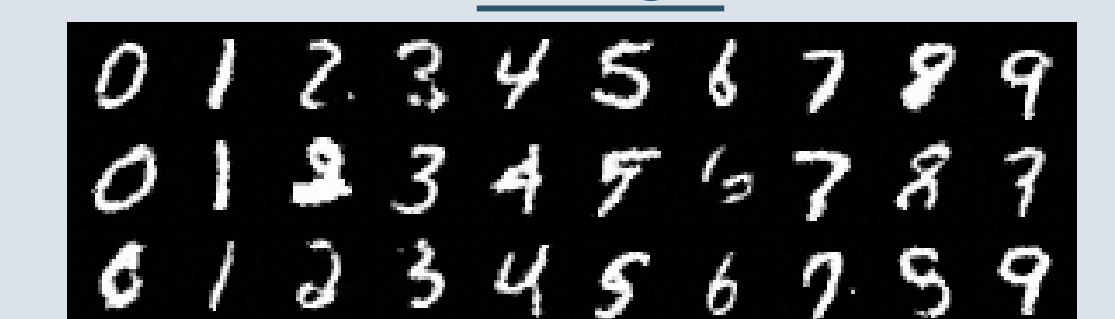


Results

Toy problems



MNIST



Density Estimation

Dataset	POWER	GAS	HEPMASS	MINIBOONE	BSDS300	MNIST
RealNVP - Dinh et al. [2017]	-0.17 ± 0.14	-8.33 ± 0.14	18.71 ± 0.02	13.55 ± 0.49	-153.28 ± 1.78	-
(a) Glow - Kingma and Dhariwal [2018]	-0.17 ± 0.01	-8.15 ± 0.40	19.92 ± 0.08	11.35 ± 0.07	-155.07 ± 0.03	-
FFJORD - Grathwohl et al. [2018]	-0.46 ± 0.01	-8.59 ± 0.12	14.92 ± 0.08	10.43 ± 0.04	-157.40 ± 0.19	-
MADE - Germain et al. [2015]	3.08 ± 0.03	-3.56 ± 0.04	20.98 ± 0.02	15.59 ± 0.50	-148.85 ± 0.28	2.04 ± 0.01
(b) MAF - Papamakarios et al. [2017]	-0.24 ± 0.01	-10.08 ± 0.02	17.70 ± 0.02	11.75 ± 0.44	-155.69 ± 0.28	1.89 ± 0.01
TAN - Oliva et al. [2018]	-0.60 ± 0.01	-12.06 ± 0.02	13.78 ± 0.02	11.01 ± 0.48	-159.80 ± 0.07	1.19
NAF - Huang et al. [2018]	-0.62 ± 0.01	-11.96 ± 0.33	15.09 ± 0.40	8.86 ± 0.15	-157.73 ± 0.30	-
(c) B-NAF - De Cao et al. [2019]	-0.61 ± 0.01	-12.06 ± 0.09	14.71 ± 0.38	8.95 ± 0.07	-157.36 ± 0.03	-
SOS - Jaini et al. [2019]	-0.60 ± 0.01	-11.99 ± 0.41	15.15 ± 0.11	8.90 ± 0.11	-157.48 ± 0.41	1.81
UMNN-MAF (ours)	-0.63 ± 0.01	-10.89 ± 0.7	13.99 ± 0.21	9.67 ± 0.13	-157.98 ± 0.01	1.13 ± 0.02

Fun Facts

1. The numerical integration is performed with static Clenshaw-Curtis method which is proven to converge for Lipschitz continuous functions.
2. The backward computation is performed by solving numerically another integral coming from the Leibnitz integral rule which leads to:

$$\nabla_\psi F(x; \psi) = \int_0^x \nabla_\psi f(t; \psi) + \nabla_\psi \beta.$$