

Conditional BRUNO: A Deep Recurrent Process for Exchangeable Labelled Data

Iryna Korshunova¹ Yarin Gal² Joni Dambre¹ Arthur Gretton³
¹Ghent University ²University of Oxford ³Gatsby Unit, University College London

Overview

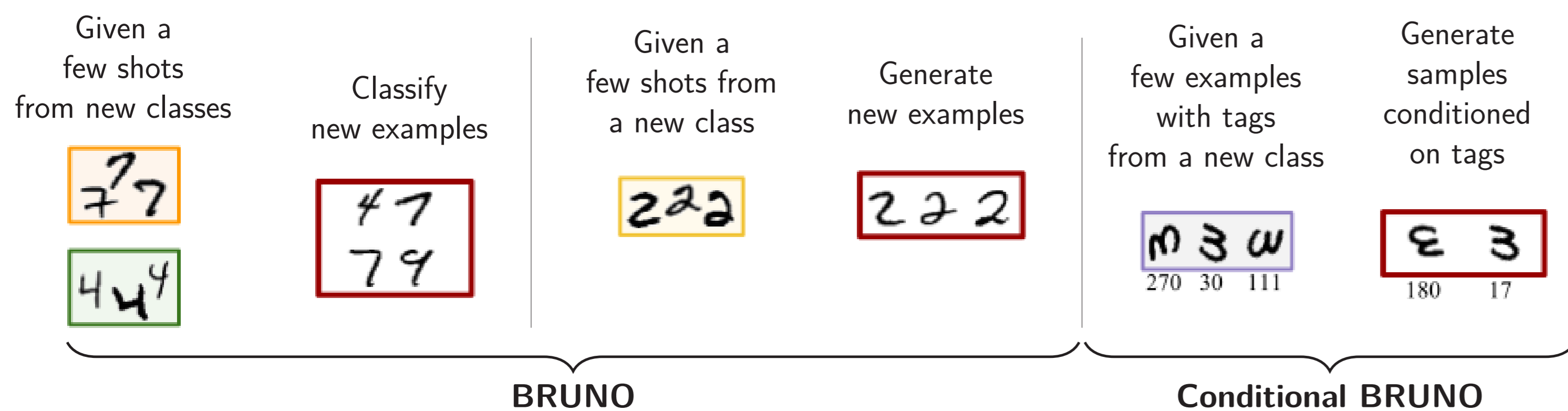
BRUNO combines the expressiveness of deep neural networks with the data-efficiency of \mathcal{GP} s to model exchangeable sequences of complex observations.

BRUNO can be extended to the **conditional** case so that it can model sequences of observations $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$ conditionally on a set of labels or tags $\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3, \dots$.

Conditional BRUNO enjoys a few properties that are desirable in practice:

- ✓ predictive distribution $p(\mathbf{x}_n | \mathbf{h}_n, \mathbf{x}_{1:n-1}, \mathbf{h}_{1:n-1})$ is fast to evaluate and to sample from
- ✓ $p(\mathbf{x}_n | \mathbf{h}_n, \mathbf{x}_{1:n-1}, \mathbf{h}_{1:n-1})$ is differentiable with respect to the model parameters
- ✓ can be trained efficiently in an RNN-like fashion

Exchangeability and meta-learning



Exchangeability and Bayesian computations

A stochastic process x_1, x_2, x_3, \dots is exchangeable if for all n and all permutations π :

$$p(x_1, \dots, x_n) = p(x_{\pi(1)}, \dots, x_{\pi(n)})$$

De Finetti's theorem says that every exchangeable process is a mixture of i.i.d. processes:

$$p(x_1, \dots, x_n) = \int p(\theta) \prod_{i=1}^n p(x_i | \theta) d\theta,$$

where θ is some parameter conditioned on which the data is i.i.d.

example

$x_1, \dots, x_n \sim \mathcal{N}(\mathbf{0}, \Sigma)$ with **compound symmetric covariance**

$$\Sigma = \begin{bmatrix} \nu & \rho & \rho & \dots & \rho \\ \rho & \nu & \rho & \dots & \rho \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \rho & \dots & \nu \end{bmatrix} \quad 0 \leq \rho < \nu$$

x_1, \dots, x_n are i.i.d. with $x_i \sim \mathcal{N}(\theta, \nu - \rho)$ conditioned on $\theta \sim \mathcal{N}(\mathbf{0}, \rho)$

De Finetti's theorem in terms of **predictive distributions**:

$$p(x_n | x_{1:n-1}) = \int \underbrace{p(x_n | \theta)}_{\text{likelihood}} \underbrace{p(\theta | x_{1:n-1})}_{\text{posterior}} d\theta$$

This gives two ways for defining models of exchangeable sequences:

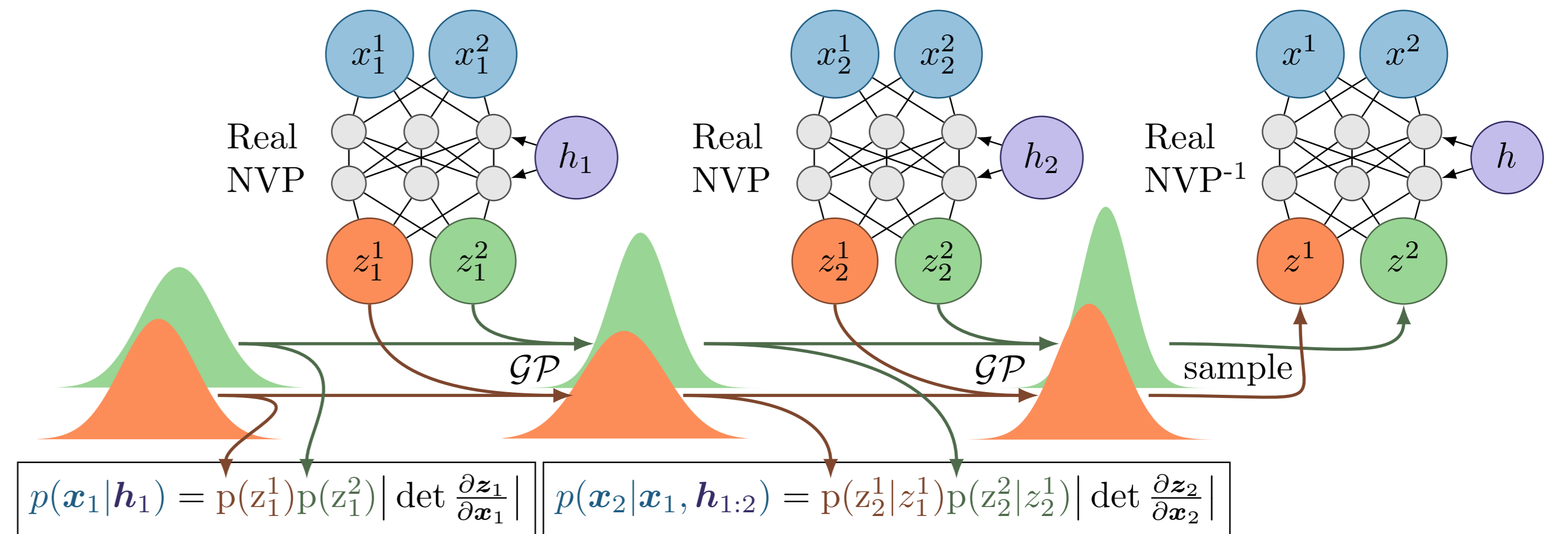
- 1) via explicit Bayesian modelling
- 2) via exchangeable processes \rightarrow BRUNO

For conditional real-valued processes, where x_1, x_2, x_3, \dots is associated with h_1, h_2, h_3, \dots , de Finetti's theorem is not proven.

conjecture The decomposition of the form $p(x_{1:n} | h_{1:n}) = \int p(\theta) \prod_{i=1}^n p(x_i | h_i, \theta) d\theta$ exists if the following conditions hold:

1. $p(x_1, \dots, x_n | h_1, \dots, h_n) = p(x_{\pi(1)}, \dots, x_{\pi(n)} | h_{\pi(1)}, \dots, h_{\pi(n)})$
2. $p(x_{1:m} | h_{1:m}) = \int p(x_{1:n} | h_{1:n}) dx_{m+1:n}$ for $1 \leq m < n$.

Conditional BRUNO



- assumptions**
- A1:** dimensions $\{z^d\}_{d=1, \dots, D}$ are independent, so $p(\mathbf{z}) = \prod_{d=1}^D p(z^d)$
- A2:** for each dimension d , we assume that $(z_1^d, \dots, z_n^d) \sim \text{MVN}_n(\mu^d \mathbf{1}, \Sigma^d)$
- mean $\mu^d \mathbf{1}$ is a $1 \times n$ vector filled with $\mu^d \in \mathbb{R}$
 - covariance $n \times n$ matrix Σ^d with $\Sigma_{ii}^d = \nu^d$ and $\Sigma_{ij, i \neq j}^d = \rho^d$ where $0 \leq \rho^d < \nu^d$

For a sequence $(\mathbf{x}_1, \mathbf{h}_1), (\mathbf{x}_2, \mathbf{h}_2), \dots, (\mathbf{x}_N, \mathbf{h}_N)$ the model is trained to maximise

$$\mathcal{L} = \sum_{n=m+1}^N \log p(\mathbf{x}_n | \mathbf{h}_n, \mathbf{x}_{1:m}, \mathbf{h}_{1:m})$$

with respect to Real NVP parameters and Σ parameters for every latent dimension.

Real NVP*

$$f: \mathcal{X} \mapsto \mathcal{Z} \text{ with } \mathcal{X} = \mathbb{R}^D \text{ and } \mathcal{Z} = \mathbb{R}^D$$

- f is bijective
- forward $\mathbf{z} = f(\mathbf{x})$ and inverse $\mathbf{x} = f^{-1}(\mathbf{z})$ mappings are equally expensive
- computing the Jacobian takes $\mathcal{O}(D)$

Coupling layer - the main building block of Real NVP:

$$\begin{cases} \mathbf{y}^{1:d} = \mathbf{x}^{1:d} \\ \mathbf{y}^{d+1:D} = \mathbf{x}^{d+1:D} \odot \exp(\mathbf{s}(\mathbf{x}^{1:d})) + \mathbf{t}(\mathbf{x}^{1:d}) \end{cases}$$

scales and translates only half of the input dimensions at a time; \mathbf{s} and \mathbf{t} are deep neural nets

For a **conditional Real NVP** mapping $\mathbf{z} = f_{\mathbf{h}}(\mathbf{x})$, we can make \mathbf{s} and \mathbf{t} depend on \mathbf{h} by adding a bias computed from the features of \mathbf{h} to every layer inside \mathbf{s} and \mathbf{t} .

Given a distribution $p(\mathbf{z})$, we can evaluate $p(\mathbf{x} | \mathbf{h})$ using the *change of variables* formula:

$$p(\mathbf{x} | \mathbf{h}) = p(\mathbf{z}) \left| \det \left(\frac{\partial f_{\mathbf{h}}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$

*L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using Real NVP. In ICLR'17

Exchangeable Gaussian processes

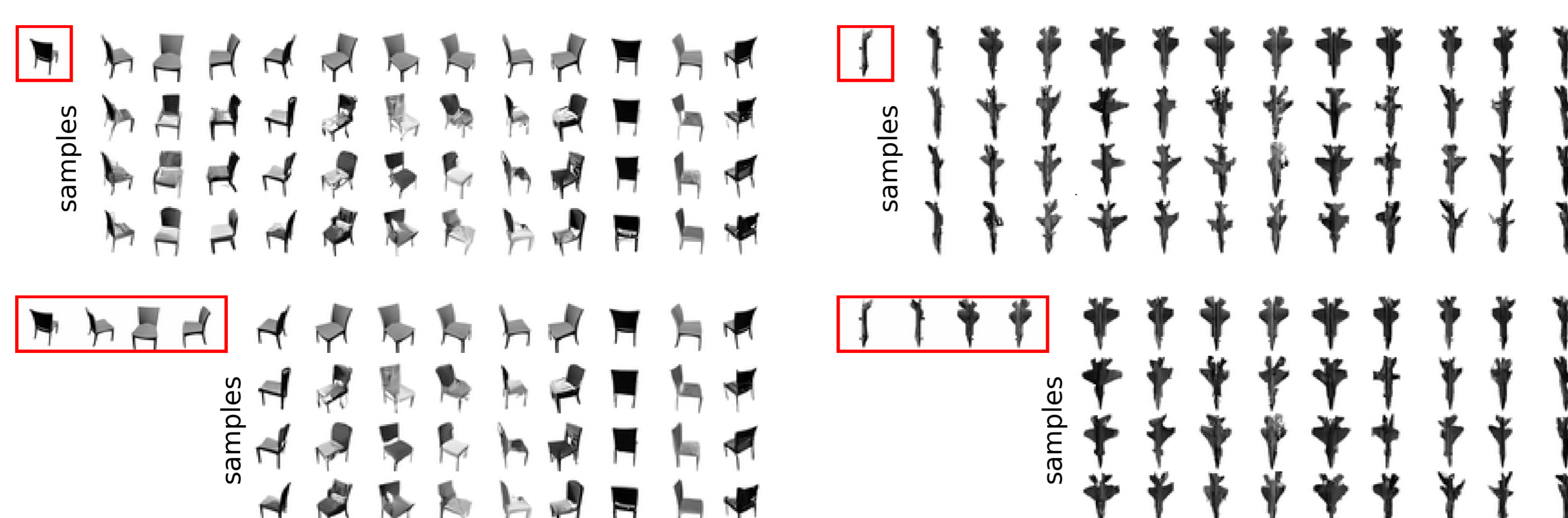
In a \mathcal{GP} , where any finite collection $(z_1, \dots, z_n) \sim \text{MVN}_n(\mu \mathbf{1}, \Sigma)$ with a compound symmetric Σ , recurrent updates for the params of $p(z_{n+1} | z_{1:n}) = \mathcal{N}(\mu_{n+1}, \nu_{n+1})$ are:

$$\begin{aligned} \mu_{n+1} &= (1 - d_n) \mu_n + d_n z_n & \text{with} & & d_n &= \rho / (\nu + \rho(n-1)) \\ \nu_{n+1} &= (1 - d_n) \nu_n + d_n (\nu - \rho) & & & \mu_1 &= \mu, \nu_1 = \nu \end{aligned}$$

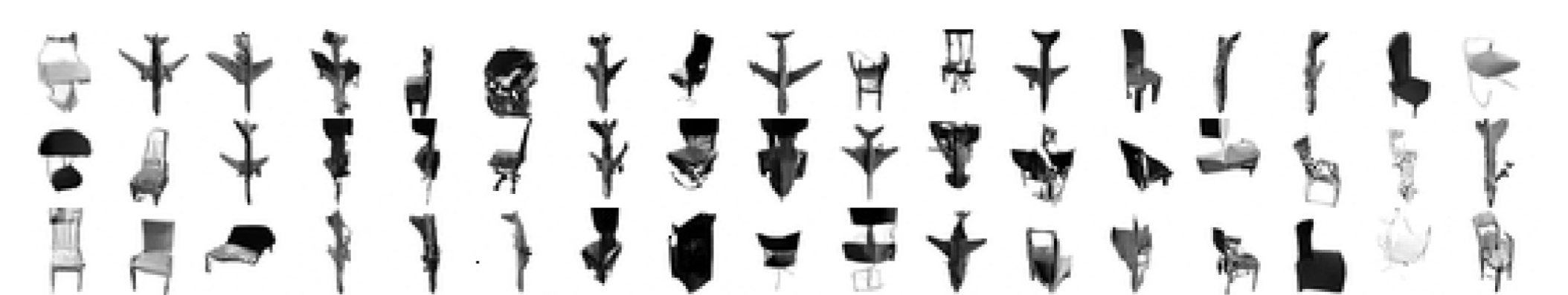
$\mathcal{O}(n)$ runtime and $\mathcal{O}(1)$ memory complexity!

Experiments

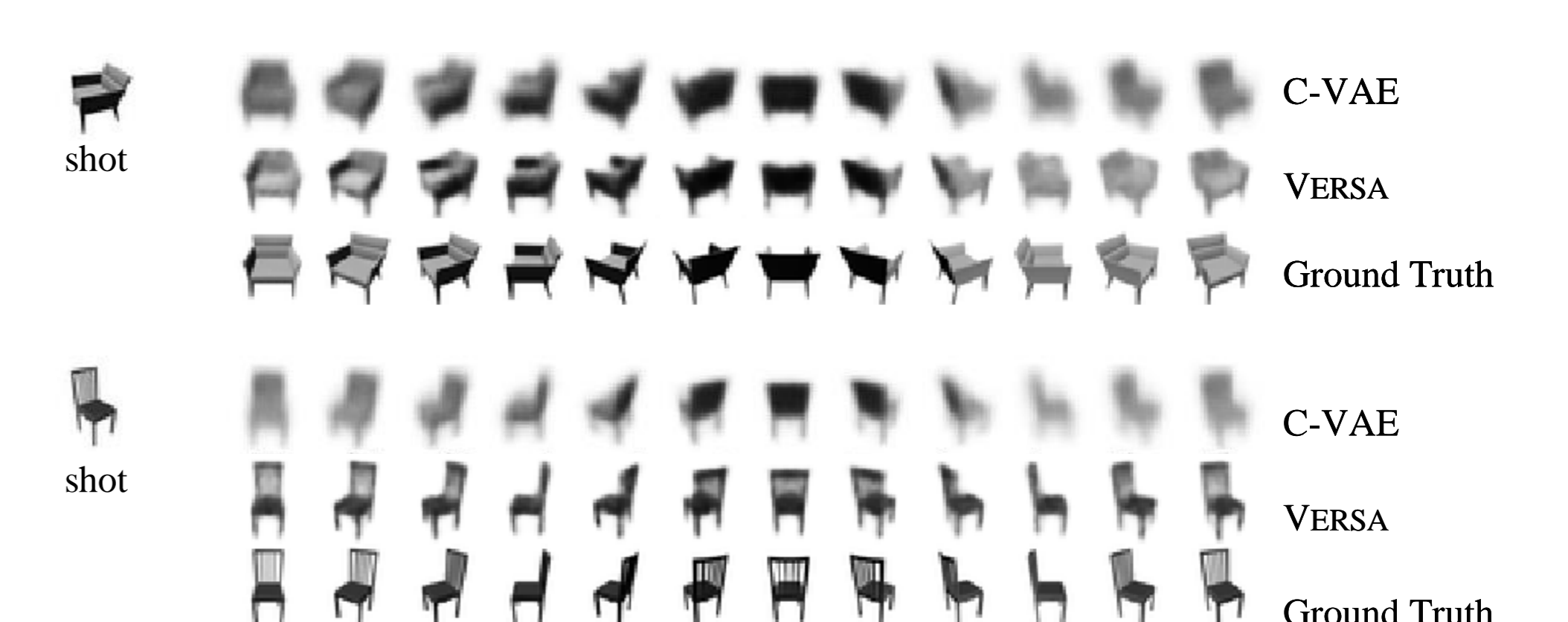
Conditional BRUNO trained on ShapeNet chairs and airplanes



Conditional BRUNO prior samples



VERSA* trained on ShapeNet chairs



*J. Gordon, J. Bronskill, M. Bauer, S. Nowozin, R. E. Turner. Decision-Theoretic Meta-Learning: Versatile and Efficient Amortization of Few-Shot Learning. arXiv:1805.09921