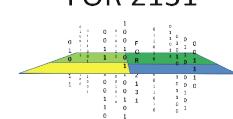




Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG



FOR 2131



International Max Planck Research School
on Earth System Modelling



09.01.2019

Deep assimilation: Adversarial variational Bayes for non-linear data assimilation

Tobias Sebastian Finn

What is the connection between my PhD thesis and machine learning?



?

Writing his PhD thesis

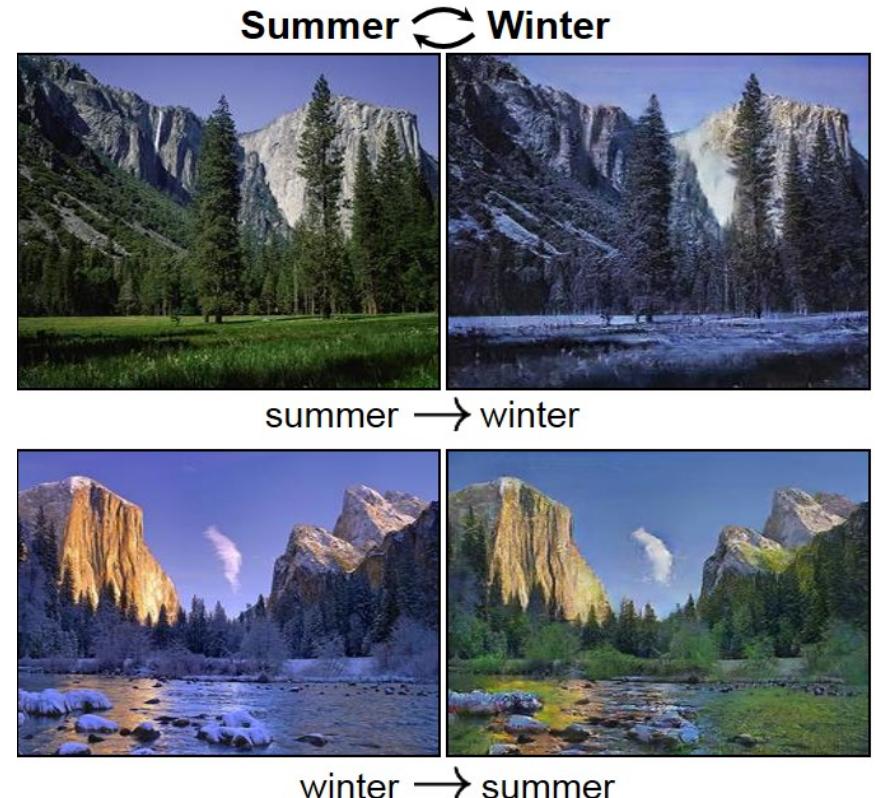


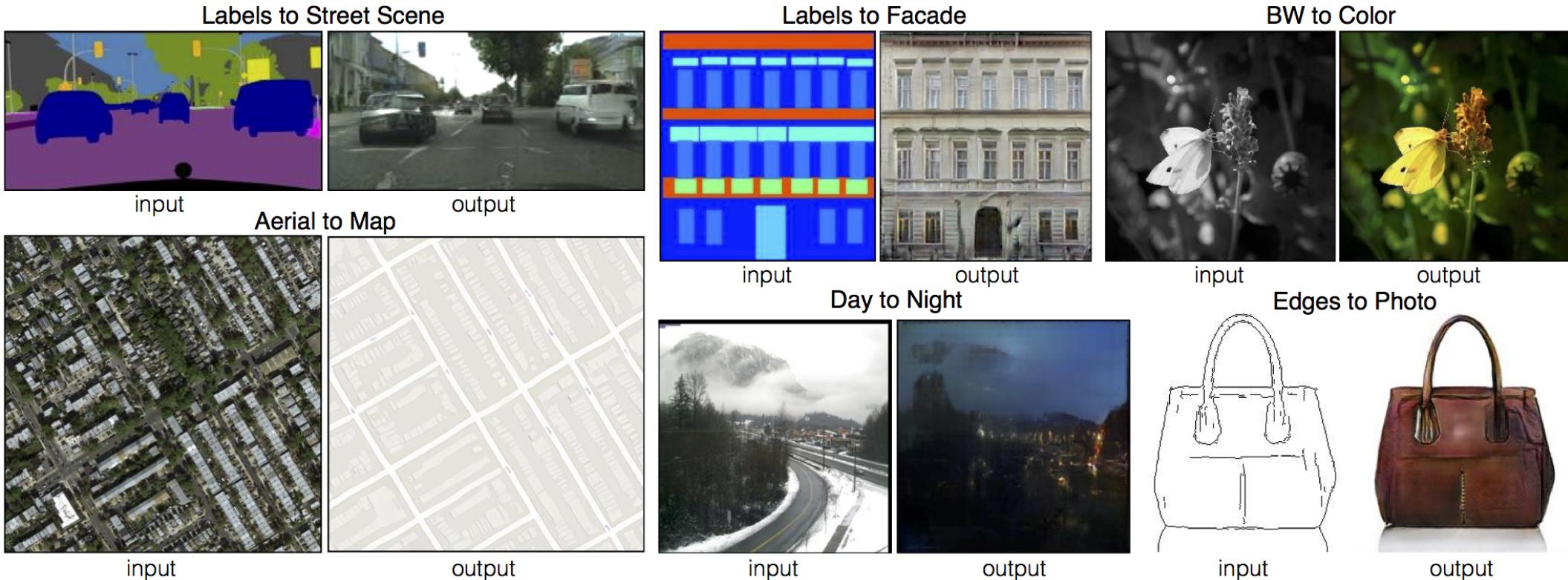
Image-to-Image translation

Sources:

<https://junyanz.github.io/CycleGAN/>

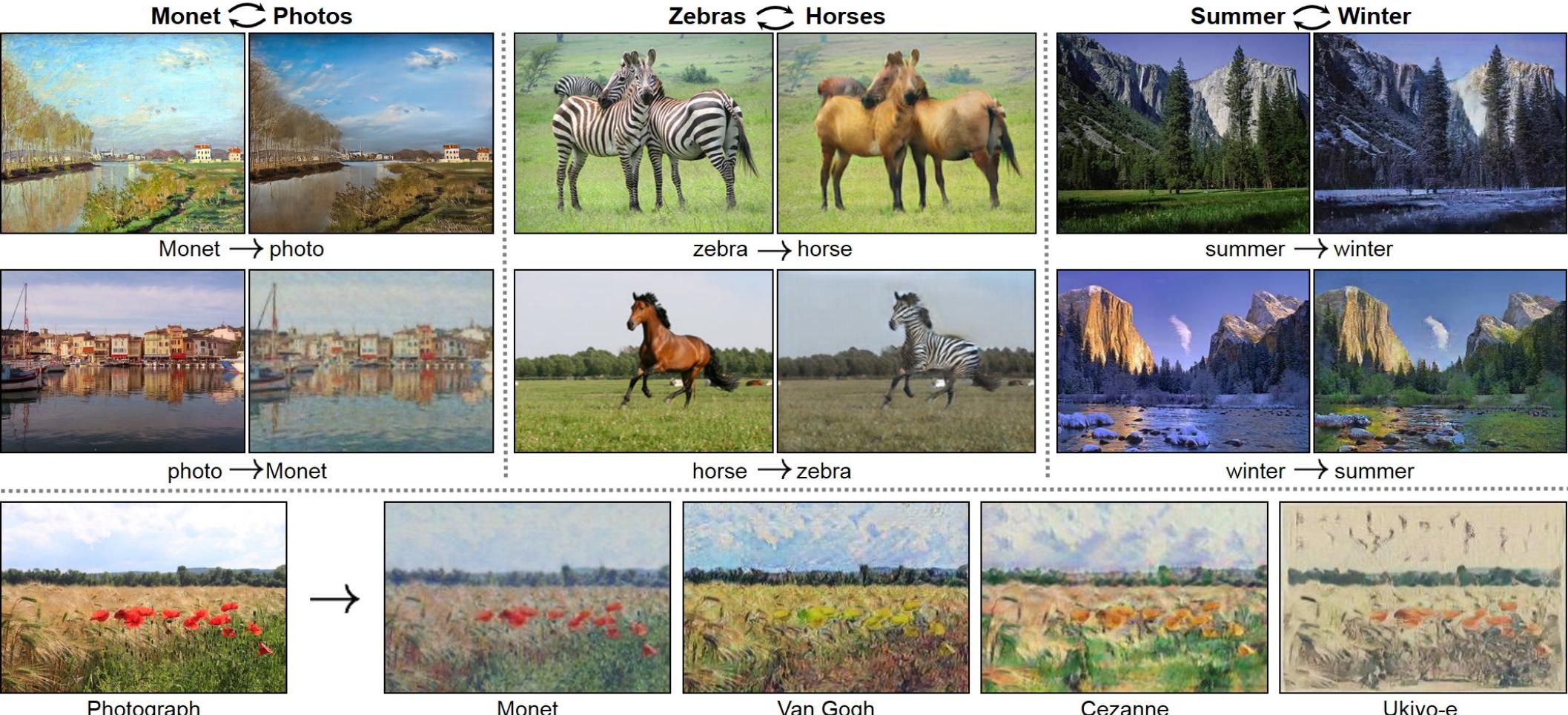
<https://magazin.spiegel.de/SP/2018/52/161498484/index.html>

How can we use machine learning for data assimilation?



Source: <https://phillipi.github.io/pix2pix/> (pix2pix)

How can we use machine learning for data assimilation?



Source: <https://junyanz.github.io/CycleGAN/> (CycleGAN)

How can we use machine learning for data assimilation?



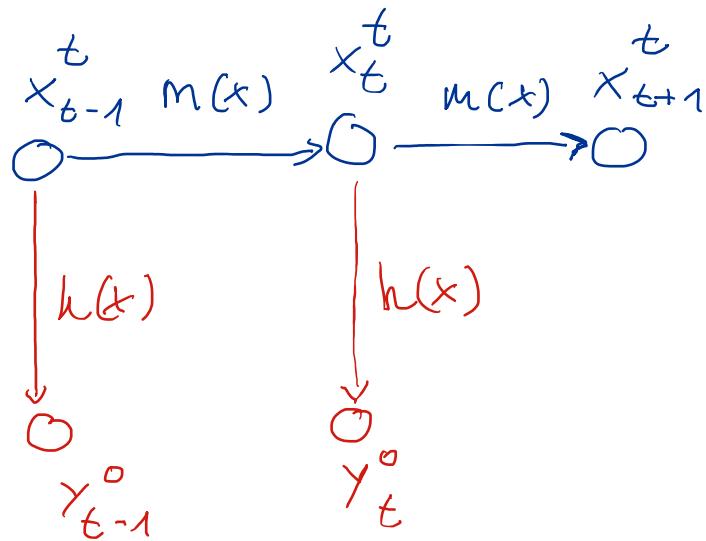
Figure 3: Street scene image translation results. For each pair, left is input and right is the translated image.

Source: <https://arxiv.org/abs/1703.00848> (UNIT)

What's new?

- First application + theoretical foundation of deep learning for non-linear data assimilation
- No assumption about state space
→ flexible distributions + non-Gaussian assimilation
- Efficient and fast inference
- Unclear: assimilation of not yet used observables
(e.g. cloud camera, radar)
- Unclear: unification of data assimilation algorithms

Data assimilation problem



How to get x_t^t based on observations $y_{0:t}^o$?

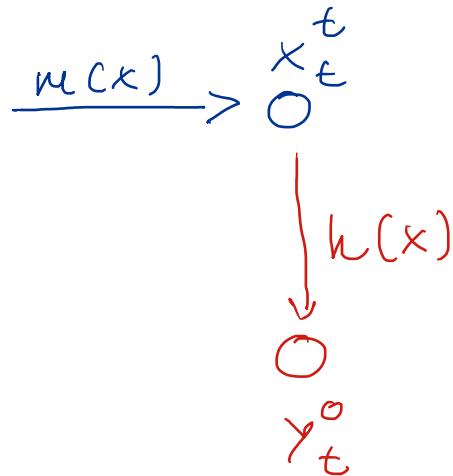
$$y_t^o = h(x_t^t) + \varepsilon^o \quad \varepsilon^o \sim p(\varepsilon^o)$$

Assumption: Markov-1 process, we know $m(x)$ and $h(x)$ up to errors ε .

↪ look at time $t \rightarrow$ update problem \rightarrow simplifies to static problem

Latent variable model for data assimilation

Let's look at time t

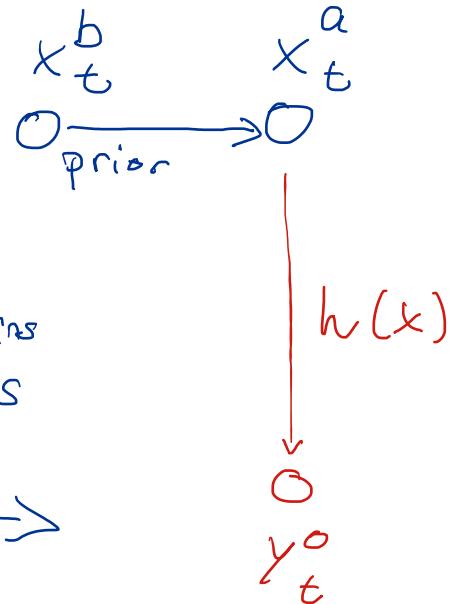


We can only estimate an approximation to x_t^t : x_t^a .

Based on an old analysis and model $m(x)$, we have a prior

x_t^b for x_t^a . This prior contains all information from previous observations.

Static



We can estimate x_t^a with Bayes' formula:

$$p(x_t^a | y_{0:t}^o) = \frac{p(y_t^o | x_t^a) p(x_t^a | y_{0:t}^o)}{p(y_{0:t}^o)} = \frac{p(y_t^o | x_t^a) \cdot p(x_t^b)}{p(y_{0:t}^o)}$$

intractable!

Let's approximate with $q_\phi(x_t^a)$

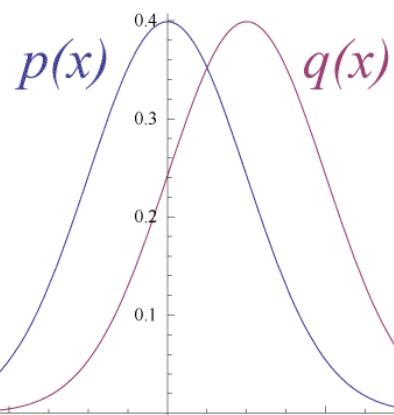
↳ we want $q_\phi(x_t^a) \approx p(x_t^a | y_{0:t}^o) \rightarrow$ Kullback-Leibler Divergence

Kullback-Leibler divergence

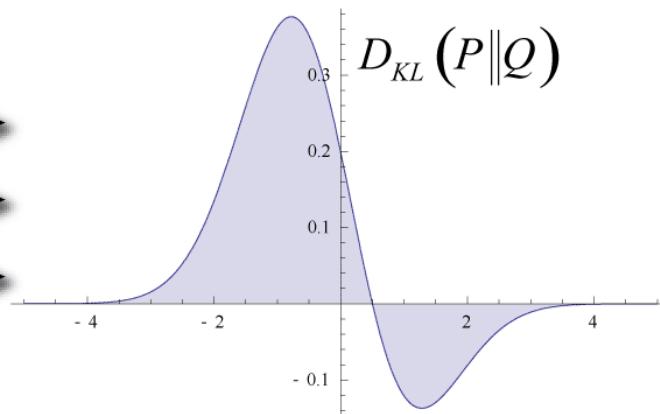
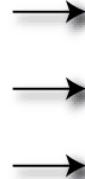
$$\text{KL}(q(x) \parallel p(x)) = - \sum_{x \in \mathcal{X}} q(x) \log\left(\frac{q(x)}{p(x)}\right) = \mathbb{E}_{q(x)} \log\left(\frac{q(x)}{p(x)}\right)$$



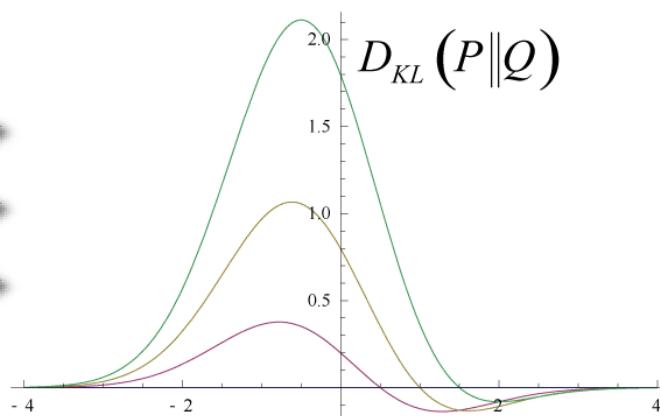
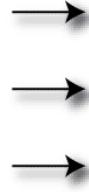
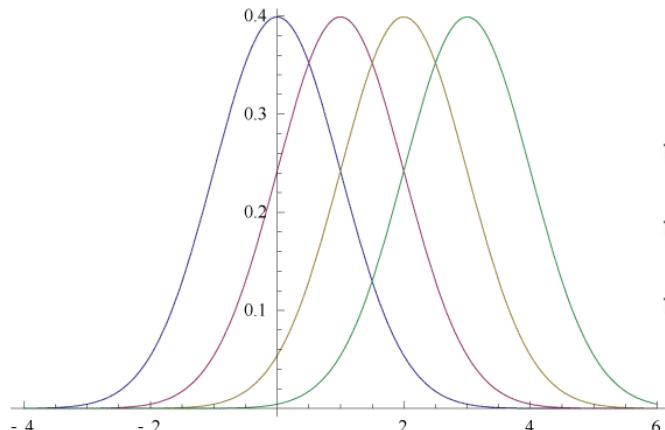
Integral under
KL-area



Original Gaussian PDF's



KL Area to be Integrated



T. Nathan Mundhenk, PhD thesis appendix C.

Variational inference

- Minimize $KL(q_{\phi_t}(\mathbf{x}_t^a) || p(\mathbf{x}_t^a | \mathbf{y}_{0:t}^o))$
- Search for approximated $q_{\phi_t}(\mathbf{x}_t^a)$ with parameters ϕ_t from variational family, e.g. Gaussian assumption
- $p(\mathbf{x}_t^a | \mathbf{y}_{0:t}^o)$ intractable, we cannot optimize!
- We need to circumvent this problem
→ rephrase the minimization into a simpler one!

↳ like done in first two slides!

Evidence lower bound (ELBO)

use of Bayes + rearrange

$$\text{KL}(q_{\phi_t}(\mathbf{x}_t^a) \parallel p(\mathbf{x}_t^a | \mathbf{y}_{o:t}^o)) \geq \text{KL}(q_{\phi_t}(\mathbf{x}_t^a) \parallel p(\mathbf{x}_t^b)) - \mathbb{E}_{\mathbf{x} \sim q_{\phi_t}(\mathbf{x}_t^a)} \log p(\mathbf{y}_t^o | \mathbf{x})$$

Similarity to first guess + Reconstruction loss
(e.g. MSE)

- Can be **optimized** by minimization
- Similar loss to standard data assimilation losses
- Computational **expensive** (can be slower than 4dVar!)
- Every new analysis has to be optimized

↳ online learning like classical DA *only local parameters ϕ_t*

e.g. Vrettas, M. D., Opper, M., & Cornford, D. (2015). Variational mean-field algorithm for efficient inference in large systems of stochastic differential equations.

Amortized inference

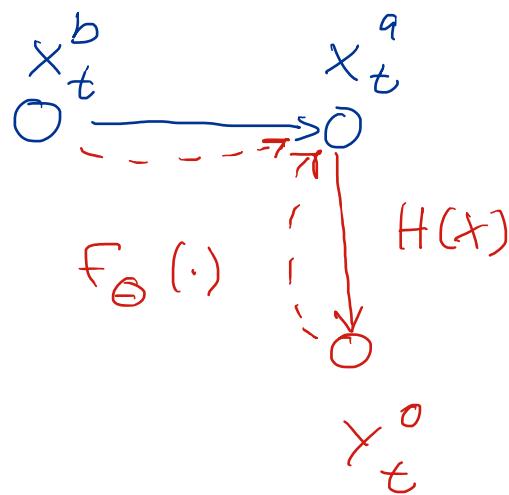
We want to get global fixed parameters for inference

$q_{\phi_t}(x_t^a)$ by making time-dependent conditional probabilities

$$q_{\theta}(x_t^a | y_{0:t}^o) \rightarrow q_{\theta}(x_t^a | y_{t-1}^o, x_t^b)$$

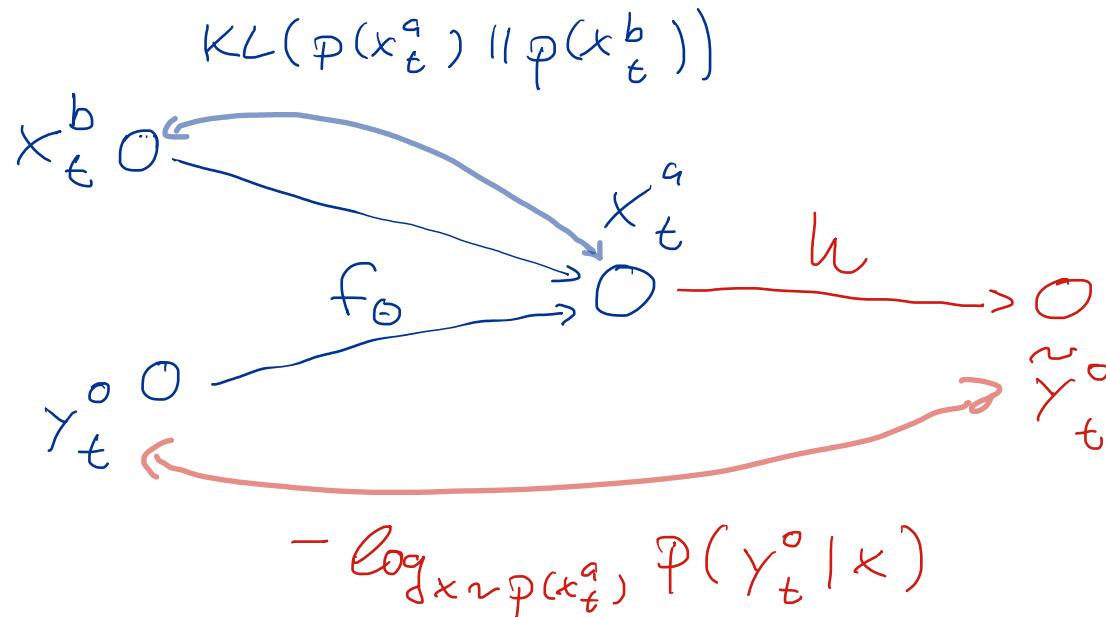
$$\Leftrightarrow x_t^a = f_{\theta}(y_{t-1}^o, x_t^b, \epsilon_t), \quad \epsilon_t \sim p(\epsilon)$$

\Rightarrow



Roweis, S., and Z. Ghahramani, 1999: A Unifying Review of Linear Gaussian Models
Gershman, S., and N. Goodman, 2014: Amortized Inference in Probabilistic Reasoning

Variational autoencoder (VAE)



\Rightarrow Normal assumption: $p(x_t^a | y_t^o, x_t^b) = N\left[f_\theta^\mu(y_t^o, x_t^b), f_\theta^\sigma(y_t^o, x_t^b)\right]$

\hookrightarrow for backpropagation: reparametrization trick

Kingma, D. P., and M. Welling, 2013: Auto-Encoding Variational Bayes

Rezende, D. J., S. Mohamed, and D. Wierstra, 2014: Stochastic Backpropagation and Approximate Inference in Deep Generative Models

Variational autoencoder (VAE)

Θ is trained + applied at inference

- Computational expensive parts are moved to training
→ Efficient at inference!

- (L)ETKF: Linearized VAE + Gaussian assumption
→ analytical solution

} move to weight space instead of model space

- Candidate to unify variational (not the same) + sequential data assimilation
- Needs closed solution of Kullback-Leibler
→ Explicit modelling of prior and posterior

Variational inference ≠ variational DA

For connection ETKF (PPCA) to VAE:
Dai, B., Y. Wang, J. Aston, G. Hua, and D. Wipf, 2017:
Hidden Talents of the Variational Autoencoder.

Variational autoencoder (VAE)

- Computational expensive parts are moved to training
→ Efficient at inference!
- (L)ETKF: Linearized VAE + Gaussian assumption
→ analytical solution
- Candidate to unify
variational (not the same) + sequential data assimilation
- Needs closed solution of Kullback-Leibler
→ Explicit modelling of prior and posterior



We want implicit modelling!

Implicit representation of density ratio

- Kullback-Leibler divergence is density ratio between $q_{\theta}(x_t^a | x_t^b, y_t^o)$ and $p(x_t^b)$
- We can sample from $p(x_t^b)$
- **Approximation via classification:**

Samples from $p(x_t^b)$: True
Samples from $q_{\theta}(x_t^a x_t^b, y_t^o)$: False
- Use another NN for classification

new classifier for
every new
sample !

$$D_t(x)$$

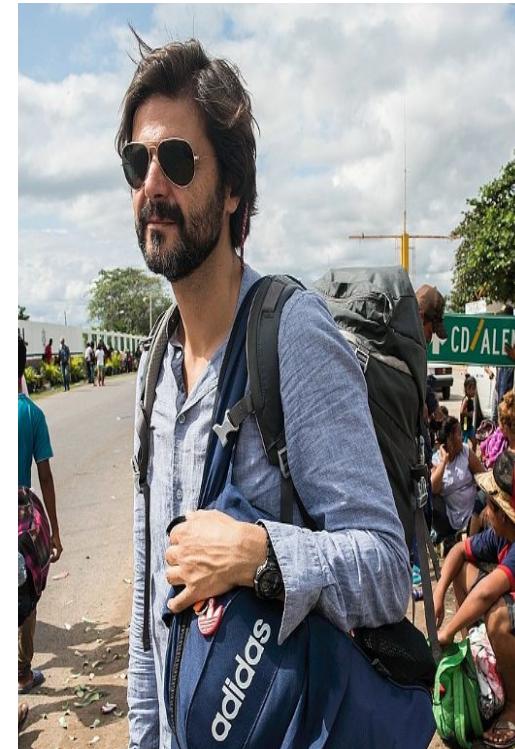
we have to
amortize too !

$$D(x, y)$$

C

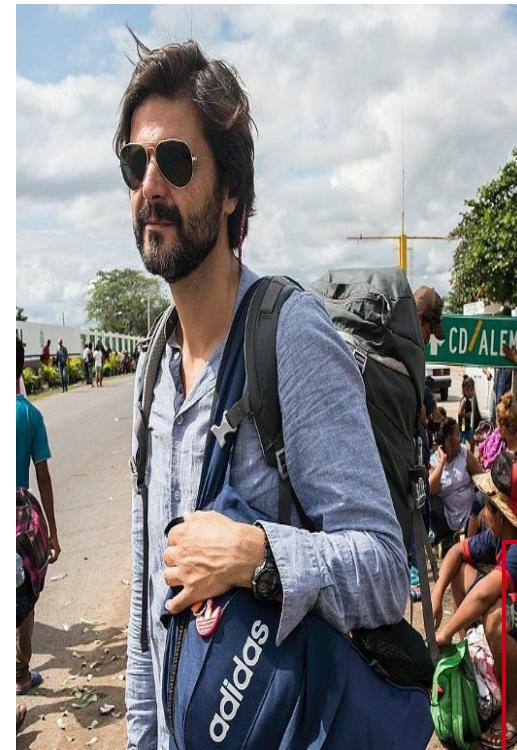
Friedman, J., T. Hastie, and R. Tibshirani, 2001: The elements of statistical learning.
Sugiyama, M., 2012: Density Ratio Estimation in Machine Learning.

Search: Find real Tobiases



Sources: Spiegel

Search: Find real Tobiases



Sources: Spiegel

Generative adversarial network

- Used as **amortized approximation** to Kullback-Leibler divergence

inference wants to
maximize disc loss

BIG
PROBLEM

- MinMax game → Nash equilibrium

Disc wants to
minimize its loss

- Unstable training (ongoing work in CS)

- Many different variants: \rightarrow standard GAN is here used, but others can be also derived by approximation
(Wasserstein GAN, boundary equilibrium GAN, relativistic GAN, gradient regularization)

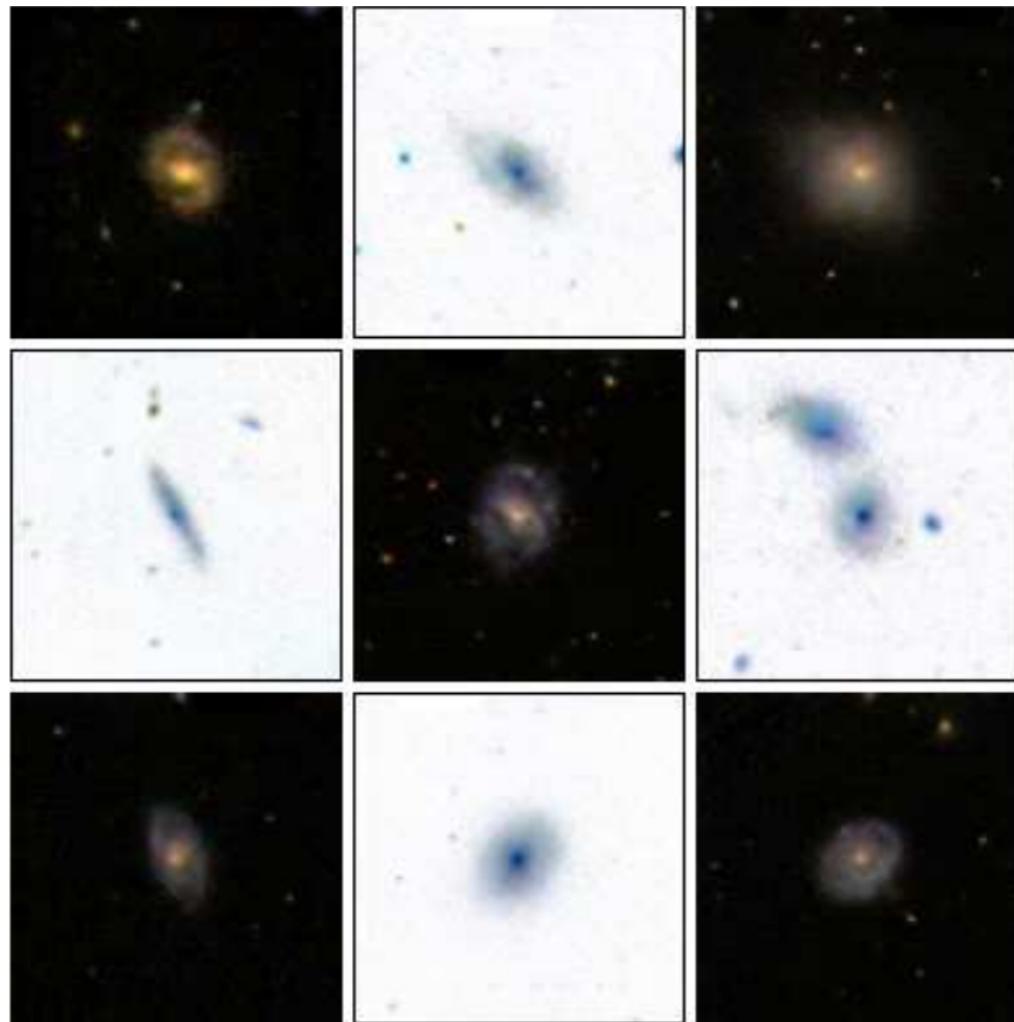
Goodfellow, I., J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, 2014:
Generative Adversarial Nets.

Examples: Generative adversarial networks (GANs)



from Karras, T., S. Laine, and T. Aila, 2018: A Style-Based Generator Architecture for Generative Adversarial Networks.

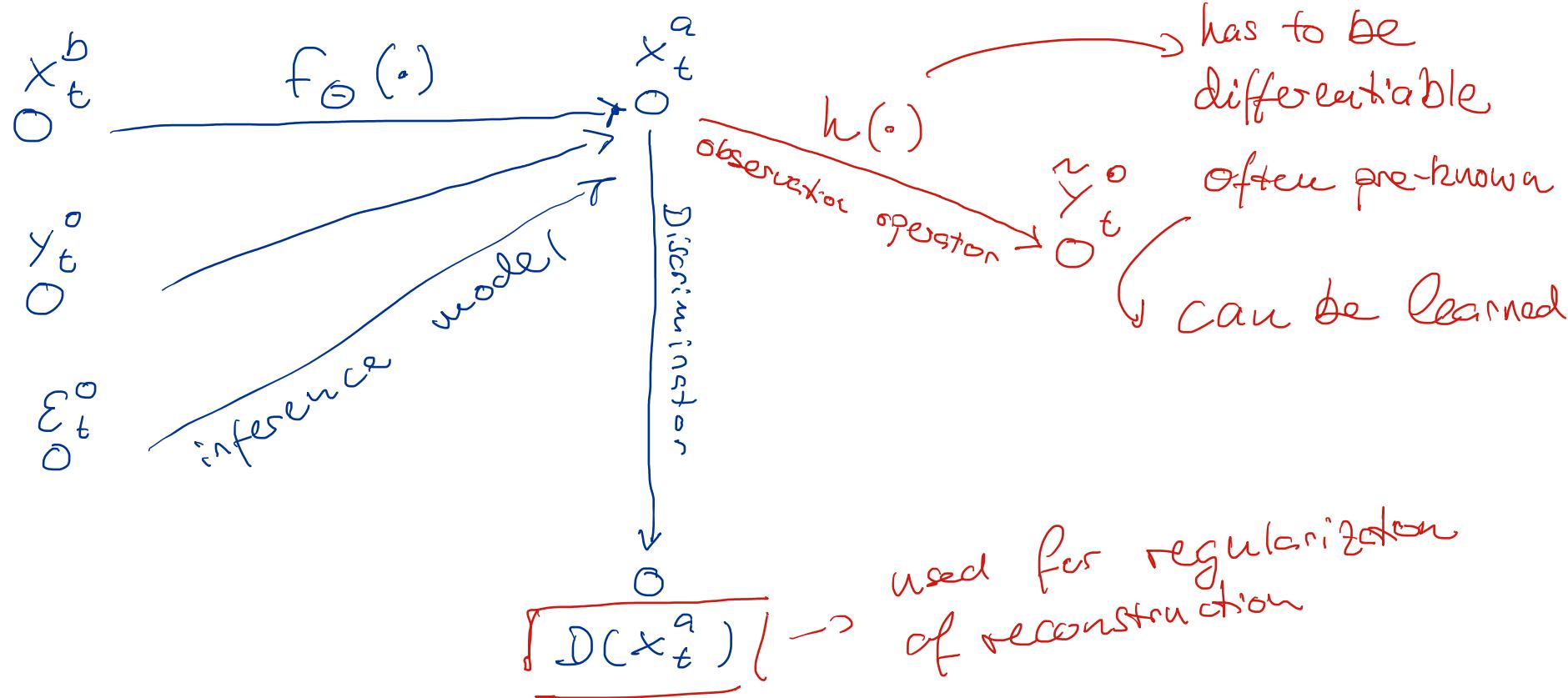
Examples: Generative adversarial networks (GANs) in physics



also
partially used
for
LHC at
CERN

from Fussell, L., and B. Moews, 2018: Forging new worlds: high-resolution synthetic galaxies with chained generative adversarial networks.

Adversarial variational Bayes for DA



⇒ Discriminator is here trained by ensemble trick:

$$x_{t,i}^a = f_\theta(y_t^o, x_{t,i}^b, \epsilon_t^o) \rightarrow \text{fake sample}$$

$$x_{t,j}^b \rightarrow \text{real sample}$$

} like adaptive contrast

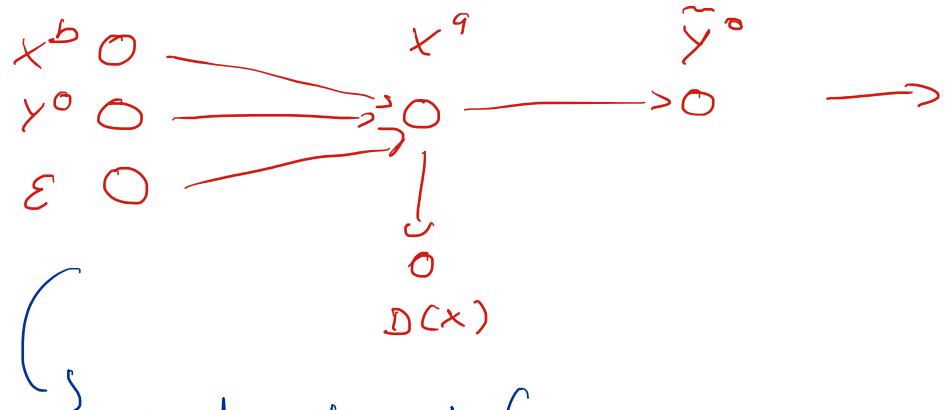
Based on – Mescheder, L., S. Nowozin, and A. Geiger, 2017:

Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks.

Experiment

- Model: Lorenz 1996, 40 grid points, RK4
- Virtual reality: $F=8$ (train), 8.1 (eval), 7.9 (test)
- Ensemble: 50 member, $F = \mathcal{N}(8, 0.25)$,
 $IC = VR + \mathcal{N}(0, 4)$
- Obs: $dt = 8$ ts, 20 observations, $VR + \mathcal{N}(0, 0.25)$
- LETKF, Gaspari-Cohn 5 grid points,
mult. inflation 1.1

Used neural network



Stochastic inference
network, no ensemble information

- prescribed obs operator
- gradient penalty for real data only (Mescheder 2018)
- Multiplicative concatenation of ϵ and (y^o, x^b)
- linear layers only
 - ↳ batch normalization
 - ↳ leaky relu
- RMS Prop

TRAINING:

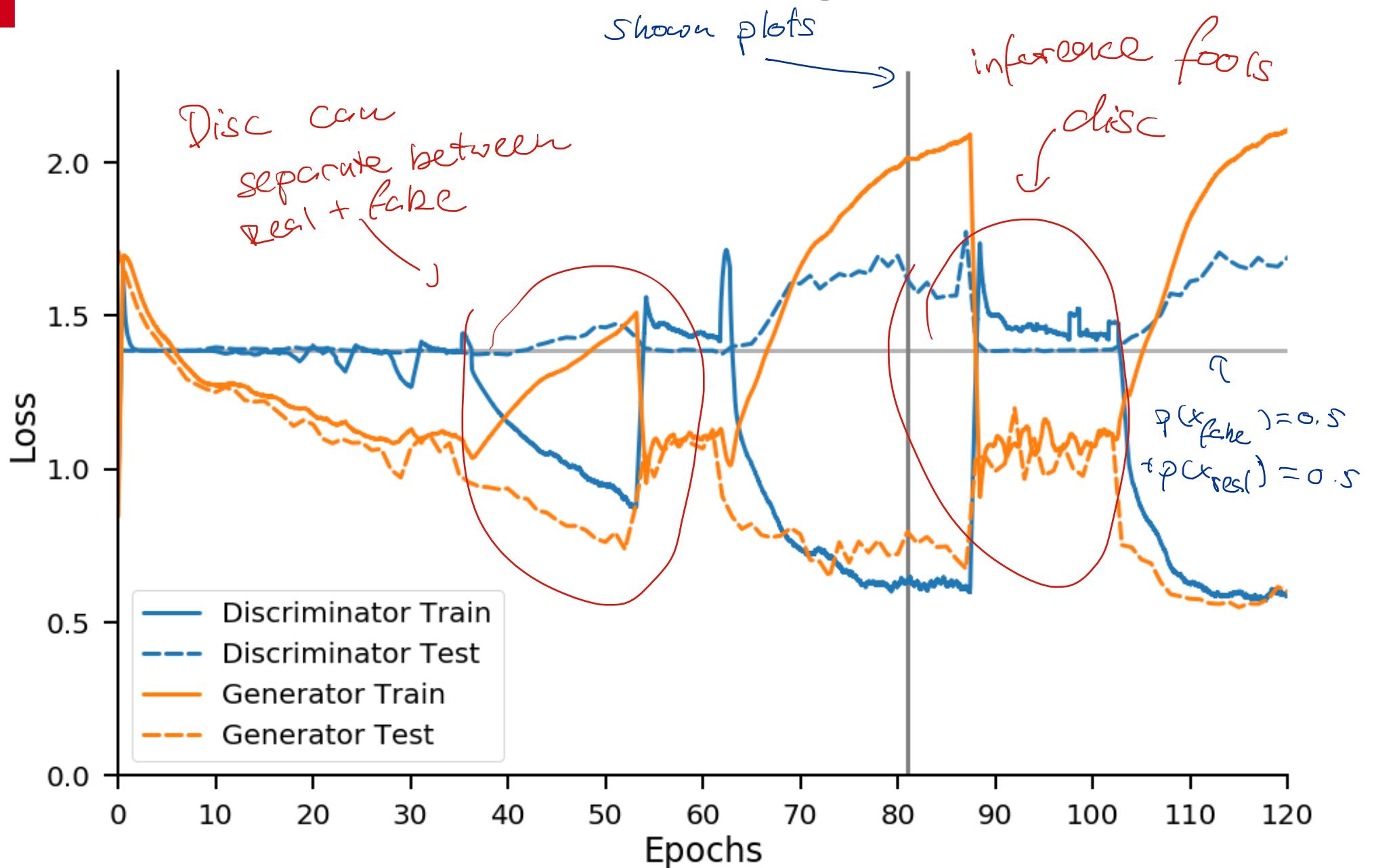
Per epoch random ensemble member selection

ASSIMILATION:

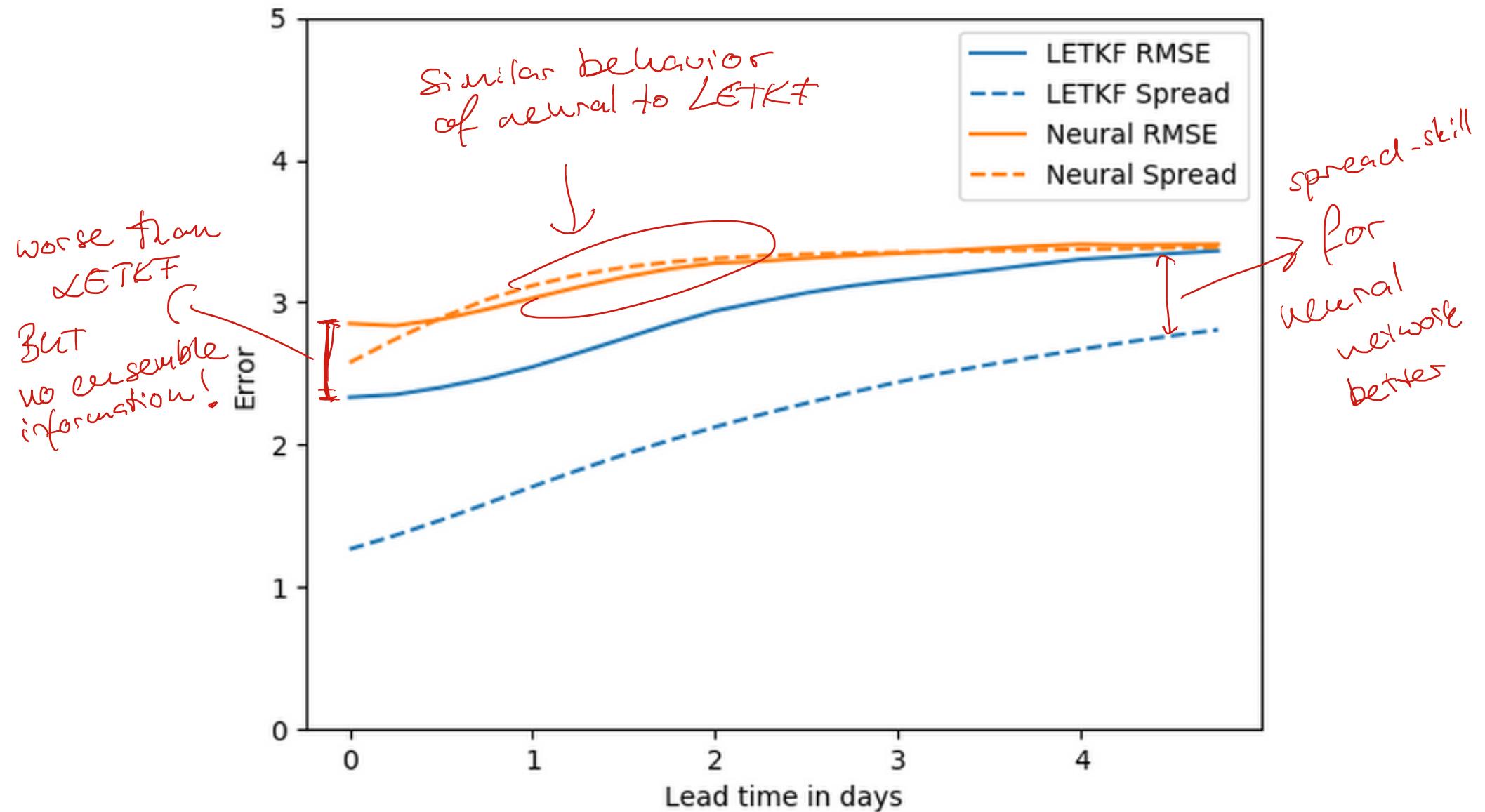
Inference for every ens member independently

like stochastic ensemble Kalman filter

Losses oscillate & diverge

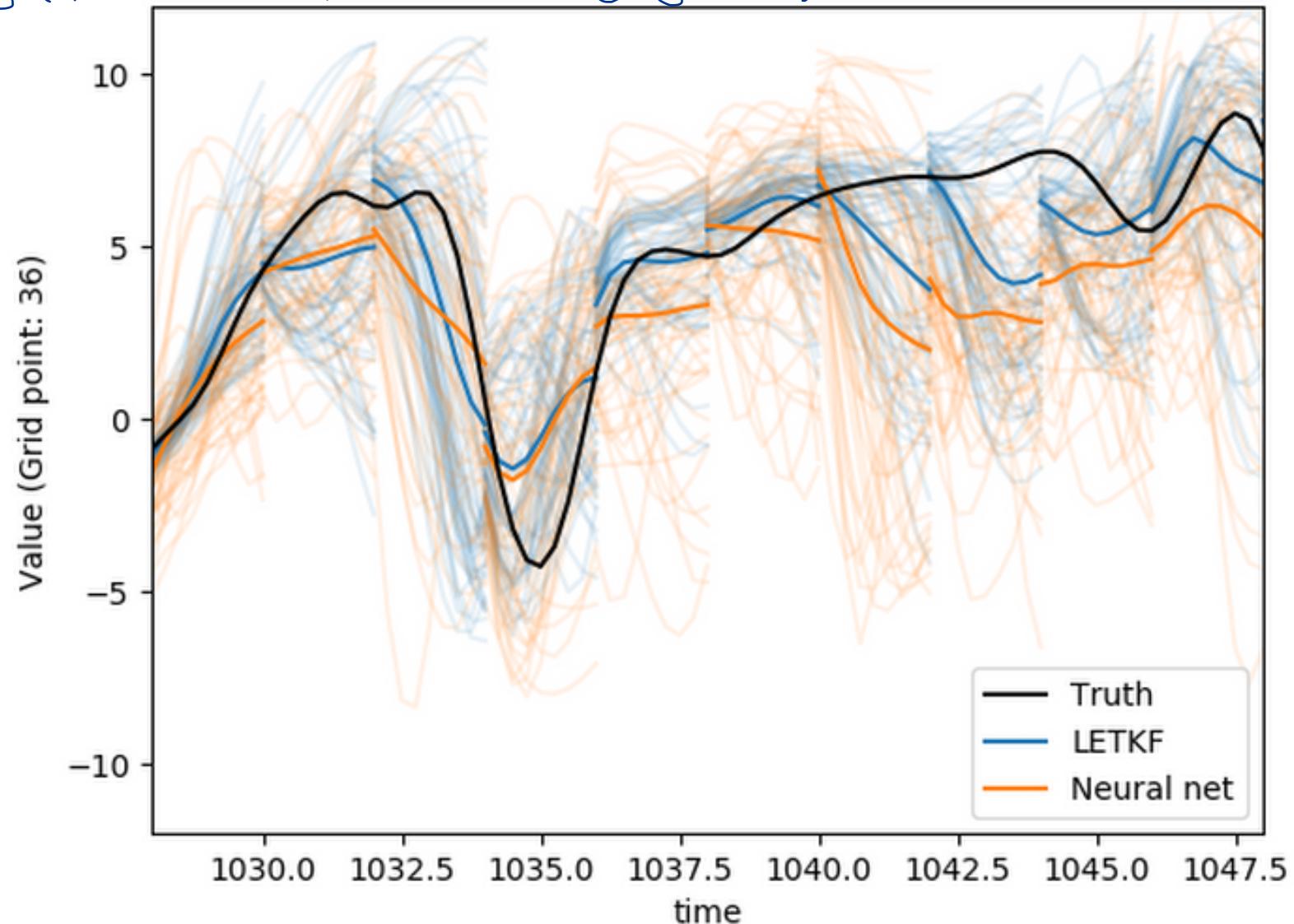


No problems with spread



Captures time evolution

Neural assimilation needs lots of crazy loss members independently
↳ in most cases it look good!



Nudged to truth



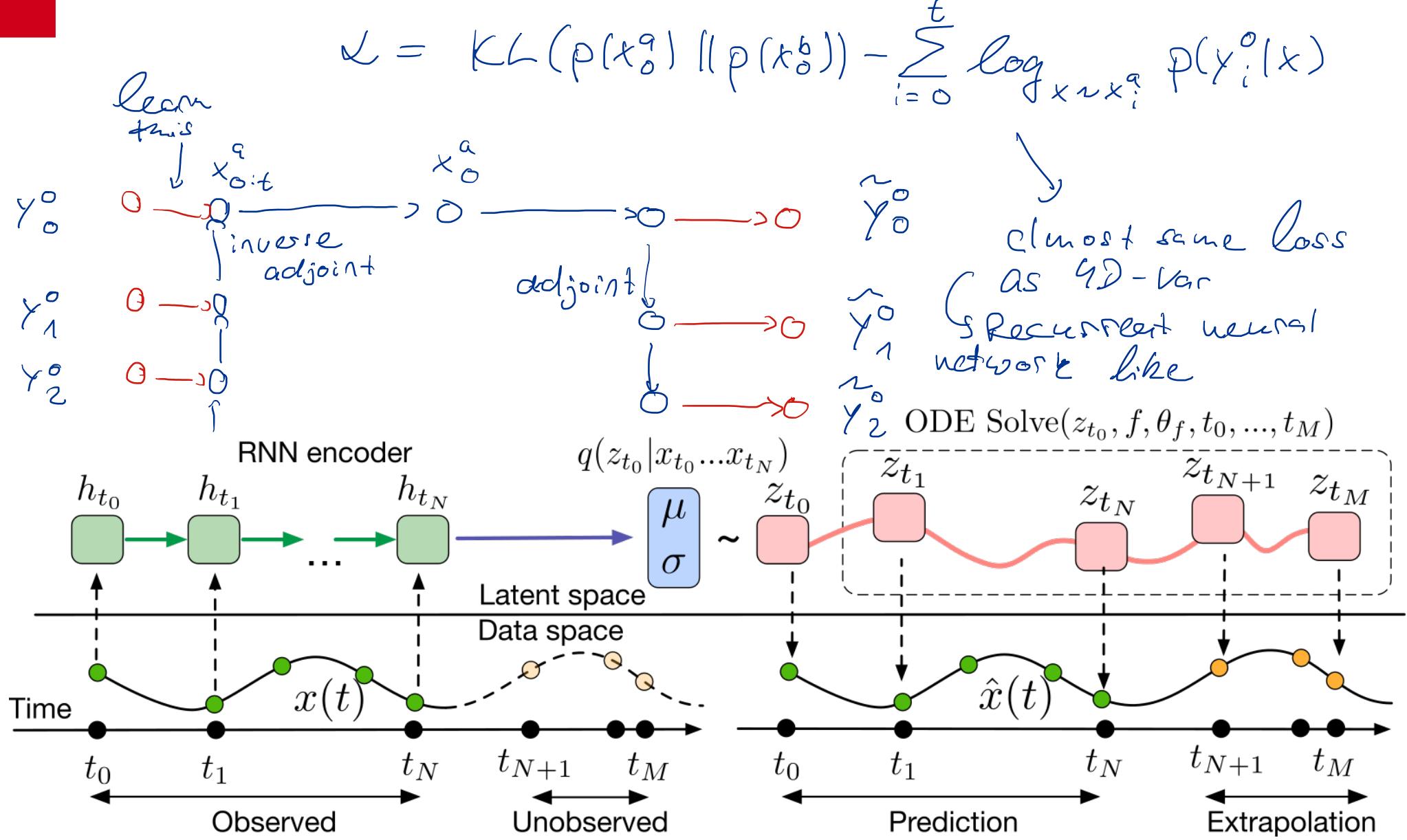
Which questions have to be answered

- How to stabilize training (open question for GANs)?
- Analytical solution for simplified problem (DiracGAN)?
Mescheder 2018
- Is unsupervised training possible?
↳ non overlapping times in samples between prior and obs during training
- What happens if observation has no impact?
↳ analysis should then be the prior
- Can we use ensemble for flow-dependent inference?
↳ include ensemble state?
- How to train in operational environment?
↳ transfer learning and recursive fitting of layers
- Can we train new observation operators?
↳ do we need to preset an observation operator?
- Can we generalize other data assimilation algorithms?
↳ this has to be proven mathematically for optimal discriminator!

Summary

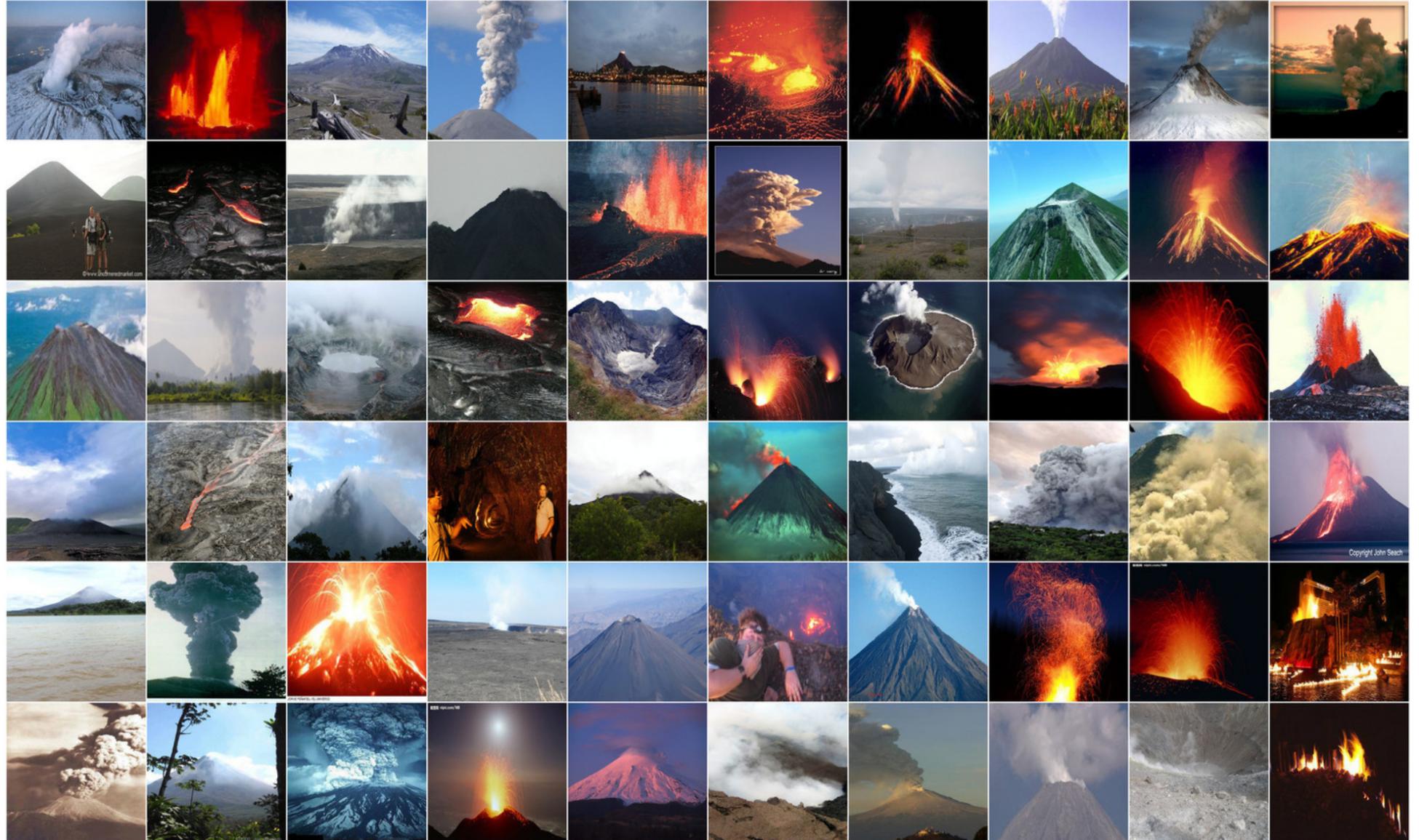
- Combined deep learning & data assimilation
- Amortized variational inference
 - + adversarial training (most problems)
- Many open questions remain
 - developments in computer science
- Can be a candidate for coupled non-linear data assimilation + to solve normal problems

4DVar as variational autoencoder



Chen, R. T. Q., Y. Rubanova, J. Bettencourt, and D. Duvenaud, 2018: Neural Ordinary Differential Equations

Generative adversarial network



Source: arxiv.org/abs/1612.00005

Generative adversarial network



Source: arxiv.org/abs/1612.00005