

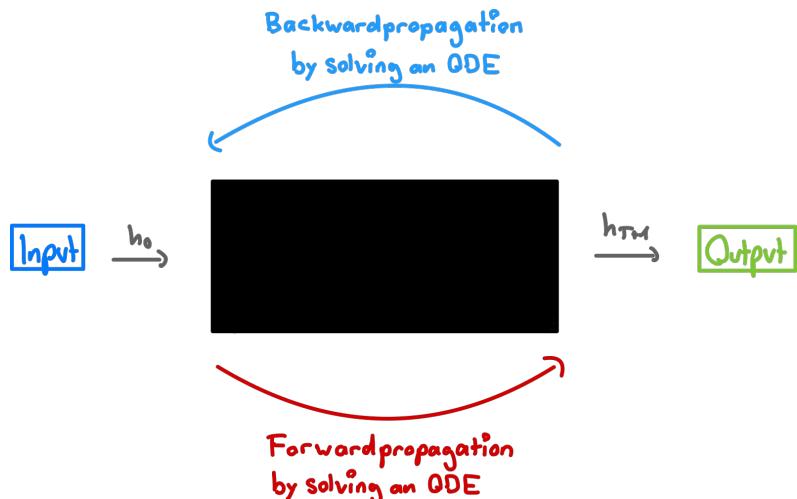
Neural ODE's

A Paper by **Ricky Chen, Julia Rubanova, Jesse Bettencourt,
David Duvenaud**

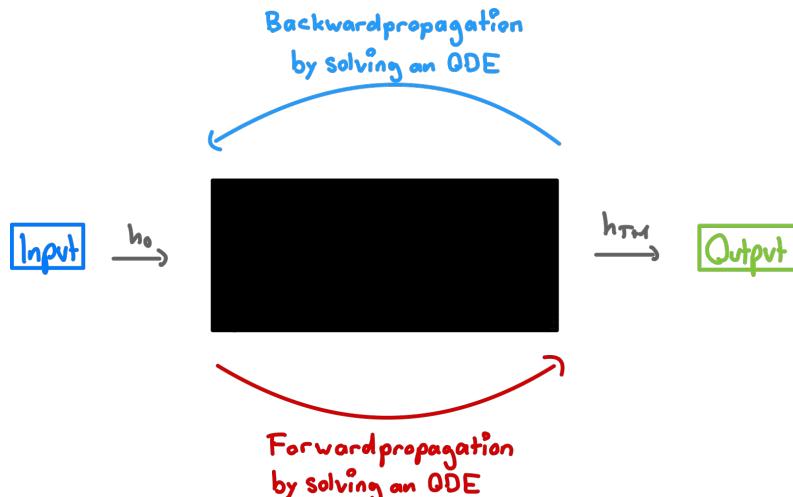
Yuri Simantob
@ysimantob



Neural ODE's: A Brief Introduction

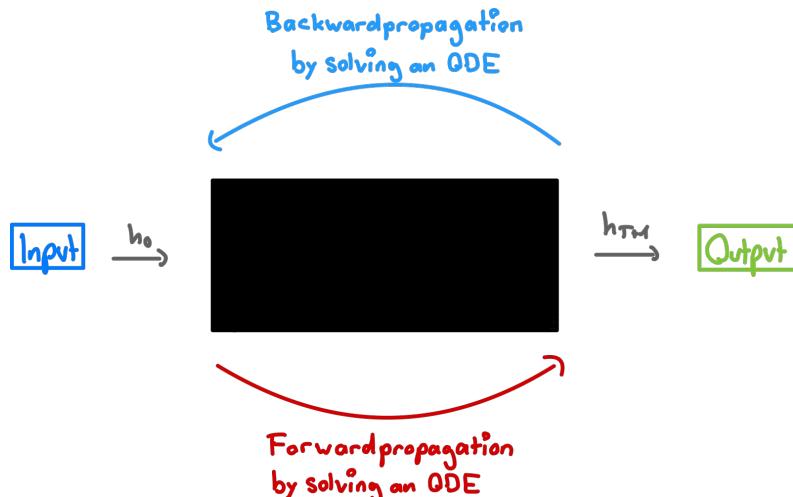


Neural ODE's: A Brief Introduction



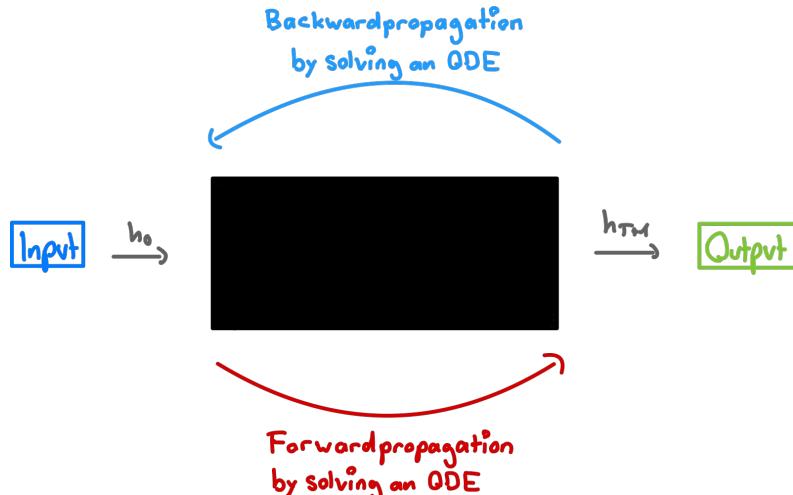
- One can trade-off accuracy with runtime

Neural ODE's: A Brief Introduction



- One can trade-off accuracy with runtime
- Simpler model

Neural ODE's: A Brief Introduction



- One can trade-off accuracy with runtime
- Simpler model
- Memory-Efficient Propagations

Overview of The Presentation

- How It Works: Forward & Backwardpropagation
- ODENets vs. ResNets
- Theoretical Application: Continuous Normalizing Flows
- Practical Application: Generative Latent Time-Series Models
- Criticism, Further Research & Applications

Recap: ODEs

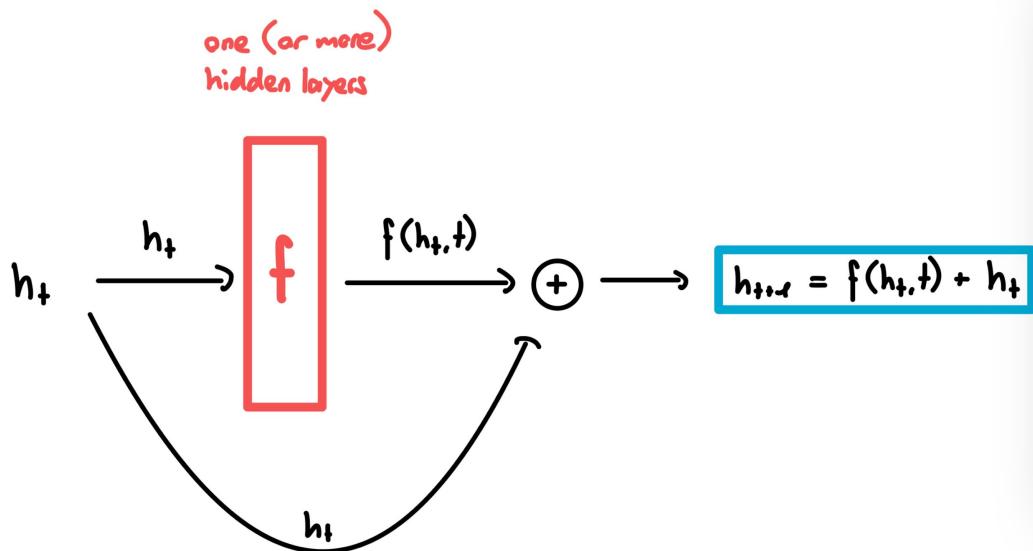
$$\frac{dx(t)}{dt} = f(x(t))$$

Recap: ODEs

$$\frac{dx(t)}{dt} = f(x(t))$$

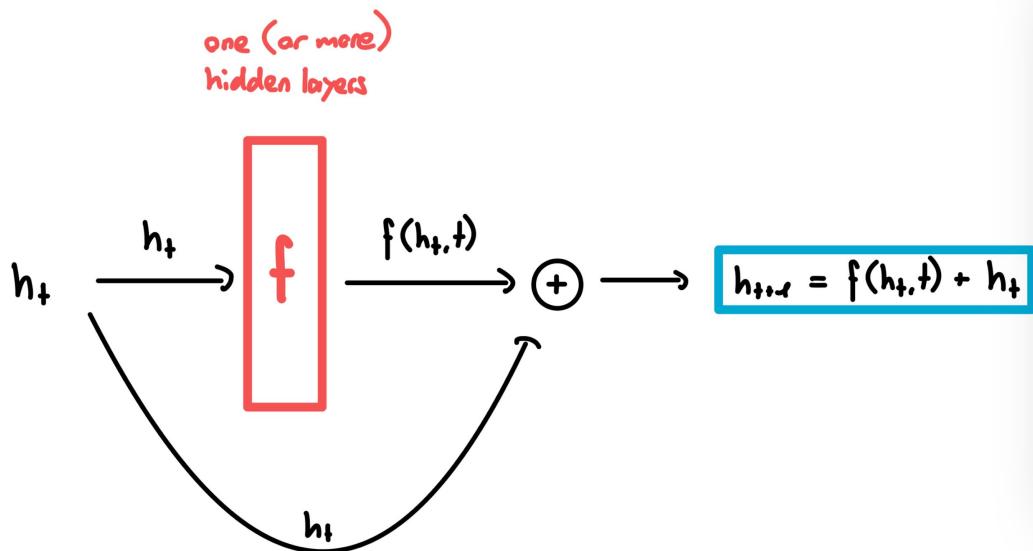


Forwardpropagation: ResNets



- addresses vanishing gradient problem
- Continues to be state-of-the-art in Computer Vision

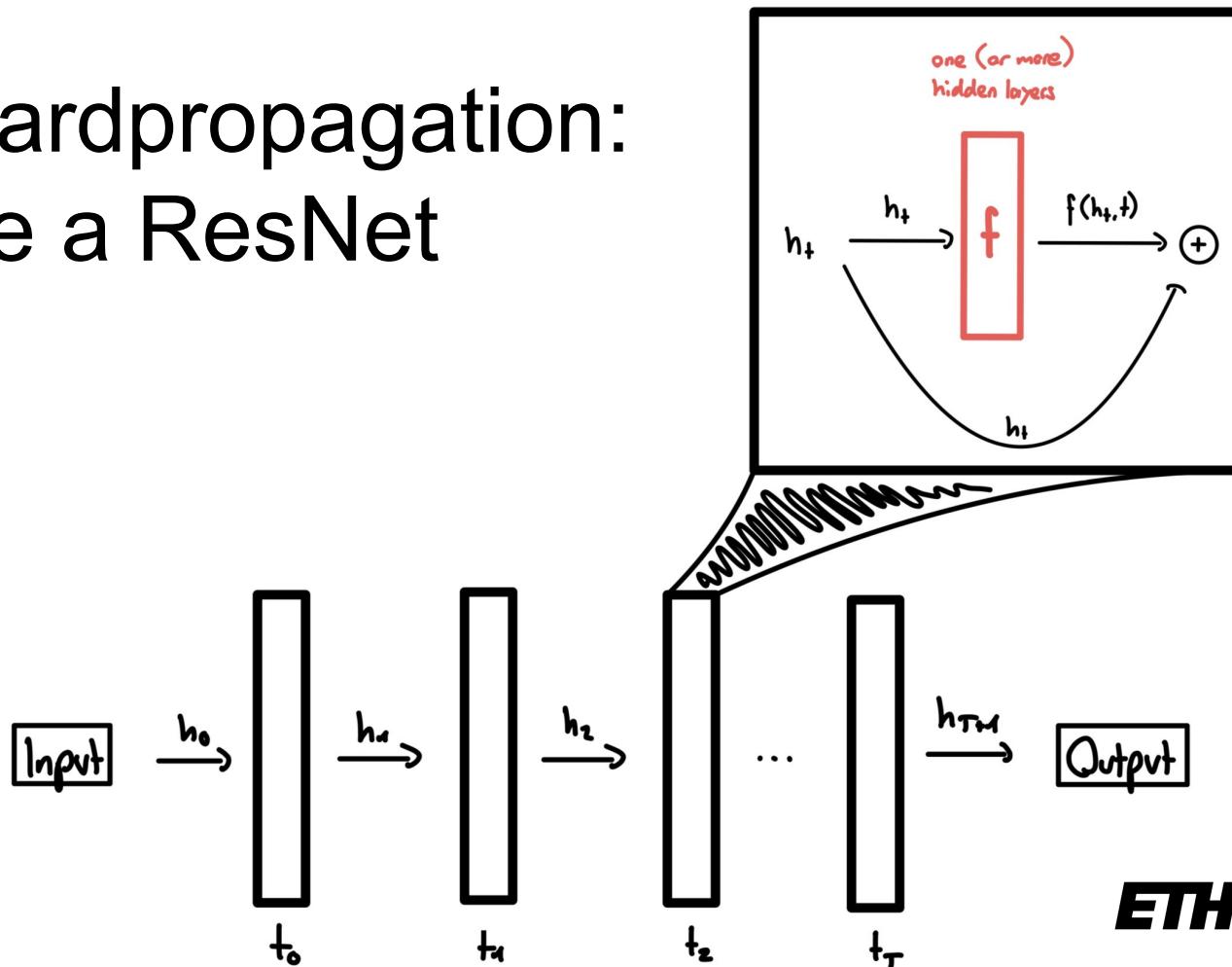
Forwardpropagation: ResNets



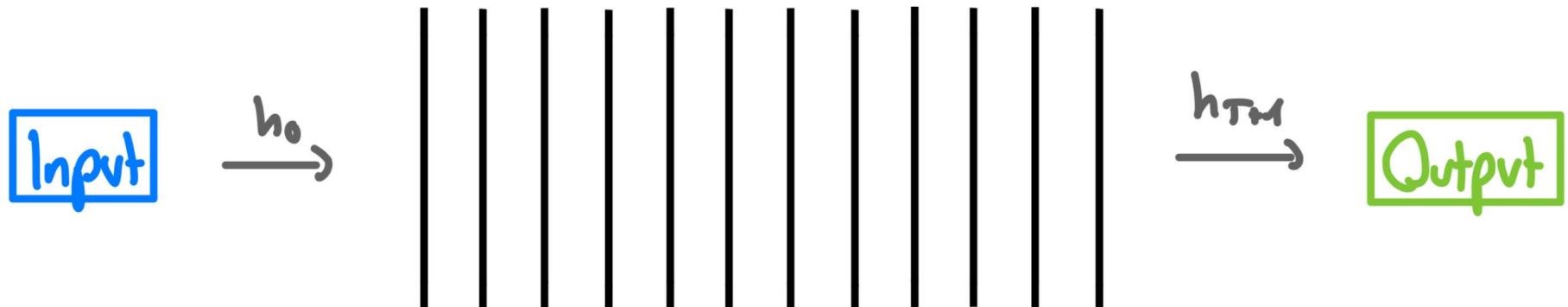
- addresses vanishing gradient problem
- Continues to be state-of-the-art in Computer Vision

$$\mathbf{h}_{t+1} = \mathbf{h}_t + f(\mathbf{h}_t, \theta_t)$$

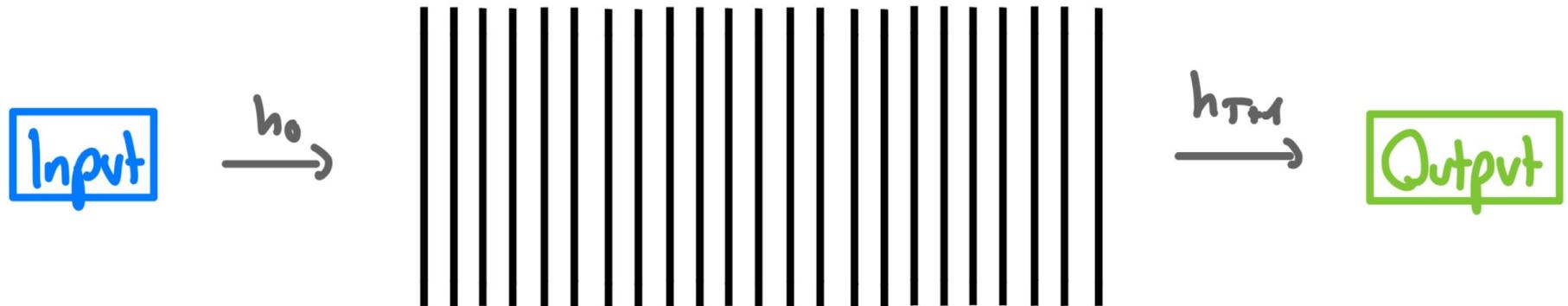
Forwardpropagation: Inside a ResNet



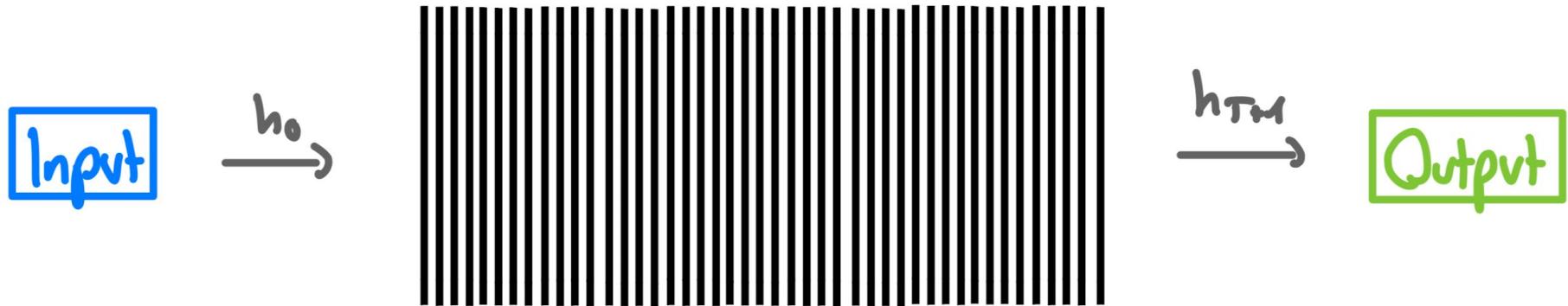
What if $T \rightarrow \infty$?



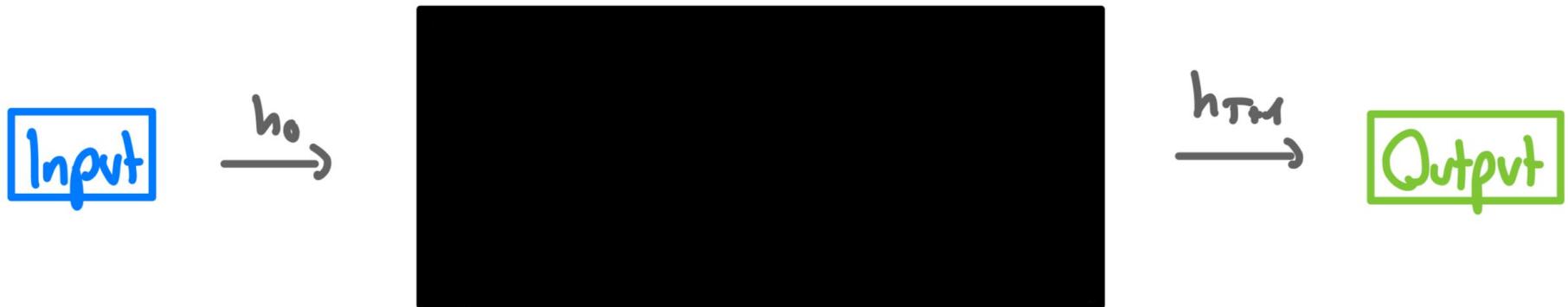
What if $T \rightarrow \infty$?



What if $T \rightarrow \infty$?

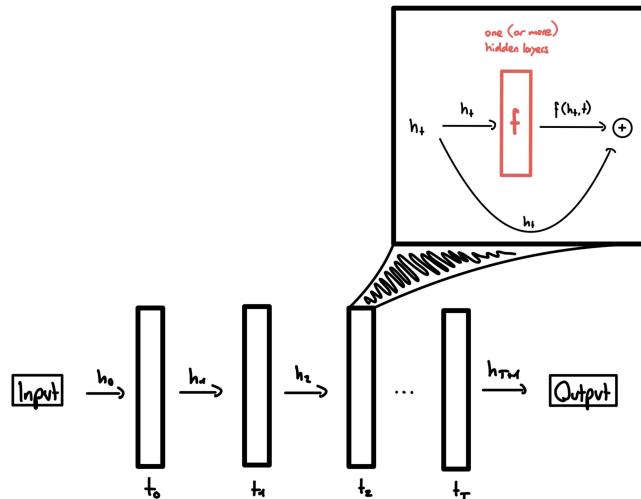


What if $T \rightarrow \infty$?



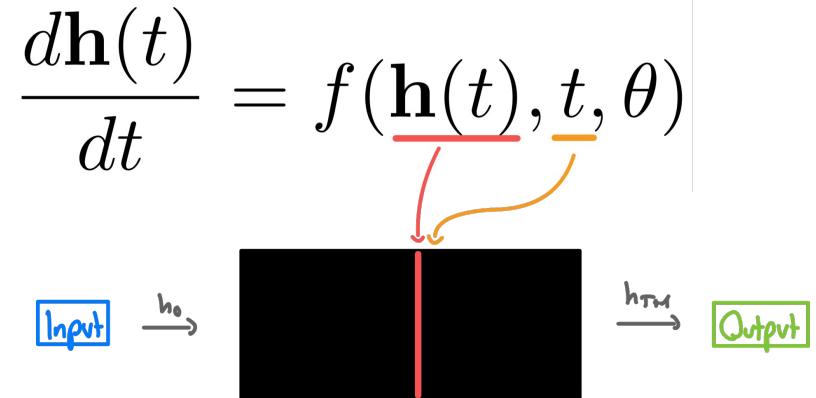
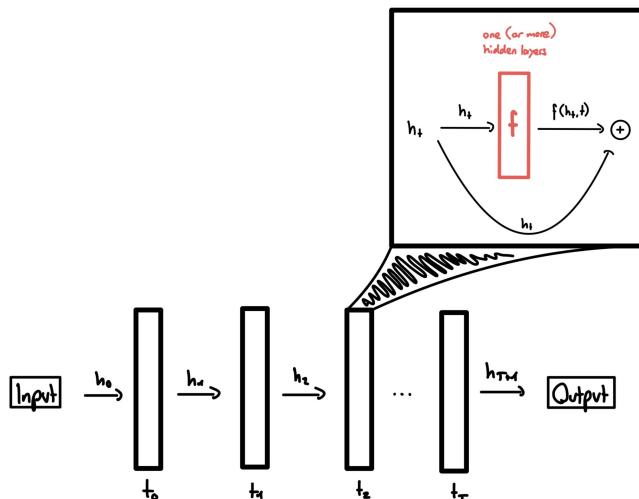
Forwardpropagation: Constructing the ODE

$$\mathbf{h}_{t+1} = \mathbf{h}_t + f(\mathbf{h}_t, \theta_t)$$



Forwardpropagation: Constructing the ODE

$$\mathbf{h}_{t+1} = \mathbf{h}_t + f(\mathbf{h}_t, \theta_t)$$



Forwardpropagation: Constructing the ODE

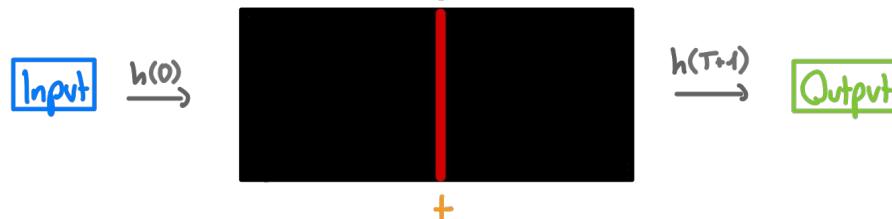
$$\frac{d\mathbf{h}(t)}{dt} = f(\underline{\mathbf{h}(t)}, \underline{t}, \theta)$$



Hidden State $\mathbf{h}(t)$
Network Depth t

Forwardpropagation: Constructing the ODE

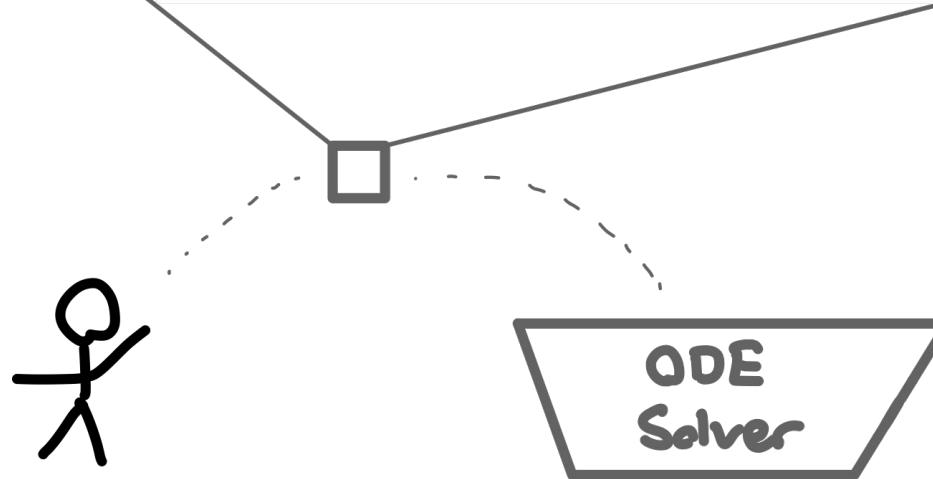
$$\frac{d\mathbf{z}(t)}{dt} = f(\underline{\mathbf{z}(t)}, \underline{t}, \theta)$$



Hidden State $\mathbf{z}(t)$
Network Depth t

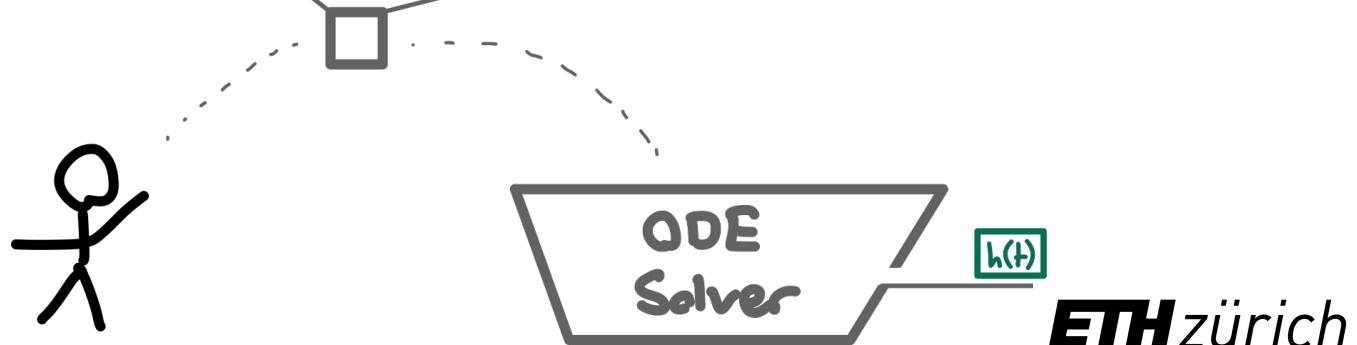
Forwardpropagation: Solving the ODE

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta)$$



Forwardpropagation: Solving the ODE

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta)$$



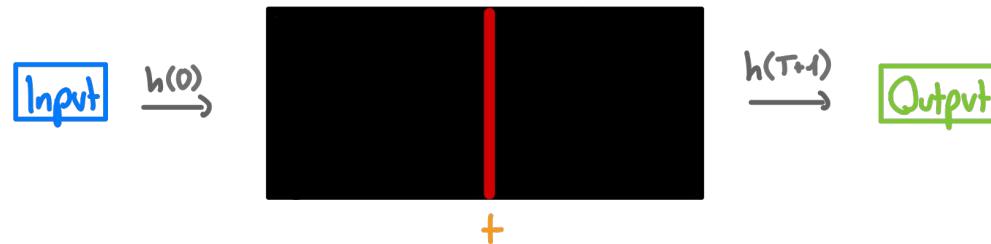
Forwardpropagation: Getting the Output

For a ODE network with depth T:

$h(0)$: Input

$h(t)$: Value of the hidden "layer" at depth t

$h(T+1)$: Output



Backpropagation

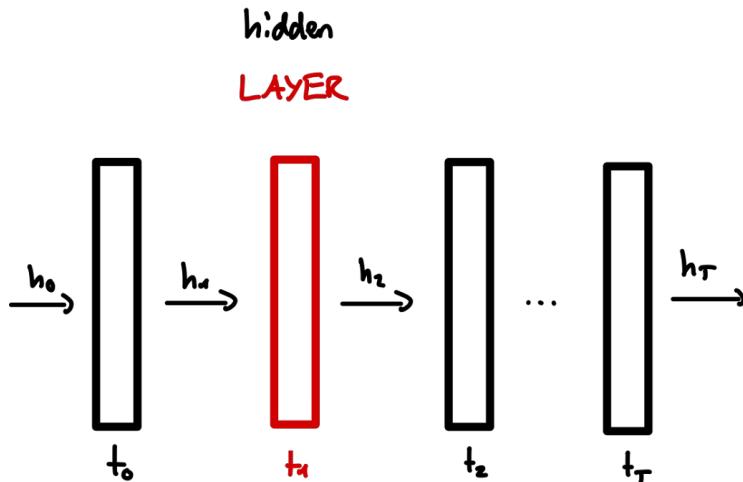
Backpropagation: Notation

Discrete

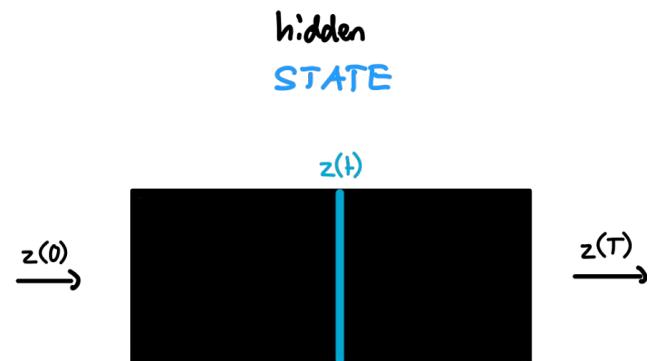
Continuous

Backpropagation: Notation

Discrete

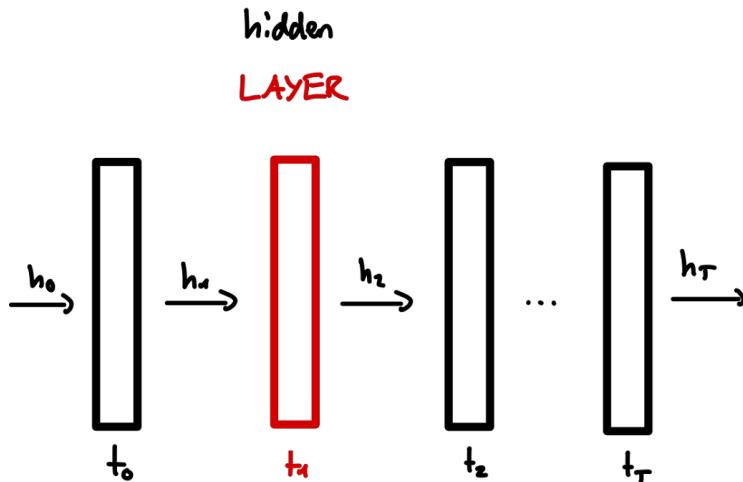


Continuous

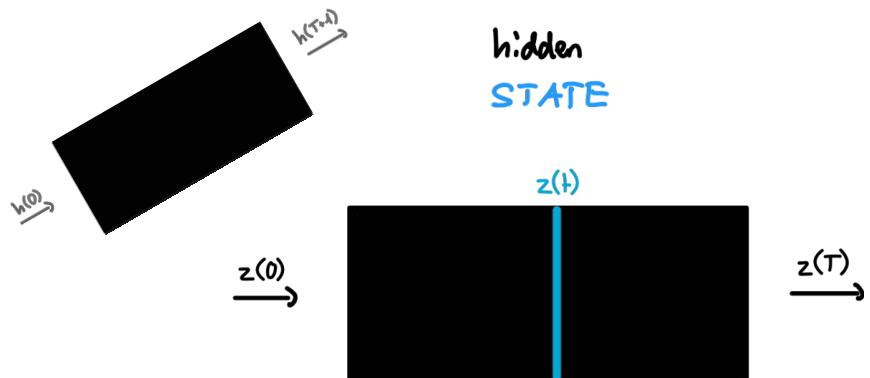


Backpropagation: Notation

Discrete



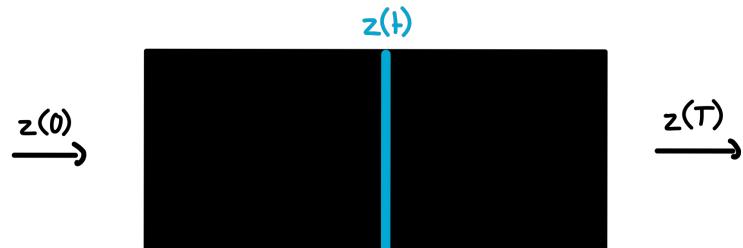
Continuous



Backpropagation: The Adjoint

$$\alpha(t) = \frac{dL}{dz(t)}$$

or how much a hidden state $z(t)$ contributes to the loss L



Backpropagation: Solving for the Adjoint

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t)^\top \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}}$$

Backpropagation: Solving for the Adjoint

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t)^\top \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}}$$



ETH zürich

Backpropagation: Weight Updates

$a(t)$



$$\frac{dL}{d\theta} = \int_{t_1}^{t_0} \mathbf{a}(t)^\top \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \theta} dt$$

Backpropagation: Weight Updates

$a(t)$



$$\frac{dL}{d\theta} = \int_{t_1}^{t_0} \mathbf{a}(t)^\top \frac{\partial f(\underline{\mathbf{z}(t)}, t, \theta)}{\partial \theta} dt$$

$\frac{d\mathbf{z}(t)}{dt} = f(\mathbf{z}(t), t, \theta)$

Backpropagation: Weight Updates

$a(t)$



Our pre-defined
model dynamics

$$\frac{d\mathbf{z}(t)}{dt} = f(\mathbf{z}(t), t, \theta)$$

$$\frac{dL}{d\theta} = \int_{t_1}^{t_0} \mathbf{a}(t)^\top \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \theta} dt$$

Backpropagation: Weight Updates

a(t)

$$\frac{dL}{d\theta} = \int_{t_1}^{t_0} \underline{\mathbf{a}(t)^\top} \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \theta} dt$$

Our pre-defined model dynamics

$$\frac{d\mathbf{z}(t)}{dt} = f(\mathbf{z}(t), t, \theta)$$
$$\frac{da(t)}{dt} = -\mathbf{a}(t)^\top \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}}$$

Advantages of Continuous Propagation

Advantages of Continuous Propagation

- Forward propagation requires one call to an ODE solver

Advantages of Continuous Propagation

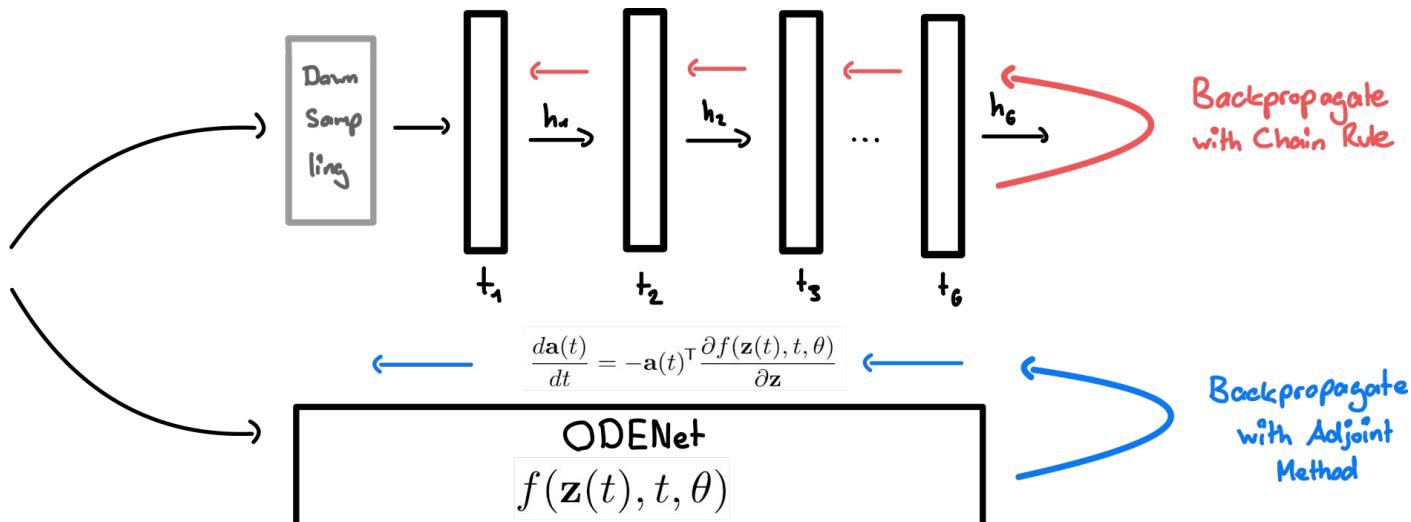
- Forward propagation requires one call to an ODE solver
- Backward propagation requires one call to an ODE solver

Advantages of Continuous Propagation

- Forward propagation requires one call to an ODE solver
- Backward propagation requires one call to an ODE solver
- Runtime & Memory independent of network depth

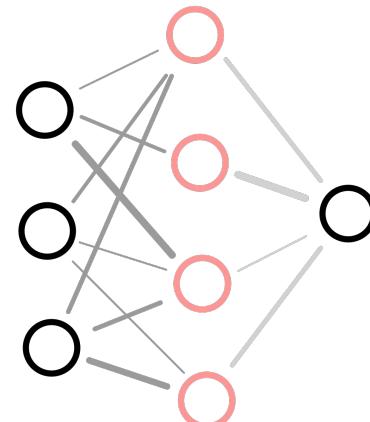
ODENets vs. ResNets on MNIST

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9



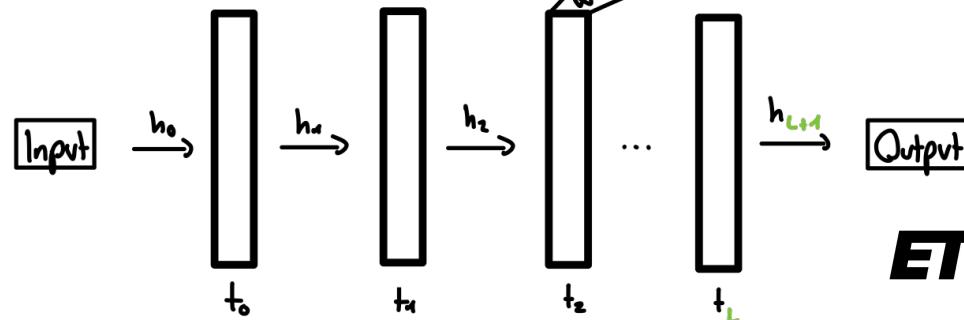
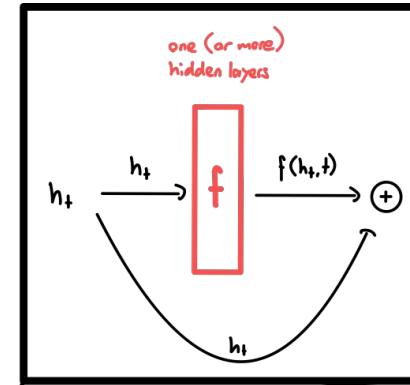
Results: 1-Layer Perceptron

	Test Error	# Params	Memory	Time
1-Layer MLP [†]	1.60%	0.24 M	-	-
ResNet	0.41%	0.60 M	$\mathcal{O}(L)$	$\mathcal{O}(L)$
ODE-Net	0.42%	0.22 M	$\mathcal{O}(1)$	$\mathcal{O}(\tilde{L})$



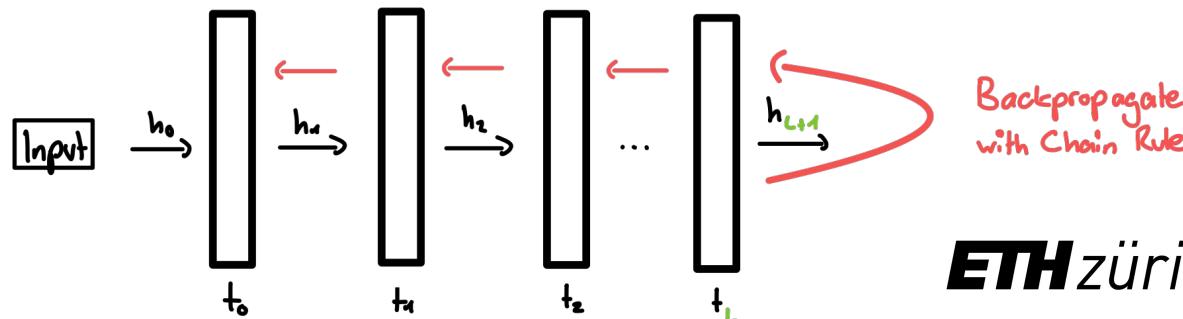
Results: ResNet

	Test Error	# Params	Memory	Time
1-Layer MLP [†]	1.60%	0.24 M	-	-
ResNet	0.41%	0.60 M	$\mathcal{O}(L)$	$\mathcal{O}(L)$
ODE-Net	0.42%	0.22 M	$\mathcal{O}(1)$	$\mathcal{O}(\bar{L})$



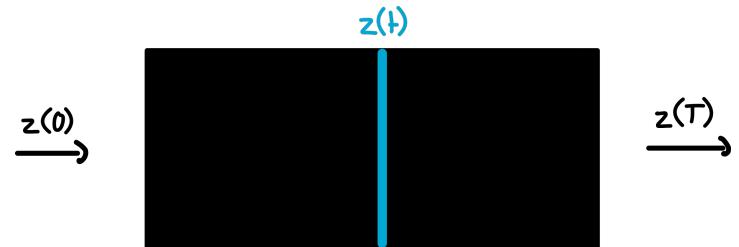
Results: ResNet

	Test Error	# Params	Memory	Time
1-Layer MLP [†]	1.60%	0.24 M	-	-
ResNet	0.41%	0.60 M	$\mathcal{O}(L)$	$\mathcal{O}(L)$
ODE-Net	0.42%	0.22 M	$\mathcal{O}(1)$	$\mathcal{O}(\tilde{L})$



Results: ResNet

	Test Error	# Params	Memory	Time
1-Layer MLP [†]	1.60%	0.24 M	-	-
ResNet	0.41%	0.60 M	$\mathcal{O}(L)$	$\mathcal{O}(L)$
ODE-Net	0.42%	0.22 M	$\mathcal{O}(1)$	$\mathcal{O}(\tilde{L})$

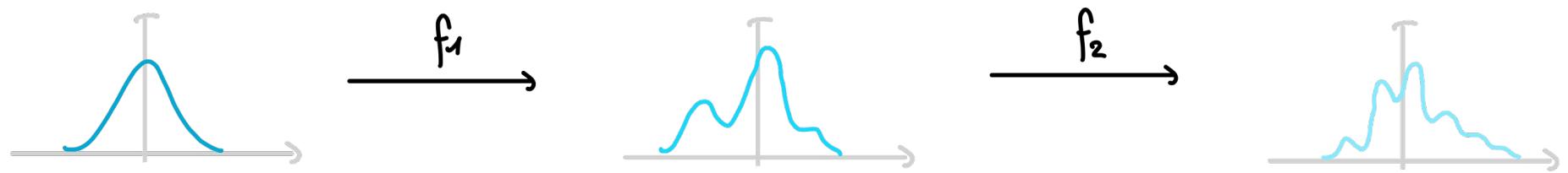


ODENets vs. ResNets: The Results

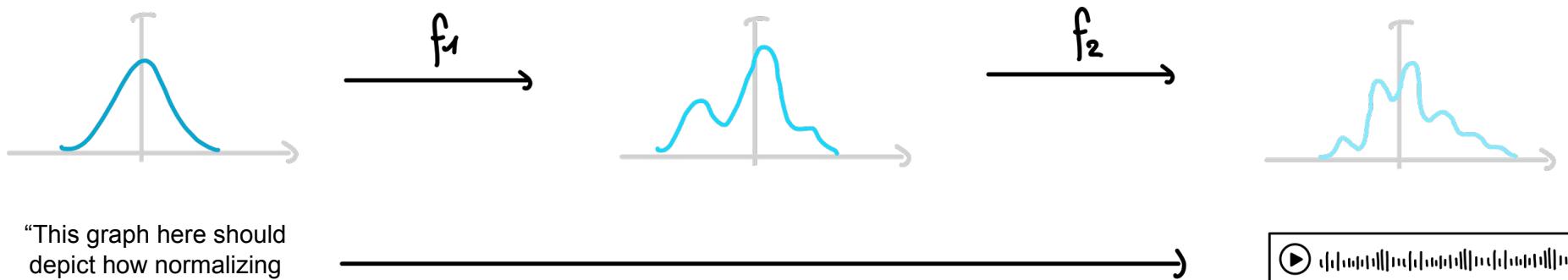
- constant memory with respect to number of layers
 - fewer parameters
 - runtime is linear with respect to the number of evaluations
-
- linear memory with respect to number of layers
 - more parameters
 - runtime is linear with respect to the number of layers

Theoretical Application: Continuous Normalizing Flows

Normalizing Flows: Transforming Distributions



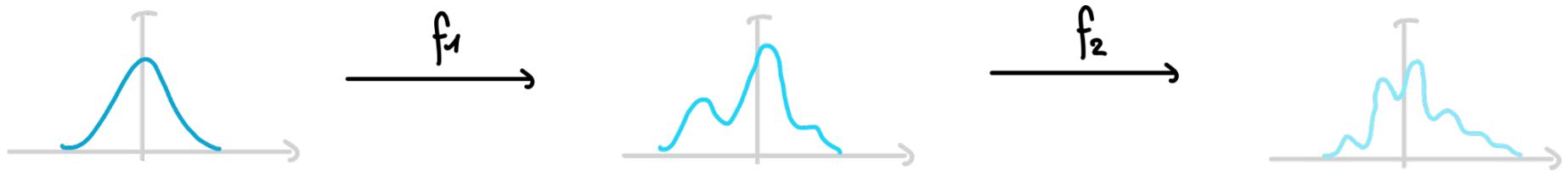
Normalizing Flows: Transforming Distributions



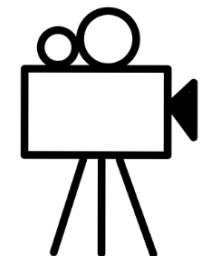
"This graph here should depict how normalizing flows create complex distributions from simple ones"



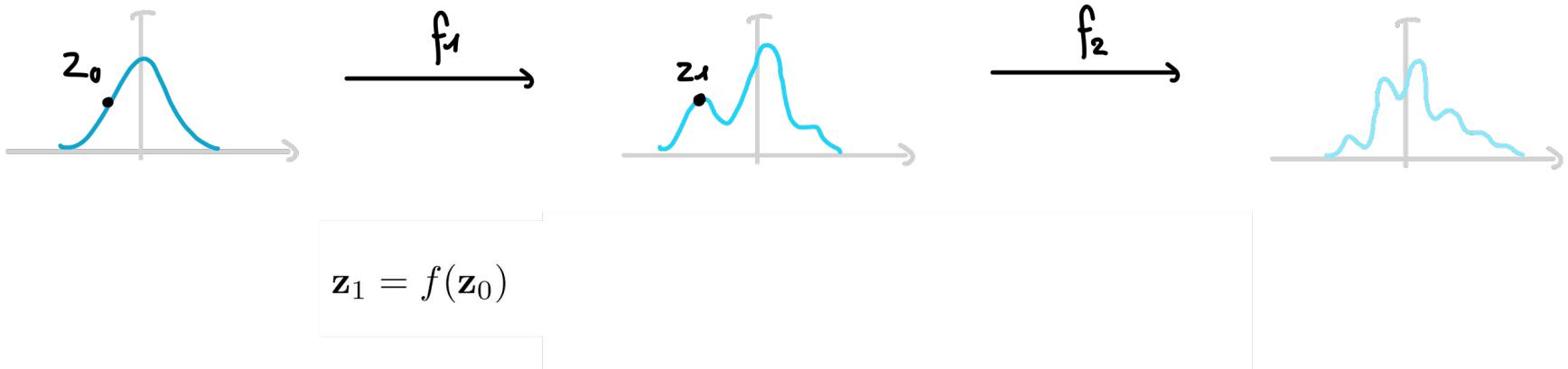
Normalizing Flows: Transforming Distributions



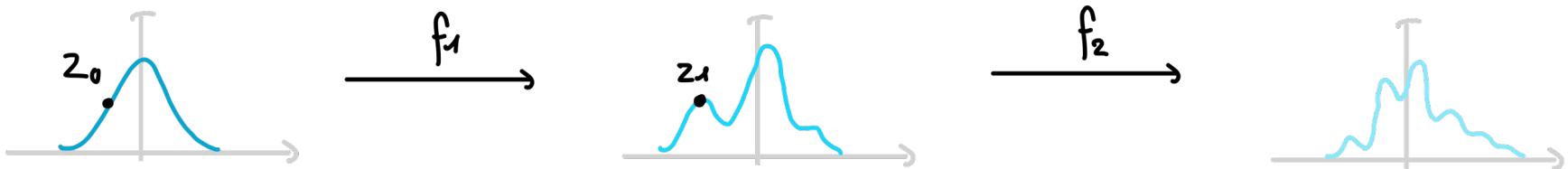
"This graph here should depict how normalizing flows create complex distributions from simple ones"



Normalizing Flows: The Main Bottleneck

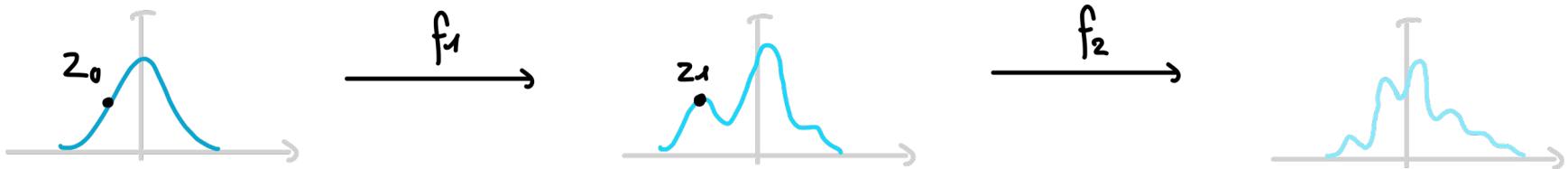


Normalizing Flows: The Main Bottleneck



$$\mathbf{z}_1 = f(\mathbf{z}_0) \implies \log p(\mathbf{z}_1) = \log p(\mathbf{z}_0) - \log \left| \det \frac{\partial f}{\partial \mathbf{z}_0} \right|$$

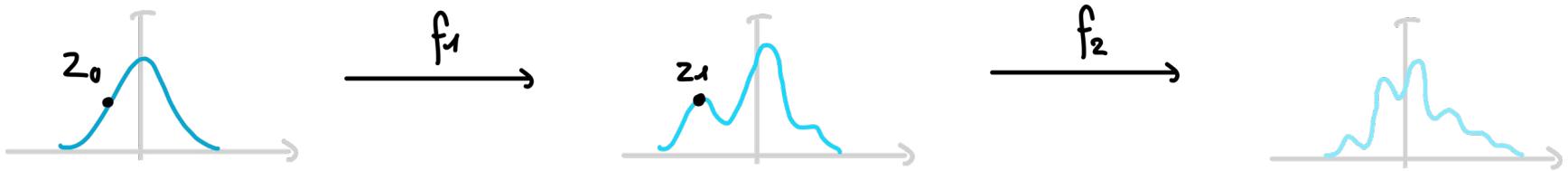
Normalizing Flows: The Main Bottleneck



$$\mathbf{z}_1 = f(\mathbf{z}_0) \implies \log p(\mathbf{z}_1) = \log p(\mathbf{z}_0) - \log \left| \det \frac{\partial f}{\partial \mathbf{z}_0} \right|$$

$\mathcal{O}(n^3)$!

Normalizing Flows: The Main Bottleneck



$$\mathbf{z}_1 = f(\mathbf{z}_0) \implies \log p(\mathbf{z}_1) = \log p(\mathbf{z}_0) - \log \left| \det \frac{\partial f}{\partial \mathbf{z}_0} \right|$$

$$\mathbf{h}_{t+1} = \mathbf{h}_t + f(\mathbf{h}_t, \theta_t)$$

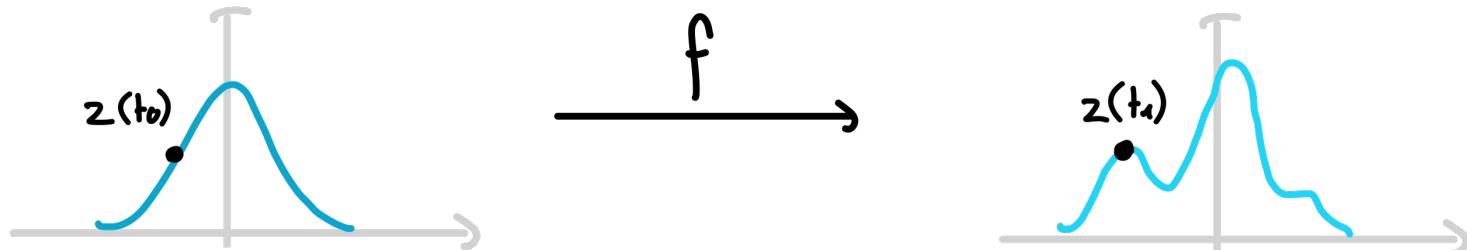
Continuous Normalizing Flows: Making It Efficient

Theorem 1 (Instantaneous Change of Variables). *Let $\mathbf{z}(t)$ be a finite continuous random variable with probability $p(\mathbf{z}(t))$ dependent on time. Let $\frac{d\mathbf{z}}{dt} = f(\mathbf{z}(t), t)$ be a differential equation describing a continuous-in-time transformation of $\mathbf{z}(t)$. Assuming that f is uniformly Lipschitz continuous in \mathbf{z} and continuous in t , then the change in log probability also follows a differential equation,*

$$\frac{\partial \log p(\mathbf{z}(t))}{\partial t} = -\text{tr} \left(\frac{df}{d\mathbf{z}(t)} \right) \quad (8)$$

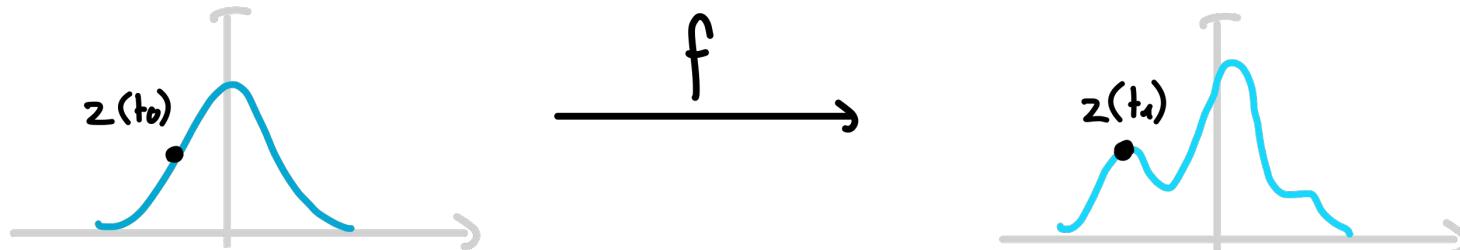
Continuous Normalizing Flows: Making It Efficient

$$\frac{\partial \log p(\mathbf{z}(t))}{\partial t} = -\text{tr} \left(\frac{df}{d\mathbf{z}(t)} \right)$$

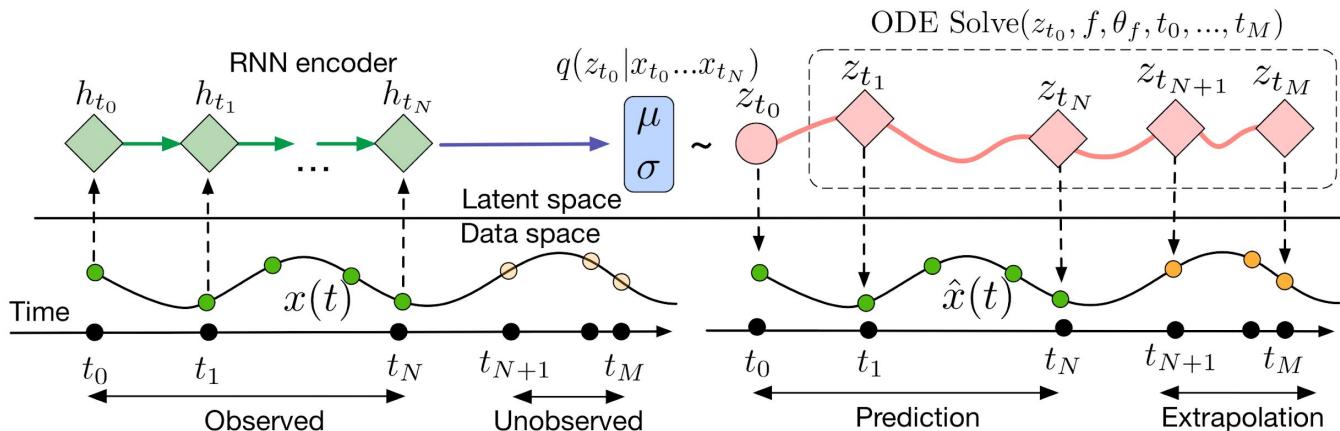


Continuous Normalizing Flows: Making It Efficient

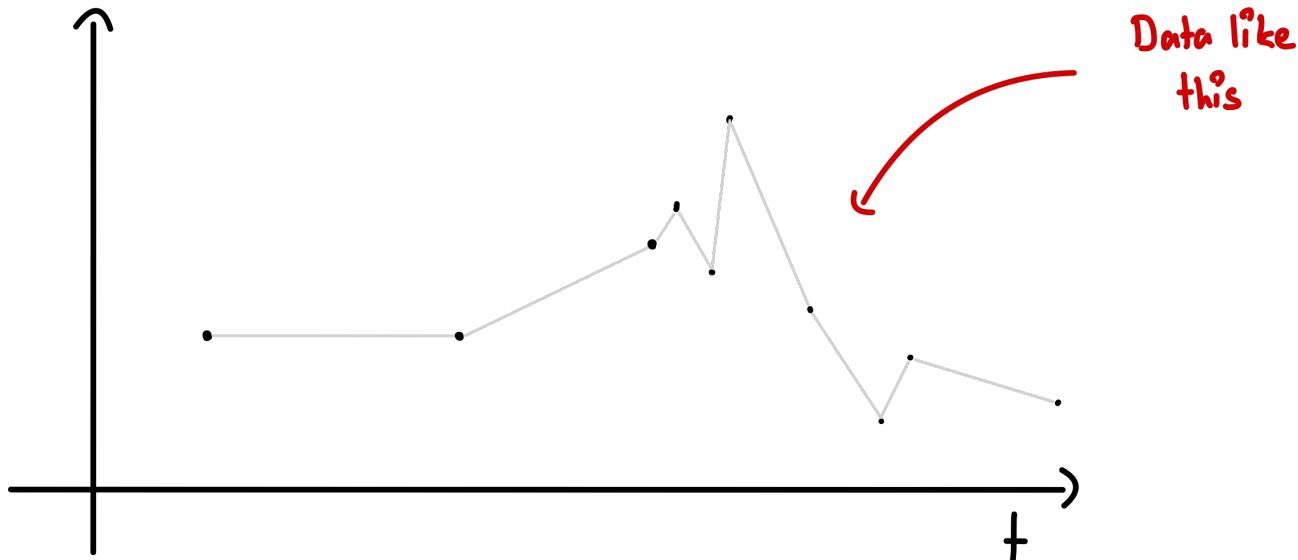
$$\frac{\partial \log p(\mathbf{z}(t))}{\partial t} = -\text{tr} \left(\frac{df}{d\mathbf{z}(t)} \right) \text{O(n)} !!$$



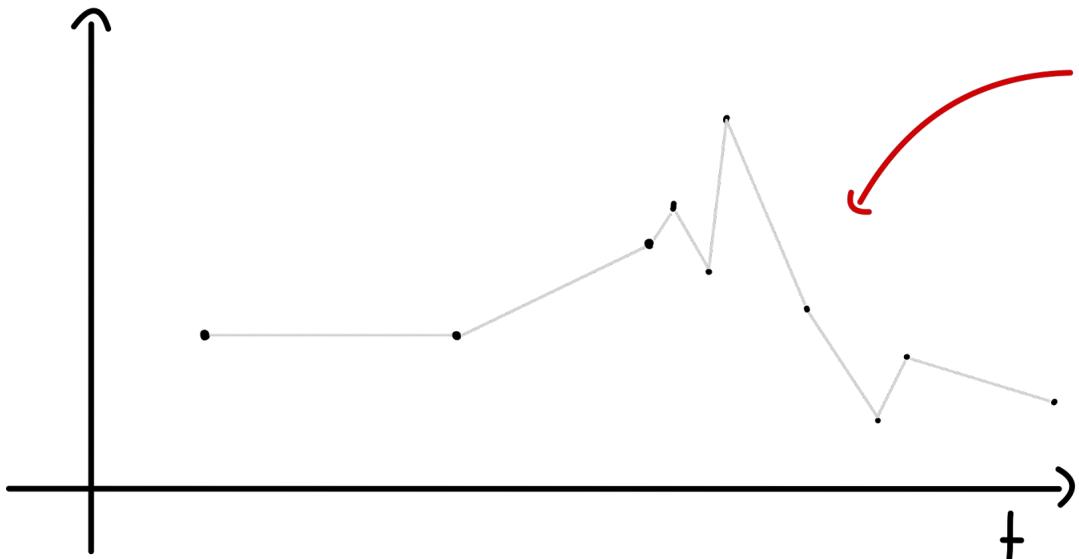
Generative Latent Time-Series Model



GM: The Problem



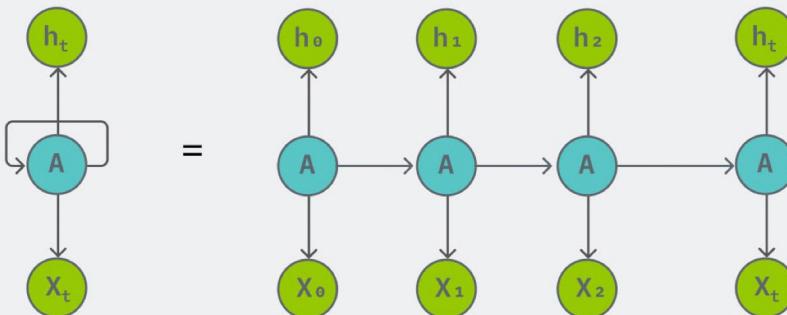
GM: The Problem



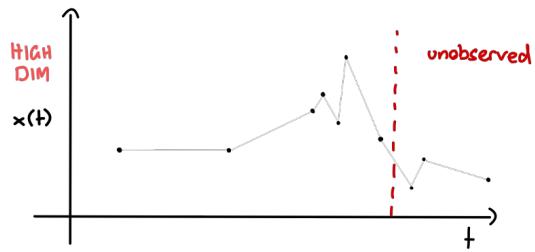
Data like
this

Irregular Time-Series
Data

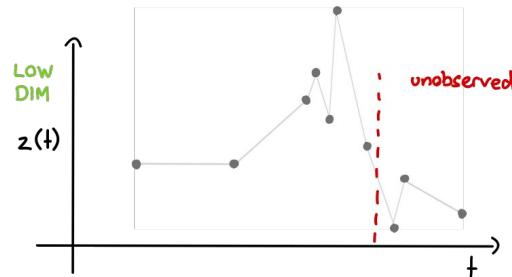
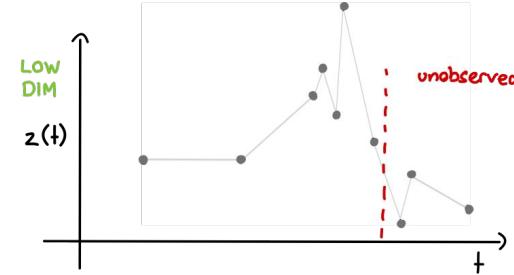
GM: The Common Solution



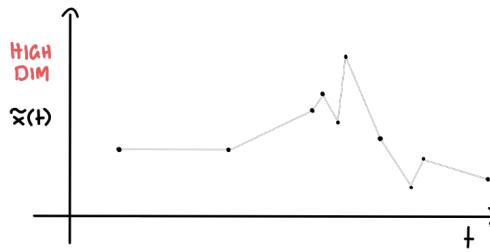
- Using a Recurrent Neural Network (RNN)

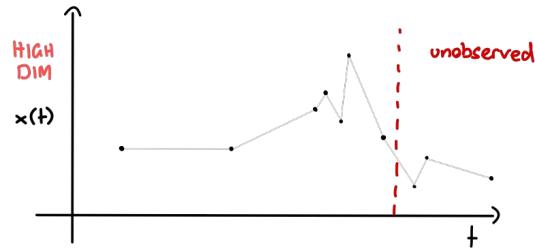


ENCODE

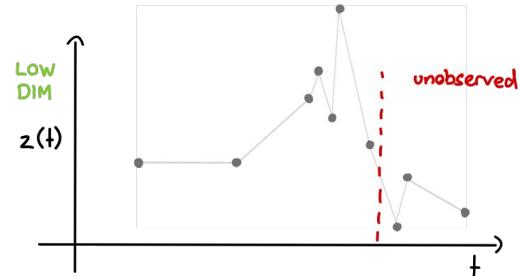
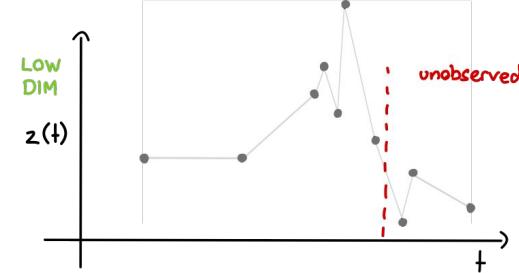


DECODE

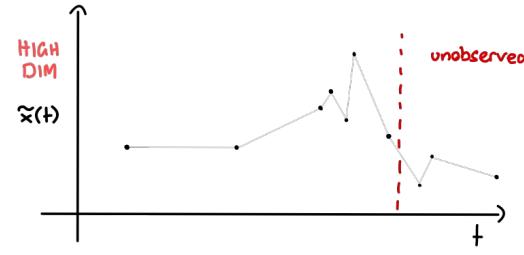




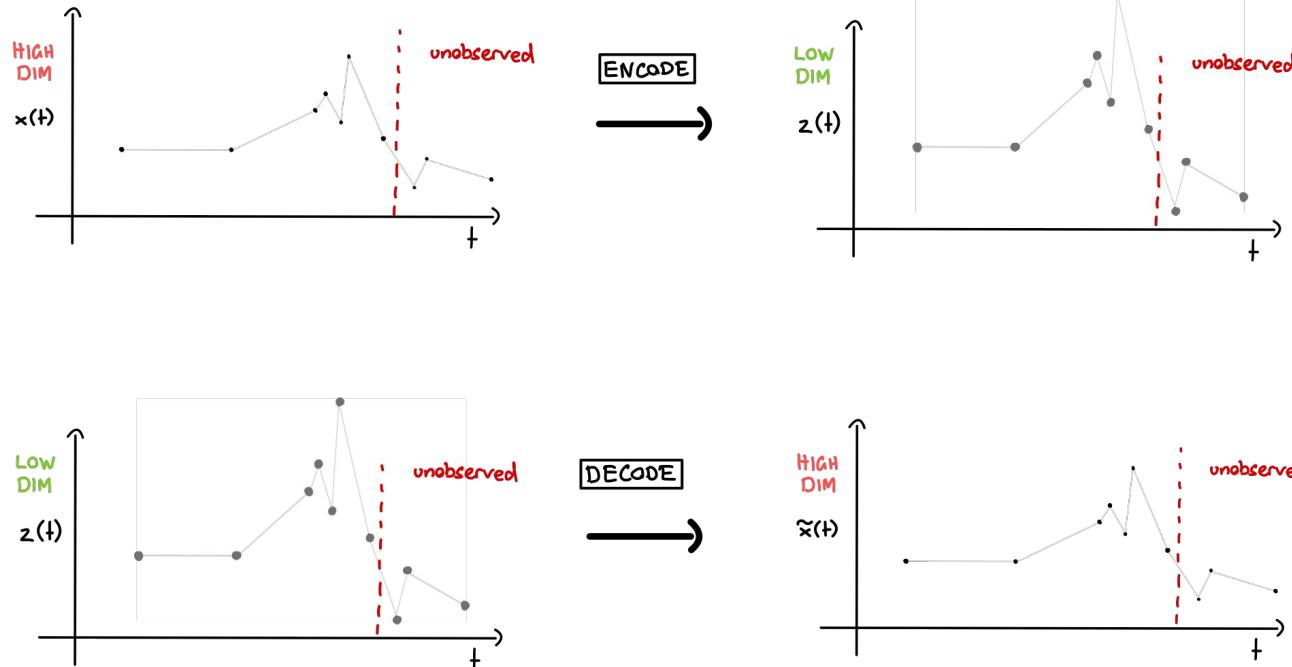
ENCODE
→

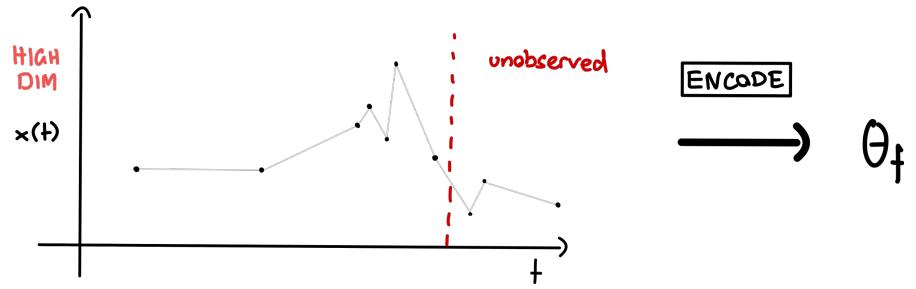


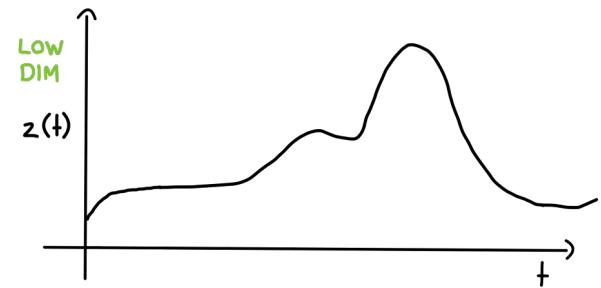
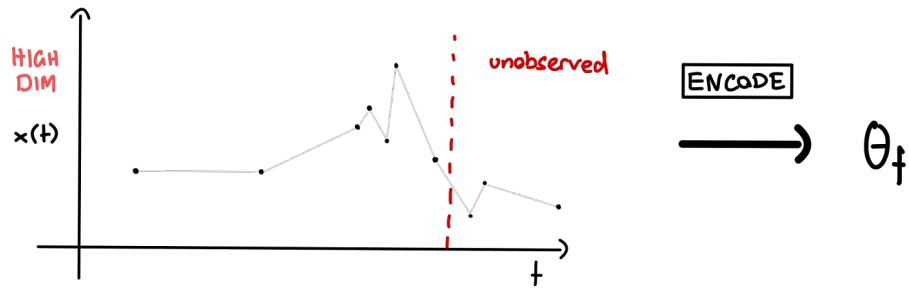
DECODE
→

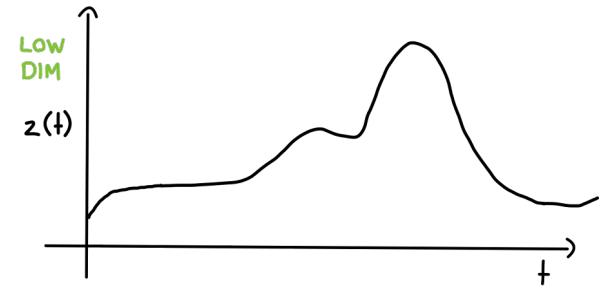
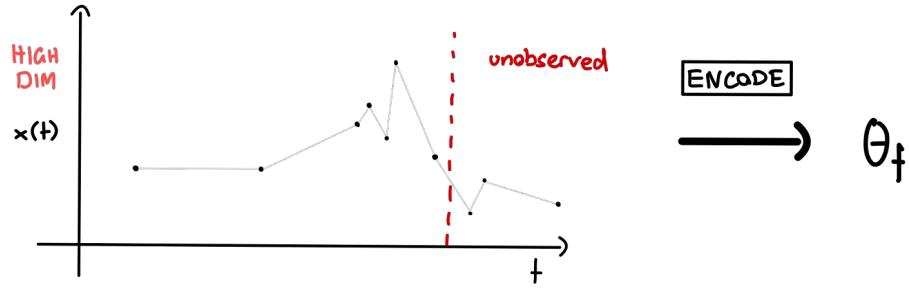


GM: Autoencoders

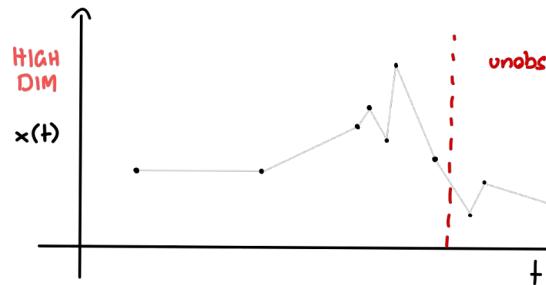






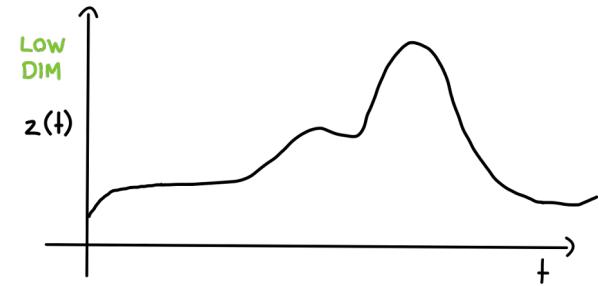


\Rightarrow described by f



ENCODE

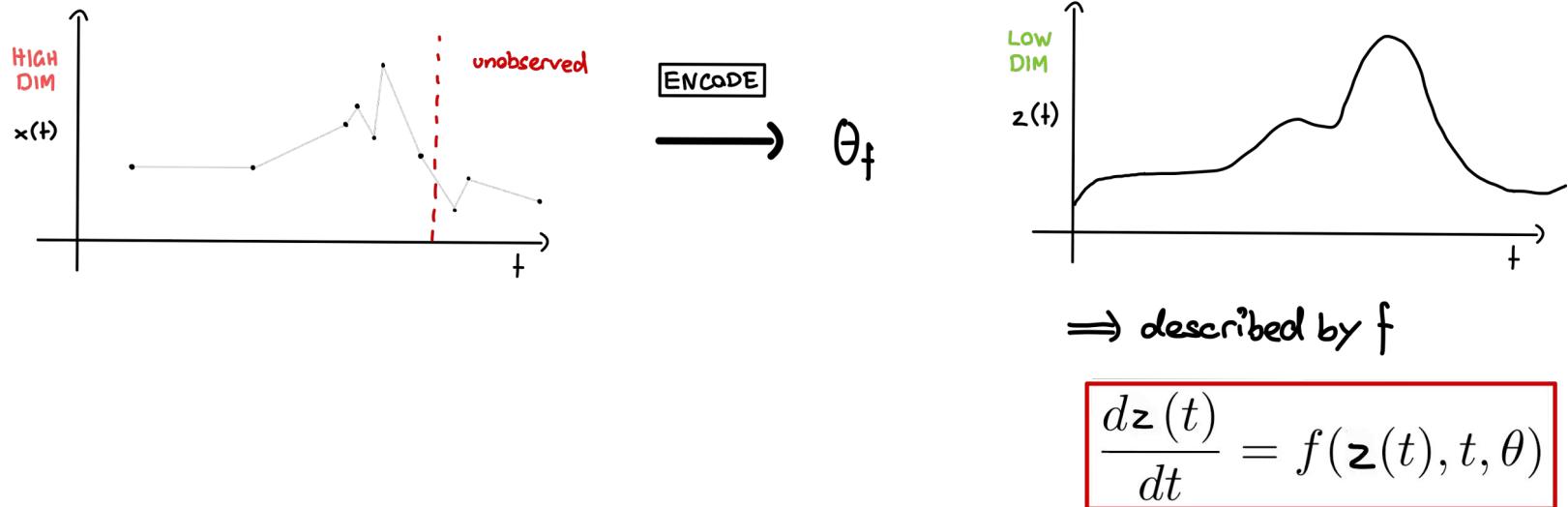
$\longrightarrow \theta_f$



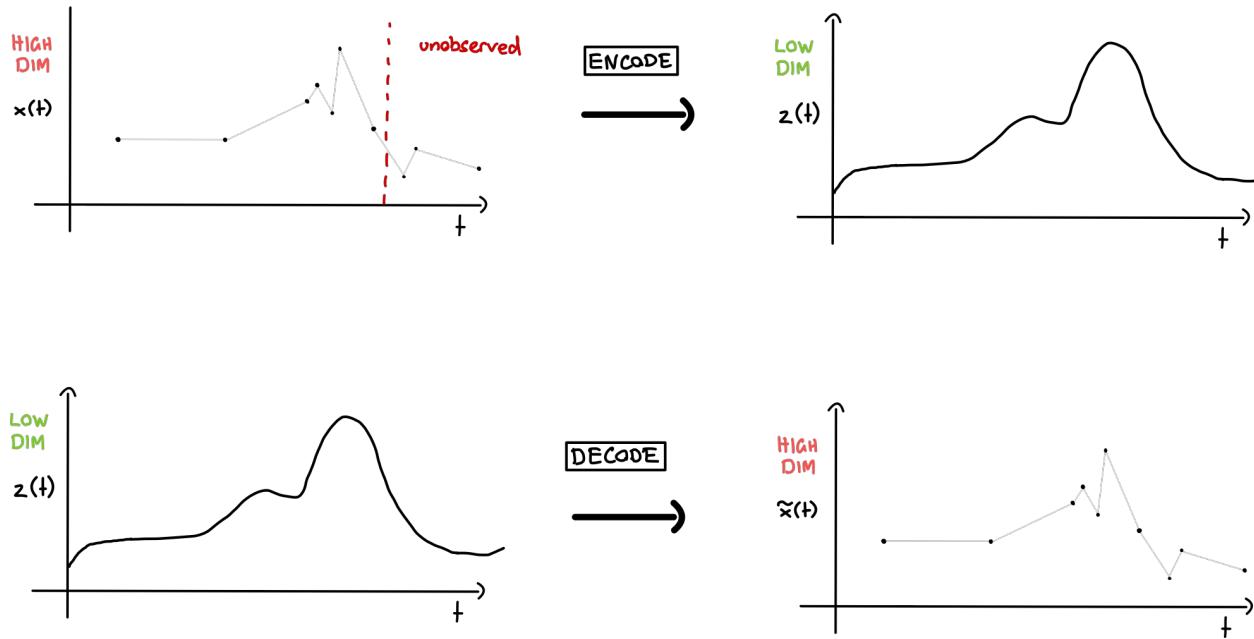
\Rightarrow described by f

$$\boxed{\frac{dz(t)}{dt} = f(z(t), t, \theta)}$$

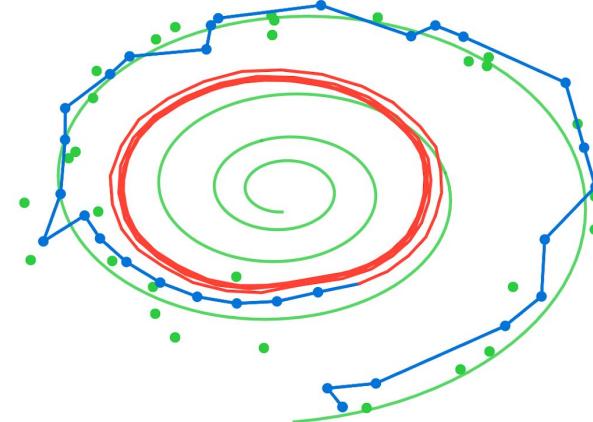
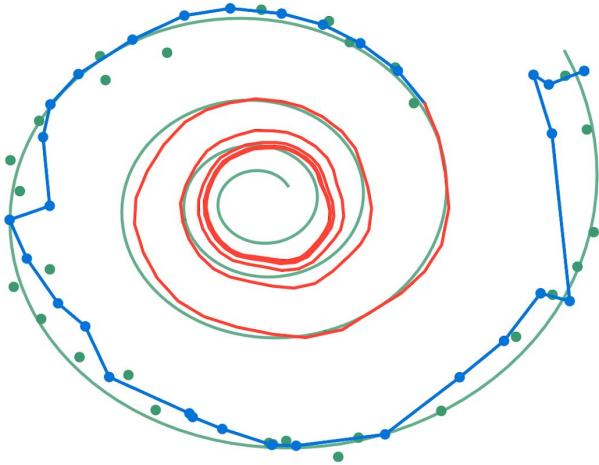
GM: Variational Autoencoders



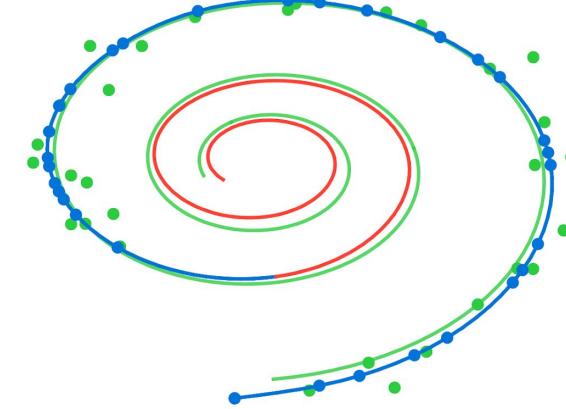
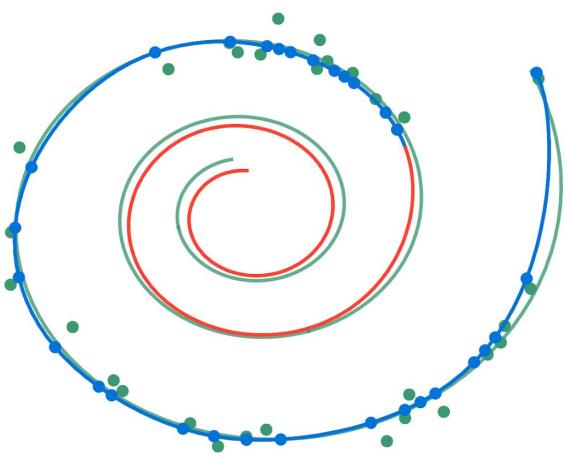
GM: Variational Autoencoders



GM: RNN Results



GM: Results with our Generative Model



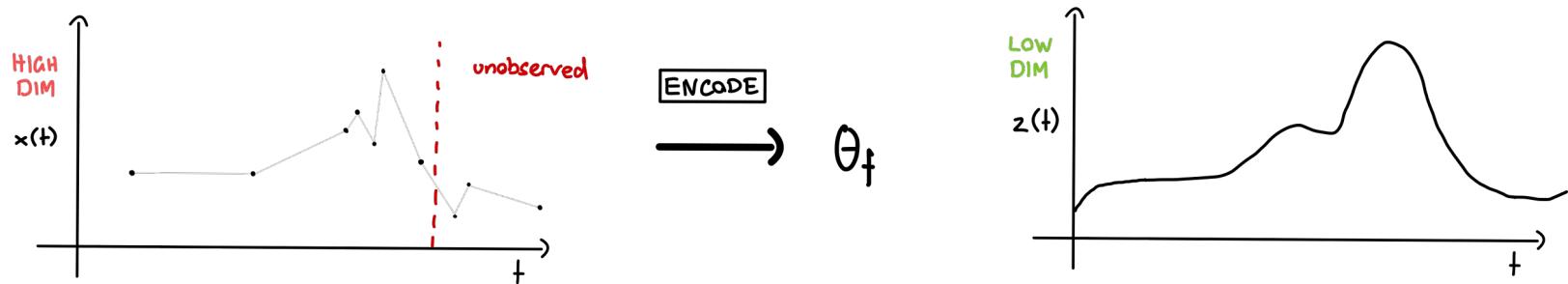
Criticism

Criticism

- Experiments are only performed on toy-ish datasets (MNIST)

Criticism

- Experiments are only performed on toy-ish datasets (MNIST)
- Using an RNN for the Encoding part of the Generative Model



Open Questions: Robustness

ON ROBUSTNESS OF NEURAL ORDINARY DIFFERENTIAL EQUATIONS

Hanshu Yan*, Jiawei Du*, Vincent Y. F. Tan & Jiashi Feng

Department of Electrical and Computer Engineering

National University of Singapore

{hanshu.yan, dujiawei}@u.nus.edu, {vtan, elefjia}@nus.edu.sg

- Shown to be more robust than CNNs
- Still unclear for high-dimensional datasets

EFFICIENT CERTIFIED TRAINING AND ROBUSTNESS VERIFICATION OF NEURAL ODES

Mustafa Zeqiri, Mark Niklas Müller, Marc Fischer & Martin Vechev

Department of Computer Science

ETH Zurich, Switzerland

mzeqiri@ethz.ch, {mark.mueller, mark.fischer, martin.vechev}@inf.ethz.ch

Further Research: Neural ODE Processes

Problem: Neural ODE's don't adapt well to new incoming data without the entire model being retrained

Solution:

NEURAL ODE PROCESSES

Alexander Norcliffe* †

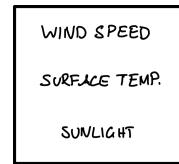
Department of Computer Science
University College London
London, United Kingdom
ucabino@ucl.ac.uk

Cristian Bodnar*, Ben Day*, Jacob Moss* & Pietro Liò

Department of Computer Science
University of Cambridge
Cambridge, United Kingdom
{cb2015, bjd39, jm2311, pl219}@cam.ac.uk

Further Research: Identification and Interventions

They introduced a method that can recover the governing ODEs from data alone



$$\frac{dS}{dt} = Q - \lambda S$$

Beyond Predictions in Neural ODEs: Identification and Interventions

Hananeh Aliee

Helmholtz Center, Munich

Fabian J. Theis

TUM, Helmholtz Center, Munich

Niki Kilbertus

Helmholtz Center, Munich

{hananeh.aliee,fabian.theis,niki.kilbertus}@helmholtz-muenchen.de

ETH zürich

Future Applications

Forecasting virus outbreaks with social media data via neural ordinary differential equations

Matías Núñez, Nadia L. Barreiro, Rafael A. Barrio, Christopher Rackauckas

doi: <https://doi.org/10.1101/2021.01.27.21250642>

July 3, 2020

Forecasting the weather with neural ODEs

#machine learning | #julia | #ode | #time series

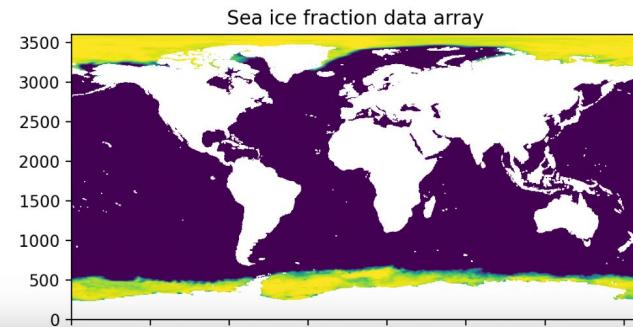
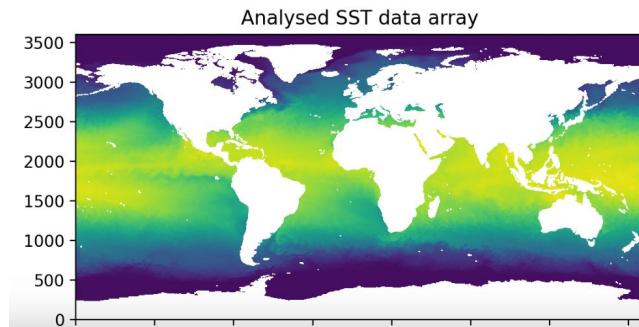
Weather forecasting is a tricky problem. Traditionally, it has been done by manually modelling weather dynamics using differential equations, but this approach is highly dependent on us getting the equations right. To avoid this problem, we can use machine learning to [directly predict the weather](#), which lets us make predictions without modelling the dynamics. However, this approach requires huge amounts of data to reach good performance. Fortunately, there is a middle ground: What if we instead use machine learning to model the *dynamics* of the weather?

Stabilized neural ordinary differential equations for long-time forecasting of dynamical systems

[Alec J. Linot^a](#)  , [Joshua W. Burby^b](#), [Qi Tang^b](#), [Prasanna Balaprakash^c](#), [Michael D. Graham^a](#), [Romit Maulik^c](#)

Future Applications: An Idea

- Use Sea Surface Temperature data to feed into a generative latent time-series model
- ESA dataset



Thank you for your
attention!