



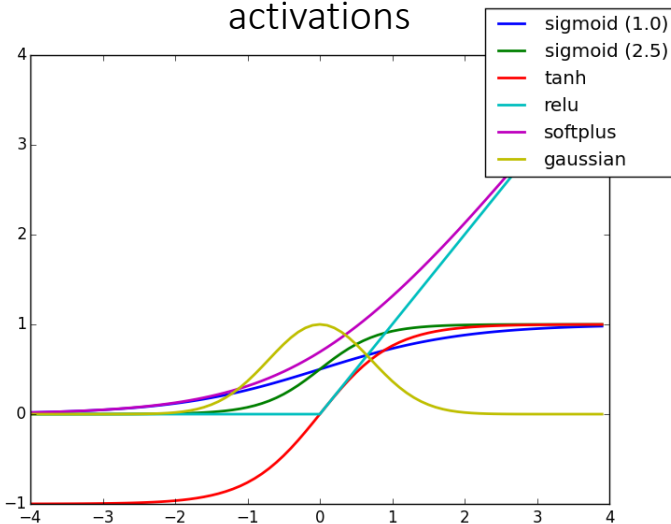
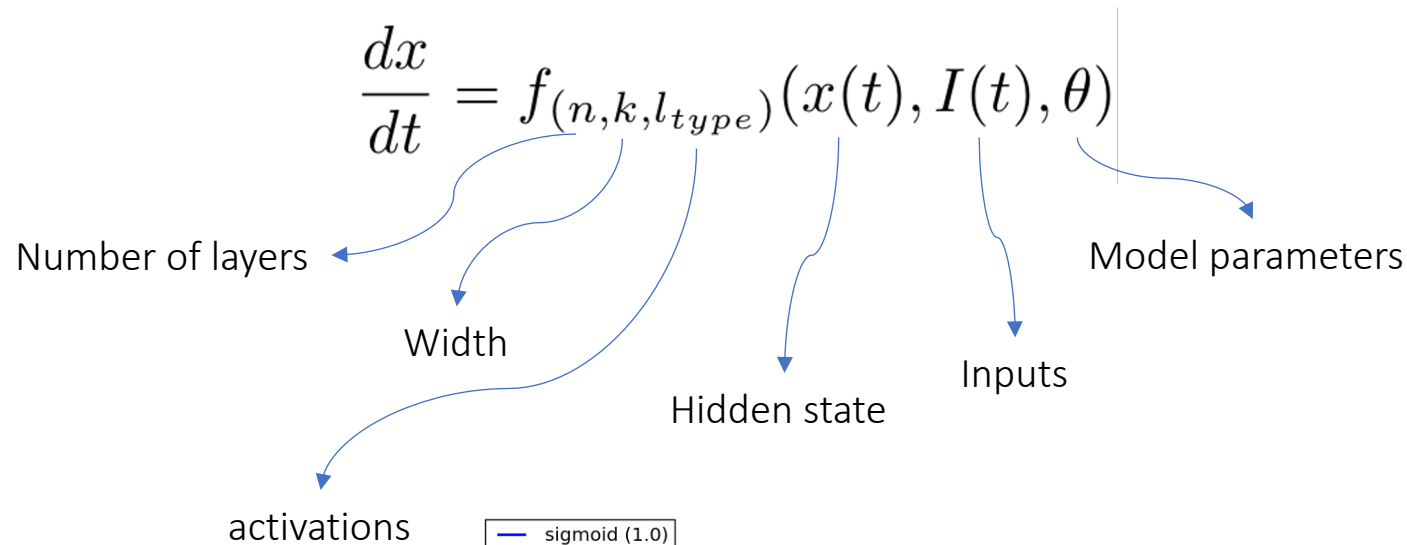
# Liquid Time-Constant Network

Ramin Hasani

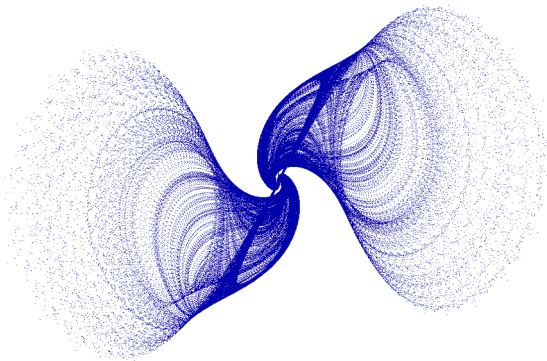
CSAIL, Massachusetts Institute of Technology, USA

Simons SMS Workshop, UC Berkeley  
March 22th, 2021

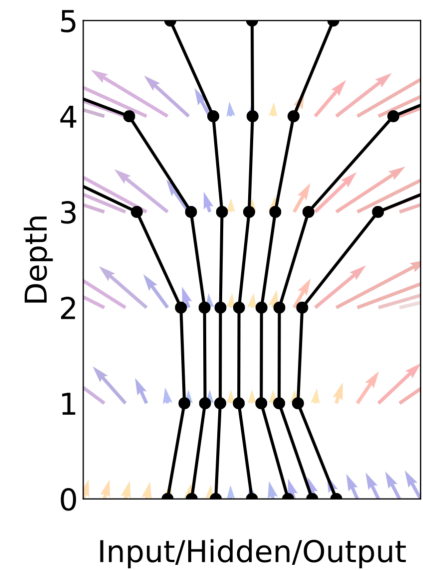
# What is a time-continuous neural network?



Dynamical systems



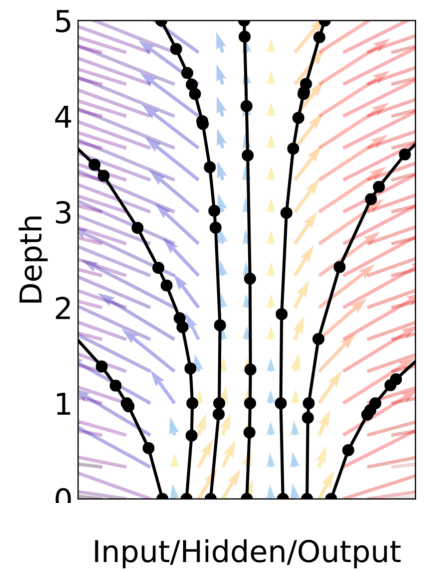
Residual Network



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

He et al.  
CVPR 2016

ODE Network



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

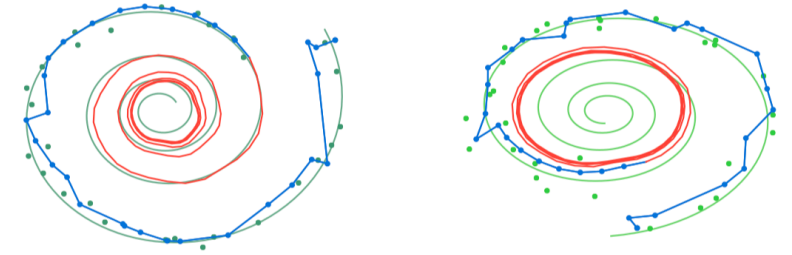
Chen et al.  
NeurIPS 2018

Figure Credit: Chen et al. NeurIPS 2018

# What is a time-continuous neural network?

Standard Recurrent  
Neural Network (RNN)  
Hopfield 1982

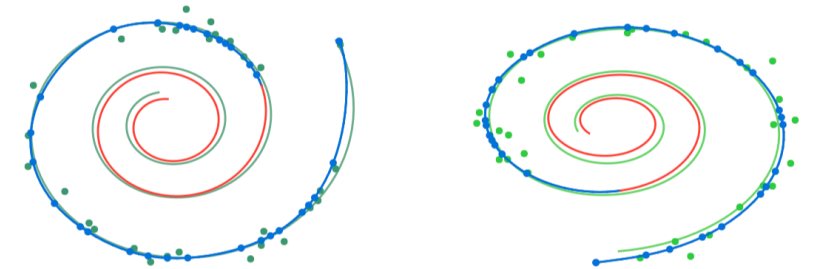
$$x(t + 1) = f(x(t), I(t), t; \theta)$$



(a) Recurrent Neural Network

Neural ODE  
Chen et al. NeurIPS, 2018

$$\frac{dx(t)}{dt} = f(x(t), I(t), t; \theta)$$



(b) Latent Neural Ordinary Differential Equation

Continuous-time  
(CT) RNN  
Funahashi et al. 1993

$$\frac{dx(t)}{dt} = -\frac{x(t)}{\tau} + f(x(t), I(t), t; \theta)$$

— Ground Truth  
● Observation  
— Prediction  
— Extrapolation

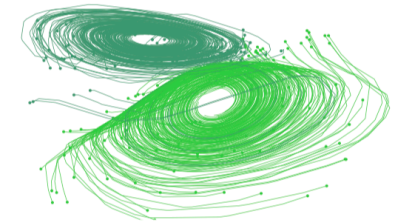


Figure Credit: Chen et al. NeurIPS 2018

# Time-continuous neural networks

How to implement them?

$$d\mathbf{x}(t)/dt = f(\mathbf{x}(t), t, \theta)$$

Numerical ODE solvers

$$\frac{d\mathbf{x}(t)}{dt} \approx \frac{\mathbf{x}(t + \delta t) - \mathbf{x}(t)}{\delta t} \approx f(\mathbf{x}(t), t, \theta)$$

Forward-pass

$$\mathbf{x}(t + \delta t) = \mathbf{x}(t) + \delta t f(\mathbf{x}(t), t, \theta)$$

Choice of the way we do an integration step determines forward pass complexity

# Time-continuous neural networks

## How to train them?

Memory Complexity  $O(1)$  Per layer of  $f$

Adjoint Sensitivity Method  
 [Pontryagin et al. 1962, Chen et al. NeurIPS, 2018]

Loss function  $L(\mathbf{x}(t_1)) = L(\text{ODESolve}(\mathbf{x}(t_0), f, t_0, t_1, \theta))$

Neural ODE  $\frac{d\mathbf{x}(t)}{dt} = f(\mathbf{x}(t), t, \theta)$

Adjoint State  $\mathbf{a}(t) = \frac{\partial L}{\partial \mathbf{x}(t)}$

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

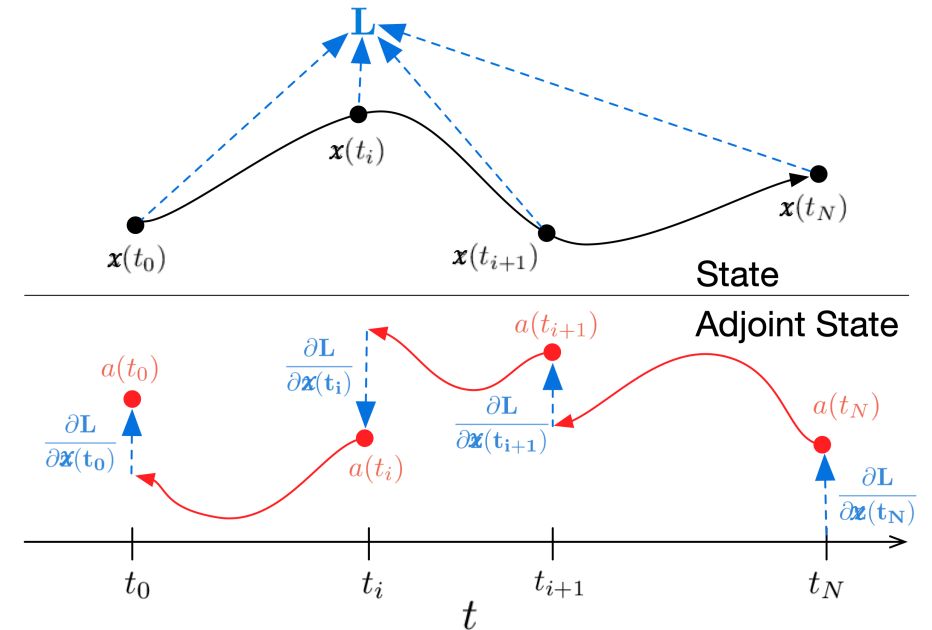


Figure Credit: Chen et al. NeurIPS 2018

# Time-continuous neural networks

## How to train them?

Memory Complexity  $O(L * T)$  Per layer of  $f$

Depth      sequence length

### Backpropagation through-time (BPTT)

[Werbos, 1990, Gholami et. al, 2019, Lechner et al. 2019, Lechner et al. 2020, Hasani et al. 2020]

Perform a forward-pass

$$\mathbf{x}(t + \delta t) = \mathbf{x}(t) + \delta t f(\mathbf{x}(t), t, \theta)$$

Compute gradients through the ODE solver

$$d\Theta = \left[ \frac{dL}{dx(t + \delta t)}, \frac{dx(t + \delta t)}{dx(t)}, \frac{dx(t + \delta t)}{df}, \frac{df}{dx(t)}, \frac{df}{dt}, \frac{df}{d\theta} \right]$$

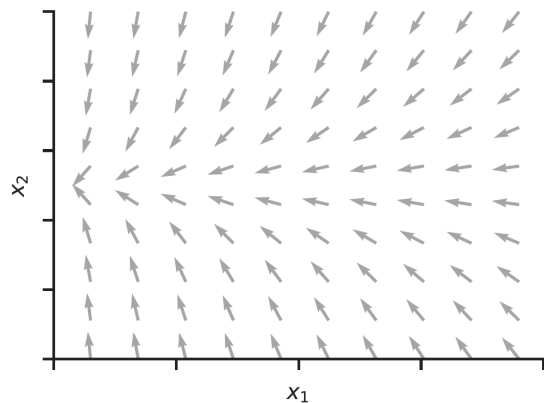
Update parameters

$$\Theta_{new} \leftarrow \Theta_{old} + \gamma d\Theta$$

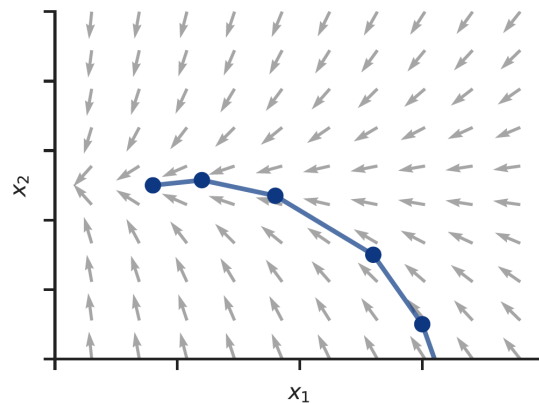
# Time-continuous neural networks

## Better Stay with BPTT

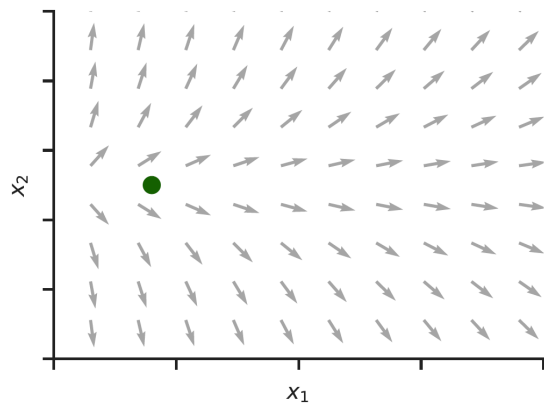
a) Continuous vector field



b) Discrete ODE-solver trajectory



c) Adjoint time-reversed vector field



d) Adjoint ODE-solver trajectory

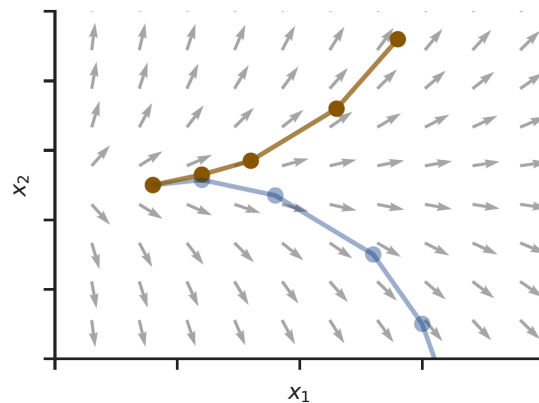


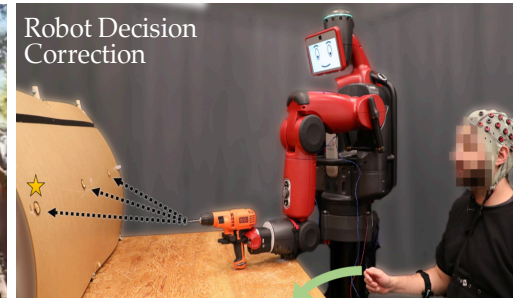
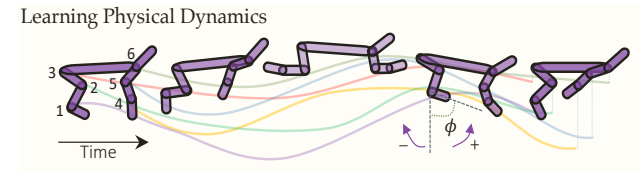
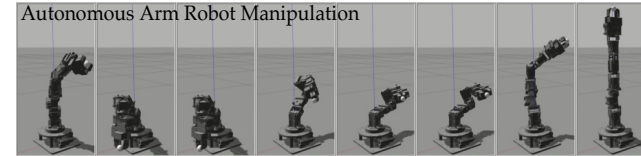
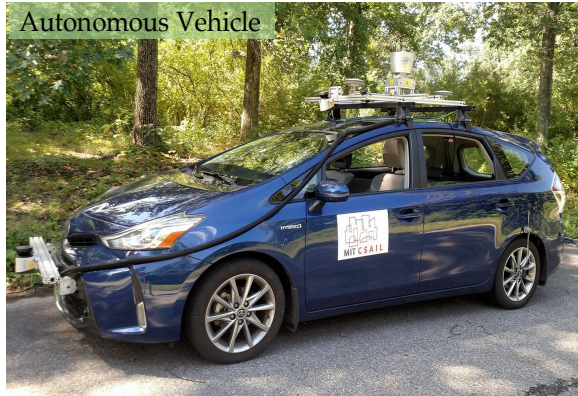
Table 1: Complexity of the vanilla BPTT algorithm compared to the adjoint method, for a single layer neural network  $f$

	<b>Vanilla BPTT</b>	<b>Adjoint</b>
Time	$O(L \times T \times 2)$	$O((L_f + L_b) \times T)$
Memory	$O(L \times T)$	<b>O(1)</b>
Depth	$O(L)$	$O(L_b)$
FWD acc	High	High
BWD acc	<b>High</b>	Low

**Note:**  $L$  = number of discretization steps,  $L_f = L$  during forward-pass.  $L_b = L$  during backward-pass.  $T$  = length of sequence, Depth = computational graph depth.

Neural ODE  
Chen et al. 2018

$$\frac{dx}{dt} = f_{(n,k,l_{type})}(x(t), I(t), \theta)$$



Can Neural ODEs be as expressive as advanced RNNs in real-world applications?



# Liquid Time-Constant Networks

$$d\mathbf{x}(t)/dt = -\mathbf{x}(t)/\tau + \mathbf{S}(t) \quad \mathbf{S}(t) \in \mathbb{R}^M$$

$$\mathbf{S}(t) = f(\mathbf{x}(t), \mathbf{I}(t), t, \theta)(A - \mathbf{x}(t))$$

$$\frac{d\mathbf{x}(t)}{dt} = - \left[ \frac{1}{\tau} + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) \right] \mathbf{x}(t) + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) A$$

“Liquid” = variable

# LTCs have stable state and time-constant

$$\frac{d\mathbf{x}(t)}{dt} = - \left[ \frac{1}{\tau} + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) \right] \mathbf{x}(t) + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) A$$

System time-constant

$$\tau_{sys} = \frac{\tau}{1 + \tau f(\mathbf{x}(t), \mathbf{I}(t), t, \theta)}$$

**Theorem 1.** *Let  $x_i$  denote the state of a neuron  $i$  within an LTC network identified by Eq. 1, and let neuron  $i$  receive  $M$  incoming connections. Then, the time-constant of the neuron,  $\tau_{sys_i}$ , is bounded to the following range:*

$$\tau_i / (1 + \tau_i W_i) \leq \tau_{sys_i} \leq \tau_i, \quad (4)$$

**Theorem 2.** *Let  $x_i$  denote the state of a neuron  $i$  within an LTC, identified by Eq. 1, and let neuron  $i$  receive  $M$  incoming connections. Then, the hidden state of any neuron  $i$ , on a finite interval  $Int \in [0, T]$ , is bounded as follows:*

$$\min(0, A_i^{min}) \leq x_i(t) \leq \max(0, A_i^{max}), \quad (5)$$

# Liquid Time-Constant Networks are Universal Approximators

$$\frac{d\mathbf{x}(t)}{dt} = - \left[ \frac{1}{\tau} + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) \right] \mathbf{x}(t) + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) A$$

**Theorem 3.** *Let  $\mathbf{x} \in \mathbb{R}^n$ ,  $S \subset \mathbb{R}^n$  and  $\dot{\mathbf{x}} = F(\mathbf{x})$  be an autonomous ODE with  $F : S \rightarrow \mathbb{R}^n$  a  $C^1$ -mapping on  $S$ . Let  $D$  denote a compact subset of  $S$  and assume that the simulation of the system is bounded in the interval  $I = [0, T]$ . Then, for a positive  $\epsilon$ , there exist an LTC network with  $N$  hidden units,  $n$  output units, and an output internal state  $\mathbf{u}(t)$ , described by Eq. 1, such that for any rollout  $\{\mathbf{x}(t) | t \in I\}$  of the system with initial value  $x(0) \in D$ , and a proper network initialization,*

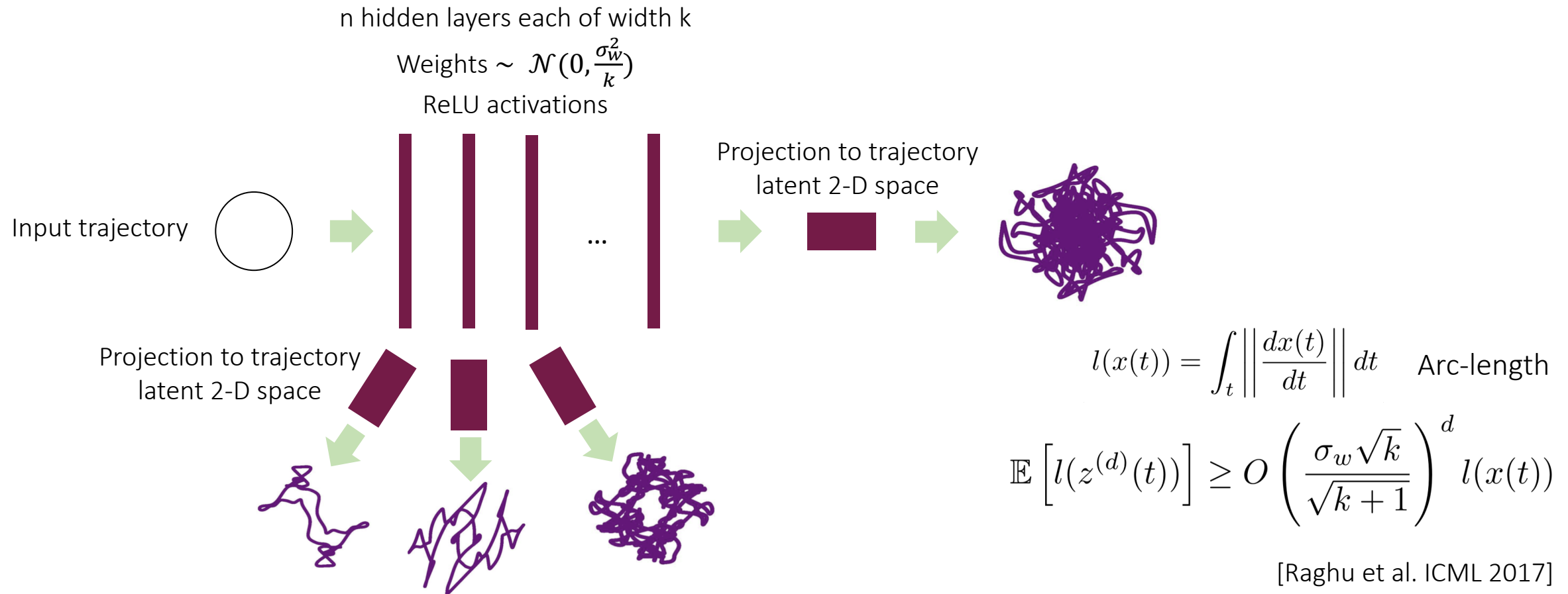
$$\max_{t \in I} |\mathbf{x}(t) - \mathbf{u}(t)| < \epsilon \quad (6)$$

# Expressivity

Defining a better measure

*Trajectory Length* as a measure of expressivity of Deep networks

[Raghu et al. ICML 2017]



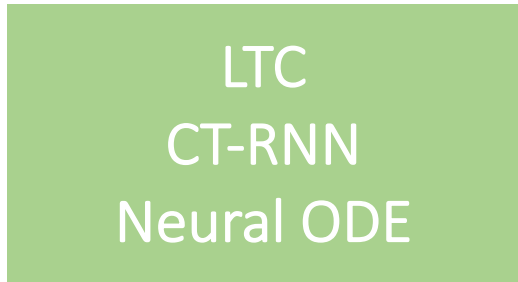
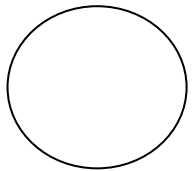
[Raghu et al. ICML 2017]

# Expressivity

Trajectory length as a measure of expressivity

Let's implement the trajectory space for time-continuous models

Input trajectory

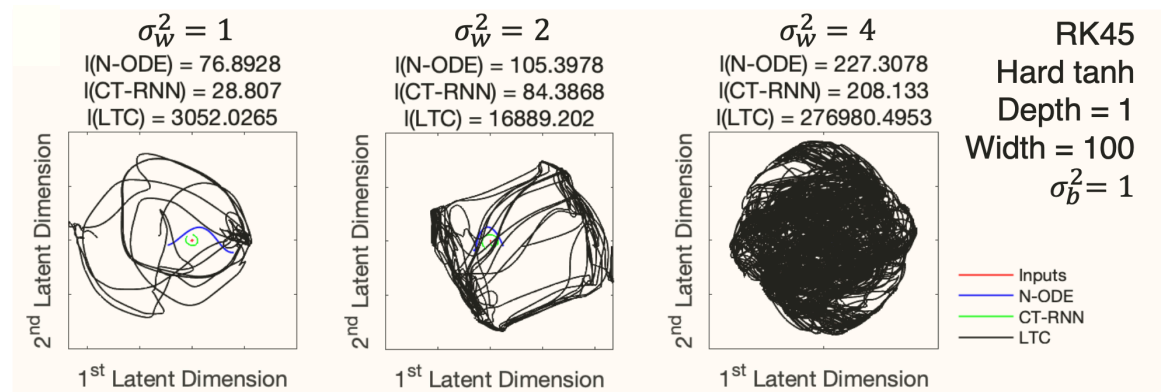
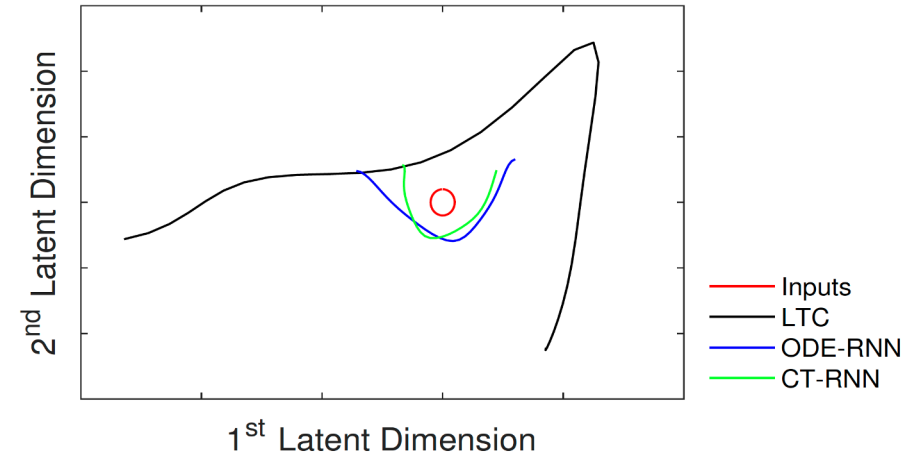


Weights  $\sim \mathcal{N}(0, \frac{\sigma_w^2}{k})$

activation functions:

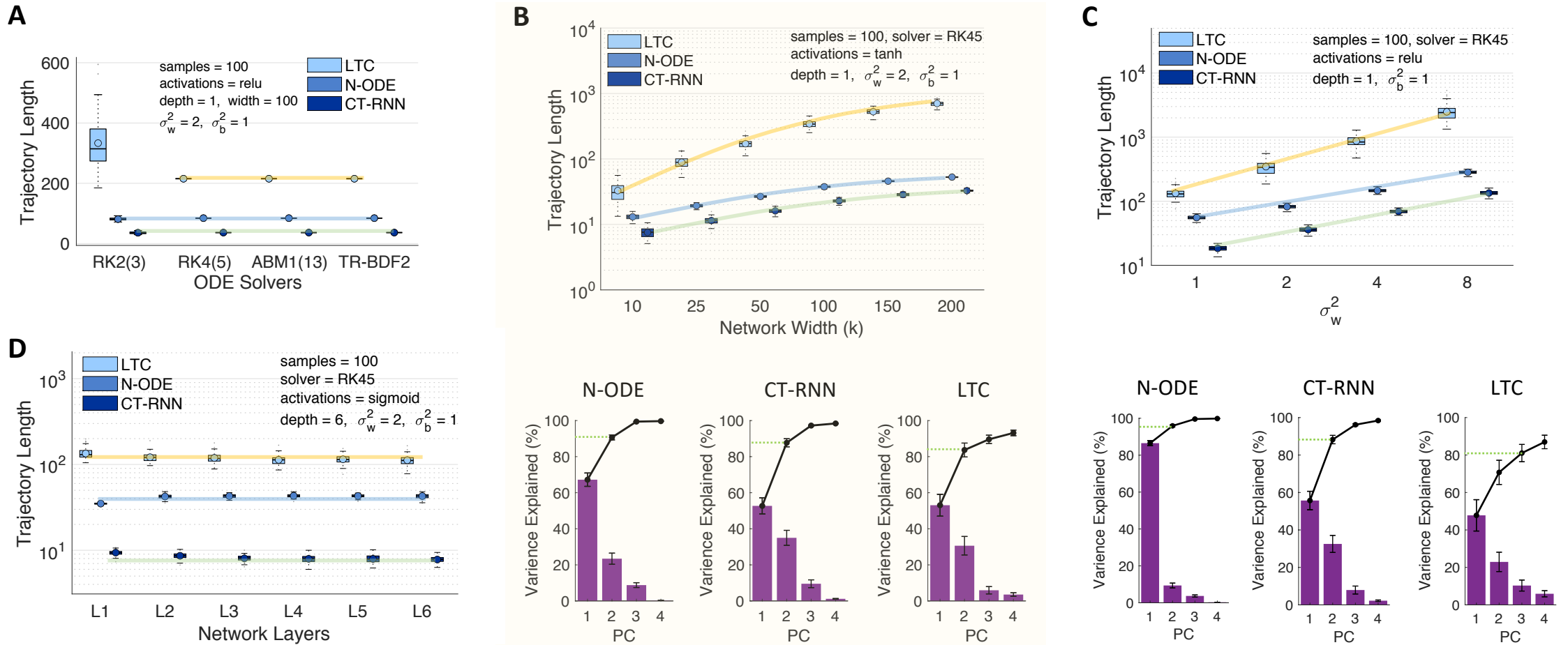
*ReLU, tanh, logistic sigmoid*

Projection to a latent trajectory 2-D space



# Expressivity

Trajectory length as a measure of expressivity



# Expressivity

Trajectory length lower bound

Weight scale      Bias scale      Width      Depth      Number of discretization steps

Neural ODE:  $\mathbb{E} \left[ l(z^{(d)}(t)) \right] \geq O \left( \frac{\sigma_w \sqrt{k}}{\sqrt{\sigma_w^2 + \sigma_b^2 + k \sqrt{\sigma_w^2 + \sigma_b^2}}} \right)^{d \times L} l(I(t))$

CT-RNN:  $\mathbb{E} \left[ l(z^{(d)}(t)) \right] \geq O \left( \frac{(\sigma_w - \sigma_b) \sqrt{k}}{\sqrt{\sigma_w^2 + \sigma_b^2 + k \sqrt{\sigma_w^2 + \sigma_b^2}}} \right)^{d \times L} l(I(t))$

LTC:  $\mathbb{E} \left[ l(z^{(d)}(t)) \right] \geq O \left( \left( \frac{\sigma_w \sqrt{k}}{\sqrt{\sigma_w^2 + \sigma_b^2 + k \sqrt{\sigma_w^2 + \sigma_b^2}}} \right)^{d \times L} \left( \sigma_w + \frac{\|z^{(d)}\|}{\min(\delta t, L)} \right) \right) l(I(t))$

System's dynamic time-scale

# Performance

LTCs in modeling physical dynamics

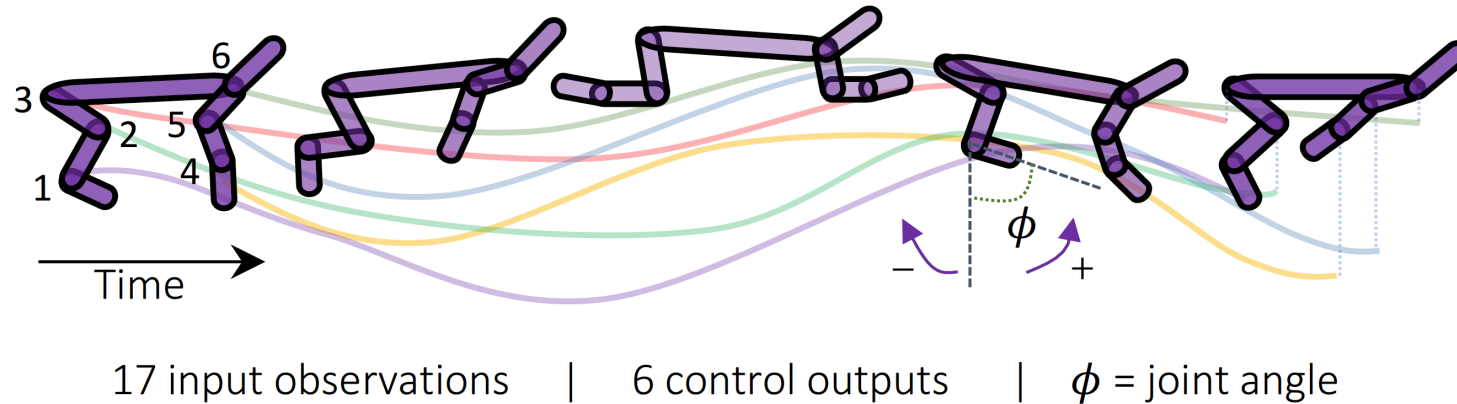


Table 6: Sequence modeling  
Half-Cheetah dynamics  $n=5$

Algorithm	MSE
LSTM	$2.500 \pm 0.140$
CT-RNN	$2.838 \pm 0.112$
Neural ODE	$3.805 \pm 0.313$
CT-GRU	$3.014 \pm 0.134$
LTC (ours)	<b><math>2.308 \pm 0.015</math></b>



# Performance

LTCs in modeling irregularly sampled data

Table 5: Person activity, 2nd setting

Algorithm	Accuracy
RNN $\Delta_t$ * [47]	$0.797 \pm 0.003$
RNN-Decay* [38]	$0.800 \pm 0.010$
RNN GRU-D* [5]	$0.806 \pm 0.007$
RNN-VAE* [47]	$0.343 \pm 0.040$
Latent ODE (D enc.)*	$0.835 \pm 0.010$
ODE-RNN *	$0.829 \pm 0.016$
Latent ODE(C enc.)*	$0.846 \pm 0.013$
-----	-----
LTC (ours)	<b><math>0.882 \pm 0.005</math></b>

**Note:** Accuracy values for algorithms indicated by \*, are taken directly from [47]. RNN  $\Delta_t$  = classic RNN + input delays. RNN-Decay = RNN with exponential decay on the hidden states. GRU-D = gated recurrent unit + exponential decay + input imputation. D-enc. = RNN encoder. C-enc = ODE encoder. n=5

[5] Che et al. Nature Scientific Reports, 2018

[38] Moser et al. Arxiv, 2017

[47] Rubanova et al. NeurIPS 2019

# Performance

LTCs in modeling real-life time series data

**Table 3: Time series prediction Mean and standard deviation, n=5**

Dataset	Metric	LSTM [28]	CT-RNN [47]	Neural ODE [6]	CT-GRU [38]	LTC (ours)
Gesture	(accuracy)	64.57% $\pm$ 0.59	59.01% $\pm$ 1.22	46.97% $\pm$ 3.03	68.31% $\pm$ 1.78	<b>69.55% <math>\pm</math> 1.13</b>
Occupancy	(accuracy)	93.18% $\pm$ 1.66	94.54% $\pm$ 0.54	90.15% $\pm$ 1.71	91.44% $\pm$ 1.67	<b>94.63% <math>\pm</math> 0.17</b>
Activity recognition	(accuracy)	95.85% $\pm$ 0.29	95.73% $\pm$ 0.47	<b>97.26% <math>\pm</math> 0.10</b>	96.16% $\pm$ 0.39	95.67% $\pm$ 0.575
Sequential MNIST	(accuracy)	<b>98.41% <math>\pm</math> 0.12</b>	96.73% $\pm$ 0.19	97.61% $\pm$ 0.14	98.27% $\pm$ 0.14	97.57% $\pm$ 0.18
Traffic	(squared error)	0.169 $\pm$ 0.004	0.224 $\pm$ 0.008	1.512 $\pm$ 0.179	0.389 $\pm$ 0.076	<b>0.099 <math>\pm</math> 0.0095</b>
Power	(squared-error)	0.628 $\pm$ 0.003	0.742 $\pm$ 0.005	1.254 $\pm$ 0.149	<b>0.586 <math>\pm</math> 0.003</b>	0.642 $\pm$ 0.021
Ozone	(F1-score)	0.284 $\pm$ 0.025	0.236 $\pm$ 0.011	0.168 $\pm$ 0.006	0.260 $\pm$ 0.024	<b>0.302 <math>\pm</math> 0.0155</b>

[28] Hochreiter et al. 1997

[6] Chen et al. NeurIPS, 2018

[38] Moser et al. Arxiv, 2017

[47] Rubanova et al. NeurIPS 2019

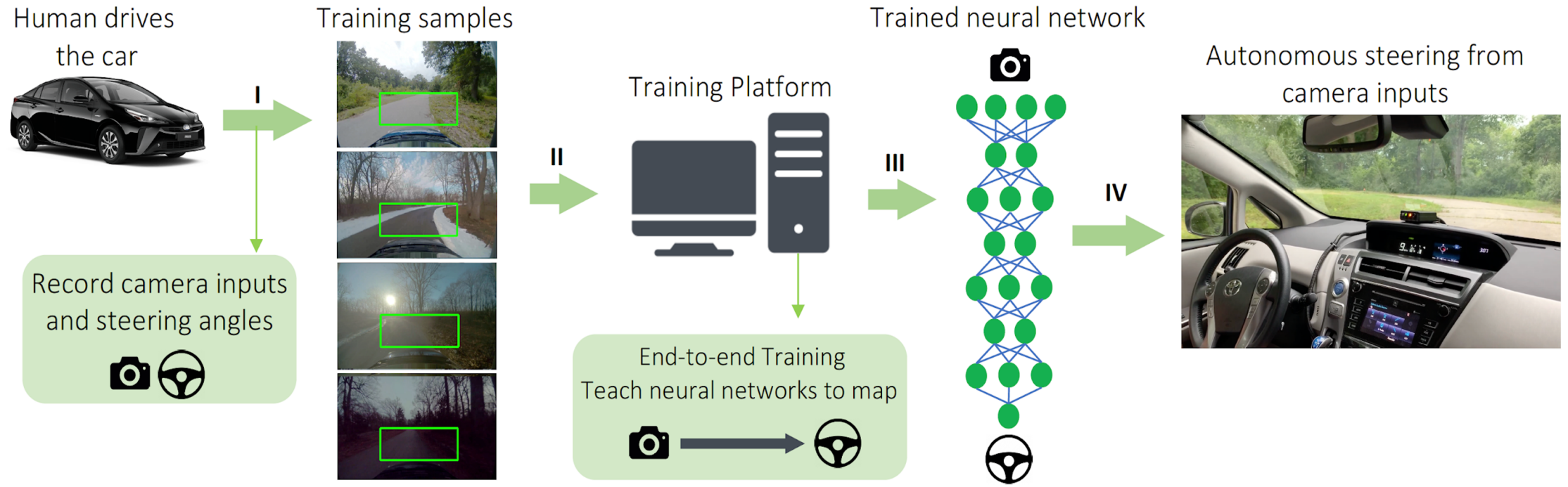
# Summary

- ✓ A novel time-continuous neural networks for efficient time-series modelling
  - LTCs are universal approximators
  - LTCs are stable dynamical systems
  - LTCs show better degrees of expressivity
  - They can vary their behavior even post-training
  - Learning irregularly-sampled data
  - Their effectiveness in modeling continuous-time processes.

What can we do with LTCs in real world?

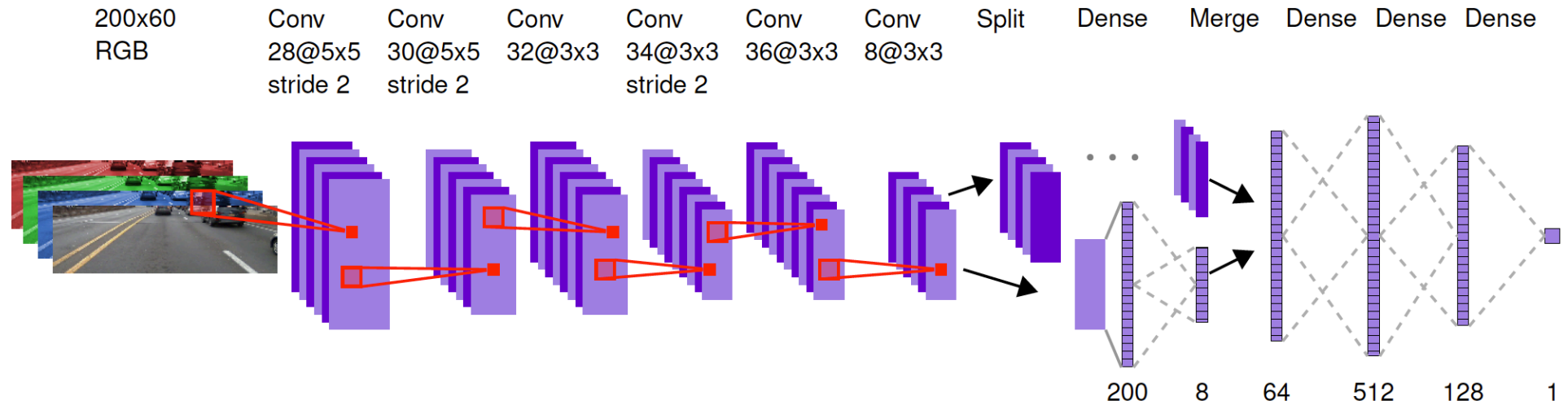
# Performance

High-fidelity autonomy by LTCs  
end-to-end learning

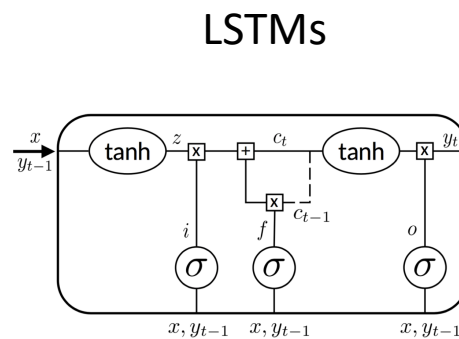
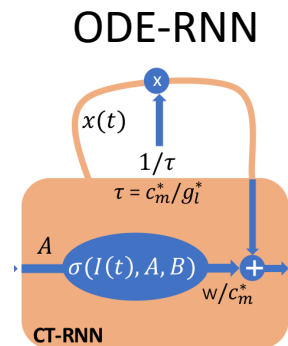


# LTCs: Performance

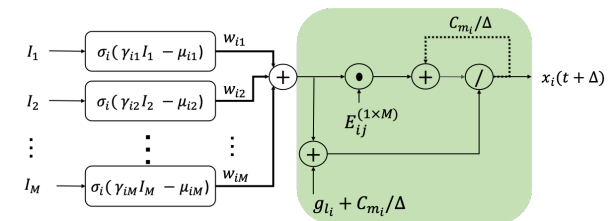
High-fidelity autonomy by LTCs - end-to-end learning



What if we replace the fully connected layers by a recurrent neural network?



LTC-based Networks?

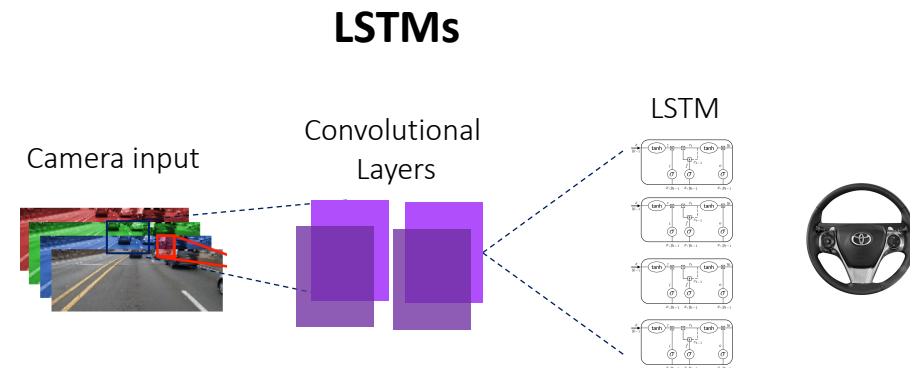
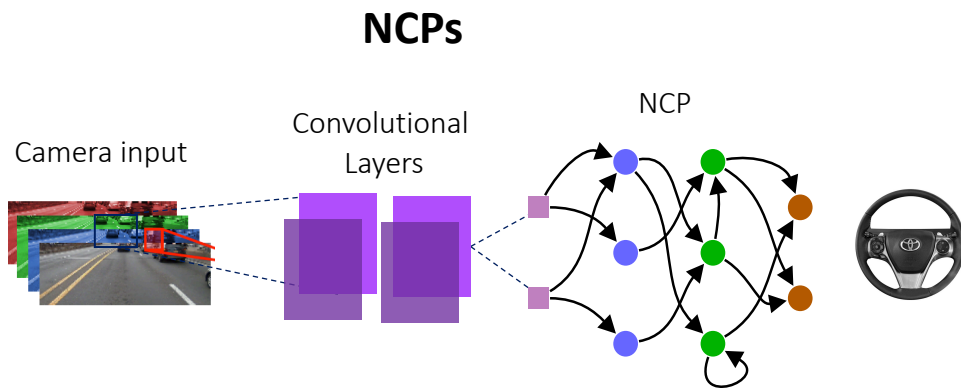


# LTCs: Performance

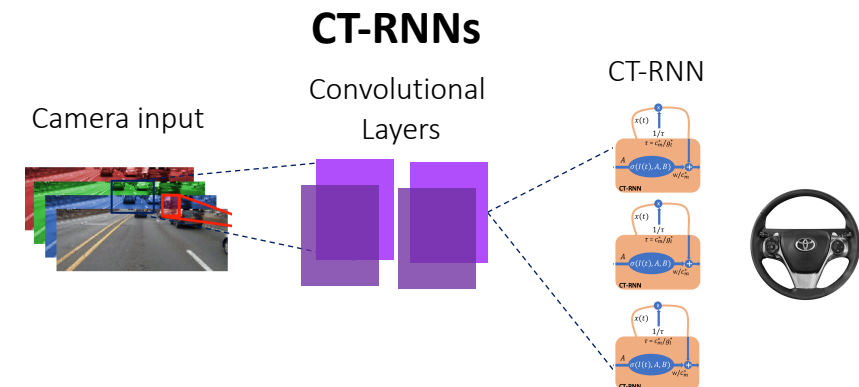
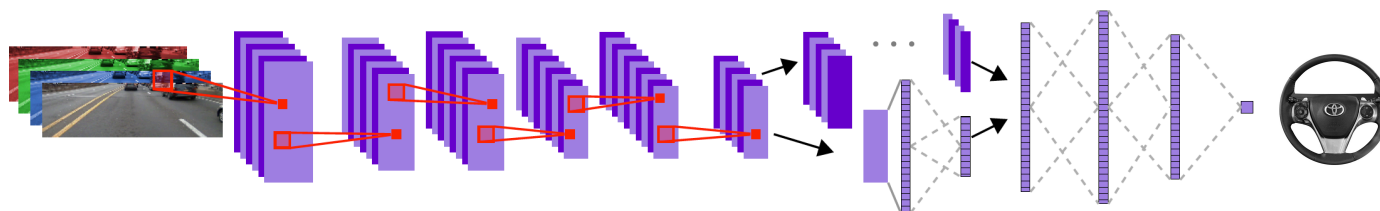
High-fidelity autonomy by LTCs

end-to-end learning of Neural Circuit Policies (NCP)

Now we compare properties of NCPs with a number of other models



## Convolutional Neural Networks (CNNs)



# LTCs: Performance

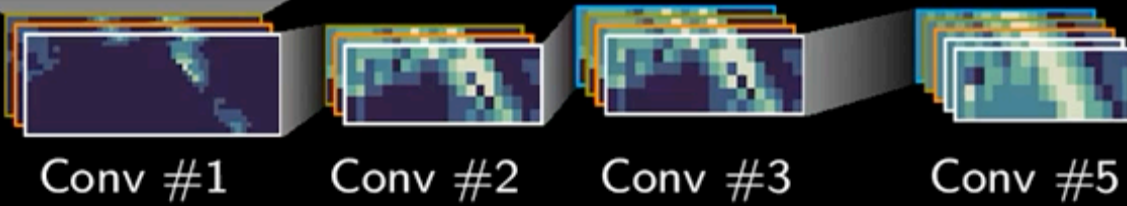
High-fidelity autonomy by LTCs

Parameter efficiency

<b>Model</b>	<b>Conv layers Param</b>	<b>RNN neurons</b>	<b>RNN synapses</b>	<b>RNN trainable param</b>
CNN	5,068,900	-	-	-
CT-RNN	79,420	64	6112	6273
LSTM	79,420	64	24640	24897
<b>NCP</b>	79,420	<b>19</b>	<b>253</b>	<b>1065</b>

# CNN driving performance

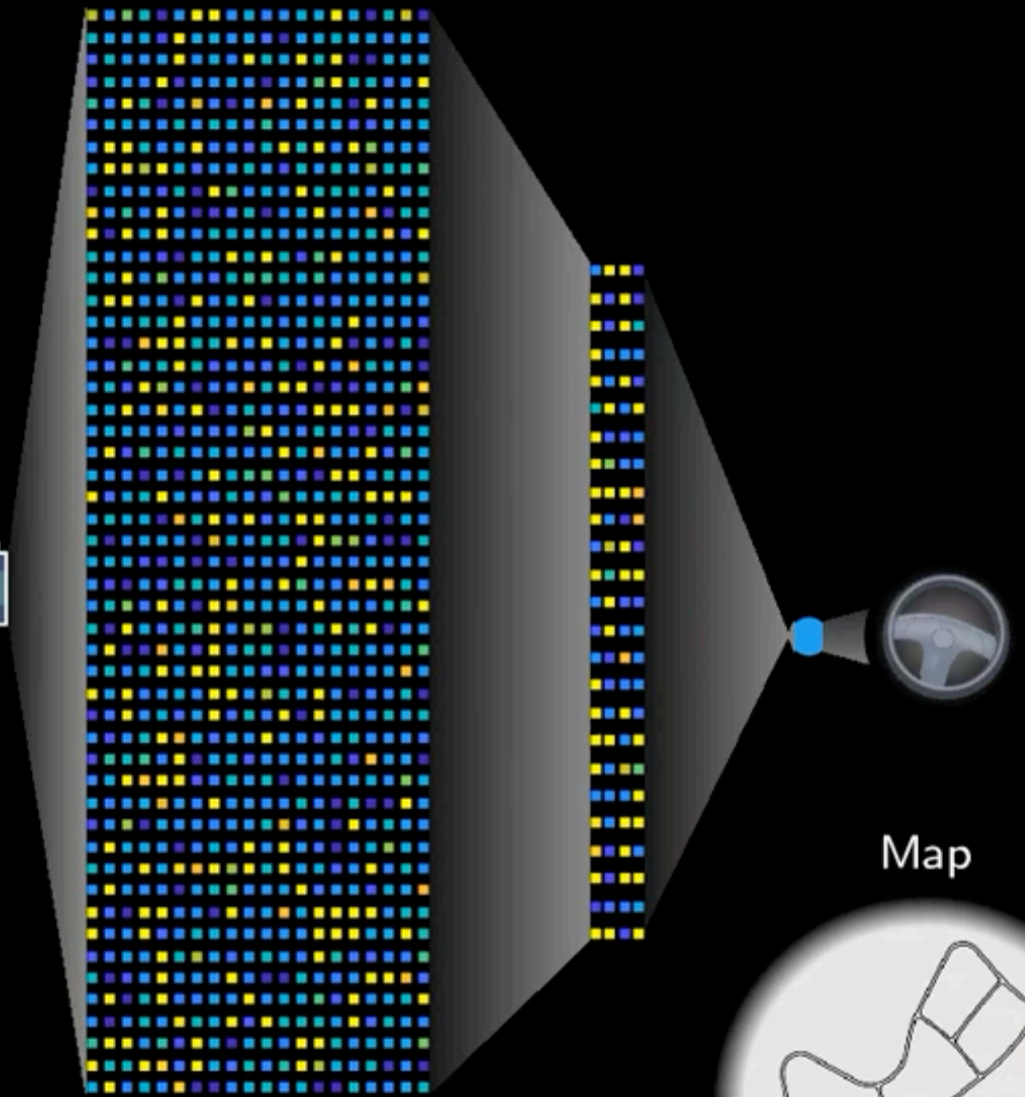
Camera input stream



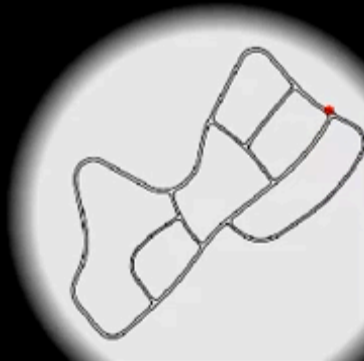
Attention map



Fully-connected layer #1   Fully-connected layer #2   Motor neurons



Map



● Mode: Manual

Normalized neural state



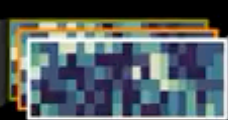


CNN driving performance  
under  $\sigma^2=0.1$  perturbation

Camera input stream



Conv #1



Conv #2

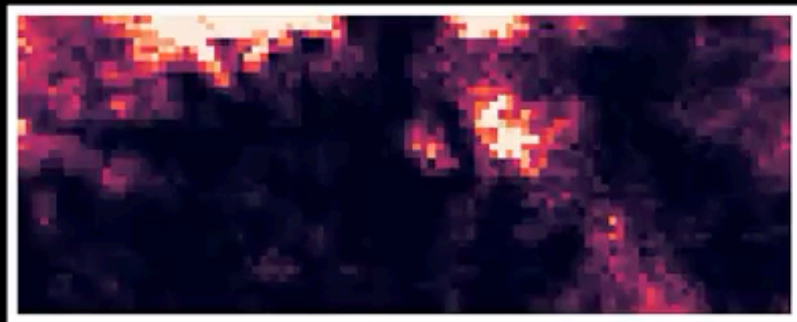


Conv #3



Conv #5

Attention map

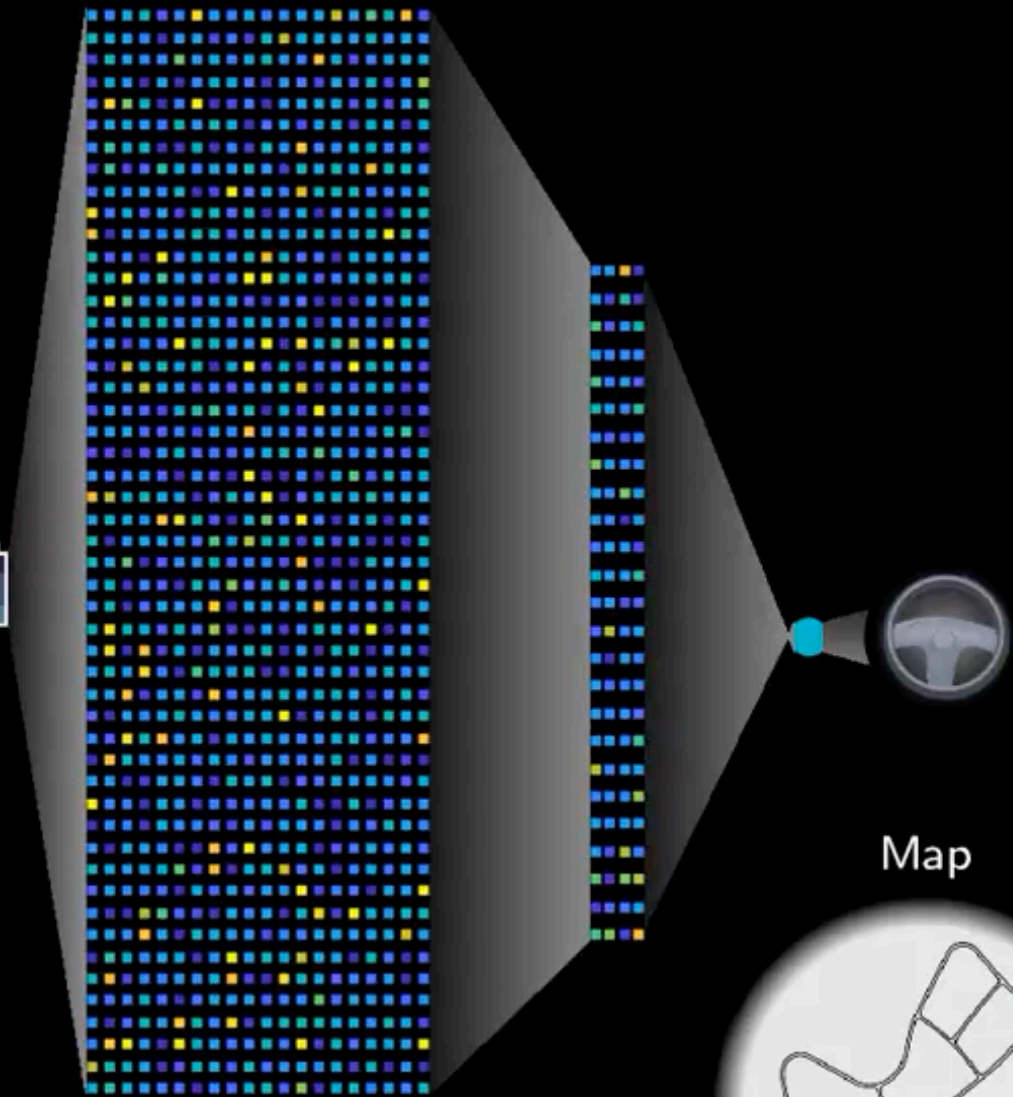


● Mode: Manual

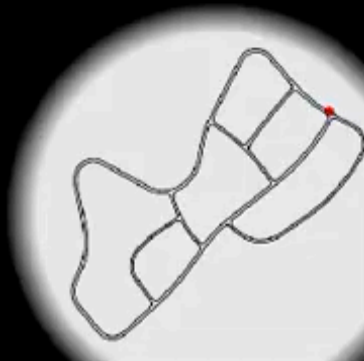
Fully-connected  
layer #1

Fully-connected  
layer #2

Motor  
neurons



Map

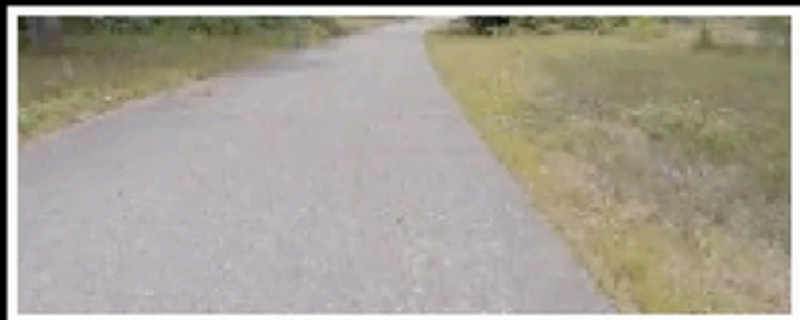


Normalized neural state



# NCP driving performance

Camera input stream



Conv #1



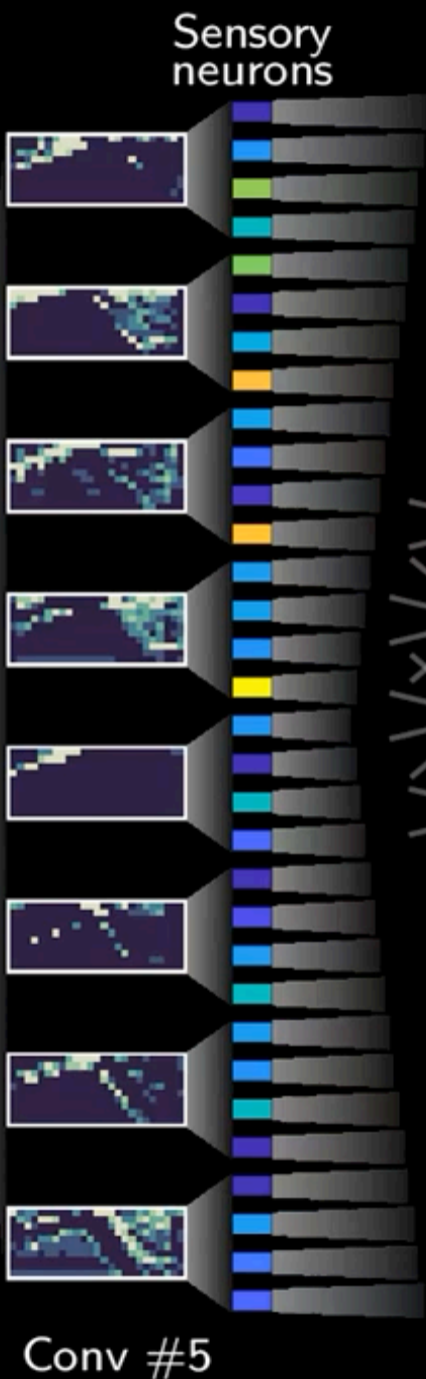
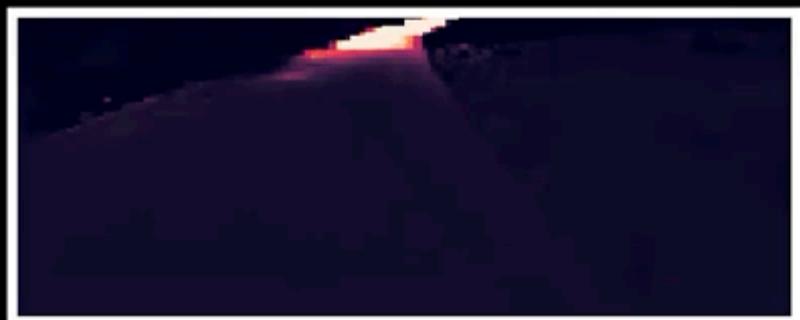
Conv #2



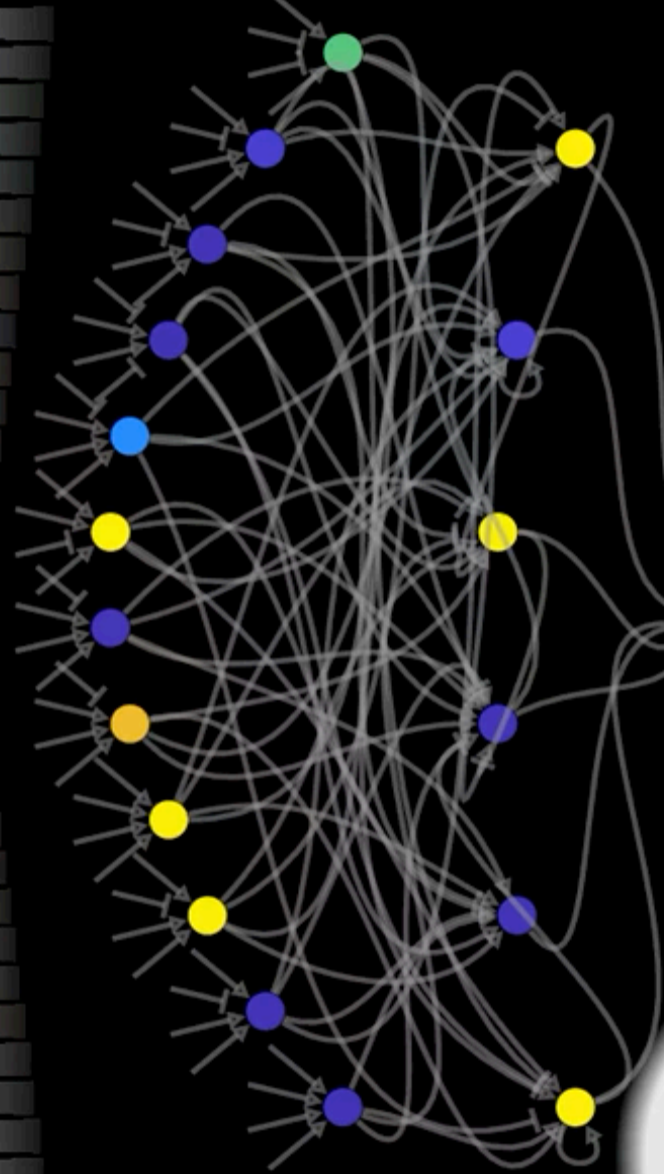
Conv #3



Attention map



Inter neurons



Command neurons

Motor neurons



Map

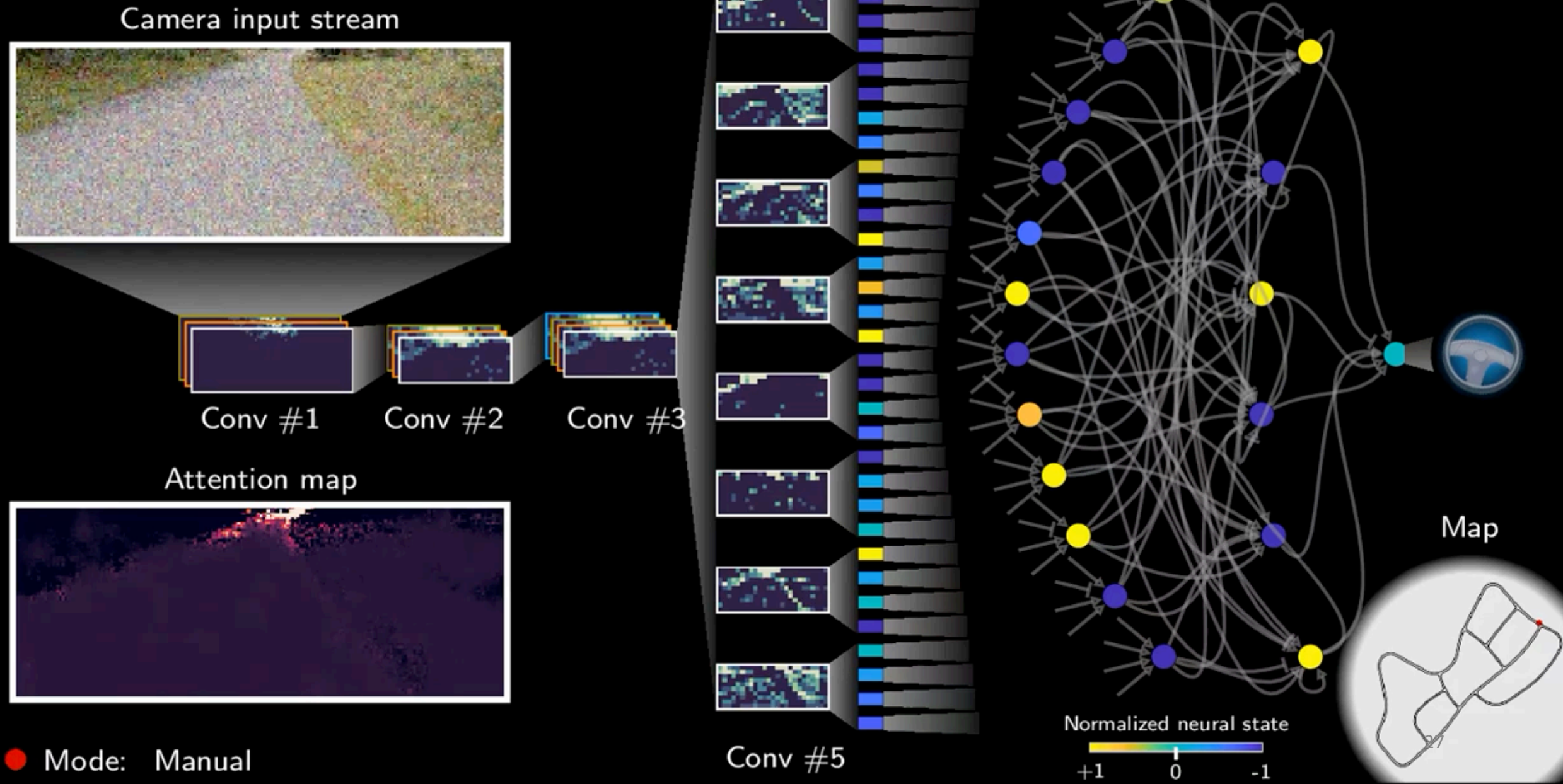


● Mode: Manual

Normalized neural state

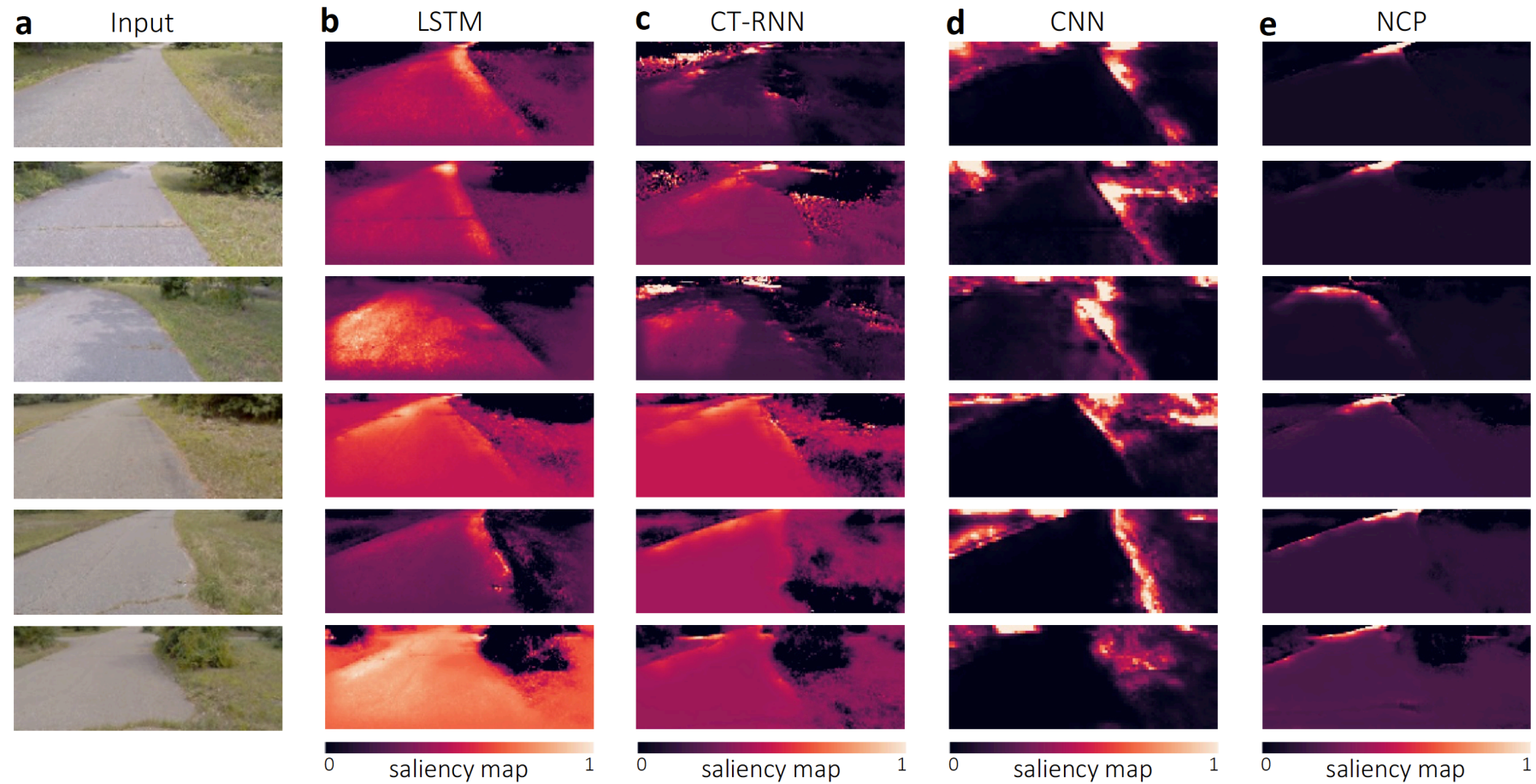


# NCP driving performance under $\sigma^2=0.1$ perturbation



# LTCs: Performance

## High-fidelity autonomy by LTCs

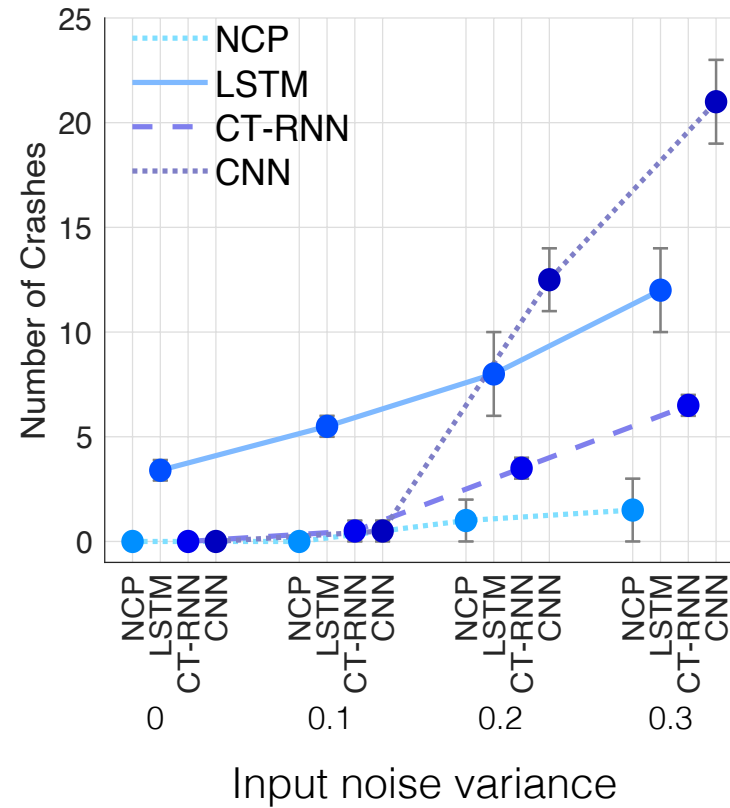


Saliency maps show where each network is learned to attend while driving

# LTCs: Performance

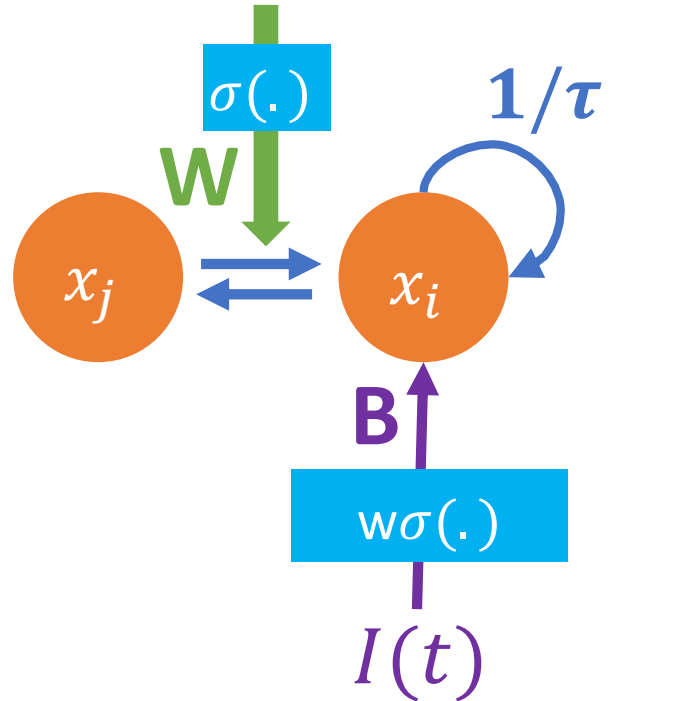
High-fidelity autonomy by LTCs – Robustness

## Noise Robustness (Interventions)



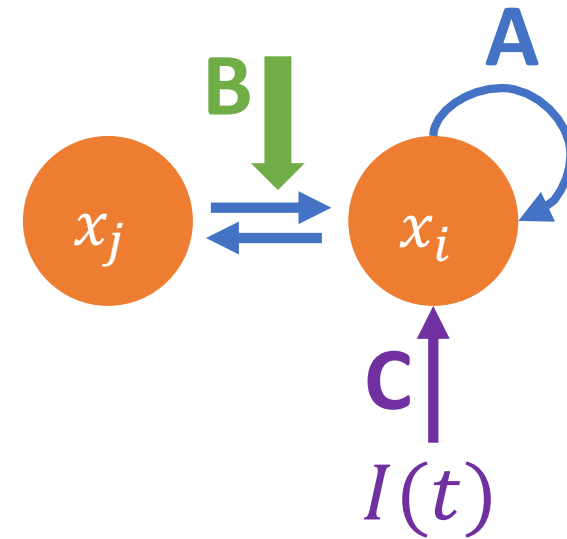
# Neural Circuit Policies are Dynamic Causal Models

Liquid time-constant RNN



- $1/\tau$  intrinsic coupling
- $W\sigma(\cdot)$  liquidity modulator
- $B$  input regulator

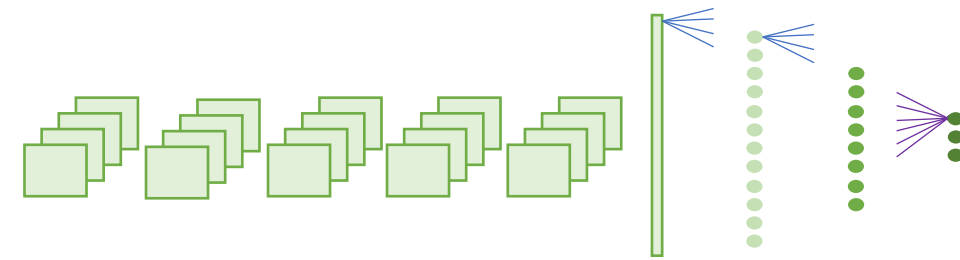
Dynamic causal model



- $A$  intrinsic coupling
- $B$  dynamic modulator
- $C$  input regulator

# Neural Circuit Policies

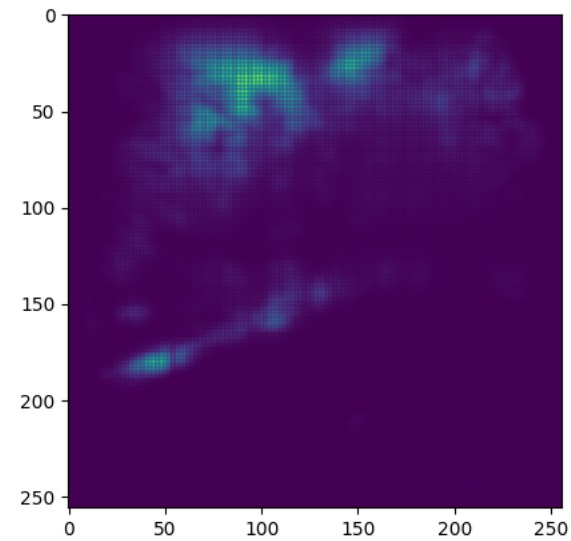
Performance – Attention



Flying Performance

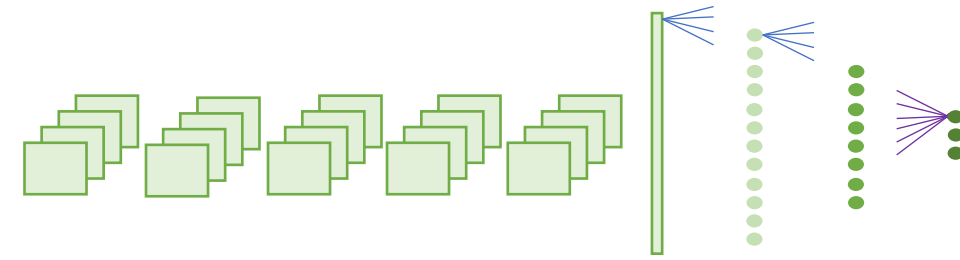


Visual Backprop Attention Map



# Neural Circuit Policies

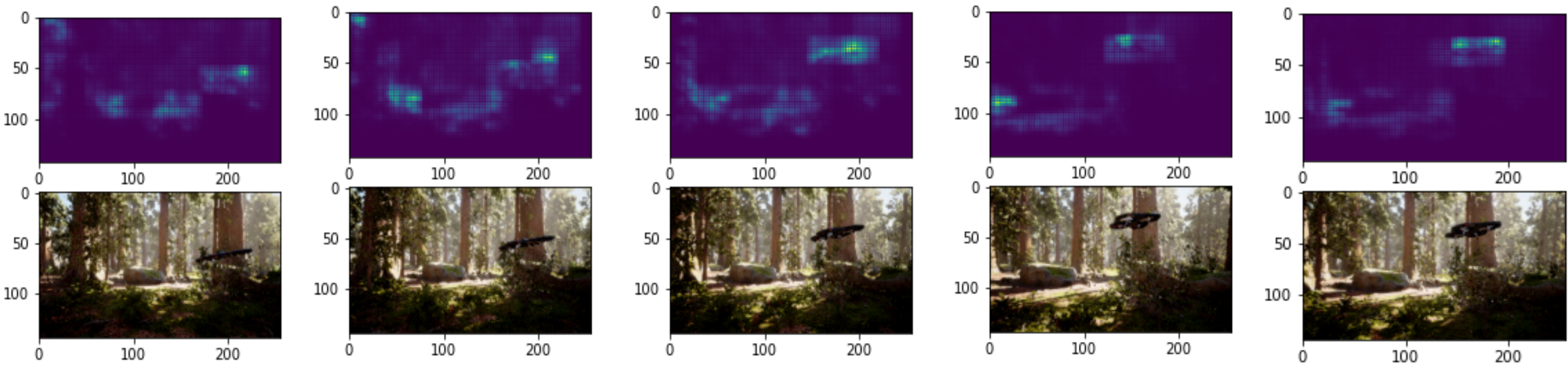
## Leader Following task



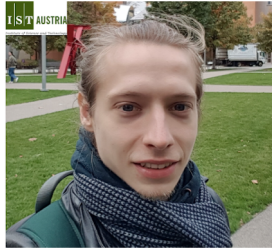


# Neural Circuit Policies

## Leader Following task



Thank you!



Mathias Lechner



Alexander Amini



Daniela Rus



Radu Grosu

Feel free to

Check out our latest repositories

[https://github.com/raminmh/liquid\\_time\\_constant\\_networks](https://github.com/raminmh/liquid_time_constant_networks)



<https://github.com/mlech261/keras-ncp>



<https://github.com/mlech261/ode-lstms>



and to reach out:  
rhasani@mit.edu