# Neural Differential Equations for Learning to Program Neural Nets Through Continuous Learning Rules

Istituto Dalle Molle di studi sull'intelligenza artificiale

IDSIA

USI/SUPSI

[1] The Swiss AI Lab, IDSIA, USI & SUPSI Switzerland
[2] AI Initiative, KAUST Saudi Arabia

Kazuki Irie[1]   Francesco Faccio[1]   Jürgen Schmidhuber[1,2]

## Neural ODEs for Sequence Processing

- Original **Neural ODEs**: continuous-depth version of (feedforward) residual nets
- There are extensions to **process sequences**, e.g.,

**Neural Controlled Differential Equations**  Kidger et al. 2020

Hidden state $h(t) \in \mathbb{R}^d$   Differentiable control signal $x(t) \in \mathbb{R}^{d_{in}}$

$$h(t) = h(t_0) + \int_{s=t_0}^{t} F_\theta(h(s))dx(s)$$

*"vanilla" RNN in the continuous-time domain*

$$= h(t_0) + \int_{s=t_0}^{t} F_\theta(h(s))x'(s)ds$$

- *Good empirical performance (outperform other ODE based sequence processors), but*
- *Scalability limitation:*  $\mathbb{R}^d \to \mathbb{R}^{d \times d_{in}}$

## Fast Weight Programmers (FWPs)

Schmidhuber 1991

- NN that learns to program other NNs by rapidly generating weight changes
- Outer product version: linear Transformer
- General purpose sequence processor

Katharopoulos et al. 2020 etc.

Example: *DeltaNet* Schlag et al. 2021

At each step $n$
- Input $x_n \in \mathbb{R}^{d_{in}}$   $\beta_n, q_n, k_n, v_n = W_{slow}x_n$
- Fast Weight Matrix $W_n = W_{n-1} + \sigma(\beta_n)(v_n - W_{n-1}\phi(k_n)) \otimes \phi(k_n)$
- Output $y_n \in \mathbb{R}^{d_{out}}$   $y_n = W_n\phi(q_n)$

**To hear more about FWPs: Visit our poster on Friday**
*"Memory in Artificial and Real Intelligence" WS*

## You only have 2 min?

*We introduce continuous-time counterparts of Fast Weight Programmers (FWP)/ linear Transformers by combining FWPs with Neural ODEs*

→ *We obtain a new type of **Neural ODE/CDE based sequence processors**, that*
→ *Conceptually **scale better** than existing Neural CDE models*
→ *Empirically **outperform** existing Neural CDE based models*

*We propose **multiple model variations**, depending on*

→ *Smoothness of input control signals, and*
→ *Different learning rule parameterisations (Hebb, Oja, Delta)*

*General Idea:*

Discrete-time Weight Update

$$\beta_n, q_n, k_n, v_n = W_{slow}x_n$$

$$W_n = W_{n-1} + \sigma(\beta_n)(v_n - W_{n-1}\phi(k_n)) \otimes \phi(k_n)$$

Continuous-time Counterpart

$$W(t) = W(t_0) + \int_{s=t_0}^{t} F_\theta(W(s), x(s))ds$$

**Code: github.com/IDSIA/neuraldiffeq-fwp**

## Continuous-Time FWPs

### NCDE FWPs

Differentiable Input Control Signal $x(t) \in \mathbb{R}^{d_{in}}$

Forward pass: ODE solver
Backward pass: Continuous adjoint method

State (Fast Weight Matrix) $$W(t) = W(t_0) + \int_{s=t_0}^{t} F_\theta(W(s), x(s), x'(s))x'(s)ds$$
$$\in \mathbb{R}^{d_{out} \times d_{key}}$$

Output
$$y(T) = \begin{cases} W(T)^\top W_q x(T) & \text{Hebb and Oja} \\ W(T)W_q x'(T) & \text{Delta} \end{cases}$$
$$\in \mathbb{R}^{d_{out}}$$

Query

Key   Value   Outer Product

$$= \sigma(\beta(s)) \begin{cases} W_k x(s) \otimes W_v x'(s) & \text{Hebb} \\ (W_k x(s) - W(s)^\top W_v x'(s)) \otimes W_v x'(s) & \text{Oja} \\ (W_v x(s) - W(s)W_k x'(s)) \otimes W_k x'(s) & \text{Delta} \end{cases}$$

Learning rate

### Key properties

- **Scalable**: outer product-based vector field
- **Expressive**: Sum all rank-1 updates in the continuous-time domain, then use the resulting weight matrix to compute the output (i.e., sum before matrix multiplication)

**vs. basic NCDEs with a rank-1 vector field: scalable but not expressive**
- Good empirical performance (Transformer!)
(rank-1 mat. used in isolation)

### Direct NODE FWPs

(piece-wise) Continuous/Bounded Input Control Signal $x(t) \in \mathbb{R}^{d_{in}}$

- Similar idea to (left) but w/o derivative of control signal
- Theoretically NCDEs are more powerful (Kidger et al. 2020), but
- With our parameterisations, performance gap is small

$$W(t) = W(t_0) + \int_{s=t_0}^{t} F_\theta(W(s), x(s))ds$$

$$q(T) = W_q x(T) \text{ Query}$$

Output $$y(T) = W(T)q(T)$$

$$F_\theta(W(s), x(s)) = \sigma(\beta(s)) \begin{cases} k(s) \otimes v(s) & \text{Hebb-style} \\ v(s) \otimes (k(s) - W(s)^\top v(s)) & \text{Oja-style} \\ (v(s) - W(s)k(s)) \otimes k(s) & \text{Delta-style} \end{cases}$$

Learning rate   Key   Value
$$[\beta(s), k(s), v(s)] = W_{slow}x(s)  \quad  W_{slow} \in \mathbb{R}^{(1+d_{key}+d_{out}) \times d_{in}}$$

## Time Series Classification Tasks

### Speech Commands & PhysioNet Sepsis

| Type | Model | Speech Commands | PhysioNet Sepsis OI | PhysioNet Sepsis no-OI |
|---|---|---|---|---|
| Direct NODE | GRU-ODE | 47.9 (2.9) | 85.2 (1.0) | 77.1 (2.4) |
| | Hebb | 82.8 (1.1) | **90.4 (0.4)** | 82.9 (0.7) |
| | Oja | **85.4 (0.9)** | 88.9 (1.4) | 82.9 (0.5) |
| | Delta | 81.5 (3.8) | 89.8 (1.0) | **84.5 (2.9)** |
| CDE | NCDE | 89.8 (2.5) | 88.0 (0.6) | 77.6 (0.9) |
| | Hebb | 89.5 (0.3) | 89.9 (0.6) | **85.7 (0.3)** |
| | Oja | 90.0 (0.7) | **91.2 (0.4)** | 85.1 (2.5) |
| | Delta | **90.2 (0.2)** | 90.9 (0.2) | 84.5 (0.7) |

- *FWPs **outperform** the existing ODE/CDE baselines*
- *No clear winner among different learning rules*

### EigenWorms

- long sequences (> 4000 timesteps)
- The delta rule outperform others

| Model | Sig-Depth | Step | Test Acc. [%] |
|---|---|---|---|
| NRDE | 2 | 4 | 83.8 (3.0) |
| Hebb | 2 | 4 | 45.6 (5.9) |
| Oja | | | 46.7 (7.5) |
| Delta | | | **87.7 (1.9)** |
| NCDE | 1 | 4 | 66.7 (11.8) |
| Hebb | 1 | 4 | 41.0 (6.5) |
| Oja | | | 49.7 (9.9) |
| Delta | | | **91.8 (3.4)** |

### Overall:
- *FWPs outperform the best existing Neural ODE/CDE based models, but*
- *There exist **discrete-time models** that perform equally well or better, e.g., LEM for Eigenworms GRU-D for PhysioNet Sepsis no-OI*

## Model-based Reinforcement Learning

*What if we need to directly work with irregularly sampled discrete inputs? FWP analogs to Latent ODE-RNNs*

$$u_n = \text{ODESolve}(f_{\theta_1}, h_{n-1}, t_{n-1}, t_n)$$
$$h_n, W_n = \text{FWP}([x_n, u_n], W_{n-1}; \theta_2)$$

- *Setting: **MuJoCo** with irregularly timed observations (semi-MDP; repeated actions)*
- *FWPs perform as well or better than Latent ODE-RNNs*



(a) Swimmer



(b) Hopper