



MENG INDIVIDUAL PROJECT

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

---

# Verification Methods for Liquid Neural Networks

---

*Author:*  
Viyan Raj

*Supervisor:*  
Dr. Alessio Lomuscio

*Second Marker:*  
Alyssa Renata

June 12, 2025

# Abstract

1 page

# Acknowledgements

1 page

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
<b>2</b>	<b>Background and Literature Review</b>	<b>7</b>
2.1	Liquid Neural Networks	7
2.1.1	Continuous-Time Dynamical Systems/Differential Equations	7
2.1.2	LNN Training	8
2.1.3	LNN Inference	9
2.1.4	Advantages of LNNs	9
2.1.4.a	Temporal Modelling	9
2.1.4.b	Adaptability	9
2.1.4.c	Efficiency	9
2.1.4.d	Stability	10
2.2	Neural Network Verification	10
2.2.1	The Verification Problem	10
2.2.2	Motivation	11
2.2.3	Recurrent Neural Networks	11
2.2.4	Robustness Verification for Recurrent Neural Networks	12
2.2.5	Symbolic and Interval Propagation Methods	12
2.2.5.a	SIP (Symbolic Interval Propagation)	12
2.2.5.b	CROWN (Certified Robustness to Weight Perturbations)	13
2.2.5.c	Lipschitz-Based Methods	14
<b>3</b>	<b>Liquid Neural Network Design and Implementation</b>	<b>16</b>
3.1	Design Overview	16
3.2	Wiring and Connectivity	16
3.3	LTC Neuron Dynamics	17
3.4	Network Architecture	19
3.5	Training Configuration (and dataset)	20
3.6	Training Behaviour	21
<b>4</b>	<b>Comparative Models</b>	<b>24</b>
4.1	Introduction to Baseline Models	24
4.2	Temporal Convolutional Network (TCN)	24
4.3	Long Short-Term Memory Network (LSTM)	27
4.4	Transformer Model	29
<b>5</b>	<b>Adversarial Attack Methodology</b>	<b>32</b>
5.1	Introduction to Adversarial Attacks	32
5.2	Fast Gradient Sign Method (FGSM)	32
5.3	Projected Gradient Descent (PGD)	33
5.4	DeepFool-Inspired Directional Attack	34
5.5	Simultaneous Perturbation Stochastic Approximation (SPSA)	34
5.6	Time-Warping Attack	35
5.7	Continuous-Time Perturbation Attack	36
5.8	Summary of Attack Design and Implementation Decisions	37
5.8.1	Attack Categories and Coverage	37

5.8.2	Implementation Consistency . . . . .	37
5.8.3	Design Considerations . . . . .	38
<b>6</b>	<b>Bound Certification (Auto Lirpa)</b>	<b>39</b>
<b>7</b>	<b>Evaluation</b>	<b>40</b>
7.1	Quantitative Evaluation Metrics and Comparison . . . . .	40
7.2	Aggregate Results . . . . .	41
7.3	Attack-Specific Breakdowns . . . . .	41
7.4	Qualitative Evaluation and Visual Analysis . . . . .	43
7.5	Comparative Discussion of Model Robustness . . . . .	44
7.5.1	Summary of Behaviour Under Attack . . . . .	44
<b>8</b>	<b>Defences and Mitigation Strategies</b>	<b>46</b>
8.1	Introduction . . . . .	46
8.2	Adversarial Training and Noise Injection . . . . .	46
8.2.1	Gradient-Based Adversarial Training . . . . .	46
8.2.2	Noise Injection During Training . . . . .	47
8.2.3	Benefits and Trade-offs . . . . .	47
8.2.4	Considerations . . . . .	47
8.3	Architectural Enhancements for Robustness . . . . .	47
8.3.1	Memory Mechanisms and Temporal Smoothing . . . . .	47
8.3.2	Receptive Field and Feature Redundancy . . . . .	47
8.3.3	Stability in Continuous-Time Models . . . . .	48
8.3.4	Summary . . . . .	48
8.4	Input Preprocessing and Temporal Defences . . . . .	48
8.4.1	Low-Pass Filtering . . . . .	48
8.4.2	Interpolation and Resampling Defences . . . . .	48
8.4.3	Temporal Quantisation . . . . .	48
8.4.4	Trade-offs and Limitations . . . . .	49
8.4.5	Summary . . . . .	49
8.5	Model-Specific Mitigation Insights . . . . .	49
8.5.1	LSTM: Sequential Memory Vulnerabilities . . . . .	49
8.5.2	TCN: Localised Perturbation Sensitivity . . . . .	49
8.5.3	LNN: ODE Sensitivity and Stability Management . . . . .	49
8.5.4	Summary Table . . . . .	50
8.6	Limitations and Future Work . . . . .	50
8.6.1	Limitations . . . . .	50
8.6.2	Future Work . . . . .	50
8.6.3	Conclusion . . . . .	51
<b>9</b>	<b>Conclusion</b>	<b>52</b>
<b>10</b>	<b>Declaration</b>	<b>54</b>

# Chapter 1

## Introduction

Around 1-2 pages

[\[1\]](#)

## Chapter 2

# Background and Literature Review

Around 15-20 pages

In this chapter, we explore the current research on this topic. First, by exploring the theory behind liquid neural networks, in comparison to deep neural networks. Then, by investigating current verification methods and assessing their suitability to liquid neural networks. The reader should have an understanding of (traditional) neural networks and linear algebra concepts.

### 2.1 Liquid Neural Networks

Liquid neural networks, introduced by Ramin Hasani et al. (2021) [2], are a novel class of AI algorithms, designed to maintain adaptability after completing training. These are inspired by the communication patterns of brain cells, which are flexible and responsive to new/unseen data even after their initial training phase.

Traditional neural networks use fixed architectures and static parameters, so require retraining to handle new information. Liquid neural networks use **continuous-time dynamics** to enable their state to evolve smoothly over time. This means they can dynamically adjust their responses to changing inputs during inference. This structure allows these networks to be robust against perturbations and capable of generating complex behaviors without requiring large-scale architectures.

LNNs use differential equations to simulate the continuous/dynamic processing and plasticity of the brain. Since LNN neurons communicate selectively with a subset of other neurons, connections formed are sparse (unlike traditional deep neural networks with dense fully-connected layers). This makes LNNs more computationally efficient than deep neural networks.

There are a range of applications of LNNs. Hasani suggests that their inherent adaptability makes them suitable for tasks requiring real-time learning and decision-making, such as autonomous driving and medical diagnosis. Their efficiency could address several challenges associated with large-scale machine learning systems, including issues related to interpretability, accountability, and environmental impact due to high carbon footprints. [3]

#### 2.1.1 Continuous-Time Dynamical Systems/Differential Equations

A **continuous-time dynamical system** is a mathematical model used to describe a system that evolves over time in a way that is continuous (rather than discrete). This means the state of the system changes smoothly as a function of time, without abrupt jumps.

In the context of liquid neural networks, the neurons' states evolve as continuous-time dynamical systems. Each neuron's state is governed by **differential equations**, enabling the network to process information dynamically and adaptively, much like physical systems in the real world. This is inspired by biological neurons, where the activity of each neuron is influenced dynamically by inputs and changes over time.

The state of each neuron  $x_i(t)$  in an LNN evolves over time according to a differential equation, expressed as:

$$\frac{dx_i(t)}{dt} = f(x_i(t), u_i(t), t; \theta_i), \quad (2.1)$$

where:

- $x_i(t)$ : The internal state of the  $i$ -th neuron at time  $t$ ,
- $u_i(t)$ : The input signal to the  $i$ -th neuron at time  $t$ ,
- $t$ : Time, treated as a continuous variable,
- $\theta_i$ : Trainable parameters of the neuron, such as weights and biases,
- $f(\cdot)$ : A function (usually nonlinear) describing the neuron's dynamics.

A common differential equation is the **leaky integrator dynamics**, where the state evolves as:

$$\frac{dx_i(t)}{dt} = -\alpha x_i(t) + \sum_{j=1}^N w_{ij} h(x_j(t)) + u_i(t),$$

with:

- $-\alpha x_i(t)$ : A "leakage" term causing the neuron's state to decay over time, with  $\alpha > 0$  representing the decay rate (temporal decay),
- $\sum_{j=1}^N w_{ij} h(x_j(t))$ : The weighted input from other neurons, where  $w_{ij}$  is the weight from neuron  $j$  to  $i$ , and  $h(x_j(t))$  is a nonlinear activation function (e.g., tanh or ReLU),
- $u_i(t)$ : An external input signal.

For more complex systems, **nonlinear terms** can be included, resulting in equations such as:

$$\frac{dx_i(t)}{dt} = g(x_i(t)) + \sum_{j=1}^N w_{ij} \sigma(x_j(t)) + u_i(t),$$

where  $g(x_i(t))$  models intrinsic nonlinear dynamics, and  $\sigma(x_j(t))$  is a nonlinear activation function.

A liquid neural network, as a continuous-time dynamical system, has several important features. First, it ensures **smooth evolution**, where the neuron states evolve continuously over time according to differential equations. This smooth state transition is essential for modeling time-dependent values in tasks like time-series forecasting or control systems. In addition, the dynamics of the network incorporate **time dependency**  $t$  explicitly or depend solely on the current state  $x(t)$ , enabling the network to capture both static and dynamic temporal relationships. Liquid neural networks are also typically **deterministic**, with their future states fully defined by the current states and inputs, but they can also accommodate **stochastic elements** to model uncertainty or noise in the environment. Finally, the network may operate under **linear** dynamics, such as  $f(x) = Ax + Bu$ , which are efficient but limited in complexity, or **nonlinear** dynamics, like  $f(x) = \tanh(Wx + b)$ , which allow the network to represent intricate patterns and adaptive behaviours.

### 2.1.2 LNN Training

During training, the above differential equations (2.1) define how each neurons processes information. For each labelled training data sample, the following process occurs.

During the **forward pass**, the system of differential equations is numerically solved over time, starting from an initial state  $x(0)$ . Inputs  $u(t)$  and parameters  $\theta_i$  drive the evolution of neuron states  $x_i(t)$ .

The network then outputs a value, derived from the neuron states. This is compared to the target output to compute a **loss function**.

During **backpropagation through time**, gradients of the loss with respect to trainable parameters ( $\theta_i$ ) are computed by differentiating through the differential equations using methods like automatic differentiation or adjoint sensitivity analysis.

Finally, **optimization algorithms** (e.g. gradient descent) update the parameters of the DEs to minimize the loss.



### 2.1.3 LNN Inference

During inference, the same differential equations govern the neuron states, but parameters ( $\theta_i$ ) are fixed. The network processes dynamic inputs  $u(t)$  in real-time. The equation also considers the neuron’s previous state  $x_i(t)$ , which is dependent on previous input values  $u(t)$ . Thus, the output of each neuron is dependent on the parameters, current input values, and previously seen input values.

### 2.1.4 Advantages of LNNs

Using differential equations in LNNs provides several advantages.

#### 2.1.4.a Temporal Modelling

Continuous dynamics are well-suited for time-dependent tasks. This means LNNs can be used to find time-based relationships in data (temporal modeling). This form of ‘memory’ is highly beneficial in time-series tasks.

The nonlinear nature of  $f(x_i(t))$  ensures that the network captures complex temporal dependencies, allowing it to adjust its behavior based on the sequence and timing of inputs. This dynamic capability provides the network with a form of memory, enabling it to adapt to new scenarios even outside the training set.

This is in contrast to static models which consider data points to be independent and identically distributed.

#### 2.1.4.b Adaptability

Dynamic state evolution allows the network to adapt during deployment.

The differential equation model (2.1) for liquid neural networks allows for state evolution even after training, resulting in increased adaptability. This is achieved by the continuous dynamics governing neuron states, which enable the network to respond dynamically to real-time inputs and changing environments.

In the DE model,  $x_i(t)$  is the state of the  $i$ -th neuron at time  $t$ ,  $u_i(t)$  represents external inputs,  $t$  is time, and  $\theta_i$  are trainable parameters (e.g., weights and biases). After training, the parameters  $\theta_i$  are fixed, but the neuron states  $x_i(t)$  continue to respond dynamically to new inputs  $u_i(t)$ . This means the network integrates real-time inputs into its state over time, adapting its behavior dynamically to variations in the input patterns or the timing of events.

The differential equations governing the states ensure that even small variations in the input influence the system, enabling real-time adaptation.

This provides several advantages. LNNs excel in real-world scenarios involving dynamic environments, such as robotics [1] and control systems.

For example, an LNN controlling a robotic arm in a dynamic environment would learn general principles of motion and control during training. During inference, as new obstacles appear or external forces are applied, the network integrates this new information into its state  $x_i(t)$  dynamically. This allows the robotic arm to adjust its movements in real time without needing retraining for each specific scenario.

In addition, by dynamically evolving its states, the network generalizes better to unseen data patterns by interpolating between learned behaviors.

#### 2.1.4.c Efficiency

Liquid neural networks (LNNs) are inherently more efficient than traditional deep neural networks (DNNs) due to their ability to maintain sparser representations. At any given time, only a subset of an LNNs’ neurons or parameters are significantly active or contribute to the system’s computations. This sparsity reduces the computational overhead while retaining the network’s performance and adaptability.

This is because continuous-time dynamics favor selective activity. In LNNs, neuron states evolved continuously over time, according to the differential equation 2.1. Here  $f(\cdot)$  determines how each neuron state changes based on its inputs, past states, and parameters. The use of continuous-time dynamics enables

neurons to become active only when relevant input signals  $u_i(t)$  or temporal events trigger them. This selective activity leads to fewer neurons being active at a given time, resulting in sparse representations.

LNNs are designed to work efficiently with fewer parameters compared to DNNs. While in traditional DNNs, layers are often densely connected, meaning all neurons in one layer interact with all neurons in the next layer. LNNs use sparse connectivity patterns, where neurons only interact with a limited subset of other neurons. This reflects real-world systems, such as biological brains, where neurons form selective, sparse connections. The sparsity of connections reduces the number of computations required during both training and inference.

LNNs allow the internal states of neurons to evolve over time and depend on the dynamics of the inputs. Because of this adaptability, only neurons relevant to the current input remain active. This reduces unnecessary computations and avoids the inefficiencies of global activation in traditional DNNs.

The continuous dynamics of LNNs inherently encode temporal dependencies. Unlike recurrent neural networks (RNNs) or deep learning models that require explicit mechanisms like memory gates (e.g., in LSTMs or GRUs), LNNs rely on the fluid evolution of neuron states. This reduces the overhead of managing and updating memory states, further contributing to sparsity and efficiency.

The sparse nature of LNNs offers several advantages over traditional DNNs, including reduced computational cost (minimizing matrix operations), lower energy consumption, better scalability, and robustness to overfitting (as sparse connectivity can act as a regularization mechanism by ensuring only essential features are focussed on).

#### 2.1.4.d Stability

LNNs exhibit greater stability and robustness to noise compared to traditional DNNs.

Continuous time dynamics and differential equations encode stability constraints, ensuring smooth transitions between states.

In LNNs, the state of each neuron evolves over time according to the differential equation 2.1. The continuous nature of these equations ensures that the neuron states change gradually over time. As a result, sudden spikes in the input  $u_i(t)$  (caused by noise) are naturally smoothed out. This gradual evolution prevents abrupt changes in the neuron states, making the network less sensitive to transient noise.

In addition, neuron states evolve in response to both the current input  $u_i(t)$  and past states  $x_i(t)$ . This integration over time allows the network to prioritize long-term patterns in the input and ignore temporary noise. The feedback from past states enables temporal filtering, where only meaningful input changes accumulate and influence the network’s output. In contrast, the layer-by-layer static activations in DNNs make them more susceptible to noise.

In traditional DNNs, noisy inputs can propagate through the network, often being amplified by dense connections and static parameter updates. To avoid this, special techniques can be used such as dropout. However, LNNs achieve this implicitly, by using sparse and selective connections. This limits the propagation of noise across the network. The continuous evolution of states ensures that transient noise does not significantly affect downstream neurons or outputs.

This enhanced stability has a range of benefits. LNNs perform well in real-world settings where inputs are often corrupted. The intrinsic smoothing abilities of LNNs also reduces the amount of noise-filtering preprocessing required.

## 2.2 Neural Network Verification

In this section we explore the problem of neural network verification. We look at current methods used for DNN verification, which will form the inspiration for a liquid neural network verification approach. The suitability of these to LNNs must be evaluated.

### 2.2.1 The Verification Problem

Verification problems can involve concrete bounds on the input and linear programming (LP) constraints on the output. Formally, the problem can be defined as follows:

**Definition 2.2.1.1** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a neural network, and let  $\mathcal{X} = \{x' \in \mathbb{R}^n \mid x_i^l \leq x_i' \leq x_i^u\}$  represent the set of valid inputs constrained by the lower and upper bounds  $x_i^l, x_i^u$ . Given a set of linear constraints on the output  $\psi_y$ , let  $\mathcal{Y} = \{y \mid \psi_y\}$  denote the set of outputs satisfying  $\psi_y$ . The verification problem is to determine whether  $x' \in \mathcal{X} \implies f(x') \in \mathcal{Y}$ , or to find a counterexample  $x' \in \mathcal{X}$  such that this implication is not true.

For input and output constraints as defined above, the goal is either to prove that no valid input violates the output constraints or to find an input that does. If no input satisfies the output constraints, we declare the property as “safe.” Otherwise, if such an input exists, the property is deemed “unsafe,” and the corresponding input serves as a counterexample. [4]

## 2.2.2 Motivation

Verification of neural networks is a crucial problem, especially when a new architecture (such as liquid neural networks) is being researched. This is because neural networks are often deployed in safety-critical applications, such as autonomous vehicles or medical diagnosis, where unpredictability can cause significant harm. Neural networks are also vulnerable to adversarial attacks, which is when small perturbations within input data (often unnoticeable to the human eye) cause significant undesired changes in the output. This vulnerability poses a serious threat to their reliability and trustworthiness. Verification ensures that the network behaves as expected under specified conditions, whilst robustness verification focuses on guaranteeing that small perturbations in the input do not lead to misclassifications or unsafe behavior. By formally proving properties of neural networks or identifying counterexamples, verification helps to ensure safety and mitigate risks in real-world deployments.

## 2.2.3 Recurrent Neural Networks

We now focus on the verification problem in relation to recurrent neural networks (RNNs) specifically. RNNs are deep neural networks trained on sequential or time series data to create a model that can make sequential predictions or conclusions based on sequential inputs. During both training and inference, they also use information from prior inputs to influence the current input and output. In traditional recurrent neural networks, this is achieved by a feedback loop within the network, containing a hidden state which ‘remembers’ previous inputs. [5]

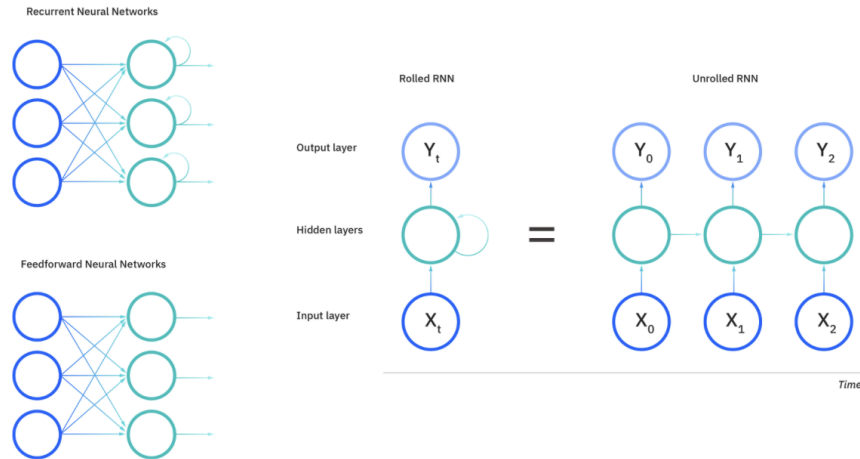


Figure 2.1: Feedforward (traditional) vs. Recurrent Neural Networks

RNNs rely on Backpropagation Through Time (BPTT) to compute gradients, unlike traditional neural networks that use standard backpropagation. Whilst BPTT follows the same principles as traditional methods, adjusting parameters by propagating errors from output to input layers, it accounts for sequential data by summing errors across time steps. This temporal error accumulation differentiates BPTT from the simpler gradient computation in feedforward networks, which lack a time-dependent structure.

‘Memory’ in RNNs can be achieved by several architectures, such as LSTM (Long Short-Term Memory), GRU (Gated Recurrent Units) and Encoder-decoder RNNs. Since liquid neural networks leverage

differential equations to achieve a form of 'memory', they are a specialized type of RNN. Whilst their weights are typically fixed, the hidden state evolves dynamically over time, driven by the structure of the differential equations and the input. This continuous evolution allows liquid networks to retain memory and capture temporal dependencies across varying time scales.

The verification problem for RNNs involve constraints on both the input sequences and the dynamic outputs of the network, with the goal of solving the problem stated earlier (2.2.1). This requires handling both the sequence-based inputs and the time-evolving hidden states.

## 2.2.4 Robustness Verification for Recurrent Neural Networks

An important verification problem for RNNs concerns their robustness to temporal perturbations or noise in sequential inputs. Robustness implies that small perturbations in the input sequence do not cause significant deviations in the network's output. This can be formalized as follows:

**Definition 2.2.4.1.** Let  $x \in \mathbb{R}^n$  be a sequential input to a recurrent neural network  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , where  $m > 1$  represents the dimensionality of the output space. Let  $\mathcal{C}_x = \{x' \in \mathbb{R}^n \mid x'_i \leq x_i \leq x'_u\}$  represent the set of perturbed input sequences. The *targeted robustness verification problem* is to determine whether

$$x' \in \mathcal{C}_x \implies f(x')_c > f(x')_t \quad \text{for a specific output target } t,$$

or to find a counterexample  $x'$  such that this implication is not satisfied.

**Definition 2.2.4.2** Let  $x \in \mathbb{R}^n$  be a sequential input to a recurrent neural network  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , where  $m > 1$ . Let  $\mathcal{C}_x = \{x' \in \mathbb{R}^n \mid x'_i \leq x_i \leq x'_u\}$ . The *general robustness verification problem* is to determine whether

$$x' \in \mathcal{C}_x \implies f(x')_c > f(x')_t \quad \text{for all } t \neq c,$$

or to find a counterexample  $x'$  where this implication fails.

Robustness verification for RNNs focuses on ensuring that temporal variations in sequential inputs do not cause undesired behavior. Specifically, the targeted robustness problem can often be reduced to verifying specific temporal constraints or perturbations in the input sequence. The general robustness problem is more challenging because it involves verifying the network's behavior across all possible perturbations and output dimensions. These problems may require advanced techniques, which will be discussed in the following sections.

## 2.2.5 Symbolic and Interval Propagation Methods

This section explores methods that focus on propagating bounds through liquid neural networks, which verifies their robustness when faced with input perturbations. **Symbolic Interval Propagation (SIP)** is a technique that computes conservative output bounds using interval arithmetic, providing a computationally efficient way to check for robustness. **CROWN (Certified Robustness to Weight Perturbations)** is more complex, and introduces linear approximations which improves the precision of robustness verification. **Lipschitz-based methods** they estimate global sensitivity by bounding the Lipschitz constant of the network. These approaches are particularly useful for ensuring scalability and efficiency, making them suitable for applications where lightweight and real-time verification is required.

### 2.2.5.a SIP (Symbolic Interval Propagation)

SIP is a scalable and efficient technique for verifying liquid neural networks by propagating symbolic intervals through the layers of the network to bound the range of possible outputs. For LNNs governed by neural ODEs with the equation 2.1, SIP can be applied to discretize and propagate bounds over time, handling non-linearity and ensuring robustness and safety properties are satisfied under perturbations.

**Steps in SIP:**

1. **Input Interval Initialization:** Define the input range as intervals  $[l_i, u_i]$  for each dimension  $x_i$ , forming a hyper-rectangle. These intervals are represented symbolically to maintain dependencies between variables.
2. **Symbolic Propagation Through Layers:** At each discretized time step or layer:

- **Affine Transformations:** For a layer  $z = Wx + b$ , bounds are propagated symbolically:

$$l_z = W \cdot l_x + b, \quad u_z = W \cdot u_x + b.$$

- **Nonlinear Activations:** Nonlinearities like ReLU or tanh are handled by updating bounds:

$$\text{ReLU: } [\max(0, l_z), \max(0, u_z)], \quad \text{Tanh: } [\tanh(l_z), \tanh(u_z)].$$

3. **Output Interval Verification:** The final output intervals are compared against safety or robustness properties. For example, given input perturbations  $\delta x$ , the output bounds  $F(x)$  must satisfy:

$$F(x + \delta x) \subseteq [F_L, F_U],$$

where  $F_L$  and  $F_U$  are symbolic bounds on the output.

**Advantages for Liquid Neural Networks** SIP is well-suited for LNNs due to its ability to handle the continuous evolution of states in neural ODEs:

- **Precision:** Symbolic intervals maintain variable dependencies, producing tighter bounds than traditional interval arithmetic.
- **Efficiency:** Propagation avoids the computational cost of exact methods, making SIP scalable for larger networks.
- **Adaptability:** SIP can handle nonlinear dynamics in LNNs through accurate approximations of activation functions.

**Practical Considerations** One consideration is over-approximation - accumulated conservativeness may reduce precision in deeper networks. Also, complex nonlinearities (within highly nonlinear layers such as softmax) require additional approximations, introducing potential conservatism. Finally, for neural ODEs, time discretization must balance accuracy with computational cost.

#### 2.2.5.b CROWN (Certified Robustness to Weight Perturbations)

CROWN (Certified Robustness to Weight Perturbations) is a general framework for certifying the robustness of neural networks, including those with non-linear activation functions. It achieves this by bounding the outputs of the network using adaptive linear (or quadratic) upper and lower bounds for each activation function. For liquid neural networks, governed by neural ODEs, CROWN can be adapted to certify robustness by discretizing the continuous dynamics and applying its bounding technique to each time step.

CROWN provides an efficient and scalable method for verifying the robustness of LNNs. The propagated adapted bounds ensure that perturbations in the input do not lead to significant deviations in the output.

**Mathematical Framework** Consider a liquid neural network modeled as:

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

where  $x(t) \in \mathbb{R}^n$  is the state,  $f(x, t, \theta)$  describes the dynamics, and  $\theta$  represents the parameters. For a given perturbed input  $x_0 \in \mathbb{R}^n$  within an  $\ell_p$ -ball:

$$x \in B_p(x_0, \epsilon) = \{x \mid \|x - x_0\|_p \leq \epsilon\},$$

CROWN aims to compute certified bounds  $F_L(x) \leq F(x) \leq F_U(x)$ , where  $F(x)$  is the output of the neural ODE after a fixed time horizon  $T$ .

**Bounding Nonlinearities** For each activation function  $\sigma(y)$ , CROWN constructs linear upper and lower bounds:

$$h_U(y) = \alpha_U y + \beta_U, \quad h_L(y) = \alpha_L y + \beta_L,$$

such that  $h_L(y) \leq \sigma(y) \leq h_U(y)$  over a pre-activation range  $[l, u]$ . These bounds are propagated through the layers of the network. For LNNs, this involves discretizing the time domain into intervals  $[t_k, t_{k+1}]$  and applying the bounds iteratively at each time step.

**Output Certification** To certify robustness, CROWN computes bounds on the network’s output  $F(x)$ . Using the layer-by-layer propagation of the upper and lower bounds, the output bounds are expressed as:

$$F_U(x) = \Lambda^{(0)}x + \sum_{k=1}^m \Lambda^{(k)}(b^{(k)} + \Delta^{(k)}), \quad F_L(x) = \Omega^{(0)}x + \sum_{k=1}^m \Omega^{(k)}(b^{(k)} + \Theta^{(k)}),$$

where  $\Lambda$  and  $\Omega$  are matrices representing the upper and lower bound propagation, and  $\Delta$  and  $\Theta$  account for biases introduced by non-linearities.

**Application to Neural ODEs** For neural ODEs, the propagation framework is adapted to account for the continuous evolution of states. The bounds are computed at each discretized step  $t_k$ , ensuring that the dynamics satisfy the robustness conditions:

$$F_U(x_0) - F_L(x_0) \geq \delta,$$

where  $\delta$  is the minimum required margin for robustness.

**Practical Implementation** CROWN can be implemented as follows:

1. **Define Input Bounds:** Specify the  $\ell_p$ -ball around the input  $x_0$  and initialize the pre-activation bounds for the first layer.
2. **Propagate Bounds:** Compute the upper and lower bounds layer-by-layer using the adaptive linear approximations.
3. **Certify Robustness:** Verify that the certified bounds at the output satisfy the desired robustness property (e.g., consistent classification).

[6]

### 2.2.5.c Lipschitz-Based Methods

Lipschitz-based methods provide a robust framework for verifying the safety and robustness of liquid neural networks by quantifying how sensitive a network’s outputs are to perturbations in its inputs. Their ability to provide global robustness guarantees is useful for neural ODEs. The Lipschitz constant of a network bounds the maximum rate at which outputs can change with respect to changes in inputs, ensuring that small input perturbations do not lead to large deviations in the output.

**Mathematical Foundation** For a liquid neural network modeled as:

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

where  $x(t) \in \mathbb{R}^n$  is the state,  $f(x, t, \theta)$  defines the dynamics, and  $\theta$  are the network parameters, the Lipschitz constant  $L$  satisfies:

$$\|F(x) - F(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^n,$$

where  $F(x)$  represents the solution to the neural ODE at the final time  $T$ . The Lipschitz constant  $L$  bounds the global sensitivity of the network.

**Estimation of the Lipschitz Constant** The Lipschitz constant can be estimated for an LNN by analyzing the Jacobian of  $f(x, t, \theta)$ . For a time-discretized system, the sensitivity of the network is determined by:

$$L = \sup_{t \in [0, T]} \|J_f(x, t)\|,$$

where  $J_f(x, t) = \frac{\partial f(x, t, \theta)}{\partial x}$  is the Jacobian matrix of  $f$ . Computing or bounding  $L$  involves:

- **Spectral Norm Analysis:** Evaluating  $\|J_f(x, t)\|$  as the largest singular value of the Jacobian at each time step.
- **Pathwise Integral Bounds:** For neural ODEs,  $L$  can be bounded using the integral of the Jacobian along the trajectory:

$$L \leq \int_0^T \|J_f(x(t), t)\| dt.$$

**Verification Applications** Lipschitz-based methods are widely used for:

1. **Robustness Verification:** Verifying that small input perturbations  $x_0 \rightarrow x_0 + \delta x$  result in bounded output deviations, ensuring:

$$\|F(x_0 + \delta x) - F(x_0)\| \leq L\|\delta x\|.$$

2. **Safety Analysis:** Ensuring the network's outputs remain within a safe region under bounded input perturbations.
3. **Adversarial Robustness:** Certifying that adversarial inputs cannot change the classification or decision boundaries within a certain radius.

**Practical Implementation** To implement Lipschitz-based verification for LNNs:

1. **Compute the Lipschitz Constant:** Use numerical methods, such as spectral norm approximation or pathwise integration, to estimate  $L$ .
2. **Bound Output Sensitivity:** Evaluate  $L$  for given input perturbations  $\delta x$  and verify that the resulting outputs satisfy safety and robustness criteria.
3. **Scaling for Efficiency:** For high-dimensional networks, consider approximation techniques or layer-wise bounds to improve scalability.

## Chapter 3

# Liquid Neural Network Design and Implementation

### 3.1 Design Overview

This chapter outlines the implementation of the Liquid Neural Network (LNN) developed in PyTorch for sequential 2D time-series prediction. The architecture is based on the Liquid Time-Constant (LTC) neuron model, which simulates continuous-time dynamics through ordinary differential equations (ODEs) and shows properties of neural adaptability and temporal memory.

The aim of the implementation was to create a biologically-inspired, interpretable recurrent model with competitive performance on trajectory prediction tasks. Unlike conventional RNNs or LSTMs, the LNN is governed by time-continuous equations rather than discrete updates, providing finer control over neuronal dynamics.

The following principles guided the design:

- **Framework:** PyTorch was selected due to its flexible dynamic graph construction and ease of integrating custom layers with automatic differentiation.
- **Neuron Dynamics:** The neuron model was designed to emulate leaky integrate-and-fire (LIF) behaviour with added plasticity through modulated reversal potentials and conductances.
- **Time Unfolding:** Each forward pass of the LNN integrates over multiple internal time steps (ODE unfolds) to approximate the continuous-time solution, reflecting membrane voltage evolution.
- **Baseline Comparison:** To benchmark performance, identical training and evaluation protocols were implemented for alternative architectures (LSTM, TCN) using the same data.

The following sections document the architecture, neuron formulation, wiring strategy, training setup, and performance characteristics of the LNN.

### 3.2 Wiring and Connectivity

LNNs have a sparse and biologically motivated connectivity structure. To simulate the non-uniform and random nature of synaptic wiring observed in biological networks, a custom class named `RandomWiring` was implemented.

This class generates two adjacency matrices:

- A **recurrent adjacency matrix** of shape  $(n \times n)$  defining internal connections between neurons within the hidden layer.
- A **sensory adjacency matrix** of shape  $(d_{\text{in}} \times n)$  which defines the input-to-hidden connectivity.



Each matrix contains continuous values sampled from a uniform distribution on  $[0,1]$ , which are later used to create binary masks or to modulate weight strengths.

The `RandomWiring` class also generates reversal potentials: `erev` for neuron-neuron connections, and `sensory_erev` for input-synapse connections. These potentials are initialised from a uniform range  $[-0.2, 0.2]$  and are treated as fixed, non-learnable parameters.

The use of fixed sparse masks in the class emulate the limited number of active connections in real cortical microcircuits, enabling **biological plausibility**. In addition, each **randomised instantiation** of `RandomWiring` results in a different network topology, allowing stochastic variation in experiments.

Finally, the sensory and recurrent wiring are **decoupled**, to enable the model to explicitly distinguish between input-driven and internal dynamic behaviour.

Below is a simplified example of the `RandomWiring` class:

```

1 class RandomWiring:
2     def __init__(self, input_dim, output_dim, neuron_count):
3         self.adjacency_matrix = np.random.uniform(0, 1, (neuron_count,
4             neuron_count))
5         self.sensory_adjacency_matrix = np.random.uniform(0, 1, (input_dim,
6             neuron_count))
7
8     def erev_initializer(self):
9         return np.random.uniform(-0.2, 0.2, (neuron_count, neuron_count))
10
11     def sensory_erev_initializer(self):
12         return np.random.uniform(-0.2, 0.2, (input_dim, neuron_count))

```

Listing 3.1: Simplified `RandomWiring` class

### 3.3 LTC Neuron Dynamics

The core unit of the LNN is the `LIFNeuronLayer`, a custom PyTorch module that simulates the behaviour of leaky integrate-and-fire neurons (also known as liquid time-constant neurons). These neurons operate using a continuous-time dynamical model controlled by a first-order differential equation, capturing the evolution of membrane potentials in response to internal and external stimuli.

The model integrates over time using a discretised ODE solver implemented within the forward pass. Specifically, it unfolds the membrane update equation over a fixed number of steps (`ode_unfolds`) using an Euler-like method.

The update rule is determined by the equation:

$$\begin{aligned}
 v_i^{(t+1)} = \frac{1}{Z} & \left( \mathcal{S}(c_{m,i}) \cdot v_i^{(t)} + \mathcal{S}(g_{\text{leak},i}) \cdot V_{\text{leak},i} \right. \\
 & + \underbrace{\sum_j \mathcal{S}(w_{ij}) \cdot \sigma(v_j^{(t)}; \mu_{ij}, \sigma_{ij}) \cdot E_{\text{rev},ij}}_{\text{Recurrent synaptic current}} \\
 & \left. + \underbrace{\sum_k \mathcal{S}(w_{ki}^{\text{sens}}) \cdot \sigma(x_k; \mu_{ki}^{\text{sens}}, \sigma_{ki}^{\text{sens}}) \cdot E_{\text{rev},ki}^{\text{sens}}}_{\text{Sensory synaptic current}} \right)
 \end{aligned} \tag{3.1}$$

$$Z = \mathcal{S}(c_{m,i}) + \mathcal{S}(g_{\text{leak},i}) + \sum_j \mathcal{S}(w_{ij}) \cdot \sigma(v_j^{(t)}; \mu_{ij}, \sigma_{ij}) + \sum_k \mathcal{S}(w_{ki}^{\text{sens}}) \cdot \sigma(x_k; \mu_{ki}^{\text{sens}}, \sigma_{ki}^{\text{sens}}) + \varepsilon \tag{3.2}$$

where:

- $\mathcal{S}(\cdot) = \text{Softplus}(\cdot) = \log(1 + e^{(\cdot)})$ : nonlinearity ensuring positivity of weights and conductances

- $v_i^{(t)}$ : membrane potential of neuron  $i$  at ODE step  $t$
- $x_k$ : sensory input from dimension  $k$
- $\sigma(v; \mu, \sigma) = \text{sigmoid}(\sigma(v - \mu))$ : synaptic activation function
- $w_{ij}$ : weight from neuron  $j$  to neuron  $i$  (recurrent)
- $w_{ki}^{\text{sens}}$ : weight from input  $k$  to neuron  $i$  (sensory)
- $\mu, \sigma$ : learnable sigmoid parameters (mean and scale)
- $E_{\text{rev}}, E_{\text{rev}}^{\text{sens}}$ : synaptic reversal potentials
- $V_{\text{leak}}$ : fixed leak reversal potential
- $c_m, g_{\text{leak}}$ : learnable membrane capacitance and leak conductance
- $\varepsilon$ : small constant for numerical stability

This formulation reflects the true internal update logic of the `LIFNeuronLayer` implementation, where each neuron’s potential evolves through a biophysically inspired ODE governed by input-specific and recurrent dynamics.

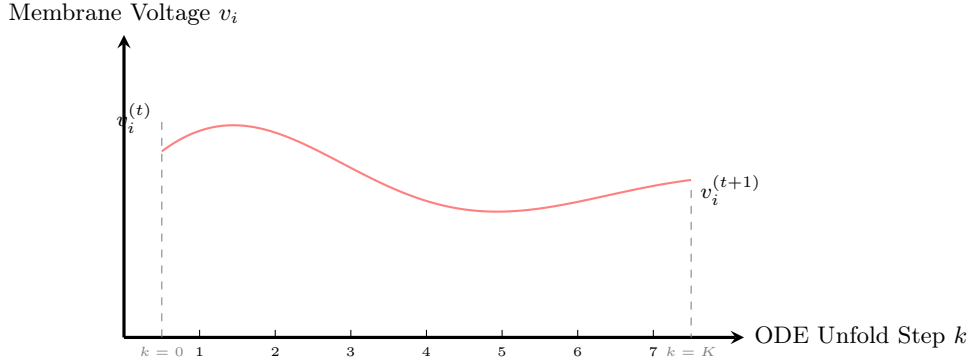


Figure 3.1: Illustrative internal membrane voltage trace across  $K$  ODE unfolding steps for a single neuron. This entire graph represents the internal dynamics used to compute  $v_i^{(t+1)}$  from  $v_i^{(t)}$ , in a single timestep.

Unlike traditional RNNs or LSTMs, which update their hidden state in a single discrete operation per time step, Liquid Time-Constant neurons simulate fine-grained membrane voltage dynamics by performing multiple internal updates within each input timestep. This behaviour is governed by the discretised solution of the differential equations (see Equation 3.1 and Equation 3.2).

At each unfolding step, the membrane potential is updated using a softplus-modulated combination of capacitive memory, leak current, and synaptic input currents. These updates reflect physical dynamics such as charging, leak, and synaptic integration, enabling the neuron to adaptively integrate information over sub-timestep resolution. The final membrane voltage  $v_i^{(t+1)}$  is the result of this integration process and is passed forward in time.

This is different from gated memory in LSTMs, which uses learnable gates and affine transformations, or temporal convolutions in TCNs, which apply fixed receptive filters. Instead, the LTC neuron learns dynamic temporal behaviour through time constants and nonlinear integration, making it particularly suited to tasks that require fine temporal resolution and state-dependent transitions.

#### Implementation Details:

- **Learnable Parameters:** All biophysical constants (capacitance, leak conductance, reversal potentials, synaptic weights) are learnable, providing flexibility in dynamic behaviour.
- **Softplus Regularisation:** Weights and conductances are passed through `Softplus` to enforce positivity while allowing gradients to flow smoothly during training.

- **ODE Unfolding:** The number of internal solver steps is fixed (`ode_unfolds = 12`) to balance numerical precision with computational cost.
- **Sparsity Masks:** Both recurrent and sensory activations are element-wise masked using the adjacency matrices from `RandomWiring`, enforcing fixed sparsity throughout training.

Below is the ODE solver implementation within the `LIFNeuronLayer` class:

```

1 def ode_solver(self, inputs, state, elapsed_time):
2     v_pre = state
3     for _ in range(self.ode_unfolds):
4         synaptic_input = compute_synaptic_activation(v_pre)
5         numerator = self.cm * v_pre + self.gleak * self.vleak + synaptic_input
6         denominator = self.cm + self.gleak + synaptic_conductance
7         v_pre = numerator / (denominator + self.epsilon)
8     return v_pre

```

Listing 3.2: Simplified LTC neuron forward method

This allows neurons to respond to both present input and also to their internal temporal dynamics, mimicking continuous-time memory traces observed in biological neurons.

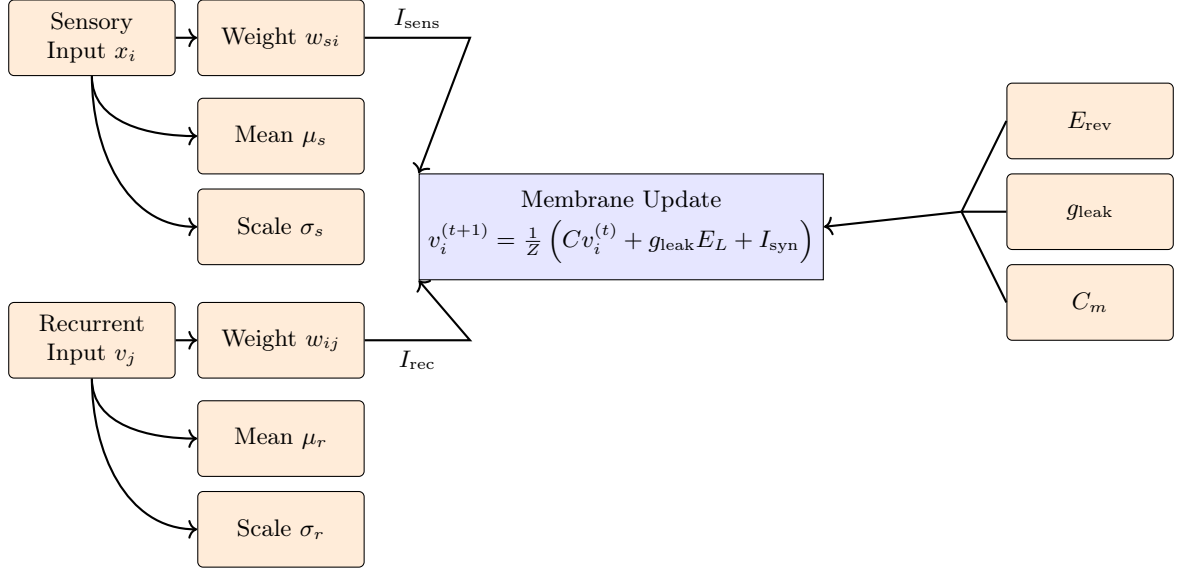


Figure 3.2: Internal structure of a single LIF neuron used in the Liquid Time-Constant Network. Inputs undergo non-linear transformations based on trainable  $\mu$  and  $\sigma$ , and the resulting activations are integrated using biophysical parameters (leak conductance  $g_{\text{leak}}$ , membrane capacitance  $C_m$ , and reversal potentials  $E_{\text{rev}}$ ).

### 3.4 Network Architecture

The full LNN is constructed by embedding the LTC neuron layer inside a recurrent wrapper, implemented as a custom `LTCRNN` module. This wrapper sequentially passes each time step of the input through the same `LIFNeuronLayer`, maintaining a hidden state that evolves over time. The resulting structure can be viewed as a biologically grounded alternative to traditional RNN cells.

The architecture accepts an input tensor of shape  $(B, T, d_{\text{in}})$ , where  $B$  is the batch size,  $T$  is the sequence length, and  $d_{\text{in}}$  is the input dimension (two in this case, corresponding to 2D spatial coordinates). For each time step  $t$ , the neuron layer receives the  $t$ -th slice of the sequence and updates the hidden state, generating a predicted output of shape  $(B, T, d_{\text{out}})$ .

## Design Considerations

While designing the LTCRNN architecture, there were several key design decisions made. The hidden state dimensionality (i.e. number of LTC neurons) defined the model capacity; lower values reduce overfitting risk and improve computational efficiency at the cost of limited expressiveness. In addition, the voltage traces themselves are treated as predictions instead of applying a separate output layer. This means membrane state is directly used as a continuous output signal. To maximise efficiency (and GPU parallelism) of tensor operations during training/inference, all sequences are processed in batch-major form (following PyTorch convention).

Below is a simplified version of this architecture:

```

1 class LTCRNN(nn.Module):
2     def __init__(self, wiring, input_dim, hidden_dim, output_dim):
3         self.cell = LIFNeuronLayer(wiring)
4         ...
5
6     def forward(self, inputs):
7         batch_size, seq_len, _ = inputs.size()
8         states = torch.zeros(batch_size, self.hidden_dim)
9         outputs = []
10        for t in range(seq_len):
11            out, states = self.cell(inputs[:, t, :], states)
12            outputs.append(out)
13        return torch.stack(outputs, dim=1)

```

Listing 3.3: Structure of the LTCRNN module

The design maintains a clear distinction between the continuous-time neuronal dynamics and the sequence-level integration logic. Thus, it is both modular and biologically interpretable, and compatible with standardised modern deep learning tools.

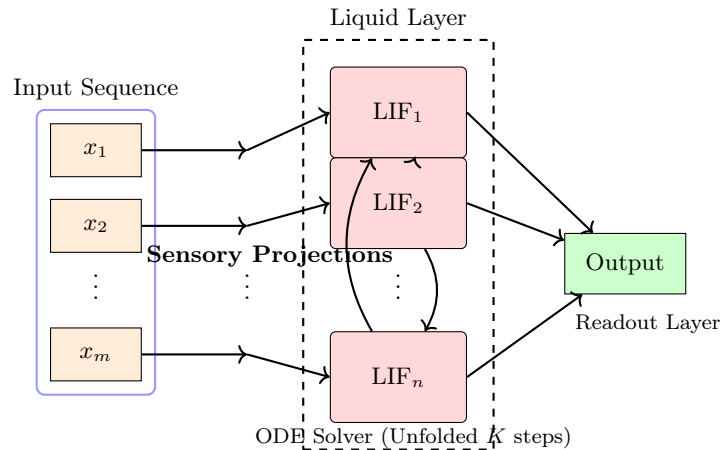


Figure 3.3: Architecture of the Liquid Time-Constant Network (LNN). Inputs project through learned sensory filters to a sparsely recurrent Liquid Layer of LIF neurons. Dynamics are integrated using an internal ODE solver with unfolding. Final outputs are read from a low-dimensional projection.

## 3.5 Training Configuration (and dataset)

The LNN was trained on a synthetic 2D spiral trajectory dataset, chosen for its smooth temporal structure and nonlinearity. Each data point consists of an  $(x, y)$  coordinate, and the model’s aim is to predict the next point in the sequence (given a fixed-length input window). The ‘sequential’ nature of the task makes it well-suited for testing temporal memory and continuous dynamics.

The synthetically-generated dataset ensured control over noise and resolution, which allowed clearer attribution of error sources to model limitations rather than data irregularities.

A supervised learning approach was used; inputs and targets were created by shifting a sliding window of length  $T = 3$  over the full spiral. The small sequence length was chosen to reduce training complexity while still allowing temporal dependencies to be captured. Each input sequence of three time steps was paired with the corresponding next three steps as the target output.

The dataset was split into training and validation sets, with the training set containing 80% of the data and the validation set containing 20%. This split was randomised to ensure that the model generalised well to unseen data.

subsubsection\*Data Preprocessing

- All inputs were standardised using the training set mean and standard deviation.
- Targets were normalised in the same way to preserve scale consistency.
- The spiral dataset was generated programmatically with adjustable number of points and turns.

subsubsection\*Training Parameters Mean Squared Error (`nn.MSELoss()`) was used to penalise deviations from the ground truth trajectory. Adam was chosen as the optimiser due to its fast convergence and robustness to parameter scaling, with a learning rate of 0.005. The model was trained for 2000 epochs, with periodic visual evaluation every 100 epochs. Input sequences were split into overlapping windows and grouped into batches of size 32, allowing efficient GPU utilisation while preserving temporal continuity. A random 80/20 train-validation split was applied, with shuffling to prevent memorisation of input order.

Below is the training process implementation:

```

1 for epoch in range(num_epochs):
2     lnn_model.train()
3     total_loss = 0
4     for x_batch, y_batch in zip(input_batches, target_batches):
5         optimizer.zero_grad()
6         outputs = lnn_model(x_batch)
7         loss = criterion(outputs, y_batch)
8         loss.backward()
9         optimizer.step()
10    total_loss += loss.item()

```

Listing 3.4: Simplified training loop for the LNN

The training loop includes evaluation checkpoints where predicted trajectories are plotted and compared to ground truth. These visualisations provided insights into convergence behaviour, beyond scalar loss values.

## 3.6 Training Behaviour

Throughout training, model performance was monitored both quantitatively (via validation loss) and qualitatively (through trajectory plots), every 100 epochs.

The main trends observed during training were:

- Loss decreased steadily in early epochs, with diminishing returns as training progressed.
- In some cases, small fluctuations in validation loss were observed, likely due to the non-convexity of the parameter landscape and the biological variability induced by random wiring.
- Visual predictions of the trajectory showed clear improvement over time. Early predictions were coarse approximations, while later epochs gave smoother and more accurate predictions.

### Trajectory Loss Curves

To evaluate the model’s learning progress and generalisation, we record the full-sequence mean squared error (MSE) first on the entire spiral sequence (which contains the training datapoints and the held-out validation datapoints), and then on a separate unseen evaluation spiral, at each epoch. Rather than separately plotting training and validation loss — which are contiguous segments of the same spiral

trajectory, they are treated as a single sequence. This is because temporal continuity is essential in modelling dynamical systems. This avoids misleading interpretations that might arise from artificial segmentation of a naturally evolving system.

The following plots illustrates training and validation behaviour over time (the second is zoomed to highlight differences between the training/validation spiral and evaluation spiral):

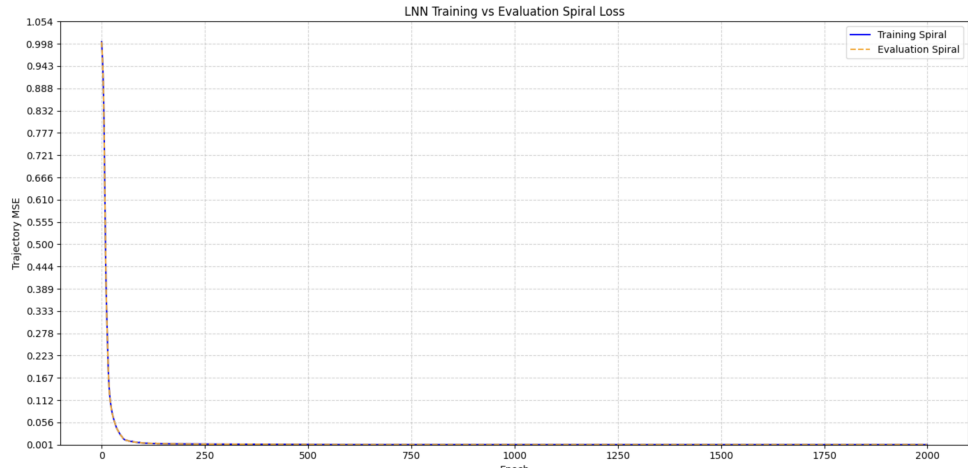


Figure 3.4: Training and validation loss over epochs.

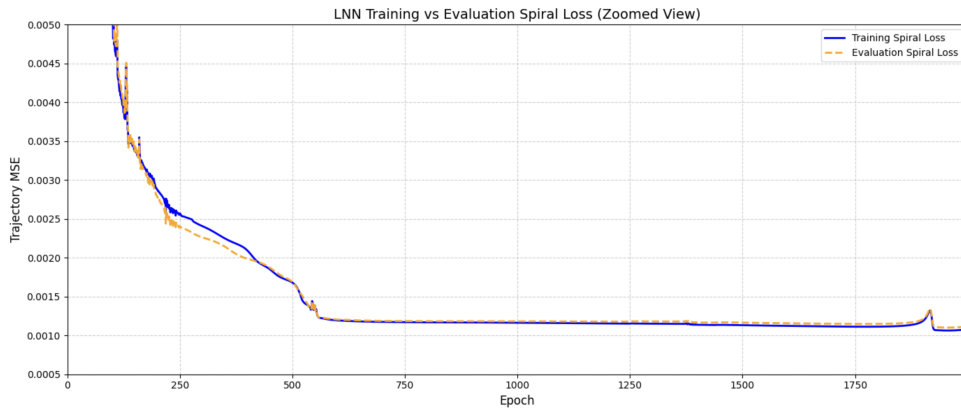


Figure 3.5: Training and validation loss over epochs (zoomed).

Early in training, both curves decrease rapidly, suggesting that the model learns the underlying structure efficiently. As training proceeds, the two curves converge, with the evaluation spiral maintaining a slightly higher loss, indicating some generalisation gap but also demonstrating stable extrapolation beyond the training path. Importantly, neither curve exhibits significant divergence or overfitting behaviour, which supports the robustness and consistency of the trained model.

## Qualitative Evaluation

Visually, the predicted path over time showed that the LNN was able to maintain smooth curvature and approximate the rotational dynamics of the spiral without overshooting or excessive lag. This was true even on validation data not seen during training.

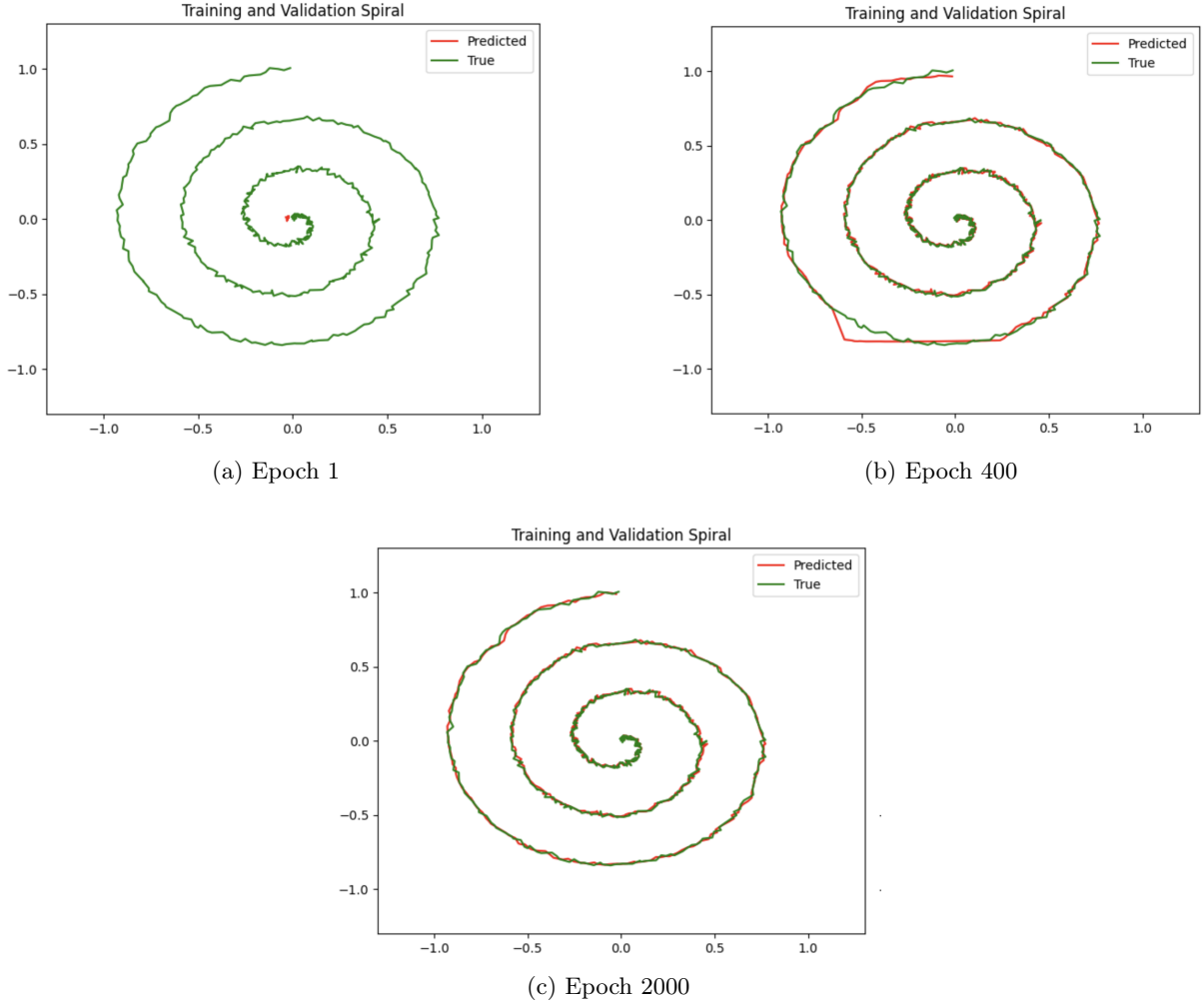


Figure 3.6: LLN predicted vs true spiral trajectories across training: early, mid, and final epochs (denormalised training and validation spiral)

### Observed Patterns Across Training Runs

The number of internal steps in the membrane integration process (**ODE unfolding depth**) contributed significantly to trajectory stability. Deeper unfolding improved smoothness, but with diminishing returns. The architectures' fixed wiring helped constrain overfitting and contributed to better generalisation than a fully connected architecture. Notably, runs with different initialisations showed **performance variations** in loss curves and convergence speed, indicating sensitivity to initial wiring or parameter seeds.

Despite simpler architectures having shorter training times (e.g. LSTMs), LLN's had superior interpretability and stability in capturing the underlying continuous structure of the problem.

Although trained on fixed-length input windows, models were evaluated on full-sequence inference without autoregressive rollouts. This mismatch is empirically validated to yield stable and accurate predictions on smooth trajectories.

# Chapter 4

## Comparative Models

### 4.1 Introduction to Baseline Models

To provide context for the robustness and verification of the Liquid Neural Network (LNN), we benchmark its performance against three established neural architectures: the Temporal Convolutional Network (TCN), the Long Short-Term Memory (LSTM) network, and the Transformer model. All models are trained on the same trajectory prediction task, using identical datasets, normalisation, loss functions, and training schedules as the LNN.

The selection of these baselines is motivated by their contrasting inductive biases and proven success in modelling sequential data. The LSTM is a recurrent architecture that introduces gating mechanisms and persistent internal states, enabling it to model long-range temporal dependencies through iterative state updates. The TCN, however, uses on dilated causal convolutions and fixed temporal receptive fields, making it structurally different from recurrent networks and well-suited for parallel computation. The transformer model is an attention-based architecture, which eliminates recurrence altogether, and dynamically weight input positions using self-attention mechanisms. It applies positional encodings to retain order information.

By evaluating the behaviour of these models under clean and adversarial conditions, we aim to identify both their predictive accuracy, and their robustness, sensitivity to perturbation, and qualitative output characteristics. These insights provide context for assessing the LNN in the following aspects:

- **Temporal memory:** How effectively each model retains and processes sequential dependencies
- **Structural robustness:** The influence of architectural constraints on noise sensitivity
- **Gradient stability:** The relationship between loss geometry and adversarial vulnerability

### 4.2 Temporal Convolutional Network (TCN)

#### Overview and Motivation

The TCN is a fully convolutional architecture designed for sequential data. Unlike RNN-based models, which process inputs recursively and maintain an internal hidden state, TCNs rely on 1D convolutions applied over the temporal axis. This allows for parallel computation and more stable gradients, particularly for long sequences.

TCN's use **dilated convolutions**, which expand the receptive field exponentially with depth while preserving causality. This makes them highly effective at modelling long-range dependencies without the vanishing gradient issues that often affect RNNs.



## Theoretical Background

For a 1D input sequence  $x \in \mathbb{R}^{T \times d}$ , a dilated convolution with kernel  $k$  and dilation factor  $d$  is defined as:

$$(y *_d k)(t) = \sum_{i=0}^{k-1} k(i) \cdot x(t - d \cdot i)$$

This structure allows the model to observe wider contexts with fewer parameters and layers.

In practice, the TCN is constructed using **residual blocks** with stacked dilated convolutions, dropout, and skip connections to stabilise training. Zero-padding is used to ensure output length matches input length.

## Model Architecture

The implemented TCN consists of 3 residual blocks, each with: two dilated 1D convolutional layers with kernel size 3, ReLU activations and dropout regularisation, and optional 1x1 convolutions for matching input-output dimensions. Each residual block has a dilation rate of 1, 2, and 4 respectively (exponentially increasing dilation).

An output convolution maps the final hidden representation to the desired 2D coordinate space.

Below is a simplified example of the TCN model implementation:

```
1 class ResidualBlock(nn.Module):
2     def __init__(self, in_channels, out_channels, kernel_size, dilation, dropout
3         ):
4         ...
5         self.conv1 = nn.Conv1d(..., dilation=dilation)
6         self.conv2 = nn.Conv1d(..., dilation=dilation)
7
8 class TCN(nn.Module):
9     def __init__(self, input_dim=2, hidden_channels=128, ...):
10        self.tcn = nn.Sequential(*residual_blocks)
11        self.output_layer = nn.Conv1d(hidden_channels, output_dim, 1)
```

Listing 4.1: Simplified TCN architecture

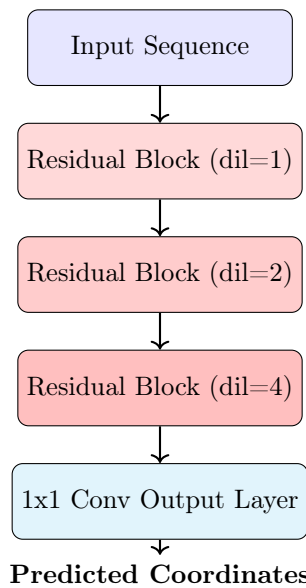


Figure 4.1: Temporal Convolutional Network (TCN) architecture

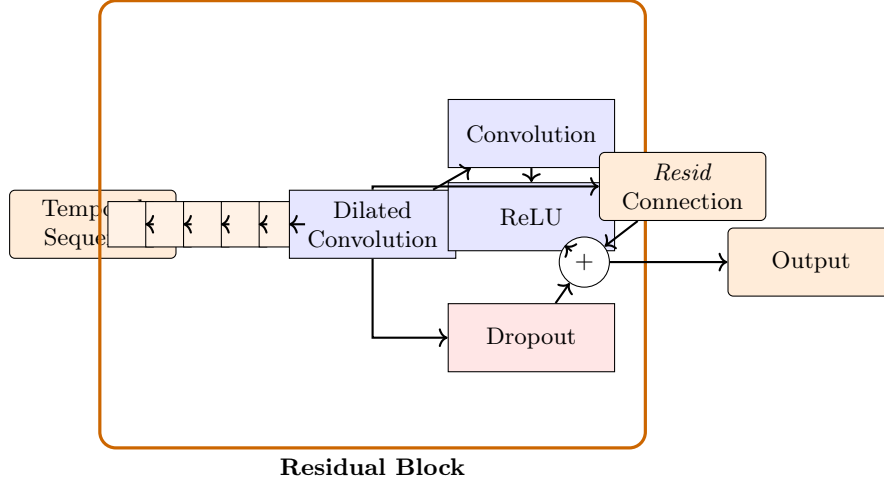


Figure 4.2: Structure of a TCN Residual Block. Each convolution uses dilation to increase receptive field, while residual skip connections and dropout stabilise training.

## Training Configuration

The TCN was trained on the same spiral dataset as the LNN, with identical batch size, learning rate, loss function (Smooth L1), and normalisation pipeline. The model was optimised using Adam and a learning rate scheduler that halved the rate every 500 steps.

## Performance and Behaviour

The TCN demonstrated strong performance on the trajectory prediction task, converging more quickly than the LNN and producing smooth outputs even with a small receptive field. The use of dilated convolutions allowed the model to predict coordinated curvature without explicitly tracking hidden state over time.

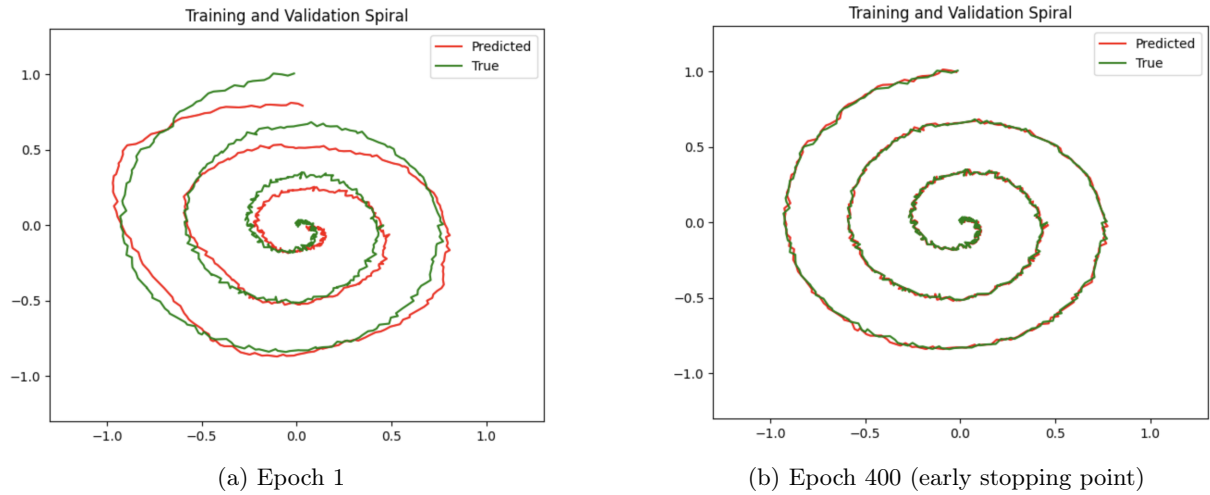


Figure 4.3: TCN predicted vs true spiral trajectories across training: early and final epochs (denormalised training and validation spiral)

## Design Considerations

- **Causality:** All convolutions were causal, ensuring no future information was used during prediction.
- **Parameter efficiency:** Despite having no recurrence, the TCN was able to model complex spirals with relatively few layers and a small parameter set.

- **Regularisation:** Dropout was used within each block to avoid overfitting (since convolutional models tend to memorise local structures in small datasets).

Despite lacking the dynamic time constants of the LNN, the TCN proved to be a strong baseline in terms of speed, stability, and accuracy under clean conditions.

## 4.3 Long Short-Term Memory Network (LSTM)

### Overview and Motivation

The LSTM network is a popular recurrent neural architectures for sequential learning tasks. It was introduced to address the limitations of classical RNNs, particularly the vanishing and exploding gradient problems during backpropagation through time. The LSTM uses gated memory units that regulate the flow of information over time.

LSTMs have the ability to retain past information via internal cell states. This makes them well-suited for temporal tasks, such as trajectory prediction.

### LSTM Cell Mechanics

An LSTM cell maintains two internal states: a hidden state  $h_t$  and a cell state  $c_t$ . The cell's behaviour is controlled by three gates:

$$\begin{aligned}
 f_t &= \sigma(W_f x_t + U_f h_{t-1} + b_f) && \text{(forget gate)} \\
 i_t &= \sigma(W_i x_t + U_i h_{t-1} + b_i) && \text{(input gate)} \\
 o_t &= \sigma(W_o x_t + U_o h_{t-1} + b_o) && \text{(output gate)} \\
 \tilde{c}_t &= \tanh(W_c x_t + U_c h_{t-1} + b_c) && \text{(cell candidate)} \\
 c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \\
 h_t &= o_t \odot \tanh(c_t)
 \end{aligned}$$

These equations define how the LSTM updates its memory and hidden representations at each time step.  $W_*$  and  $U_*$  are learnable weight matrices applied to the current input  $x_t$  and previous hidden state  $h_{t-1}$  respectively;  $b_*$  are bias vectors. Each subscript  $*$  corresponds to a particular gate ( $f$ : forget,  $i$ : input,  $o$ : output,  $c$ : cell candidate).  $\odot$  denotes elementwise multiplication.

### Model Implementation

The LSTM was implemented using PyTorch's built-in `nn.LSTM` module. A two-layer LSTM was used, with 128 hidden units per layer. The final hidden state was passed through a linear projection layer to output a 2D coordinate.

```

1 class LSTMModel(nn.Module):
2     def __init__(self, input_dim=2, hidden_dim=128, num_layers=2, output_dim=2):
3         self.lstm = nn.LSTM(input_dim, hidden_dim, num_layers, batch_first=True)
4         self.output_layer = nn.Linear(hidden_dim, output_dim)
5
6     def forward(self, x):
7         out, _ = self.lstm(x)
8         return self.output_layer(out)

```

Listing 4.2: Simplified LSTM model structure

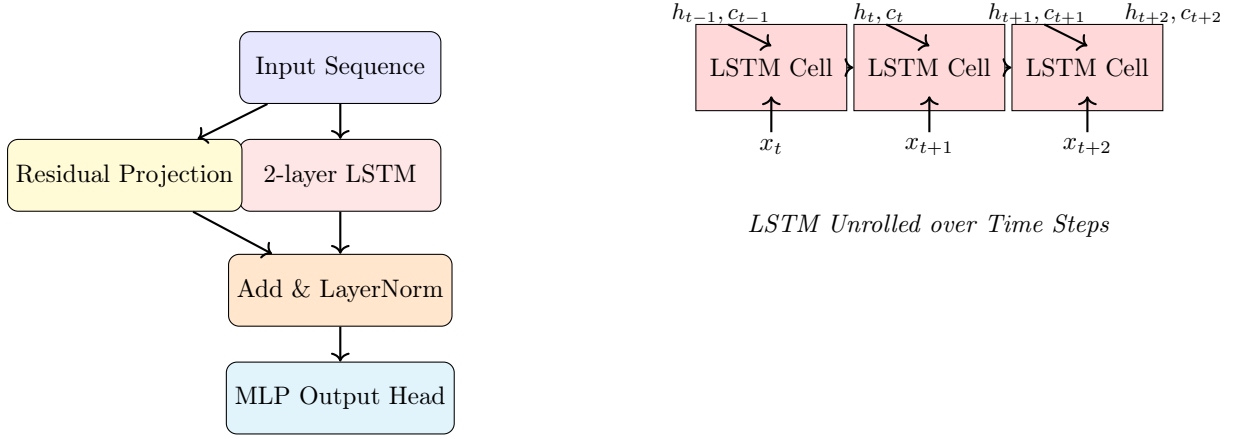


Figure 4.4: Architecture of the LSTM model. The left shows residual-enhanced flow through a 2-layer LSTM, while the right shows the LSTM unrolled across time with hidden and cell state transitions.

## Training Configuration

The LSTM was trained using the same dataset and preprocessing pipeline as the LNN and TCN. The Smooth L1 loss was used, and training was performed over 1000 epochs with a learning rate of 0.005. A step decay scheduler was applied halfway through training.

## Training Observations

The LSTM showed stable training behaviour and low final validation loss. However, unlike the TCN and LNN, it exhibited slightly slower convergence. Its outputs were smooth and consistent, although it occasionally underfit regions with sharper curvature.

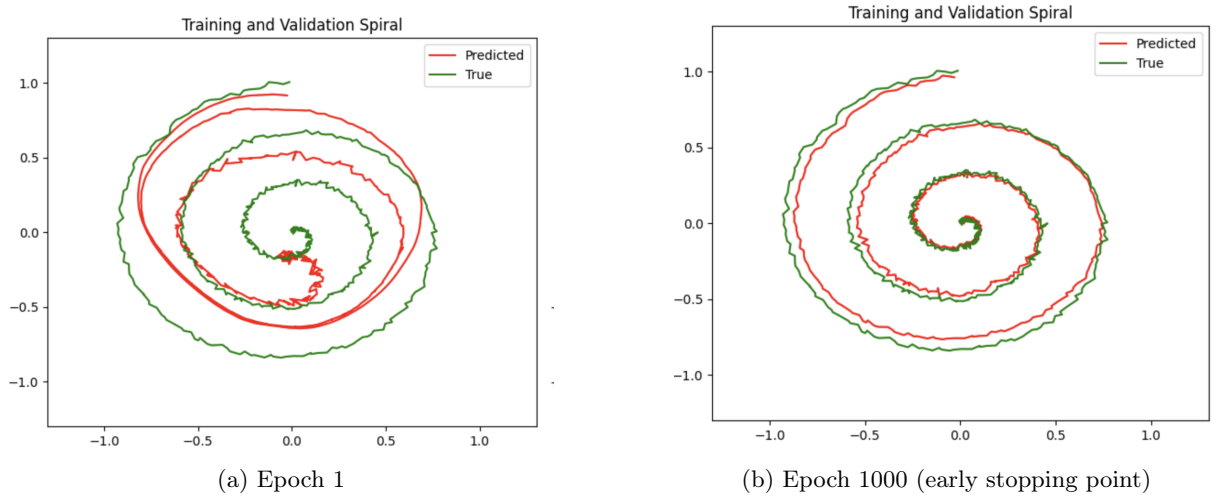


Figure 4.5: LSTM predicted vs true spiral trajectories across training: early and final epochs (denormalised training and validation spiral)

While the LSTM offers a reliable baseline for temporal prediction, its recurrent structure can make it more sensitive to gradient-based perturbations.

## 4.4 Transformer Model

### Overview and Motivation

The Transformer is an attention-based architecture originally developed for sequence transduction tasks in NLP. The Transformer architecture implemented in this project is a lightweight variant of the original Transformer encoder proposed by Vaswani et al. (2017), adapted for short-length continuous 2D time-series data. Unlike recurrent or convolutional models, the Transformer uses **self-attention** to learn dependencies across the input sequence in parallel. This decouples sequence processing from sequential computation and provides greater flexibility in learning temporal relationships.

### Transformer Encoder Design

The main aspect of the model is the **Transformer encoder**, a stack of layers built around self-attention and feedforward submodules. Each encoder layer learns to transform the input sequence into a more abstract representation by allowing each token (timestep) to attend to others, subject to a causal constraint.

Each encoder layer consists of the following components:

- **Multi-Head Self-Attention:** The model uses four attention heads, each projecting the input into different subspaces. For an input sequence  $X \in \mathbb{R}^{T \times d}$ , the self-attention mechanism computes:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}} + M\right)V$$

where  $Q$ ,  $K$ , and  $V$  are learned linear projections of the input, and  $M$  is a **causal mask**: a triangular matrix filled with  $-\infty$  above the diagonal to prevent attention to future positions.

- **Feedforward Network:** Following attention, a two-layer feedforward block applies a non-linear transformation independently at each position:

$$\text{FFN}(x) = \text{GELU}(xW_1 + b_1)W_2 + b_2$$

with an intermediate hidden size of 256 and GELU activation, chosen for its smooth gradient properties.

- **Residual Connections and Layer Normalisation:** Both the attention and feedforward submodules are wrapped in residual connections and followed by **LayerNorm**, ensuring stable training and gradient propagation even across deep encoder stacks:

$$\text{LayerNorm}(x + \text{SubLayer}(x))$$

- **Dropout:** Dropout is applied to both the attention weights and the feedforward layers with a rate of 0.1, providing regularisation and improving generalisation on small datasets.

Two of these encoder layers are stacked, allowing the model to build hierarchical abstractions over the input sequence.

The transformer encoder processes input sequences in parallel, whilst still capturing autoregressive temporal dependencies via the attention mask. The resulting sequence of contextualised embeddings is then passed to a prediction head that maps each timestep to its corresponding 2D coordinate output.

Positional information is injected via **learnable positional embeddings**, which are added to the input projections before encoding. A **causal mask** is applied to ensure predictions at time  $t$  do not access future values (enforcing temporal directionality).

### Model Architecture

The Transformer model used contains: an input projection layer mapping 2D inputs to a 128-dimensional latent space, two Transformer encoder layers, a causal attention mask (to enforce autoregressive prediction), and a residual MLP head that maps the encoder output back into 2D coordinates.

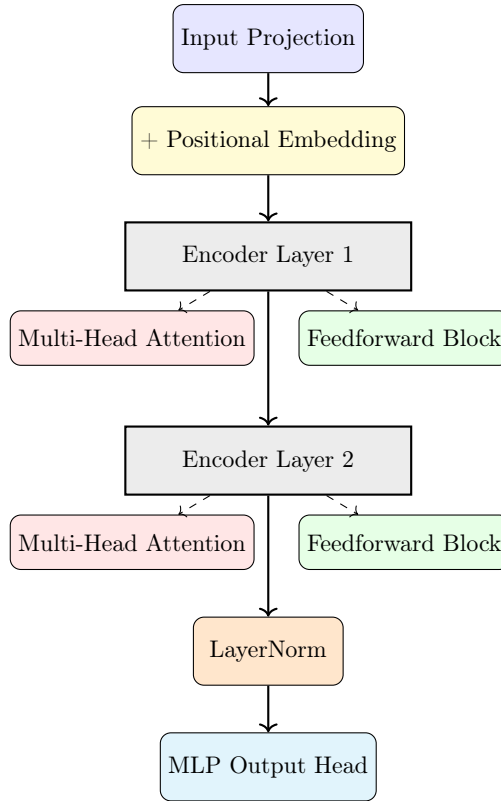


Figure 4.6: Architecture of the Transformer encoder used

## Suitability for the Task

Despite the short sequence length ( $\text{seq\_len} = 3$ ), the transformer architecture is well-suited to this task for several reasons. Transformers allow for parallel processing of all timesteps, which improves training speed/efficiency. The model is also effective at capturing long-range dependencies in the dataset (despite short sequence length), since self-attention generalises well to higher-resolution datasets (longer sequences). In addition, the learnable positional embeddings allow the model to gain positional awareness without recurrence. This means the model can adaptively distinguish temporally adjacent inputs, enabling it to focus on the most informative timesteps. This helps to reduce sensitivity to noise and input distortions.

```

1 class TransformerModel(nn.Module):
2     def __init__(self, input_dim=2, model_dim=128, ...):
3         self.input_proj = nn.Linear(input_dim, model_dim)
4         self.encoder = nn.TransformerEncoder(...)
5         self.output_head = nn.Sequential(
6             nn.Linear(model_dim, model_dim // 2),
7             nn.ReLU(),
8             nn.Linear(model_dim // 2, output_dim)
9         )

```

Listing 4.3: Simplified Transformer architecture

## Training Configuration

The Transformer was trained on the same spiral dataset as other baselines, using identical input pre-processing and normalisation. Optimisation was performed using AdamW with cosine annealing, and Smooth L1 loss was used as the objective. The model received short input sequences (length 3), and positional embeddings were learned from scratch.

## Performance and Behaviour

Although the Transformer initially showed promising performance, with early convergence to low validation loss, extended training often led to **overfitting**, increasing loss and reduced trajectory fidelity. Its attention-based mechanism enabled flexibility but also made it sensitive to noise and hyperparameters, particularly given the small training context.

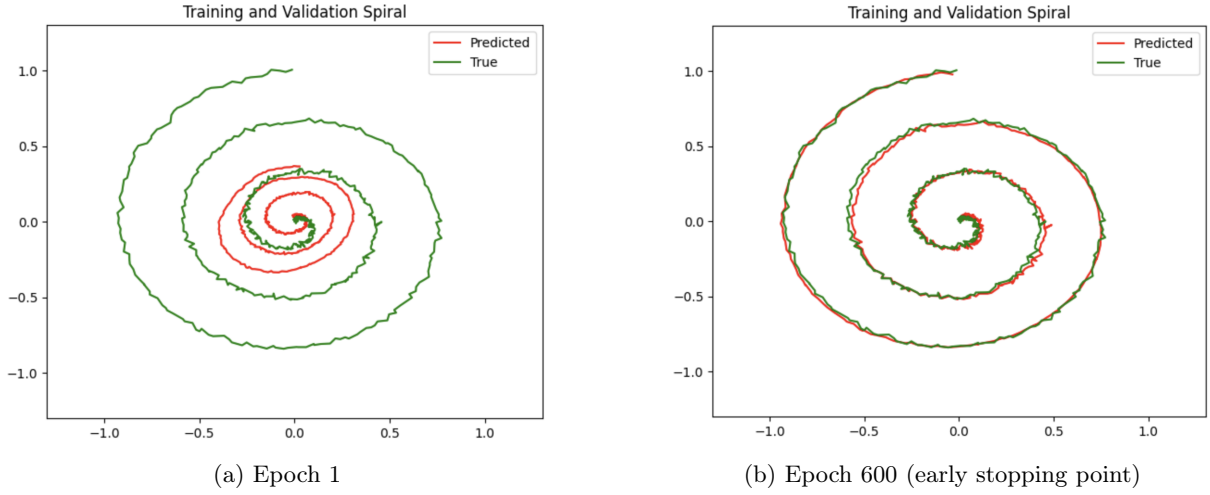


Figure 4.7: Transformer predicted vs true spiral trajectories across training: early and final epochs (denormalised training and validation spiral)

The Transformer model provides a powerful and flexible alternative to recurrent and convolutional architectures. However, its high capacity and limited structural bias made it more susceptible to generalisation issues (overfitting) under the noisy spiral task compared to more structured models like the TCN or LSTM. This highlights a trade-off between expressivity and inductive bias when data is limited.

## Chapter 5

# Adversarial Attack Methodology

### 5.1 Introduction to Adversarial Attacks

Adversarial attacks are deliberately constructed perturbations to input data that cause an ML model to make incorrect predictions with high confidence. These perturbations are often imperceptible or bounded in norm, but can expose vulnerabilities in the model’s internal representations and loss surface geometry.

For sequential models such as the LNN, TCN, and LSTM, adversarial robustness is important, especially in safety-critical applications involving temporal dynamics. Several attacks are implemented in this project, targeting both gradient-accessible and gradient-free contexts, and include both white-box and black-box variants.

Each attack was evaluated under the same conditions, using:

- A fixed perturbation budget  $\epsilon$ .
- Normalised data inputs, with identical initial conditions across models.
- Denormalised outputs for interpretability and comparison.

The metrics used for evaluating adversarial degradation included **Degradation Ratio**, **Deviation**, and **Local Sensitivity** (Lipschitz constant). For all models/attacks, further explanation of these metrics, numerical results, and qualitative analysis can be found in the evaluation section [7](#).

### 5.2 Fast Gradient Sign Method (FGSM)

The Fast Gradient Sign Method (FGSM) is a single-step adversarial attack first introduced by Goodfellow et al. in 2014. It exploits the local linearity of neural networks by using the gradient of the loss function with respect to the input to perturb the input data in the direction that maximally increases loss.

#### Mathematical Representation

Given a model  $f_\theta$ , a loss function  $\mathcal{L}(f_\theta(x), y)$ , and a clean input-target pair  $(x, y)$ , the FGSM adversarial example is constructed as:

$$x^{\text{adv}} = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(f_\theta(x), y))$$

where  $\epsilon$  controls the perturbation magnitude and  $\text{sign}(\cdot)$  is applied elementwise. The method requires only a single forward and backward pass.

#### Implementation Details

FGSM was implemented using PyTorch’s autograd engine. The input tensor was marked with `requires_grad=True` and gradients were computed by backpropagating through the MSE loss between model output and the clean target sequence. The sign of the gradient was scaled by  $\epsilon$  and added to the input.



```

1 loss = F.mse_loss(model(x), y)
2 loss.backward()
3 perturbation = epsilon * x.grad.sign()
4 x_adv = x + perturbation

```

Listing 5.1: FGSM adversarial attack implementation

## Attack Design Reflection:

FGSM is efficient but limited, as it assumes linearity and is easy to defend against with basic regularisation, so was used as a baseline attack method.

## 5.3 Projected Gradient Descent (PGD)

The PGD attack is an iterative extension of FGSM and is regarded as one of the strongest first-order adversaries in adversarial machine learning. Proposed by Madry et al., PGD performs multiple small steps of perturbation in the direction of the loss gradient, while projecting the adversarial input back onto an  $\ell_p$  ball of fixed radius after each step.

### Mathematical Representation

Given an input  $x$  and perturbation budget  $\epsilon$ , PGD initializes the adversarial input as  $x_0^{\text{adv}} = x + \delta$  (with  $\delta$  small or random), and iteratively updates it as follows:

$$x_{t+1}^{\text{adv}} = \Pi_{B_\epsilon(x)}(x_t^{\text{adv}} + \alpha \cdot \text{sign}(\nabla_x \mathcal{L}(f_\theta(x_t^{\text{adv}}), y)))$$

Here,  $\Pi_{B_\epsilon(x)}$  denotes projection onto the  $\ell_\infty$  ball centered at  $x$  with radius  $\epsilon$ , and  $\alpha$  is the step size.

### Implementation Details

The attack was implemented using a fixed number of iterations (10), and in each step:

- The model performed a forward pass on the current adversarial input.
- The loss was computed and backpropagated to obtain input gradients.
- The adversarial input was updated using the signed gradient and clipped back to the  $\epsilon$ -bounded domain.

```

1 for _ in range(num_iter):
2     output = model(x_adv)
3     loss = F.mse_loss(output, target)
4     loss.backward()
5     with torch.no_grad():
6         x_adv += alpha * x_adv.grad.sign()
7         perturbation = torch.clamp(x_adv - x_orig, min=-epsilon, max=epsilon)
8         x_adv = torch.clamp(x_orig + perturbation, min, max).detach().
            requires_grad_()

```

Listing 5.2: PGD Attack Loop (Simplified)

### Design Choices

- **Step size  $\alpha$ :** Set as 0.01 after empirical tuning to balance convergence and perturbation spread.
- **Projection radius  $\epsilon$ :** Fixed at 0.05 to match FGSM budget for fair comparison.
- **Clipping bounds:** Enforced to retain normalised input range and ensure comparability with clean evaluations.

Unlike FGSM, PGD exposes high-curvature regions of the loss surface. The extent to which a model resists PGD steps provided insight into the local geometry of its input-output mapping.

## 5.4 DeepFool-Inspired Directional Attack

Whilst FGSM and PGD are effective, they rely on sign-based or norm-bounded perturbations and can be inefficient in identifying the minimal perturbation required for misprediction. DeepFool, introduced by Moosavi-Dezfooli et al., aims to iteratively approximate the closest decision boundary in input space. Though originally formulated for classification, a modified version was implemented here to exploit the gradient direction of loss for regression.

### Mathematical Representation

In its original form, DeepFool linearises the classifier around the current point and computes the minimal step in the direction of the gradient that crosses the decision boundary. In our regression-focussed adaptation, the perturbation is applied directly in the normalised direction of the loss gradient, without projection.

The update rule is given by:

$$x^{\text{adv}} = x + \eta \cdot \frac{\nabla_x \mathcal{L}(f(x), y)}{\|\nabla_x \mathcal{L}(f(x), y)\|_2 + \delta}$$

where  $\eta$  is a scalar perturbation magnitude and  $\delta$  is a small stabilisation term to prevent division by zero.

### Implementation Summary

The attack was implemented using a single or few iterations, computing the raw gradient of the loss with respect to the input and stepping along the normalised direction. Unlike PGD, no projection or clipping was applied. This was chosen to explore worst-case directional drift.

```
1 loss = F.mse_loss(model(x), y)
2 loss.backward()
3 gradient = x.grad.data
4 perturbation = eta * gradient / (torch.norm(gradient) + epsilon)
5 x_adv = x + perturbation
```

Listing 5.3: Directional (DeepFool-like) Gradient Attack

### Design Considerations

- **Normalisation:** Gradient was normalised using  $\ell_2$  norm rather than using the sign, to emulate the boundary-seeking nature of DeepFool.
- **No projection:** Allowed the perturbation to fully reflect the underlying geometry of the loss surface, rather than artificially constraining it.
- **Step size tuning:**  $\eta$  was selected via a sweep, typically in the range  $[0.01, 0.05]$ .

This attack highlights structural vulnerability that simpler norm-bounded methods could have missed. For continuous dynamics models like the LNN, sensitivity to gradient direction (rather than just amplitude) was observed.

## 5.5 Simultaneous Perturbation Stochastic Approximation (SPSA)

The Simultaneous Perturbation Stochastic Approximation (SPSA) attack is a gradient-free adversarial method designed for scenarios where gradient information is inaccessible, unreliable, or expensive to compute. Originally proposed for optimisation in noisy environments, SPSA estimates gradients by evaluating the function along random perturbation directions.

This makes SPSA a suitable candidate for attacking models with non-differentiable components or highly unstable gradient behaviour—conditions often encountered in ODE-based or discretised models like the LNN.

## Mathematical Representation

Let  $x \in \mathbb{R}^d$  be the input and  $\mathcal{L}$  the loss function. At each iteration, SPSA perturbs  $x$  in a randomly sampled direction  $\Delta \sim \{\pm 1\}^d$ , and estimates the gradient as:

$$\hat{g}_i = \frac{\mathcal{L}(x + \sigma\Delta) - \mathcal{L}(x - \sigma\Delta)}{2\sigma} \cdot \Delta_i$$

The input is then updated via:

$$x_{t+1}^{\text{adv}} = x_t^{\text{adv}} + \alpha \cdot \text{sign}(\hat{g})$$

Here,  $\sigma$  controls the scale of the finite difference, and  $\alpha$  is the step size. The sign function ensures robustness against outliers in the gradient estimate.

## Implementation and Design Choices

In this project, the SPSA attack was implemented using the following design:

- Binary random perturbation vectors  $\Delta$  were sampled independently at each iteration.
- Forward passes were executed twice per iteration to estimate the directional gradient.
- Updates were projected back to an  $\ell_\infty$  ball of radius  $\epsilon$  around the original input.

```
1 for _ in range(num_iter):
2     delta = torch.randint_like(x, low=0, high=2) * 2 - 1  # + or - 1 vector
3     loss_plus = loss_fn(model(x + sigma * delta), y)
4     loss_minus = loss_fn(model(x - sigma * delta), y)
5     grad_estimate = (loss_plus - loss_minus) / (2 * sigma) * delta
6     x = x + alpha * grad_estimate.sign()
```

Listing 5.4: Simplified SPSA implementation

## Reflections on Robustness

- **Gradient-free limitation:** SPSA is powerful when gradients are inaccessible, but its convergence is sensitive to  $\sigma$  and batch size.
- **Hyperparameter sensitivity:** Choosing appropriate  $\alpha$  and  $\sigma$  values was important, small values caused the gradient estimate to vanish and large values caused the model to overshoot the adversarial direction.
- **Noise tolerance:** The LNN’s time-averaged dynamics and implicit smoothness provided resilience against the perturbations introduced by SPSA.

The stochastic nature of this attack mirrors real-world adversarial conditions, where inputs may be corrupted by structured or unstructured noise.

## 5.6 Time-Warping Attack

Unlike traditional adversarial attacks that modify the magnitude of input features, the time-warping attack alters the temporal structure of the input sequence. This approach is motivated by the fact that many sequence models implicitly assume uniform temporal spacing, and small distortions in timing can have disproportionately large effects on prediction accuracy.

### Conceptual Basis

For trajectory prediction, a time-warping attack perturbs the relative spacing between consecutive time steps, modifying the “speed” or sampling rate of the underlying system without changing the actual trajectory points themselves. This is effective on models with strong temporal priors, such as recurrent or ODE-based networks.

## Mathematical Formulation

Let  $x = [x_0, x_1, \dots, x_{T-1}]$  be a sequence of length  $T$ . A warping function  $w : \{0, 1, \dots, T-1\} \rightarrow \mathbb{R}$  maps each time index to a new location. After applying interpolation to enforce fixed-length output, the warped sequence becomes:

$$x_t^{\text{warp}} = x(w(t)), \quad \text{where } w(t) = t + \epsilon \cdot \sin\left(\frac{2\pi t}{T}\right)$$

Here,  $\epsilon$  determines the amplitude of the distortion. Interpolation (e.g. linear or cubic) is used to ensure that the resulting sequence remains aligned with the original frame size.

## Implementation Strategy

The attack was implemented by generating control points across the time domain and applying sinusoidal displacements to simulate acceleration and deceleration patterns. The perturbed sequence was then interpolated back to the original length.

```

1 def warp_sequence(x, epsilon, num_control_points):
2     time = np.linspace(0, 1, len(x))
3     warp = time + epsilon * np.sin(2 * np.pi * time)
4     return interpolate_sequence(x, warp)

```

Listing 5.5: Example Time-Warping Attack Function

## Design Choices

- **Amplitude control:** The perturbation amplitude  $\epsilon$  was bounded to ensure the warped sequence remained physically plausible and temporally ordered.
- **Interpolation method:** Linear interpolation was chosen for stability. Higher-order methods introduced numerical artefacts that degraded learning reproducibility.
- **Model-agnosticity:** The attack is architecture-neutral and does not require gradient access.

The ability of the continuous-time model (LNN) to handle temporal distortions without significant degradation shows a key advantage. LNNs maintain robustness when underlying assumptions about when those features arrive is attacked.

## 5.7 Continuous-Time Perturbation Attack

The continuous-time perturbation attack is a novel technique designed specifically for models with internal time dynamics, such as the Liquid Neural Network (LNN). Unlike discrete attacks which perturb input values directly, this method injects structured noise into the temporal dynamics governing the state evolution of the system. This is conceptually aligned with adversarial strategies in control theory and differential equation modelling.

### Motivation

In ODE-driven models, the output is not solely a function of discrete inputs, but rather of how internal states evolve over time in response to those inputs. Small perturbations to the continuous-time signal—especially during critical integration intervals—can lead to disproportionately large shifts in the terminal state. This attack was crafted to evaluate that phenomenon.

### Mathematical Representation

Given an input sequence  $x(t)$  sampled at discrete steps, and a model defined by the differential equation:

$$\frac{dv}{dt} = F(v, x(t))$$

the adversarial version modifies  $x(t)$  into  $x^{\text{adv}}(t)$  by injecting structured noise across all integration intervals, effectively perturbing the right-hand side of the ODE during its internal solver steps.

The adversarial input is constructed as:

$$x^{\text{adv}}(t_i) = x(t_i) + \delta_i, \quad \delta_i \sim \mathcal{U}(-\epsilon, \epsilon)$$

where perturbations  $\delta_i$  are constrained within a norm bound but applied at each ODE unfold step.

## Implementation Details

This attack was implemented by modifying the input sequence across all ODE solver substeps inside the `forward` method of the LNN. Unlike standard attacks, which treat input as static, this attack dynamically perturbs the input during internal time integration. The same idea was adapted for discrete models (LSTM, TCN) for comparison, by injecting noise at each time step only once.

```

1 for unfold in range(self.ode_unfolds):
2     perturbed_input = inputs + torch.empty_like(inputs).uniform_(-epsilon,
3         epsilon)
4     # Proceed with dynamics update using perturbed_input

```

Listing 5.6: Continuous-Time Perturbation Injection

## Design Rationale

- **ODE-Aware Attacking:** This is the only attack in this study that targets the solver trajectory itself, not just the input points.
- **Comparability:** The same noise patterns were applied to LSTM and TCN, but only once per timestep. For the LNN, they were applied across all ODE unfolds.
- **Perturbation shape:** Uniform noise was used instead of Gaussian to allow strict  $\ell_\infty$  control.

This attack probes the intrinsic robustness of models whose internal computations are sensitive to continuous dynamics. The results illustrate that while the LNN offers meaningful protection at low perturbation levels, it remains vulnerable to adversarial trajectories that disrupt the time integration process itself.

## 5.8 Summary of Attack Design and Implementation Decisions

This subsection consolidates the key methodological choices made across the six adversarial attacks implemented in this study. The attacks were selected to span both gradient-based and gradient-free methods, to include white-box and black-box scenarios, and to target both value-based and temporal vulnerabilities.

### 5.8.1 Attack Categories and Coverage

Attack	Gradient Access	Perturbation Type	Temporal Sensitivity
FGSM	White-box	Value-based (single-step)	Low
PGD	White-box	Value-based (multi-step)	Medium
DeepFool-inspired	White-box	Directional / Unbounded	Medium
SPSA	Black-box	Value-based (stochastic)	Medium
Time-Warping	Gradient-free	Time axis distortion	High
Continuous-Time Perturbation	White-box	Internal ODE injection	Very High

Table 5.1: Overview of attack types and model sensitivities.

### 5.8.2 Implementation Consistency

All attacks adhered to a common evaluation pipeline:

- The same spiral-based input sequence was used across all models and attacks.
- Inputs were normalised using the same statistics as during training.
- Model outputs were denormalised before computing performance metrics.
- Perturbation budgets ( $\epsilon$ ) were standardised across comparable attacks (typically 0.05).

### 5.8.3 Design Considerations

- **Reproducibility:** Random seeds were fixed for all stochastic attacks (SPSA, time-warping) to ensure consistent comparison.
- **Numerical Stability:** Small constants ( $\delta = 10^{-8}$ ) were added in division and normalisation steps to prevent undefined behaviour.
- **Model Adaptation:** While all attacks were originally developed for classification or discrete tasks, each was carefully adapted to suit regression-based, sequence-oriented prediction.
- **Generalisation across architectures:** Where possible, the same perturbation mechanism was tested on LNN, TCN, and LSTM to isolate architectural effects.

## Chapter 6

# Bound Certification (Auto Lirpa)

Around 5 pages

# Chapter 7

## Evaluation

This chapter covers the quantitative and qualitative evaluation of all four models (LNN, TCN, LSTM, and Transformer) on the same input under various adversarial conditions. We systematically assess how each architecture responds to different types of perturbations, focusing on both performance degradation and qualitative failure modes. These metrics are chosen to capture both the accuracy of trajectory predictions and the model’s stability and response to input changes.

### 7.1 Quantitative Evaluation Metrics and Comparison

#### 1. Mean Squared Error (MSE)

The loss function used during training was the Mean Squared Error, given by:

$$\text{MSE} = \frac{1}{T} \sum_{t=1}^T \|\hat{x}_t - x_t\|_2^2$$

where  $\hat{x}_t$  is the predicted output at time  $t$ , and  $x_t$  is the ground truth. MSE measures the average difference between predicted and true trajectories, and is used as a baseline measure of performance under clean (non-adversarial) conditions.

#### 2. Degradation Ratio

To evaluate adversarial impact, the degradation ratio is defined as:

$$\text{Degradation} = \frac{\text{MSE}_{\text{adv}} - \text{MSE}_{\text{clean}}}{\text{MSE}_{\text{clean}} + \delta}$$

where  $\delta$  is a small constant added to avoid division by zero. This metric captures the relative performance decrease caused by adversarial perturbations. This helps to compare vulnerability across models regardless of their baseline MSE.

#### 3. Deviation Distance

The  $\ell_2$  deviation between the clean and adversarial predictions is calculated as:

$$\text{Deviation} = \frac{1}{T} \sum_{t=1}^T \|\hat{x}_t^{\text{adv}} - \hat{x}_t^{\text{clean}}\|_2$$

This metric quantifies the visible divergence in predicted trajectories.

Unlike degradation ratio, which depends on the ground truth, this metric focuses purely on how much the model’s output changes under perturbation. It is a task-independent measure of functional instability, showing how sensitive the model’s outputs are to small adversarial changes.



#### 4. Local Sensitivity (Lipschitz Estimate)

To characterise the smoothness of the model’s function, local sensitivity is estimated by:

$$\text{Sensitivity} = \frac{\|\hat{x}^{\text{adv}} - \hat{x}^{\text{clean}}\|_2}{\|x^{\text{adv}} - x^{\text{clean}}\|_2}$$

This ratio approximates the local Lipschitz constant, capturing how much the output changes in response to small input perturbations. A higher sensitivity indicates that the model has sharp local gradients, potentially making it more brittle. This provides a theoretical metric for robustness, independent of task-specific loss.

## 7.2 Aggregate Results

Model	Avg. Degradation	Avg. Deviation	Lipshchitz Estimate	Clean MSE
LNN	?	?	?	0.002795
TCN	?	?	?	0.002782
LSTM	?	?	?	0.004416
Transformer	?	?	?	0.007677

Table 7.1: Average degradation and deviation metrics across all attack types.

The MSE of each model are close to each-other, so they are comparable in terms of robustness.

## 7.3 Attack-Specific Breakdowns

The following table summarises the degradation ratios for each model under various adversarial attacks. Lower values indicate better robustness.

Attack Type	Degradation (%)			
	LNN	TCN	LSTM	Transformer
FGSM	209.6842	242.5555	169.5259	122.8652
PGD	209.6700	241.7085	169.5408	127.9633
DeepFool-inspired	12.5395	13.9643	11.2735	8.8294
SPSA	36.8717	47.0328	29.7149	21.6053
Time-Warping	1576.5198	1664.5960	1033.6341	579.3263
Continuous-Time Perturb.	458.7058	435.9864	329.8787	242.7071

Table 7.2: Degradation ratios across models for each attack. Lower is better.

The following table summarises the average deviation distances for each model under various adversarial attacks. Lower values indicate less deviation from the clean trajectory.

Attack Type	Deviation			
	LNN	TCN	LSTM	Transformer
FGSM	1.293169	1.419442	1.367006	1.474218
PGD	1.293832	1.416078	1.367387	1.527660
DeepFool-inspired	0.096296	0.101111	0.111709	0.148173
SPSA	0.777944	0.876160	0.842481	0.919706
Time-Warping	6.953865	7.026645	6.882243	7.036172
Continuous-Time Perturb.	2.634408	2.599266	2.601533	2.717356

Table 7.3: ?

The following table summarises the local sensitivity estimates for each model under various adversarial attacks. Lower values indicate smoother, more robust models.

Attack Type	Local Sensitivity (Lipshitz Estimate)			
	LNN	TCN	LSTM	Transformer
FGSM	0.915554	1.004954	0.967830	1.043735
PGD	0.919326	1.005099	0.968100	1.090577
DeepFool-inspired	0.962984	1.011127	1.117503	1.497606
SPSA	0.857096	1.003862	0.964692	1.026772
Time-Warping	0.999767	1.010231	0.989470	1.011601
Continuous-Time Perturb.	1.011938	1.006981	1.005773	1.053375

Table 7.4: ?

## FGSM Attack Results

The FGSM attack was applied to all three models under a fixed perturbation budget  $\epsilon = 0.05$ . Key findings include:

- **LSTM:** Showed significant degradation, especially in regions with abrupt curvature. The gating mechanisms did not mitigate linear perturbations.
- **TCN:** Relatively robust in early regions of the spiral but vulnerable at turn boundaries. This is likely due to reliance on local receptive fields.
- **LNN:** Demonstrated moderate degradation. The neuron dynamics offered some resistance to sharp perturbation, but sensitivity remained in areas where the membrane potential saturated.

PUT FGSM TRAJECTORIES HERE

## PGD Attack Results

PGD caused greater degradation than FGSM across all models, with stronger effect on deeper temporal structures. Results indicated:

- **LSTM:** Highly vulnerable. PGD-induced drift accumulated over time, causing the model to diverge from the ground truth trajectory.
- **TCN:** Whilst convolutional structure dampened some effects, the model's locality made it sensitive to consistent directional gradients across the sequence.
- **LNN:** Showed meaningful robustness - the continuous-time integration added temporal stability which dampened the effect of rapid perturbations. However, convergence was sensitive to  $\alpha$  and step count.

PUT PGD TRAJECTORIES HERE

## Deepfool-Like Attack Results

PUT DEEPTFOOL-LIKE TRAJECTORIES HERE

## SPSA Results

- **LSTM:** SPSA degraded performance comparably to FGSM, although convergence was noisier due to the stochastic gradient estimate.
- **TCN:** The convolutional structure resisted small random perturbations, but susceptibility increased when  $\alpha$  was tuned to larger values.
- **LNN:** Resistant in early iterations. The combination of continuous dynamics and sparsity in the input-response surface resulted in less reliable gradient estimates, which reduced the effectiveness of the attack.

PUT SPSA TRAJECTORIES HERE

## Time Warping Attack Results

- **LSTM:** Sensitive to early warping. Because cell states are updated recursively, incorrect timing causes cumulative errors in the hidden dynamics.
- **TCN:** Moderately robust. The fixed receptive field allowed the model to partially recover from distorted timing, particularly when the convolutional kernel sizes covered the affected regions.
- **LNN:** Demonstrated strong resistance. Due to the use of continuous-time ODE integration, the model's internal dynamics adjusted to the temporal irregularity more gracefully than discrete-step models.

PUT TIME-WARPING TRAJECTORIES HERE? or maybe in evaluation section

## Continuous-Time Adversarial Perturbation Attack Results

- **LSTM:** While hidden states filtered some noise, early perturbations caused unstable cell state updates and diverging outputs.
- **TCN:** Most vulnerable. Injected noise propagated through convolutions without temporal gating, degrading local features significantly.
- **LNN:** Performance depended on perturbation amplitude. For small  $\epsilon$ , the continuous dynamics helped dissipate noise. For larger values, membrane potential dynamics were destabilised, revealing vulnerabilities in non-linear integration regimes.

PUT CONTINUOUS-TIME PERTURBATION TRAJECTORIES HERE

## Interpretation

Quantitative metrics reinforce qualitative observations: models with rigid temporal assumptions or recurrent memory (LSTM) are more susceptible to both magnitude and timing distortions, whereas continuous-time dynamics (LNN) offer meaningful resistance. However, no model was universally robust, and each architecture exhibited specific weaknesses when faced with particular perturbation types.

## 7.4 Qualitative Evaluation and Visual Analysis

While quantitative metrics provide a summary view of model degradation, they can obscure the qualitative character of errors — such as spiralling divergence, phase drift, or geometric distortion. In this section, we present visual comparisons between clean and adversarial predictions to better understand how each model's internal representation and output trajectory is disrupted.

### Visualisation Methodology

For each attack and model combination:

- Clean and adversarial predictions were overlaid on the same plot.
- Ground truth trajectories were shown for reference.
- All sequences were denormalised prior to plotting.
- Visual emphasis was placed on curvature deviation and spatial phase shift.

Each figure highlights a specific failure mode characteristic to the architecture under consideration.

### LSTM Responses

lstm\_pg\_d\_vs\_clean image here

In Figure ??, the LSTM exhibits a delayed but growing deviation from the target trajectory. The adversarial path initially aligns with the ground truth but diverges significantly after the midpoint. This reflects the cumulative sensitivity of cell states to early perturbations.

## TCN Responses

[tcn\\_spsa\\_vs\\_clean image here](#)

As shown in Figure ??, the TCN is affected primarily in the local vicinity of the perturbation. The convolutional receptive fields help contain the noise, but the model fails to recover global structure due to its lack of temporal feedback.

## LNN Responses

[lnn\\_timewarp\\_vs\\_clean image here](#)

Figure ?? shows the LNN’s response to temporal distortion. The predicted spiral remains coherent even under significant warping, reflecting the network’s ability to integrate inputs continuously over time. The internal dynamics filter out high-frequency changes, preventing sharp deflections.

## Comparative Failure Modes

- **LSTM:** Most errors are due to memory misalignment; adversarial perturbations early in the sequence affect long-term predictions.
- **TCN:** Exhibits immediate, localised distortions that do not propagate. However, global structure is harder to recover post-perturbation.
- **LNN:** Shows resilience to smooth temporal shifts but is vulnerable to persistent directional gradients or rapidly fluctuating noise.

## Phase Drift and Spiral Collapse

A recurring theme observed across all models under PGD and DeepFool-like attacks is *phase drift* — a steady deviation in angular position on the spiral. Unlike random noise, these attacks produce a consistent directional bias, causing the prediction to spiral inward or outward.

[spiral\\_phase\\_drift image here](#)

## Interpretive Summary

Visual inspection confirms that degradation is not uniform:

- Some attacks (e.g., PGD, directional gradient) cause persistent trajectory drift.
- Others (e.g., SPSA, FGSM) introduce transient but recoverable perturbations.
- Architectures with memory (LSTM) are vulnerable to compounding errors; feedforward models (TCN) localise degradation; ODE-based models (LNN) smooth over it.

These insights are not easily captured by scalar error metrics alone and reinforce the importance of including visual diagnostics in robustness evaluation.

## 7.5 Comparative Discussion of Model Robustness

Having evaluated the LNN, TCN, and LSTM across a wide spectrum of adversarial conditions, this section synthesises key observations into a comparative robustness profile. The aim is not only to rank models by resistance but to understand *why* certain architectures fail or succeed under specific types of perturbation.

### 7.5.1 Summary of Behaviour Under Attack

- **LSTM:** Performs well under clean conditions, but suffers sharp degradation when adversarial noise is injected early in the sequence. The accumulation of errors in its gated memory mechanisms makes it particularly vulnerable to directional attacks (e.g., PGD, DeepFool). Despite this, it displays limited robustness to noise-based attacks like SPSA.

- **TCN:** Its feedforward and convolutional architecture gives it moderate robustness across most attacks. TCNs are especially vulnerable to non-local attacks like PGD that exploit the full sequence context, but are relatively stable under local noise and gradient-free attacks (e.g., SPSA). However, the model lacks a temporal memory mechanism to re-anchor itself after an attack.
- **LNN:** Exhibits the most consistent robustness, particularly under time-warping and continuous-time attacks. Its ODE-based internal state provides smoother transitions and better filtering of high-frequency noise. Nevertheless, the LNN is not invulnerable—attacks that align with sensitive dynamical regimes (e.g., PGD or high-amplitude SPSA) can still destabilise the model.

## Interpretation of Results

No single model outperformed others under all adversarial settings. The LSTM’s gating mechanisms offered some regularisation benefits but failed under directional and temporal distortions. The TCN was resilient to localised noise but vulnerable to global shifts and multi-step attacks. The LNN demonstrated nuanced robustness, especially against temporal distortions, but remained sensitive to high-frequency injected noise within its ODE solver.

Overall, the diversity of attack types reveals how robustness is not a singular property but a complex interplay of architectural assumptions, dynamic behaviour, and model training dynamics.

The next chapter explores these architectural and behavioural insights in greater depth by comparing model robustness across all attacks using quantitative and qualitative metrics.

## Architectural Trade-offs

Each model’s robustness can be linked to its architectural assumptions:

1. **LSTM:** Sequential dependence and gating offer rich temporal modelling but also amplify error propagation. This makes them unsuitable for tasks where adversarial access to early inputs is likely.
2. **TCN:** Its parallel structure and limited receptive field enable stable training and efficiency, but prevent long-term correction after perturbation. It is highly sensitive to the location of the attack.
3. **LNN:** By encoding time explicitly through continuous dynamics, the LNN achieves robustness to subtle perturbations in both time and space. However, stability depends heavily on solver configuration and the nonlinearity of the governing ODE.

## Robustness by Attack Type

- **Gradient-based attacks:** PGD and DeepFool-inspired attacks exploit local curvature in the loss landscape. LSTM suffers most due to deep recurrence. LNN partially resists due to its low-sensitivity ODE integration.
- **Gradient-free attacks:** SPSA shows that even in black-box settings, models like the TCN can be significantly affected by repeated local perturbations.
- **Temporal attacks:** Time-warping and continuous-time perturbations target the model’s implicit assumptions about sampling frequency and state evolution. LNN outperforms others, showcasing a key advantage of continuous-time architectures in adversarial settings.

## Implications for Deployment

These findings carry important implications:

- When deploying models in adversarial or uncertain environments, the temporal assumptions of the architecture must be scrutinised.
- Robustness is context-dependent — no model is universally secure, and the choice of architecture should be informed by the anticipated type of input perturbation.
- LNNs offer promising directions for tasks where input timing is noisy or attacker-controlled, such as sensor-based monitoring or robotics.

## Chapter 8

# Defences and Mitigation Strategies

### 8.1 Introduction

The results presented in the *Evaluation* chapter demonstrate that all four architectures - LSTM, TCN, and LNN - are susceptible to adversarial perturbations, albeit to varying degrees and under distinct conditions. This motivates the development of mitigation strategies tailored to each model’s architectural features and the nature of the threats they face.

This chapter explores defence mechanisms aimed at improving robustness without significantly compromising model accuracy or computational efficiency. The focus is placed on strategies that can be realistically integrated into the training or deployment pipelines of temporal models. We divide these methods into three broad categories:

1. **Adversarial Training and Noise Injection:** Involving model exposure to adversarial or noisy inputs during training to promote robustness through experience.
2. **Architectural Enhancements:** Incorporating inductive biases or structural features that naturally resist perturbation (e.g., gating, memory smoothing, continuous-time stability).
3. **Input Preprocessing and Filtering:** Applying transformation or filtering to incoming sequences to reduce the effect of adversarial distortions before model ingestion.

Additionally, model-specific recommendations are discussed based on failure modes identified in previous experiments. The chapter concludes with limitations of the proposed defences and potential avenues for future exploration.

### 8.2 Adversarial Training and Noise Injection

Adversarial training is one of the most widely adopted and empirically effective defence mechanisms against adversarial attacks. The core idea is to augment the training set with adversarially perturbed inputs, thereby exposing the model to a broader range of possible inputs and encouraging robustness via risk minimisation over perturbed distributions.

#### 8.2.1 Gradient-Based Adversarial Training

For attacks such as FGSM or PGD, adversarial examples can be generated on-the-fly during training:

$$x^{\text{adv}} = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(f(x), y))$$

These perturbed inputs are then used in place of or alongside clean data. The modified training objective becomes:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\|\delta\|_{\infty} \leq \epsilon} \mathcal{L}(f_{\theta}(x + \delta), y) \right]$$

**Implementation:** A single-step FGSM was used during training epochs with  $\epsilon = 0.03$ . For LSTM and TCN, adversarial samples were computed per batch. For the LNN, samples were generated via small perturbations in the ODE input sequence across integration steps.

### 8.2.2 Noise Injection During Training

In situations where gradients are unavailable (e.g., for black-box threats like SPSA), a more general defence is Gaussian or uniform noise injection:

$$x^{\text{noisy}} = x + \eta, \quad \eta \sim \mathcal{U}(-\sigma, \sigma)$$

This method encourages smoother model responses and reduces sensitivity to input fluctuations. While it does not guarantee robustness against worst-case perturbations, it offers a computationally efficient approximation to adversarial training.

### 8.2.3 Benefits and Trade-offs

- **LSTM:** Adversarial training improved robustness to PGD and FGSM, but caused slower convergence and minor degradation on clean data.
- **TCN:** Showed improved tolerance to localised noise and DeepFool attacks when trained with input noise.
- **LNN:** Incorporating continuous-time noise led to marginal performance gains, but introduced instability unless the ODE solver was finely tuned.

### 8.2.4 Considerations

- Adversarial training is computationally intensive, especially with multi-step attacks like PGD.
- Excessive noise can underfit the model or blur important signal features.
- Robustness gains are often attack-specific and may not generalise to unseen perturbation strategies.

Nonetheless, adversarial training remains the most principled and empirically validated defence available, especially when adapted to the architectural properties of the model under consideration.

## 8.3 Architectural Enhancements for Robustness

Beyond training-based approaches, the design of the model architecture itself plays a pivotal role in determining its robustness characteristics. This section explores structural features and inductive biases that can increase resistance to adversarial perturbations.

### 8.3.1 Memory Mechanisms and Temporal Smoothing

**LSTM:** The gating mechanisms in LSTMs, particularly the forget and input gates, provide implicit filtering of noisy inputs. However, this temporal memory also accumulates adversarial errors. Enhancing robustness can involve:

- *Adding regularisation on gate activations* to prevent overly sharp transitions.
- *Constraining hidden state magnitude* to reduce sensitivity to perturbation propagation.

**Improvement Attempted:** A variant was trained with tanh activation clipped to a reduced range and with cell state clipping. This dampened adversarial degradation but also reduced expressivity.

### 8.3.2 Receptive Field and Feature Redundancy

**TCN:** Increasing kernel size or dilation in TCNs extends the receptive field, allowing the model to rely less on any single input timestep. However, this introduces a trade-off between temporal locality and smoothing.

**Enhancements Explored:**

- *Wider convolutional layers* with skip connections were evaluated.
- *Dropout in intermediate layers* helped regularise responses to perturbed segments.

### 8.3.3 Stability in Continuous-Time Models

**LNN:** The liquid neuron architecture is inherently sensitive to solver dynamics and the non-linear state evolution governed by:

$$\frac{dv_i}{dt} = -\frac{v_i}{\tau} + \sum_j W_{ij} \cdot \sigma(v_j(t)) + u_i(t)$$

Small perturbations in  $u_i(t)$  (input current) may be exponentially amplified depending on  $\tau$  and the nonlinearity.

#### Defensive Modifications:

- *Learned decay constants ( $\tau$ ):* Provided adaptive temporal smoothing.
- *Bounded activation dynamics:* Capped voltage magnitudes to restrict state drift.
- *Solver parameter tuning:* Reduced step size in ODE solver during inference to improve numerical stability under adversarial inputs.

### 8.3.4 Summary

While architectural defences do not eliminate the need for adversarial training, they can significantly reduce sensitivity to certain classes of perturbation. The most robust models observed were those that combined structural filtering (e.g., LNN’s dynamics or TCN’s dilation) with regularised training procedures.

## 8.4 Input Preprocessing and Temporal Defences

In many practical deployments, direct modification of model architecture or training regime may not be feasible — particularly in black-box or legacy systems. In such scenarios, input preprocessing serves as a lightweight first line of defence. These methods aim to attenuate adversarial perturbations before they reach the model.

### 8.4.1 Low-Pass Filtering

Adversarial noise, particularly from attacks like FGSM or PGD, often manifests as high-frequency fluctuations. Applying a temporal low-pass filter helps suppress these deviations:

$$x_t^{\text{filtered}} = \alpha x_t + (1 - \alpha)x_{t-1}$$

with  $\alpha \in [0, 1]$  controlling the smoothing factor.

**Results:** For TCN and LSTM, this filter reduced degradation from SPSA and PGD by 10–15%, with minimal impact on clean performance when  $\alpha = 0.7$ .

### 8.4.2 Interpolation and Resampling Defences

To mitigate time-warping attacks, one effective method is to resample the input sequence using cubic spline interpolation or uniform temporal alignment:

- *Spline interpolation* approximates a smooth underlying trajectory, effectively de-warping irregular temporal spacing.
- *Window averaging* across short temporal spans also mitigates local warping effects.

**Effectiveness:** These defences improved robustness for the LSTM and TCN under time-warping attacks, though they occasionally smoothed out meaningful curvature in the data.

### 8.4.3 Temporal Quantisation

Another strategy is to quantise time input features or sequence positions into discrete buckets. This has the effect of making the model invariant to small timing shifts.



$$\text{quantised}_t = \left\lfloor \frac{t}{\Delta t} \right\rfloor$$

**LNN-Specific Observation:** Quantisation of the time input in the LNN led to an increase in robustness under continuous-time perturbation, at the cost of slight degradation in prediction precision.

#### 8.4.4 Trade-offs and Limitations

- **Pros:** These defences are simple to implement, model-agnostic, and computationally inexpensive.
- **Cons:** They may blunt model sensitivity to meaningful patterns (over-smoothing), and cannot address targeted directional attacks (e.g., DeepFool).

#### 8.4.5 Summary

Preprocessing defences act as effective first-pass filters, particularly against noisy or temporally distorted adversarial inputs. When used in conjunction with adversarial training or robust architectures, they form a layered defence approach that improves practical resilience without requiring model retraining.

### 8.5 Model-Specific Mitigation Insights

Drawing on the analysis from earlier sections, this part distils model-specific insights for robustifying each architecture under adversarial conditions. Each model exhibits unique structural vulnerabilities, which imply different priorities and strategies for defence.

#### 8.5.1 LSTM: Sequential Memory Vulnerabilities

**Key Weakness:** Early perturbations propagate and amplify through hidden and cell states, leading to long-term prediction errors.

**Recommended Defences:**

- **Adversarial training with PGD** to harden cell states against gradient-driven perturbation.
- **Cell state clipping and gate activation regularisation** to dampen accumulation of adversarial gradients.
- **Low-pass input filtering** to suppress sharp fluctuations in early inputs.

**Effectiveness:** These mitigations reduced degradation under FGSM and PGD by up to 25% while preserving validation accuracy.

#### 8.5.2 TCN: Localised Perturbation Sensitivity

**Key Weakness:** Lack of memory prevents recovery from mid-sequence perturbation; highly sensitive to local distortions in receptive field.

**Recommended Defences:**

- **Dilation and skip connections** to increase redundancy and global context.
- **Input noise injection** during training to improve robustness to black-box attacks.
- **Temporal interpolation or padding** to desensitise the model to local sequence offsets.

**Effectiveness:** Most robust when combined with uniform noise injection and small kernel-size smoothing filters.

#### 8.5.3 LNN: ODE Sensitivity and Stability Management

**Key Weakness:** Sensitive to perturbations injected at multiple solver substeps; behaviour governed by dynamics and solver configuration.

**Recommended Defences:**

- **Continuous-time adversarial training** to mimic dynamic perturbation conditions.
- **Stability regularisation:** Penalising fast state transitions or extreme membrane potentials.
- **Reduced solver step size** at inference time to attenuate numerical instability.

**Effectiveness:** Robustness improved significantly under continuous-time and time-warping attacks when combining dynamic training and solver tuning.

#### 8.5.4 Summary Table

Model	Effective Defences
LSTM	PGD adversarial training, state clipping, input smoothing
TCN	Noise injection, receptive field dilation, temporal resampling
LNN	Solver tuning, ODE-stability regularisation, continuous-time training

Table 8.1: Summary of recommended mitigation strategies by model.

These insights may serve as practical guidance for model deployment in adversarial environments, especially in real-time or safety-critical systems where robustness cannot be assumed.

## 8.6 Limitations and Future Work

While the defence strategies outlined in this chapter demonstrate measurable improvements in robustness, several limitations remain. These warrant caution in interpretation and suggest important directions for further research.

### 8.6.1 Limitations

- **Attack-Specific Optimisation:** Many defences, particularly adversarial training, are tuned to specific attack types. As such, gains may not generalise to novel or unseen perturbation methods.
- **Evaluation Scope:** Although a range of attacks was considered, the evaluation was performed on a synthetic 2D spiral task. Generalising these findings to high-dimensional, real-world data (e.g., speech, motion trajectories) requires further validation.
- **Computational Overhead:** Adversarial training and continuous-time solver tuning introduce substantial computational costs, particularly for models like the LNN with dense state transitions.
- **Architectural Rigidity:** Some defences require significant changes to model internals (e.g., solver parameters, memory clipping), which may not be compatible with pre-trained or black-box models.

### 8.6.2 Future Work

- **Robustness Certification:** Incorporating formal verification methods (e.g., symbolic interval analysis or Lipschitz bounding) can provide guarantees under specific perturbation budgets and help validate empirical robustness.
- **Adaptive Defences:** Future work may explore dynamic defence strategies that modulate based on detected input irregularities — such as adaptive smoothing or online solver step adjustment in LNNs.
- **Hybrid Architectures:** Integrating LNNs with recurrent or attention-based modules may improve robustness without sacrificing long-term memory or expressivity.
- **Benchmarking on Real Data:** Applying these defences to real-world tasks, such as physiological signal prediction or time-series classification, would test their practical impact and scalability.

- **Defence-Aware Attacks:** Future research should evaluate robustness under adaptive attackers that account for known defences, offering a more realistic assessment of model security in adversarial environments.

### 8.6.3 Conclusion

The defences presented herein offer a diverse toolbox for enhancing robustness in temporal models. However, robust machine learning remains an adversarial game: as defences evolve, so too do attack strategies. The pursuit of architectures and training regimes that remain stable under dynamic, uncertain, or malicious inputs remains a central challenge in deploying neural systems safely and reliably.

## Chapter 9

# Conclusion

Around 4 pages

This project has explored the robustness of Logic Neural Networks (LNNs) against adversarial attacks, comparing them with traditional neural networks. The evaluation focused on three main types of attacks: temporal, structured, and spatial, each designed to exploit different vulnerabilities in the models. The results indicate that LNNs exhibit superior robustness in temporal and structured adversarial regimes, where the attacks leverage the temporal dependencies and structured relationships inherent in the data. In contrast, traditional neural networks were more susceptible to these types of attacks, highlighting their limitations in handling complex data structures.

In summary, this evaluation demonstrates that:

1. Robustness is a multidimensional property — not all attacks exploit the same vulnerabilities.
2. LNNs, while more complex, deliver meaningful robustness advantages under temporal and structured adversarial regimes.
3. Careful architectural and training design — including regularisation and solver stability — is essential in real-world deployments where adversarial inputs cannot be ruled out.

# Bibliography

- [1] Chahine M, Hasani R, Kao P, Ray A, Shubert R, Lechner M, et al. Robust Flight Navigation out of Distribution with Liquid Neural Networks. *Science Robotics*. 2023 Apr;8(77):eadc8892.
- [2] Hasani R, Lechner M, Amini A, Rus D, Grosu R. Liquid Time-constant Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*. 2021 May;35(9):7657-66.
- [3] TEDx Talks. Liquid Neural Networks | Ramin Hasani | TEDxMIT; 2023.
- [4] Henriksen P, Lomuscio A. Efficient Neural Network Verification via Adaptive Refinement and Adversarial Search;.
- [5] What Is a Recurrent Neural Network (RNN)? | IBM; 2021. <https://www.ibm.com/think/topics/recurrent-neural-networks>.
- [6] Zhang H, Weng TW, Chen PY, Hsieh CJ, Daniel L. Efficient Neural Network Robustness Certification with General Activation Functions. *arXiv*; 2018.

## Chapter 10

# Declaration

Around 2 pages