ROBUST LURIE NETWORKS WITH CONTROLLABLE CONVERGENT DYNAMICS

Anonymous authors

Paper under double-blind review

ABSTRACT

The Lurie Network is proposed as a unifying architecture for modelling timeinvariant nonlinear dynamical systems. Many existing continuous-time models including Recurrent Neural Networks and Neural Oscillators are special cases of the Lurie Network when applied to this domain. Motivated by the need for a general inductive bias, shared by many systems, this paper proposes an approach to enable network weights and biases to be trained in such a manner that a generalised concept of stability is guaranteed. This generalised stability measure is that of k-contraction which enables global convergence to a point, line or plane in the neural state-space. This result is leveraged to construct a Graph Lurie Network (GLN) satisfying the same convergence properties. Unconstrained parametrisations of these conditions are derived allowing the models to be trained using standard optimisation algorithms, whilst limiting the search space to solutions satisfying the k-contraction constraints. Empirical results show significant improvement in terms of prediction accuracy, generalisation and robustness compared to other unconstrained and stability-constrained models. Furthermore, both models consistently learnt representations which respected the convergence behaviour of the dynamics.

025 026 027

003

006

009 010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

029 A Lurie¹ system is a class of nonlinear ordinary differential equations (ODE) comprising a linear time-invariant (LTI) component interconnected with a, potentially time-varying, nonlinearity. Such systems are ubiquitous throughout the sciences and engineering, including machine learning (ML) 031 and neuroscience Pauli et al. (2021a); Lessard et al. (2016); Lanthaler et al. (2024); Wilson & Cowan (1972). When modelling time-invariant dynamical systems, many ML models including Linear State 033 Space Models (LSSM), Recurrent Neural Networks (RNN) and some Graph Neural Networks (GNN) 034 are special cases of a Lurie system $(\S3.1)$. Such models have proven to be highly expressive as demonstrated by their successful application on a wide range of tasks such as sequential processing Kozachkov et al. (2022a); Erichson et al. (2020); Chang et al. (2019); Gu et al. (2021), computer vision Erichson et al. (2020); Chang et al. (2019); Gu et al. (2021), language modelling Gu et al. 037 (2021) and computational chemistry Rusch et al. (2022). 038

Consider dynamical systems in the field of neuroscience: convergence and stability of a latent state 040 are crucial for learning, information propagation Kozachkov et al. (2020); Centorrino et al. (2022); 041 Pascanu et al. (2013); Vogt et al. (2022), memory storage Kozachkov et al. (2022b); Krotov & Hopfield (2020); Ramsauer et al. (2020) and robustness Manchester et al. (2021); Pauli et al. (2021b; 042 2023). For example, multi-stable and orbitally stable systems are, respectively, observed in associative 043 memory Kozachkov et al. (2023) and working memory Kozachkov et al. (2022b). Furthermore, these 044 processes occur in different regions of the brain interconnected through a graph structure. A graph structure and these convergent properties are shared with many other dynamical systems such as 046 chemical processes Ofir et al. (2023), opinion dynamics and power systems Ofir et al. (2024). 047

The convergence and stability analysis of dynamical systems has been well-studied in the control theory literature. A pertinent example is the absolute stability problem where the nonlinearity of the Lurie system is unknown, but assumed to be sector-bounded or slope-restricted. The goal is to find conditions on the model parameters which ensure the trajectories of all Lurie systems, with nonlinearities in the assumed class, uphold a chosen definition of convergence. Approaches to this problem can be classified as Lyapunov analysis Khalil (2002); Park (2002; 1997), Zames-Falb

¹Named after Anatolii Isakovich Lurie and sometimes spelt Lur'e or Lurye.

multipliers Zames & Falb (1968); Turner & Drummond (2019); Carrasco et al. (2016); Drummond et al. (2024) or k-contraction analysis Zhang & Cui (2013); Ofir et al. (2023; 2024). Lyapunov and Zames-Falb multiplier methods are primarily designed to analyse the convergence to an equilibrium point, whereas k-contraction methods analyse a variety of global convergence behaviours including convergence to points, lines and planes. As many ML models are examples of Lurie systems and many activation functions are sector-bounded or slope-restricted (Drummond et al., 2022, Table 1), much of the literature on the absolute stability problem is applicable to systems involving neural networks Pauli et al. (2021a); Richardson et al. (2023; 2024); Fazlyab et al. (2020).

Although designing networks with convergent dynamics is well motivated as an inductive bias, 062 ensuring such a property requires constraints on the network parameters, which can be detrimental if 063 too restrictive. With this in mind, this work focuses on: (i) using k-contraction analysis to derive mild 064 constraints on the weights of the Lurie Network which ensure global convergence to a point, line or 065 plane in the neural state space, for all Lurie Networks with slope-restricted activation functions; (ii) 066 establishing unconstrained parametrisations of these conditions which allows the Lurie Network to be 067 trained using gradient based optimisation algorithms whilst limiting the search space to weights which 068 satisfy the convergence conditions; (iii) constructing a Graph Lurie Network (GLN) from individual Lurie Networks and deriving constraints on the graph coupling matrix which ensure the k-contraction 069 property is maintained. Similar unconstrained parametrisations are also derived. We compare the proposed models against other continuous-time architectures and show the proposed approach leads 071 to more accurate, generalisable and robust models for a range of time-invariant dynamical systems. 072

073

075

090

2 PRELIMINARIES

076 2.1 NOTATION

077 For two integers i < j, we define $[i, j] := \{i, i + 1, \dots, j\}$. The set of non-negative real numbers is 078 denoted by \Re_+ . Symmetric matrices of dimension n are denoted by S^n with the positive definite 079 subset denoted by S_{+}^{n} . All other positive definite subsets are denoted by a + subscript. Square diagonal matrices are denoted by \mathcal{D}^n and $n \times m$ diagonal matrices are symbolised by $\mathcal{D}^{n\bar{m}}$. A positive definite (semi-definite) matrix P is sometimes indicated by $P \succ 0$ ($P \succeq 0$). Negative 081 definite (semi-definite) matrices are indicated analogously. The set of $n \times n$ orthogonal and skewsymmetric matrices are respectively denoted by $\mathcal{SO}(n)$ and so(n). For $W \in \Re^{n \times m}$, the ordered singular values are represented by $\sigma_1(W) \ge \cdots \ge \sigma_{\min(n,m)}(W) \ge 0$ and for $W \in \Re^{n \times n}$, the ordered eigenvalues are denoted by $\lambda_1(W) \ge \cdots \ge \lambda_n(W)$. The k-multiplicative and k-additive 084 085 compound matrices of W are respectively denoted by $W^{(k)}$ and $W^{[k]}$. The Jacobian of a function f(t,x) is denoted by $J_f(t,x)$. The scaled 2-norm of a vector $x \in \Re^n$ with respect to (w.r.t) an 087 invertible scaling matrix $\Theta \in \Re^{n \times n}$ is defined by $|x|_{2,\Theta} := |\Theta x|_2$, and the matrix measure induced by the scaled 2-norm is

$$\mu_{2,\Theta}(W) := \mu_2(\Theta W \Theta^{-1}) = \lambda_1 \left(\frac{\Theta W \Theta^{-1} + (\Theta W \Theta^{-1})^T}{2}\right)$$

2.2 k-CONTRACTION ANALYSIS

In this work, we leverage k-contraction analysis Wu et al. (2022); Muldowney (1990), the geometrical generalisation of contraction analysis Lohmiller & Slotine (1998), as a tool for controlling 094 convergence in the neural state space. Intuitively, k-contraction implies the volume of k-dimensional bodies exponentially converges to zero when governed by the system dynamics. Alternatively, this 096 could be thought of as exponential convergence to a (k-1)-dimensional subspace. When k = 1, this 097 reduces to standard contraction Lohmiller & Slotine (1998), which implies that all trajectories expo-098 nentially converge to a single trajectory. For a general time-varying dynamical system, satisfying the k-contraction property does not guarantee stability. However, for time-invariant dynamical systems, 100 it has been shown that for every bounded solution: 1-contraction implies global convergence to a unique equilibrium point Lohmiller & Slotine (1998), 2-contraction implies global convergence to an equilibrium point, which is not necessarily unique but must be connected along a line Muldowney 102 (1990), and 3-contraction, under certain assumptions, implies convergence to a non-unique attractor 103 Cecilia et al. (2023). Three examples of k-contracting dynamics are presented in Figure 1. 104

Time-invariant dynamical systems which satisfy the *k*-contraction property for $k \in \{1, 2, 3\}$ have several desirable properties for ML models. They can exhibit a wide range of complex convergent behaviours such as multi-stable and orbitally stable systems Zoboli et al. (2024). This suggests that a model can be *expressive* whilst satisfying the *k*-contraction conditions, particularly for higher values



Figure 1: Trajectories from three dynamical systems satisfying the *k*-contraction property. Crosses denote the initial condition and stars denote equilibirum points.

of k where the constraints are less restrictive, as highlighted in §3.2. The k-contraction property also implies an inherent *robustness* as the trajectories can only converge to a finite number of long term behaviours. Next, we present the fundamental k-contraction result from Wu et al. (2022).

Theorem 1 Fix $k \in [1, n]$ and consider the nonlinear system $\dot{x} = f(t, x)$ with $f : \Re_+ \times \Re^n \to \Re^n$ continuously differentiable. If there exists $\eta > 0$ and an invertible matrix $\Theta \in \Re^{n \times n}$ such that

$$\mu_{2,\Theta^{(k)}}\left(J_f^{[k]}(t,x)\right) \le -\eta \quad \forall \ x \in \Re^n \ and \ t \in \Re_+$$

$$\tag{1}$$

then the nonlinear system is k-contracting in the 2-norm w.r.t the metric $P := \Theta^{\top} \Theta$.

130 This result has two features: (i) it requires the existence of an invertible matrix Θ . In the simplest 131 case, one can expect a solution $\Theta = pI_n$ to exist. For other systems, such simple solutions will 132 not exist and more general matrices such as $\Theta \in S^n$ will be required, making the proofs more 133 difficult; (ii) it requires the use of compound matrices. For a matrix $W \in \Re^{n \times m}$, the matrix $W^{[k]}$ 134 with $k \in [1, \min(n, m)]$ will have the size $\binom{n}{k} \times \binom{m}{k}$ which is typically much larger and more 135 computationally difficult to work with. A more technical introduction to k-contraction analysis and 136 compound matrices is presented in A. In 3.2 we derive results which verify (1) for the special case of nonlinear systems of the form (2). 137

3 LURIE NETWORK

141 A Lurie Network is defined by (2) with weights $A \in \Re^{n \times n}, B \in \Re^{n \times m}, C \in \Re^{m \times n}$ and biases 142 $b_x \in \Re^n, b_y \in \Re^m$.

$$\dot{x}(t) = Ax(t) + B\Phi(y(t)) + b_x \qquad y(t) = Cx(t) + b_y \qquad x(0) = x_0 \qquad (2)$$

The model has a biased linear component interconnected with a nonlinearity of the form $\Phi(y) :=$ 145 $[\phi_1(y_1) \dots \phi_m(y_m)]'$ where $\phi_i(y_i)$ is assumed to be slope-restricted with an upper bound g > 0, such 146 that $0 \leq J_{\Phi}(y) \leq gI_m$. This separation of the linear and nonlinear components is useful for analysis. 147 Activation functions which satisfy this slope-restricted assumption include the hyperbolic tangent 148 (tanh) and the rectified linear unit (ReLU). For simplicity, we assume the same scalar nonlinearity is 149 applied element-wise and drop the subscript. The proposed Lurie Network is a very general model as 150 highlighted by the special cases in §3.1 and its relationship to deep feedforward neural networks as 151 outlined in §B.1. Finally, it is important to observe that the model is time-invariant, this implies that 152 if Theorem 1 is satisfied, the model will inherit the appealing convergence and robustness properties 153 stated in §2.2.

154

119

120

124

125

126 127 128

129

138 139

140

143 144

3.1 EXAMPLE LURIE NETWORKS

156

When applied to time-invariant dynamical systems, many models from the ML literature become special cases of (2). In this setting, the time-varying external inputs are replaced with trainable biases. A subset of examples are presented next with further examples included in §B.2. As the results in §3.2 apply to any model of the form (2), they can also be applied to these special cases.

161 **Lipschitz RNN:** A stability constrained RNN Erichson et al. (2020) where the parameters A, C are expressed as a weighted sum of symmetric and skew-symmetric terms in order to control the

171

172

177

178

179

180 181 182

200

215

eigenvalues of the Jacobian. The remaining components are $B = I_n$, $b_x = 0$ and $\phi(\cdot) \equiv \tanh(\cdot)$.

$$A, C \in \{(1 - \beta)(W + W^{\top}) + \beta(W - W^{\top}) - \gamma I_n | W \in \Re^{n \times n}, 0.5 \le \beta \le 1, \gamma > 0\}$$

165 Antisymmetric RNN: A constrained RNN Chang et al. (2019) with the same motivations as above. 166 Related to (2) by A = 0, $B = I_n$, $C \in \{W - \gamma I_n | W \in so(n), \gamma > 0\}$, $b_x = 0$ and $\phi(\cdot) \equiv \tanh(\cdot)$. 167 BNN: A constrained RNN Chang et al. (2019) with the same motivations as above. 168 Antisymmetric RNN: A constrained RNN Chang et al. (2019) with the same motivations as above. 169 Related to (2) by A = 0, $B = I_n$, $C \in \{W - \gamma I_n | W \in so(n), \gamma > 0\}$, $b_x = 0$ and $\phi(\cdot) \equiv \tanh(\cdot)$.

RNNs of RNNs: A 1-contracting graph coupled RNN Kozachkov et al. (2022a). A special case of (2) with $A = L - aI_n$, $B \in \Re^{n \times n}$, $C = I_n$ and $b_y = 0$. The matrix B is block diagonal and contains the synaptic weights of the individual RNNs. The matrix L is the graph coupling matrix.

3.2 *k*-contraction Analysis of Lurie Networks

Two sufficient results which satisfy Theorem 1 and guarantee (2) is k-contracting are presented next. Conditions were derived in (Ofir et al., 2024, Theorem 2) which verify Theorem 1 for a Lurie Network with $A \in \mathcal{D}^n$ and $b_y = 0$. Theorem 2 extends them to account for $A \in \Re^{n \times n}$ and $b_y \neq 0$. Refer to §C.1 for the proof.

Theorem 2 Consider the Lurie Network (2) with $\Phi(y) := [\phi_1(y_1) \dots \phi_m(y_m)]'$ being sloperestricted such that $0 \leq J_{\Phi}(y) \leq gI_m$. Fix $k \in [1, n]$ and define $\alpha_k := (2k)^{-1} \sum_{i=1}^k \lambda_i (A + A^{\top})$. If $\alpha_k < 0$ and

$$g^2 \sum_{i=1}^k \sigma_i^2(B) \sigma_i^2(C) < \alpha_k^2 k \tag{3}$$

then (2) is k-contracting in the 2-norm w.r.t the metric $P = -\alpha_k^{-1}I_n$.

185 The additional freedom permitted by k-contraction over standard contraction is highlighted by the summation of the eigenvalues and singular values. In 1-contraction, Theorem 2 requires the largest eigenvalue of the symmetric component of A to be negative whereas for $k \in [2, n]$, this condition on 187 A becomes incrementally more relaxed as k is increased. Equation (3) illustrates a similar relaxation 188 of the constraints on B and C. Theorem 2 has several appealing features: (i) it does not require 189 the computation of the troublesome compound matrices; (ii) it provides a way of embedding the 190 k-contraction property into the structure of a Lurie Network based on fairly simple unconstrained 191 parametrisations of the weights, as shown in §4; (iii) the biases are not present in the condition, 192 so are naturally unconstrained. The limitation of the result is that only Lurie Networks which are 193 *k*-contracting in a scalar metric can be verified.

We now present a second result which addresses the scalar metric drawback; however, it comes at the cost of strong constraints on the weights B and C. See §C.2 for the proof.

Theorem 3 Consider the Lurie Network (2) with $\Phi(y) := [\phi_1(y_1) \dots \phi_n(y_n)]'$ being sloperestricted such that $0 \leq J_{\Phi}(y) \leq gI_n$. Fix $k \in [1, n]$. If $B \in \mathcal{D}^n$, $C = B^{-1}$ and

$$P^{(k)}A^{[k]} + (A^{[k]})^{\top}P^{(k)} + 2kqP^{(k)} \prec 0$$
(4)

then (2) is k-contracting in the 2-norm w.r.t the metric $P \in \mathcal{D}^n_+$.

Theorem 3 improves upon Theorem 2 in the sense that Lurie Networks which are k-contracting in a diagonal metric can now be verified; however, only when $C = B^{-1}$. Due to this constraint, Theorem 3 has very limited practical use; however, it it may prove to be theoretically insightful for addressing the scalar metric drawback. Finally, it is important to highlight that Theorem 2 and Theorem 3 apply to the class of slope-restricted nonlinearities, so these results address the absolute stability problem for the k-contraction property.

209 210 3.3 GRAPH LURIE NETWORKS

Many larger scale dynamical systems such as molecular, social, biological, and financial networks
 Hamilton et al. (2017) naturally have a graph structure. To make the Lurie Network more applicable
 to these problems, a graph coupling term is introduced to model a set of *q* interacting Lurie Networks.
 To illustrate this, *q* independent Lurie Networks can be modelled by

$$\dot{x}(t) = A_G x(t) + B_G \Phi(y(t)) + b_x \qquad y(t) = C_G x(t) + b_y \qquad x(0) = x_0 \qquad (5)$$

216 where the weights are $A_G := \text{blockdiag}(A_1, \ldots, A_q) \in \Re^{qn \times qn}, B_G := \text{blockdiag}(B_1, \ldots, B_q) \in$ 217 $\Re^{qn \times qm}$, $C_G := \text{blockdiag}(C_1, \ldots, C_q) \in \Re^{qm \times qn}$ and the biases are $b_x \in \Re^{qn}$, $b_y \in \Re^{qm}$. The 218 Graph Lurie Network (GLN) is then defined by 219

$$\dot{x}(t) = (A_G + L)x(t) + B_G \Phi(y(t)) + b_x \qquad y(t) = C_G x(t) + b_y \qquad x(0) = x_0 \qquad (6)$$

where $L := [L_{il}]$ is a block matrix with block $L_{il} \in \Re^{n \times n}$ connecting Lurie Network l to Lurie 221 Network j. The state and nonlinearity of the GLN are defined by $x \in \Re^{qn}$ and $\Phi : \Re^{qm} \to \Re^{qm}$ 222 where the states of the independent Lurie Networks have been stacked into a single state. Interestingly, 223 the GLN (6) is actually a special case of a Lurie Network; however, as the networks get larger, so 224 does the search space. Thus, imposing any assumptions which respect the structure of the problem 225 can reduce the search space and lead to more robust and generalisable models. Any further prior 226 knowledge about the graph can be encoded through constraints on the graph coupling matrix.

227 228 229

220

3.4 k-contraction Analysis of Graph Lurie Networks

In this section, we assume that q independent Lurie Networks (2) are k-contracting in the 2-norm 230 w.r.t metrics P_1, \ldots, P_q . This is equivalent to (5) k-contracting in the 2-norm w.r.t the metric 231 $P = \text{blockdiag}(P_1, \dots, P_q)$. Theorem 2 and Theorem 3 provide two results which can, respectively, 232 verify this w.r.t a scalar metric $P_j = p_j I_n$ and a diagonal metric $P_j \in \mathcal{D}^n_+$ for $j \in [1,q]$. Other 233 results may be used providing they apply to systems of the form (2). Under this assumption, Theorem 234 4 provides a constraint on the graph coupling term which ensures the GLN is k-contracting when 235 constructed from q independently k-contracting Lurie Networks. The proof is detailed in §C.3. 236

Theorem 4 Fix $k \in [1, n]$. Consider a GLN (6) where the q independent Lurie Networks are collectively defined by the weights $A_G := \text{blockdiag}(A_1, \ldots, A_q)$, $B_G := \text{blockdiag}(B_1, \ldots, B_q)$, $C_G := \text{blockdiag}(C_1, \ldots, C_q)$ and biases $b_x \in \Re^{qn}$, $b_y \in \Re^{qm}$ and are k-contracting in the 2-norm w.r.t the metric $P := \text{blockdiag}(P_1, \ldots, P_q)$. If the graph coupling matrix $L \in \Re^{qn \times qn}$ satisfies $P^{(k)}L^{[k]} + (L^{[k]})^{\top}P^{(k)} \prec 0$

240 241

243

244

245

246

247 248

249

237

238 239

242

then (6) is also k-contracting in the 2-norm w.r.t the metric P.

Remark 1 It should be noted that Theorem 4 could be applied more generally for constructing k-contracting graph coupled systems. Providing the q subsystems are k-contracting in the 2-norm w.r.t some metric $P = \text{blockdiag}(P_1, \ldots, P_q)$, then Theorem 4 is applicable when coupling them through the graph coupling term. It is not necessary for the q subsystems to be Lurie Networks.

PARAMETRISATION OF k-CONTRACTING LURIE NETWORKS 4

250 To train a k-contracting Lurie Network using gradient based optimisers, parametrisations which 251 express the constrained weights in terms of unconstrained variables must be found. To formalise this idea, we define the sets $\Omega_2(g,k)$ and $\Omega_4(k,P)$. As the biases do not appear in these sets, they are 253 naturally unconstrained. The set $\Omega_2(g,k)$ contains all the weights of the Lurie Network which satisfy 254 Theorem 2.

$$\Omega_2(g,k) := \left\{ (\bar{A}, \bar{B}, \bar{C}) \mid \alpha_k < 0, \ z := g^2 \sum_{i=1}^{\kappa} \sigma_i^2(\bar{B}) \sigma_i^2(\bar{C}) < \alpha_k^2 k \right\}$$
(8a)

(7)

266 267 268

255

$$\Omega_4(k,P) := \left\{ \bar{L} \in \Re^{qn \times qn} \mid P^{(k)} \bar{L}^{[k]} + (\bar{L}^{[k]})^T P^{(k)} \preceq 0 \right\}$$
(8b)

259 The set $\Omega_4(k, P)$ contains the graph coupling matrices which satisfy Theorem 4. A parametrisation 260 associated with Theorem 3 was not established due to its limitations mentioned earlier. The next 261 results present two different parametrisations of the set $\Omega_2(g, k)$. See §C.4 for the proofs which 262 leverage the eigenvalue and singular value decompositions. 263

Theorem 5 Given g > 0, $k \in [1, n]$, $U_A, U_B, V_C \in \mathcal{SO}(n)$, $V_B, U_C \in \mathcal{SO}(m)$, $\Sigma_B \in \mathcal{D}^{nm}_+$, 264 $\Sigma_C \in \mathcal{D}^{mn}_+, Y_A \in so(n), G_A \in \mathcal{D}^n_+$ and define 265

$$A := \frac{1}{2} U_A \Sigma_A U_A^\top + \frac{1}{2} Y_A \qquad \qquad \Sigma_A := -\sqrt{\frac{4z}{k}} I_n - G_A \qquad (9a)$$

then $(A, B, C) \in \Omega_2(g, k)$.

Theorem 6 Given g > 0, $k \in [1, n]$, $U_A, U_B, V_C \in \mathcal{SO}(n)$, $V_B, U_C \in \mathcal{SO}(m)$, $\Sigma_B \in \mathcal{D}^{nm}_+$, $\Sigma_C \in \mathcal{D}^{mn}_+$, $Y_A \in so(n)$, $\Sigma_{A1} \in \mathcal{D}^{k-1}_+$, $G_{A2} > 0$, $G_{A3} \in \mathcal{D}^{n-k}_+$ and define

$$A := \frac{1}{2} U_A \Sigma_A U_A^\top + \frac{1}{2} Y_A \qquad \qquad B := U_B \Sigma_B V_B^\top \qquad \qquad C := U_C \Sigma_C V_C^\top \quad (10a)$$

$$\Sigma_A := \text{blockdiag}(\Sigma_{A1}, \Sigma_{A2}, \Sigma_{A3}) \qquad \Sigma_{A1} \in \mathcal{D}^{k-1}$$
(10b)

$$\Sigma_{A2} := -\sqrt{4kz} - \sum_{i}^{n-1} (\Sigma_{A1})_{ii} - G_{A2} \quad \Sigma_{A3} := \min(\Sigma_{A1}, \Sigma_{A2}) I_{n-k} - G_{A3}$$
(10c)

277 278 279

273

274 275 276

then
$$(A, B, C) \in \Omega_2(g, k)$$

Remark 2 In both results, a mapping between the sets $so(\cdot)$ and $SO(\cdot)$ was exploited to express the orthogonal variables Lezcano-Casado & Martinez-Rubio (2019). The remaining variables are unconstrained or simply require positive elements, which can be obtained by taking the absolute value of unconstrained variables.

285 In Theorem 5 and Theorem 6, the variables which express the B and C matrices are completely 286 unconstrained. Furthermore, the variables which represent the singular values of B and C are directly 287 used to upper bound α_k since PyTorch provides a differentiable implementation of a sort function. 288 The only source of conservatism in the parametrisations is introduced through the definition of Σ_A . 289 In Theorem 5, the constraint on α_k is replaced by a uniform negative constraint on the individual eigenvalues of $(A + A^{\top})$. If this assumption is true, this significantly speeds up the learning 290 process; however, it can be prohibitive if not. Theorem 6 allows the largest (k-1) eigenvalues to 291 be unconstrained, meaning non-Hurwitz A matrices are encapsulated in the parametrisation. The 292 definition of Σ_A is split into one unconstrained block for the largest (k-1) eigenvalues, a block to 293 ensure Theorem 2 is satisfied, and a block for the remaining eigenvalues which must be less than the other k. The next result provides an unconstrained parametrisation of the set $\Omega_4(k, P)$ in (8b). 295 Coupling this with the earlier parametrisations allows the k-contracting GLN to be trained with 296 gradient based optimisers. See §C.5 for the proof.

Theorem 7 Given $k \in [1, n]$, $G_{L1}, G_{L2} \in \Re^{qn \times qn}$, $\Theta = \text{blockdiag}(\Theta_1, \dots, \Theta_q)$ where $\Theta_j \in S^n_+$ for $j \in [1, q]$, $P = \Theta^\top \Theta$ and define

299 300 301

297

298

$$L := G_{L1} - P^{-1} G_{L1}^{\top} P + \Theta^{-1} G_{L2} \Theta$$
(11)

302 where $(G_{L2} + G_{L2}^{\top})^{[k]} \leq 0$, then $L \in \Omega_4(k, P)$.

Theorem 7 provides flexibility in the choice of G_{L1} , G_{L2} to define different graph structures, where the constraint on G_{L2} is equally a constraint on the sum of it's k-largest eigenvalues (Fact 7). This constraint could be parametrised in a similar manner to the symmetric component of A in Theorem 5 and Theorem 6; in which case, Theorem 7 is satisfied and the graph has a directed all-to-all structure, including self-loops.

Remark 3 If $G_{L2} = 0$ along with the main diagonal blocks of G_{L1} , then (11) is a boundary condition of (7). Thus, Theorem 7 is satisfied and the self-loops are removed from the graph structure.

311 One consideration when constructing k-contracting Lurie Networks is the number of additional 312 parameters required for the parametrisation. A Lurie Network, as defined in (2), has a total parameter count $N_L = n^2 + 2nm + n + m$; whereas a k-contracting Lurie network, constructed according to Theorem 5 or Theorem 6, has a total parameter count $N_K \le 2n^2 + m^2$. The number of parameters only increases significantly when m is large compared to n; however, throughout the literature, many 313 314 315 special cases of a Lurie Network set n = m (see §3.1); in which case, N_K becomes marginally 316 smaller than N_L . Furthermore, the all-to-all graph coupling term has $N_C = 2(qn)^2$ parameters or 317 $N_C = qn^2(q-1)$ if the self-loops are removed according to Remark 3. This analysis highlights that 318 ensuring the Lurie Network and GLN are k-contracting comes at a minimal computational expense. 319

320 321

- 5 EMPIRICAL EVALUATION
- 323 The aim of this section was to test the significance of the theoretical results presented earlier in the paper and the additional expressivity of the Lurie Network compared to special cases from the

327 Model Train. Time Mean Squared Error (mean \pm std, best) 328 Params. Opinion Hopfield Attractor -5, $(1.5 \pm 1) \times 10^{-2},$ $(3.5 \pm 1) \times 10^{-3}.$ k-Lurie 87 8.1 mins $(8.0\pm3)\times10^{-1}$ 330 5.1×10^{-5} $\mathbf{2.6}\times\mathbf{10^{-3}}$ $\mathbf{1.7}\times\mathbf{10^{-3}}$ Net. 331 $(3.7 \pm 4) \times 10^{-5}$ $(3.6\pm2)\times10^{-1},$ $(5.1\pm5)\times10^{-1}$ Lurie 33 8.3 mins 332 5.3×10^{-4} 3.9×10^{-2} 5.7×10^{-2} Net. 333 $(2.0 \pm 2) \times 10^{-5}$ $(2.5 \pm 1) \times 10^{-1}$ $(2.0 \pm 1) \times 10^{-1}$ Neural 983 91.2 mins 334 ODE $4.3 imes 10^{-5}$ 1.5×10^{-2} 1.0×10^{-2} 335 $(3.9 \pm 3) \times 10^{-1}$ -2 $(2.9 \pm 2) \times 10^{-1}$ 1.48 Lipschitz 237.8 mins \pm 2.07, $\underline{1.1}\times 10^{-2}$ 8.8×10^{-3} 3.0×10^{-2} 336 RNN $\begin{array}{ccc} 0.29 & \pm & 0.06, \\ 2.1 \times 10^{-1} \end{array}$ SVD $(8.6 \pm 3) \times 10^{-5}$ 3.12 \pm 337 406.2 mins0.51, 5.7×10^{-4} Combo 2.74338 Antisym. 12 \pm 0.002, 0.43 ± 6.8812 mins0.280.02, \pm 0.17,339 RNN 0.280.416.74340

Table 1: MSE on a test set of 100 trajectories. The mean and standard deviation were calculated after training each model N = 3 times on a single T4 GPU (Google Colab). The lowest MSE for each of the N repetitions is also presented alongside the training time from one run on the opinion dataset.

literature. The importance of these were evaluated using the prediction accuracy, generalisation
 and robustness of the *k*-contracting Lurie Network and GLN on a range of datasets generated by
 time-invariant dynamical systems.

344 5.1 DATASETS AND TRAINING345

We first consider three time-invariant dynamical systems: (i) an opinion dynamics model of a social 346 network where all opinions agree and thus converge to a unique equilibrium point; (ii) a Hopfield 347 network of associative memory with two stable equilibrium points and one unstable; (iii) a generic 348 simple attractor which could be used to model the stored patterns in working memory. Each system 349 has n = 3 states allowing for easy visualisation of the ground truth and predictions. For each 350 dynamical system, we have a dataset including 1,000 trajectories, sampled every 0.01s over a 351 20s interval. The test sets were formed by holding out 100 trajectories. The input to each model 352 was the initial condition sampled from a uniform distribution with the domain $(-1, +1)^3$ for the 353 opinion/Hopfield datasets and $(-3, +3)^3$ for the simple attractor. The full trajectory was then used as the target to train the model. An illustration of these datasets can be seen in Figure 1 and further 354 details about the data generation can be found in Appendix D. 355

356 To test the out of distribution generalisation and robustness, two additional datasets were generated 357 for each of the opinion, Hopfield and attractor tasks. These differ from the training datasets in the 358 following ways: (i) they include 100 trajectories over a 30s interval; (ii) the initial conditions were sampled from a uniform distribution over the intervals $1 < |x_i(0)| < 4$ for the opinion/Hopfield 359 datasets and $3 < |x_i(0)| < 6$ for the simple attractor, where $i \in [1, 2, 3]$. These trajectories are 10s 360 longer than the training data with initial conditions also sampled outside the training distribution. 361 Finally, to test the robustness, a third dataset was generated using the same trajectories as the targets 362 but with noise sampled from the standard normal distribution added to the initial conditions. 363

For the second set of experiments, we consider two 30 dimensional dynamical systems. The first is a graph-coupled (GC) Hopfield network, formed by connecting 10 previously described Hopfield networks through a graph coupling matrix. To ensure the convergence property was preserved, the matrix was expressed by (11) with G_{L1} sampled from a uniform distribution and $G_{L2} = 0$. The second system was a graph-coupled attractor, constructed in the same way. The datasets were generated as described in the previous two paragraphs; however, they each included 30,000 trajectories. The study of interacting memory networks is the motivation behind these examples.

The training settings used are explicitly detailed in Appendix D. For all models and all datasets, the
 mean squared error (MSE) loss was used alongside the Adam optimiser. All code was implemented
 in PyTorch.

373

5.2 *k*-contracting Lurie Network

375

In this section, we compare the k-contracting Lurie Network against five other continuous-time models: (i) the unconstrained Lurie Network, for testing the importance of the k-contraction constraints; (ii) an unconstrained Neural ODE with two hidden layers, each comprised of 20 neurons and ReLU Table 2: MSE on a new test set of 100 trajectories where: (i) the initial conditions were sampled outside the range used for training and the trajectories were 10s longer than the training set; (ii) additionally, noise sampled from the standard normal distribution was added to the initial conditions.

382	Model	Me	an Squared Er	ror	Mean Squared Error (noisy inputs)				
383		Opinion	Hopfield	Attractor	Opinion	Hopfield	Attractor		
204	k-Lurie Net.	$2.9 imes10^{-3}$	$5.6 imes10^{-2}$	$2.3 imes10^{-1}$	$2.0 imes10^{-2}$	$3.2 imes10^{-1}$	1.28		
304	Lurie Net.	6.4×10^{-2}	1.8×10^{-1}	5.96	2.7×10^{-1}	4.4×10^{-1}	6.79		
385	Neural ODE	1.2×10^{-2}	1.09	2.31	3.9×10^{-2}	1.63	4.76		
386	Lipschitz RNN	2.3×10^{-1}	7.3×10^{-1}	6.2×10^{-1}	2.9×10^{-1}	9.7×10^{-1}	1.71		
387	SVD Combo	1.0×10^{-2}	2.38	20.9	$3.3 imes 10^{-2}$	7.97	30.30		
388	Antisym. RNN	6.43	5.25	52.1	7.29	6.33	52.9		

activations, as detailed in Xia et al. (2021); (iii) three other constrained continuous-time models: the Lipschitz RNN Erichson et al. (2020), Antisymmetric RNN Chang et al. (2019) and a single node from the 1-contracting SVD Combo Network Kozachkov et al. (2022a). For each model, the weights were initialised using the default PyTorch settings and the biases were initialised as zero. The numerical integration was performed using the Euler method with step size $\delta = 1 \times 10^{-2}$. Besides the Neural ODE, each model had tanh activations, which is slope-restricted with g = 1. The parameter count of each model, excluding the Neural ODE, is determined by the dimension of the state as they do not explicitly have any hidden layers. This can be seen from (2) and Section 3.1.

For the opinion and Hopfield datasets, the *k*-contracting Lurie Network was constructed according to Theorem 5 whereas Theorem 6 was used for the simple attractor. We set k = 1 for the opinion dynamics, k = 2 for the mulit-stable Hopfield network and k = 3 for the simple attractor.

400 Table 1 compares the MSE on the test set of each task. The k-contracting Lurie Network achieved 401 the best MSE on two out of three examples. The importance of the k-contraction conditions is particularly clear when comparing the mean and standard deviation with that of the unconstrained 402 Lurie Network. These conditions clearly reduce the search space to a tractable region to optimise 403 over as the MSE of the unconstrained Lurie Network is at least an order of magnitude worse than 404 its k-contracting counterpart. The other models perform as one would expect. The Neural ODE 405 demonstrates strong accuracy across all tasks, albeit at the expense of considerably longer training 406 time. The SVD Combo performs well on the opinion dataset, where the 1-contraction assumption 407 is valid, but struggles on the others. The Lipschitz RNN struggles on the attractor dataset whilst 408 the Antisymmetric RNN struggles across the board due to the A matrix being fixed at zero and the eigenvalues of the C matrix being fixed to almost purely imaginary values. Figures 4, 7, 10 show a 409 random sample of trajectories from each test set, along with the predictions of each model. 410

411 Table 2 compares the generalisation and robustness of the models on each task. No new models 412 were trained, instead the best models from Table 1 were directly applied to these out of distribution and noisy datasets (Section 5.1). The k-contracting Lurie Network performs the best on all of these 413 datasets and for some, it still demonstrates a MSE of an order of magnitude lower than the next best 414 model. The MSE of the unconstrained Lurie Network and Neural ODE tended to drop off for these 415 datasets whereas the MSE of the constrained models, excluding the Antisymmetric RNN, tended to 416 stay fairly consistent when their assumptions were valid. Figures 5, 8, 11 show a random sample of 417 noise-free trajectories and predictions for each dataset, whilst Figures 6, 9, 12 repeat the same for the 418 noisy inputs. The k-contracting Lurie Network was the only model which converged to the correct 419 long-term behaviour under all conditions; even when noise was added, the error was only present 420 during the initial transient.

421 422

423

5.3 *k*-contracting Graph Lurie Network

This section repeats the same experiments as the previous section, but for two 30-dimensional graphcoupled (GC) dynamical systems: the GC Hopfield network and the GC simple attractor (Section
5.1). The state of these datasets is significantly larger than those used in other dynamical systems
datasets such as: (i) the LASA dataset Lemme et al. (2015) where the 2-d trajectories are typically
stacked to form 4-d or 8-d trajectories; (ii) simulated datasets of the 2, 4 or 8 link pendulums which,
respectively, have 4, 8 or 16 dimension trajectories.

For both datasets, the GLN was constructed according to Theorem 7 and Remark 3 where n = m = 3and q = 10. The individual Lurie Networks were constructed according to Theorem 5 for the GC Hopfield network, with k = 2 and Theorem 6 for the GC attractor, with k = 3. The Neural ODE

	-	-	_	-
Model	Params.	Train. Time	Mean Squared Erro	or (mean \pm std, best)
			GC Hopfield	GC Attractor
k-Lurie Net.	8046	24.8 mins	$0.238 \pm 0.0011, 0.237$	$0.737 \pm 0.0108, 0.723$
Lurie Net.	2760	21.5 mins	$2.537 \pm 1.3327, 1.157$	$291.8 \pm 191.84, 21.83$
Neural ODE	20903	126.5 mins	$0.138 \pm 0.0291, 0.114$	$3.445 \pm 0.3626, 2.942$
Lipschitz RNN	1832	17.6 mins	$0.124 \pm 0.0161, 0.105$	$0.658 \pm 0.1604, 0.433$
Antisym. RNN	930	25 mins	$0.448 \pm 0.0023, 0.444$	$6.386 \pm 0.0066, 6.380$
GC SVD Combo	1300	18.4 mins	$0.339 \pm 0.1363, 0.229$	$3.024 \pm 1.1240, 1.435$
GLN	1770	25.5 mins	$0.016 \pm 0.0003, 0.016$	$0.293 \pm 0.1969, 0.015$

Table 3: MSE on a test set of 100 trajectories. The mean and standard deviation were calculated after 433 training each model N = 3 times on a single A100 GPU (Google Colab). The lowest MSE for each 434 of the N repetitions is also presented alongside the training time from one run on the opinion dataset. 435

Table 4: MSE on a new test set of 100 trajectories where: (i) the initial conditions were sampled outside the range used for training and the trajectories were 10s longer than the training set; (ii) additionally, noise sampled from the standard normal distribution was added to the initial conditions.

Model	Mean	Squared Error	Mean Square	Mean Squared Error (noisy inputs)			
	GC Hopfield	GC Attractor	GC Hopfield	GC Attractor			
k-Lurie Net.	0.77	6.10	1.05	6.90			
Lurie Net.	20350	5638	25481	6368			
Neural ODE	2.12	24.93	2.76	25.11			
Lipschitz RN	NN 3.58	6.17	4.25	7.84			
Antisym. RN	NN 5.32	50.68	6.21	52.10			
GC SVD Co	ombo 0.84	11.40	1.07	12.30			
GLN	0.08	1.67	0.26	2.85			

455 was formed using a two layers with 100 neurons and ReLU activations. The SVD Combo leveraged a similar graph structure to the one used in this paper. The other models were the same as in the 456 previous section but with a 30-dimensional state. 457

458 Table 3 shows the GLN had a lower MSE than all other models by a factor of 10. This included the 20,903 parameter Neural ODE. Comparing the GLN, k-contracting Lurie Network and the 459 unconstrained Lurie Network highlights the improvements due to the k-contraction conditions and 460 the graph structure. Table 4 also indicates that the GLN generalised and remained robust to noise, 461 even in these high-dimensional systems. The same can be said for the k-contracting Lurie Network 462 which achieved the second lowest MSE on the generalisation and robustness tests. The inherent graph 463 structure of the SVD Combo may be the reason behind its improved ranking in the GC Hopfield 464 tasks whereas the Neural ODE particularly struggled to generalise for the GC attractor. Possible 465 explanations behind the poor performance of the constrained benchmark models are suggested next.

RELATED WORK 6

469 Several constrained continuous-time RNN models exist in the ML literature. The model structure of 470 three notable examples were presented in §3.1 as they happen to be special cases of a Lurie Network, 471 when modelling time-invariant dynamical systems. The Antisymmetric RNN Chang et al. (2019) 472 and the Lipschitz RNN Erichson et al. (2020) were designed to address the exploding and vanishing gradient problem Pascanu et al. (2013). The Antisymmetric RNN did so by parametrising the RNN 473 such that the real eigenvalues of the Jacobian were zero. It achieved this by setting A = 0 and 474 restricting C to being skew-symmetric. Whilst this does prevent the gradients from exploding and 475 vanishing, it restricts the dynamics which the model can learn to purely oscillatory behaviour. The 476 Lipschitz RNN has a more relaxed parametrisation. This model constructs the A and C matrices such 477 that they are both convex combinations of symmetric and skew-symmetric matrices. However, the 478 weight of the symmetric matrix can only vary between 0 and 0.5, whereas the weight of the skew-479 symmetric matrix can vary between 0.5 and 1. Again, this addresses the vanishing and exploding gradient problem, but the model cannot encode dynamics which are predominately decaying or 480 growing. Finally, the RNN proposed in Kozachkov et al. (2022a) has biological motivations and 481 encodes 1-contracting dynamics. This implies that all possible trajectories will exponentially converge, 482 making the model robust to input disturbances. 483

Whilst the models above address a variety of problems, each model is quite limited in the range of 484 dynamics they can learn, which is a problem when trying to design a general model for learning 485 time-invariant dynamical systems. With respect to (2), the Lurie Network has more flexibility than all

444

445

466 467

of these models. Firstly, it includes all three weight matrices (A, B, C) and biases (b_x, b_y) , whereas the models mentioned above fix at least one of these parameters. Secondly, the constraints imposed, and the corresponding parametrisations, allow the model to learn a variety of dynamics including, but not limited to, those mentioned above. The only limitation is that the dynamics must converge in some way.

The Lurie Network is also related to a class of feed-forward models named *Implicit* or *Equilibrium Networks*. These models use an implicit equation to express the relationship between the model output, layer outputs and model input in a compact vectorised form El Ghaoui et al. (2021). Like the Lurie Network, these models can be represented by the interconnection of a linear time-invariant system and a nonlinearity. This makes analysis tools from Control Theory, such as Lipschitz bounds, applicable to these models Fazlyab et al. (2019). An additional connection is that the solution to the implicit equations correspond to equilibrium points of a Lurie system Revay et al. (2020).

As mentioned in the introduction, the k-contraction constraints used in this paper have an interesting 498 connection to some properties observed in biological learning systems. A 2-contracting model 499 can replicate the behaviour of associative memory, where every stored pattern corresponds to an 500 equilibrium Kozachkov et al. (2023). Furthermore, a 3-contracting model can replicate the dynamics 501 of working memory, where patterns are retained as attractor states Kozachkov et al. (2022b). The 502 conditions developed in this paper could be of interest to Neuroscientists and ML researchers interested in memory storage and retrieval Ramsauer et al. (2020); Hopfield (1984); Krotov & Hopfield 504 (2020). The formation of the GLN also has connections to biology. For example, constructing larger 505 systems from smaller modules can be motivated by evolutionary biology Simon (1962) where the name facilitated variation Gerhart & Kirschner (2007) is used to describe the development of traits 506 in response to adaptation of the regulatory elements that connect modules, rather than the core 507 components themselves. This was investigated in Kozachkov et al. (2022a) where the weights of the 508 individual RNNs satisfied a 1-contraction condition but were fixed. Only the graph coupling weights 509 were learnt during training. 510

The relationship between properties guaranteed by k-contraction analysis and those observed in 511 associative and working memory suggest the k-contracting Lurie Network possesses a number of 512 appealing properties for an ML model; hence, it may be suitable for a wider class of ML problems. 513 This proposition is supported by the successful application of the constrained RNN models (special 514 cases of the Lurie Network) on a wide range of ML tasks. Since the Lurie Network is a more 515 structured, time-invariant example of a Neural ODE Chen et al. (2018), it will be applicable to a 516 similar array of tasks. Beyond modelling time-invariant dynamical systems, this includes image classification and continuous normalising flows Kidger (2022). The only limiting requirement is 517 that the input must be passed in through the initial condition which in some cases, such as image 518 classification, may require a pre-processing layer. 519

- 520
- 521
- 522 523

7 CONCLUSION

524 525

526 The Lurie Network was presented as a novel and unifying architecture for modelling time-invariant 527 dynamical systems, with more flexibility than comparable methods. Many dynamical systems of 528 interest exhibit convergent behaviour in some form; this inductive bias was built into the Lurie 529 Network through the use of k-contraction analysis. Furthermore, a principled approach was proposed 530 for constructing k-contracting Graph Lurie Networks out of k-contracting modules. Both the kcontracting Lurie Network and GLN demonstrated improved prediction accuracy, out of distribution 531 generalisation and robustness on a range of examples. Furthermore, they were the only models 532 to consistently and accurately predict the convergence behaviour of the dynamics in both the high 533 and low dimensional datasets. Future theoretical work will try to expand the class of systems the 534 k-contracting Lurie Network can be optimised over, by obtaining similar results for systems which are k-contracting in a diagonal metric. It would also be interesting to find conditions on the inputs of a 536 time-varying Lurie Network for which the convergence properties of the time-invariant k-contracting Lurie Network are also upheld, making the Lurie Network applicable to sequential processing tasks. Finally, we would like to empirically investigate the performance of the model on a wider range of 538 applications. The recently proposed working memory benchmark Sikarwar & Zhang (2024) is of particular interest due to its relationship with 3-contracting dynamics.

540 REFERENCES

565 566

- Eyal Bar-Shalom, Omri Dalin, and Michael Margaliot. Compound matrices in systems and control theory: a tutorial. *Mathematics of Control, Signals, and Systems*, pp. 1–55, 2023.
- Joaquin Carrasco, Matthew C Turner, and William P Heath. Zames-Falb multipliers for absolute
 stability: From O'Shea's contribution to convex searches. *European Journal of Control*, 28:1–19, 2016.
- Andreu Cecilia, Samuele Zoboli, Daniele Astolfi, Ulysse Serres, and Vincent Andrieu. Generalized
 Lyapunov conditions for k-contraction: analysis and feedback design. 2023.
- Veronica Centorrino, Francesco Bullo, and Giovanni Russo. Contraction analysis of Hopfield neural networks with Hebbian learning. In 2022 IEEE 61st Conference on Decision and Control (CDC), pp. 622–627. IEEE, 2022.
- Bo Chang, Minmin Chen, Eldad Haber, and Ed H Chi. Antisymmetricrnn: A dynamical system view
 on recurrent neural networks. *arXiv preprint arXiv:1902.09689*, 2019.
- 556 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- Ross Drummond, Matthew C Turner, and Stephen R Duncan. Reduced-order neural network synthesis
 with robustness guarantees. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- Ross Drummond, Chris Guiver, and Matthew C Turner. Exponential input-to-state stability for
 Lur'e systems via integral quadratic constraints and Zames–Falb multipliers. *IMA Journal of Mathematical Control and Information*, pp. dnae003, 2024.
 - Laurent El Ghaoui, Fangda Gu, Bertrand Travacca, Armin Askari, and Alicia Tsai. Implicit deep learning. *SIAM Journal on Mathematics of Data Science*, 3(3):930–958, 2021.
 - N Benjamin Erichson, Omri Azencot, Alejandro Queiruga, Liam Hodgkinson, and Michael W Mahoney. Lipschitz recurrent neural networks. *arXiv preprint arXiv:2006.12070*, 2020.
- Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas. Efficient and accurate estimation of lipschitz constants for deep neural networks. *Advances in neural information processing systems*, 32, 2019.
- Mahyar Fazlyab, Manfred Morari, and George J Pappas. Safety verification and robustness analysis
 of neural networks via quadratic constraints and semidefinite programming. *IEEE Transactions on Automatic Control*, 67(1):1–15, 2020.
- John Gerhart and Marc Kirschner. The theory of facilitated variation. *Proceedings of the National Academy of Sciences*, 104(suppl_1):8582–8589, 2007.
- Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. *Advances in neural information processing systems*, 33: 1474–1487, 2020.
- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modelling long sequences with structured
 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *arXiv preprint arXiv:1709.05584*, 2017.
- John J Hopfield. Neurons with graded response have collective computational properties like those of
 two-state neurons. *Proceedings of the national academy of sciences*, 81(10):3088–3092, 1984.
- ⁵⁸⁸ Roger A Horn and Charles R Johnson. *Topics in matrix analysis*. Cambridge university press, 1994.
- 590 Hassan K Khalil. Nonlinear systems. *Patience Hall*, 115, 2002.
- ⁵⁹¹ Patrick Kidger. On neural differential equations. *arXiv preprint arXiv:2202.02435*, 2022.
- 593 Leo Kozachkov, Mikael Lundqvist, Jean-Jacques Slotine, and Earl K Miller. Achieving stable dynamics in neural circuits. *PLoS computational biology*, 16(8):e1007659, 2020.

594 595 596	Leo Kozachkov, Michaela Ennis, and Jean-Jacques Slotine. RNNs of RNNs: Recursive construction of stable assemblies of recurrent neural networks. <i>Advances in Neural Information Processing Systems</i> , 35:30512–30527, 2022a.
598 599 600	Leo Kozachkov, John Tauber, Mikael Lundqvist, Scott L Brincat, Jean-Jacques Slotine, and Earl K Miller. Robust and brain-like working memory through short-term synaptic plasticity. <i>PLOS Computational Biology</i> , 18(12):e1010776, 2022b.
601 602	Leo Kozachkov, Jean-Jacques Slotine, and Dmitry Krotov. Neuron-astrocyte associative memory. <i>arXiv preprint arXiv:2311.08135</i> , 2023.
603 604 605	Dmitry Krotov and John Hopfield. Large associative memory problem in neurobiology and machine learning. <i>arXiv preprint arXiv:2008.06996</i> , 2020.
606 607	Samuel Lanthaler, T Konstantin Rusch, and Siddhartha Mishra. Neural oscillators are universal. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
608 609 610	Andre Lemme, Yaron Meirovitch, M Khansari-Zadeh, Tamar Flash, Aude Billard, and Jochen J Steil. Open-source benchmarking for learned reaching motion generation in robotics. <i>Paladyn, Journal</i> <i>of Behavioral Robotics</i> , 6(1):000010151520150002, 2015.
611 612 613	Laurent Lessard, Benjamin Recht, and Andrew Packard. Analysis and design of optimization algorithms via integral quadratic constraints. <i>SIAM Journal on Optimization</i> , 26(1):57–95, 2016.
614 615 616	Mario Lezcano-Casado and David Martinez-Rubio. Cheap orthogonal constraints in neural networks: A simple parametrization of the orthogonal and unitary group. In <i>International Conference on Machine Learning</i> , pp. 3794–3803. PMLR, 2019.
617 618 619	Winfried Lohmiller and Jean-Jacques E Slotine. On contraction analysis for non-linear systems. <i>Automatica</i> , 34(6):683–696, 1998.
620 621 622	Ian R Manchester, Max Revay, and Ruigang Wang. Contraction-based methods for stable identifi- cation and robust machine learning: a tutorial. In 2021 60th IEEE Conference on Decision and Control (CDC), pp. 2955–2962. IEEE, 2021.
623 624	James S Muldowney. Compound matrices and ordinary differential equations. <i>The Rocky Mountain Journal of Mathematics</i> , pp. 857–872, 1990.
625 626 627	Ron Ofir, Jean-Jacques Slotine, and Michael Margaliot. <i>k</i> -contraction in a generalized Lurie system. <i>arXiv preprint arXiv:2309.07514</i> , 2023.
628 629	Ron Ofir, Alexander Ovseevich, and Michael Margaliot. Contraction and k-contraction in Lurie systems with applications to networked systems. <i>Automatica</i> , 159:111341, 2024.
630 631	Poogyeon Park. A revisited Popov criterion for nonlinear Lur'e systems with sector-restrictions. International Journal of Control, 68(3):461–470, 1997.
633 634	PooGyeon Park. Stability criteria of sector-and slope-restricted Lur'e systems. <i>IEEE Transactions on Automatic Control</i> , 47(2):308–313, 2002.
635 636 637	Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In <i>International conference on machine learning</i> , pp. 1310–1318. Pmlr, 2013.
638 639 640	Patricia Pauli, Dennis Gramlich, Julian Berberich, and Frank Allgöwer. Linear systems with neural network nonlinearities: Improved stability analysis via acausal Zames-Falb multipliers. In 2021 60th IEEE Conference on Decision and Control (CDC), pp. 3611–3618. IEEE, 2021a.
641 642	Patricia Pauli, Anne Koch, Julian Berberich, Paul Kohler, and Frank Allgöwer. Training robust neural networks using Lipschitz bounds. <i>IEEE Control Systems Letters</i> , 6:121–126, 2021b.
643 644 645 646	Patricia Pauli, Dennis Gramlich, and Frank Allgöwer. Lipschitz constant estimation for 1d convolu- tional neural networks. In <i>Learning for Dynamics and Control Conference</i> , pp. 1321–1332. PMLR, 2023.

647 Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.

- 648 Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas Adler, 649 Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, et al. Hopfield networks 650 is all you need. arXiv preprint arXiv:2008.02217, 2020. 651 Max Revay, Ruigang Wang, and Ian R Manchester. Lipschitz bounded equilibrium networks. arXiv 652 preprint arXiv:2010.01732, 2020. 653 654 Carl Richardson, Matthew Turner, Steve Gunn, and Ross Drummond. Strengthened stability analysis of discrete-time Lurie systems involving ReLU neural networks. In 6th Annual Learning for 655 Dynamics & Control Conference, pp. 209-221. PMLR, 2024. 656 657 Carl R Richardson, Matthew C Turner, and Steve R Gunn. Strengthened Circle and Popov Criteria 658 for the stability analysis of feedback systems with ReLU neural networks. IEEE Control Systems 659 Letters, 2023. 660 T Konstantin Rusch, Ben Chamberlain, James Rowbottom, Siddhartha Mishra, and Michael Bronstein. 661 Graph-coupled oscillator networks. In International Conference on Machine Learning, pp. 18888– 662 18909. PMLR, 2022. 663 Ankur Sikarwar and Mengmi Zhang. Decoding the enigma: benchmarking humans and AIs on the 664 many facets of working memory. Advances in Neural Information Processing Systems, 36, 2024. 665 666 Herbert A Simon. The architecture of complexity. Proceedings of the American philosophical society, 667 106(6):467-482, 1962. 668 Matthew C Turner and Ross Drummond. Analysis of MIMO Lurie systems with slope restricted 669 nonlinearities using concepts of external positivity. In 2019 IEEE 58th Conference on Decision 670 and Control (CDC), pp. 163–168. IEEE, 2019. 671 672 Mathukumalli Vidyasagar. Nonlinear systems analysis. SIAM, 2002. 673 Ryan Vogt, Maximilian Puelma Touzel, Eli Shlizerman, and Guillaume Lajoie. On Lyapunov 674 exponents for RNNs: Understanding information propagation using dynamical systems tools. 675 Frontiers in Applied Mathematics and Statistics, 8:818799, 2022. 676 677 Hugh R Wilson and Jack D Cowan. Excitatory and inhibitory interactions in localized populations of model neurons. *Biophysical journal*, 12(1):1–24, 1972. 678 679 Chengshuai Wu, Ilya Kanevskiy, and Michael Margaliot. k-contraction: Theory and applications. 680 Automatica, 136:110048, 2022. 681 Hedi Xia, Vai Suliafu, Hangjie Ji, Tan Nguyen, Andrea Bertozzi, Stanley Osher, and Bao Wang. 682 Heavy ball neural ordinary differential equations. Advances in Neural Information Processing 683 Systems, 34:18646-18659, 2021. 684 685 George Zames and PL Falb. Stability conditions for systems with monotone and slope-restricted 686 nonlinearities. SIAM Journal on Control, 6(1):89-108, 1968. 687 Xiaojiao Zhang and Baotong Cui. Synchronization of Lurie system based on contraction analysis. 688 Applied Mathematics and Computation, 223:180–190, 2013. 689 690 Samuele Zoboli, Andreu Cecilia, and Sophie Tarbouriech. Quadratic abstractions for k-contraction. 2024. 691 692 693 696 697 699
- 700
- 701

702 A EXTENDED PRELIMINARIES

As many of the tools used in this paper are not well-known in the machine learning community, a condensed background is presented here, based on results from Bar-Shalom et al. (2023); Wu et al. (2022). The first section is on compound matrices, an important algebraic tool needed for generalising contraction analysis to k-contraction analysis. Following that, a geometric interpretation of the k-compound matrix is stated through its relationship to the volume of a parametrised set. Finally, the results in these two sections are utilised to define and provide intuition for the convergence of k-contracting dynamics.

711 712

724 725

727

728

729

730 731

738 739 740

741

742 743

744

745 746

747 748

749

752

753

755

A.1 COMPOUND MATRICES

In this section, we document several known definitions and algebraic results related to compound matrices. The results are included without proof; the interested reader should refer to Bar-Shalom et al. (2023) for a more detailed tutorial on the topic.

T16 T17 T18 Let n be a positive integer and fix $k \in [1, n]$. The ordered set of *increasing* sequences of k integers from [1, n] is denoted by Q(k, n). For example: $Q(3, 4) = \{(1, 2, 3), (1, 2, 4), (1, 3, 4), (2, 3, 4)\}.$

Now consider a matrix $W \in \Re^{n \times m}$. For $\alpha \in Q(k, n)$ and $\beta \in Q(k, m)$, the matrix $W[\alpha|\beta]$ denotes the $k \times k$ sub-matrix obtained by taking the entries of W along the rows indexed by α and columns indexed by β . As an example, if k = 2 and n = m = 4, then $Q(2, 4) = \{(1, 2), (1, 3), (1, 4), (2, 3), (2, 4), (3, 4)\}$. The sub-matrix W[(1, 2)|(3, 4)] would then be given by

$$W[(1,2)|(3,4)] = \begin{bmatrix} w_{13} & w_{14} \\ w_{23} & w_{24} \end{bmatrix}$$

The k-minors of the matrix W are defined as $W(\alpha|\beta) := \det(W[\alpha|\beta])$.

Definition 1 (*k*-multiplicative compound) Let $W \in \Re^{n \times m}$ and fix $k \in [1, \min(n, m)]$. The *k*-multiplicative compound of W, denoted $W^{(k)}$, is the $\binom{n}{k} \times \binom{m}{k}$ matrix containing all the *k*-minors of W ordered lexicographically.

For example, if we have n = m = 3 and k = 2 then $\alpha, \beta \in Q(2,3) = \{(1,2), (1,3), (2,3)\}$ and

$$W^{(2)} = \begin{bmatrix} W((1,2)|(1,2)) & W((1,2)|(1,3)) & W((1,2)|(2,3)) \\ W((1,3)|(1,2)) & W((1,3)|(1,3)) & W((1,3)|(2,3)) \\ W((2,3)|(1,2)) & W((2,3)|(1,3)) & W((2,3)|(2,3)) \end{bmatrix}$$

Some important special cases include

$$W^{(1)} = W \qquad W^{(n)} = \det(W) \qquad (pI_n)^{(k)} = p^k I_s \qquad W \in \mathcal{D}^n \to W^{(k)} \in \mathcal{D}^s$$
(12)

with $s := \binom{n}{k}$. Next, we present a series of algebraic results concerned with the k-multiplicative compound.

Fact 1 (Cauchy-Binet Formula) If $U \in \Re^{n \times m}$, $V \in \Re^{m \times p}$ and $k \in [1, \min(n, m, p)]$, then $(UV)^{(k)} = U^{(k)}V^{(k)}$

Fact 2 Fix $k \in [1, \min(n, m)]$. As a consequence of Definition 1, if $W \in \Re^{n \times m}$ then

 $(W^{\top})^{(k)} = (W^{(k)})^{\top}$

Fact 3 Fix $k \in [1, n]$. If $W \in \mathbb{R}^{n \times n}$ is non-singular, then by Theorem 1

$$(W^{-1})^{(k)} = (W^{(k)})^{-1}$$

Fact 4 Fix $k \in [1, \min(n, m, p)]$. If $W \in \mathbb{R}^{n \times n}$, $U \in \mathbb{R}^{p \times n}$ and $V \in \mathbb{R}^{n \times p}$, then by Theorem 1

 $(UWV)^{(k)} = U^{(k)}W^{(k)}V^{(k)}$

Fact 5 Fix $k \in [1, n]$. An implication of Theorem 1 is that if $W \in \Re^{n \times n}$ with eigenvalues $\lambda_1, \ldots, \lambda_n$, then the eigenvalues of $W^{(k)}$ are the $\binom{n}{k}$ products

$$\left\{\prod_{l=1}^k \lambda_{i_l} : 1 \le i_1 < \dots < i_k \le n\right\}$$

We now introduce the definition of a second compound matrix, the k-additive compound, and a set of algebraic results related to it.

Definition 2 (*k*-additive compound) Let $W \in \Re^{n \times n}$ and $k \in [1, n]$. The *k*-additive compound of W is the $\binom{n}{k} \times \binom{n}{k}$ matrix defined by

$$W^{[k]} := \frac{d}{d\epsilon} \left(I_n + \epsilon W \right)^{(k)}|_{\epsilon = 0}$$

Special cases include

$$W^{[1]} = W \qquad W^{[n]} = tr(W) \qquad (pI_n)^{[k]} = kpI_s \qquad W \in \mathcal{D}^n \to W^{[k]} \in \mathcal{D}^s$$
(13)

with $s := \binom{n}{k}$. Like before, we now present some useful algebraic results related to the k-additive compound.

Fact 6 If $W \in \Re^{n \times n}$ and $k \in [1, n]$, then as a consequence of Definition 2

$$(W^{\top})^{[k]} = (W^{[k]})^{\top}$$

Fact 7 Fix $k \in [1, n]$. For $W \in \Re^{n \times n}$ with eigenvalues $\lambda_1, \ldots, \lambda_n$, the eigenvalues of $W^{[k]}$ are the $\binom{n}{k}$ sums

$$\{\sum_{l=1}^k \lambda_{i_l} : 1 \le i_1 < \dots < i_k \le n\}$$

An important consequence of Fact 7 is that if W is positive definite (semi-definite), then this property is upheld by $W^{[k]}$. Opposite conclusions can be drawn if W is negative definite (semi-definite).

Fact 8 Fix $k \in [1, n]$. If $U, V \in \Re^{n \times n}$, then

$$(U+V)^{[k]} = U^{[k]} + V^{[k]}$$

Fact 9 Fix $k \in [1, \min(n, p)]$. If $W \in \mathbb{R}^{n \times n}$, $U \in \mathbb{R}^{p \times n}$, $V \in \mathbb{R}^{n \times p}$ and $UV = I_p$, then

$$(UWV)^{[k]} = U^{(k)}W^{[k]}(U^{(k)})^{-1}$$

This section aims to provide a clear geometric interpretation of the k-multiplicative compound. The section begins by defining a k-set (the codomain of a function dependent on k variables) before presenting Theorem 8, the key result which exposes the relationship between the volume of a k-set and the k-multiplicative compound of the Jacobian. A k-parallelotope is then shown as an example. Like before, these are existing results, so are presented without proof. Refer to Wu et al. (2022) for more information.

Definition 3 (*k*-sets) Consider a compact set $\mathcal{D} \subset \Re^k$ and a continuous differentiable map $\Psi : \mathcal{D} \to \mathbb{C}$ \Re^n , with $k \in [1, n]$. The codomain of Ψ is given by the parametrised set

$$\Psi(D) := \{\Psi(r) : r \in \mathcal{D}\} \subseteq \Re^n \tag{14}$$

Since \mathcal{D} is compact and $\Psi(\cdot)$ is continuous, $\Psi(\mathcal{D})$ is a closed set.



Figure 2: The 3-parallelotope with vertices $x_1, x_2, x_3 \in \Re^3$ parametrised by the unit cube, \mathcal{D} .

Theorem 8 (Volume of k-sets) Fix $k \in [1, n]$. Consider a compact set $\mathcal{D} \subset \Re^k$ and a continuously differentiable map $\Psi : \mathcal{D} \to \Re^n$. The volume of the parametrised set (14) is given by

$$\operatorname{vol}\left(\Psi(D)\right) = \int_{\mathcal{D}} \left| J_{\Psi}^{(k)}(r) \right| dr$$

where $J_{\Psi}(r) = \begin{bmatrix} \frac{\partial \Psi(r)}{\partial r_1} & \dots & \frac{\partial \Psi(r)}{\partial r_k} \end{bmatrix}$ is the Jacobian of Ψ . Note that as $J_{\Psi} : \mathcal{D} \to \Re^{n \times k}$, it implies $J_{\Psi}^{(k)}(r) : \mathcal{D} \to \Re^s$ with $s = \binom{n}{k}$.

Theorem 8 states that the volume of a k-set is governed by the k-multiplicative compound of the Jacobian. An important geometrical feature, is that for $k \in \{1, 2, 3\}$ the volume of the k-set is equivalent to the standard notions of length, area and volume. We now consider the k-parallelotope as an example of a k-set.

Definition 4 (*k*-parallelotope) Fix $k \in [1, n]$ and let vectors $x_1, \ldots, x_k \in \Re^n$. The parallelotope generated by these vectors (and the zero vertex) is the set given by

$$P(x_1, \dots, x_k) := \left\{ \sum_{i=1}^k r_i x_i : r_i \in [0, 1] \ \forall \ i \right\}$$

Based on Definition 4, the k-parallelotope is a k-set with the following compact domain \mathcal{D} and continuous differentiable function $\Psi(r; x_1, \ldots, x_k)$, as illustrated for k = n = 3 in Figure 2.

$$\mathcal{D} := \left\{ r \in \Re^k : r_i \in [0, 1] \; \forall \; i \in [1, k] \right\} \qquad \Psi(r; x_1, \dots, x_k) := \sum_{i=1}^k r_i x_i$$

The Jacobian of the k-parallelotope is $J_{\Psi}(r) = X$ and from Theorem 8, the volume of the k-parallelotope is given by $|X^{(k)}|$.

A.3 k-CONTRACTION ANALYSIS

852853Consider the time-varying nonlinear system

$$\dot{x} = f(t, x) \tag{15}$$

where $f : \Re_+ \times \Re^n \to \Re^n$. It is assumed throughout that f is continuously differentiable w.r.t x. Fix $k \in [1, n]$ and let S^k denote the unit simplex.

$$\mathcal{S}^k := \{r \in \Re^k : r_i \ge 0 \text{ and } r_1 + \dots + r_k \le 1\}$$

The convex combination of a set of initial conditions $x_1, \ldots, x_{k+1} \in \Re^n$ is defined by $h : S^k \to \Re^n$.

862
863
$$h(r; x_1, \dots, x_{k+1}) := \sum_{i=1}^k r_i x_i + \left(1 - \sum_{i=1}^k r_i\right) x_{k+1}$$

The set $h(S^k)$ is a k-set and can be thought of as a k-dimensional body of states representing initial conditions of (15). We now define $w_i(t, r)$ as a measure of the sensitivity of a solution to (15) at time t, to a change in the initial condition h(r), caused by a change in r_i .

$$w_i(t,r) := \frac{\partial x(t,h(r))}{\partial r_i} \quad \text{where } w_i(0,r) = \frac{\partial h(r)}{\partial r_i} = x_i - x_{k+1} \quad \text{for all } i \in [1,k]$$

We are now ready to present the definition of k-contraction, followed by its geometric interpretation.

Definition 5 (k-contraction) Fix $k \in [1, n]$. The nonlinear system (15) is k-contracting if there exists an $\eta > 0$ and a vector norm $|\cdot|$ such that for any $x_1, \ldots, x_{k+1} \in \Re^n$ and any $r \in S^k$, the mapping $W: \Re_+ \times S^k \to \Re^{n \times k}$ defined by $W(t,r) := [w_1(t,r) \dots w_k(t,r)]$ satisfies

$$|W^{(k)}(t,r)| \le \exp(-\eta t)|W^{(k)}(0,r)| \quad \forall t \in \Re_+$$

To explain the geometric meaning of this definition, pick a domain $\mathcal{D} \subseteq \mathcal{S}^k$ and recall that $h(\mathcal{D})$ is a k-set representing k-dimensional bodies of initial conditions for (15); thus, $x(t, h(\mathcal{D})) :=$ $\{x(t, h(r)) : r \in \mathcal{D}\}\$ is a k-set describing how k-dimensional bodies evolve over time. We now leverage Theorem 8, to show how the volume of these bodies evolves over time when governed by (15).

$$vol\left(x(t,h(\mathcal{D}))\right) = \int_{\mathcal{D}} \left|J_x(t,h(r))^{(k)}\right| dr$$

$$= \int_{\mathcal{D}} \left| \begin{bmatrix} \frac{\partial x(t,h(r))}{\partial r_1} & \dots & \frac{\partial x(t,h(r))}{\partial r_k} \end{bmatrix}^{(k)} \right| dr$$
$$= \int_{\mathcal{D}} \left| W^{(k)}(t,r) \right| dr$$

When (15) is k-contracting, the volume of these bodies is upper bounded by the initial volume scaled by an exponentially decaying term.

$$vol\left(x(t,h(\mathcal{D}))\right) \leq exp(-\eta t) \int_{\mathcal{D}} |W^{(k)}(0,r)| dr$$
$$= exp(-\eta t) \left| \left[(x_1 - x_{k+1}) \quad \dots \quad (x_k - x_{k+1}) \right]^{(k)} \right| \int_{\mathcal{D}} dr$$

Therefore, k-contraction of (15) implies the volume of k-dimensional bodies $x(t, h(\mathcal{D}))$ converges to zero at an exponential rate. This can also be interpreted as the volume of k-dimensional bodies is contracting or converging to a (k-1)-dimensional subspace. Figure 1 in the main paper gives an illustration of (1, 2, 3)-contracting systems where the 1-contracting system converges to a point, the 2-contracting system converges to a line and the 3-contracting system converges to a plane.

Many existing k-contraction results, including Theorem 1, are expressed in terms of matrix measures. An overview of their definitions and properties can be found in (Vidyasagar, 2002, Section 2.2). Theorem 1 provides a sufficient condition for verifying k-contraction in the 2-norm w.r.t a metric P. The 2-norm was chosen for this work due to its relationship with the eigenvalues of its argument, but other norms could be chosen. Furthermore, one may apply an invertible linear transformation Θ to $w_i(t,r)$ and the k-contraction analysis may be performed in this new domain whilst implying the same property holds in the original domain. This idea is made clear in Lohmiller & Slotine (1998) and translates analogously to the k-contraction case. When using such an invertible transformation, the system is said to be k-contracting w.r.t the metric $P = \Theta^{\top} \Theta$.

В LURIE NETWORK

B.1 RELATIONSHIP BETWEEN LURIE NETWORKS AND DEEP FEEDFORWARD MODELS

Consider the vector field $\dot{z} = f(z)$ defined by the following L-layer feedforward network

924		
925	$\dot{z} = W_L \Phi(u_{L-1}) + b_L$	
926	$u_{L-1} = W_{L-1}\Phi(u_{L-2}) + b_{L-1}$	
927	$u_{L-2} = W_{L-2}\Phi(u_{L-3}) + b_{L-2}$	
928		(16)
929	:	
930	$u_2 = W_2 \Phi(u_1) + b_2$	
931	$u_1 = W_1 z + b_1$	
932		

To illustrate the superior expressivity of the Lurie Network, we wish to show that a special case of (2) can approximate the deep feedforward network (16). An alternative expression for (16) is

$$\dot{z} = 0z + W_L \Phi(u_{L-1}) + b_L$$

$$\epsilon \dot{u}_{L-1} = -u_{L-1} + W_{L-1} \Phi(u_{L-2}) + b_{L-1}$$

$$\epsilon \dot{u}_{L-2} = -u_{L-2} + W_{L-2} \Phi(u_{L-3}) + b_{L-2}$$

$$\vdots$$

$$\epsilon \dot{u}_2 = -u_2 + W_2 \Phi(u_1) + b_2$$

$$u_1 = W_1 z + b_1$$

(17)

where $\epsilon \to 0$. Defining a new state $x := \begin{bmatrix} z & u_{L-1} & u_{L-2} & \dots & u_3 & u_2 \end{bmatrix}^{\top}$ and an output vector $y := \begin{bmatrix} u_{L-1} & u_{L-2} & \dots & u_2 & u_1 \end{bmatrix}$ it is clear that (17) is a special case of the Lurie Network (2) with state-space matrices and biases defined by the sparse structures below.

949		0	0	0		0	0]		Γ0	W_L	0		0	0
950		0	-I	0		0	0		0	0	W_{L-1}		0	0
951	$A = \epsilon^{-1}$	0	0	-I		0	0	$B = \epsilon^{-1}$	L 0	0	0	•••	0	0
952		:	:	:	۰.	:	;		:	:	:	٠.	:	:
953		$\dot{0}$	0	0		•	-I		0	0	0		0	\dot{W}_2
954	Г	0	Ţ	0	0	01	-		ГЬ	, T		гол		
955		0	0	I.	$ \begin{array}{ccc} $				b_T	$\begin{bmatrix} L \\ -1 \end{bmatrix}$				
956	C -	Õ	Õ	0.	0	0		$h = e^{-\frac{1}{2}}$	$ b_L $	$\begin{bmatrix} -1 \\ -2 \end{bmatrix}$	h —	0		
957	C –	•	•	• •				$v_x = \epsilon$		-	$v_y =$			
958		: 177	:	:	·· :				1					
959	L	<i>vv</i> 1	0	υ.	0	0]			LU	2		$\lfloor v_1 \rfloor$		

This is just one realisation of a Lurie Network which approximates a feedforward network. Other permutations of the state would result in different realisations of the Lurie Networks weights and biases. Finally, due to the division by ϵ , it would not be possible to train a Lurie Network to have the exact form of a feedforward network; however, this analysis shows that it is possible to approximate the strucuture of a deep feedforward network with a Lurie Network.

972 B.2 FURTHER EXAMPLES

Neural Oscillators: This example is from the graph ML literature Lanthaler et al. (2024). The state 975 of the general neural oscillator is governed by a second order ODE; however, it's equivalent first order 976 representation takes the form (2) with one possible realisation given by $C_{21} \in \Re^{n \times n}$, $b_x = 0$ and

$A = \begin{bmatrix} 0 & I \\ 0 & 0 \end{bmatrix}$	$B = \begin{bmatrix} 0 & 0 \\ 0 & I \end{bmatrix}$	$C = \begin{bmatrix} 0 & 0\\ C_{21} & 0 \end{bmatrix}$	$b_y = \begin{bmatrix} 0\\b_{y2} \end{bmatrix}$
----------------------------------------------------	----------------------------------------------------	--------------------------------------------------------	-------------------------------------------------

980 The solution is then passed through an affine readout layer.

Graph Coupled Oscillators: Another example of a second order ODE from the graph ML literature Rusch et al. (2022). The state is defined as a matrix, but this can simply be recast in a vectorised form which relates to (2) if a linear coupling function is chosen. This requires the weights to have a block matrix form, where one possible realisation is defined by $C_{21} \in \Re^{n \times n}$, $b_x = b_y = 0$ and

$$A = \begin{bmatrix} 0 & I \\ -\gamma I & -\alpha I \end{bmatrix} \qquad \qquad B = \begin{bmatrix} 0 & 0 \\ 0 & I \end{bmatrix} \qquad \qquad C = \begin{bmatrix} 0 & 0 \\ C_{21} & 0 \end{bmatrix}$$

LSSM: When the external input is replaced by nonlinear output feedback (i.e., $u(t) \equiv \Phi(y)$) and D = 0, the linear state-space layer used in S4 Gu et al. (2021) and Hippo Gu et al. (2020) is related to (2) with A being a lower triangular Hippo matrix, $B \in \Re^{n \times m}$, $C \in \Re^{m \times n}$ and $b_x = b_y = 0$.

PROOFS С

C.1 PROOF OF THEOREM 2

We aim to verify Theorem 1 for the particular case where the nonlinear system is described by the Lurie Network (2). Our proof begins with Theorem 9, which restates (Ofir et al., 2024, Theorem 1). This result is sufficient to satisfy Theorem 1 for systems of the form (18).

$$\dot{x} = \bar{A}x(t) - \bar{B}\Psi(t, y) \qquad \qquad y = \bar{C}x \tag{18}$$

If there exists $\eta_1, \eta_2 > 0$ and an invertible $\Theta \in \Re^{n \times n}$ such that

Theorem 9 Fix $k \in [1, n]$ and consider the system below.

$$P^{(k)}\bar{A}^{[k]} + (\bar{A}^{[k]})^{\top}P^{(k)} + \Theta^{(k)} \Big((\Theta \bar{B}\bar{B}^{\top}\Theta)^{[k]} + (\Theta^{-1}\bar{C}^{\top}\bar{C}\Theta^{-1})^{[k]} \Big) \Theta^{(k)} \preceq -\eta_1 P^{(k)}$$
(19)

and

$$\left(\Theta^{-1}\bar{C}^{\top}(J_{\Psi}^{\top}(t,y)J_{\Psi}(t,y)-I_m)\bar{C}\Theta^{-1}\right)^{[k]} \preceq -\eta_2 I_s \quad \forall \ t \in \Re_+ \ and \ y \in \Re^m$$
(20)

where $s = \binom{n}{k}$, then (18) is k-contracting in the 2-norm w.r.t the metric $P := \Theta^{\top} \Theta$.

We first need to express the Lurie Network in the form (18). By (3) there exists $\gamma < 0$ satisfying

$$0 < \gamma^2 < \alpha_k^2 \quad \text{and} \quad g^2 \sum_{i=1}^k \sigma_i^2(B) \sigma_i^2(C) < \gamma^2 k \tag{21}$$

Using γ , we can express (2) in the form (18) through the definitions below, where the dependence on t has been dropped from Ψ .

$$\bar{A} := A \qquad \bar{B} := \gamma I_n \qquad \bar{C} := I_n \qquad \Psi(x) := -\gamma^{-1} B \Phi(Cx + b_y) - \gamma^{-1} b_x \tag{22}$$

The next step is to verify (19). Subbing (22) into the left hand side of (19) and assuming $\Theta = \Theta^{\perp}$ results in the first equality. Setting $P := pI_n$ with p > 0 results in the second. Now we must leverage some of the facts presented in A.1. Using the relevant special cases from (12) and (13) leads to equality three and consequently applying Fact 6 and Fact 8 results in equality four. Re-applying (13) and Fact 8 results in the final equality.

$$= P^{(k)}A^{[k]} + (A^{[k]})^{\top}P^{(k)} + \Theta^{(k)}\left((\gamma^{2}P)^{[k]} + (P^{-1})^{[k]}\right)\Theta^{(k)}$$

$$= (pI_{n})^{(k)}A^{[k]} + (A^{[k]})^{\top}(pI_{n})^{(k)} + (p^{\frac{1}{2}}I_{n})^{(k)}\left((\gamma^{2}pI_{n})^{[k]} + (p^{-1}I_{n})^{[k]}\right)(p^{\frac{1}{2}}I_{n})^{(k)}$$

$$= (pI_{n})^{(k)}A^{[k]} + (A^{[k]})^{\top}(pI_{n})^{(k)} + (p^{\frac{1}{2}}I_{n})^{(k)}\left((\gamma^{2}pI_{n})^{[k]} + (p^{-1}I_{n})^{[k]}\right)(p^{\frac{1}{2}}I_{n})^{(k)}$$

$$= p^k (A^{[k]} + (A^{[k]})^{\top}) + k(\gamma^2 p + p^{-1})p^k I_s$$

= $p^k ((A + A^{\top})^{[k]} + k(\gamma^2 p + p^{-1})I_s)$

$$= p^{n} \left((A + A^{+})^{n} + k(\gamma^{2}p + p) \right)$$

$$= p^{k} \left(A + A^{+} + (\gamma^{2} p + p^{-1}) I_{n} \right)^{[\kappa]}$$

If the matrix above is negative definite, then (19) is satisfied for some suitably chosen $\eta_1 > 0$. This is true when

$$(A + A^{\top} + (\gamma^2 p + p^{-1})I_n)^{[k]} \prec 0$$

By Fact 7, the inequality above can be equivalently expressed as a condition on the sum of the klargest eigenvalues of the matrix inside the k-compound operator. Leveraging (Petersen et al., 2008, Eq. 285) allows us to separate p from the eigenvalues of the symmetric component of A, resulting in the equality below.

1075
1076
1077
$$\sum_{i=1}^{k} \lambda_i \left(A + A^\top + (\gamma^2 p + p^{-1}) I_n \right) = k(\gamma^2 p + p^{-1}) + \sum_{i=1}^{k} \lambda_i (A + A^\top) < 0$$
1077

By the definition of α_k in Theorem 2, this simplifies to

$$\gamma^2 p^2 + 2\alpha_k p + 1 < 0$$

For for γ satisfying (21), the quadratic inequality always emits at least one solution $p = -\alpha_k^{-1}$. The final step is to verify (20). The Jacobian of $\Psi_{\rm c}$ as defined in (22), is

The final step is to verify (20). The Jacobian of
$$\Psi$$
, as defined in (22),

$$J_{\Psi}(x) = -\gamma^{-1}BJ_{\Phi}C$$

For $\Theta = p^{\frac{1}{2}} I_n$ and the definitions from (22), the left hand side of (20) reduces to

$$= \left(p^{-1} J_{\Psi}^{\top} J_{\Psi} - p^{-1} I_n \right)^{[k]}$$

If the matrix above is negative definite, then (20) is satisfied for some suitably chosen $\eta_2 > 0$. Repeating the same steps as above, this negative definite requirement reduces to the inequality below.

$$\sum_{i=1}^{k} \lambda_i (p^{-1} J_{\Psi}^{\top} J_{\Psi}) - k p^{-1} = p^{-1} \sum_{i=1}^{k} \sigma_i^2 (J_{\Psi}) - k p^{-1} < 0$$

Subbing in the definition of J_{Ψ} and applying the well-known property of singular values (Horn & Johnson, 1994, Theorem 3.3.14), then (20) is true if

$$\gamma^{-2} \sum_{i=1}^k \sigma_i^2(B) \sigma_i^2(J_\Phi) \sigma_i^2(C) < k$$

By the assumption made on the slope of Φ , this inequality will always be satisfied if (3) holds.

C.2 PROOF OF THEOREM 3

For this proof, we aim to directly verify Theorem 1 for $\Theta \in \mathcal{D}^n$ where f represents the Lurie Network (2). We start by substituting the Jacobian of the Lurie Network into the left hand side of (1), followed by the application of Fact 8 to obtain the second equality. The subadditivity property of the matrix measure μ_2 was then leveraged to split the terms (Vidyasagar, 2002, Section 2.2). As the second term is difficult to manipulate, we rely on the simplifying assumption $C = B^{-1}$ in order to apply Fact 9. As $\Theta, B, J_{\Phi} \in \mathcal{D}^n$, so are both of their k-compound counterparts (12) (13), which means the Θ, B terms cancel out. We then apply the property $\mu_2(\cdot) \leq ||\cdot||_2$ (Vidyasagar, 2002, Theorem 16) and Fact 7, which allows us to leverage the slope restricted assumption on Φ . By the relevant special case of the k-additive compound (13), we can directly calculate the 2-norm. Finally, kg can be incorporated in μ_2 as shown in the final inequality.

....

$$\mu_{2,\Theta^{(k)}}(J_f^{[k]}(t,x)) = \mu_{2,\Theta^{(k)}}\left((A + BJ_{\Phi}C)^{[k]}\right)$$

1117
$$= \mu_{2,\Theta^{(k)}} \left(A^{[k]} + (BJ_{\Phi}(y)C)^{[k]} \right)$$

1119
$$\leq \mu_{2,\Theta^{(k)}}(A^{[k]}) + \mu_{2,\Theta^{(k)}}((BJ_{\Phi}C)^{[k]})$$

1120
$$= \mu_{2,\Theta^{(k)}}(A^{[k]}) + \mu_2(\Theta^{(k)}B^{(k)}J^{[k]}_{\Phi}B^{-(k)}\Theta^{-(k)})$$

1121
1122
$$= \mu_{2,\Theta^{(k)}}(A^{[k]}) + \mu_2(J_{\Phi}^{[k]})$$
(1122 (112)

1123
$$\leq \mu_{2,\Theta^{(k)}}(A^{[k]}) + ||(gI_n)^{[k]}||_2$$

1124
1125
$$= \mu_{2,\Theta^{(k)}}(A^{[k]}) + kg$$

1126
$$= \mu_2(\Theta^{(k)}A^{[k]}\Theta^{-(k)} + kgI_n)$$

If the final inequality is negative, then Theorem 1 will be satisfied for some suitably chosen η . This is equivalent to the matrix inequality below.

1131
1132

$$\frac{1}{2} \left(\Theta^{(k)} A^{[k]} \Theta^{-(k)} + \Theta^{-(k)} (A^{[k]})^\top \Theta^{(k)} + 2kgI_n \right) \prec 0$$

Multiplying on the left by $2\Theta^{(k)}$ and on the right by $\Theta^{(k)}$ results in (4).

1134 C.3 PROOF OF THEOREM 4 1135

1136 The aim of this proof is to verify Theorem 1 when f is the GLN (6). We begin by expressing the Jacobian as a sum of the Jacobian of the q independent Lurie Networks, J_{indep} , and Jacobian of the 1137 coupling term J_{couple} 1138

 $J_f(x) = J_{indep}(x) + J_{couple}(x)$ 1139 where $J_{indep}(x) = A_G + B_G J_{\Phi} C_G$ and $J_{couple}(x) = L$. Subbing J_f into the left hand side of (1) 1140 and applying the subadditivity property of μ_2 results in 1141

1142
1143
$$\mu_{2,\Theta^{(k)}}(J_f^{[k]}) \le \mu_{2,\Theta^{(k)}}(J_{\text{indep}}^{[k]}) + \mu_{2,\Theta^{(k)}}(J_{\text{couple}}^{[k]})$$

As we assume the q independent Lurie Networks are k-contracting in the 2-norm w.r.t the metric 1144 $P = \text{blockdiag}(P_1, \dots, P_q)$, we know that $\mu_{2,\Theta^{(k)}}(J_{\text{indep}}^{[k]}) < 0$. Under this assumption, Theorem 1 1145 1146 is satisfied if $\mu_{2,\Theta^{(k)}}(J_{\text{couple}}^{[k]}) \le 0$ 1147

1150 1151 1152

1169 1170 1171

1181 1182

This is equivalent to the matrix inequality below, where J_{couple} has been subbed in. 1149

$$\frac{1}{2} \left(\Theta^{(k)} L^{[k]} \Theta^{-(k)} + \Theta^{-(k)} (L^{[k]})^T \Theta^{(k)} \right) \leq 0$$

1153 Multiplying on the left by $2\Theta^{(k)}$ and on the right by $\Theta^{(k)}$ results in (7). 1154

1155 C.4 PROOF OF THEOREM 5 AND THEOREM 6 1156

1157 The aim of this proof is to express the weights of the Lurie Network (2) such that Theorem 2 is always satisfied. More formally, this requires $A, B, C \in \Omega_2(g, k)$ to always hold. As both theorems share 1158 the same proof until the final step, where Σ_A is defined, they have been combined into one proof. 1159

1160 To expose the singular values of B, C, we leverage the singular value decomposition, as in (9b) and (10a). This requires the matrices U_B, U_C, V_B, V_C to be orthogonal. We can immediately use 1161 the unconstrained parametrisation of the orthogonal class from Lezcano-Casado & Martinez-Rubio 1162 (2019) to express these matrices as unconstrained symmetric matrices. The matrices Σ_B, Σ_C contain 1163 the singular values of the respective matrix, hence $\Sigma_B \in \mathcal{D}^{nm}_+$ and $\Sigma_C \in \mathcal{D}^{mn}_+$. We also treat these as 1164 unconstrained sets since any element can be obtained by taking the absolute value of an unconstrained 1165 diagonal matrix with the same shape. 1166

To verify Theorem 2, we can combine $\alpha_k < 0$ and (3) into one inequality representing the intersection 1167 of the two sets. 1168

$$\sum_{i=1}^{k} \lambda_i (A + A^T) < -2k \sqrt{\frac{z}{k}} \quad \text{where } z := g^2 \sum_{i=1}^{k} \sigma_i^2(B) \sigma_i^2(C)$$
(23)

1172 Thanks to the definition of B and C, the right hand side is a function of the hyperparameters g, k and elements of the parameters Σ_B, Σ_C , so can be easily computed using sort and sum functions. 1173

1174 To impose this constraint directly on the eigenvalues of the symmetric component of A, we express A1175 as a sum of symmetric and skew-symmetric matrices. The skew-symmetric matrix is unconstrained, 1176 so this can be left as it is. Finally, we apply the eigenvalue decomposition of a symmetric matrix to obtain (9a) (10a). This expression allows us to directly place the constraint above on the diagonal 1177 matrix Σ_A . 1178

1179 This is where Theorem 5 and Theorem 6 differ. Defining Σ_A as in (9a) ensures 1180

$$\lambda_i(A + A^{\top}) < -2\sqrt{\frac{z}{k}} \quad \text{for all } i \in [1, n]$$

1183 which guarantees both conditions of Theorem 2 will be satisfied; however, conservatism is introduced 1184 as all the eigenvalues must be negative. Theorem 6 was established to address this issue. Defining 1185 Σ_A as in (10b) guarantees both conditions of Theorem 2 will be satisfied via (23). The definition of Σ_A is split into one unconstrained block for the first (k-1)-eigenvalues, a block for $\lambda_k (A + A^{\top})$ 1186 which is defined to ensure (23) holds, and finally a block for the remaining eigenvalues which must 1187 be defined to ensure the k eigenvalues involved in (23) are the largest.

C.5 PROOF OF THEOREM 7

This proof aims to show that Theorem 4 is satisfied when L is defined as in (11). First, multiply (7) on the left and right by $\Theta^{-(k)}$. Sequentially applying Fact 6, Fact 9 and Fact 8 results in

Subbing in the definition of L from (11) and recalling that $P = \Theta^\top \Theta$ leads to

$$\left(\Theta L \Theta^{-1} + \Theta^{-1} L^T \Theta\right)^{[k]} \preceq 0$$

 $\left(G_{L2} + G_{L2}^{\top}\right)^{[k]} \preceq 0$

1196	
1197	
1198	
1199	which is assumed to hold.
1200	
1201	
1202	
1203	
1204	
1205	
1206	
1207	
1208	
1209	
1210	
1211	
1212	
1213	
1214	
1215	
1216	
1217	
1218	
1219	
1220	
1221	
1222	
1223	
1224	
1225	
1226	
1227	
1228	
1229	
1230	
1231	
1232	
1233	
1234	
1235	
1236	
1237	
1238	
1239	
1240	

1242 D EXTENDED EMPIRICAL EVALUATION

1244 D.1 Dата

The data was synthetically generated by numerically integrating over the analytical models of each dynamical system. The integration was performed using the Euler method with step size $\delta = 1 \times 10^{-2}$.

1249 D.1.1 OPINION DYNAMICS

This model was presented in Ofir et al. (2024). It has the following state space equations

$$\dot{x} = -1.5I_3 + 0.5\Phi(Cx) + b$$

1253 where 1254

$C = \begin{bmatrix} +1\\ -1\\ 0 \end{bmatrix}$	-1 + 1 - 1	$\begin{bmatrix} 0\\ -1\\ +1 \end{bmatrix}$	b =	$\begin{bmatrix} +0.2\\0\\-0.2\end{bmatrix}$
ΓU	-1	+1]		$\lfloor -0.2 \rfloor$

and the tanh function is applied element-wise. The model is 1-contracting and has a unique equilibrium point at b.

1261 D.1.2 HOPFIELD NETWORK

This model is a variation of the Hopfield network presented in Ofir et al. (2024). It has the following state space equations

$$\dot{x} = -2.5I_3 + B\Phi(x)$$

1265 where

and the tanh function is also applied element-wise. The model is 2-contracting and has two stable equilibrium points: $e_1 = \begin{bmatrix} 0.79 & 0.79 & 0.79 \end{bmatrix}^{\top}$, $e_2 = -e_1$; and an unstable equilibrium point $e_3 = 0$.

 $B = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$

1273 D.1.3 SIMPLE ATTRACTOR

¹²⁷⁴ This model was presented in Cecilia et al. (2023). It has the following state space equations

$$\dot{x} = Ax + B\Phi(Cx)$$

where

ΓO 1	-2]	Г	0	0	[0		ΓO	0	[0
$A = \begin{bmatrix} -1 & 0 \end{bmatrix}$	-1	B =	0	0	0	C =	0	0	0
$\begin{bmatrix} 0.5 & 0 \end{bmatrix}$	-0.5	L	-0.5	0	0		$\lfloor 1$	0	0

and $\phi(z) = z^3$ is the nonlinearity applied element-wise. This function is not slope-restricted, so the simple attractor does not satisfy the assumptions of the Lurie Network. The model is 3-contracting and has several attractor states.

1296 D.2 TRAINING

The default training settings common to the opinion, Hopfield and attractor datasets are presented in Table 5. The only parameters which varied between models were the learning rate and epoch which it was cut. Deviations from the default settings are detailed in Table 6. These values were chosen based on observations during training. No hyperparameter sweep was performed. The same details are also presented for the graph-coupled Hopfield network and graph-coupled attractor datasets in Table 7 and Table 8. When training the models for the opinion, Hopfield and attractor datasets, a single T4 GPU (accessed through Google Colab) was used. A single A100 GPU (also accessed through Google Colab) was used for training the models on the graph-coupled Hopfield and graph-coupled attractor datasets. Training curves from the opinion, Hopfield and attractor datasets are also plotted to compare the convergence time and variance.

1	3	0	7
1	3	0	8

Parameter	Value	
Batches	10	
Batch size	100	
Test split	0.1	
Epochs	100	
Criterion	Mean squared error	
Optimiser	Adam (default settings)	
Learning rate	1×10^{-2} (no cuts)	

Table 6: Deviations from the default settings for opinion, Hopfield and attractor datasets.

Dataset	Model	Learning Rate	Cut (decay, epoch)
Opinion Dynamics	Lurie Network	5×10^{-3}	N/A
Opinion Dynamics	Antisymmetric RNN	5×10^{-3}	N/A
Hopfield Network	k-Lurie Network	5×10^{-3}	N/A
Hopfield Network	Neural ODE	1×10^{-3}	0.1,75
Hopfield Network	Antisymmetric RNN	5×10^{-3}	N/A
Simple Attractor	Neural ODE	1×10^{-3}	N/A
Simple Attractor	Antisymmetric RNN	5×10^{-3}	N/A

Table 7: Default training settings for the graph-coupled Hopfield and attractor datasets.



1351Table 8: Deviations from the default settings for the graph-coupled Hopfield and attractor datasets.

Figure 3: Mean \pm two standard deviations of the training loss (blue) and test loss (orange) across N = 3 training runs.

1404 D.3 FURTHER RESULTS





Figure 7: Random sample of trajectories from the Hopfield network test set. Predictions are made by the model associated with the best MSE in Table 1. Crosses denote initial conditions, stars denote equilibrium points.



Figure 9: Random sample of 30s trajectories from the noisy, out of distribution Hopfield network dataset. Predictions are made by the model associated with the best MSE in Table 1. Crosses denote initial conditions, stars denote equilibrium points.

1504



1616 initial conditions.



Figure 12: Random sample of 30s trajectories from the noisy, out of distribution simple attractor dataset. Predictions are made by the model associated with the best MSE in Table 1. Crosses denote initial conditions.