

A unified theory of feature learning in RNNs and DNNs

Jan Bauer¹ · Kirsten Fischer² · Moritz Helias² · Agostina Palmigiano¹

¹Gatsby Computational Neuroscience Unit, UCL, London, UK ²IAS-6, Forschungszentrum Jülich, Germany {jan.bauer, a.palmigiano}@ucl.ac.uk

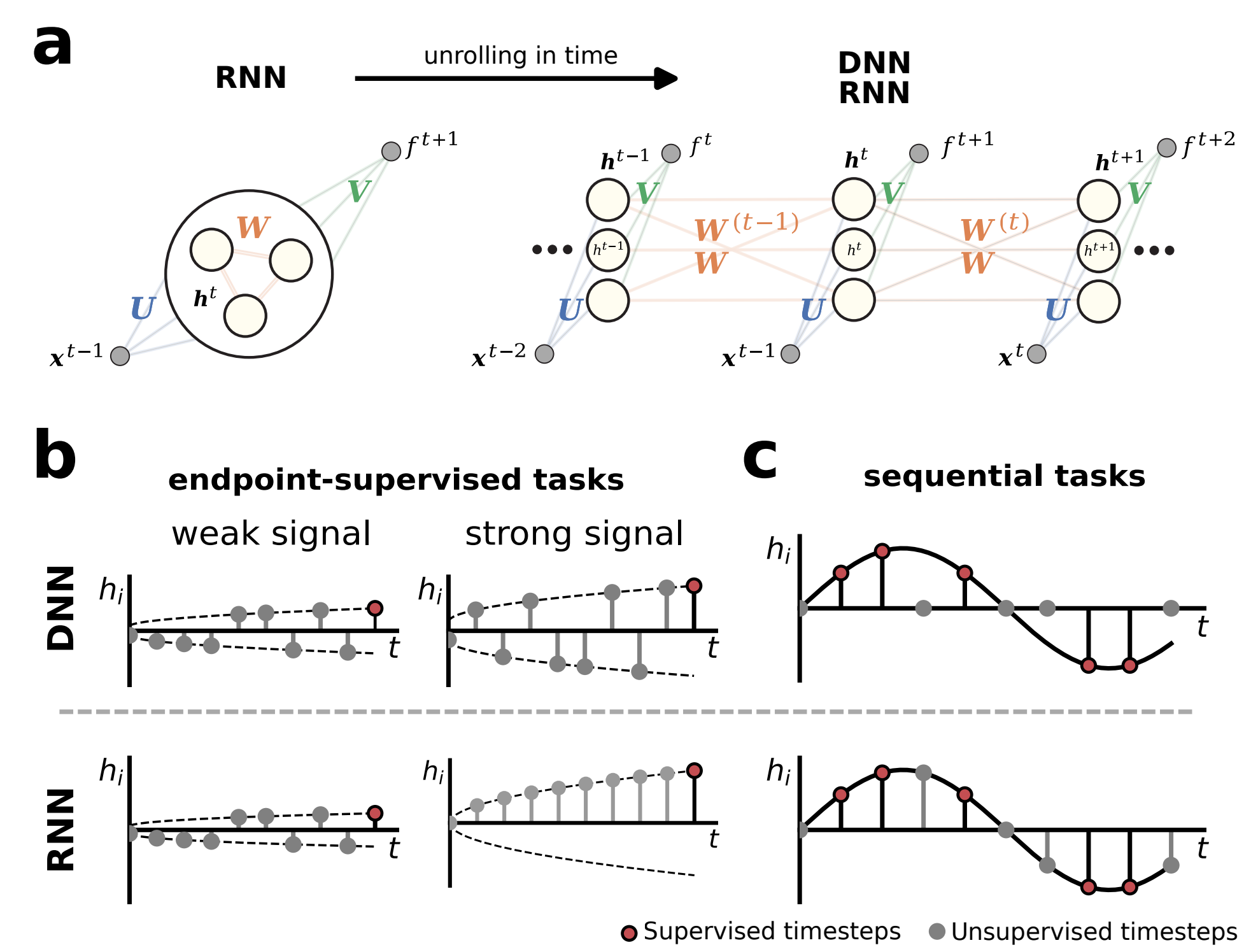
TL;DR: We derive one mean-field theory for trained RNNs and DNNs. Weight sharing, their only difference, becomes a **masking of temporal correlations**: it triggers a **phase transition** to temporally coherent representations and acts as an **inductive bias** that makes RNNs sample-efficient sequence learners.

1 · RNNs are DNNs with shared weights

Unrolled in time, an RNN is a deep network whose layers are timesteps, except that all layers share the same weights \mathbf{W} . Are the DNN's extra parameters an advantage, or is weight sharing a useful inductive bias?

$$h^t = \mathbf{W}^{(t-1)}\phi(h^{t-1}) + \mathbf{U}\mathbf{x}^{t-1}, \quad f^{t+1} = \mathbf{V}\phi(h^t),$$

$$\mathbf{W}^{(t)} = \begin{cases} \mathbf{W}, & \text{RNN (shared)} \\ \mathbf{W}^{(t)}, & \text{DNN (untied)} \end{cases}$$



a) Unrolling-in-time. b) Endpoint supervision: only the RNN turns coherent in time. c) Sequential tasks: weight sharing interpolates unsupervised timesteps.

2 · Training = Bayesian inference

Noisy gradient descent (SGLD: gradients + weight decay K/G_θ + isotropic noise $\sqrt{2K}$) converges to a stationary posterior over weights,

$$\theta \sim P(\theta|y, \mathbf{x}) \propto \underbrace{e^{-P|\mathcal{T}|\mathcal{L}(\theta; y, \mathbf{x})/K}}_{\text{likelihood (loss)}} \underbrace{\mathcal{N}(\theta|0, G_\theta)}_{\text{prior (weight decay)}}$$

Requiring $\mathcal{O}(1)$ kernels as $N \rightarrow \infty$ recovers μP feature-learning scaling, $(G_U, G_W, G_V) = (u/D, w/N, v/N^2)$.

3 · One kernel theory, two architectures

Weights enter the forward pass linearly \Rightarrow integrate them out exactly; activations follow a posterior governed by the kernel

$$\Phi_{pp'}^{tt'} = \frac{1}{N} \sum_i \phi(h_{i,p}^t) \phi(h_{i,p'}^{t'}) \approx \langle \phi(h_p^t) \phi(h_{p'}^{t'}) \rangle:$$

$$P(h|y, \mathbf{x}) \propto \exp\left\{-\frac{1}{2}\text{tr}[\mathbb{Y}\mathcal{T}(v\Phi\mathcal{T} + \kappa)^{-1}] - \frac{1}{2}h^T(w[\Phi^-] + u\mathbb{X}^-)^{-1}h + \phi(h)^T \tilde{\Phi} \phi(h)\right\}$$

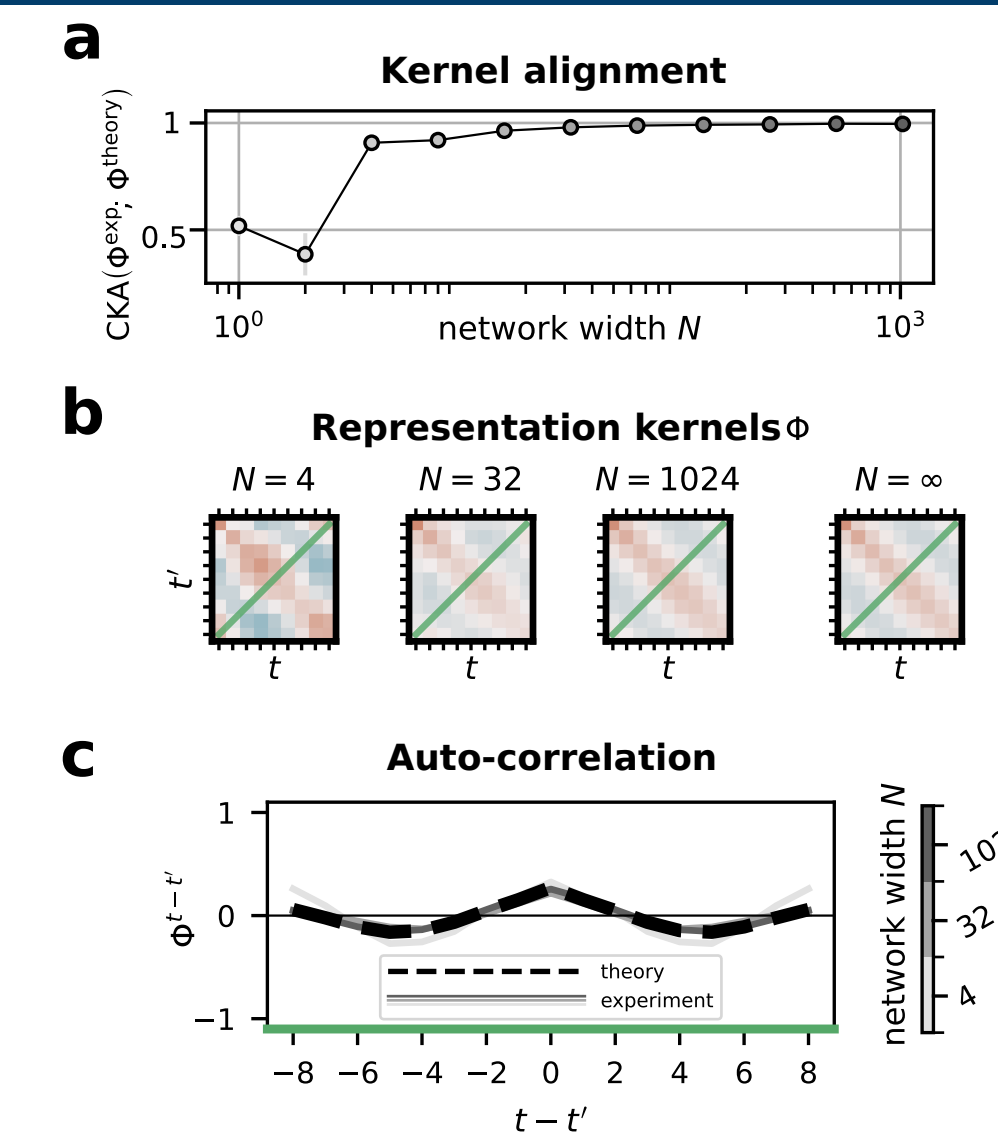
Architecture enters **only** through a masking of temporal correlations,

$$[\Phi^-] = \begin{cases} \Phi^-, & \text{RNN} \\ \text{diag}(\Phi^-), & \text{DNN} \end{cases}$$

shared weights are perfectly correlated in time \Rightarrow coherence can build up; independent weights force it to zero.

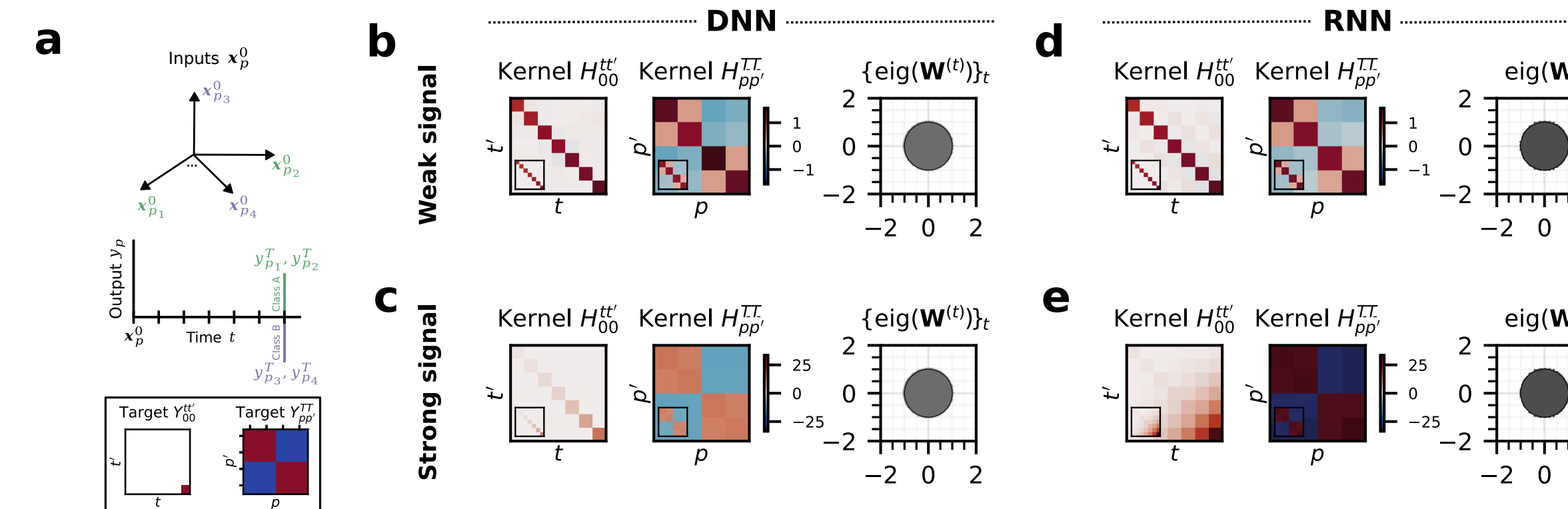
4 · The theory describes trained networks

Kernels from SGLD-trained RNNs converge to the theory's prediction as the width N grows (CKA $\rightarrow 1$); temporal coherence and autocorrelation are captured already at finite N .



5 · Endpoint supervision: a phase transition to temporal coherence

Input at $t = 0$, labels only at the last timestep (the classic DNN setting). **Both** architectures learn the task: $\mathbb{H}^{T-T} \simeq \mathbb{Y}$, giving similar predictors f .



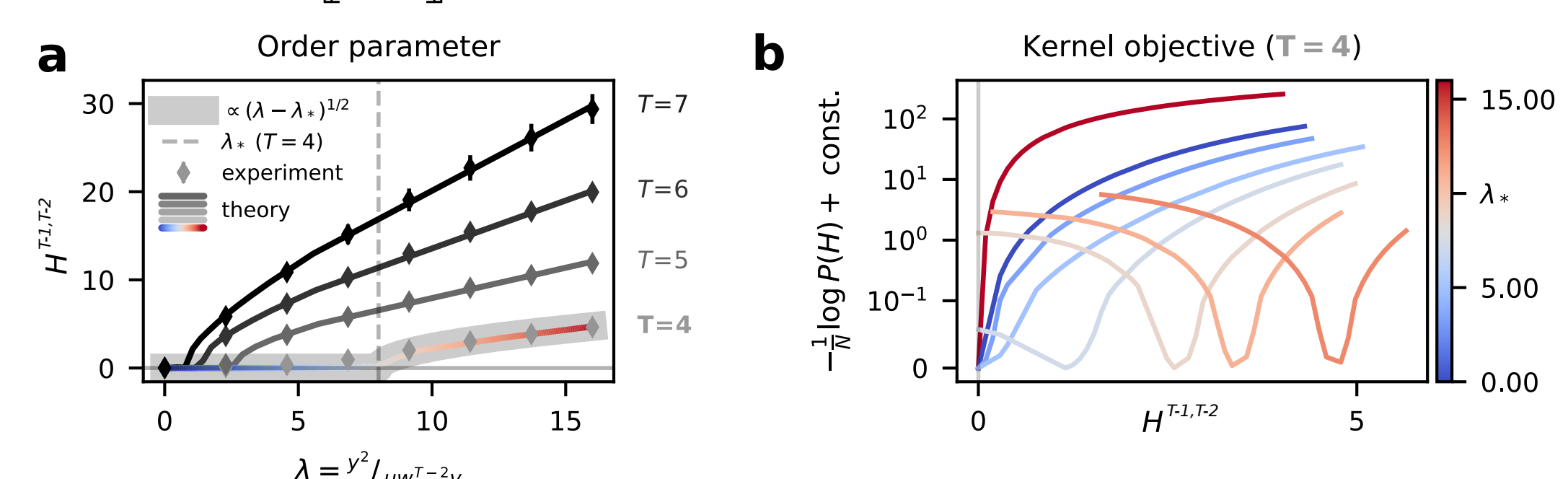
Weak signal (top): RNN and DNN representations coincide. Strong signal (bottom): only the RNN develops temporal coherence, with an outlier eigenvalue in \mathbf{W} .

Beyond a critical learning signal $\lambda = \|\mathbf{y}\|^2/vwT^{-2}u$, the ratio of label strength to initialization variance, only the RNN crosses a **second-order phase transition** (critical exponent $\frac{1}{2}$, analytic at $T = 4$) into a temporally coherent phase, reminiscent of a BBP transition. For linear activation, the kernel posterior is explicit:

$$-\ln P(\mathbb{H}|y, \mathbf{x})/N = \frac{1}{2}\text{tr}[\mathbb{Y}\mathcal{T}(v\mathbb{H}\mathcal{T} + \kappa)^{-1}] + \frac{1}{2}\text{tr}[\mathbb{H}(w[\mathbb{H}^-] + u\mathbb{X}^-)^{-1}] - \frac{1}{2}\ln \frac{|\mathbb{H}|}{|w[\mathbb{H}^-] + u\mathbb{X}^-|}$$

D_{KL} to the random-weight forward pass (entropy)

Aligning \mathbb{H}^{T-T} to the labels lowers the energy and costs entropy; the RNN can pay through off-diagonals in $[\mathbb{H}^-]$, the DNN cannot.

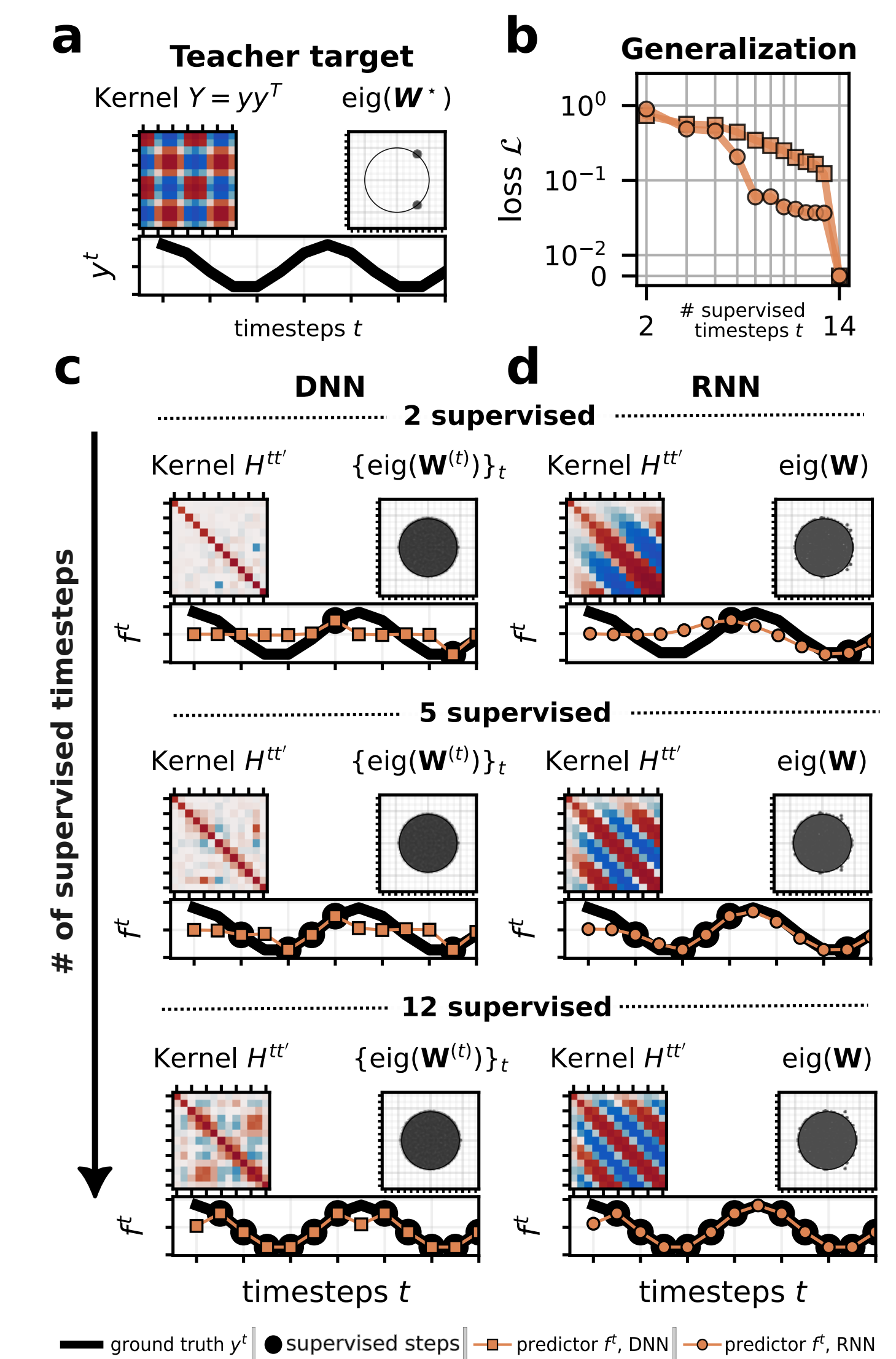


Kernel off-diagonal \mathbb{H}^{T-T-2} as order parameter: theory (lines) vs. weight-space SGLD (diamonds); the kernel landscape $-\ln P(\mathbb{H}|y, \mathbf{x})$ develops a coherent minimum.

Funding: Gatsby Charitable Foundation (GAT3850), Simons Foundation (1156607), DFG (368482240/GRK2416, SPP 2205 (533396241), HE 9032/4-1), Helmholtz ACA (SO-092); compute on JURECA, FZ Jülich (JINB33).

6 · Sequential tasks: weight sharing is an inductive bias

Sequence regression $y^t = \mathbf{V}^*(\mathbf{W}^*)^t \mathbf{U}^* \mathbf{x}^0$ (teacher rotation), supervising only a subset of timesteps:



RNN: kernel matches \mathbb{Y} from few supervised steps; \mathbf{W} recovers the teacher eigenvalues. DNN: posterior stays diagonal \Rightarrow falls back to the prior mean $f^t = 0$ outside supervision, despite larger expressivity.

Mechanism (perturbation in \mathbb{Y}): kernel inverses act as *propagators* that pass label messages across unsupervised timesteps,

$$\Delta_{2,\text{RNN}}^{tt'} \propto \frac{1}{h_0^{t'}} \left(\sum_{t''} \mathbb{Y}_{+tt''} \frac{1}{h_0^{t''}} \mathbb{Y}_{+t''t'} \right) \frac{1}{h_0^t}$$

interpolating along paths $t_3 \leftarrow t_2 \leftarrow t_1$. The DNN's masking cuts these paths \Rightarrow learning needs more samples.

7 · Takeaways

- ▶ One Bayesian kernel theory covers trained RNNs **and** DNNs; the architecture enters as a masking of temporal correlations.
- ▶ Strong learning signals drive RNNs, and only RNNs, through a phase transition to temporally coherent representations.
- ▶ In sequential tasks, weight sharing induces **task-model alignment**: the learned kernel matches the task structure \mathbb{Y} , giving better generalization from fewer labels. **Structure** \rightarrow **function**.



Paper & code