

Reviewed Preprint

v1 • August 20, 2025

Not revised

Reviewed Preprint

v2 • April 24, 2026

Revised by authors

✉ For correspondence:

pgribble@uwo.ca

† co-senior authors

‡ Present address: Dept. Organismic and Evolutionary Biology, Harvard University, Boston, United States

Competing interests: No competing interests declared

Funding: See page 18

Reviewing editor: Rui Ponte Costa, University of Oxford, United Kingdom

© 2025, Shahbazi et al. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

A Context-Free Model of Savings in Motor Learning

Mahdiyar Shahbazi^{1,†}, Olivier Codol^{2,3}, Jonathan A Michaels^{4,5,‡}, Paul L Gribble^{1,4,‡} ✉

¹Dept. Psychology, Western University, London, Canada • ²Mila–Québec Artificial Intelligence Institute, Montréal, Canada • ³Dept. Neuroscience, Université de Montréal, Montréal, Canada • ⁴Dept. Physiology & Pharmacology, Schulich School of Medicine & Dentistry, London, Canada • ⁵School of Kinesiology and Health Science, Faculty of Health, York University, Toronto, Canada

eLife Assessment

This study presents **valuable** computational findings on the neural basis of learning new motor memories and the savings using recurrent neural networks. The evidence supporting the claims of the authors is **solid**, but it would benefit from more detailed discussion on the specific conditions under which savings emerges from purely implicit mechanisms. This work will be of interest to computational and experimental neuroscientists working in motor learning.

<https://doi.org/10.7554/eLife.107423.2.sa2>

Abstract

Learning to adapt voluntary movements to an external perturbation, whether mechanical or visual, is faster during a second encounter than during the first. The mechanisms underlying this phenomenon, known as savings, remain unclear. Recent studies propose that the high dimensionality of neural control enables the retention of learning traces that may facilitate savings. To test this idea we used MotorNet, a framework for training recurrent neural networks (RNNs) to control biomechanical models of the human upper limb. RNNs were trained to perform reaching movements with a velocity-dependent force field (FF) and without (NF) in the sequence NF1 (baseline), FF1 (adaptation), NF2 (washout), and FF2 (re-adaptation). RNNs showed behavioural signatures of savings in the absence of any explicit contextual input signalling the presence or absence of the FF. Savings was more robust in RNNs with larger numbers of units. We identified a component of RNN activity associated with savings—a shift in preparatory activity that persisted even after washout. Displacing this preparatory activity in the direction of the shift enhanced savings, whereas perturbations in the opposite direction reduced or eliminated savings. These findings suggest a potential neural basis for motor memory retention underlying savings that is reliant on the high dimensionality of neural circuits for control, and is independent of cognitive or strategic learning.

Introduction

In studies of motor learning, *savings* commonly refers to a phenomenon in which learning is superior after previous exposure to a motor task. Savings has been demonstrated in the context of voluntary reaching movements for adaptation to novel visuomotor perturbations and for learning to counter novel mechanical environments such as velocity-dependent curl force fields (FF) (Haith et al., 2015 [↗](#); Morehead et al., 2015 [↗](#); Leow et al., 2016 [↗](#); Nguyen et al., 2019 [↗](#); Yin and Wei, 2020 [↗](#); Hadjiiosif et al., 2023 [↗](#); Coltman et al., 2019 [↗](#)). In a typical experiment, participants are initially exposed to a novel perturbation environment and they practice reaching to targets until an asymptotic level of performance is reached, for example recovery of approximately straight-line hand paths. Following this initial learning the perturbation is removed, and participants practice again until behavioural performance in this “washout” phase returns to pre-learning

baseline levels. After washout participants are re-exposed to the perturbing environment. Savings is observed as a faster learning rate when re-exposed to the perturbing environment compared to initial learning, and sometimes also as a superior initial level of performance compared to when the perturbing environment was first encountered (Coltman et al., 2019 [↗](#); Herzfeld et al., 2014 [↗](#)).

There is some debate about the mechanisms that may be responsible for savings (Leow et al., 2016 [↗](#); Yin and Wei, 2020 [↗](#)). Some propose that savings is produced by explicit, cognitive or strategic processes such as a conscious memory of the action selection strategy (Morehead et al., 2015 [↗](#)), a contextual signal associated with the perturbing environment (Heald et al., 2021 [↗](#)), meta-learning of adaptation parameters such as learning rate (Zarahn et al., 2008 [↗](#); McDougle et al., 2015 [↗](#); Albert and Shadmehr, 2018 [↗](#)), or reinforcement-based memories of successful execution (Huang et al., 2011 [↗](#)). Others have proposed that savings may arise from implicit learning processes that are not based on conscious, strategic mechanisms, including an increased sensitivity to previously experienced errors (Herzfeld et al., 2014 [↗](#)), use-dependent plasticity (Diedrichsen et al., 2010 [↗](#)), or implicit updating of internal models that predict the sensory consequences of action (Wolpert et al., 1995 [↗](#)).

Recent advances have been made in the ability to record from large numbers of neurons during motor learning tasks (Trautmann et al., 2025 [↗](#)) and this has resulted in new approaches to understanding the relationship between high dimensional neural population activity and sensory, motor and task parameters during, and even prior to, voluntary movement (Kobak et al., 2016 [↗](#); Dubreuil et al., 2022 [↗](#)). Sun, O’Shea, and colleagues recorded neural activity in primary motor cortex of rhesus macaques during a FF reaching task and identified a neural subspace of network activity during the preparatory period prior to movement that shifted after learning (Sun et al., 2022 [↗](#)). This “uniform shift” persisted even after behavioural washout of FF learning. The authors proposed that this neural trace of prior learning could facilitate subsequent savings (Sun et al., 2022 [↗](#)). Losey and colleagues used a brain-computer interface learning paradigm to study how neural population activity in the primary motor cortex of monkeys supports motor learning of multiple tasks (Losey et al., 2024 [↗](#)). They identified a change in a subspace of neural population activity that supported behavioural performance of a new task without interfering with a previously learned task. They proposed that the high dimensionality of neural activity in primary motor cortex allows for the formation of memory traces that supports multiple behaviours without interference.

In the present paper we used artificial recurrent neural network (RNN) models of upper limb motor control to test the idea that high dimensional neural control facilitates the encoding of multiple novel motor memories, and that neural traces of previous learning underlie subsequent savings, without the need for contextual signals. We used MotorNet (Codol et al., 2024b [↗](#)) to train RNNs to control a mathematical model of the upper limb neuromuscular system (Kistemaker et al., 2010 [↗](#)) in the context of a simulated FF reaching task (Shadmehr and Mussa-Ivaldi, 1994 [↗](#); Conditt et al., 1997 [↗](#)). Even without any explicit contextual cue signalling the presence of absence of FFs, RNNs showed behavioural signatures of savings. In addition savings was more robust as the number of units in RNNs increased. Using approaches similar to those described in previous studies we identified a learning-related shift in neural activity in the preparatory period prior to movement on-set (Sun et al., 2022 [↗](#); Losey et al., 2024 [↗](#)). We established a causal relationship between this neural shift and savings by perturbing neural activity along the direction of this neural shift. When RNN hidden unit activity was shifted in the direction of the neural shift, savings was enhanced, whereas neural perturbations in the opposite direction reduced or abolished savings. Our findings support the hypothesis that a neural basis of motor memory retention underlies savings, one that could be independent of cognitive or strategic learning components and that depends upon the high dimensionality of neural population activity.

Results

We trained 40 recurrent neural networks (RNNs) with 128 fully connected gated recurrent units (GRUs) to control a mathematical model of the human upper limb (Codol et al., 2024b [↗](#)) (Figure 1a,b [↗](#)). Task inputs to the RNN are the Cartesian coordinates of the movement target (x, y) and a

binary go signal (0 or 1) indicating when to initiate movement (Figure 1c). The RNN also receives time-delayed feedback signals corresponding to the length and velocity of each muscle, and the Cartesian coordinates of the endpoint of the limb. The output of the RNN is time-varying muscle stimulation commands to each of 6 upper limb muscles (Figure 1c). Muscle stimulation commands range between 0 and 1 and act on a musculoskeletal model of the upper limb which includes multi-joint limb dynamics and a hill-type muscle model (Kistemaker et al., 2010).

The networks were initially trained to produce point to point reaching movements between targets located in random locations throughout the limb's workspace. No perturbing FF was applied during this initial "growing up" training phase. We refer to the absence of a perturbing FF as a "null field" (NF). RNN weights were updated using backpropagation through time (Werbos, 1990), using the Adam optimizer (Kingma and Ba, 2014) implemented in PyTorch (Paszke et al., 2019). The loss function for optimization was based on minimizing the difference between hand position and target position, and also included regularization terms that encouraged the network to produce smooth, human-like kinematics, phasic muscle commands, and stable hidden unit activity (Michaels et al., 2020; Sussillo et al., 2015) (see Figure 1e and Methods for details).

After training, the networks produced reaching movements with human-like characteristics including smooth, relatively straight hand paths with bell-shaped velocity profiles, and phasic activity in agonist and antagonist muscles spanning the shoulder and elbow (Figure 1c).

Figure 1d shows examples of reaching trajectories for reaches to targets located randomly across the limb's workspace. Models produced human-like reaches both when tested on reaches with random starting points and targets (Figure 1d) and when tested on centre-out reaches toward 8 equidistant targets (Figure 1e). Consistent with electrophysiological recordings in monkeys, RNN hidden units showed activity patterns that were relatively stable over time, and distinct for different movement targets during the delay period prior to the go signal (time 0 in Figure 1e). RNN hidden unit activity during movement was similarly distinct for different movement directions, and showed oscillatory activity consistent with that seen in recordings from motor cortex in non-human primates (Churchland et al., 2012; Churchland and Shenoy, 2024) (also see Figure 6).

Force field adaptation

After the networks were trained to perform point to point reaches in a NF, we implemented a relatively standard experimental sequence common in studies of human motor learning. We trained networks on a centre-out reaching task either in the absence of perturbing forces (null field, NF) or in the presence of a velocity-dependent curl force field (FF) (see Methods). First, networks were exposed to a NF (NF1) to characterize baseline performance. Following this, networks were trained to produce reaches in a clockwise FF (FF1, 3200 batches of training). After initial FF learning networks were re-trained in a NF (NF2, 10000 batches). Following this "washout" phase networks were re-trained in the FF (FF2, 3200 batches) (see Figure 2a).

We characterized behavioural performance of the networks by measuring the maximum deviation of each hand trajectory from a straight line connecting initial and final target positions. During centre-out reaching in the initial NF baseline tests (NF1, Figure 2b) the network produced straight hand paths with very little lateral deviation.

We then trained the networks to compensate for the effects of a clockwise force field (FF1). The first time networks encountered the FF (batch 0), they exhibited large lateral deviations from a straight line trajectory (Figure 2a,b). Over the course of training the network's hidden weights were modified so that the networks produced different muscle stimulation commands that compensated for the forces produced by the FF (Figure 2b). Only the hidden (recurrent) weights were modified after the "growing up" phase, and not the input/output weights. By the end of FF1, relatively straight hand paths were recovered, and lateral deviation was near that in the NF baseline tests (NF1).

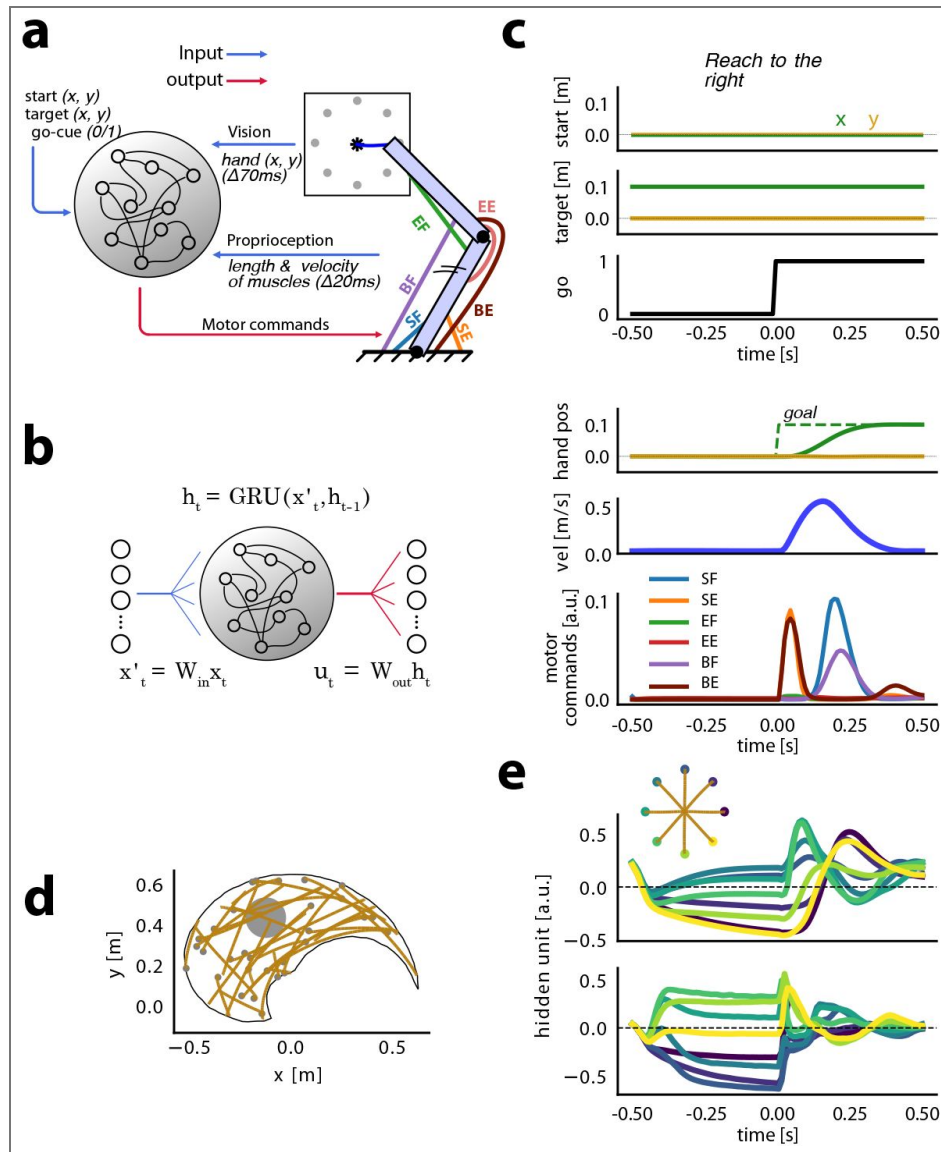


Figure 1. Recurrent neural network model inputs and outputs.

(a) RNNs receive a 17-dimensional input signal consisting of the location of the movement target in Cartesian coordinates, a “visual” feedback signal giving the arm’s endpoint position delivered with a 70 ms delay, a “proprioceptive” feedback signal consisting of the length and velocity of each of the 6 limb muscles delivered with a 20 ms delay, and a binary go cue. RNNs output 6 motor stimulation commands (between 0 and 1) to drive each muscle: SF (shoulder flexors), SE (shoulder extensors), EE (elbow extensors), EF (elbow flexors), BE (bi-articular extensors), and BF (bi-articular flexors). (b) The 17-dimensional input signal was mapped to the recurrent network using linear weights W_{in} . RNN output was transformed into motor commands by linear weights W_{out} . The vector h_t is the activity of hidden units at time t . (c) Task-related RNN inputs for a reach in a null field toward the rightmost target depicted in (a). For the purpose of illustration in this Figure, we translated the starting and target positions such that the start position is at the coordinates (0, 0). The simulation duration was 1 s, with 10 ms time steps. The goal (dashed lines) was set to the hand’s starting position before the go signal changed to 1, to the movement target position after that. (d) Sample endpoint trajectories after training RNNs on reaches to random target locations. Reaching trajectories are indicated in orange, and small gray dots show target positions. The large gray circle indicates the position of the centre-out reaches within the workspace. (e) Reaching trajectories and hidden unit activity from two example hidden units over time at the end of training in the centre-out task. colours indicate each of 8 targets. The go-cue switches from 0 to 1 at time $t = 0$.

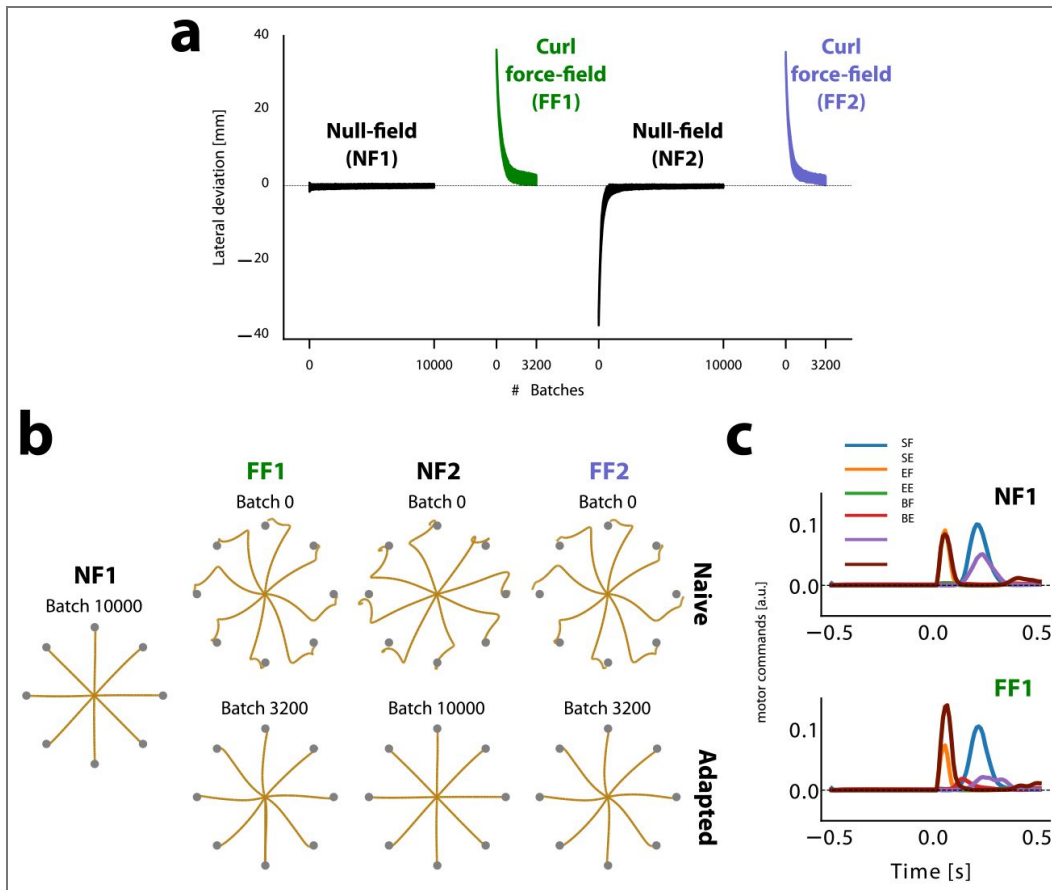


Figure 2. Networks learn to compensate for curl force fields without any contextual input.

(a) Lateral deviation averaged over 8 centre-out reaches for each batch. Black indicates null-field phases (NF1 and NF2), green indicates the first phase of the force field (FF1), and purple indicates the second phase of the force field (FF2). Positive values indicate deviation in the direction of the force field, which is clockwise relative to the line connecting the starting and target positions. (b) Simulated reaching trajectories at the beginning and end of each phase, grouped in different columns. (c) Motor commands during reaching toward the rightmost target for NF1 and FF1.

Importantly, at no time during training did the networks receive any contextual signal related to the presence or absence of the FF. Adaptation occurred during FF training because as the simulated limb is perturbed by the FF, hand position deviates away from the target, and the loss function increases. Over training the values of the RNN hidden unit weights are changed to minimize the loss function, and in turn, recover straight hand paths.

Washout

After force field adaptation (FF1) we trained the networks in a washout phase (NF2) in which we removed the simulated FF that the networks trained on in FF1. As is the case in empirical studies of FF learning, the networks initially showed an after-effect in movement kinematics in the opposite direction of the force field (Figure 2a,b), indicating that the networks prepared muscle commands to compensate for the (now absent) FF. After training in NF2, networks recovered straight line hand paths that were very similar to the performance in the NF1 baseline phase (only 0.1 mm difference in lateral deviation).

Re-adaptation

Following washout we re-trained networks in the same curl field (FF2) that they had been exposed to previously in the FF1 block. We observed that when networks initially encountered the FF in FF2, they exhibited smaller lateral deviations than those observed in the beginning of FF1, when they were first exposed to the simulated FF ($t(39) = 10.940$, $p = 1.9e-13$) (Figure 3a). In addition, networks adapted to the force field in FF2 faster than they did in FF1, as measured by an increased learning rate based on an exponential fit to the learning curve (see Methods; $t(39) = 9.284$, $p = 2.0e-11$). This pattern of improved performance in FF2 is seen in both human and monkey studies of motor learning and is typically referred to as savings.

The difference in performance in the first exposure to FF1 versus the first exposure to FF2 indicates that the network weights changed in such a way that facilitated improved initial performance and faster learning rate in FF2 as compared to FF1, while also not interfering with the performance in NF2. In other words, after training FF1, and washout in NF2, the networks retained enough information about the force field to improve re-learning in FF2, and this retained information was stored in such a way that it did not interfere with the ability of the networks to perform in NF2 just as they had done in NF1, before any FF learning. The pattern of savings observed here in our RNNs occurred despite the absence of any explicit contextual signal indicating the presence or absence of the simulated FF. When contextual signals are present, neural network models often create separate representations for two different tasks (Driscoll et al., 2022). We hypothesize that in our study the high dimensionality of the RNNs allows them to develop representations that effectively serve the ongoing task while also preserving some information about previously learned tasks.

We tested this hypothesis by repeating all simulations using RNNs with different numbers of hidden units. We found that as the number of hidden units increases, the likelihood of networks expressing savings also increases. For models with 256 hidden units, there was an 80% chance of expressing savings based on both the learning rate (if the learning rate is faster in FF2 than in FF1) and lateral deviation (if the initial lateral deviation in FF2 is smaller than when FF1 is first encountered). This probability could drop to nearly chance level if the number of hidden units is smaller than 32. This supports the idea that savings may depend upon on the dimensionality of the network's weight space.

As an additional control we trained networks after the growing up phase on an opposing force field (CCW) and then as above, exposed the networks to a NF washout phase, and then to a CW force field. In this case no savings was observed in the CW force field, either for initial lateral deviation, or for learning rate (Figure 4). In fact, we observed that initial lateral deviation is larger for the novel force field ($t(39) = -4.918$, $p = 1.6e-5$). This observation is in line with the finding that learning opposing force fields sequentially results in interference (Sun et al., 2022).

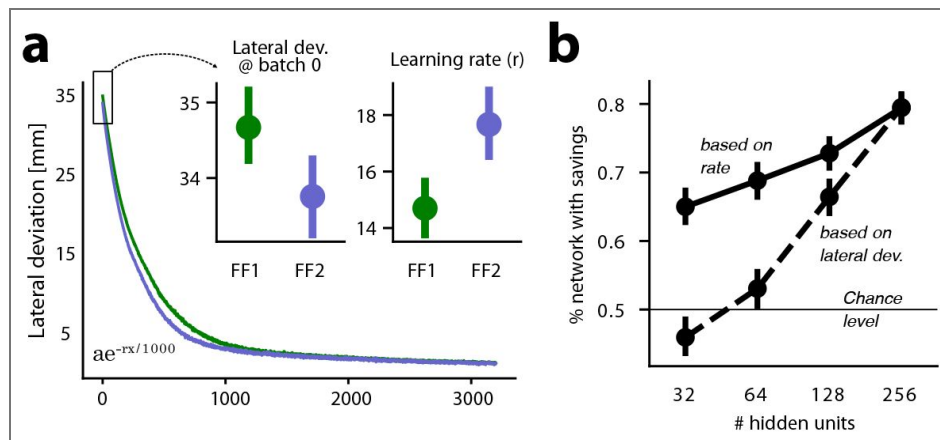


Figure 3. RNN models exhibit behavioural characteristics of savings.

(a) Mean learning curves averaged over RNN models. Sub-panels show (left) lateral deviation at batch 0 (before training in the corresponding phase) and (right) learning rate r after fitting an exponential curve to lateral deviations over training for each network. (b) The percentage of networks ($n=40$ total) with savings is plotted against the number of RNN hidden units. The dashed line indicates the percentage of networks with savings, defined as FF2 learning rate greater than FF1 learning rate. The solid line shows percentage of networks showing savings (lateral deviation at batch 0 smaller in FF2 than FF1). Error bars indicate 95% confidence interval.

The results of these control simulations underscore that the savings effect observed in our main study was learning-specific—it was due to prior learning of the CCW force field, and not a general effect of learning any novel dynamics.

Learning related changes in hidden unit activity

To characterize changes in hidden unit activity after learning we followed a similar approach to that described by Sun, O’Shea, and colleagues (Sun et al., 2022), in which they investigated how neural population activity in motor cortex changed after monkeys learned to adapt to FFs in an upper limb reaching task. The focus is on hidden unit activity during the preparatory phase, prior to the go signal, as this is the primary determinant of the feed-forward motor commands to muscles (Churchland and Shenoy, 2024).

We examined changes in a subspace defined by the relationship between hidden unit activity in the preparatory period, prior to the go cue, and movement-related force at the initial acceleration phase of the movement (hereafter referred to as the TDR subspace, see Methods for further details). We identified the TDR subspace by linearly regressing the preparatory hidden unit activity (340 ms before the go-cue) onto the endpoint (hand) force early during execution (90 ms after the go-cue; see Methods). Projecting the preparatory hidden unit activity associated with all 8 centre-out reaches onto this subspace revealed a ring formation (Figure 5a). This circular pattern has also been observed in empirical studies of adaptation for motor cortical neurons (Sun et al., 2022). This ring rotated in a counter-clockwise (CCW) following adaptation to a clockwise (CW) FF. This is consistent with the idea that after training in the CWFF the preparatory activity of the RNN hidden units is tuned to facilitate the production of forces in the direction opposite to the upcoming FF during movement.

After washout (NF2), this rotation reverted, leaving little or no residual part of the initial CCW rotation in the preparatory hidden unit activity (Figure 5a). The preparatory hidden unit activity again rotated in a CCW direction after adaptation to the CWFF in FF2.

The neural trajectories for preparation and for movement can be visualized in principal component space. Figure 6 shows trajectories during planning and early execution for initial FF1 and FF2 exposure. Hidden unit activity was subjected to a principal components analysis, and neural trajectories within the first three PCs are shown for movements to each of the eight movement targets. Filled circles indicate neural state 200 ms prior to the go cue. During the preparatory period trajectories travel along PC1 and then disperse across PC2 and PC3 into the circular pattern indicated by the filled stars, which indicate time of the go cue (also see Figure 5A). After the go cue neural trajectories shift back along PC1 and rotate along oscillatory patterns characteristic of populations of motor cortical neurons in non-human primates during movement (Churchland and Shenoy, 2024).

We probed the learning-related changes in RNN hidden unit activity that occurred outside of the TDR subspace. To do this, we calculated the extent to which the centroid of the preparatory hidden unit activity for 8 centre-out targets shifted following FF learning. Following the procedure used in Sun, O’Shea et al., to isolate learning-related changes in hidden unit activity common to all reach directions we calculated the difference between the mean preparatory activity (340 ms before go cue) of NF1 and FF1 and then orthogonalized this with respect to the TDR subspace (Sun et al., 2022). The result is referred to as a “uniform shift” (Figure 5b). After projecting the mean preparatory activity of all experimental phases onto this uniform shift and normalizing the projection values such that the NF1 projection is 0 and the FF1 projection is 1 (see Methods), we observed that at the end of the washout phase (NF2), the RNN hidden unit activity still showed a projection onto this direction that was significantly different than zero ($t(39)=7.484$, $p=4.6e-9$; Figure 5c). This indicates that the mean hidden unit activity during the preparatory period did not fully revert to pre-adaptation levels, despite full behavioural washout by the end of the washout phase (Figure 2a). This result is consistent with findings in monkey motor cortex (Sun et al., 2022), and the idea that this residual hidden unit activity captures information about the previously learned FF, and can be linked to subsequent savings.

Figure 4. No savings observed for novel force field in FF2.

(left) the lateral deviation at batch 0 (before training in the corresponding phase) and (right) the learning rate r after fitting an exponential curve to the lateral deviations over training for each network.

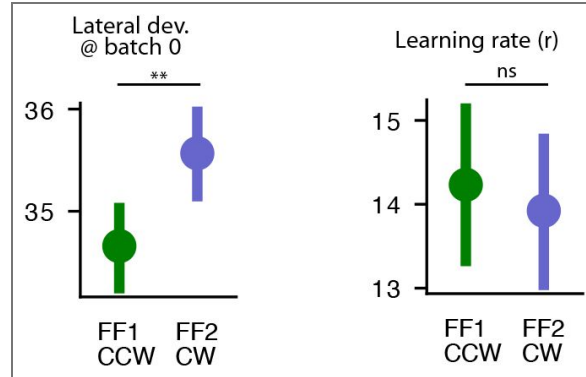
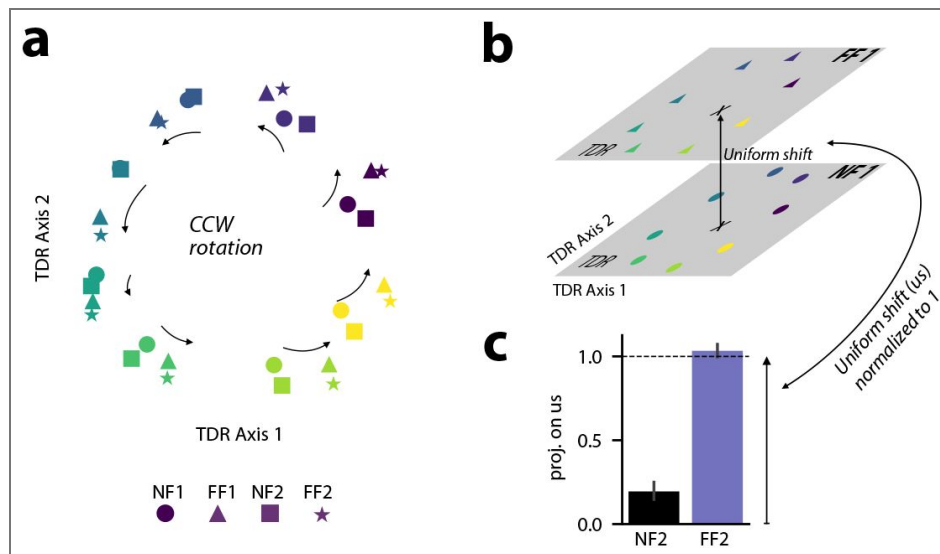


Figure 5. Changes in the preparatory activity of RNN hidden units following FF learning.

a: Projection of the hidden preparatory activity (340 ms before the go-cue) of an example trained model performing centre-out reaches on the force-predictive subspace acquired with targeted dimensionality reduction (TDR). Different reaching targets are indicated with different colours, and different adaptation phases are indicated with different shapes: circle for NF1, triangle for FF1, square for NF2, and star for FF2. **b:** A schematic illustration of the uniform shift. Each cross indicates the centre of the hidden preparatory activity for NF1 and FF1, and the arrow indicates the uniform shift. **c:** Projection of the hidden preparatory activity of all phases onto the uniform shift after orthogonalizing the uniform shift with respect to TDR. The data are scaled so that the projection of NF1 onto the uniform shift is zero and the projection of FF1 is one.



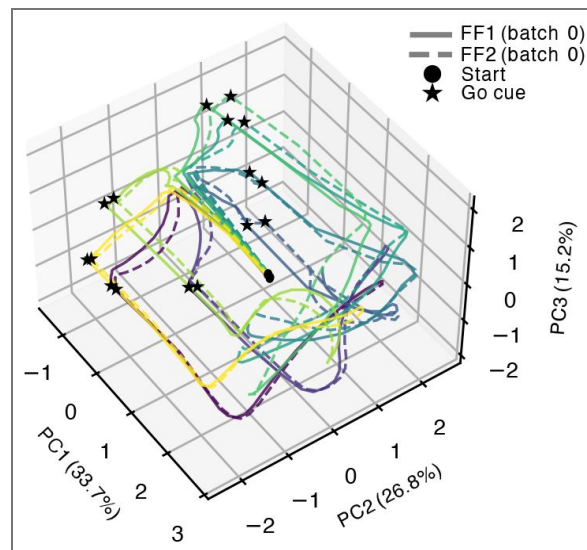


Figure 6. Neural trajectories during initial FF1 and FF2.

Trajectories for eight movement targets starting at the go cue and ending 200 ms into movement execution. Colors show movement directions. PC1–3 represent the first three principal components of neural activity variance.

Perturbing preparatory hidden unit activity along the uniform shift

The fact that the preparatory activity of RNN hidden units did not fully revert to the NF1 levels in the uniform shift direction suggests that this residual activity might underlie savings (Sun et al., 2022 [↗](#); Losey et al., 2024 [↗](#)). We tested this idea directly by perturbing hidden unit activity along the direction of the uniform shift, and we probed the effect of these perturbations on behavioural measures of savings.

We added a proportion of the uniform shift to the preparatory hidden unit activity of networks performing centre-out reaches and we measured the resulting changes in the lateral deviation of simulated hand trajectories. Importantly, the perturbations to preparatory hidden unit activity resulted in little or no changes in muscle activity prior to movement, and during movement itself no perturbations were delivered. We used models at batch 0 of the FF2 phase, when they had not yet been trained on FF adaptation a second time. We have already shown that in FF2, hand trajectory lateral deviation is smaller than in FF1, and we used this as a metric to characterize savings (Figure 3a [↗](#)). We examined how much the lateral deviation at batch 0 would change following perturbations to RNN hidden unit activity in the direction of the uniform shift. Details about how the perturbations to hidden unit activity were implemented are found in Methods.

When we perturbed the preparatory hidden unit activity in the opposite direction of the uniform shift (negative magnitudes in Figure 7a [↗](#)), lateral deviation of movement-related hand trajectories increased. By increasing the magnitude of perturbation further, we could effectively abolish savings altogether. In contrast, if we perturbed the RNN hidden unit preparatory activity in the same direction as the uniform shift, lateral deviation of hand trajectories was reduced, thus enhancing savings. Example trajectories toward the right-most target are shown in Figure 7a [↗](#) for each uniform shift perturbation. These results support the idea that the hidden unit activity in the direction of the uniform shift that remains after washout represents a neural trace of the initial FF learning, one that supports subsequent savings when the models adapt to the FF a second time.

The perturbations changed the activity of hidden units (Figure 7c [↗](#)). These changes were large early after the delivery of the perturbation, and then they reached a steady state. However, it is important to note that these changes in the preparatory hidden unit activity did not result in substantive changes in the motor commands (Figure 7b [↗](#)), which emphasizes that the uniform shift resides in the null space of motor output.

In summary, the activity of networks along the direction of the uniform shift did not revert fully to pre-adaptation levels after the washout phase, despite the behaviour (hand trajectories and muscle activity) being fully washed out, the same as pre-adaptation baseline performance. Perturbing RNN hidden unit activity during the preparatory period along the direction of the uniform shift enhanced savings while perturbations which reduced this residual hidden unit activity reduced or eliminated savings. RNN models with higher dimensionality (more hidden units) were more likely to exhibit evidence of savings while smaller models were less likely to show savings.

Discussion

We have shown here that RNNs trained to control a computational model of the upper limb show behavioural signatures of savings when learning multiple FFs in succession. We also found that savings is enhanced when the dimensionality of the control network is higher. Following the approach described in a recent electrophysiological study of monkey motor cortex (Sun et al., 2022 [↗](#)), we identified a subspace of RNN hidden unit activity during the preparatory period after target presentation but prior to movement that shifts after initial FF learning, and subsequently retreats after behavioural washout in a NF. Importantly, this “uniform shift” did not retreat all the way back to pre-FF baseline levels after washout (see Figure 5c [↗](#)). Despite this residual trace of FF

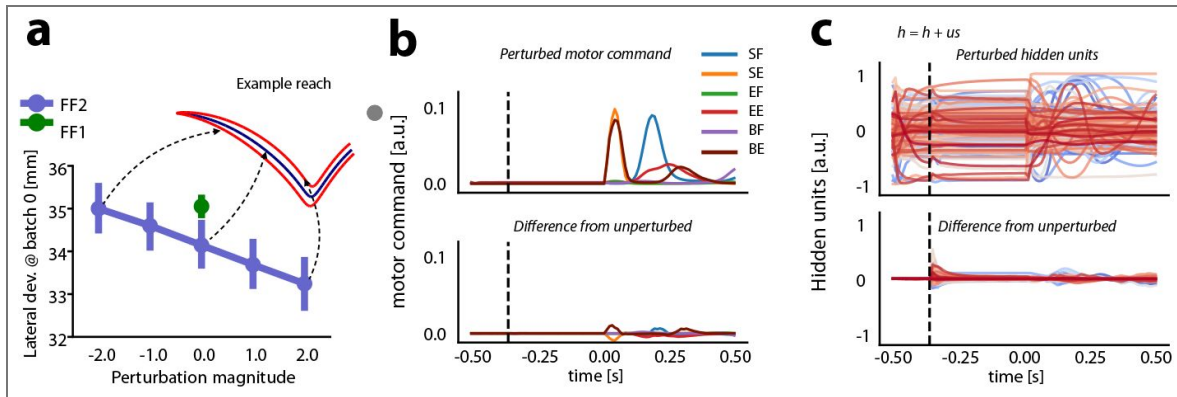


Figure 7. Uniform shift in RNN hidden unit activity is related to savings.

A: Lateral deviation in FF2 (purple) when hidden preparatory activity was perturbed in the positive (+) and negative (-) directions of the uniform shift with different magnitudes. Arrows indicate trajectories of an example model when hidden unit activity was perturbed (red) or not (blue). Lateral deviation of hand trajectories for FF1 are shown in green. **B:** Motor commands when hidden unit activity was perturbed with a magnitude of 2.0. Vertical dashed line indicates when the perturbation was delivered. The lower sub-panel shows the difference from the unperturbed motor command. **C:** Activity of 128 RNN hidden units after perturbation (dashed line indicates time of perturbation along the uniform shift). Lower sub-panel shows the difference between the unperturbed and perturbed RNN hidden unit activity.

learning after washout, in a NF the RNN produced reaches to targets that were the same as those produced prior to initial learning. Importantly, alternating network training on opposing fields (CW and CCW) did not produce savings.

We interpret this residual hidden unit activity as a neural trace of the initial learning that remains after washout, and this is consistent with the idea that this persistent signal contributes to savings (Sun et al., 2022 [↗](#)). We tested this hypothesis directly by perturbing RNN hidden unit activity during the preparatory period along the direction of the uniform shift—an approach that is possible in a computational model but is not presently feasible in a biological neural network. After NF washout when we re-exposed networks to a previously learned FF, perturbations that amplified the residual trace of learning (increasing activity along the uniform shift) resulted in increased savings. When we perturbed hidden unit activity in the other direction to reduce activity along the uniform shift, savings was reduced and in some cases abolished altogether. These results provide evidence of a causal link between the identified neural trace of learning that remains after washout and subsequent savings when RNNs are re-exposed to the previously learned FF.

Our results are compatible with the proposal by Sun, O’Shea, et al. that the activity of neurons in primary motor cortex during the preparatory period prior to movement contains a component that tracks previously learned motor behaviour (Sun et al., 2022 [↗](#)). They propose that these residual traces of prior learned behaviour are encoded in a way that separates the associated motor memories in neural state space and facilitates recall of the appropriate control policies. Losey et al. describe a similar account in the context of a brain-machine interface in which new motor learning is encoded in a neighbouring region of neural state space such that it solves the motor control task but doesn’t interfere with prior learning (Losey et al., 2024 [↗](#)). In our RNNs this was achieved through learning-related changes in the recurrent weights, such that after NF washout the component of the uniform shift that remained didn’t interfere with NF motor behaviour, but did produce savings when networks were re-exposed to the previously learned FF.

Our finding that higher-dimensional RNNs are more likely to produce savings supports the idea that encoding newly learned control policies so that they do not interfere with previously learned motor memories depends upon the availability of adequate dimensionality in neural state space. This implies that multiple motor memories can be encoded in neural subspaces so that they do not interfere with each other. A number of recent empirical studies support this idea. Kim et al. recorded from anterior lateral motor cortex of mice over several months and tracked how neurons represented different learned motor tasks (Kim et al., 2025 [↗](#)). They found that learning produced new neural representations that did not modify existing representations, and re-exposure to a previously learned motor task re-activated the previous neural activity patterns. In a recent paper Bernardi et al. recorded neural activity in prefrontal cortex and hippocampus of monkeys during cognitive tasks, and found that multiple abstract task-related variables were encoded in neural state space using a geometry that allowed separability using linear classifiers (Bernardi et al., 2020 [↗](#)). Neural recordings in monkey motor cortex show that this kind of task representation emerges prior to movement, in preparatory activity after a movement target is shown but before a go signal instructs the animal to begin a movement (Churchland et al., 2012 [↗](#); Sun et al., 2022 [↗](#)). A similar time course emerged in our RNN simulations here (Figure 1e [↗](#), Figure 5a [↗](#)).

In primates presumably high-level contextual cues can aid in indexing the appropriate previously learned control policy by activating populations of neurons in a neural direction appropriate for the previously learned task. Indeed a number of theoretical accounts exist that position contextual cues as a driver of motor memory encoding and selection (Wolpert and Kawato, 1998 [↗](#); Haruno et al., 2001 [↗](#); Heald et al., 2021 [↗](#)). Similar indexing is likely occurring in accounts where the learning of FFs that would normally interfere is avoided by associating each with the planning (Sheahan et al., 2016 [↗](#)) or imagination (Sheahan et al., 2018 [↗](#)) of different follow-through movements. In our RNNs no such contextual signal was provided, and so the question arises, how are residual traces of previously learned FFs activated? One possibility is that the error signals encountered when our RNNs are re-exposed to a previously learned FF activate neurons within

the previously learned subspace. This kind of scheme in which re-exposure to previously encountered errors produces savings is consistent with accounts in which a history of errors or corrections to errors plays a role in motor memory formation (Herzfeld et al., 2014; Leow et al., 2016).

In a recent computational modelling study Dubreuil et al. proposed that non-random neural population connectivity structures involving multiple subpopulations that play functionally distinct roles encode multiple tasks better than random connectivity structures (Dubreuil et al., 2022). This seems compatible with the idea presented here and in other recent work that low-dimensional recurrent subspaces embedded within a high-dimensional neural control space are used to encode features of movements such as target directions (Churchland et al., 2012), anticipated sensory consequences of perturbations (Michaels et al., 2024), and motor skills (Sun et al., 2022; Losey et al., 2024). The idea that motor memories are encoded in a distributed sensorimotor network and that features of motor adaptation emerge as a result of the dynamical properties of recurrent neural circuits have also been discussed in other computational modelling studies using recurrent neural networks. Ajemian et al. proposed a theoretical framework in which error signals prompt a reorganization of synaptic connectivity to encode motor memories within a high dimensional neural space (Ajemian et al., 2010, 2013).

The present results are based on RNNs trained in an error-based approach using backpropagation through time (Werbos, 1990) using the Adam optimizer (Kingma and Ba, 2014). Other techniques for training RNNs have been proposed including the FORCE algorithm (Sussillo and Abbott, 2009). In addition, several recent reports have demonstrated success using reinforcement learning approaches to train neural networks in the context of sensorimotor control tasks (Lillicrap et al., 2015; Codol et al., 2024a). An interesting avenue for future work is to determine how the present results may or may not generalize to different neural network architectures and learning rules.

The work presented here adds to the growing literature documenting how complex features of motor behaviour can arise due to the dynamics of preparatory neural activity in motor cortex. Fuelner et al. show that a number of features of motor adaptation emerge as a result of the dynamical properties of recurrent neural circuits in which sensory feedback modulates motor output (Fuelner et al., 2025). In a recent paper Smoulder et al. show that neural activity in monkey motor cortex scales with reward magnitude, and that this reward signal interacts with movement preparation signals in such a way that high rewards disrupt movement preparation and result in poor motor performance compared to moderate rewards (Smoulder et al., 2024).

The phenomenon of savings in motor learning implies that motor memory associated with an initial bout of training leads to faster subsequent relearning. The nature of the memory that is stored and how it influences subsequent relearning has been a topic of some debate in the recent literature. One account of savings emphasizes the effect of explicit, strategic components of motor learning (Huberdeau et al., 2015; Morehead et al., 2015; Avraham et al., 2021). Another line of work focuses on the idea that faster relearning may also be driven by implicit learning processes that result from previously experienced errors (Leow et al., 2016; Coltman et al., 2019; Yin and Wei, 2020; Coltman et al., 2021).

The simulations described here do not constitute evidence that savings in motor learning tasks is exclusively implicit in animals and humans. The purely context-free learning implemented in our simulations is not meant to be a full model of biological learning, as in biological systems some form of contextual information is invariably available. Indeed, computational models of motor learning that incorporate contextual effects already exist, e.g. (Heald et al., 2021). Nevertheless, our simulations provide a useful window into what the context-free component of savings may look like. This approach offers a powerful means of probing the context-free component of savings in isolation—something that is not readily achievable in animal or human experiments.

Recent empirical work suggests that relearning after washout of implicit adaptation can be attenuated rather than facilitated, a phenomenon attributed to anterograde interference from the washout phase (Leow et al., 2020; Yin and Wei, 2020; Hamel et al., 2021, 2022; Avraham et al., 2021).

al., 2021 [↗](#); Hadjiosif et al., 2023 [↗](#); Wang and Ivry, 2025 [↗](#)). The savings observed in our simulations differs from these behavioral findings. Crucially, our model excludes both contextual interference (since no cues signal which force field is present) and explicit-implicit interactions (since context-driven explicit learning is absent). Our goal was not to model a complete explicit-implicit system, but rather to probe how savings may emerge from a purely implicit mechanism and to compare the underlying neural geometry to monkey electrophysiology data. Our results suggest that high-dimensional neural circuits possess an intrinsic capacity for savings via persistent preparatory traces. How and when this capacity may be masked by interference or explicit-implicit interactions in biological systems remains an open question for future work.

The results of our work here with RNNs and the related electrophysiological studies of populations of motor cortical neurons of non-human primates point to a neural basis of savings (Sun et al., 2022 [↗](#); Losey et al., 2024 [↗](#)). We showed here that by increasing the number of hidden units in our RNNs, and hence increasing the dimensionality of the control space, RNNs were more likely to produce behavioural signatures of savings (Figure 3b [↗](#)). The high dimensionality of neural space enables a new motor memory to be encoded in such a way that it doesn't interfere with other previously learned information, while still facilitating savings when the network is re-exposed to the newly learned skill. This neural basis of savings can be characterized as implicit, since in our RNN simulations we did not provide the network with any contextual input that would signal the presence or absence of any given FF.

Methods

RNN model

Recurrent neural networks (RNNs) were trained to control movements of a simulated two degree of freedom arm that included rotations about a shoulder joint and an elbow joint in a horizontal plane. The model includes six rigid-tendon Hill-type muscle actuators comprising mono-articular flexors and extensors spanning the shoulder and elbow, as well as a pair of bi-articular muscles producing flexion or extension forces about both shoulder and elbow joints (Kistemaker et al., 2010 [↗](#)). RNN models are implemented in PyTorch (Paszke et al., 2019 [↗](#)) and receive a 17-dimensional input signal to a linear input layer, which is fully connected to a recurrent layer consisting of gated recurrent units (GRUs) (Cho et al., 2014 [↗](#)). The GRU layer is connected to a 6-dimensional linear output layer which provides stimulation commands over time to each of the 6 muscles in the arm model (see Figure 1 [↗](#)). Simulations were carried out in Python 3.10 using our open-source MotorNet toolbox (Codol et al., 2024b [↗](#)).

Input-hidden (W_{in}) and hidden-hidden recurrent (W_r) weights (see Figure 1 [↗](#)) were initialized using Glorot initialization (Glorot and Bengio, 2010 [↗](#)) and orthogonal initialization (Hu et al., 2020 [↗](#)), respectively, with biases set to 0. The output layer used a sigmoid activation function. The hidden-output weights (W_{out}) were initialized with the Glorot initialization scheme, and its biases were set to -5.0. The sigmoid activation function ensured the controller's output remained close to 0 at the start of training, promoting a stable initialization state. We set the initial hidden unit activity of the network (h_0) as a learnable parameter and initialized it to 0.

The RNN models received a 17-dimensional input vector consisting of task-related inputs along with time-delayed feedback representing visual and proprioceptive signals. The task-related input consisted of a 2-element vector of (x, y) Cartesian coordinates for the target position, and a binary go-cue signal that switched from 0 to 1 when the movement should be initiated. The visual feedback was a Cartesian coordinate of the arm's endpoint (x, y) , and the proprioceptive feedback was the lengths and velocities of all 6 muscles. The time step for simulations was set to 10 ms, the visual delay (Δ_v) was 70 ms, and the proprioceptive delay (Δ_p) was 20 ms. We also treated the go cue as a visual signal, meaning that at each time step the network received the 70 ms time delayed value. At each time step the RNN model transformed the above described 17-dimensional input into a 6-dimensional muscle stimulation command.

Growing up training phase

During an initial “growing up” phase new initialized RNN models were trained to move the arm from random starting positions to random target positions, both drawn from a uniform distribution across the joint space of the arm model. Note that due to muscle lengths and joint geometry, only a subset of the workspace was reachable for the model (Figure 1d). In 50% of simulations, no go-cue was provided (a catch trial). This was done to ensure that the network avoided producing anticipatory muscle stimulation commands. In the other 50% of cases the time of the go-cue switch from 0 to 1 was drawn from a random uniform distribution between 100 ms and 300 ms after the start of each simulated trial.

The loss function for training was mainly comprised of position loss, the Euclidean distance between the arm endpoint position \mathbf{x}_t and the desired position \mathbf{x}_t^* . The desired position was set to be equal to the starting position of the limb’s endpoint before the go cue, and after that the target position. We also included terms in the loss function that punished large and oscillatory hidden and muscle activity, and the jerk (the second derivative of acceleration) of the endpoint trajectory (Flash and Hogan, 1985). The full form of the loss function is shown in Equation 1:

$$\begin{aligned}
 L &= \frac{\sum_{t=1}^N L_t}{N} \\
 L_t &= 10^3 L_t^p + 10^5 L_t^j + 10^{-1} L_t^m + 10^{-5} L_t^h \\
 L_t^p &= \|\mathbf{x}_t^* - \mathbf{x}_t\|_1 \\
 L_t^j &= \ddot{\mathbf{x}}_t^T \ddot{\mathbf{x}}_t \\
 L_t^m &= \mathbf{f}_t^T \mathbf{f}_t + 3 \times 10^{-3} \dot{\mathbf{f}}_t^T \dot{\mathbf{f}}_t \\
 L_t^h &= \mathbf{h}_t^T \mathbf{h}_t + 10^2 \dot{\mathbf{h}}_t^T \dot{\mathbf{h}}_t
 \end{aligned} \tag{1}$$

where the subscript t indicates time step, N is the total number of time steps in the simulation, T is the transpose operation, and $\|\cdot\|_1$ is the L_1 norm. L_t^p , L_t^j , L_t^m , and L_t^h indicate position, jerk, muscle, and hidden loss, respectively. \mathbf{h}_t is a n -element vector of hidden unit activity, and $\dot{\mathbf{h}}_t$ is its derivative. \mathbf{f}_t is a 6-element vector of muscle forces, and $\dot{\mathbf{f}}_t$ is its derivative. Note that muscle forces are different from muscle stimulation commands (Figure 1A,C), which are inputs to the Ordinary Differential Equation that produces muscle forces (Kistemaker et al., 2010).

The RNN models were initially trained on 20,000 batches with a batch size of 32 on simulations of 1 s (100 time steps). RNN weights were adjusted using the Adam optimization scheme (Kingma and Ba, 2014) with a learning rate $lr = 0.003$.

Motor learning phases

Once the networks were trained to perform reaches to random targets in a null-field (NF), we fixed the input-to-hidden weights (W_{in}) and the hidden-to-output weights (W_{out}) and their biases. This allowed us to isolate subsequent learning-related changes resulting from our experimental manipulations to the recurrent connectivity of hidden units (Feulner et al., 2025).

We trained networks on centre-out reaches from a start position to 8 equidistant targets (0.1 m) around the circumference of a circle (see Figure 1). The start position corresponded to external joint angles of 60 degrees at the shoulder and 90 degrees at the elbow. We exposed models sequentially to FFs or NFs and in each phase we continued to adjust hidden recurrent weights to optimize the loss function described in Equation 1. During these experimental phases we used a stochastic gradient descent optimization scheme with learning rate parameter $lr = 0.005$ (Sutskever et al., 2013). This ensures batch-local learning, and thus provides greater control and transparency over the course of learning. It also results in gradual learning over batches, better resembling learning curves from empirical studies of force field learning in humans and non-human primates.

In the NF1 experimental phase the models were trained on centre-out reaches only. We trained the models for 10,000 batches of size 32 (4 repetitions of each of the 8 targets). As in the growing-up phase, 50% of trials were catch-trials in which the go-cue did not change from 0 to 1. After NF1 training, we continued to train the models to perform centre-out reaches but we introduced a clockwise curl force field (CWFF) for all movements (FF1). The external force (F_x, F_y) applied at the arm's endpoint that produced a CWFF is described by the following equation:

$$\begin{bmatrix} F_x \\ F_y \end{bmatrix} = b \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix} \quad (2)$$

where \dot{x} and \dot{y} are the velocity of the arm's endpoint in Cartesian coordinates and $b = 8 \text{ Ns/m}$ is a scalar constant defining the strength of the FF. In the null field (NF), $b = 0$. We trained the models for 3,200 batches of size 32 in the FF1 experimental phase, with 50% catch-trials.

Following FF1, the models were trained again in a null field (NF2), using the same procedures as in NF1 (10,000 batches of size 32). After NF2, the models were again trained in the presence of a CWFF (FF2), exactly as in FF1.

Lateral deviation

We evaluated the behavioural performance of the models during NF1, FF1, NF2, and FF2 by calculating the maximum lateral deviation of the endpoint trajectory from straight lines connecting the starting and target positions. We will refer to this measure as “lateral deviation”, and it was considered positive if it was in the clockwise direction from the straight line, and negative otherwise. For each batch of training we calculated the mean lateral deviations across all 8 targets.

For each model, we characterized the learning rate during FF1 and FF2 by fitting an exponential function of the following form to the mean lateral deviation across training batches:

$$\hat{y} = \alpha e^{-rx_n/1000} \quad (3)$$

where \hat{y} is the modelled lateral deviation at batch number x_n , α is a scaling factor that determines the initial lateral deviation, and r is the rate at which lateral deviation decays over time (indicating the learning rate). Before fitting we smoothed learning rates over batches by window-averaging with a kernel size of 5 batches.

Targeted dimensionality reduction

Following the procedure described in (Sun et al., 2022 [↗](#)) we identified a subspace of RNN hidden unit activity that predicts the arm's endpoint initial forces based on the preparatory activity of hidden units (hidden unit activity before the go-cue). To do this we applied targeted dimensionality reduction (TDR) using model data at the end of the NF1 experimental phase. The subspace is defined as:

$$\mathbf{H}_{-340 \text{ ms}}^{\text{NF1}} = [\mathbf{F}_{+90 \text{ ms}}^{\text{NF1}} \quad \mathbb{1}] \mathbf{W} \quad (4)$$

where $\mathbf{H}_{-340 \text{ ms}}^{\text{NF1}}$ is the matrix of hidden unit activity of size (targets \times units) 340 ms before the go-cue, $\mathbf{F}_{+90 \text{ ms}}^{\text{NF1}}$ is the matrix of endpoint forces of size (targets \times 2) 90 ms after go-cue, $\mathbb{1}$ is the targets-element vector concatenated to the force matrix, and \mathbf{W} is the matrix of size (3 \times units). For $\mathbf{F}_{+90 \text{ ms}}^{\text{NF1}}$ the parameter 90 ms was chosen because it coincides with peak acceleration. For the $\mathbf{H}_{-340 \text{ ms}}^{\text{NF1}}$ parameter, 340 ms was chosen because it ensures that hidden unit activity is stabilized.

We then calculated the pseudo-inverse of \mathbf{W} , resulting in a (units \times 3) matrix \mathbf{W}^+ . To find the force-predictive subspace, we took the first two columns of \mathbf{W}^+ (ignoring the intercept) and orthogonalized them using the Gram-Schmidt orthogonalization scheme, resulting in $\hat{\mathbf{W}}^+$. We projected

the preparatory hidden unit activity of all experimental phases ($\mathbf{H}_{-340 \text{ ms}}^{\text{NF1}}$, $\mathbf{H}_{-340 \text{ ms}}^{\text{FF1}}$, $\mathbf{H}_{-340 \text{ ms}}^{\text{NF2}}$, $\mathbf{H}_{-340 \text{ ms}}^{\text{FF2}}$) onto $\hat{\mathbf{W}}^+$ after removing their global mean.

Uniform shift

Following the experimental phases FF1 and FF2 we calculated the direction in which the RNN hidden unit activity during the preparatory period shifted, averaged across movement directions. We averaged the preparatory hidden unit activity over the 8 targets in FF1 and NF1, and then calculated the difference. Following the convention in (Sun et al., 2022) this shift is referred to as a “uniform shift” (us):

$$\mathbf{us} = \bar{\mathbf{H}}_{-340 \text{ ms}}^{\text{FF1}} - \bar{\mathbf{H}}_{-340 \text{ ms}}^{\text{NF1}} \quad (5)$$

where the bar indicates averaging over 8 movement directions. We orthogonalized the uniform shift with respect to $\hat{\mathbf{W}}^+$. This allowed us to test for changes outside the force-predictive subspace. We then projected the preparatory hidden unit activity (340 ms before the go-cue) of all experimental phases after removing the global mean. Next, we normalized the projection values for each model so that the projection of $\mathbf{H}_{-340 \text{ ms}}^{\text{NF1}}$ onto the uniform shift is zero, and the projection of $\mathbf{H}_{-340 \text{ ms}}^{\text{FF1}}$ onto the uniform shift is one.

Perturbing hidden unit activity

To conduct a causal test of the idea that the non-zero uniform shift activity that remained following NF2 is related to savings, we perturbed the activity of hidden units at the end of the preparatory period by adding to each hidden unit a proportion (-2, -1, 0, 1, 2) of the projection of that unit onto the uniform shift. We conducted these perturbations separately for each movement direction. The perturbation was applied 340 ms prior to the go cue and the duration of the perturbation was one simulation time step.

Data availability

Python code to reproduce the simulations and analyses described here is available on GitHub at the following repository: <https://github.com/mshahbazi1997/MotorSavingModel.git>.

Acknowledgements

This work was supported by the Natural Sciences and Engineering Research Council of Canada through a Discovery Grant RGPIN/05458-2018 to P.L.G., and a FRQNT Strategic Clusters Program grant to O.C. J.A.M. was supported by a Banting Postdoctoral Fellowship and a BrainsCAN Postdoctoral Fellowship, and by Canadian Institutes of Health Research grant PJT-175010 to Dr. Andrew Pruszynski. The authors wish to thank Mehrdad Kashefi for useful discussions about this project.

Additional information

Author contributions

M.S., O.C., J.A.M., and P.L.G. conception and design of research; M.S. performed simulations; M.S. and P.L.G. analyzed data; M.S., P.L.G., and J.A.M. interpreted results of experiments; M.S. prepared figures; M.S., P.L.G., and J.A.M. drafted manuscript; M.S., P.L.G., and J.A.M. edited and revised manuscript; M.S., O.C., P.L.G., and J.A.M. approved final version of manuscript.

Funding

Funder	Grant reference number	Author
Natural Sciences and Engineering Research Council of Canada (NSERC)	RGPIN/05458-2018	Paul L Gribble
Fonds de recherche du Québec (FRQ)		Olivier Codol
Banting Research Foundation (BRF)		Jonathan A Michaels
Canada First Research Excellence Fund (CFREF)		Jonathan A Michaels

Author ORCID iDs

Mahdiyar Shahbazi: <https://orcid.org/0000-0002-4883-4376>

Olivier Codol: <https://orcid.org/0000-0003-0796-5457>

Jonathan A Michaels: <https://orcid.org/0000-0002-5179-3181>

Paul L Gribble: <https://orcid.org/0000-0002-1368-032X>

References

Ajemian R, D'Ausilio A, Moorman H, Bizzi E. (2010) Why professional athletes need a prolonged period of warm-up and other peculiarities of human motor learning. *J Mot Behav* **42**:381-388

<https://doi.org/10.1080/00222895.2010.528262> | PubMed

Ajemian R, D'Ausilio A, Moorman H, Bizzi E. (2013) A theory for how sensorimotor skills are learned and retained in noisy and nonstationary neural circuits. *Proc Natl Acad Sci U S A* **110**:E5078-87

<https://doi.org/10.1073/pnas.1320116110> | PubMed

Albert ST, Shadmehr R. (2018) Estimating properties of the fast and slow adaptive processes during sensorimotor adaptation. *J Neurophysiol* **119**:1367-1393

<https://doi.org/10.1152/jn.00197.2017> | PubMed

Avraham G, Morehead JR, Kim HE, Ivry RB (2021) Reexposure to a sensorimotor perturbation produces opposite effects on explicit and implicit learning processes. *PLoS Biol* **19**:e3001147

<https://doi.org/10.1371/journal.pbio.3001147> | PubMed

Bernardi S, Benna MK, Rigotti M, Munuera J, Fusi S, Salzman CD (2020) The geometry of abstraction in the hippocampus and prefrontal cortex. *Cell* **183**:954-967.e21.

<https://doi.org/10.1016/j.cell.2020.09.031> | PubMed

Cho K, van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, Bengio Y. (2014) Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv*

<https://doi.org/10.48550/arxiv.1406.1078>

Churchland MM, Cunningham JP, Kaufman MT, Foster JD, Nuyujukian P, Ryu SI, Shenoy KV (2012) Neural population dynamics during reaching. *Nature* **487**:51-56

<https://doi.org/10.1038/nature11129> | PubMed

Churchland MM, Shenoy KV (2024) Preparatory activity and the expansive null-space. *Nat Rev Neurosci* **25**:213-236

<https://doi.org/10.1038/s41583-024-00796-z> | PubMed

Codol O, Krishna NH, Lajoie G, Perich MG (2024) Brain-like neural dynamics for behavioral control develop through reinforcement learning. *bioRxiv*

<https://doi.org/10.1101/2024.10.04.616712>

Codol O, Michaels JA, Kashefi M, Pruszynski JA, Gribble PL (2024) MotorNet: a Python toolbox for controlling differentiable biomechanical effectors with artificial neural networks. *eLife*

<https://doi.org/10.7554/eLife.88591> | PubMed

Coltman SK, van Beers RJ, Medendorp WP, Gribble PL (2021) Sensitivity to error during visuomotor adaptation is similarly modulated by abrupt, gradual, and random perturbation schedules. *J Neurophysiol* **126**:934-945

<https://doi.org/10.1152/jn.00269.2021> | PubMed

- Coltman SK, Cashaback JGA, Gribble PL (2019) Both fast and slow learning processes contribute to savings following sensorimotor adaptation. *J Neurophysiol* **121**:1575-1583 <https://doi.org/10.1152/jn.00794.2018> | PubMed
- Conditt MA, Gandolfo F, Mussa-Ivaldi FA (1997) The motor system does not learn the dynamics of the arm by rote memorization of past experience. *J Neurophysiol* **78**:554-560 <https://doi.org/10.1152/jn.1997.78.1.554> | PubMed
- Diedrichsen J, White O, Newman D, Lally N. (2010) Use-dependent and error-based learning of motor behaviors. *J Neurosci* **30**:5159-5166 <https://doi.org/10.1523/jneurosci.5406-09.2010> | PubMed
- Driscoll L, Shenoy K, Sussillo D. (2022) Flexible multitask computation in recurrent networks utilizes shared dynamical motifs. *bioRxiv* <https://doi.org/10.1101/2022.08.15.503870>
- Dubreuil A, Valente A, Beiran M, Mastrogiuseppe F, Ostojic S. (2022) The role of population structure in computations through neural dynamics. *Nat Neurosci* **25**:783-794 <https://doi.org/10.1038/s41593-022-01088-4> | PubMed
- Feulner B, Perich MG, Miller LE, Clopath C, Gallego JA (2025) A neural implementation model of feedback-based motor learning. *Nat Commun* **16**:1-14 <https://doi.org/10.1038/s41467-024-54738-5> | PubMed
- Flash T, Hogan N. (1985) The coordination of arm movements: an experimentally confirmed mathematical model. *J Neurosci* **5**:1688-1703 <https://doi.org/10.1523/jneurosci.05-07-01688.1985> | PubMed
- Glorot X, Bengio Y. (2010) Understanding the difficulty of training deep feedforward neural networks. In: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics. Sardinia, Italy. pp. 249-256
- Hadjosif AM, Morehead JR, Smith MA (2023) A double dissociation between savings and long-term memory in motor learning. *PLoS Biol* **21**:e3001799 <https://doi.org/10.1371/journal.pbio.3001799> | PubMed
- Haith AM, Huberdeau DM, Krakauer JW (2015) The influence of movement preparation time on the expression of visuomotor learning and savings. *J Neurosci* **35**:5109-5117 <https://doi.org/10.1523/jneurosci.3869-14.2015> | PubMed
- Hamel R, Dallaire-Jean L, Fontaine De La, Lepage JF, Bernier PM (2021) Learning the same motor task twice impairs its retention in a time and dose-dependent manner. *Proc Biol Sci* **288**:20202556 <https://doi.org/10.1098/rspb.2020.2556> | PubMed
- Hamel R, Lepage JF, Bernier PM (2022) Anterograde interference emerges along a gradient as a function of task similarity: A behavioural study. *Eur J Neurosci* **55**:49-66 <https://doi.org/10.1111/ejn.15561> | PubMed
- Haruno M, Wolpert DM, Kawato M. (2001) Mosaic model for sensorimotor learning and control. *Neural Comput* **13**:2201-2220 <https://doi.org/10.1162/089976601750541778> | PubMed
- Heald JB, Lengyel M, Wolpert DM (2021) Contextual inference underlies the learning of sensorimotor repertoires. *Nature* **600**:489-493 <https://doi.org/10.1038/s41586-021-04129-3> | PubMed
- Herzfeld DJ, Vaswani PA, Marko MK, Shadmehr R. (2014) A memory of errors in sensorimotor learning. *Science* **345**:1349-1353 <https://doi.org/10.1126/science.1253138> | PubMed
- Hu W, Xiao L, Pennington J. (2020) Provable benefit of orthogonal initialization in optimizing deep linear networks. *arXiv* <https://doi.org/10.48550/arXiv.2001.05992>
- Huang VS, Haith A, Mazzoni P, Krakauer JW (2011) Rethinking motor learning and savings in adaptation paradigms: model-free memory for successful actions combines with internal models. *Neuron* **70**:787-801 <https://doi.org/10.1016/j.neuron.2011.04.012> | PubMed
- Huberdeau DM, Haith AM, Krakauer JW (2015) Formation of a long-term memory for visuomotor adaptation following only a few trials of practice. *J Neurophysiol* **114**:969-977 <https://doi.org/10.1152/jn.00369.2015> | PubMed

- Kim JH, Daie K, Li N. (2025) A combinatorial neural code for long-term motor memory. *Nature* **637**:663-672 <https://doi.org/10.1038/s41586-024-08193-3> | [PubMed](#)
- Kingma DP, Ba J. (2014) Adam: A method for stochastic optimization. *arXiv* <https://doi.org/10.48550/arXiv.1412.6980>
- Kistemaker DA, Wong JD, Gribble PL (2010) The central nervous system does not minimize energy cost in arm movements. *J Neurophysiol* **104**:2985-2994 <https://doi.org/10.1152/jn.00483.2010> | [PubMed](#)
- Kobak D, Brendel W, Constantinidis C, Feierstein CE, Kepecs A, Mainen ZF, Qi XL, Romo R, Uchida N, Machens CK (2016) Demixed principal component analysis of neural population data. *eLife* **5**:9424 <https://doi.org/10.7554/eLife.10989> | [PubMed](#)
- Leow LA, Marinovic W, de Rugy A, Carroll TJ (2020) Task errors drive memories that improve sensorimotor adaptation. *J Neurosci* **40**:3075-3088 <https://doi.org/10.1523/JNEUROSCI.1506-19.2020> | [PubMed](#)
- Leow LA, de Rugy A, Marinovic W, Riek S, Carroll TJ (2016) Savings for visuomotor adaptation require prior history of error, not prior repetition of successful actions. *J Neurophysiol* **116**:1603-1614 <https://doi.org/10.1152/jn.01055.2015> | [PubMed](#)
- Lillicrap TP, Hunt JJ, Pritzel A, Heess N, Erez T, Tassa Y, Silver D, Wierstra D. (2015) Continuous control with deep reinforcement learning. *arXiv* <https://doi.org/10.48550/arXiv.1509.02971>
- Losey DM, Hennig JA, Oby ER, Golub MD, Sadtler PT, Quick KM, Ryu SI, Tyler-Kabara EC, Batista AP, Yu BM, *et al.* (2024) Learning leaves a memory trace in motor cortex. *Current Biology* **34**:1519-1531.e4. <https://doi.org/10.1016/j.cub.2024.03.003> | [PubMed](#)
- McDougle SD, Bond KM, Taylor JA (2015) Explicit and Implicit Processes Constitute the Fast and Slow Processes of Sensorimotor Learning. *J Neurosci* **35**:9568-9579 <https://doi.org/10.1523/jneurosci.5061-14.2015> | [PubMed](#)
- Michaels JA, Kashefi M, Zheng J, Codol O, Weiler J, Kersten R, Gribble PL, Diedrichsen J, Pruszynski JA (2024) Sensory expectations shape neural population dynamics in motor circuits. *bioRxiv* <https://doi.org/10.1101/2024.12.22.629295>
- Michaels JA, Schaffelhofer S, Agudelo-Toro A, Scherberger H. (2020) A goal-driven modular neural network predicts parietofrontal neural dynamics during grasping. *Proc Natl Acad Sci U S A* **117**:32124-32135 <https://doi.org/10.1073/pnas.2005087117> | [PubMed](#)
- Morehead JR, Qasim SE, Crossley MJ, Ivry R. (2015) Savings upon Re-Aiming in Visuomotor Adaptation. *J Neurosci* **35**:14386-14396 <https://doi.org/10.1523/jneurosci.1046-15.2015> | [PubMed](#)
- Nguyen KP, Zhou W, McKenna E, Colucci-Chang K, Bray LCJ, Hosseini EA, Alhussein L, Rezazad M, Joiner WM (2019) The 24-h savings of adaptation to novel movement dynamics initially reflects the recall of previous performance. *J Neurophysiol* **122**:933-946 <https://doi.org/10.1152/jn.00569.2018> | [PubMed](#)
- Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, *et al.* (2019) PyTorch: An imperative style, high-performance deep learning library. *arXiv* <https://doi.org/10.48550/arxiv.1912.01703>
- Shadmehr R, Mussa-Ivaldi FA (1994) Adaptive representation of dynamics during learning of a motor task. *J Neurosci* **14**:3208-3224 <https://doi.org/10.1523/jneurosci.14-05-03208.1994> | [PubMed](#)
- Sheahan HR, Franklin DW, Wolpert DM (2016) Motor Planning, Not Execution, Separates Motor Memories. *Neuron* **92**:773-779 <https://doi.org/10.1016/j.neuron.2016.10.017> | [PubMed](#)
- Sheahan HR, Ingram JN, Žalalytė GM, Wolpert DM (2018) Imagery of movements immediately following performance allows learning of motor skills that interfere. *Scientific reports* **8**:14330 <https://doi.org/10.1038/s41598-018-32606-9> | [PubMed](#)
- Smoulder AL, Marino PJ, Oby ER, Snyder SE, Miyata H, Pavlovsky NP, Bishop WE, Yu BM, Chase SM, Batista AP (2024) A neural basis of choking under pressure. *Neuron* **112** <https://doi.org/10.1016/j.neuron.2024.08.012> | [PubMed](#)

- Sun X, OShea DJ, Golub MD, Trautmann EM, Vyas S, Ryu SI, Shenoy KV (2022) Cortical preparatory activity indexes learned motor memories. *Nature* **602**:274-279 <https://doi.org/10.1038/s41586-021-04329-x> | PubMed
- Sussillo D, Abbott LF (2009) Generating coherent patterns of activity from chaotic neural networks. *Neuron* **63**:544-557 <https://doi.org/10.1016/j.neuron.2009.07.018> | PubMed
- Sussillo D, Churchland MM, Kaufman MT, Shenoy KV (2015) A neural network that finds a naturalistic solution for the production of muscle activity. *Nat Neurosci* **18**:1025-1033 <https://doi.org/10.1038/nn.4042> | PubMed
- Sutskever I, Martens J, Dahl G, Hinton G. (2013) On the importance of initialization and momentum in deep learning. In: Proceedings of the 30th International Conference on Machine Learning. Georgia, United States: Pmlr. pp. 1139-1147
- Trautmann EM, Hesse JK, Stine GM, Xia R, Zhu S, O'Shea DJ, Karsh B, Colonell J, Lanfranchi FF, Vyas S, et al. (2025) Large-scale high-density brain-wide neural recording in nonhuman primates. *Nat Neurosci* **28**:1562-1575 <https://doi.org/10.1038/s41593-025-01976-5> | PubMed
- Wang T, Ivry RB (2025) Contextual effects during sensorimotor adaptation are an emergent property of population coding in a cerebellar-inspired model. *Sci Adv* **11**:eadr4540 <https://doi.org/10.1126/sciadv.adr4540> | PubMed
- Werbos PJ (1990) Backpropagation through time: what it does and how to do it. *Proc IEEE Inst Electr Electron Eng* **78**:1550-1560 <https://doi.org/10.1109/5.58337>
- Wolpert D, Kawato M. (1998) Multiple paired forward and inverse models for motor control. *Neural Netw* **11**:1317-1329 [https://doi.org/10.1016/S0893-6080\(98\)00066-5](https://doi.org/10.1016/S0893-6080(98)00066-5) | PubMed
- Wolpert DM, Ghahramani Z, Jordan MI (1995) An internal model for sensorimotor integration. *Science* **269**:1880-1882 <https://doi.org/10.1126/science.7569931> | PubMed
- Yin C, Wei K. (2020) Savings in sensorimotor adaptation without explicit strategy. *J Neurophysiol* <https://doi.org/10.1152/jn.00524.2019> | PubMed
- Zarahn E, Weston GD, Liang J, Mazzoni P, Krakauer JW (2008) Explaining savings for visuomotor adaptation: linear time-invariant state-space models are not sufficient. *J Neurophysiol* **100**:2537-2548 <https://doi.org/10.1152/jn.90529.2008> | PubMed

Peer reviews

Reviewer #2 (Public review):

Summary:

Shahbazi et al. trained recurrent neural networks (RNNs) to simulate human upper limb movement during adaptation to a force field perturbation. They demonstrated that throughout adaptation, the pattern of motor commands to the muscles of the simulated arm changed, allowing the perturbed movements to regain their typical, perturbation-free straight-line paths. After this initial learning block (FF1), the network encountered null-fields to wash out the adaptation, before re-experiencing the force in a second learning block (FF2). Upon re-exposure, the network learned faster than during initial learning, consistent with the savings observed in behavioral studies of adaptation. They also found that as the number of hidden units in the RNN increased, so did the probability of exhibiting savings. The authors concluded that these results propose a neural basis for savings that is independent of context and strategic processes.

Strengths:

The paper addresses an important and controversial topic in motor adaptation: the mechanism underlying motor memory. The RNN simulation reproduces behavioral hallmarks of adaptation, and it provides a useful illustration of the pattern of muscle activity

underlying human-like movements under both normal and perturbing conditions. While the savings effect produced by the network, though significant, appears somewhat small, the simulation demonstrating an increase in savings with a greater number of hidden units is particularly intriguing.

Main weakness:

The introduction details the ongoing debate in the literature regarding the mechanisms underlying savings, particularly whether it stems from explicit or implicit learning processes. However, it remains unclear how the current work addresses this debate. There is already a considerable body of research, particularly in visuomotor adaptation, demonstrating that savings is predominantly driven by explicit strategies (e.g., Morehead et al. 2015, Haith et al., 2015; Huberdeau et al., 2019; Avraham et al., 2021). Furthermore, there have been multiple reports that implicit adaptation exhibits attenuation upon relearning (Avraham et al., 2021, Leow et al., 2020; Yin and Wei, 2020; Hamel et al., 2021; Hamel et al., 2022; Wang and Ivry, 2023; Hadjiosif et al., 2023). In the discussion, the authors acknowledge that their goal was not to model a complete explicit-implicit system, but rather to probe how savings may emerge from a purely implicit mechanism. Given the central debate introduced by the authors, the manuscript would benefit from a more detailed discussion explaining how their findings elucidate the specific conditions under which savings emerge from purely implicit mechanisms versus when cognitive strategies predominate.

<https://doi.org/10.7554/eLife.107423.2.sa1>

Author response:

The following is the authors' response to the original reviews.

Public Reviews:

Reviewer #1 (Public review):

Summary:

Shahbazi et al used a recurrent neural network model trained to control a musculoskeletal model of the arm to investigate how neural populations accommodate activity patterns underpinning savings. The paper draws upon the recent finding of a "uniform shift" in preparatory activity in monkey motor cortex associated with savings, and leverages full access to a computational model to establish causality.

Strengths:

The paper is well written, and the figures are clearly presented. The key finding that the uniform shift first reported based on neural recordings by Sun et al. emerges in artificial neural networks performing a similar task is interesting and well-backed by their analyses. Manipulating this uniform shift to show that it drives behavioural savings is an important causal confirmation of the proposal by Sun et al.

Weaknesses / Comments:

As mentioned earlier, the core results are well backed by the analyses. Most of my comments relate to adding more controls and additional questions that could be explored with the model to strengthen the paper.

(1) Savings are quantified as more rapid relearning of the FF upon re-exposure (e.g., Figure 3). This finding is based on backpropagation through time, but would this hold when using a different optimiser, e.g., FORCE?

This is an interesting question, and indeed, there are an increasing number of studies addressing how different neural network learning rules may affect the kinds of representations that arise after learning (Codol et al., 2024). However the focus of the present paper is not on which neural network approaches or which specific optimisers produce savings, rather, the focus is on the basis and neural geometry of savings when it emerges.

We have added a short paragraph to the Discussion section [lines 349-355] to address this:

“The present results are based on RNNs trained in an error-based approach using backpropagation through time (Werbos, 1990) using the Adam optimizer (Kingma and Ba, 2014). Other techniques for training RNNs have been proposed including the FORCE algorithm (Sussillo and Abbott, 2009). In addition, several recent reports have demonstrated success using reinforcement learning approaches to train neural networks in the context of sensorimotor control tasks (Lillicrap et al., 2015; Codol et al., 2024a). An interesting avenue for future work is to determine how the present results may or may not generalize to different neural network architectures and learning rules.”

(2) The authors should include a "null model" showing that training on a different reaching task following NF, as opposed to FF2, won't show something akin to a uniform shift during preparation due to the adoption of TDR and having similar targets.

This is a critical point. Training on a different reaching task other than FF2 (e.g. a different force field) will indeed result in a uniform shift, but critically, a shift in a different direction in neural state space than the uniform shift associated with FF2. The central focus of the present paper is to show that when there remains a non-zero projection of preparatory neural activity along the direction of the uniform shift associated with a given learning task, this residual projection underlies savings when networks are subsequently re-exposed to the same task.

In the Results section we had included a short paragraph to describe control simulations that we performed that address this concept. We have expanded this text and added a Figure and the results of statistical tests to better describe this control [lines 179-187]:

“As an additional control we trained networks after the growing up phase on an opposing force field (CCW) and then as above, exposed the networks to a NF washout phase, and then to a CW force field. In this case no savings was observed in the CW force field, either for initial lateral deviation, or for learning rate (Figure 3). In fact, we observed that initial lateral deviation is larger for the novel force field ($t(39)=-4.918$, $p=1.6e-5$). This observation is in line with the finding that learning opposing force fields sequentially results in interference (Sun et al., 2022). The results of these control simulations underscore that the savings effect observed in our main study was learning-specific—it was due to prior learning of the CCW force field, and not a general effect of learning any novel dynamics.”

(3) The analyses of network activity during movement preparation (Figure 4) nicely replicate the key finding in Sun et al, but I think the authors could leverage the full access to their network and go further, e.g., by examining changes (or the lack of) during execution in FF2 with respect to FF (and perhaps in a future NF2 with respect to NF), including whether execution activity lives also lives in parallel hyperplanes, etc.

We agree that a visualization of the neural activity during movement would be beneficial to the reader. To address this we have added a new Figure (Fig. 6) and associated text [lines 210-219]. The Figure shows the neural trajectories when the RNNs are first exposed to the FF1 and when they are first exposed to FF2 (after NF2 washout). Trajectories are plotted in 3D corresponding to the first 3 principal components, starting at the go cue and ending 200 ms into the movement, for each of the 8 movement targets.

“The neural trajectories for preparation and for movement can be visualized in principal component space. Figure 6 shows trajectories during planning and early execution for initial FF1 and FF2 exposure. Hidden unit activity was subjected to a principal components analysis, and neural trajectories within the first three PCs are shown for movements to each of the eight movement targets. Filled circles indicate neural state 200 ms prior to the go cue. During the preparatory period trajectories travel along PC1 and then disperse across PC2 and PC3 into the circular pattern indicated by the filled stars, which indicate time of the go cue (also see Figure 5A). After the go cue neural trajectories shift back along PC1 and rotate along oscillatory patterns characteristic of populations of motor cortical neurons in non-human primates during movement (Churchland and Shenoy, 2024).”

(4) Related to the above, while the results are interesting and the paper is well done, I kept wishing that the authors had done "more" with their model. This could be one or two final sections on "predictions" that would nicely complement their "validation" of the uniform shift, and that, in my opinion, would greatly increase the impact of the paper. In particular:

(a) What would be the effect of learning more "tasks"? For example, is there a limit on how many fields can be learned? (You show something related by manipulating network size, but this is slightly different.)

These are interesting questions and to some extent they are already addressed in the paper. Of course, the number of tasks that a network is able to learn, will be related to how much those tasks overlap in a control space. Indeed, this idea goes back to early theoretical accounts of connectionist models such as Hopfield nets and capacity for representing information (Hopfield, 1982; Hopfield et al., 1983). The control simulations that we described in the paper [lines 179-187 and Figure 4] are a test of one extreme version of this, in which two tasks are in direct opposition to each other (opposite force fields), and in this situation no savings emerges. We believe it is an interesting question, but beyond the scope of the present paper to undertake a comprehensive exploration of the nature of task-overlap in upper limb reaching learning tasks.

(b) Figure 5 is a nice causal demonstration that the uniform shift is related to savings. However, and related to comment #3, it'd be interesting to see more details about how the behaviour and the network activity changes as preparatory activity shifts along this axis, in particular regarding how moving the preparatory states affect the organisation and dynamics of upcoming execution activity -these are the kind of intuitions that modelling studies like this one can provide.

This has been addressed above by the changes we made to address the reviewer's comment #3.

(c) The authors focus on a task design that spans baseline, FF, NF, FF2 to replicate the original study by Sun et al. However, it would be interesting if they generated predictions for neural changes to other types of tasks that have been studied behaviourally. These could include, for example: (i) modelling a visuomotor rotation or a mirror reversal task; (ii) having to adapt to a FF in the opposite direction; (iii) investigating the role of adding an explicit context and having the networks learn multiple FF; and (iv) trying to learn FF fields in opposite directions, perhaps restricted to specific targets. As the authors know, all these questions and more have been studied with similar behavioural paradigms, and it would be nice to see what neural predictions are generated by this model.

See responses above e.g. to comment 4. We have clarified the text and provided a new Figure to illustrate our opposite FF control simulations. The other suggestions about visuomotor rotations, and contextual cues, are interesting and potentially important questions that we

are working on, but we believe are beyond the scope of the current paper which is focused specifically around the question of savings in FF learning.

(5) On the Discussion: When extrapolating from neural network results to animals, the fact that your networks can learn implicitly doesn't mean that animals do learn implicitly. Indeed, I think the consensus view is that different perturbations may lead to the expression of different types of savings (e.g., FF vs VR, which seems to be more explicit). Besides, these different mechanisms may be primarily implemented by brain regions less directly tied to motor control (e.g., cerebellum, parietal cortex?), which are not directly implemented in the authors' model.

Of course the reviewer is correct that our simulations are not evidence that savings in motor tasks learned by animals is only implicit, and we do not make any such claims in the paper. The model we describe in the present paper is not meant to be a comprehensive model of motor learning in humans/animals. Indeed, the pure “context free” type of learning that we implement in our simulations basically cannot occur in animals, because there is always some information that provides contextual information. Indeed there are computational models of motor learning that include these effects, e.g. the COIN model (Heald et al., 2021). Our model however provides a useful window into what the context-free component of savings may look like. The approach we describe in the present paper is a powerful way to probe the context-free component of savings in isolation in a way that is not possible (at least not readily) in animals/humans. We have modified the text in the Discussion [lines 372-379] to better articulate this point.

“The simulations described here do not constitute evidence that savings in motor learning tasks is exclusively implicit in animals and humans. The purely context-free learning implemented in our simulations is highly unrealistic, as some form of contextual information is invariably available. Indeed, computational models of motor learning that incorporate contextual effects already exist, e.g. (Heald et al. 2021). Nevertheless, our simulations provide a useful window into what the context-free component of savings may look like. This approach offers a powerful means of probing the context-free component of savings in isolation—something that is not readily achievable in animal or human experiments.”

Reviewer #2 (Public review):

Summary:

Shahbazi et al. trained recurrent neural networks (RNNs) to simulate human upper limb movement during adaptation to a force field perturbation. They demonstrated that throughout adaptation, the pattern of motor commands to the muscles of the simulated arm changed, allowing the perturbed movements to regain their typical, perturbation-free straight-line paths. After this initial learning block (FF1), the network encountered null-fields to wash out the adaptation, before re-experiencing the force in a second learning block (FF2). Upon re-exposure, the network learned faster than during initial learning, consistent with the savings observed in behavioral studies of adaptation. They also found that as the number of hidden units in the RNN increased, so did the probability of exhibiting savings. The authors concluded that these results propose a neural basis for savings that is independent of context and strategic processes.

Strengths:

The paper addresses an important and controversial topic in motor adaptation: the mechanism underlying motor memory. The RNN simulation reproduces behavioral hallmarks of adaptation, and it provides a useful illustration of the pattern of muscle activity underlying human-like movements under both normal and perturbing conditions. While the savings effect produced by the network, though significant, appears

somewhat small, the simulation demonstrating an increase in savings with a greater number of hidden units is particularly intriguing.

Weaknesses:

(1) To be transparent, savings in motor adaptation have been a primary focus of my own research. Some core findings presented in this paper are at odds with the ideas I and others have previously put forward. While I don't want to impose my agenda on the authors of this paper, I do think the authors should address these issues.

(a) The authors acknowledge the ongoing debate in the literature regarding the mechanisms underlying savings, particularly whether it stems from explicit or implicit learning processes. However, it remains unclear how the current work addresses this debate. There is already a considerable body of research, particularly in visuomotor adaptation, demonstrating that savings is predominantly driven by explicit strategies. For example, when people are asked to report their strategy, they recall a strategy that was useful during the first learning block (Morehead et al. 2015). Furthermore, savings are abolished under experimental manipulations designed to eliminate strategic contributions (e.g., Haith et al., 2015; Huberdeau et al., 2019; Avraham et al., 2021). The authors briefly state that their findings support the hypothesis that a neural basis of memory retention underlying savings can be independent of cognitive or strategic learning components, and that savings can be characterized as implicit. While these statements may be true, it is not clear how this work substantiates these claims.

We have addressed a similar point raised by Reviewer 1, see point #5 above. Our work represents an example of how savings can occur from implicit mechanisms in the absence of explicit contextual cues. Our goal is not to resolve the debate about how this occurs in humans/animals. Rather, our model provides a useful window into what the context-free component of savings may look like. Our approach is a powerful way to probe the context-free component of savings in isolation in a way that is not possible (at least not readily) in animals/humans. We have modified the text in the Discussion [lines 372-379] to better articulate this point.

“The simulations described here do not constitute evidence that savings in motor learning tasks is exclusively implicit in animals and humans. The purely context-free learning implemented in our simulations is not meant to be a full model of biological learning, as in biological systems some form of contextual information is invariably available. Indeed, computational models of motor learning that incorporate contextual effects already exist, e.g. (Heald et al. 2021). Nevertheless, our simulations provide a useful window into what the context-free component of savings may look like. This approach offers a powerful means of probing the context-free component of savings in isolation—something that is not readily achievable in animal or human experiments.”

(b) Our research has also demonstrated that if implicit adaptation is completely washed out after the initial learning block, it not only fails to exhibit savings but is actually attenuated relative to the first learning block (Avraham et al., 2021). This phenomenon of attenuation upon relearning can also be seen in other studies of visuomotor adaptation (e.g., Leow et al., 2020; Yin and Wei, 2020; Hamel et al., 2021; Hamel et al., 2022; Wang and Ivry, 2023; Hadjiosif et al., 2023). More recently, we have shown that this attenuation is due to anterograde interference arising from the experience with the washout block experience (Avraham and Ivry, 2025). We illustrated that the implicit system is highly susceptible to interference; it doesn't require exposure to salient opposite errors and can occur even following prolonged exposure to veridical feedback. The central thesis of this paper, namely that implicit savings can emerge through RNNs, is at odds with these empirical results. The authors should address this discrepancy.

These empirical results are interesting and intriguing, and we agree that they are relevant in the context of the debate about the relative contributions and interactions between explicit and implicit learning systems and savings. Importantly, contextual interference is impossible in our model, since there are no contextual cues about which force field is present or absent. Interactions between an explicit system and an implicit learning system are also impossible in our model, since there is no possibility of context-driven explicit learning or memory. The approach we have taken in the present paper is not to model a full explicit plus implicit learning system but rather to probe how savings may emerge from a purely implicit learning mechanism alone and to compare the neural geometry underlying this implicit-drive savings to the neural recording results from monkey electrophysiology studies. Nevertheless we have added some text to the Discussion [lines 380-391] to situate our findings in the context of the studies mentioned above by the reviewer.

“Recent empirical work suggests that relearning after washout of implicit adaptation can be attenuated rather than facilitated, a phenomenon attributed to anterograde interference from the washout phase (Avraham et al., 2021; Hadjiosif et al., 2023; Hamel et al., 2022, 2021; Leow et al., 2020; Wang and Ivry, 2025; Yin and Wei, 2020). The savings observed in our simulations differs from these behavioral findings. Crucially, our model excludes both contextual interference (since no cues signal which force field is present) and explicit-implicit interactions (since context-driven explicit learning is absent). Our goal was not to model a complete explicit-implicit system, but rather to probe how savings may emerge from a purely implicit mechanism and to compare the underlying neural geometry to monkey electrophysiology data. Our results suggest that high-dimensional neural circuits possess an intrinsic capacity for savings via persistent preparatory traces. How and when this capacity may be masked by interference or explicit-implicit interactions in biological systems remains an open question for future work.”

(2) This brings me to the question about neural correlates: The results are linked to activity in the primary motor cortex. How does that align with the well-established role of the cerebellum in implicit motor adaptation? And with the studies showing that savings are due to explicit strategies, which are generally associated with prefrontal regions?

The modeling approach we use in the present paper is area agnostic, and we do not include different neural modules to represent specific brain areas such as cerebellum or prefrontal regions. In the current approach we specifically exclude explicit strategies, as a way to specifically probe implicit mechanisms alone. Also see response to reviewer 1 comment 5 above.

(3) The analysis on the complexity of the neural network (i.e., the number of hidden units) and its relationship to savings is very interesting. It makes sense to me that more complex networks would show more savings. I'm not sure I follow the author's explanation, but my understanding is that increased network complexity makes it more difficult to override the formed memory through interference (e.g., from the experience with NF2). Also, the results indicate that a network with 32 units led to a less-than-chance level of networks exhibiting savings (Figure 3b). What behavioral output does this configuration produce? Could this behavior manifest as attenuation upon relearning? Furthermore, if one were to examine an even smaller, simpler network (perhaps one more closely reflecting cerebellar circuits), would such a model predict attenuation rather than savings?

These are interesting questions, and are potentially important, for future work to explore. Our interpretation of the results of smaller networks is that these small RNNs fail to show savings presumably because the learned FF behavior is 'erased' during washout because of the limited capacity to retain the FF learning in a distinct neighborhood in neural state space. Our paper is focused specifically on the relationship between savings, implicit learning, and

neural capacity via network size, in the context of the monkey electrophysiology results in motor cortex. It would be interesting in future work to explore a cerebellar-like modeling approach.

(4) The authors emphasize that their network did not receive any explicit contextual signals related to the presence or absence of the force field (FF), thus operating in a 'context-free' manner. From my understanding, some existing models of context's role in motor memories (e.g., Oh and Schweighofer, 2019; Heald et al., 2021) propose that memory-related changes can be observed even without explicit contextual information, as contextual changes can be inferred from sudden or significant environmental shifts (e.g., the introduction or removal of perturbations). Given this, could the observed savings in the current simulation be explained by some form of contextual retrieval, inferred by the network from the re-representation of the perturbation in FF2?

It is important to note that this is not possible in the context of the modeling approach described in the present paper. For example, in trial 1 of FF2, because the network has no contextual cue signaling the FF's presence, the network has no information before movement begins that a FF will be present during movement (recall that the FF is velocity-dependent, and so is zero before movement begins). Once the network encounters the FF during movement, some component of its response I suppose could be described as contextual inference derived from effector state (similar to the account described in the COIN model), but strictly speaking the model is only responding to what it encounters in the moment. Any change in behaviour due to prior learning (e.g. savings) is due to the interaction between the residual learning-related neural state (e.g. the uniform shift), the effector state in the moment, and the errors encountered during movement. We don't interpret this as "inference" in the traditional sense of an explicit learning system.

(5) If there is residual hidden unit activity related to the FF at the end of the NF2 phase, how does the simulated movement revert back to baseline? Are there any differences in the movement trajectory, beyond just lateral deviation, between NF1 and NF2? The authors state that "changes in the preparatory hidden unit activity did not result in substantive changes in the motor commands (Figure 5b), which emphasizes that the uniform shift resides in the null space of motor output." However, Figure 5b appears to show visible changes in hidden unit activity. Don't these changes reflect a pattern of muscle activity that is the basis for behavior? These changes are indeed small, but it seems that so is the effect size for savings (Figure 3a). Could this suggest that there is not, in fact, a complete washout of initial learning during NF2 within the network?

This is precisely the point of the paper, i.e. to show that neural activity during the preparatory period before movement onset is different, even though the behaviour during the preparatory period is the same (i.e. no muscle activity and no movement). This recapitulates the empirical findings from the neural data reported in the Sun et al. (2022) paper.

The reviewer asks "Don't these changes reflect a pattern of muscle activity that is the basis for behavior?" Yes indeed they do, but not during the NF and not during the preparatory activity prior to movement onset.

The reviewer asks "Could this suggest that there is not, in fact, a complete washout of initial learning during NF2 within the network?" We addressed this in the paper (Results/Washout) by comparing kinematics after washout to that prior to FF learning; e.g. any differences in lateral deviation of the hand path for the entire reach trajectory was in the range of 0.1 mm, which is less than 0.25 % of the lateral deviation encountered in the FF and only 0.1 % of the reach distance (10 cm).

Recommendations for the authors:

Reviewer #2 (Recommendations for the authors):

(1) *Figure 1c, lower panel: Is this from the early or late stage of FF1?*

This is an example movement after learning in a null field (NF). We have clarified this in the Figure caption.

(2) *Please clarify what the two panels in Figure 1e represent.*

We have clarified in the Figure caption that these are activity from two example hidden units.

(3) *If Figure 2c is intended to illustrate the changes in motor commands for individual muscles, consider reorganizing the plots by muscle to more clearly show the change for each muscle from NF1 to FF1.*

The point here is not to make fine-grained comparisons between specific muscles, rather to show a general example of how muscle activity is different. For the sake of visual simplicity in a Figure that already has many components we have decided to keep Figure 2c the same.

(4) *The text mentions that no savings were observed when the network was trained on CCW followed by CW perturbations. However, no data or statistical analysis is presented to support this claim. I wonder if the authors would expect attenuated learning when exposed to the CW perturbation, given a memory of the opposite perturbation.*

We have added a Figure to provide data for the FF opposite control.

(5) *The relevance of the discussion on choking under pressure to the paper wasn't clear.*

We have modified the relevant text in the Discussion section [lines 356-363] to clarify the relevance of the present work to other recent work on how complex features of motor behaviour can arise due to the dynamics of preparatory neural activity in motor cortex.

References

- Avraham G, Morehead JR, Kim HE, Ivry RB. 2021. Reexposure to a sensorimotor perturbation produces opposite effects on explicit and implicit learning processes. *PLoS Biol* 19:e3001147. doi:10.1371/journal.pbio.3001147
- Codol O, Krishna NH, Lajoie G, Perich MG. 2024. Brain-like neural dynamics for behavioral control develop through reinforcement learning. *bioRxiv*. doi:10.1101/2024.10.04.616712
- Hadjiosif AM, Morehead JR, Smith MA. 2023. A double dissociation between savings and long-term memory in motor learning. *PLoS Biol* 21:e3001799. doi:10.1371/journal.pbio.3001799
- Hamel R, Dallaire-Jean L, De La Fontaine É, Lepage JF, Bernier PM. 2021. Learning the same motor task twice impairs its retention in a time- and dose-dependent manner. *Proc Biol Sci* 288:20202556. doi:10.1098/rspb.2020.2556
- Hamel R, Lepage J-F, Bernier P-M. 2022. Anterograde interference emerges along a gradient as a function of task similarity: A behavioural study. *Eur J Neurosci* 55:49–66. doi:10.1111/ejn.15561
- Heald JB, Lengyel M, Wolpert DM. 2021. Contextual inference underlies the learning of sensorimotor repertoires. *Nature* 600:489–493. doi:10.1038/s41586-021-04129-3
- Hopfield JJ. 1982. Neural networks and physical systems with emergent collective computational abilities. *Proc Natl Acad Sci U S A* 79:2554–2558. doi:10.1073/pnas.79.8.2554

Hopfield JJ, Feinstein DI, Palmer RG. 1983. “Unlearning” has a stabilizing effect in collective memories. *Nature* 304:158–159. doi:10.1038/304158a0

Leow L-A, Marinovic W, de Rugy A, Carroll TJ. 2020. Task errors drive memories that improve sensorimotor adaptation. *J Neurosci* 40:3075–3088. doi:10.1523/JNEUROSCI.1506-19.2020

Wang T, Ivry RB. 2025. Contextual effects during sensorimotor adaptation are an emergent property of population coding in a cerebellar-inspired model. *Sci Adv* 11:eadr4540. doi:10.1126/sciadv.adr4540

Yin C, Wei K. 2020. Savings in sensorimotor adaptation without an explicit strategy. *J Neurophysiol* 123:1180–1192. doi:10.1152/jn.00524.2019

<https://doi.org/10.7554/eLife.107423.2.sa0>