

In the format provided by the authors and unedited.

Learning by neural reassociation

Matthew D. Golub ^{1,2,3}, Patrick T. Sadtler^{2,4,5}, Emily R. Oby^{2,4,5}, Kristin M. Quick^{2,4,5}, Stephen I. Ryu^{3,6}, Elizabeth C. Tyler-Kabara ^{4,7,8}, Aaron P. Batista^{2,4,5}, Steven M. Chase ^{2,9,10*} and Byron M. Yu ^{1,2,9,10*}

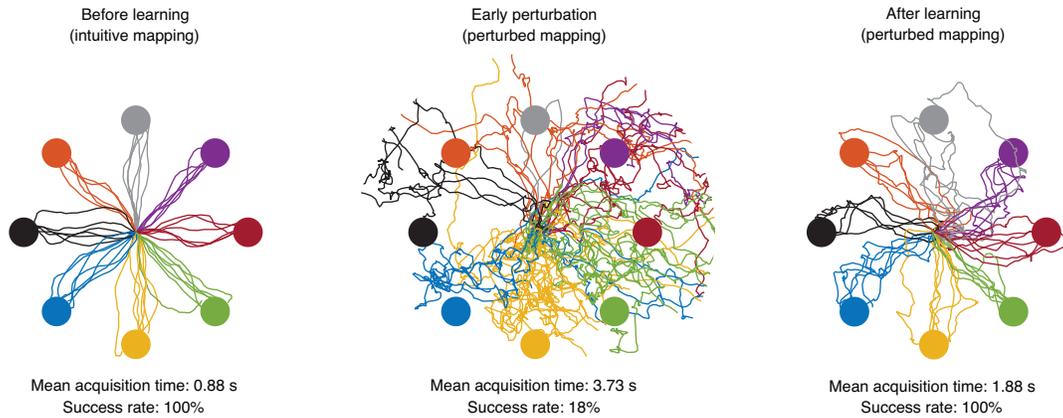
¹Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, USA. ²Center for the Neural Basis of Cognition, Carnegie Mellon University, Pittsburgh, PA, USA. ³Department of Electrical Engineering, Stanford University, Stanford, CA, USA. ⁴Department of Bioengineering, University of Pittsburgh, Pittsburgh, PA, USA. ⁵Systems Neuroscience Institute, University of Pittsburgh, Pittsburgh, PA, USA. ⁶Department of Neurosurgery, Palo Alto Medical Foundation, Palo Alto, CA, USA. ⁷Department of Physical Medicine and Rehabilitation, University of Pittsburgh, Pittsburgh, PA, USA. ⁸Department of Neurological Surgery, University of Pittsburgh, Pittsburgh, PA, USA. ⁹Department of Biomedical Engineering, Carnegie Mellon University, Pittsburgh, PA, USA. ¹⁰These authors contributed equally: Steven M. Chase, Byron M. Yu. *e-mail: schase@cmu.edu; byronyu@cmu.edu

SUPPLEMENTARY FIGURES

Supplementary Figure 1

Cursor movements from an example experiment.

Movements are from example experiment J20120525. “Before learning” trials are the last 50 trials under the intuitive BCI mapping (left gray window in **Fig. 1c**). “Early perturbation” trials are the first 50 trials under the perturbed BCI mapping (not analyzed in this work). “After learning” trials are the 50 consecutive trials of peak behavioral performance under the perturbed mapping (right gray window in **Fig. 1c**). Mean acquisition times are given for successful trials only. Cursor positions exceeding a cutoff distance from the workspace center are omitted from view in this figure.



Supplementary Figure 2

Selection and characterization of perturbed BCI mappings.

To select the perturbed BCI mapping for a given experiment, neural activity from intuitive trials was used to predict the behavioral effects of millions of candidate perturbed mappings. To generate these predictions, we computed per-target average population activity patterns and passed them through every candidate perturbed mapping, resulting in per-target open-loop cursor velocity predictions about how the animal would move the cursor in the absence of visual feedback and before any learning had taken place. As a basis for comparison, we also computed open-loop velocities using the intuitive BCI mapping. We compared these open-loop velocities from the intuitive mapping to those from the candidate perturbed mappings. Then, we excluded candidate mappings predicted to be too easy or too difficult to learn. Exclusion criteria were adjusted such that 10-100 candidate mappings remained. From these, we arbitrarily selected one mapping to be used in the experiment. See the Supplementary Math Note for further details.

(a) Example experiment J20120628. Open-loop velocities based on the intuitive mapping (circles) aligned with the cursor-to-target directions (lines), indicating proficient performance under the intuitive mapping. Exclusion criteria eliminated all but 14 candidate perturbed mappings (open-loop velocities indicated by small squares), from which one was chosen for use during the experiment. Open-loop velocities from the chosen perturbed mapping (large squares) predict complex behavioral effects, which vary across target directions (e.g., counter-clockwise rotations for directions indicated by black and blue, but clockwise rotations for the other directions; nonuniform speed scalings). These behavioral effects are in contrast with the uniform effects of a visuomotor rotation (i.e., a uniform rotation across movement directions; unity speed scaling). Note that, because open-loop velocities do not include a contribution from a previous timestep (i.e., $\mathbf{A}\mathbf{v}_{t-1}$ in equation (1)), they tend to be slower than cursor velocities from the actual experiments (which used the full equation (1)).

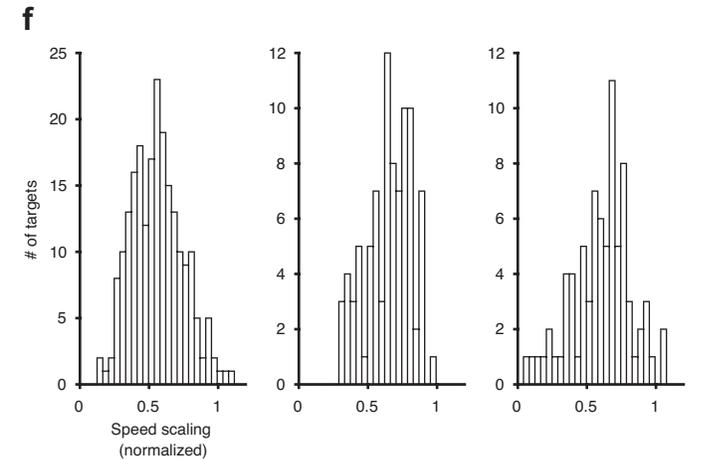
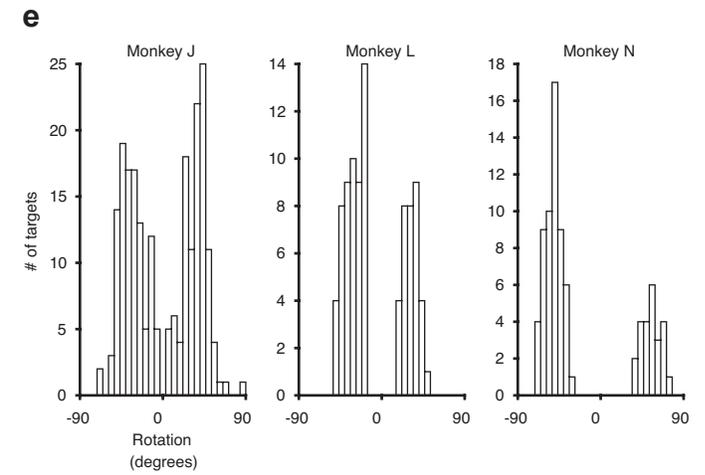
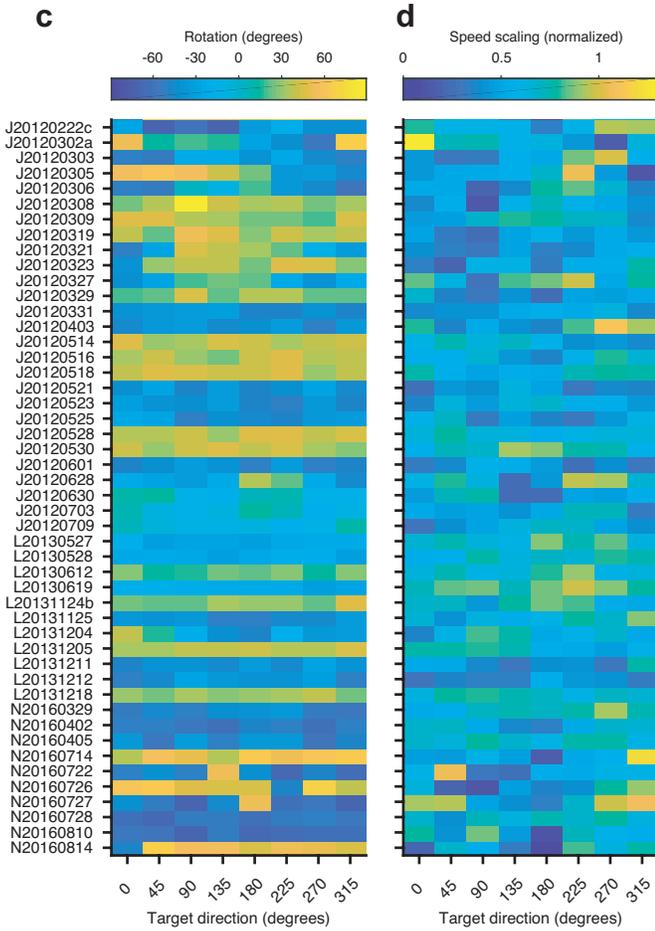
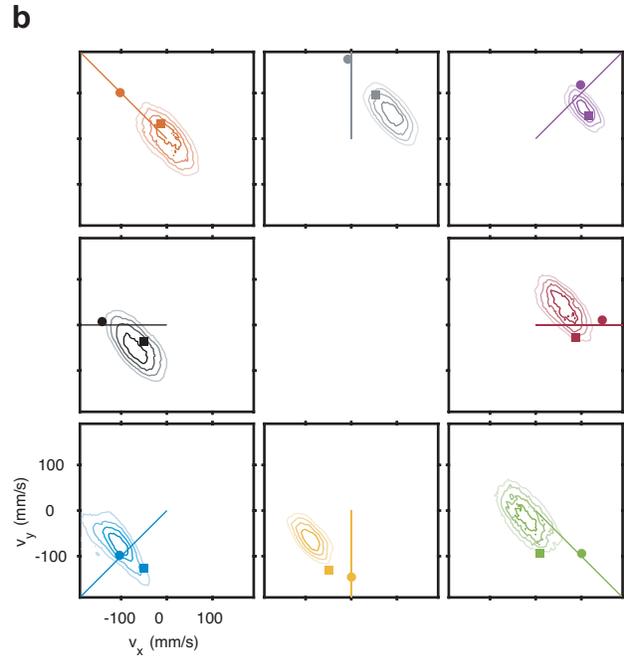
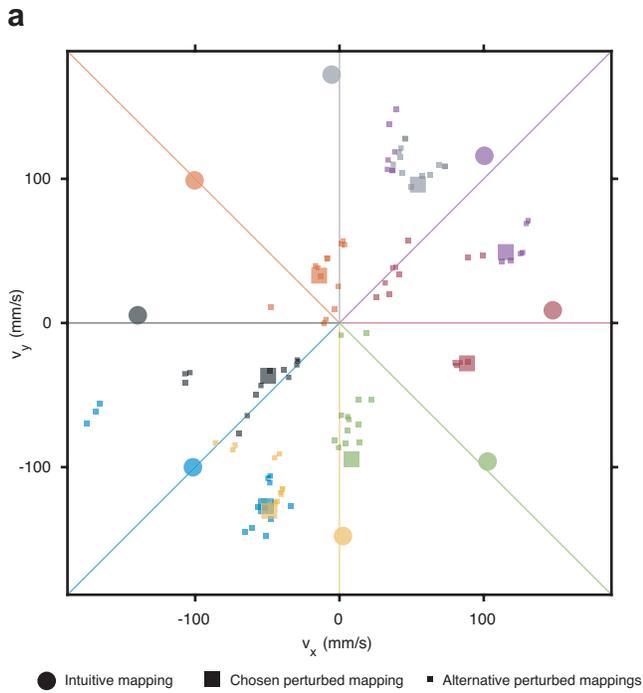
(b) Distributions of per-target open-loop velocities from all 3,628,800 candidate perturbed mappings (i.e., prior to applying exclusion criteria). Each sub-panel uses only the neural activity from trials to a particular target (corresponding to the colors in **a**). Each candidate perturbed mapping jointly specifies a set of open-loop velocities (i.e., one per sub-panel). Contour lines indicate uniformly spaced density levels, and darker lines represent regions of higher density. Open-loop velocities from the intuitive mapping (circles) and chosen perturbed mapping (squares), along with target directions (lines), are replicated from **a** for reference.

(c) Predicted per-target open-loop rotations for the chosen perturbed mapping from each experiment. For a given experiment (row) and target direction (column), the rotation was computed as the signed angle between the intuitive and perturbed open-loop velocity. Note the non-uniformity of these rotations across target directions in each experiment (colors vary across each row).

(d) Predicted per-target open-loop speed scalings for the chosen perturbed mapping from each experiment. For a given experiment (row) and target direction (column), the speed scaling was computed as the speed of the perturbed open-loop velocity (e.g., distance between a large square and the origin in **a**) divided by the speed of the intuitive open-loop velocity (e.g., distance between a circle and the origin in **a**). Note the non-uniformity of these speed scalings across target directions in each experiment (colors vary across each row).

(e) Distribution of per-target open-loop rotations for all chosen perturbed mappings (i.e., histogram of all values in **c**; 1 count per target per analyzed experiment).

(f) Distribution of per-target open-loop speed scalings for all chosen perturbed mappings (i.e., histogram of all values in **d**; 1 count per target per analyzed experiment).



Supplementary Figure 3

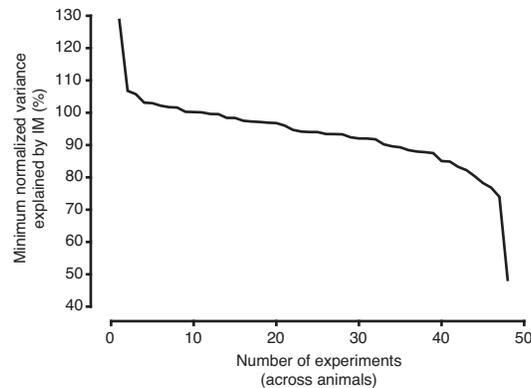
Neural activity remains consistent with the intrinsic manifold after learning.

Experiments were designed using a 10-dimensional FA model to describe the intrinsic manifold. These 10 dimensions were determined during calibration trials at the beginning of each experiment. Here we quantify the extent to which after-learning neural activity remained consistent with those 10 dimensions.

For each experiment, we computed the fraction of variance in the after-learning neural activity explained by the FA model:

$$\frac{\text{trace}(\mathbf{U}^T \mathbf{S} \mathbf{U})}{\text{trace}(\mathbf{S})}$$

where $\mathbf{U} \in \mathcal{R}^{q \times 10}$, from equation (16), describes the dimensions spanned by the FA model, and $\mathbf{S} \in \mathcal{S}_+^{q \times q}$ is the empirical covariance of the after-learning z-scored spike count vectors. To establish a baseline for interpreting these values, we also computed the fraction of variance explained by each FA model in the before-learning trials (these trials were not used to fit the FA model). We then defined the normalized variance explained (NVE) by the intrinsic manifold (vertical axis) to be the after-learning fraction of variance explained divided by the before-learning fraction of variance explained. An NVE value of 100% indicates that the same fraction of neural variability was explained by the fixed FA model after learning relative to before learning. Values above 100% indicate that neural activity preferentially migrated into the FA dimensions, and values below 100% indicate net migrations out of those dimensions. We found that normalized variance explained by the fixed FA model was at least 90% in 33/48 of experiments (across animals). In 47/48 experiments, NVE was at least 74.0%.

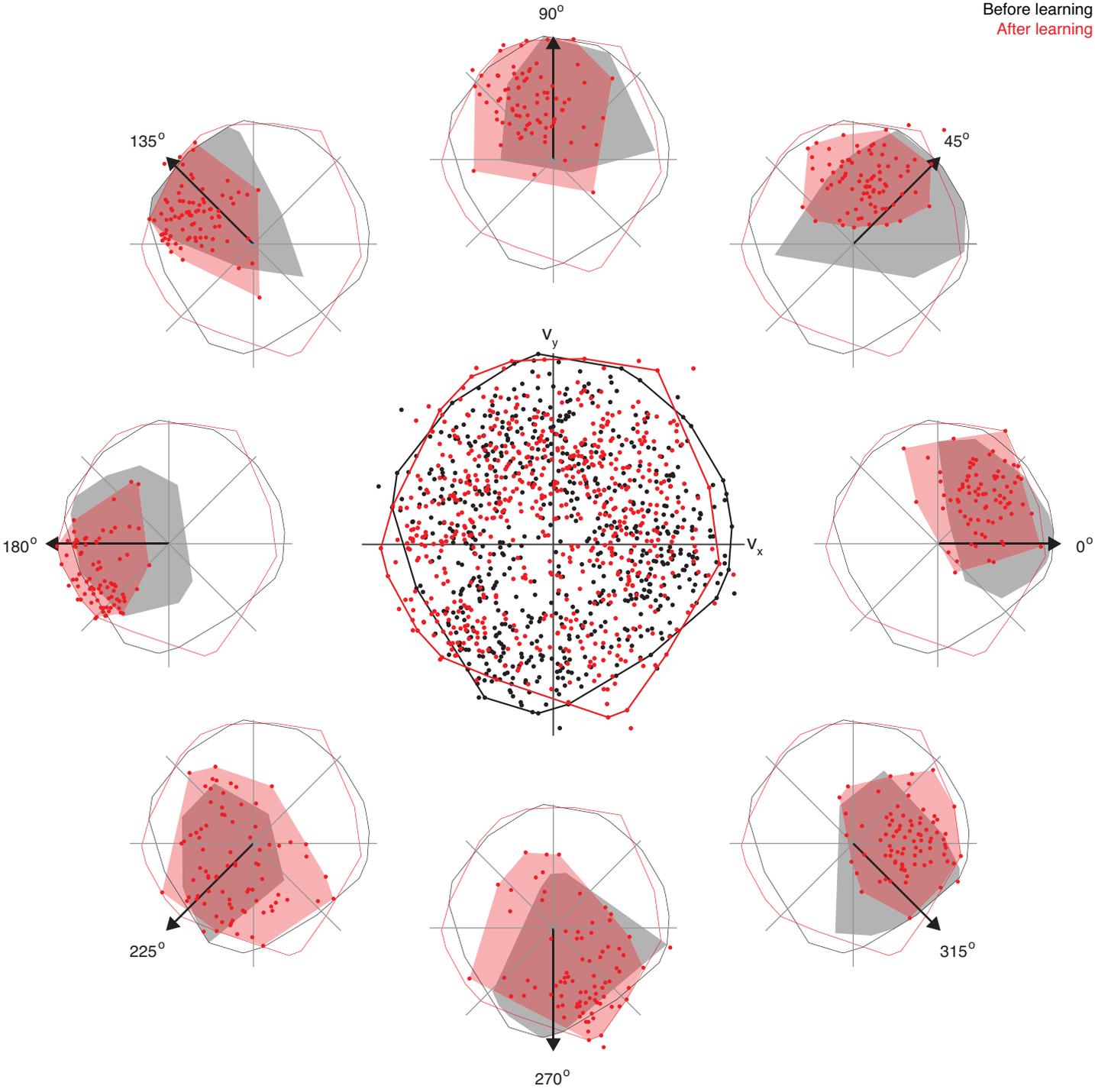


Supplementary Figure 4

Population activity patterns viewed through the intuitive BCI mapping.

Activity patterns are the same as those in **Figure 3**, and are presented using the same formatting conventions (black: before learning; red: after learning). Here 10-D activity patterns shown as their 2-D outputs through the intuitive BCI mapping, whereas in **Figure 3** patterns were shown as their outputs through the perturbed BCI mapping.

Before learning
After learning



Supplementary Figure 5

Predicted population activity patterns.

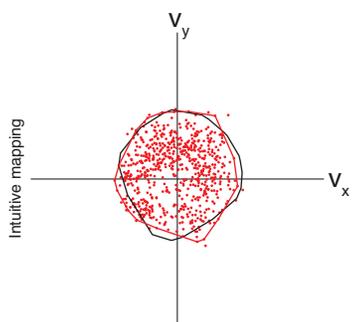
(a) Empirically observed population activity patterns, visualized as their outputs through the intuitive (top) and perturbed (bottom) BCI mappings (for reference; replicated from **Figure 3** and **Supplementary Figure 4**). Black and red outlines encapsulate 98% of before- and after-learning patterns, respectively.

(b) Realignment-predicted activity patterns (red). Black outlines are replicated from **a** for reference. Formatting conventions match that of **Figure 3** and **Supplementary Figure 4**. When viewed through the perturbed BCI mapping (bottom panel), realignment-predicted repertoire expansion is indicated by the many activity patterns (red points) that lie outside the empirical before-learning repertoire (black outline). Predicted repertoire change is quantified in **Figure 4**. This repertoire expansion is an explicit goal of realignment, as it enables high-speed movements. Notably, realignment predicts higher movement speeds (indicated by the distance of each red point from the origin) than were empirically observed under any context (black and red in **a**, top and bottom). A quantification of predicted movement speeds is given in **Supplementary Figure 7a**. The spread of activity patterns is quantified in **Figure 5**.

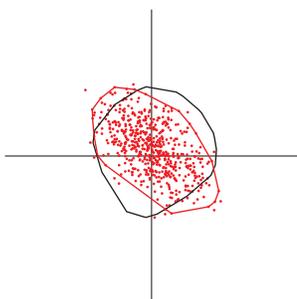
(c) Rescaling-predicted activity patterns. Here, repertoire expansion is indicated when viewing through the intuitive BCI mapping (top panel), which shows many predicted activity patterns (red points) that lie outside of the domain of the before-learning repertoire (black outline). Because the dimensions spanned by the intuitive mapping contribute more to movement under the intuitive mapping than they do under the perturbed mapping, rescaling predicts that those dimensions increase their variance to restore the magnitude of contribution they had toward movement prior to the perturbation. This expansion of variances is seen as the increased spread of Rescale-predicted patterns relative to the before-learning data (**Fig. 5c**).

(d) Reassociation-predicted activity patterns. Reassociation-predicted activity patterns match the data when visualizing through either BCI mapping (top and bottom panels). Repertoire preservation is indicated by the overwhelming majority of predicted activity patterns (red points) that lie within the domain of the before-learning repertoire (black outlines).

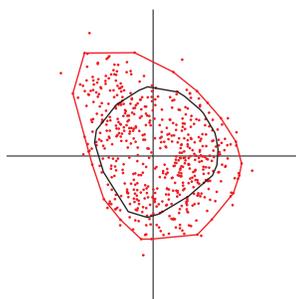
a Data, before learning
Data, after learning



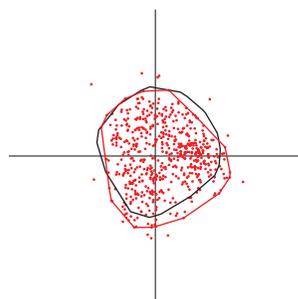
b Realignment



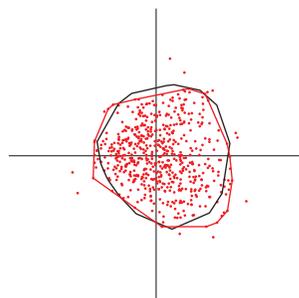
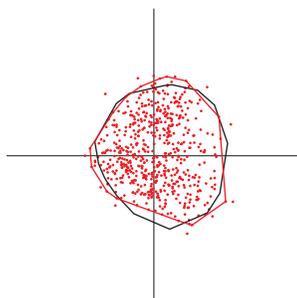
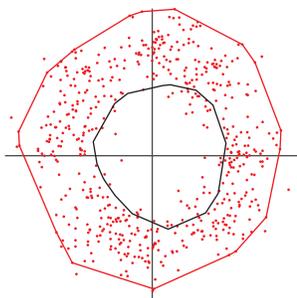
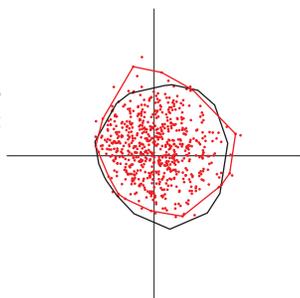
c Rescaling



d Reassociation



Perturbed mapping



Supplementary Figure 6

Characterization of the repertoire change metric.

To quantify changes in the neural repertoire, we devised a metric based on distances between each after-learning population activity pattern and its nearest neighbors among the before-learning activity patterns (see Online Methods; **Fig. 4b**, **Fig. 8c**, **Supplementary Fig. 9**, and **Supplementary Fig. 10**). Here we characterize this repertoire change metric by applying it to data generated from simple and systematic changes in the distribution of population activity patterns. For this characterization, we generated before-learning activity patterns from a standard multivariate Gaussian distribution (i.e., mean = 0 and standard deviation = 1 for each dimension). We then generated after-learning activity patterns from distributions that were scaled (**a-c**) or shifted (**d-f**) versions of the before-learning distribution. To build intuition, we first generated from 2-D distributions (**a-b** and **d-e**). Then, to match the data analyzed throughout this work, we generated from 10-D distributions (**c** and **f**). In both cases, we matched the number of generated activity patterns to the data presented in **Figure 4b** (i.e., we simulated 48 independent experiments, and for each simulated experiment, we generated the same number of activity patterns as was analyzed from the animal's overall neural repertoire in the corresponding actual experiment).

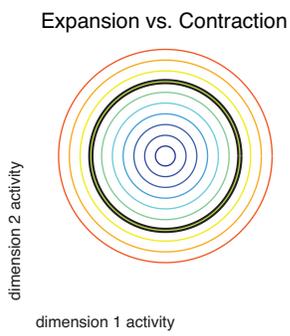
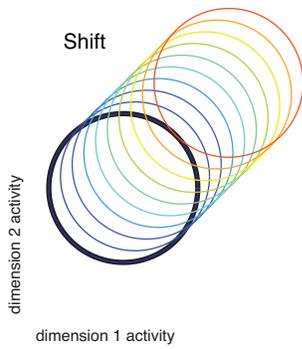
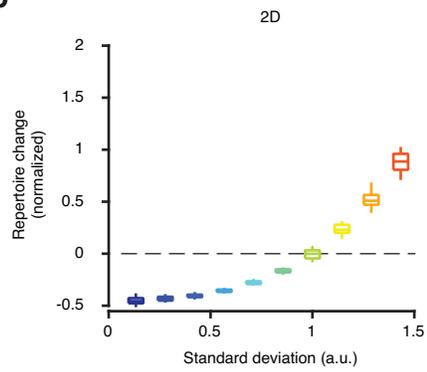
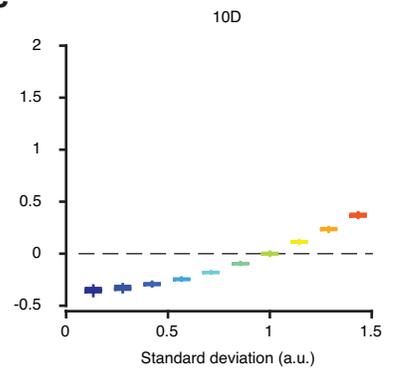
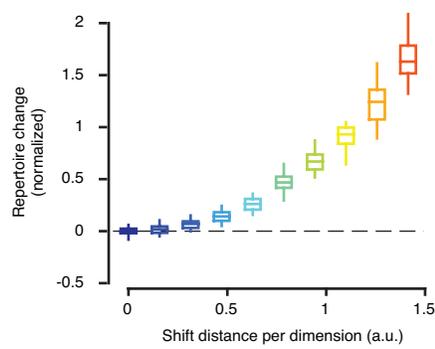
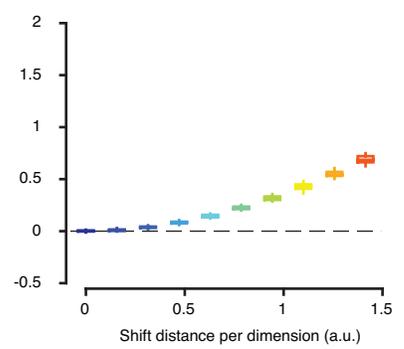
(a) Systematic scalings of the after-learning activity distributions. We illustrate the before-learning (black) and after-learning (blue through red) distributions as 2-D covariance ellipses (radii = 1 standard deviation).

(b) Repertoire change measured in the 2-D after-learning distributions from **a** relative to the before-learning distribution. Contractions (i.e., standard deviation < 1) are indicated by negative repertoire change values. Expansions (i.e., standard deviation > 1) are indicated by positive repertoire change values. Repertoire change values were first averaged within each simulated experiment. On each box, the central line indicates the median of these averaged values, the bottom and top edges indicate the 25th and 75th percentiles, respectively, and the whiskers extend to the 5th to the 95th percentiles ($n = 48$).

(c) Repertoire change measured in 10-D versions of the distributions from **a**. Measurements follow the same trends as in **b**, with negative values indicating contractions and positive values indicating expansions. The analyses of recorded population activity presented throughout this paper were performed in a 10-D space (representing the intrinsic manifold). Boxes and whiskers follow the same conventions as in **b**.

(d-f) Same format as **a-c**, but for systematic shifts in the after-learning activity distributions. Each after-learning distribution had its mean shifted by the same distance along each dimension. All dimensions of all before- and after-learning distributions had standard deviation = 1. Shifts are indicated by positive repertoire change values.

Contractions, expansions, and shifts are three qualitative descriptions of changes that would manifest quantitatively in the repertoire change metric. We illustrate these three here and because they describe changes we observed visually either in the empirical data (**Fig. 3** and **Supplementary Fig. 4**) or in the hypotheses' predictions (**Supplementary Fig. 5**). We have labeled the vertical axes in **Figure 4b** and **Figure 8c** with these qualitative descriptors ("contract", "shift / expand") to convey these corresponding intuitions. However, this repertoire change metric would be sensitive to many other types of changes not explored here (e.g., scaling of variability only along particular dimensions, as predicted by rescaling and realignment). While values near zero are consistent with repertoire preservation (as measured in the data and as predicted by reassociation), one can construct cases in which the neural repertoire changes in a manner that is not detected by this metric (e.g., a combination of a contraction and a shift could result in a value near zero; if activity is organized about an annulus, contractions could result in positive rather than negative values). To protect against potential non-identifiability, we were careful to visualize the empirical and predicted data (**Fig. 3**, **Supplementary Fig. 4**, and **Supplementary Fig. 5**), and we applied a suite of metrics tailored to identify a diversity of features in the data (**Figs. 4-8**).

a**d****b****c****e****f**

Supplementary Figure 7

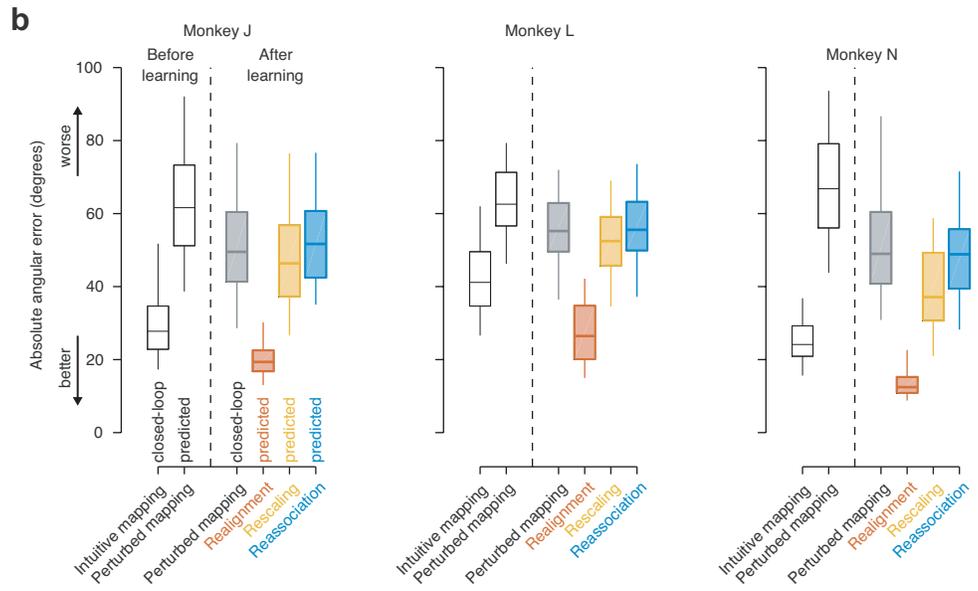
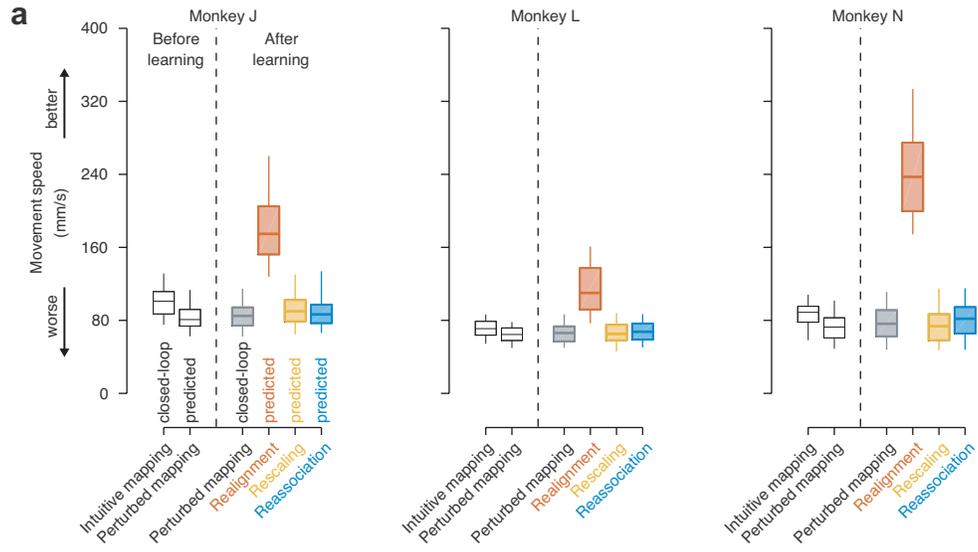
Contributions of movement speed and directional accuracy to acquisition time improvements.

The predicted acquisition times presented in **Figure 7** reflect both the speed and directional accuracy of movements. Here we breakdown those acquisition times in terms of these features. Format matches that of **Figure 7**.

(a) Timestep-by-timestep movement speeds did not increase appreciably during learning (Before learning: “Perturbed mapping” vs. After learning: “Perturbed mapping”), which is consistent with the predictions of rescaling and reassociation, but not realignment. Despite the qualitative agreement with the data for rescaling and reassociation, differences were significant for both comparisons ($p \leq 0.001$, two-sided paired Wilcoxon signed-rank test, $n = 384$: 48 experiments across animals \times 8 movement conditions). Differences were more substantial between realignment predictions and the data ($p < 10^{-10}$). On each box, the central line indicates the median, the bottom and top edges indicate the 25th and 75th percentiles of the data, respectively, and the whiskers extend to the 5th to the 95th percentiles of the data (monkey J: $n = 216$; monkey L: $n = 88$; monkey N: $n = 80$).

(b) Timestep-by-timestep directional accuracy (relative to cursor-to-target direction) improved during learning (Before learning: “Perturbed mapping” vs. After learning: “Perturbed mapping”). Reassociation-predicted angular errors were not significantly different from those in the data ($p = 0.59$, two-sided paired Wilcoxon signed-rank test, $n = 384$). Rescaling- and realignment-predicted angular errors were significantly different from those in the data ($p < 10^{-10}$). Boxes and whiskers follow the same conventions as in **a**.

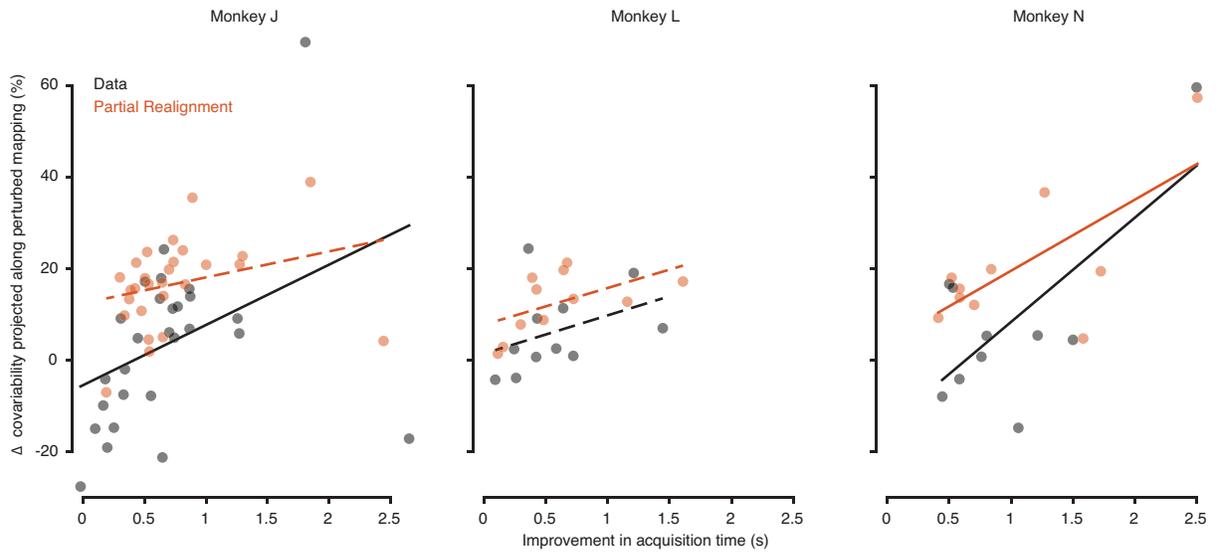
Taken together, empirical improvements in acquisition time reported in **Figure 7** (and consistent with reassociation and rescaling) are largely due to improvements in directional accuracy, rather than increases in timestep-by-timestep speeds.



Supplementary Figure 8

The neural activity shows subtle hints of realignment.

Although the covariability of the after-learning neural activity was consistent with reassociation (i.e., no substantial changes, see **Figs. 5-6**), we did find a subtle realignment-like effect. Behavioral improvements under realignment are due to expanding the covariability of the neural activity along the dimensions spanned by the perturbed BCI mapping. The data (black) show a weak, but consistently positive trend between increases in covariability along the perturbed BCI mapping and behavioral improvements after learning (F-test for nonzero slope; solid lines: $p \leq 0.05$; dashed lines: $p > 0.05$; monkey J: $p = 0.05$, $n = 27$ experiments; monkey L: $p = 0.25$, $n = 11$ experiments; monkey N: $p = 0.02$, $n = 10$ experiments). This trend matches that predicted by partial realignment (red; monkey J: $p = 0.15$; monkey L: $p = 0.09$; monkey N: $p = 0.03$). We show predictions for partial realignment rather than complete realignment so predicted behavioral improvements would approximately match those in the data (i.e., so values on the horizontal axes would be comparable between the data and the predictions).



Supplementary Figure 9

Learning strategy was not influenced by training history.

Here we address whether the animals' learning strategies might have been affected by accumulated experience controlling intuitive BCI mappings. By construction, intuitive mappings could be effectively controlled using within-repertoire neural activity patterns. An intriguing possibility is that accumulated experience with intuitive mappings (i.e., across many weeks or months) might reinforce the neural repertoire, perhaps making it increasingly difficult for animals to learn to produce out-of-repertoire activity patterns (such as those predicted by realignment or rescaling). Animals accrued experience controlling intuitive BCI mappings both prior to and during the experiments analyzed in this study (see Online Methods). Here we searched for effects of this accrued experience (horizontal axes) on features of population activity (vertical axes) that we have shown can disambiguate between different neural strategies of learning (**Figs. 4-6**).

(a) If animals became increasingly reliant on a fixed neural repertoire with accumulated experience controlling intuitive BCI mappings, we would expect to see larger positive values of repertoire change (as described in **Fig. 4**) in earlier experiments (indicating more flexibility in the neural repertoire; consistent with realignment or rescaling) and values closer to zero in later experiments (indicating more reliance on a fixed repertoire; consistent with reassociation). There was a weak trend in this direction, but this effect was not statistically significant (F-test for nonzero slope; monkey J: $p = 0.38$, $n = 27$ experiments; monkey L: $p = 0.58$, $n = 11$ experiments; monkey N: $p = 0.53$, $n = 10$ experiments).

(b) If the animals' ability to learn via realignment diminished with increased experience with intuitive BCI mappings, we would expect to see increases in population covariability along the perturbed mapping (i.e., positive values along the vertical axis in **Fig. 5c**) that decrease in magnitude throughout the course of the experiments. The data did not reveal any significant trend (monkey J: $p = 0.96$; monkey L: $p = 0.91$; monkey N: $p = 0.87$).

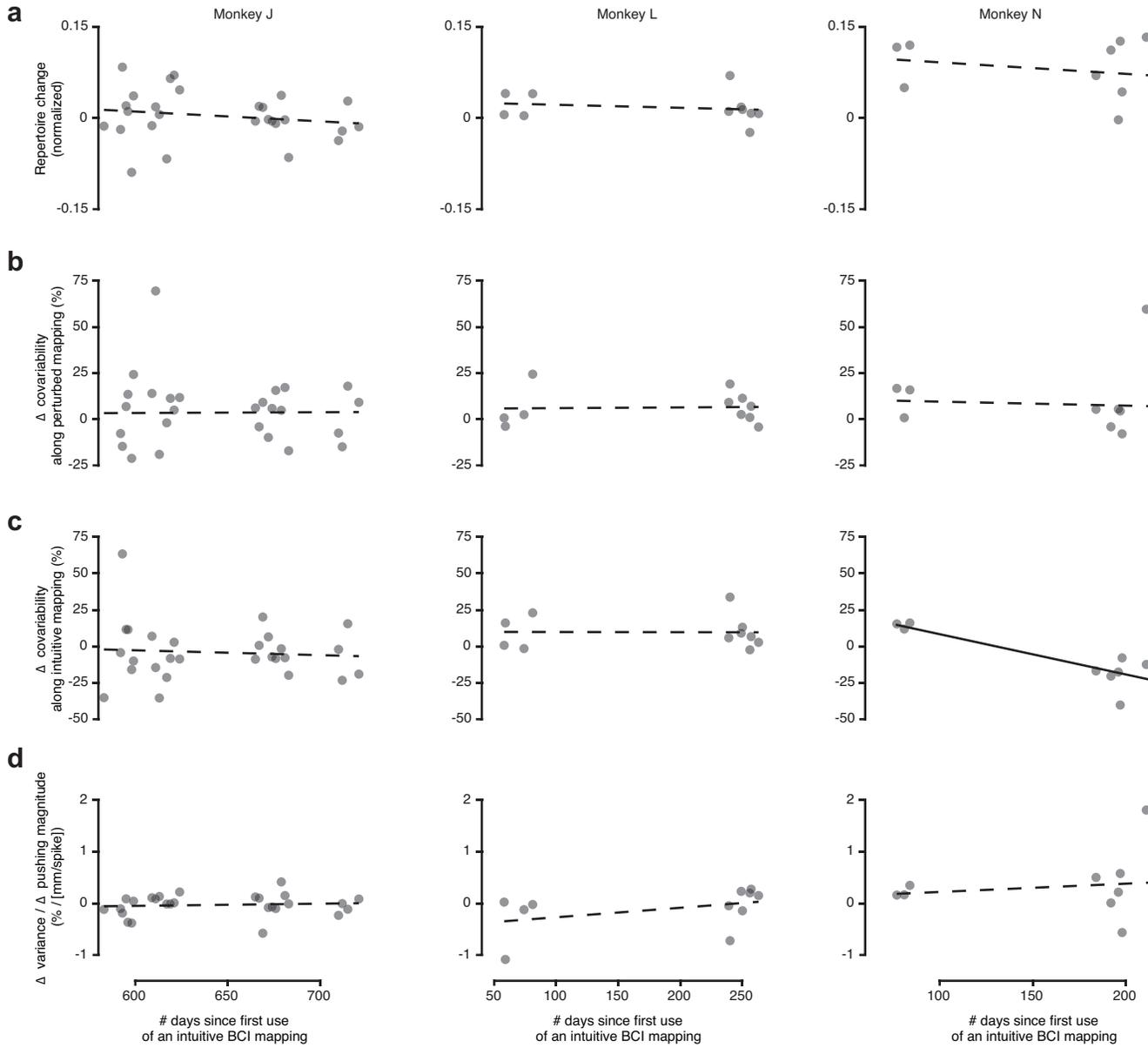
(c) If the animals' ability to learn via rescaling diminished with increased experience with intuitive BCI mappings, we would expect to see increases in population covariability along the intuitive mapping (i.e., positive values along the horizontal axis in **Fig. 5c**) that decrease in magnitude throughout the course of the experiments. The data did not consistently reveal such a trend (although see monkey N, solid line; monkey J: $p = 0.70$; monkey L: $p = 0.97$; monkey N: $p = 1.04 \times 10^{-3}$).

(d) Realignment predicts an increasing relationship between perturbation-induced dimension-by-dimension changes in pushing magnitude and changes in population variability along those dimensions (i.e., positive slopes as described in **Fig. 6f,g**). Rescaling predicts the opposite trend (i.e., negative slopes in **Fig. 6f,g**). If the animals' ability to learn via these strategies diminished with increased experience with intuitive BCI mappings, we would expect to evidence of these relationships that becomes weaker throughout the course of the experiments (i.e., for realignment: positive values along the vertical axis that decrease with time; for rescaling: negative values that increase with time). The data did not reveal either of these trends (monkey J: $p = 0.64$; monkey L: $p = 0.22$; monkey N: $p = 0.67$).

The lack of the evidence for the trends described above indicates that the animals' learning strategies were largely consistent with reassociation early, late, and throughout the entire course of the experiments analyzed in this study. Further, note that monkeys L and N had substantially less prior experience with intuitive BCI mappings than did monkey J (see starting values on horizontal axes), yet monkeys L and N do not show more evidence of realignment or rescaling. Taken together, these consistent observations suggest that the animals' learning strategies were not influenced by accumulated experience controlling intuitive BCI mappings during this study.

One possible interpretation of this finding is that animals relied on a fixed neural repertoire that was consolidated prior to these experiments. We believe that such a setting is a highly relevant regime for studying learning (e.g., motor learning) because the appropriate neural repertoire for a given class of behaviors (e.g., arm movements) is likely consolidated after years of use (e.g., in adult subjects). This is analogous to an animal having experience with an intuitive BCI mapping prior to the experiments analyzed in this work.

We did not combine data across animals in these analyses because doing so would result in distinct clustering of values along the horizontal axes. Such clustering effectively reduces the degrees of freedom in the data, and fitting a line to such clustered data could lead to spurious detection of trends intended to generalize across animals but based only on data from a single animal in the regimes about each cluster of values.



Supplementary Figure 10

Learning strategy was not influenced by pressure to change the neural repertoire.

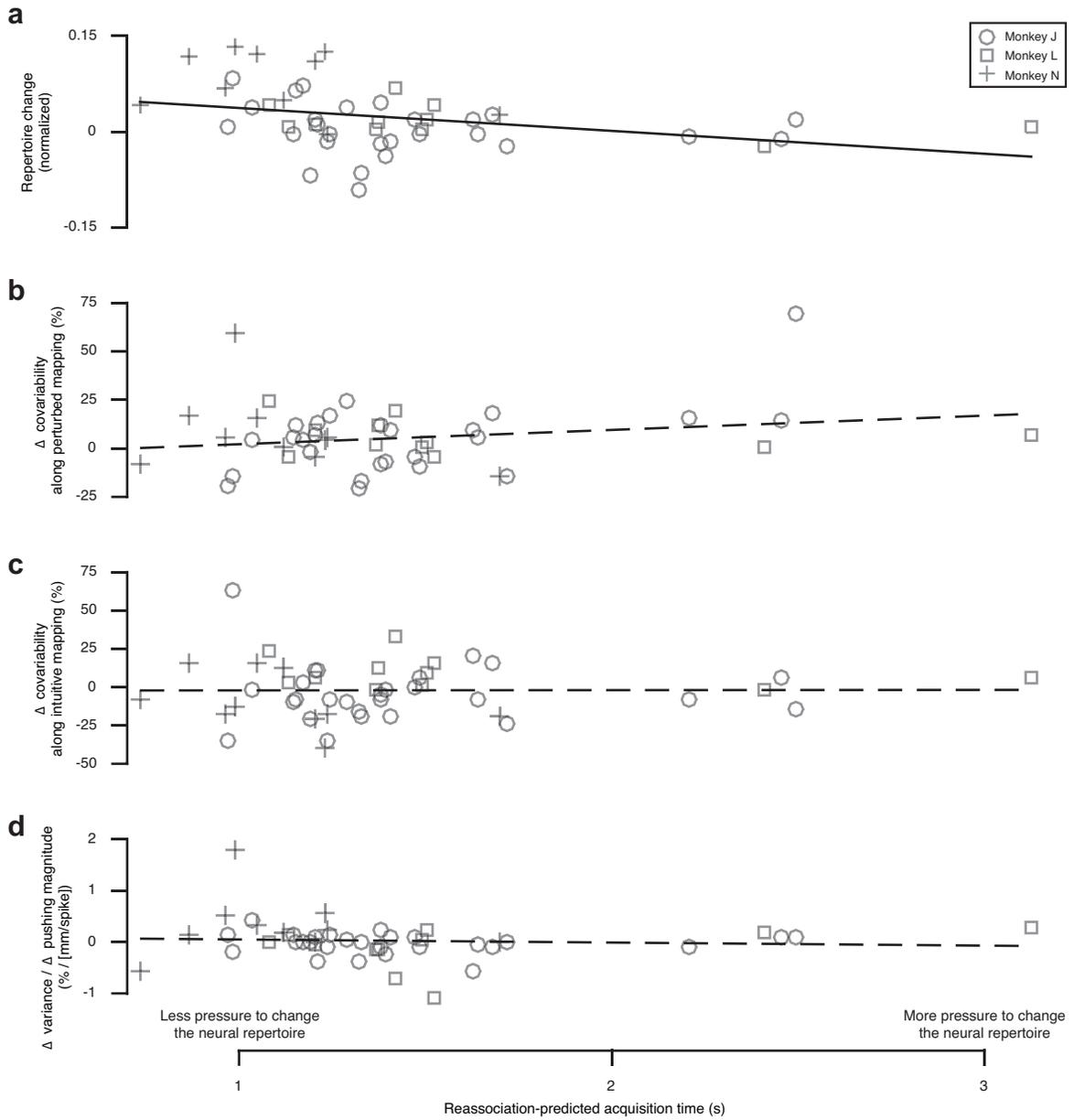
Because perturbations were chosen randomly, some perturbations could be better learned by reassociation than others. For perturbations that can be well-learned by reassociation (shorter reassociation-predicted acquisition times), there is relatively little pressure to change the neural repertoire. By contrast, perturbations that cannot be as well-learned by reassociation (longer reassociation-predicted acquisition times) exert more pressure to change the neural repertoire. Because realignment offers the largest potential for behavioral improvement (**Fig. 7**), here we ask whether animals increasingly introduced a realignment-like strategy as a function of this learning pressure. Although the potential for behavioral improvements due to rescaling is substantially lower than that for realignment (**Fig. 7**), for completeness we address both realignment and rescaling here. Overall, we did not find evidence that such perturbation-specific learning pressures influenced the animals' learning strategy.

(a) Had animals differentially engaged learning by realignment or rescaling, we would expect increases in learning pressure to be accompanied by increasing amounts of repertoire change (as described in **Fig. 4**). The data do not show this trend, but rather show a subtle but significant trend in the opposite direction ($p = 0.02$; F-test for nonzero slope, $n = 48$ experiments across animals).

(b-c) Had animals differentially engaged learning strategies, we would expect increases in learning pressure to be accompanied by larger increases in population covariability (as described in **Fig. 5c**) along the perturbed mapping for realignment (**b**) or along the intuitive mapping for rescaling (**c**), respectively. The data do not show significant trends in either of these metrics (perturbed BCI mapping: $p = 0.19$; intuitive BCI mapping: $p = 0.98$).

(d) Had animals differentially engaged learning by realignment, we would expect increases in learning pressure to be accompanied by larger increases in dimension-by-dimension population variability per unit of perturbation-induced change in pushing magnitude (slopes as described in **Fig. 6f,g**) as a function of learning pressure. Had animals differentially engaged learning by rescaling, we would expect to see the opposite trend: increasing learning pressure should yield larger decreases in variability per unit change in pushing magnitude. The data show no significant trend ($p = 0.64$).

Taken together, the lack of dose-dependent effects on these neural signatures of realignment or rescaling (**a-d**) suggests that the reason we do not find evidence in support of these learning strategies is not trivially due to insufficient behavioral pressure to change the neural repertoire. Rather, these analyses provide additional corroboration that constraints on the neural repertoire are not readily overcome on the timescale of these experiments, even when there is behavioral pressure to do so. In addition to parsing these neural signatures as a function of absolute reassociation-predicted acquisition times, we also looked at reassociation-predicted percent recovery of intuitive-level acquisition times, which yielded qualitatively similar results.



Supplementary Figure 11

Introducing Poisson variability does not violate hypotheses’ physiological plausibility.

Here we show that the physiological plausibility that we ensured in the 10-D space of the intrinsic manifold is also conferred to corresponding high-dimensional spike count vectors containing Poisson variability. Physiological plausibility is a critical feature of the hypotheses we have considered in this work. Specifically, we ensured that all hypotheses predicted neural activity patterns that correspond to spike counts for individual neural units that do not exceed the range of empirically observed spike counts for those same neural units in the before-learning trials. To ensure that none of the hypotheses relied on outside-manifold activity (which we have shown to be difficult to learn on the timescale of these experiments²²), we formulated all hypotheses’ predictions in the 10-D space of the intrinsic manifold. Because these 10-D predictions were based on the factors extracted by FA, which represent variance that is shared across units, these predictions do not include variability that is independent to each individual neural unit.

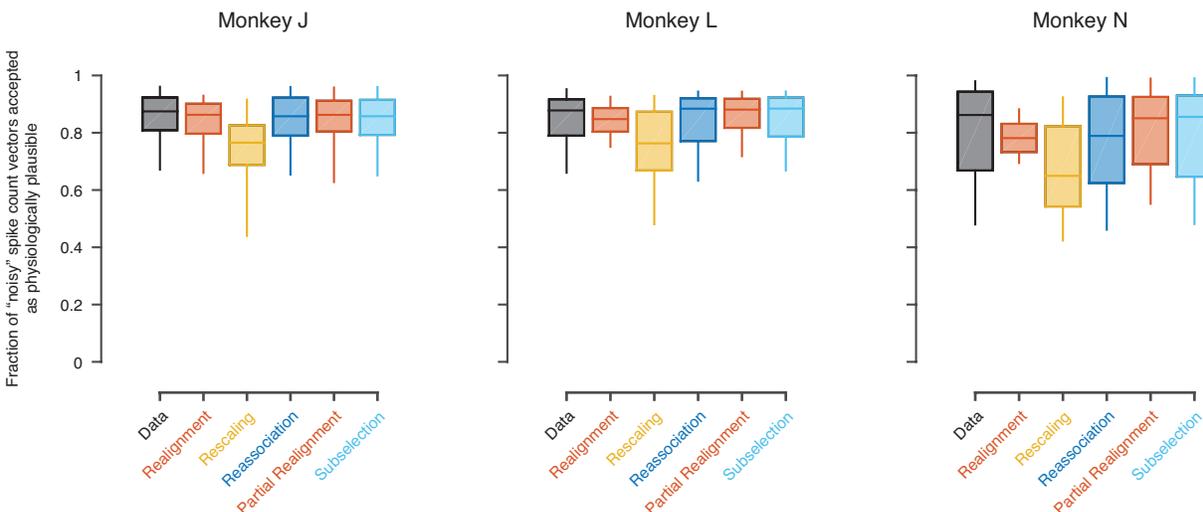
A potential concern is that, even if a neural activity pattern is deemed physiologically plausible in the 10-D space of the intrinsic manifold, it might not be truly physiologically plausible if this “denoised” 10-D representation cannot be obtained by a sufficient number of “noisy” instantiations of that pattern in the high-dimensional space of spike count vectors. Here we performed a post-hoc analysis to show that, even when introducing independent noise across neural units, all hypotheses’ predictions remain physiologically plausible.

We used the following procedure to obtain “noisy” spike count vectors. For each activity pattern in the before-learning data and for each predicted activity pattern from each hypothesis, we computed an expected firing rate for each neural unit i :

$$r_{i,t} = \max(0, E[u_{i,t} | \mathbf{z}_t])$$

where $E[u_{i,t} | \mathbf{z}_t]$ is the expected spike count for unit i given the 10-D factors[†], and the \max operation ensures that all firing rates are non-negative. We then sampled 1000 “noisy” spike count vectors using Poisson distributions with these rate parameters $r_{i,t}$. Finally, we asked what fraction of these sampled spike count vectors were physiologically plausible at the level of individual neural unit spike counts. We found that the vast majority of these samples were physiologically plausible, and we did not see qualitative differences across hypotheses or between the hypotheses and the before-learning data.

Fractions were first computed for each movement condition in each experiment. On each box, the central line indicates the median of these fractions, the bottom and top edges indicate the 25th and 75th percentiles of the data, respectively, and the whiskers extend to the 5th to the 95th percentiles of the data (monkey J: $n = 27$ experiments \times 8 movement conditions = 216; monkey L: $n = 88$; monkey N: $n = 80$).



[†]This expectation corresponds to i -th element of the mean of the distribution in equation (12), but adjusted to take into account the orthonormalization of the factors (equation (15)) and the z-scoring of spike counts.

SUPPLEMENTARY MATH NOTE

Identifying the intrinsic manifold

We used factor analysis (FA) to identify the intrinsic manifold and to summarize the recorded neural population activity at each moment in terms of a set of low-dimensional factors, \mathbf{z}_t . The probabilistic model for FA is defined as follows:

$$\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (11)$$

$$\mathbf{u}_t | \mathbf{z}_t \sim \mathcal{N}(\Lambda \mathbf{z}_t + \bar{\mathbf{u}}, \Psi) \quad (12)$$

where equation (11) describes the prior distribution of the factors, $\mathbf{z}_t \in \mathcal{R}^{10}$, and equation (12) describes the relationship between the factors, \mathbf{z}_t , and the vector of z-scored spike counts, $\mathbf{u}_t \in \mathcal{R}^q$, across q neural units (z-scoring was performed separately for each neural unit). The intrinsic manifold is defined by the column space of Λ . The FA parameters, $\Lambda \in \mathcal{R}^{q \times 10}$, $\bar{\mathbf{u}} \in \mathcal{R}^q$, and diagonal $\Psi \in \mathcal{S}_+^q$ (a $q \times q$ covariance matrix with all off-diagonal entries set to 0), were fit using the Expectation Maximization algorithm using the spike count vectors recorded during the calibration trials at the beginning of each experiment. Because spike count vectors were z-scored, each element of $\bar{\mathbf{u}}$ was 0.

After model fitting, the factors are computed as:

$$\hat{\mathbf{z}}_t = \mathbb{E}[\mathbf{z}_t | \mathbf{u}_t] = \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \mathbf{u}_t \quad (13)$$

To ensure that the perturbations (see *Perturbed BCI mappings*) would not require animals to produce factor activity beyond physiological ranges during learning, cursor movements during the closed-loop experiments were determined based on z-scored versions of these factors:

$$\hat{\mathbf{z}}_t^{\text{z-score}} = \mathbf{S}_z^{-1} \hat{\mathbf{z}}_t \quad (14)$$

where $\mathbf{S}_z \in \mathcal{S}_+^{10}$ is a diagonal matrix containing the standard deviations of the factors $\hat{\mathbf{z}}_t$ as measured from the calibration trials. Because FA was performed on z-scored spike counts from the calibration trials (\mathbf{u}_t in equations (12) and (13)), the calibration $\hat{\mathbf{z}}_t$ are already zero mean and thus an additional offset is not required in equation (14).

All offline analyses of *population activity patterns* were performed on orthonormalized factors:

$$\hat{\mathbf{z}}_t^{\text{orth}} = \mathbf{D}\mathbf{V}^T \hat{\mathbf{z}}_t \quad (15)$$

where premultiplication by $\mathbf{D}\mathbf{V}^T$ orthonormalizes the dimensions of the $\hat{\mathbf{z}}_t$, preserving the space spanned by the population activity patterns, but reorganizing the representation such that dimension 1 of $\hat{\mathbf{z}}_t^{\text{orth}}$ explains the most shared variance across the population, dimension 2 is orthogonal to dimension 1 and explains the next most shared variance, and so on³⁴. These parameters are obtained using the singular value decomposition:

$$\Lambda = \mathbf{U}\mathbf{D}\mathbf{V}^T \quad (16)$$

where the columns of $\mathbf{U} \in \mathcal{R}^{q \times 10}$ are orthonormal, $\mathbf{D} \in \mathcal{S}_+^{10}$ is a diagonal matrix containing the singular values of Λ , sorted from largest to smallest, and the columns of $\mathbf{V} \in \mathcal{R}^{10 \times 10}$ are orthonormal. Because all analyses were performed on these orthonormalized factors, we refer to $\hat{\mathbf{z}}_t^{\text{orth}}$ as \mathbf{z}_t in the main text and Online Methods for notational simplicity.

Intuitive BCI mappings

Intuitive BCI mappings translated the population activity, represented as 10-D factors into 2-D cursor velocities using a Kalman filter^{47,48}. The Kalman filter combines knowledge from a state model about how intended cursor velocity tends to evolve over time with knowledge from an observation model about how observed population activity relates to intended cursor velocity. When both of these models are linear-Gaussian, the Kalman filter provides minimum-mean-squared-error estimates of intended cursor velocity given the population activity recorded up to the current time t . The linear-Gaussian state model is written as:

$$\mathbf{v}_t \mid \mathbf{v}_{t-1} \sim \mathcal{N}(\mathbf{H}\mathbf{v}_{t-1} + \mathbf{b}, \mathbf{Q}) \quad (17)$$

The linear-Gaussian observation model is written as:

$$\mathbf{z}_t \mid \mathbf{v}_t \sim \mathcal{N}(\mathbf{C}\mathbf{v}_t + \mathbf{d}, \mathbf{R}) \quad (18)$$

We defined the state model in equation (17) to be a random walk model by setting $\mathbf{H} = \mathbf{I}_{2 \times 2}$ and $\mathbf{b} = \mathbf{0}$. The matrix $\mathbf{Q} \in \mathcal{S}_+^2$ controls the smoothness of cursor velocities over time, and fitting \mathbf{Q} to data does not always result in the BCI mapping that provides the animal with the highest performance closed-loop control. Rather than fitting to data, we tested various settings of \mathbf{Q} during closed-loop control with monkey J and found the best performance was achieved with a setting of $\mathbf{Q} = 2 \text{ (m}^2/\text{s}^2) \times \mathbf{I}_{2 \times 2}$. We used this setting of \mathbf{Q} for all experiments in all animals. We fit the remaining model parameters $\mathbf{C} \in \mathcal{R}^{10 \times 2}$, $\mathbf{d} \in \mathcal{R}^{10}$, and $\mathbf{R} \in \mathcal{S}_+^{10}$ with maximum likelihood using the factors ($\hat{\mathbf{z}}_t^{\text{z-score}}$ from equation (14)) and intended cursor velocities from the calibration trials at the start of each experiment (see *Calibrating the intuitive BCI mappings*).

The intuitive BCI mapping was the Kalman filter estimate of the animal's intended cursor velocity:

$$\hat{\mathbf{v}}_t = \mathbf{A}\hat{\mathbf{v}}_{t-1} + \mathbf{B}\hat{\mathbf{z}}_t^{\text{z-score}} + \mathbf{c} \quad (19)$$

$$\mathbf{A} = \mathbf{H} - \mathbf{KCH} \quad \mathbf{B} = \mathbf{K} \quad \mathbf{c} = -\mathbf{Kd} \quad (20)$$

where \mathbf{K} is the steady-state Kalman gain. Here we have written equation (19) in terms of the z-scored factors, $\hat{\mathbf{z}}_t^{\text{z-score}}$ from equation (14). Equivalently, we can write the intuitive BCI mapping in terms of the orthonormalized factors, $\hat{\mathbf{z}}_t^{\text{orth}}$ from equation (15):

$$\hat{\mathbf{v}}_t = \mathbf{A}\hat{\mathbf{v}}_{t-1} + \mathbf{B}^{\text{orth}}\hat{\mathbf{z}}_t^{\text{orth}} + \mathbf{c} \quad (21)$$

where \mathbf{A} and \mathbf{c} are unchanged from equation (20), and $\mathbf{B}^{\text{orth}} = \mathbf{KS}_z^{-1}\mathbf{VD}^{-1}$ follows from equations (14) and (15). Because these two representations of the BCI mapping are mathematically equivalent, in the main text we expressed the BCI mapping using a generic form (equation (1)), omitting the $\hat{\cdot}$, $^{\text{z-score}}$, and $^{\text{orth}}$ notations for clarity of presentation.

Perturbed BCI mappings

Perturbed BCI mappings were computed by permuting the factors (i.e., the elements of $\hat{\mathbf{z}}_t^{\text{z-score}}$) before they are passed into the Kalman filter. The perturbed BCI mapping can be written as:

$$\hat{\mathbf{v}}_t = \mathbf{A}\hat{\mathbf{v}}_{t-1} + \mathbf{B}^{\text{pert}}\hat{\mathbf{z}}_t^{\text{z-score}} + \mathbf{c} \quad \mathbf{B}^{\text{pert}} = \mathbf{K}\boldsymbol{\eta} \quad (22)$$

where \mathbf{A} , \mathbf{c} , and \mathbf{K} are unchanged from the intuitive BCI mapping in equation (19), and $\boldsymbol{\eta} \in \mathcal{R}^{10 \times 10}$ is a permutation matrix where exactly one entry in each row and each column is 1 and all other entries are 0. The effect of multiplying by $\boldsymbol{\eta}$ is to shuffle the elements of the 10-D factors, $\hat{\mathbf{z}}^{\text{z-score}\dagger}$.

For each experiment, there were $10! = 3,628,800$ unique permutations (i.e., settings of $\boldsymbol{\eta}$) from which we selected one as the perturbed BCI mapping. With each experiment, our aim was to select a candidate perturbation that would be difficult enough that substantial learning would be required to restore proficient control, but not so difficult as to discourage the animal from participating in the task (**Supplementary Figure 2**). To inform this selection, we predicted the cursor velocities that the animal would produce under each of the $10!$ candidate perturbed BCI mappings in the absence of visual feedback and before any learning had taken place. We term these predictions *open-loop cursor velocities*. To form these predictions, we first computed average per-target population activity patterns (based on equation (14)), $\bar{\mathbf{z}}_{\Theta}$, from the intuitive trials, Θ (monkey J: using 300ms to 1,300ms after the start of the first 200 intuitive trials; monkeys L and N: using 300ms to 1,100ms after the start of the first 150 intuitive trials). Similar to equation (4), we defined the open-loop cursor velocities as

$$\mathbf{v}_{\Theta}^{\text{OL}} = \mathbf{B}^{\text{pert}} \bar{\mathbf{z}}_{\Theta} + \mathbf{c} \quad (23)$$

where \mathbf{B}^{pert} and \mathbf{c} are the parameters of one of the $10!$ candidate perturbed mappings (from equation (22)). For each candidate perturbed mapping, we compared the open-loop velocities on a per-target basis to the set of open-loop velocities computed using the intuitive BCI mapping (i.e., replacing \mathbf{B}^{pert} in equation (23) with \mathbf{B} from equation (2)). Specifically, for each target, Θ , we measured the angle, ω_{Θ} , between the perturbed and intuitive $\mathbf{v}_{\Theta}^{\text{OL}}$. We also computed the speed of each open-loop velocity, $s_{\Theta} = \|\mathbf{v}_{\Theta}^{\text{OL}}\|_2$. These angles, ω_{Θ} , and speeds, s_{Θ} , serve to approximate the per-target velocity rotations and speed scalings, respectively, as described in the previous paragraph. We used these angles and speeds to eliminate candidate perturbed mappings deemed too difficult or not difficult enough. A candidate perturbed mapping was eliminated unless it satisfied all four of the following requirements:

1. All ω_{Θ} must be greater than an experiment-specific minimum (monkey J: $21.3^{\circ} \pm 8.8^{\circ}$; monkey L: $18.6^{\circ} \pm 5.1^{\circ}$; monkey N: $37.6^{\circ} \pm 6.0^{\circ}$)
2. All ω_{Θ} must be less than an experiment-specific maximum (monkey J: $45.0^{\circ} \pm 12.2^{\circ}$; monkey L: $41.8^{\circ} \pm 9.4^{\circ}$; monkey N: $68.9^{\circ} \pm 7.4^{\circ}$)
3. All s_{Θ} must be greater than an experiment-specific minimum (monkey J: $32.4 \text{ mm/s} \pm 20.7 \text{ mm/s}$; monkey L: $45.9 \text{ mm/s} \pm 17.9 \text{ mm/s}$; monkey N: $21.0 \text{ mm/s} \pm 17.7 \text{ mm/s}$)
4. All s_{Θ} must be less than an experiment-specific maximum (monkey J: $227.6 \text{ mm/s} \pm 174.7 \text{ mm/s}$; monkey L: $213.6 \text{ mm/s} \pm 74.1 \text{ mm/s}$; monkey N: $359.0 \text{ mm/s} \pm 3.0 \text{ mm/s}$).

We adjusted these thresholds from one experiment to the next (resulting in the thresholds' stated variabilities) to ensure that at least one candidate perturbed mapping satisfied all of the above requirements. In a typical experiment, approximately 50 candidate mappings satisfied all requirements, and we arbitrarily selected one of those as the perturbed BCI mapping for that experiment.

Predicting population activity after learning

Here we describe the procedures used to generate the population activity patterns predicted by realignment, partial-realignment, rescaling, reassociation, and subselection. For a given experiment, each hypothesis specifies a *movement-specific mean*, $\boldsymbol{\mu}_{\Theta} \in \mathcal{R}^{10}$, and a *movement-specific covariance*, $\boldsymbol{\Sigma}_{\Theta} \in \mathcal{S}_+^{10}$, of the population

[†]The effect of permuting the factors is mathematically equivalent to preserving the order of the factors, but permuting the columns of \mathbf{K} . Thus, the columns of \mathbf{B}^{pert} are a permutation of the columns of \mathbf{B} . However, all analyses of population activity patterns were performed in the space of the orthonormalized factors, $\hat{\mathbf{z}}_t^{\text{orth}}$, and thus this permutation is not visually apparent in the pushing vectors in **Figure 6a-d**. Rather, those pushing vectors are the columns of $\mathbf{B}^{\text{pert,orth}} = \mathbf{K}\boldsymbol{\eta}\mathbf{S}_z^{-1}\mathbf{V}\mathbf{D}^{-1}$, which in general is not a column-permuted version of \mathbf{B}^{orth} from equation (21).

activity patterns predicted in the after-learning movement-specific cloud for each of the 8 intended movement directions, Θ . We defined the intended movement direction at time t to be the cursor-to-target direction relative to the cursor position at time t . For realignment, rescaling, and reassociation, these parameters (μ_Θ and Σ_Θ) were determined based only on the before-learning activity patterns, the intuitive and perturbed BCI mappings, and the assumptions of realignment, rescaling, and reassociation. For partial-realignment and subselection, these parameters were also determined by the after-learning levels of behavioral performance, which was evaluated based on empirical closed-loop cursor movements. Predicted activity patterns for each movement-specific cloud were then sampled from a Gaussian distribution with mean μ_Θ and covariance Σ_Θ , ensuring that all sampled patterns correspond to physiologically plausible spike counts (see below). For each movement-specific cloud, we predicted N activity patterns to match the number of analyzed after-learning activity patterns in each empirical movement-specific cloud. For notational simplicity, all references to population activity patterns \mathbf{z}_t in this section refer to the $\hat{\mathbf{z}}_t^{\text{orth}}$ from equation (15).

Ensuring physiologically plausible spike counts

We defined each hypothesis to predict population activity patterns that correspond to physiologically plausible spike counts, which we defined independently for each neural unit as the range between the minimum and maximum observed spike count for that unit during the before-learning trials. This statistical control ensures that differences between the empirical and predicted activity patterns are not trivially due to a hypothesis predicting firing rates that are physiologically implausible. We implemented this control using the following procedure, which was common across all hypotheses' predictions.

1. Constrain each predicted movement-specific mean. FA provides a mapping between the high-dimensional space of spike count vectors \mathbf{u}_t and the 10-D population activity space of the factors \mathbf{z}_t . We constrained each movement-specific mean population activity pattern (which is defined in the 10-D population activity space of \mathbf{z}_t) to map to a high-dimensional mean spike count vector, $\Gamma_\Theta \in \mathcal{R}^q$, in which each neural unit's mean spike count lies within a physiologically plausible range. Here, since we are describing a *mean* spike count vector, we define this physiological range relative to empirical movement-specific *mean* spike counts. For each intended movement direction Θ , we found each unit's before-learning movement-specific mean spike count, $\bar{u}_{i,\Theta}^{\text{before}} \in \mathcal{R}$ for unit i . We then found each unit's minimum and maximum mean spike count across the 8 intended movement directions:

$$\bar{u}_{i,\min}^{\text{before}} = \min_{\Theta'} \bar{u}_{i,\Theta'}^{\text{before}} \quad \bar{u}_{i,\max}^{\text{before}} = \max_{\Theta'} \bar{u}_{i,\Theta'}^{\text{before}} \quad (24)$$

For all hypotheses, we used these empirical values to constrain each predicted movement-specific mean population activity pattern, μ_Θ in the 10-D space of \mathbf{z}_t , to map to a predicted mean spike count vector, Γ_Θ in the q -dimensional space of \mathbf{u}_t such that

$$\bar{\mathbf{u}}_{\min}^{\text{before}} \leq \Gamma_\Theta \leq \bar{\mathbf{u}}_{\max}^{\text{before}} \quad (25)$$

where $\bar{\mathbf{u}}_{\min} \in \mathcal{R}^q$ and $\bar{\mathbf{u}}_{\max} \in \mathcal{R}^q$ are vectors with $\bar{u}_{i,\min}^{\text{before}}$ and $\bar{u}_{i,\max}^{\text{before}}$ as their i -th entries, respectively. Predicted mean spike count vectors, $\Gamma_\Theta \in \mathcal{R}^q$, related to predicted mean population activity patterns, $\mu_\Theta \in \mathcal{R}^{10}$, based on FA (equation (13)) and its related orthonormalization procedure (equation (15)):

$$\mu_\Theta = \beta \Gamma_\Theta \quad \beta = \mathbf{D}\mathbf{V}^T \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \quad (26)$$

2. Match each predicted movement-specific covariance to the before-learning data. To ensure realistic levels of movement-specific neural variability, predicted movement-specific covariances were constrained in a

hypothesis-specific manner relative to the before-learning empirical movement-specific covariances, $\mathbf{S}_\Theta^{\text{before}}$:

$$\mathbf{S}_\Theta^{\text{before}} = \frac{1}{N} \sum_{t \in \mathcal{T}_\Theta^{\text{before}}} (\mathbf{z}_t - \bar{\mathbf{z}}_\Theta^{\text{before}})(\mathbf{z}_t - \bar{\mathbf{z}}_\Theta^{\text{before}})^\top \quad (27)$$

where $\mathcal{T}_\Theta^{\text{before}}$ is the set of timesteps that define the empirical before-learning movement- Θ cloud, and $\bar{\mathbf{z}}_\Theta^{\text{before}}$ is the movement- Θ mean population activity pattern:

$$\bar{\mathbf{z}}_\Theta^{\text{before}} = \frac{1}{N} \sum_{t \in \mathcal{T}_\Theta^{\text{before}}} \mathbf{z}_t \quad (28)$$

3. Sample predicted population activity patterns, rejecting those corresponding to non-physiological spike counts. For each hypothesis, we obtained N predicted activity patterns, $\tilde{\mathbf{z}}_\Theta^{(1)}, \dots, \tilde{\mathbf{z}}_\Theta^{(N)}$, for each intended movement direction, Θ , by sampling from the 10-D Gaussian specified by the predicted movement-specific mean and covariance:

$$\tilde{\mathbf{z}}_\Theta^{(n)} \sim \mathcal{N}(\boldsymbol{\mu}_\Theta, \boldsymbol{\Sigma}_\Theta) \quad (29)$$

Constraining predicted movement-specific means (Step 1) and covariances (Step 2) encourages sampling from this Gaussian to yield patterns with physiologically plausible spike counts. To guarantee that no predicted patterns involved non-physiological spike counts, we rejected samples corresponding to non-physiological spike counts. To do so, we asked whether each sampled population activity pattern, $\tilde{\mathbf{z}}_\Theta^{(n)}$, corresponded to a spike count for each neural unit (via FA) that was within the empirical before-learning dynamic range for that unit. This was accomplished by solving the following linear feasibility problem³⁸:

$$\begin{aligned} &\text{find} && \tilde{\mathbf{u}}_\Theta^{(n)} \\ &\text{subject to} && \tilde{\mathbf{z}}_\Theta^{(n)} = \boldsymbol{\beta} \tilde{\mathbf{u}}_\Theta^{(n)} \\ &&& \mathbf{u}_{\min}^{\text{before}} \leq \tilde{\mathbf{u}}_\Theta^{(n)} \leq \mathbf{u}_{\max}^{\text{before}} \end{aligned} \quad (30)$$

where

$$u_{i,\min}^{\text{before}} = \min_{t \in \mathcal{T}^{\text{before}}} u_{i,t} \quad u_{i,\max}^{\text{before}} = \max_{t \in \mathcal{T}^{\text{before}}} u_{i,t} \quad (31)$$

are the minimum and maximum empirically observed before-learning spikes counts for neural unit i across all intended movement directions, and $\mathbf{u}_{\min}^{\text{before}}$ and $\mathbf{u}_{\max}^{\text{before}}$ are q -dimensional vectors with $u_{i,\min}^{\text{before}}$ and $u_{i,\max}^{\text{before}}$ as their i -th entries, respectively. We solved this linear feasibility problem independently for each sampled population activity pattern, $\tilde{\mathbf{z}}_\Theta^{(n)}$. If a feasible solution does not exist for a particular sample, that indicates that the sampled activity pattern could not have been extracted (via FA) by any physiologically plausible spike count. In this case, we rejected that sample, and resampled using equation (29) until obtaining a sample with a feasible solution.

This process of rejecting non-physiological samples effectively changes the distribution of the sampled activity patterns (i.e., a movement-specific cloud might not correspond to the distribution specified in equation (29)). However, sampling from the optimized (i.e., pre-rejection) distribution specified by equation (29) is not the primary goal in our prediction procedure. Rather, the primary goal is to obtain a physiologically plausible sample distribution that is close to the optimized distribution. This rejection procedure was designed to meet this goal. Because of the constraints used in determining movement-specific means (Step 1), the large majority of activity patterns sampled using equation (29) were physiologically plausible and thus were not rejected (exact rates of rejection are provided in the prediction procedures described below).

Realignment

Under realignment, the repertoire of activity patterns changes in the manner that maximizes behavioral performance subject only to the constraints described above. Realignment-predicted movement-specific covariances were defined to match the mean movement-specific covariance from the before-learning trials. Finally, we generated realignment-predicted activity patterns from a Gaussian distribution with the specified movement-specific mean and covariance, rejecting patterns that did not correspond to high-dimensional spike counts within the empirically observed range of before-learning spike counts.

The following steps describe the details of the prediction procedure, which was done independently for each experimental session and for each intended movement direction, Θ :

1. Find the movement-specific mean that maximizes behavioral performance subject to physiological constraints. Obtaining the Realigned movement-specific mean, $\mu_{\Theta}^{\text{realign}}$, for intended movement direction Θ amounts to solving the following convex optimization problem³⁸:

$$\underset{\mu_{\Theta}^{\text{realign}}, \Gamma_{\Theta}^{\text{realign}}}{\text{maximize}} \quad P_{\text{pert}}(\mu_{\Theta}^{\text{realign}}, \Theta) \quad (32)$$

$$\text{subject to} \quad \bar{\mathbf{u}}_{\min}^{\text{before}} \leq \Gamma_{\Theta}^{\text{realign}} \leq \bar{\mathbf{u}}_{\max}^{\text{before}} \quad (33)$$

$$\mu_{\Theta}^{\text{realign}} = \beta \Gamma_{\Theta}^{\text{realign}} \quad (34)$$

The objective function in equation (32) is the cursor progress in direction Θ due to the activity pattern $\mu_{\Theta}^{\text{realign}}$ (equation (10)). The constraints in equations (33) and (34) ensure that the predicted mean activity pattern corresponds to a physiologically plausible mean spike count vector (following equations (25) and (26)).

2. Match the movement-specific covariance to the before-learning data. Predicted movement-specific covariances, $\Sigma_{\Theta}^{\text{realign}}$, for all intended movement directions Θ were defined to be the mean empirical movement-specific covariance across the 8 intended movement directions during the before-learning trials:

$$\Sigma_{\Theta}^{\text{realign}} = \frac{1}{8} \sum_{i=1}^8 \mathbf{S}_{\Theta_i}^{\text{before}} \quad (35)$$

where $\mathbf{S}_{\Theta_i}^{\text{before}}$ is empirical movement-specific covariance for movement Θ_i (equations (27) and (28)).

3. Generate realignment-predicted population activity patterns. We obtained N predicted activity patterns by sampling from the 10-D Gaussian specified by $\mu_{\Theta}^{\text{realign}}$ and $\Sigma_{\Theta}^{\text{realign}}$ (equation (29)), rejecting any samples that correspond to non-physiologically plausible spike counts (equation (30)); $1.84\% \pm 3.24\%$ of all samples were rejected for monkey J; $4.68\% \pm 6.84\%$ for monkey L; and $0.27\% \pm 0.31\%$ for monkey N).

Partial-realignment

Under partial realignment, the population activity patterns produced after learning for a given movement are intermediate between the before-learning cloud of activity patterns and realignment-predicted cloud of patterns (referred to as *complete* realignment in this section, to clearly disambiguate from *partial* realignment). For each intended movement direction, we determined how close the partial-realignment cloud is to the before-learning cloud versus the complete-realignment cloud by matching the empirical after-learning behavioral performance. To match behavior, the partial realignment clouds required, on average, a 15.2% migration from the before-learning to the complete-realignment clouds (monkey J: $14.9\% \pm 12.2\%$; monkey L: $19.4\% \pm 16.0\%$; monkey N: $11.3\% \pm 8.78\%$). The following steps describe the prediction procedure, which was done independently for each

experimental session and for each intended movement direction, Θ :

1. Determine the mean population activity pattern predicted by complete realignment by solving the optimization problem from equations (32)-(34). Find the individual-movement covariance under complete realignment using equation (35). This mean and covariance will be referred to as $\mu_{\Theta}^{\text{realign}}$ and $\Sigma_{\Theta}^{\text{realign}}$, respectively.

2. Quantify predicted behavioral performance under complete realignment by computing the cursor progress, $P^{\text{realign}} = P_{\text{pert}}(\mu_{\Theta}^{\text{realign}}, \Theta)$, using equation (10).

3. Determine the empirical mean and covariance of the before-learning activity patterns, $\bar{z}_{\Theta}^{\text{before}}$ and $S_{\Theta}^{\text{before}}$, respectively, using equations (27)-(28).

4. Quantify empirical behavioral performance of the before- and after-learning activity patterns through the perturbed BCI mapping by computing mean cursor progress values, $\bar{P}^{\text{before}} = P_{\text{pert}}(\bar{z}_{\Theta}^{\text{before}}, \Theta)$ and $\bar{P}^{\text{after}} = P_{\text{pert}}(\bar{z}_{\Theta}^{\text{after}}, \Theta)$, respectively, using equation (10).

5. Find the movement-specific mean activity pattern to match after-learning behavioral performance. We chose this activity pattern, $\mu_{\Theta}^{\text{partial}}$, to be that which best reproduces the empirical mean after-learning cursor progress out of all patterns that lie on the line between the empirical before-learning mean activity pattern, $\bar{z}_{\Theta}^{\text{before}}$, and the complete-realignment-predicted mean population activity pattern, $\mu_{\Theta}^{\text{realign}}$.

Because cursor progress is a linear function of activity pattern (equation (10)), interpolating between these mean activity patterns ($\bar{z}_{\Theta}^{\text{before}}$ and $\mu_{\Theta}^{\text{realign}}$) simultaneously interpolates between the corresponding progress values (\bar{P}^{before} and P^{realign}). To reproduce the empirical after-learning cursor progress from the before-learning and complete-realignment progress values, we solved the following equation for interpolation value, α :

$$\bar{P}^{\text{after}} = \alpha P^{\text{realign}} + (1 - \alpha) \bar{P}^{\text{before}} \quad (36)$$

We then applied the interpolation value α to obtain the partial-realignment mean population activity pattern:

$$\mu_{\Theta}^{\text{partial}} = \alpha \mu_{\Theta}^{\text{realign}} + (1 - \alpha) \bar{z}_{\Theta}^{\text{before}} \quad (37)$$

which results in predicted partial-realignment cursor progress, $P_{\Theta}^{\text{partial}} = \bar{P}^{\text{after}}$.

In the rare event that cursor progress decreased after learning for the given movement direction (i.e., solving equation (36) yields $\alpha < 0$), we set $\mu_{\Theta}^{\text{partial}} = \bar{z}_{\Theta}^{\text{before}}$ (i.e., $\alpha = 0$). It was never the case that empirical after learning cursor progress exceeded complete-realignment predicted cursor progress (i.e., $\alpha > 1$).

6. Find the movement-specific covariance of population activity pattern for partial realignment.

$$\Sigma_{\Theta}^{\text{partial}} = \alpha \Sigma_{\Theta}^{\text{realign}} + (1 - \alpha) S_{\Theta}^{\text{before}} \quad (38)$$

This corresponds to linearly interpolating between the empirical before-learning covariance and the covariance predicted by complete realignment, just as was done for the mean activity patterns in Step 5.

7. Generate partial-realignment-predicted population activity patterns. We obtained N predicted activity patterns by sampling from the 10-D Gaussian specified by $\mu_{\Theta}^{\text{partial}}$ and $\Sigma_{\Theta}^{\text{partial}}$ (equation (29)), rejecting any samples that correspond to non-physiologically plausible spike counts (equation (30); 1.08% \pm 1.89% of all samples were rejected for monkey J; 3.33% \pm 6.13% for monkey L; and 0.26% \pm 0.26% for monkey N).

Reassociation

Under reassociation, the animal relies on population activity patterns from the before-learning neural repertoire, but reassociates patterns with potentially different movements after learning. The following steps describe the prediction procedure, which was done independently for each experimental session and for each intended movement direction, Θ :

1. Define the search space for the movement-specific mean within the before-learning neural repertoire. Under reassociation, each predicted movement-specific cloud must be contained within the before-learning overall neural repertoire. This requires each movement-specific mean to be within an interior region of the overall neural repertoire, as opposed to on the boundary (if a movement-specific mean is on the boundary, then a large fraction of predictions will be outside of the overall repertoire). We chose this interior region to be the convex hull of the before-learning empirical movement-specific means (i.e., the smallest region that encloses all of those means and the line segments connecting them). To allow as much flexibility as possible for reassociating movement intents with parts of the overall neural repertoire, we defined this region with respect to movements re-discretized to cursor-to-target angles with a 3° angular spacing, $\Phi \in \{0^\circ, 3^\circ, \dots, 357^\circ\}$, rather than the 45° spacing in the set of intended movement directions $\Theta \in \{0^\circ, 45^\circ, \dots, 315^\circ\}$. To re-discretize for a given movement, Φ , we identified the N before-learning timesteps during which the cursor-to-target angle was closest to Φ . We denote the set of these timesteps as $\mathcal{T}_\Phi^{\text{before}}$. We then computed the mean, $\bar{\mathbf{z}}_\Phi^{\text{before}}$, and covariance, $\mathbf{S}_\Phi^{\text{before}}$, of the population activity patterns recorded at those timesteps, following equations (27) and (28). This re-discretization procedure exactly replicates the original discretization at the 45° spacing (see Online Methods).

2. Find the movement-specific mean within the before-learning repertoire that maximizes behavioral performance. We chose the mean population activity pattern, $\boldsymbol{\mu}_\Theta^{\text{reassociate}}$, to be the pattern from the search space defined in Step 1 that maximizes behavioral performance under the perturbed BCI mapping, as measured by cursor progress (equation (10)). This activity pattern was identified by solving the following convex optimization problem³⁸:

$$\underset{\boldsymbol{\mu}_\Theta, \boldsymbol{\Gamma}_\Theta, \mathbf{w}}{\text{maximize}} \quad P_{\text{pert}}(\boldsymbol{\mu}_\Theta, \Theta) \quad (39)$$

$$\text{subject to} \quad \boldsymbol{\mu}_\Theta = \sum_{\Phi \in \{0^\circ, 3^\circ, \dots, 357^\circ\}} w_\Phi \bar{\mathbf{z}}_\Phi^{\text{before}} \quad (40)$$

$$\mathbf{1}^T \mathbf{w} = 1, \quad \mathbf{w} \geq \mathbf{0} \quad (41)$$

$$\bar{\mathbf{u}}_{\min}^{\text{before}} \leq \boldsymbol{\Gamma}_\Theta \leq \bar{\mathbf{u}}_{\max}^{\text{before}} \quad (42)$$

$$\boldsymbol{\mu}_\Theta = \beta \boldsymbol{\Gamma}_\Theta \quad (43)$$

The objective function in equation (39) is the cursor progress in direction Θ due to the activity pattern $\boldsymbol{\mu}_\Theta$ (equation (10)). The constraints in equations (40) and (41) ensure that the predicted mean activity pattern is in the convex hull of the before-learning empirical movement-specific means. This is achieved by requiring that the optimal activity pattern $\boldsymbol{\mu}_\Theta^*$ is a weighted sum of the $\bar{\mathbf{z}}_\Phi^{\text{before}}$ (equation (40)), where the weights w_Φ are non-negative and sum to 1 (equation (41)), and \mathbf{w} is the weight vector with elements w_Φ . The constraints in equations (42) and (43) ensure that the predicted mean activity pattern corresponds to a physiologically plausible mean spike count vector (following equations (25) and (26)).

It turns out that solving this optimization problem (equations (39)-(43)) gives an optimal movement-specific mean, $\boldsymbol{\mu}_\Theta^*$, that is simply the re-discretized empirical mean pattern with the highest cursor progress:

$$\boldsymbol{\mu}_\Theta^* = \bar{\mathbf{z}}_{\Phi^*}^{\text{before}} \quad (44)$$

where

$$\Phi^* = \underset{\Phi \in \{0^\circ, 3^\circ, \dots, 357^\circ\}}{\operatorname{argmax}} P_{\text{pert}}(\bar{\mathbf{z}}_\Phi^{\text{before}}, \Theta) \quad (45)$$

In other words, the predicted mean activity pattern for movement Θ is an empirical mean activity pattern from a potentially different movement, Φ^* , before learning. [†]

3. Find the movement-specific covariance. To ensure that predicted population activity patterns resembled patterns from the before-learning neural repertoire and had realistic levels of within movement variability, we chose the covariance corresponding to the movement Φ^* from equation (45): $\Sigma_\Theta^{\text{reassociate}} = \mathbf{S}_{\Phi^*}^{\text{before}}$.

4. Generate reassociation-predicted population activity patterns. We obtained N predicted activity patterns by sampling from the 10-D Gaussian specified by $\mu_\Theta^{\text{reassociate}}$ and $\Sigma_\Theta^{\text{reassociate}}$ (equation (29)), rejecting any samples that correspond to non-physiologically plausible spike counts (equation (30)); $1.19\% \pm 2.49\%$ of all samples were rejected for monkey J; $3.18\% \pm 4.88\%$ for monkey L; and $0.40\% \pm 0.58\%$ for monkey N).

Subselection

Under subselection, the animal improves performance for a given movement by producing only the before-learning population activity patterns that are appropriate for that same movement through the perturbed BCI mapping. In this sense, the after-learning movement-specific neural repertoire is a subset of what it was before learning. The following steps describe the prediction procedure, which was done independently for each experimental session and for each intended movement direction, Θ :

1. Determine the empirical mean and covariance of the before-learning activity patterns, $\bar{\mathbf{z}}_\Theta^{\text{before}}$ and $\mathbf{S}_\Theta^{\text{before}}$, respectively, using equations (27)-(28).

2. Quantify empirical behavioral performance in the after-learning trials by computing mean cursor progress, $\bar{P}^{\text{after}} = P_{\text{pert}}(\bar{\mathbf{z}}_\Theta^{\text{after}}, \Theta)$, using equation (10).

3. Generate before-learning population activity patterns. We obtained N before-learning samples, $\tilde{\mathbf{z}}_{\Theta,1}^{\text{before}}, \dots, \tilde{\mathbf{z}}_{\Theta,N}^{\text{before}}$, by sampling from the 10-D Gaussian specified by $\bar{\mathbf{z}}_\Theta^{\text{before}}$ and $\mathbf{S}_\Theta^{\text{before}}$ (equation (29)), rejecting any samples that correspond to non-physiologically plausible spike counts (equation (30)).

4. Quantify the behavioral performance of the sampled activity patterns through the perturbed BCI mapping by computing mean cursor progress, $\bar{P}^{\text{sampled}} = \frac{1}{N} \sum_{n=1}^N P_{\text{pert}}(\tilde{\mathbf{z}}_{\Theta,n}^{\text{before}}, \Theta)$, using equation (10).

5. Subselect sampled activity patterns to match after-learning behavioral performance. We dropped the k samples with the worst cursor progress, choosing k such that the mean cursor progress of the remaining samples best matched the empirical mean after-learning cursor progress, \bar{P}^{after} . In the rare event that cursor progress decreased after learning for a particular intended movement direction (i.e., $\bar{P}^{\text{after}} < \bar{P}^{\text{sampled}}$), we instead dropped the k samples with the best cursor progress. In practice, this subselection procedure preserved a large majority of the before-learning activity patterns (the percentage of samples dropped using this subselection procedure, i.e.,

[†]This solution emerges because by definition all of the re-discretized empirical mean patterns, $\bar{\mathbf{z}}_\Phi^{\text{before}}$, correspond to physiological spike count vectors, and because the cursor-progress objective function in equation (39) is linear in μ_Θ . Thus, if μ_Θ moves from $\bar{\mathbf{z}}_{\Phi^*}^{\text{before}}$ (as defined in equation (44)) toward any other $\bar{\mathbf{z}}_\Phi^{\text{before}}$, cursor progress will change linearly from $P_{\text{pert}}(\bar{\mathbf{z}}_{\Phi^*}^{\text{before}}, \Theta)$ toward $P_{\text{pert}}(\bar{\mathbf{z}}_\Phi^{\text{before}}, \Theta)$. This change is necessarily a decrease because of equation (45). A similar argument can be made for the case of μ_Θ moving from $\bar{\mathbf{z}}_{\Phi^*}^{\text{before}}$ toward any other point in the convex hull of the $\bar{\mathbf{z}}_\Phi^{\text{before}}$. Therefore no point in the convex hull of the $\bar{\mathbf{z}}_\Phi^{\text{before}}$ has a larger cursor progress than that of $\bar{\mathbf{z}}_{\Phi^*}^{\text{before}}$.

$100 \frac{k}{N}$, was $23.9\% \pm 7.03\%$ for monkey J; $14.6\% \pm 4.76\%$ for monkey L; and $32.3\% \pm 11.6\%$ for monkey N).

6. Generate subselection-predicted population activity patterns. To be consistent with prediction procedures for all other hypotheses, we need to predict N patterns. The remaining patterns from Step 5 form the first $N - k$ of these predictions. To obtain the remaining k predictions, we fit a mean and covariance to the $N - k$ samples from Step 5, sampled from the corresponding 10-D Gaussian, and only accepted samples corresponding to physiologically plausible spike counts (equation (30)); $0.22\% \pm 0.36\%$ of all samples were rejected for monkey J; $0.23\% \pm 0.31\%$ for monkey L; and $0.09\% \pm 0.12\%$ for monkey N).

Rescaling

Under rescaling, the repertoire of activity patterns rescales along each dimension of the intrinsic manifold to restore the influence that each dimension had on movement prior to the perturbation. The animal then learns to produce the activity patterns within this rescaled neural repertoire that maximize behavioral performance subject to physiological constraints on the spike counts for each neural unit. The following steps describe the prediction procedure, which was done independently for each experimental session and for each intended movement direction, Θ :

1. Define the rescaled neural repertoire. After learning via rescaling, the *influence* of a particular dimension of the population activity on movement is restored to its pre-perturbation level. We define the *influence* for dimension i under the intuitive BCI mapping as

$$\text{influence}_i = \sigma_i \|\mathbf{b}_i\|_2 \quad (46)$$

where σ_i is the standard deviation of population activity along dimension i , and \mathbf{b}_i is the i -th column in \mathbf{B} of the intuitive BCI mapping (equation (1)). The key concept is that a dimension with high variability and a large pushing magnitude will have more influence on cursor movement than a dimension with low variability and the same pushing magnitude or a dimension with the same variability and a small pushing magnitude. To restore influence once the perturbed BCI mapping is in effect, the animal should rescale the activity along dimension i such that

$$\sigma_i^{\text{before}} \|\mathbf{b}_i\|_2 = \sigma_i^{\text{rescaled}} \|\mathbf{b}_i^{\text{pert}}\|_2 \quad (47)$$

To achieve this restored influence, we define the rescaled neural repertoire as the across-movement collection of rescaled before-learning activity patterns,

$$\mathbf{z}_t^{\text{rescaled}} = \begin{bmatrix} f_1 z_{1,t} \\ \vdots \\ f_{10} z_{10,t} \end{bmatrix} \quad (48)$$

where the scale factor for dimension i , f_i is

$$f_i = \frac{\|\mathbf{b}_i\|_2}{\|\mathbf{b}_i^{\text{pert}}\|_2} \quad (49)$$

2. Define the search space for the movement-specific mean within the rescaled neural repertoire. We required the predicted movement-specific mean to be within an interior region of the rescaled neural repertoire (see Step 1 of the reassociation prediction procedure). We defined this interior region by the convex hull of the rescaled movement-specific means, $\bar{\mathbf{z}}_{\Phi}^{\text{rescaled}}$, and we required the predicted mean to correspond to a physiologically plausible mean spike count vector (as in equation (25)). The rescaled movement-specific means, $\bar{\mathbf{z}}_{\Phi}^{\text{rescaled}}$, were computed

using equation (28), but replacing \mathbf{z}_t with $\mathbf{z}_t^{\text{rescaled}}$ of equation (48).

3. Find the movement-specific mean. We chose the mean population activity pattern, $\boldsymbol{\mu}_\Theta^{\text{rescale}}$, to be the pattern from the search space defined in Step 2 that maximizes behavioral performance under the perturbed BCI mapping. This corresponds to solving the optimization problem for reassociation, but over the rescaled neural repertoire rather than over the fixed before-learning repertoire (i.e., equations (39)-(43) with $\bar{\mathbf{z}}_\Phi^{\text{before}}$ replaced by $\bar{\mathbf{z}}_\Phi^{\text{rescaled}}$ in equation (40)). Note that, although the solution to this problem in the case of reassociation, $\boldsymbol{\mu}_\Theta^* = \bar{\mathbf{z}}_{\Phi^*}^{\text{before}}$, corresponds to the before-learning empirical mean pattern for some movement Φ^* (i.e., the optimal weights w_Φ^* are 0 except for $w_{\Phi^*}^* = 1$), in general under rescaling $\boldsymbol{\mu}_\Theta^* \neq \bar{\mathbf{z}}_\Phi^{\text{rescaled}}$ for any Φ . Rather, $\boldsymbol{\mu}_\Theta^*$ is a combination of multiple $\bar{\mathbf{z}}_\Phi^{\text{rescaled}}$ (i.e., multiple weights w_Φ^* are nonzero) because in general the $\bar{\mathbf{z}}_\Phi^{\text{rescaled}}$ might not all correspond to physiologically plausible spike counts, and it may be necessary to mix among multiple $\bar{\mathbf{z}}_\Phi^{\text{rescaled}}$ to satisfy the physiological constraints in equations (42) and (43). This mixing of multiple $\bar{\mathbf{z}}_\Phi^{\text{rescaled}}$ has the effect of choosing a $\boldsymbol{\mu}_\Theta^*$ that is on the interior of the convex hull of the $\bar{\mathbf{z}}_\Phi^{\text{rescaled}}$, where the activity patterns correspond to physiological spikes counts, rather than along the boundary where patterns might not be physiological.

4. Find the movement-specific covariance. We chose the movement-specific covariance, $\Sigma_\Theta^{\text{rescale}}$, to be the weighted combination of rescaled movement-specific covariances:

$$\Sigma_\Theta^{\text{rescale}} = \sum_{\Phi \in \{0^\circ, 3^\circ, \dots, 357^\circ\}} w_\Phi^* \mathbf{S}_\Phi^{\text{rescaled}} \quad (50)$$

where the weights w_Φ^* are those obtained in Step 3, and the rescaled movement-specific covariances, $\mathbf{S}_\Phi^{\text{rescaled}}$, are obtained using equation (27), but replacing \mathbf{z}_t with $\mathbf{z}_t^{\text{rescaled}}$ and replacing $\bar{\mathbf{z}}_\Phi^{\text{before}}$ with $\bar{\mathbf{z}}_\Phi^{\text{rescaled}}$.

5. Generate rescaling-predicted population activity patterns. We obtained N predicted activity patterns by sampling from the 10-D Gaussian specified by $\boldsymbol{\mu}_\Theta^{\text{rescale}}$ and $\Sigma_\Theta^{\text{rescale}}$ (equation (29)), rejecting any samples that correspond to non-physiologically plausible spike counts (equation (30); $10.8\% \pm 9.69\%$ of all samples were rejected for monkey J; $14.9\% \pm 11.2\%$ for monkey L; and $6.84\% \pm 4.15\%$ for monkey N). Although these rejected samples represent a small minority of all samples generated, rescaling did incur substantially higher rejection rates relative to the other hypotheses we considered. This difference is due to a modest amount of activity specified by the rescaled movement-specific covariances (from Step 4) that corresponds to non-physiological firing rates. If we replace these rescaled movement-specific covariances with the average movement-specific covariance from the before-learning data (as in equation (35) for realignment), rejection rates are on par with those from the other hypotheses we considered ($2.02\% \pm 3.13\%$ for monkey J; $4.64\% \pm 4.57\%$ for monkey L; and $1.67\% \pm 0.91\%$ for monkey N), and the key features of the predicted neural activity (i.e., as in **Figs. 4-7**) are qualitatively unchanged.

On the range of predicted behavioral performance

The range of predicted behavioral performance levels across strategies (**Fig. 7**) stems from the neural constraints imposed by each strategy. All strategies we considered predict individual-neuron firing rates within physiological ranges (equations (25), (26), and (30)) that were identified empirically from the before-learning data. All strategies also predict population activity patterns that lie within the intrinsic manifold. Realignment does not impose any additional neural constraints. As a result, realignment predicts movements that are both faster and more accurate than the observed movements and those predicted by reassociation and rescaling (**Supplementary Fig. 7**). Rescaling and reassociation impose additional neural constraints, which restrict behavioral performance. reassociation constrains the after-learning neural repertoire to match the before-learning repertoire (equations (40) and (41)). Rescaling constrains the after-learning repertoire in a manner that restores the intuitive-level behavioral influence of activity along each dimension of the population activity (equations (46)-(49)).