

1 **Predicting individual skill learning, a cautionary tale**

2 Dekel Abeles¹, Jasmine Hertzage², Moni Shahar³, Nitzan Censor^{1,2}

3 1. School of Psychological Sciences, Tel Aviv University, Tel Aviv 69978, Israel

4 2. Sagol School of Neuroscience, Tel Aviv University, Tel Aviv 69978, Israel

5 3. AI and Data Science Center of Tel-Aviv University (TAD), Tel Aviv 69978, Israel

6

7 **Abstract**

8 People show vast variability in skill learning. What determines a person's individual learning
9 ability? In this study we explored the possibility to predict participants' future learning, based
10 on their behavior during initial skill acquisition. We recruited a large online multi-session
11 sample of participants performing a sequential tapping skill learning task. We trained machine
12 learning models to predict future skill learning from raw data acquired during initial skill
13 acquisition, and from engineered features calculated from the raw data. While the models did
14 not explain learning, strong correlations were observed between initial and final performance.
15 In addition, the results suggest that in correspondence with other empirical fields testing
16 human behavior, canonical experimental tasks developed and selected to detect average
17 effects may constrain insights regarding individual variability, relevant for real-life scenarios.
18 Overall, implementing machine learning tools on large-scale data sets may provide a powerful
19 approach towards revealing what differentiates between high and low innate learning abilities,
20 paving the way for learning optimization techniques which may generalize beyond motor skill
21 learning to broad learning abilities.

22

23 Keywords: human skill acquisition, motor learning, individual differences, sequence learning,
24 machine learning

25

26

27

28 **Introduction**

29 People vary substantively in their ability to execute daily skills. What are the sources of such
30 variability? Most studies have focused on initial and online task performance, known to vary
31 between individuals (Anderson, Lohse, Lopes, & Williams, 2021). Thus, with no prior practice,
32 some individuals might exhibit outstanding performance, while others might express slow and
33 inaccurate performance. Importantly, people vary greatly in their ability to learn new skills as
34 well, with the range of possible improvement differing between individuals. Predicting learning
35 based on early skill acquisition offers an abundance of benefits and may be useful for effective
36 adjustment of training regimes in daily life and for neurorehabilitation. What determines
37 individual differences in learning abilities? Here, we aimed to investigate individual differences
38 in skill learning by predicting the amount of learning an individual will exhibit across different
39 time intervals, based on information extracted from performance at an early session.

40 Investigating individual differences with complex statistical modeling requires a large pool of
41 participants. Therefore to address this question, we leveraged online platforms enabling
42 crowdsourced recruitment producing large-scale data sets (Chandler & Shapiro, 2016; Ranard
43 et al., 2014). Furthermore, the combination of such online platforms along the recent rise of
44 machine learning models as means to understand rich data sets in neuroscience (Richards et al.,
45 2019), provides a unique opportunity to investigate individual differences in skill learning.

46 To predict the extent of learning from skill acquisition characteristics, we utilized a common
47 motor sequence learning task, widely used to model human skill acquisition (Brown &
48 Robertson, 2007; Cohen, Pascual-Leone, Press, & Robertson, 2005; Genzel et al., 2012; Karni et
49 al., 1998; Muellbacher et al., 2002; Perez et al., 2007; Reis et al., 2009; Robertson, Pascual-
50 Leone, & Press, 2004; Wiestler & Diedrichsen, 2013; Wu, Srinivasan, Kaur, & Cramer, 2014).
51 Thus, we conducted a large-scale crowdsourced experiment, recruiting online participants to
52 take part in 3 learning sessions, with a retention session following one week, and an additional
53 long-term retention session following 2-4 months. First, we validated that online participation
54 demonstrates common learning rates within each session as well as between sessions offline
55 gains (Karni et al., 1995; Lugassy, Herszage, Pilo, Brosh, & Censor, 2018; Robertson et al., 2004).
56 Next, we applied a wide array of machine learning models based on engineered features
57 derived from existing literature of motor skill learning, as well as models based on raw data,
58 using machine extracted features with no involvement of prior knowledge.

59

60 **Methods**

61 *Participants*

62 Participants were recruited online from the Amazon Mechanical Turk platform
63 (<https://www.mturk.com>). Qualifications for registered MTurk workers to participate in the first
64 session of the experiment were: above 95% approval rate in previous MTurk assignments,
65 currently located in the USA, right-handed, and did not previously participate in a sequential
66 tapping task from our lab. Each of the following sessions were made available to qualified
67 participants according to the predefined scheduling scheme and was available for 12 hours.
68 Data were collected using non overlapping batches of participants – session 1 of the
69 experiment was made available on a Monday and the next sessions accordingly. This resulted in
70 the following number of participants per session: Session 1: 571 participants, Session 2: 334,
71 Session 3: 273, Session 4: 195, Session 5: 103. Additional exclusion criteria were enforced to
72 make sure the remaining sample of participants were all attentive and complied with
73 instructions (see below). This resulted in the final sample of: session 1: N=460; 274 Female;
74 Mean age = 43.35, Std = 12.99; session 2: N=254; 154 Female; Mean age = 43.29, Std = 12.83;
75 session 3: N=203; 116 Female; Mean age = 44.07, Std = 12.72; session 4: N=134; 75 Female;
76 Mean age = 46.08, Std = 13.00; session 5: N=75; 39 Female; Mean age = 47.48, Std = 12.47. All
77 participants used a button press to sign an online informed consent form presented at the
78 beginning of each session. The payment scheme for all sessions was visible in the experiment
79 page on the Mturk platform. To minimize dropouts, the compensation increased as sessions
80 progressed (1.5\$, 2\$, 2.5\$, 2\$ for the shorter 4th Retention session, and 5\$ for the final long-
81 term Retention session).

82 *Task*

83 Participants performed a procedural motor task - the sequence tapping task (Karni et al., 1995),
84 a highly common task used in numerous motor learning studies (Albouy et al., 2012; Bönstrup,
85 Iturrate, Hebart, Censor, & Cohen, 2020; Herszage, Sharon, & Censor, 2021; Rickard, Cai, Rieth,
86 Jones, & Ard, 2008). Participants were instructed (using illustrative slides) to place their non-
87 dominant left hand on their keyboard in a one-to-one correspondence between fingers and
88 digit-numbers; pinky – #1, ring finger – #2, middle finger – #3, index finger – #4. They were
89 instructed to repeatedly tap the requested pattern (4-1-3-2-4) as fast and as accurate as
90 possible using their left hand for the entire trial duration (10 seconds). A 10 second count-down
91 screen preceded each trial and served as a break. Feedback was provided in the form of dots,
92 with each keypress adding an additional dot to the display, regardless of correctness. Except for
93 the sequence itself, this was the only visible item on the screen during the trial. The experiment
94 was programmed in Psychopy (Peirce et al., 2019) and was hosted on Pavlovia servers
95 (<https://pavlovia.org/>).

96 *Experimental procedure*

97 Before the first session, participants reported their age, gender, education level, time of weekly
98 engagement with musical instruments and time engaged in physical activities. Additionally, at
99 the beginning of each session, participants were asked to report the duration and the quality of
100 sleep on the night preceding that session. At the end of each session, as a simple attention
101 check, participants were asked to report the hand they used to perform the task. The study
102 initially comprised of 4 sessions - each consisting of 36 trials except for the Retention session
103 (4th session) containing 9 trials. A fifth session, the long-term Retention session, was made
104 available 2-4 months after the completion of the Retention session, and comprised of 36 trials,
105 identical to the first 3 sessions (figure 1a).

106 *Data analysis and machine learning feature engineering*

107 All analyses were performed using custom code written in python (Van Rossum & Drake Jr,
108 1995). Data preprocessing and handling was done using the Numpy (Harris et al., 2020) and
109 Pandas (McKinney, 2010) package. The machine learning pipeline was defined using Scikit-learn
110 (Pedregosa et al., 2011) and Pytorch (Paszke et al., 2019). The Matplotlib (Hunter, 2007) and
111 Seaborn (Waskom, 2021) libraries were used for data visualization. Statistical analysis was
112 conducted using Pinguin (Vallat, 2018).

113 Participants were qualified to continue to the next session if they did not end the experiment
114 mid-session and averaged at least 9 input characters per trial. Additionally, to validate
115 participants' attention to the task, data were discarded from all sessions if participants were
116 too slow to start the trial following a break (first input exceeded 2 seconds) or failed to respond
117 in more than 5 trials per session. Next, if the reported sleep duration was outside of the
118 acceptable range of 6-12 hours, the data from that session and all following sessions were
119 discarded.

120 Performance was defined as the overall number of correct keypresses in a trial (Censor,
121 Horovitz, & Cohen, 2014; de Beukelaar, Woolley, & Wenderoth, 2014; Herszage et al., 2021;
122 Korman et al., 2007). Keypresses were deemed correct if they were part of the complete
123 requested pattern (4-1-3-2-4). If the trial ended mid-pattern, all keypresses from the start of
124 that pattern were also considered correct. To minimize the effects of fatigue, *learning* was
125 defined as the difference between the average of the 3 best trials in each session.

126 The following statistics were extracted from each session for each participant: *start*
127 *performance* was defined as the average of trials number 2 and 3 (trial 1 not included due to
128 warm-up decrements) (Adams, 1952; Rickard et al., 2008). *End performance* was defined as the
129 mean of the last 3 trials in a session. *Maximal and minimal performance* were defined as the
130 mean of the 3 trials with highest/lowest performance within each session. *Offline gains* were
131 defined as the difference between consecutive sessions i.e., the *start performance* in session
132 n+1 was deduced from the *end performance*. *Continuity* was defined as the average of the

133 longest consecutive correct keypress of each trial across an entire session (Herszage et al.,
134 2021). The *mean accuracy* was also computed for each participant in each session based on the
135 average accuracies in all trials within the session. Additionally, the average response time of the
136 first keypress of each trial across the session was defined as the *mean first RTs* and used as a
137 proxy for estimating the level of attentiveness during the trial.

138 *Session dynamics*. Session performance, defined as number of correct keypresses per trial
139 within a session, was fitted with a learning curve according to the following equation:

$$140 \quad T_n = T_1 n^{-l(n)}, l(n) = l + f_p + 1 - \exp(f_p(n^{f_p} - 1))$$

141 where T – the amount of correct keypresses, l – learning rate, f_p –
142 fatigue paramter (Asadayoobi, Jaber, & Taghipour, 2021), n – trial number.

143 *Scipy.optimize.curve_fit* (initial guess for parameters (0.5,0.2,0) all bounded between [0-1]) was
144 used to find the optimal Parameters f_p , l and T_1 for each participant and session.

145 *End of session slopes*. A regression line (intercept and slope) was fitted for the number of
146 correct trials for the last 15 trials in the session separately for each participant and session
147 (excluding session 4, which included only 9 trials).

148 *Locally weighted scatterplot smoothing (lowess) features*. For each participant, the correct
149 number of keypresses per trial were smoothed across the session using a non-parametric local
150 regression (*statsmodel.api.nonparamateric.lowess*, *frac* = 0.5). Several features were extracted
151 from the smoothed curve. First, we defined the regions of plateau on the curve as the longest
152 streak of consecutive trials in which the derivative was below 0.25, meaning that the smoothed
153 improvement between trials was less than a quarter of a keypress. The start and end of the
154 plateau were defined as the first and last trials within this streak and the streak count was their
155 difference. Additionally, the maximum of the smoothed curve and its index within the session
156 (the trial in which it was achieved) were also extracted per participant and session.

157 *Within sequence consistency dynamics*. To derive a representation of within sequence dynamics
158 we first extracted the response time of the last sequence in a correct pattern (the 5th input per
159 sequence) in relation to the first input of the same sequence. This resulted in a vector of last
160 keypress durations (locked to the first input of the sequence) for all correct sequences in the
161 order of execution. To examine the consistency of this input over time we calculated the
162 standard deviation over a running window of 10 consecutive inputs (*running RT consistency*).
163 This running estimate was then fitted with a 3rd degree polynomial (using the *numpy.polyfit*
164 function). The coefficients of this polynomial and the fit prediction error (root mean square
165 error) were used as additional hand-crafted features which capture the pattern dynamics
166 across the session for each participant.

167 *Pattern consistency trend.* To examine the amount of monotonicity apparent in the *running RT*
168 *consistency* estimate, we used Spearman correlation with the corresponding vector of window
169 number within the session. A high negative correlation suggests that a participant's strategy
170 gradually converged to a stable pattern. A high positive correlation on the other hand, suggests
171 a diverged strategy, entering correct sequences less consistently as time progressed.

172

173 *Machine learning modeling*

174 To test the predictive power of the behavior observed during initial training (session 1) on
175 future learning induced by subsequent training sessions, three time intervals were examined: a)
176 change in performance from the 1st session to the 2nd session. b) change in performance from
177 the 1st session to the 3rd session. c) change in performance from the 1st session to the 4th
178 retention session. Two additional time intervals were used to predict skill retention a) one week
179 retention interval (from the 3rd session to the 4th) and b) a long-term retention interval (2-4
180 months) (from the 4th session to the 5th). Note that the number of participants decreases as the
181 experiment reached later sessions, hence the number of observations available for modeling of
182 later intervals is smaller. Accordingly, different modeling approaches were used, as detailed
183 below.

184 The first approach utilized the engineered features as predictors and examined a wide range of
185 machine learning techniques. Specifically, we tested: two tree-based models: Random Forest
186 regression (Ho, 1995) and Sequential Regression Trees using gradient boosting (Xgboost) (J. H.
187 Friedman, 2001). Regularized regression (Elastic net (Zou & Hastie, 2005)) and a multi-layer
188 perceptron (MLP (Haykin, 1994)). Due to the large number of potential predictors, and to avoid
189 over-fitting of the training set, we tested these pipelines both with and without an additional
190 preprocessing step of principle components analysis (PCA)-based dimensionality reduction.
191 Each modeling pipeline started with a standard scaler, transforming the feature values into z-
192 scores. We used grid search for hyper-parameters tuning of the algorithms and regularization
193 parameters. Each set of hyper-parameters was optimized separately for each type of algorithm,
194 predictors step and time interval. The best model was selected based on the average 5-fold
195 cross validation (CV) score. For each model type and time interval, the model selection was
196 done in stages. In each stage an additional set of predictors was introduced based on their
197 complexity, starting with high level features (i.e., session dynamic parameters) and ending with
198 the simplest features (performance per trial). Initially, only non-behavioral features were
199 included (i.e., Age and Gender). Next, predictors were introduced in steps. In the 1st step
200 parameters from the learning curve were introduced. The 2nd step included the parameters
201 extracted to capture *Within sequence consistency dynamics* and the *pattern consistency trend*.
202 The 3rd step included *Lowess based features*. The 4th step included *session statistics*. The 5th
203 step included the micro-offline and micro-online features of the first 5 trials (Bönstrup et al.,

204 2019). And the 6th and final step, included the performance per trial for all trials in the session.
205 For prediction purposes, normalization was done using the means and standard deviations of
206 the variables in the training set. Additionally, we tested a recurrent Long Short-Term Memory
207 (LSTM) network architecture in which the input was the most common end-point measure (de
208 Beukelaar et al., 2014; Herszage & Censor, 2017; Herszage et al., 2021; Karni et al., 1995) of the
209 task - the number of correct keypresses for each trial in the first session.

210 The second approach examined the prediction of future learning, based on all previous
211 sessions. We used a linear regression model with correlation-based feature selection,
212 introducing all available predictors at once and running a hyperparameters grid search on the
213 number of selected features.

214 In the third approach, models were trained directly on raw data from the first session,
215 predicting learning between the first and second training sessions. Task performance was
216 represented as a binary image of size 4 x 7200, where rows represent the key identity (1-4) and
217 columns represent the time where the key was pressed (in 50ms bins). For example, a key press
218 on the key “3” performed 250ms after trial start, will have a value of 1 in the coordinate (3,5).
219 We then trained a convolutional neural network to predict learning. Hyper parameters of the
220 topology and the optimization parameters were tuned manually. Similarly, a convolution
221 encoder-decoder based method was built using the above binary session image as input,
222 geared to reproduce the same image with a compact embedding layer which is then used as
223 features in a regression analysis.

224 *Model evaluation*

225 The parameters that resulted in the best performance on the training-set for each model type
226 and prediction interval were used to re-train the model on the entire training set and examine
227 it on the 20% of hold-out data that was not accessible during training. The final score is thus the
228 reported explained variance (R^2) of the hold-out dataset.

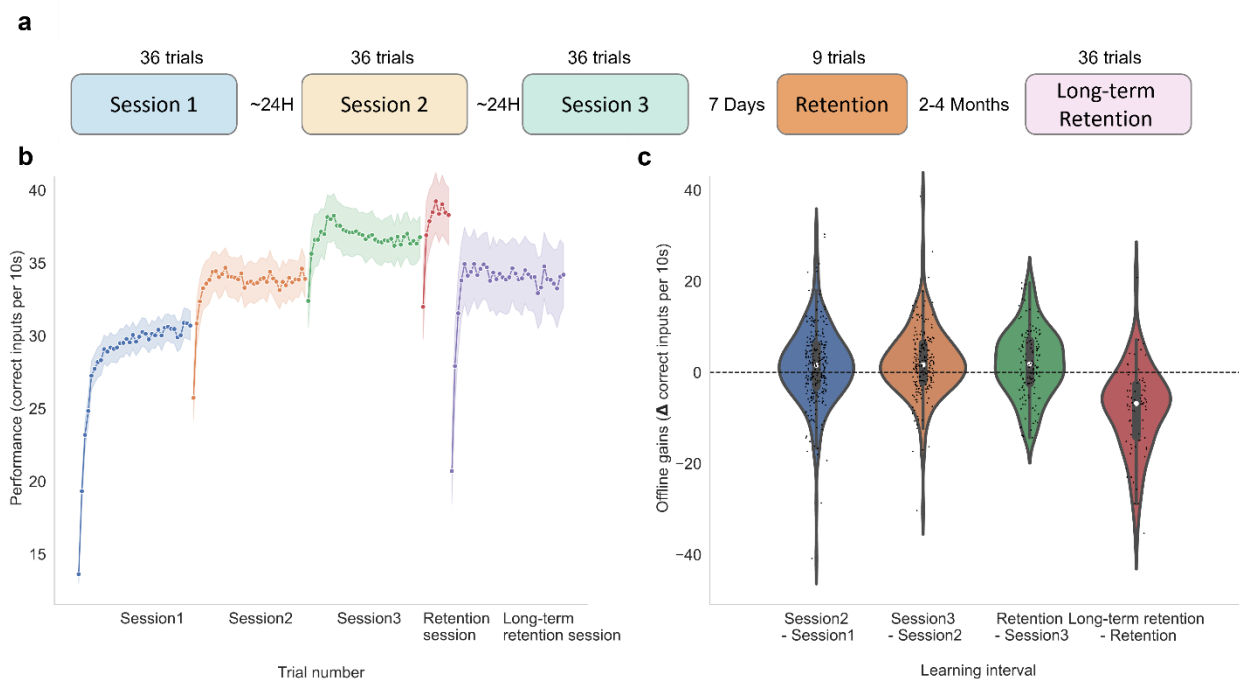
229 *Statistical analysis*

230 One sample t-tests were used to examine the statistical significance of the offline gains analysis.
231 Correlational analyses were conducted using Pearson or Spearman correlation.

232

233 Results

234 We first validated that performance was consistent with previous studies employing the same
235 task in laboratory settings (de Beukelaar et al., 2014; Herszage & Censor, 2017; Karni et al.,
236 1998; Korman et al., 2007). Indeed, participants displayed typical learning curves (figure 1b),
237 with significant learning expressed both within-session, and between-sessions as offline gains
238 (Karni et al., 1998; Press, Casement, Pascual-Leone, & Robertson, 2005; Walker, Brakefield,
239 Morgan, Hobson, & Stickgold, 2002) (figure 1c). Specifically, there were significant offline gains
240 between sessions 1 and 2 ($t(253) = 2.639, p = 0.009, \text{Cohen's } d = 0.126, CI = [0.36 \text{ } 2.45]$), and
241 between sessions 2 and 3 ($t(202) = 4.008, p < 0.001, \text{Cohen's } d = 0.191, CI = [1.08 \text{ } 3.16]$).
242 Interestingly, even when the skill memory was tested following one week, additional offline
243 gains were evident, with a significant improvement between session 3 and Retention session 4
244 ($t(133) = 3.154, p = 0.002, \text{Cohen's } d = 0.183, CI = [0.75 \text{ } 3.28]$). In addition, during the long term
245 retention interval, lasting between 2-4 months (see *Methods*) a significant reduction in
246 performance was observed (difference from Retention (4th session) to Long-term Retention (5th
247 session): $t(74) = -7.661, p < 0.001, \text{Cohen's } d = 0.722, CI = [-10.32 \text{ } -6.06]$), indicating a decay of
248 the memory trace over a period of months. Overall, these results validate typical within and
249 between session motor skill learning.



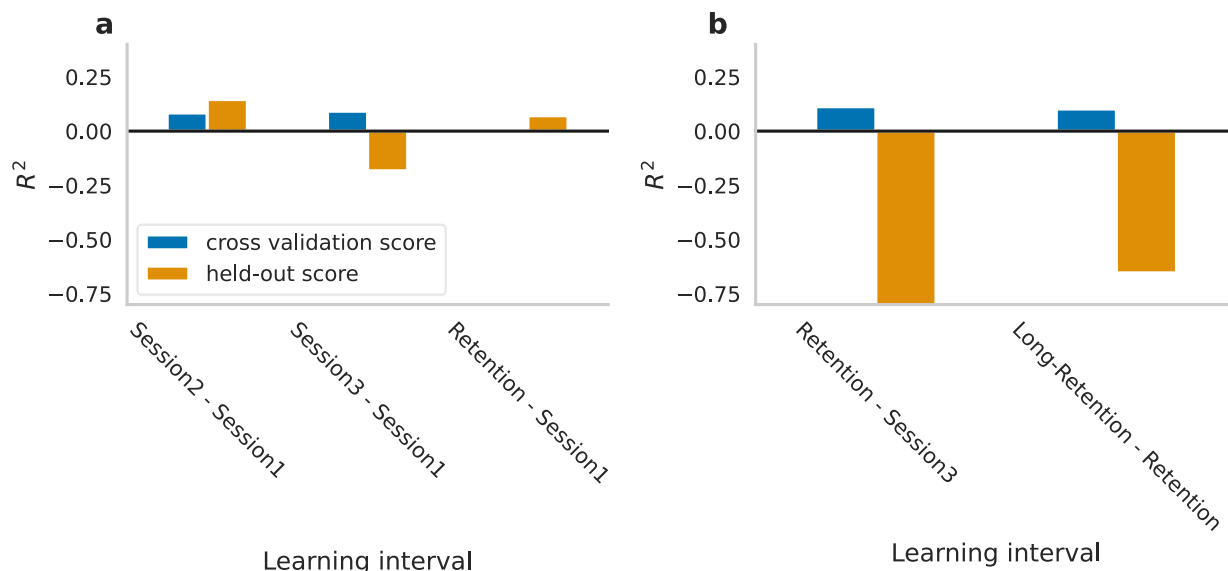
250

251 **Figure 1: Task performance within and between sessions.** a) Experimental design. b) Learning curves across all five
252 sessions (session 1 – blue, session 2 – yellow, session 3 – green, Retention session – orange, Long Term Retention
253 session – pink), the shaded area represents the 95% confidence interval. c) Offline gains between consecutive
254 sessions. Data points in the violin plots represent offline gains for each participant. The white dot represents the
255 median.

256 How could machine learning tools be applied to predict future learning? We first used ML with
257 engineered features (see *Methods*), training discriminative algorithms to predict learning based
258 on performance in the first session. To that end, our goal was to predict the improvements
259 between performance in session 1 and performance in each of the subsequent sessions 2-4. To
260 minimize within session effects of warm-up and fatigue (Adams, 1952; Rickard et al., 2008),
261 between-session learning was quantified based on maximal performance in each session (see
262 *Methods*). Potential predictors were introduced in steps with diminishing feature complexity,
263 ranging from whole session dynamics descriptors, to the number of correct keypresses in each
264 trial. The best performing model was selected based on its mean cross validation and tested on
265 a predetermined hold-out set. Models did not predict learning in the hold-out set (*session2 -*
266 *session1*: $R^2_{mean_cv_score} = 0.08$, $R^2_{test} = 0.15$; *session3 - session1*: $R^2_{mean_cv_score} = 0.09$, $R^2_{test} = -$
267 0.18 ; *Retention session 4 - session1*: $R^2_{mean_cv_score} = 0.01$, $R^2_{test} = 0.07$) (Figure 2a). Of note, a
268 negative R^2 score indicates that model predictions do not explain any variance in the dependent
269 variable.

270 Is behavior at initial stages of skill acquisition indicative of skill retention? To address this
271 question, models were trained to predict the performance change during the short (from
272 session 3 to Retention session) and long retention intervals (from Retention to Long-term
273 retention), based on performance in either the first or all 3 prior sessions. The change in
274 performance over both retention intervals was not predicted by the best performing model
275 (highest cross validation score) as reflected in the negative R^2 in the hold-out set (*Retention*
276 *session - session3*: $R^2_{mean_cv_score} = 0.11$, $R^2_{test} = -0.84$; *Long-retention - Retention session*:
277 $R^2_{mean_cv_score} = 0.10$, $R^2_{test} = -0.65$, figure 2b). Since the long-term retention interval showed
278 negative changes in performance, further investigation of the data revealed that maximum
279 performance in the Retention session was the best predictor for the subsequent long-term
280 retention interval (Pearson's $r(73) = -0.49$, $p < 0.001$, $CI = [-0.65, -0.30]$). Considering that
281 maximum performance in the Retention session reflects both innate abilities and the overall
282 benefit of training throughout the experiment, we examined the correlation between total
283 learning and retention. Pearson correlation confirmed that the amount of total learning
284 throughout the experiment (performance differences between session 1 and Retention session
285 4) was even a stronger predictor of the change in performance (Pearson's $r(73) = -0.58$,
286 $p < 0.001$, $CI = [-0.71, -0.40]$), suggesting that participants exhibit long-term decay of their own
287 learning before the retention interval.

288



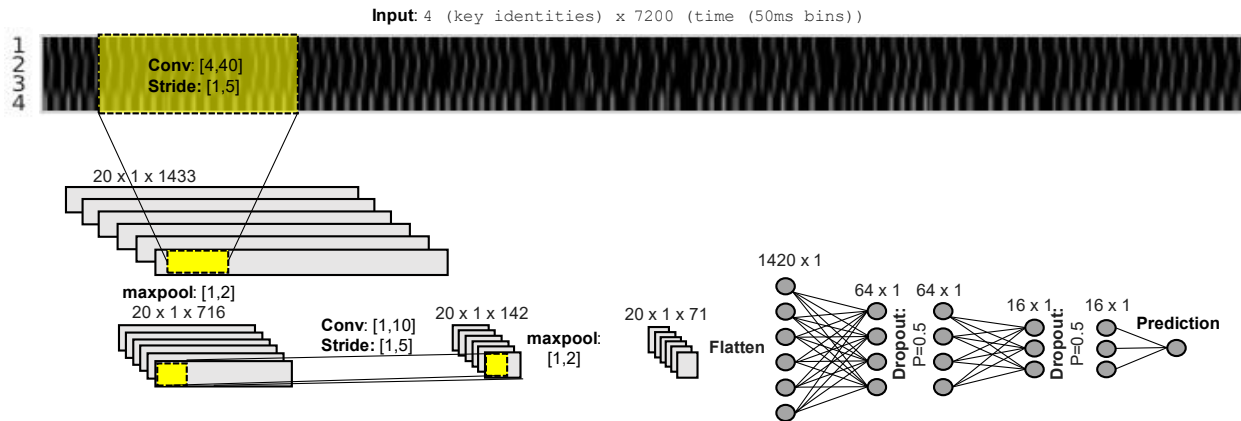
289

290 **Figure 2: Model performance with engineered features.** a) maximum mean cross-validation R^2 scores (blue) and
291 the corresponding hold-out R^2 scores (orange) for each learning interval (X axis). b) Maximum mean cross-
292 validation R^2 scores (blue) and the corresponding hold out R^2 (orange) for the two retention intervals (x axis).

293

294 Next, we tested whether a different approach of machine learning models, avoiding feature
295 selection based on prior assumptions, will achieve better prediction of future learning. To
296 further investigate prediction in that direction, we trained a convolutional neuronal network on
297 data from session 1, represented as a binary matrix of size 4 x 7200, where rows represent key
298 identity and columns represent keypress time within the session in 50ms time bins (Figure 3).
299 This representation reflects the available raw data, without imposing any definition of key
300 correctness. This analysis was focused on the prediction of learning between the first and the
301 second session, which includes the largest pool of participants. Additionally, to better utilize all
302 available data, evaluation of model performance was based solely on cross validation. The best
303 model resulted in mean cross validation $R^2_{test} = -0.049$, std = 0.053 performance. Consistent with
304 this result, two additional models, using a convolution based encoder-decoder and LSTM
305 architectures (see *Methods*), did not show predictive power.

306



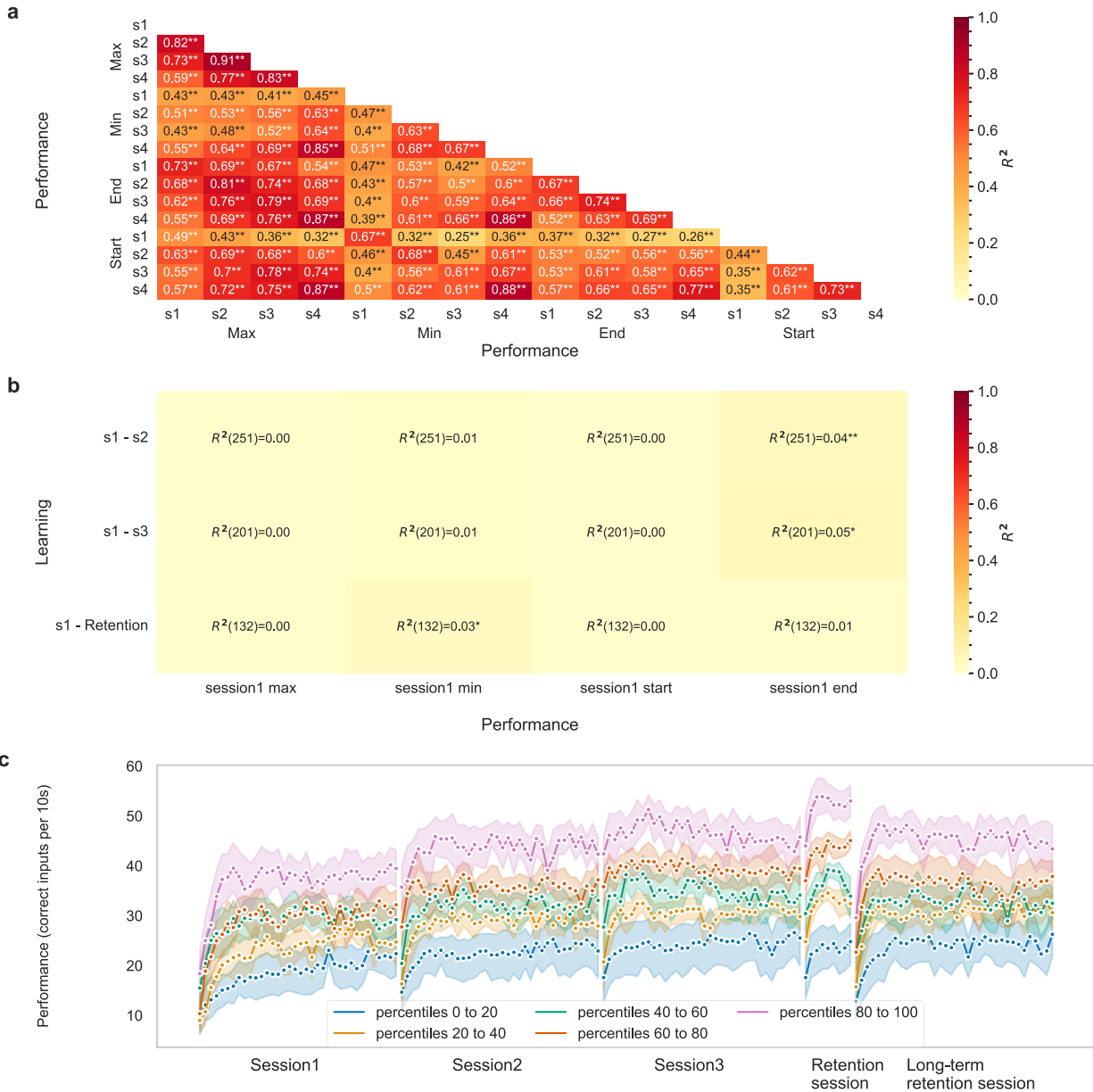
307

308 **Figure 3. Convolution based neural network architecture.** Input was represented as a 4 x 7200 binary matrix, where
309 rows represent key identity (1-4) and columns represent time within the session (in 50ms time bins). The network
310 architecture consists of two convolution layers, each followed by a pooling operation which is followed by 3 fully
311 connected layers. The Rectified linear unit (Relu) was the selected activation function.

312

313 To further investigate the above results, we assessed the consistency of simple performance
314 metrics in each session and between-session learning, using Pearson correlations. Performance
315 in each session explained a large portion of the variance in Performance scores across the 3
316 sessions and Retention session (R^2 range = [0.25-0.91], all $p < 0.001$; see figure 4a), indicating
317 high test-retest reliability and thus a stable measure of individual performance. However,
318 performance hardly explained any portion of the variance in learning (R^2 range = [0.00,0.05];
319 figure 4b). While these results suggest that variability in performance can be explained by
320 performance in previous sessions, variability in learning can hardly be explained. To further
321 illustrate this point, participants were separated into 5 quantile ranges (each spanning
322 20%)(Stafford & Dewar, 2014) based on their maximum performance in the Retention session,
323 plotted throughout the experiment (figure 4c). The plotted curves show that participant's relative
324 performance remained stable throughout the experiment.

-12-



325

326 **Figure 4: Performance was consistent across sessions but did not predict learning.** a) Performance in all sessions
 327 explains a large portion of the variability in performance (R^2 range = [0.25, 0.91]). b) Performance hardly explains
 328 the variability in learning (R^2 range = [0, 0.05]). c) Performance throughout the experiment separated according to
 329 the performance quantile in the Retention session (colors), showing that participants' relative performance rank
 330 remains stable across sessions. Shaded areas represent the 95% confidence interval. Statistical significance is
 331 marked with * for $p < 0.05$ and with ** for $p < 0.001$

332

333

334

335 **Discussion**

336 The goal of this study was to identify what determines an individual's skill learning ability, based
337 on their initial behavior during skill acquisition. Learning was measured at different intervals,
338 using large-scale crowdsourced data. Results showed that performance in early sessions did not
339 predict subsequent learning, while variability in performance was explained by performance in
340 previous sessions. In addition, participants exhibited long-term skill memory decay, bound by
341 their own learning before the retention interval.

342 Machine learning techniques were leveraged to predict learning, utilizing several families of
343 algorithms relying both on manually engineered features and on raw data representations.
344 First, we extracted various features from the observed behavior in the task, ranging from high
345 level features such as the parameters of the learning curve, to simple features such as the
346 correct number of keypress in a trial. The applied models cover a wide array of approaches:
347 Random Forest regression and Xgboost use an ensemble of weak learners and aggregate their
348 predictions either based on consensus (random forest regression) or in a sequential manner.
349 Multi-layered Perceptron (MLP), on the other hand, is a simple deep learning architecture
350 consisting only of fully connected layers. The main advantage of these algorithms is their ability
351 to capture interactions and other non-linear effects between predictors without explicitly
352 modeling them by creating new variables. Two linear regression techniques were also examined
353 due to their straightforward interpretability. Specifically, ElasticNet uses both L1 (Lasso) and L2
354 (Ridge) regularization penalties to limit model complexity while maintaining the linear relation
355 between features and target. Finally, more sophisticated deep learning techniques were
356 examined due to their ability to extract useful features from the data, without relying on expert
357 knowledge and feature engineering.

358 A prerequisite of successful prediction of individual differences is a reliable test-retest metric
359 for prediction (Spearman, 1961). This concept was demonstrated in other fields, such as the
360 field of attentional control, where many canonical tasks, including Stroop (Stroop, 1935),
361 Flanker (Eriksen & Eriksen, 1974), and Navon (Navon, 1977) result in robust between-conditions
362 experimental effects, but in unreliable estimates of individual effects (Hedge, Powell, &
363 Sumner, 2018), thus limiting insights regarding individual differences. Spearman and colleagues
364 attributed this limitation to the calculation of a composite score as the difference between two
365 measurements for the same individual (affecting test-retest reliability, Cronbach & Furby, 1970;
366 Spearman, 1961). Critically, such differences between two measurements are the key outcome
367 for evaluating skill learning. Therefore, while skill learning tasks have extensively shown robust
368 and replicable results when examined between conditions (de Beukelaar et al., 2014; Gabitov et
369 al., 2017; Herszage & Censor, 2017; Herszage et al., 2021; Korman et al., 2007), insights into
370 individual differences may be limited. Accordingly, while large sample sizes may reduce
371 standard errors and enable to detect average between-conditions effects, they do not

372 necessarily improve the reliability of individual effects. This issue could be addressed by
373 increasing the number of repeated measures or trials for each participant, as done for example
374 in studies of perceptual learning (Sagi, 2011).

375 Furthermore, our analysis revealed that separating participants into 5 groups based on their
376 performance in the Retention session, resulted in a visible, consistent classification throughout
377 all sessions, suggesting that future learning may be too small to change participants' rank.
378 Participants showing higher performance at the beginning, will also result in better
379 performance at the end of the experiment. These results are consistent with previous findings
380 of a large online sample of participants playing a complex online shooter game (Stafford &
381 Dewar, 2014). When participants were split into 5 quantile ranges based on their best
382 performance the curves remained separated from the very beginning of the task. Development
383 of novel model motor skill tasks with high variability in between-session learning, and in which
384 future performance is not determined by initial performance, may overcome the above
385 constraints and provide further insights regarding learning variability, important for real-life
386 scenarios. These may be combined with potentially useful predictors from other domains
387 (Ackerman, 1987; Anderson et al., 2021; Chen, Gully, Whiteman, & Kilcullen, 2000), functional
388 and anatomical neuroimaging information (Tomassini et al., 2011), or high-resolution kinematic
389 inputs (Friedman & Korman, 2012).

390 In correspondence with other empirical fields testing human behavior, canonical experimental
391 tasks developed and selected to detect average effects may constrain insights regarding
392 individual variability, relevant for real-life scenarios. Accordingly, development of novel tasks
393 with high test-retest reliability which model real-life learning, may shed light on the underlying
394 mechanisms of individual differences in skill learning and promote personalized learning
395 regimes geared to enhance human performance. Consequently, collecting large online datasets
396 of behaving participants combined with advanced machine learning approaches, holds great
397 potential for modeling future learning based on easily observable behavior during initial
398 training. In turn, this may allow efficient resource allocation and enhancement of training
399 regimes tailored to each person according to their innate abilities.

400

401 *Data and code availability*

402 All collected data and the code for analysis are available upon request.

403

404 **References**

- 405 Ackerman, P. L. (1987). Individual differences in skill learning: An integration of psychometric
406 and information processing perspectives. *Psychological Bulletin*, *102*(1), 3.
- 407 Adams, J. A. (1952). Warm-up decrement in performance on the pursuit-rotor. *The American*
408 *Journal of Psychology*, *65*(3), 404–414.
- 409 Albouy, G., Sterpenich, V., Vandewalle, G., Darsaud, A., Gais, S., Rauchs, G., ... Maquet, P.
410 (2012). Neural correlates of performance variability during motor sequence acquisition.
411 *NeuroImage*, *60*(1), 324–331. <https://doi.org/10.1016/j.neuroimage.2011.12.049>
- 412 Anderson, D. I., Lohse, K. R., Lopes, T. C. V., & Williams, A. M. (2021). Individual differences in
413 motor skill learning: Past, present and future. *Human Movement Science*, *78*, 102818.
- 414 Asadayoobi, N., Jaber, M. Y., & Taghipour, S. (2021). A new learning curve with fatigue-
415 dependent learning rate. *Applied Mathematical Modelling*, *93*, 644–656.
416 <https://doi.org/10.1016/j.apm.2020.12.005>
- 417 Bönstrup, M., Iturrate, I., Hebart, M. N., Censor, N., & Cohen, L. G. (2020). Mechanisms of
418 offline motor learning at a microscale of seconds in large-scale crowdsourced data. *Npj*
419 *Science of Learning*, *5*(1), 1–10. <https://doi.org/10.1038/s41539-020-0066-9>
- 420 Bönstrup, M., Iturrate, I., Thompson, R., Cruciani, G., Censor, N., & Cohen, L. G. (2019). A rapid
421 form of offline consolidation in skill learning. *Current Biology*, *29*(8), 1346–1351.
- 422 Brown, R. M., & Robertson, E. M. (2007). Inducing motor skill improvements with a declarative
423 task. *Nature Neuroscience*, *10*(2), 148–149. <https://doi.org/10.1038/nn1836>
- 424 Censor, N., Horowitz, S. G., & Cohen, L. G. (2014). Interference with Existing Memories Alters
425 Offline Intrinsic Functional Brain Connectivity. *Neuron*, *81*(1), 69–76.
426 <https://doi.org/10.1016/j.neuron.2013.10.042>
- 427 Chandler, J., & Shapiro, D. (2016). Conducting clinical research using crowdsourced convenience
428 samples. *Annual Review of Clinical Psychology*, *12*, 53–81.
- 429 Chen, G., Gully, S. M., Whiteman, J.-A., & Kilcullen, R. N. (2000). Examination of relationships
430 among trait-like individual differences, state-like individual differences, and learning
431 performance. *Journal of Applied Psychology*, *85*(6), 835.
- 432 Cohen, D. A., Pascual-Leone, A., Press, D. Z., & Robertson, E. M. (2005). Off-line learning of
433 motor skill memory: a double dissociation of goal and movement. *Proceedings of the*
434 *National Academy of Sciences of the United States of America*, *102*(50), 18237–18241.
- 435 Cronbach, L. J., & Furby, L. (1970). How we should measure "change": Or should we?

- 436 *Psychological Bulletin*, 74(1), 68.
- 437 de Beukelaar, T. T., Woolley, D. G., & Wenderoth, N. (2014). Gone for 60 seconds: Reactivation
438 length determines motor memory degradation during reconsolidation. *Cortex*, 59, 138–
439 145. <https://doi.org/10.1016/j.cortex.2014.07.008>
- 440 Eriksen, B., & Eriksen, C. W. (1974). Effects of noise letters upon the identification of a target
441 letter in a nonsearch task*. *Perception & Psychophysics*, 16(1), 143–149.
- 442 Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of*
443 *Statistics*, 1189–1232.
- 444 Friedman, J., & Korman, M. (2012). Kinematic strategies underlying improvement in the
445 acquisition of a sequential finger task with self-generated vs. cued repetition training. *PLoS*
446 *One*, 7(12), e52063.
- 447 Gabitov, E., Boutin, A., Pinsard, B., Censor, N., Fogel, S. M., Albouy, G., ... Doyon, J. (2017). Re-
448 stepping into the same river: competition problem rather than a reconsolidation failure in
449 an established motor skill. *Scientific Reports*, 7(1), 9406. <https://doi.org/10.1038/s41598-017-09677-1>
- 451 Genzel, L., Quack, A., Jager, E., Konrad, B., Steiger, A., & Dresler, M. (2012). Complex motor
452 sequence skills profit from sleep. *Neuropsychobiology*, 66(4), 237–243.
453 <https://doi.org/10.1159/000341878>
- 454 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ...
455 Oliphant, T. E. (2020). Array programming with {NumPy}. *Nature*, 585(7825), 357–362.
456 <https://doi.org/10.1038/s41586-020-2649-2>
- 457 Haykin, S. (1994). *Neural networks: a comprehensive foundation*. Prentice Hall PTR.
- 458 Hedge, C., Powell, G., & Sumner, P. (2018). The reliability paradox: Why robust cognitive tasks
459 do not produce reliable individual differences. *Behavior Research Methods*, 50(3), 1166–
460 1186. <https://doi.org/10.3758/s13428-017-0935-1>
- 461 Herszage, J., & Censor, N. (2017). Memory Reactivation Enables Long-Term Prevention of
462 Interference. *Current Biology : CB*, 27(10), 1529-1534.e2.
463 <https://doi.org/10.1016/j.cub.2017.04.025>
- 464 Herszage, J., Sharon, H., & Censor, N. (2021). Reactivation-induced motor skill learning.
465 *Proceedings of the National Academy of Sciences of the United States of America*, 118(23),
466 1–6. <https://doi.org/10.1073/PNAS.2102242118>
- 467 Ho, T. K. (1995). Random decision forests. In *Proceedings of 3rd international conference on*

- 468 *document analysis and recognition* (Vol. 1, pp. 278–282).
- 469 Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science &*
470 *Engineering*, 9(3), 90–95. <https://doi.org/10.1109/MCSE.2007.55>
- 471 Karni, A., Meyer, G., Jezard, P., Adams, M. M., Turner, R., & Ungerleider, L. G. (1995).
472 Functional MRI evidence for adult motor cortex plasticity during motor skill learning.
473 *Nature*, 377(6545), 155–158. <https://doi.org/10.1038/377155a0>
- 474 Karni, A., Meyer, G., Rey-Hipolito, C., Jezard, P., Adams, M. M., Turner, R., & Ungerleider, L. G.
475 (1998). The acquisition of skilled motor performance: fast and slow experience-driven
476 changes in primary motor cortex. *Proceedings of the National Academy of Sciences of the*
477 *United States of America*, 95(3), 861–868. <https://doi.org/10.1073/pnas.95.3.861>
- 478 Korman, M., Doyon, J., Doljansky, J., Carrier, J., Dagan, Y., & Karni, A. (2007). Daytime sleep
479 condenses the time course of motor memory consolidation. *Nature Neuroscience*, 10(9),
480 1206–1213. <https://doi.org/10.1038/nn1959>
- 481 Lugassy, D., Herszage, J., Pilo, R., Brosh, T., & Censor, N. (2018). Consolidation of complex motor
482 skill learning: Evidence for a delayed offline process. *Sleep*, 41(9), zsy123.
483 <https://doi.org/10.1093/sleep/zsy123>
- 484 McKinney, W. (2010). Data structures for statistical computing in Python. In S. van der
485 Walt & J. Millman (Eds.), *Proceedings of the 9th Python in Science Conference* (pp.
486 56–61). <https://doi.org/10.25080/Majora-92bf1922-00a>
- 487 Muellbacher, W., Ziemann, U., Wissel, J., Dang, N., Kofler, M., Facchini, S., ... Hallett, M. (2002).
488 Early consolidation in human primary motor cortex 3. *Nature*, 415(0028-0836 (Print)),
489 640–644. <https://doi.org/10.1038/nature712>
- 490 Navon, D. (1977). Forest before trees: The precedence of global features in visual perception.
491 *Cognitive Psychology*, 9(3), 353–383. [https://doi.org/10.1016/0010-0285\(77\)90012-3](https://doi.org/10.1016/0010-0285(77)90012-3)
- 492 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... Chintala, S. (2019).
493 PyTorch: An imperative style, high-performance deep learning library. (H. Wallach, H.
494 Larochelle, A. Beygelzimer, F. d\textquotesingle Alché-Buc, E. Fox, & R. Garnett, Eds.),
495 *Advances in Neural Information Processing Systems*. Curran Associates, Inc. Retrieved from
496 [http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
497 [deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
- 498 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay, É.
499 (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*,
500 12(85), 2825–2830. Retrieved from <http://jmlr.org/papers/v12/pedregosa11a.html>

- 501 Peirce, J., Gray, J. R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., ... Lindeløv, J. K.
502 (2019). PsychoPy2: Experiments in behavior made easy. *Behavior Research Methods*, 51(1),
503 195–203. <https://doi.org/10.3758/s13428-018-01193-y>
- 504 Perez, M. A., Tanaka, S., Wise, S. P., Sadato, N., Tanabe, H. C., Willingham, D. T., & Cohen, L. G.
505 (2007). Neural Substrates of Intermanual Transfer of a Newly Acquired Motor Skill. *Current*
506 *Biology*, 17(21), 1896–1902. <https://doi.org/10.1016/j.cub.2007.09.058>
- 507 Press, D. Z., Casement, M. D., Pascual-Leone, A., & Robertson, E. M. (2005). The time course of
508 off-line motor sequence learning. *Cognitive Brain Research*, 25(1), 375–378.
- 509 Ranard, B. L., Ha, Y. P., Meisel, Z. F., Asch, D. A., Hill, S. S., Becker, L. B., ... Merchant, R. M.
510 (2014). Crowdsourcing—harnessing the masses to advance health and medicine, a
511 systematic review. *Journal of General Internal Medicine*, 29(1), 187–203.
- 512 Reis, J., Schambra, H. M., Cohen, L. G., Buch, E. R., Fritsch, B., Zarahn, E., ... Krakauer, J. W.
513 (2009). Noninvasive cortical stimulation enhances motor skill acquisition over multiple
514 days through an effect on consolidation. *Proceedings of the National Academy of Sciences*
515 *of the United States of America*, 106(5), 1590–1595.
516 <https://doi.org/10.1073/pnas.0805413106>
- 517 Richards, B. A., Lillicrap, T. P., Beaudoin, P., Bengio, Y., Bogacz, R., Christensen, A., ... Ganguli, S.
518 (2019). A deep learning framework for neuroscience. *Nature Neuroscience*, 22(11), 1761–
519 1770.
- 520 Rickard, T. C., Cai, D. J., Rieth, C. A., Jones, J., & Ard, M. C. (2008). Sleep Does Not Enhance
521 Motor Sequence Learning. <https://doi.org/10.1037/0278-7393.34.4.834>
- 522 Robertson, E. M., Pascual-Leone, A., & Press, D. Z. (2004). Awareness Modifies the Skill-Learning
523 Benefits of Sleep. *Current Biology*, 14(3), 208–212. [https://doi.org/10.1016/S0960-](https://doi.org/10.1016/S0960-9822(04)00039-9)
524 [9822\(04\)00039-9](https://doi.org/10.1016/S0960-9822(04)00039-9)
- 525 Sagi, D. (2011). Perceptual learning in Vision Research. *Vision Research*, 51(13), 1552–1566.
526 <https://doi.org/https://doi.org/10.1016/j.visres.2010.10.019>
- 527 Spearman, C. (1961). The proof and measurement of association between two things.
- 528 Stafford, T., & Dewar, M. (2014). Tracing the Trajectory of Skill Learning With a Very Large
529 Sample of Online Game Players. *Psychological Science*, 25(2), 511–518.
530 <https://doi.org/10.1177/0956797613511466>
- 531 Stroop, J. R. (1935). Studies of interference in serial verbal reactions. *Journal of Experimental*
532 *Psychology*, 18(6), 643–662. <https://doi.org/10.1037/h0054651>

- 533 Tomassini, V., Jbabdi, S., Kincses, Z. T., Bosnell, R., Douaud, G., Pozzilli, C., ... Johansen-Berg, H.
534 (2011). Structural and functional bases for individual differences in motor learning. *Human*
535 *Brain Mapping*, 32(3), 494–508. <https://doi.org/10.1002/hbm.21037>
- 536 Vallat, R. (2018). Pingouin: statistics in Python. *The Journal of Open Source Software*, 3(31),
537 1026.
- 538 Van Rossum, G., & Drake Jr, F. L. (1995). *Python reference manual*. Centrum voor Wiskunde en
539 Informatica Amsterdam.
- 540 Walker, M. P., Brakefield, T., Morgan, A., Hobson, J. A., & Stickgold, R. (2002). Practice with
541 sleep makes perfect: Sleep-dependent motor skill learning. *Neuron*, 35(1), 205–211.
542 [https://doi.org/10.1016/S0896-6273\(02\)00746-8](https://doi.org/10.1016/S0896-6273(02)00746-8)
- 543 Waskom, M. L. (2021). seaborn: statistical data visualization. *Journal of Open Source Software*,
544 6(60), 3021. <https://doi.org/10.21105/joss.03021>
- 545 Wiestler, T., & Diedrichsen, J. (2013). Skill learning strengthens cortical representations of
546 motor sequences. *ELife*, 2013(2), 1–20. <https://doi.org/10.7554/eLife.00801>
- 547 Wu, J., Srinivasan, R., Kaur, A., & Cramer, S. C. (2014). Resting-state cortical connectivity
548 predicts motor skill acquisition. *NeuroImage*, 91, 84–90.
549 <https://doi.org/10.1016/j.neuroimage.2014.01.026>
- 550 Zou, H., & Hastie, T. (2005). Regularization and variable selection via the Elastic Net. *Journal of*
551 *the Royal Statistical Society, Series B*, 67, 301–320.
- 552
- 553