1    **An emergent temporal basis set robustly supports cerebellar time-series learning**

2    Jesse I. Gilmer[1,2], Michael A. Farries[3], Zachary Kilpatrick[4], Ioannis Delis[5], and Abigail L. Person[2]

3

4    1.    Neuroscience Graduate Program, University of Colorado School of Medicine, Aurora CO
5    2.    Dept. Physiology and Biophysics, University of Colorado School of Medicine, Aurora CO
6    3.    Knoebel Institute for Healthy Aging, University of Denver, Denver CO
7    4.    Department of Applied Mathematics, University of Colorado Boulder, Boulder CO
8    5.    School of Biomedical Sciences, University of Leeds, Leeds UK

9    **Abstract**
10    Learning plays a key role in the function of many neural circuits. The cerebellum is considered a 'learning
11    machine' essential for time interval estimation underlying motor coordination and other behaviors.
12    Theoretical work has proposed that the cerebellum's input recipient structure, the granule cell layer
13    (GCL), performs pattern separation of inputs that facilitates learning in Purkinje cells (P-cells). However,
14    the relationship between input reformatting and learning outcomes has remained debated, with roles
15    emphasized for pattern separation features from sparsification to decorrelation. We took a novel approach
16    by training a minimalist model of the cerebellar cortex to learn complex time-series data from naturalistic
17    inputs, in contrast to traditional classification tasks. The model robustly produced temporal basis sets
18    from naturalistic inputs, and the resultant GCL output supported learning of temporally complex target
19    functions. Learning favored surprisingly dense granule cell activity, yet the key statistical features in GCL
20    population activity that drove learning differed from those seen previously for classification tasks.
21    Moreover, different cerebellar tasks were supported by diverse pattern separation features that matched
22    the demands of the tasks. These findings advance testable hypotheses for mechanisms of temporal basis
23    set formation and predict that population statistics of granule cell activity may differ across cerebellar
24    regions to support distinct behaviors.
25

26    **Introduction**
27    The cerebellum refines movement and maintains calibrated sensorimotor transformations by learning to
28    predict outcomes of behaviors through error-based feedback (Ito, 1972; Herzfeld et al., 2015; Medina
29    2000; Mauk and Buonomano, 2004; Raymond et al., 1996). A major site of cerebellar learning is in the
30    cerebellar cortex, where Purkinje cells (P-cells) receive sensorimotor information from parallel fibers
31    (Huang et al. 2013) whose synaptic strengths are modified by the conjunction of presynaptic (parallel
32    fiber) activity and climbing fiber inputs to P-cells thought to convey instructive feedback (McCormick et
33    al., 1982; Yang and Lisberger, 2014; Mauk et al., 1986; De Zeeuw et al., 1998). P-cell activity is
34    characterized by rich temporal dynamics during movements, representing putative computations of
35    internal models of the body and the physics of the environment (Wolpert et al., 1998; Shadmehr and
36    Mussa-Ivaldi 1994). Parallel fibers are the axons of cerebellar granule cells (GCs), a huge neuronal
37    population (comprising roughly half of the neurons in the entire brain; Herculano-Houzel 2010), which
38    are the major recipient of extrinsic inputs to the cerebellum. Thus, understanding the output of the GCL is
39    key in determining the encoding capacity and information load of incoming activity projected to the
40    cerebellum. Inputs to GCs arise from mossy fibers (MFs), which convey sensorimotor information used
41    by the cerebellum to predict the consequences of motor commands (Rancz et al., 2007; Ishikawa et al.,
42    2015). There are massively more GCs than MFs and each GC typically receives input from just 4 MFs
43    (Palkovits et al., 1971), such that the information carried by each MF is spread among many GCs but each
44    GC samples from only a tiny fraction of total MFs (Jakab and Hamori 1988; Eccles et al., 1967).

45

46    The GCL has been the focus of theoretical work spanning decades that has explored the computational
47    advantages of the unique architecture of the structure. Notably, early studies of the cerebellar circuit by

48    Marr (1969) and Albus (1971) proposed that a key component of the cerebellar algorithm is the sparse
49    representation of MF inputs by GCs. In this view, the cerebellum often must discriminate between
50    overlapping, highly correlated patterns of MF activity with only subtle differences distinguishing them
51    (Bengsston and Jorntell 2009). Sparse recoding of MF activity in a much larger population of GCs
52    ("expansion recoding") increases the dimensionality of population representation and transforms
53    correlated MF activity into independent activity patterns among a subset of GCs (Litwin-Kumar et al.,
54    2017; Cayco-Gajic et al., 2017; Gilmer and Person 2018). These decorrelated activity patterns are easier
55    to distinguish by learning algorithms operating in P-cells, leading to better associative learning and credit
56    assignment (Cayco-Gajic et al., 2017; Sanger et al., 2020).
57
58    The machine learning perspective of Marr-Albus theory tends to assume that the cerebellum is presented
59    with a series of static input patterns that must be distinguished and categorized. However, neuronal
60    population dynamics are hardly ever static and precise timing of circuit inputs to the cerebellum remains
61    an essential part of cerebellar function. Mauk and Buonomano (2004) revisited cerebellar expansion
62    recoding in the context of delayed eyeblink conditioning, a cerebellum-dependent learning task where the
63    subject hears a tone followed by an aversive air puff to the eye at a fixed delay from tone onset and must
64    learn to initiate an eyeblink at the correct delay to protect the eye. They proposed that a static activity
65    pattern in MFs (representing the tone) could be recoded in the GC layer as a temporally evolving set of
66    distinct activity patterns. P-cells could learn to recognize the GC activity pattern present at the correct
67    delay and initiate an eyeblink to avert the "error" signal representing the air puff to the eye. In other
68    words, P-cells would select from a "temporal basis set" for correct error prediction and learning adaptive
69    behavior.
70
71    Expansion recoding creates the possibility of representing a single MF pattern as a series of distinct GC
72    patterns (a "temporal basis set"; Albus 1975; Zhou et al., 2020; Tyrrell and Willshaw 1992; Liu et al.,
73    2019; Kalmbach et al., 2011). The existence of this predicted temporal basis set within the cerebellum has
74    been supported experimentally in electric fish, where GCs represent the duration of mimicked electric
75    organ discharge through a range of onsets (Kennedy et al., 2014). Although these studies have been
76    highly influential, little is known about how the GCL would produce a temporally diverse basis set from
77    static input data. Local inhibition, short-term synaptic plasticity, and varying GC excitability all may
78    work together to diversify time-invariant input (Chabrol et al., 2015; Duguid et al, 2012; Crowley et al.,
79    2009; Rudolph et al., 2015; Buonomano and Mauk 1994; Kanichay and Silver 2008; Simat et al., 2007;
80    Mapelli et al., 2009; Rossi et al., 1996; Gall et al., 2005; Armano et al., 2000; Rizwan et al. 2016; Tabuchi
81    et al., 2019; D'Angelo and De Zeeuw 2009). However, the assumption that MFs ever provide truly static
82    input to the cerebellum is probably unrealistic; even a static stimulus like a tone will generate time-
83    varying activity patterns in the auditory brainstem as units undergo adaptation (Eriksson and Robert
84    1999). Moreover, most of the input signals that the cerebellum must process are intrinsically dynamic
85    (Bengsston and Jorntell 2009; Chabrol et al., 2015). We seek to explore how expansion recoding of
86    dynamic, naturalistic input activity assists cerebellar function.
87
88    To test how expansion recoding of naturalistic input contributes to learning, we developed a simple model
89    of the GCL and a time-series prediction task to explore the effect of putative GCL filtering mechanisms
90    on expansion recoding and learning (Fig. 1A). Similar to previous models, this simplified model made
91    GC activity sparser relative to MF inputs (Marr 1969; Albus 1971) and increased the dimensionality of
92    the input activity (Litwin-Kumar et al., 2017) while preserving information (Billings et al., 2014). That
93    these features of GCL function were achieved using only basic approximations of GC physiology
94    suggests that the crystalline connectivity and feedforward inhibition of the cerebellum incorporated in our
95    model are sufficient to produce pattern separation of naturalistic time-varying inputs. This model
96    demonstrates greatly enhanced learning accuracy and speed by P-cells on a difficult time series prediction
97    task when compared to MF inputs alone.  Although we observed robust sparsening of input activity by
98    GCL output, the relationship between pattern separation metrics and the observed learning was dependent

99    upon the task being performed, suggesting that GCL output covers a span of modalities supporting
100   flexible feature selection by P-cells to meet the needs of particular learning targets. These findings
101   reinforce the ideas explored previously that the GCL balances input sparsening against information loss to
102   optimize learning (Cayco-Gajic et al., 2017; Cayco-Gajic and Silver 2019), and that the balance between
103   these features of GCL output can be functionally controlled through adjustments in the strength of local
104   inhibition. We conclude by showing that muscle activity during reaching movements (Delis, et al. 2018),
105   a proxy for time-varying efference copy signals received by the cerebellum, gives rise to information-
106   preserving sparseness that supports time-series predictions, suggesting that physiological input sources to
107   the GCL, like the spinocerebellar pathways, are sufficient to drive learning. Together, these results
108   suggest that the cerebellar GCL provides a rich basis for learning in downstream Purkinje cells, providing
109   a mixture of lossless representation (Billings et al., 2014) and enhanced spatiotemporal representation
110   (Litwin-Kumar et al. 2017) that are selected for by associative learning to support the learning of diverse
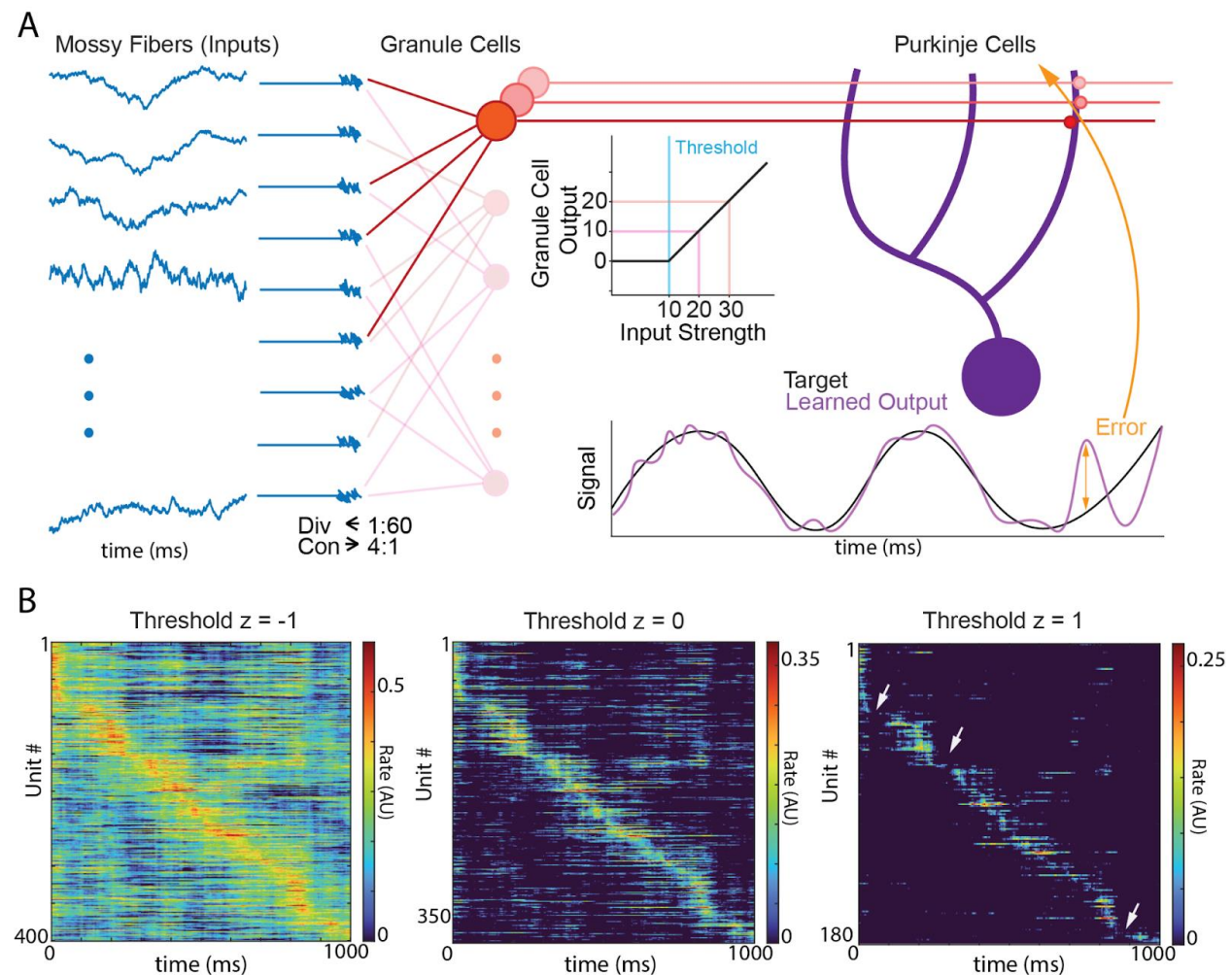111   outputs that support adaptive outputs in a variety of tasks (Fujita 1982; Dean and Porrill 2008).
112

113 **Results**
114 **Temporal basis set formation as emergent property of GCL filtering of physiological-like inputs**
115   The cerebellar granule cell layer (GCL) is theorized to convert spatiotemporally dense inputs into discrete
116   representations through coincidence detection and feedforward and feedback inhibition-mediated
117   thresholding (Marr 1969; Solinas et al., 2010). How the GCL expands spatiotemporal representation has
118   been the subject of debate and scientific inquiry for decades. While cellular and circuit mechanisms have
119   been proposed to expand time invariant signals such as tones (Mauk and Buonomano 2004; Medina
120   2000), naturalistic cerebellar inputs are typically time varying by virtue of dynamic sensorimotor
121   interactions with the environment (Rancz et al., 2007; Eriksson and Robert 1999). Moreover, cerebellar
122   learning is thought to sculpt complex time-varying outputs in Purkinje cells (P-cells) that reflect
123   behavioral adaptations. This observation raises the question of how GCL output supports time series
124   learning, a divergence from traditional classification tasks used in cerebellar models. To address this, we
125   investigated how such naturalistic input patterns were transformed by the GCL to support learning time-
126   varying output patterns, such as those required for generating and correcting movements, or for producing
127   predictions of sensory events (Fig. 1; Izawa et al. 2012).
128

129   We created a simple model capturing the dominant circuit features of the GCL: sparse sampling of mossy
130   fiber (MFs) inputs by postsynaptic granule cells (GCs) and coincidence detection regulated by cellular
131   excitability and local feedforward inhibition (Figure 1A; Eq.1,2; Marr 1969; Albus 1971; Palkovits et al.,
132   1971; Chabrol et al., 2015). MF inputs are represented as smooth time-varying functions, i.e., as variable
133   firing rates rather than spike trains. GC output is generated by summing MF inputs and thresholding the
134   resultant sum; anything below threshold is set to zero while suprathreshold summed activity is passed on
135   (minus the threshold) as GC output (Fig. 1A, center). The GC threshold level represents both intrinsic
136   excitability and the effect of local feedforward inhibition on regulating GC activity. To model MF activity
137   patterns, we sought a statistical ensemble that was rich enough to capture the dynamic nature of
138   naturalistic inputs while remaining analytically tractable and easily parameterized. We chose to utilize the
139   Ornstein-Uhlenbeck (OU) stochastic process, whose output is Gaussian and varies over an adjustable
140   timescale. The statistics of an OU process can be fully characterized by just three parameters: mean,
141   standard deviation, and correlation time; samples drawn from an OU process are shown in Fig. 1A (left,
142   blue). Since the input to GCs is Gaussian in our model, the summed activity that is thresholded is
143   Gaussian as well. For that reason, we found it convenient to define the GC threshold in terms of z-scores.
144   Thus a GC with a threshold of "zero" would have its threshold set at the mean value of its MF inputs;
145   such a GC would be silent 50% of the time on average because the Gaussian presynaptic input would be
146   below the mean value half the time. This makes it possible to discuss functionally similar thresholds
147   across varying network architectures (e.g., a GC with a threshold of zero would discard half of its input
148   on average regardless of whether it received 2 or 8 MF inputs). Via this simple mechanism, our model
149   GCL generates temporally sparse activity that could support learning by downstream P-cells (Fig. 1A,
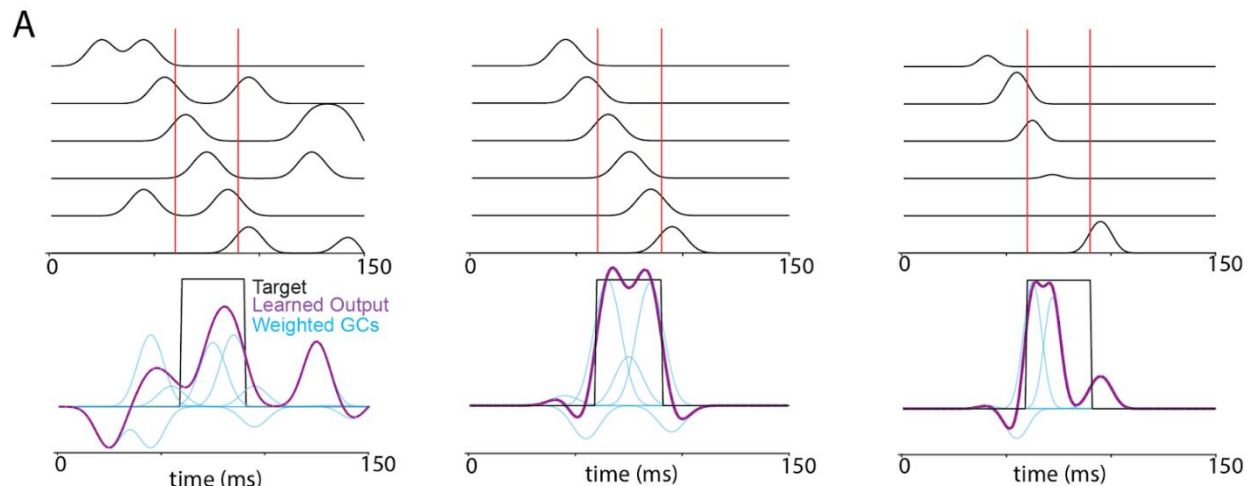
150 right). Indeed, when subjected to this form of filtering, the resultant representation in the GCL population
151 became spatiotemporally distinct at moderate thresholding levels (near 0, Fig. 1B, center). However, too
152 little thresholding resulted in dense representation (Fig. 1B, left) while too much thresholding resulted in
153 over-sparsening, leading to loss of representation in the temporal domain (Fig. 1B, right, arrows indicate
154 loss of representation). The emergence of sparse spatiotemporal representation under the simplistic
155 constraints of the model suggests that the cerebellum's intrinsic circuitry is sufficient to produce
156 spatiotemporal separation when given sufficiently time-varying inputs.
157
158



159
160 *Figure 1: Model architecture and effects of thresholding on GCL population activity.*
161 *A. Diagram of algorithm implmentation. Left shows Ornstein-Uhlenbeck processes (see Methods) as*
162 *proxies for mossy fiber (MFs, blue) inputs to granule cell (GCs, red) units, with convergence and*
163 *divergence of MFs to GCs noted beneath MFs. GCs employ threshold-linear filtering shown beneath the*
164 *red parallel fibers. GC outputs are then transmitted to downstream Purkinje cells (P-cells). P-cells learn*
165 *to predict target functions based on summation of weighted GC inputs and differences between the*
166 *prediction and true target are transmitted as an 'error', which determines the updates to the weights*
167 *between GCs and P-cells. B. Example unit GC population rates when threshold is -1.0, 0 and 1.0*
168 *showing the gradual sparsening of GCL output. Arrows on 1.0 plot indicate regions of gaps in*
169 *representation (lossiness) by the GCL population due to over-sparsening.*
170

A



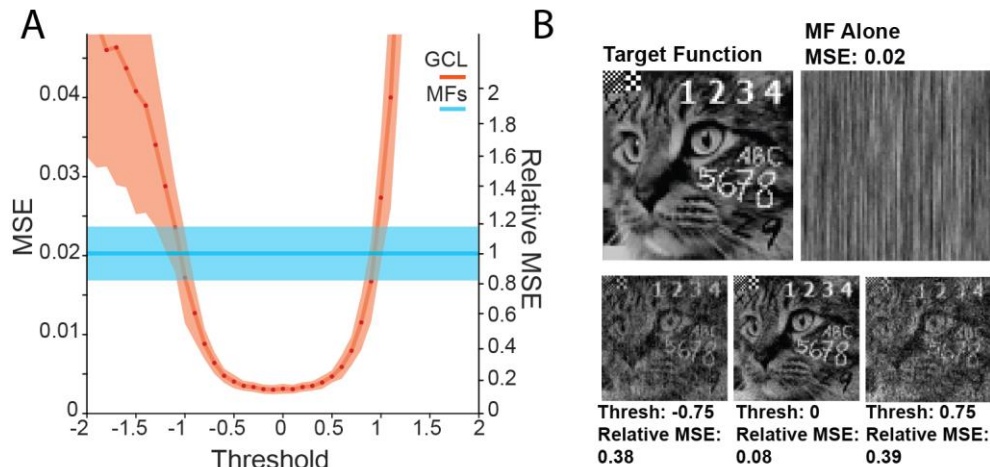*Figure 1, figure supplement 1: Example of basis set utility in learning.*
*A. Diagram relating fictive GC activity (top) with resultant learning (bottom) using those fictive signals*
*as the basis for learning. Target functions are shown in black and learned outputs with minimized error*
*are shown in purple. Note that the best learning occurs with uniform, minimally overlapping GCs, tiling*
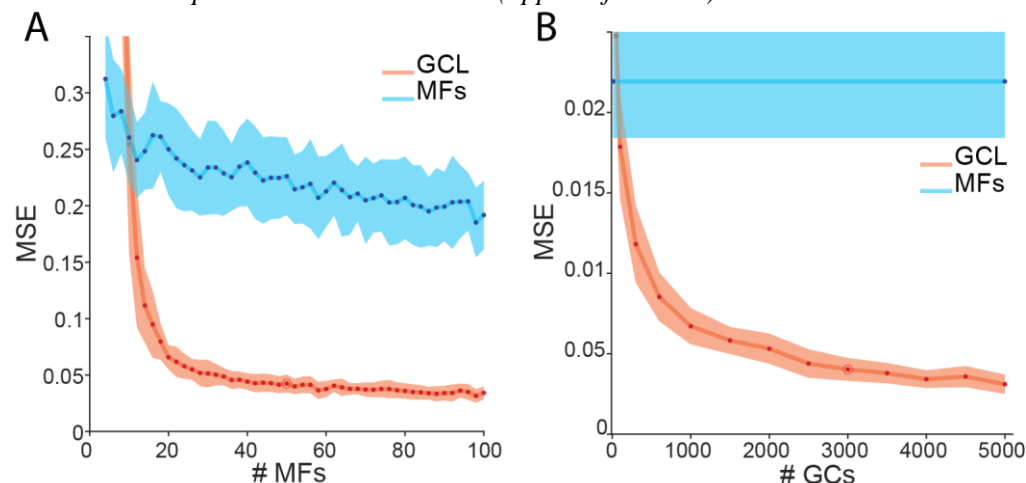*the epoch when the target signal is active (red lines; middle panel).*

**GCL improves time series learning accuracy**

Next, we tested whether GCL population activity seen above assisted learning. We devised a learning task
where P-cells learned to generate a specific time-varying activity pattern in response to the dynamic
activity patterns generated by MFs, which better represents the tasks performed by the cerebellum than
pattern classification. The target patterns that P-cells were tasked with generating were drawn from an OU
process with an autocorrelation time of 10 ms (see Methods). P-cells initially produced output very unlike
the target, but over repeated trials their output converged towards the target function (see Fig. 3A for
example progression of learning). We compared this convergence of P-cell output to target when input
activity was filtered through the GCL to performance the case when MF activity is sent directly to P-cells
("MFs alone"). The GCL enhanced convergence to target at thresholds between –1 and 1 (Fig. 2A),
achieving a minimized mean squared error (MSE) of roughly 0.005 compared to 0.02 when using MFs
alone. It may seem that the performance with MFs alone was still quite good, if slightly quantitatively
inferior, when compared to the range of the target function (normalized to a range of [0,1]). Thus, to
establish intuition into the practical difference of this range of MSEs, we tasked the model with
recapitulating a complex image with an identical range of target function values (with identical range of
[0,1], Fig. 2B). Importantly, the model GCL generated a recognizable image, with an MSE of 0.002 while
experiments using MF alone generated an unrecognizable image with an MSE of 0.02. (The relative
MSE, i.e. the ratio of GCL MSE to MFs alone MSE, was 0.08). Thus, this MSE range represented the
difference between noise and easily recognizable images and text (Fig. 2B top right vs three thresholds,
bottom). This principle was qualitatively true of abstract target functions used in OU input experiments as
well (Fig. 3A for example target functions and estimations). Thus, the inclusion of the GCL in the
filtering process greatly improves learning of complex functions by P-cells in this task, supporting an
order of magnitude improvement in MSE of learned target functions compared to MFs alone. Importantly,
this was not a consequence of the large population expansion between MFs and GCs, as increasing the
number of MFs alone did not improve performance to the levels observed in the model GCL (Supp. Fig.
2A), but a sufficiently large GCL population is required to improve learning (Supp. Fig. 2B).

*Figure 2: Enhanced time series learning using GCL model.*
*A: Relation between mean squared error (MSE) and threshold in a 50 MF, 3000 GC system, showing a significant reduction in error between a threshold of –1 and 1 for the learning model using GCL output (orange) compared to mossy fibers alone (blue). Transparent bounds represent standard deviation of learning outcome. Relative MSE of the GCL is shown on the right margin and represents the ratio of MSE for the GCL compared to MF alone. Values less than 1 indicate GCL outperforming MFs alone. **B**. An intuitive demonstration of the difference in the small MSE change produced by the MF-direct task, and the much clearer MSE produced by the GCL model used as input to P-cells. Panels show the outcomes of the same task with the target function being an image of a cat, with both handwritten and typeface text, and a 1- and 2-pixel width checkerboard (upper left corner).*



*Figure 2, Figure Supplement 1: Effects of input and output population sizes on learning.*
*A. Relation between the number of mossy fiber inputs and the resultant MSE, with MFs either inputted directly to P-cells (blue) or fed through 3000 GC unit model (red). **B.** MSE as a function of GC number compared to 50 MFs alone (blue). GC threshold fixed at 0 for these simulations.*
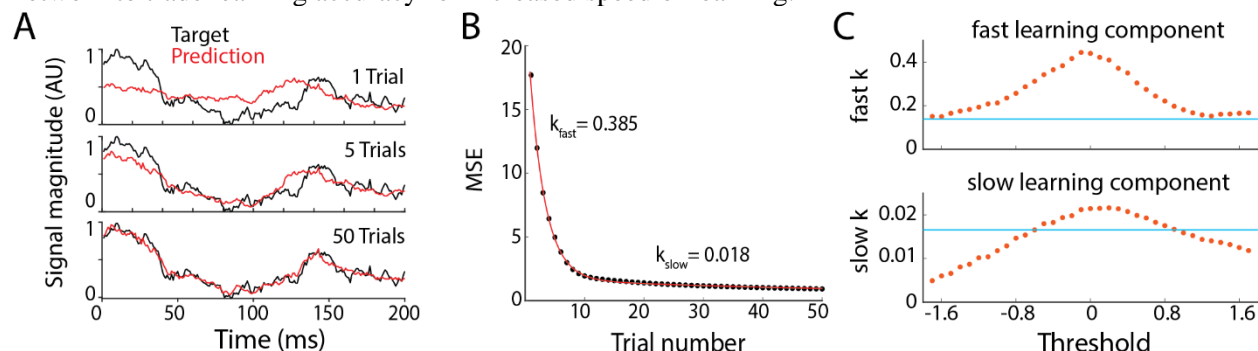
**GCL model speeds time series learning**

Having found that the GCL improves the match between predicted output and target output over a range of thresholds, we next examined whether the structure also increased the speed of convergence. Examining the MSE between output and target on each trial as training progresses (Fig. 3C, *red circles*), we found that output usually converged rapidly at first then more slowly in later stages of training (Fig. 3A). The reduction in MSE over training in our model was reasonably well fit by a double exponential (Fig. 3B, *red curve*), of the form

230
231 $$MSE(n) = A_1 \, e^{(-k_1 \, n)} + A_2 \, e^{(-k_2 \, n)} + C$$
232
233 where $n$ is the trial number. We measured the convergence speed of a simulation by the rate constants $k_1$
234 and $k_2$. In the vast majority cases, one of these rate constants was 5-50 times larger than the other; we
235 denote the larger constant $k_{fast}$ and the other $k_{slow}$. For most parameter values, $k_{fast}$ accounts for more than
236 80% of learning.
237
238 We next examined the influence of key parameters on convergence speed. First, we looked at the effect of
239 the GC threshold. Learning was fastest for GCL thresholds near zero (Fig. 3C, *red circles*), the level that
240 filters out half of the input received by a GC. Convergence in networks that lack a GCL (MFs directly
241 innervating P-cells) was consistently slower (Fig. 3C, *blue line*) than networks with a GCL. Convergence
242 can also be sped up by increasing the size of the parameter jumps in synaptic weight space during
243 gradient descent (the "step size"), but only to a limited degree (Supp. Fig. 3A). Indeed, at a GCL
244 threshold of 0, convergence speed *decreased* as the step size size was increased beyond ~$10^{-6}$ (au). We
245 speculated that this trade-off was a consequence of a failure to converge in a subset of simulations. To test
246 this, we looked at the fraction of simulations that converged towards a low MSE as a function of the
247 update magnitude. We found that the fraction of simulations that converged ("fraction successful")
248 decreased with increasing step size (Supp. Fig. 3B); in simulations that did not converge, the MSE
249 increased explosively and synaptic weights diverged. In such cases, we assume the large weight updates
250 made it impossible to descend the MSE gradient; each network weight update drastically changed the cost
251 function such that local MSE minima were overshot. When larger step sizes did permit convergence,
252 progress was nevertheless slowed, likely because the relatively large learning rates led to inefficient
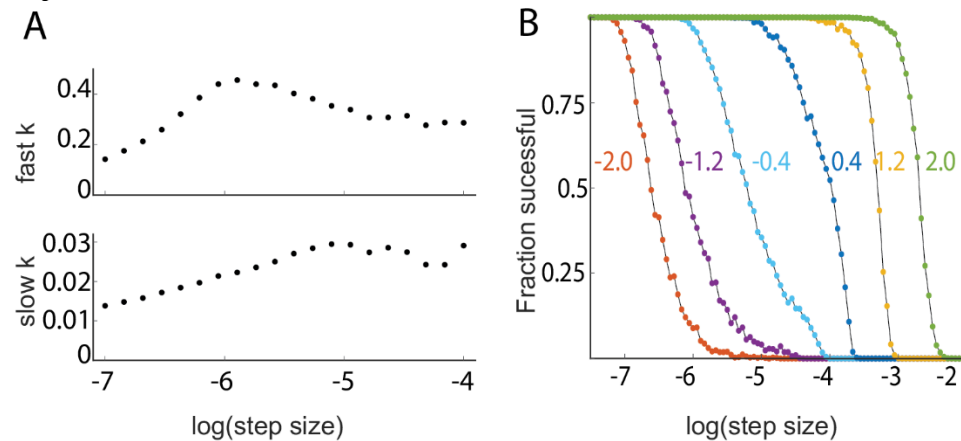253 progress towards the MSE minimum.
254
255 Although larger step sizes eventually cause learning to slow and then fail entirely at a given GCL
256 threshold, higher thresholds permitted larger step sizes before failures predominated (Supp. Fig. 3B).
257 Since higher thresholds permit larger step sizes before convergence failure sets in, convergence speed
258 might be maximized by jointly optimizing step size and GCL threshold. We tested this by systematically
259 raising step sizes at each threshold until convergence success fell to 50%. We defined the "maximum
260 convergence rate" for a given threshold as the maximum convergence rate (derived from fitting the MSE
261 trajectory with a double exponential) yielding successful convergence at least 50% of the time. We found
262 that the threshold giving the fastest convergence was indeed higher when step size was also optimized
263 (Supp. Fig. 3B) than when step size was fixed (Fig. 3C). Thus, increased GCL thresholding can allow the
264 network to trade learning accuracy for increased speed of learning.



265
266 *Figure 3. Learning speed increases with GCL.*
267 *A. Example of learned predictions after 1,5, and 50 trials of learning, with predictions in red and target*
268 *function in black. B. Example learning trajectory of MSE fit with a double exponential. Black circles:*
269 *MSE of network output on each trial. Red line: double exponential fit MSE during learning. Here, step*
270 *size was $10^{-6}$ and z-scored GCL threshold was 0. We use the exponents k from the exponential fit to*

271  *measure learning speed. **C**. Learning speed as a function of GCL threshold (red dots). Blue line: learning*
272  *speed in networks lacking GCL, i.e. mossy fibers directly innervate output Purkinje unit, gradient descent*
273  *step size was $10^{-6}$.*



274
275  ***Figure 3, figure supplement 1: Effects of gradient descent step size on learning speed.***
276  ***A**. Learning speeds (exponential time constant) for different simulations using varying gradient descent*
277  *step sizes, showing differentially maximized learning speeds occurring at different step sizes. **B**. Fraction*
278  *of simulations that converge to asymptotic MSE values as a function of gradient descent step sizes for*
279  *different values of GCL threshold (colors denote threshold values). Note that larger step sizes and faster*
280  *learning are supported in models with higher thresholds.*

281
282
283
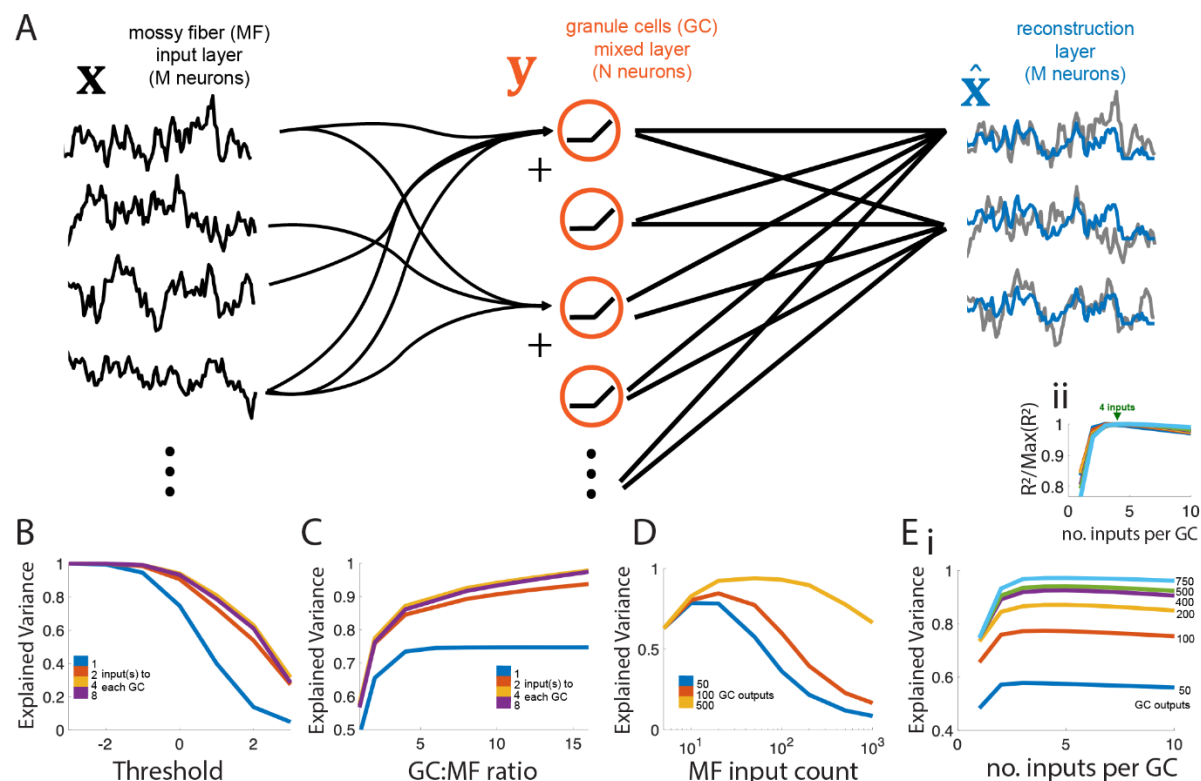284  **Recovering GCL input from GCL output**
285  Having established a framework for studying GCL processing of naturalistic inputs, we wanted to
286  understand to what extent thresholding GCL activity led to the loss of information supplied by MF inputs,
287  which potentially contains useful features for learning. In other words, would Purkinje neurons be
288  deprived of behaviorally relevant mossy fiber information if these inputs are severely filtered by the
289  GCL?  To assess this issue, we used a metric of information preservation called *explained variance*,
290  (Achen 1982); however, in this special case, we use the term '*variance retained*', because this metric
291  represents the preservation of information about the input after being subjected to filtering in the GCL
292  layer. Let $x_t$ denote the MF input at time *t*. If the GCL activity preserves the information present in $x_t$, then
293  it should be possible to reconstruct the activity of MFs from GCL activity (see Methods for details on
294  how this reconstruction was performed). The variance retained is then the mean squared error between the
295  actual MF input $x_t$ and the reconstructed input, normalized by the MF input variance:
296

297  $$R^2 = 1 - \frac{\sum_{t=1}^{T} \langle (\hat{x}_t - x_t)^2 \rangle}{\sum_{t=1}^{T} Var[x_t]}$$

298
299  Our primary finding is that the GCL transmits nearly all of the information present in the MF inputs even
300  at fairly high thresholds, but only if the GCL is sufficiently large relative to the MF population. The
301  threshold, feedforward architecture, and relative balance of MF inputs and GC outputs all affect the
302  quality of the reconstruction. Variance retained by the reconstruction layer decreased with the GC layer
303  threshold, since it masked some subthreshold input values (Fig. 4B). Allowing more MF inputs per GC
304  recovered some of this masked information, since some subthreshold values are revealed through
305  summing with sufficiently suprathreshold values. However, these gains cease beyond a few MF inputs
306  per GC, since the exponential growth of MF combinations rapidly exceeds the number that the GCs can
307  represent (Marr 1969; Gilmer and Person 2017).

308
309 To disentangle the information contained in the summed inputs, many different combinations of inputs
310 must be represented to disambiguate the contributions of each MF input. Increasing the number of GCs
311 generally increases the variance retained, since more combinations of MF inputs are represented, and
312 reveal subthreshold input values (Fig. 4C). Interestingly, variance retained by the network varied non-
313 monotonically with the number of MF inputs (M) when the number of GCs (N) was fixed. This is because
314 having too few MF inputs means there may not be a sufficient number of combinations so that
315 subthreshold values can be revealed (by summing them with suprathreshold inputs) but having too many
316 saturates the information load of the GC layer (Fig. 4D). Lastly, when fixing the number of MF inputs
317 and GCs, there is an optimal number of MF inputs to each GC, which aligns with the anatomical
318 convergence factor of 4 MF/GC (Fig. 4E), related to previous findings that suggest the best way to
319 maximize dimensionality in the GC output layer is to provide sparse input from the mossy fibers (Litwin-
320 Kumar et al., 2017; Cayco-Gajic et al., 2017). Thus, there are two key features that shape the information
321 transferred to the GCL from the MF inputs. First, the way in which MF inputs are combined to form the
322 total input to each GC determines how much information about subthreshold inputs can be transferred
323 through the nonlinearity. Second, the total number of GC outputs determines how many MF input
324 combinations can be represented, so that, ultimately, the random sums of MFs can be disentangled by the
325 downstream reconstruction layer. Together, information transfer requires a combined summation and
326 downstream decorrelation process accomplished by the three layer network.
327



328
329 *Figure 4: Recovering inputs with an optimal linear readout.*
330 *A. Network model schematic. Granule cell (GC, red, center) layer thresholds the sum of (4 here)*
331 *randomly chosen mossy fiber (MF, black, left) inputs, which are then fed into a reconstruction layer*
332 *which uses the optimal weighting from all N GCs to approximate each of the M inputs (compare blue*
333 *readouts to grey inputs). B. Increasing the threshold of the GC layer (N=500 outputs) decreases the*

334     *explained variance (i.e. variance retained) of the best reconstruction layer (M=50), but the effect is*
335     *reduced with an intermediate number of MF inputs per GC. **C.** Variance retained increases with the ratio*
336     *of GCs per MF but gains from increasing the number of inputs to each GC are limited (max at 4 inputs).*
337     *Here there are M=50 MF inputs at the threshold = 0. **D.** For a fixed number of GC outputs N, there is an*
338     *optimal number of MF inputs (M) for which the variance retained of the reconstruction layer is*
339     *maximized. **E. i.** For a fixed number of GC outputs N and MF inputs M=50, there is an optimal number of*
340     *inputs per G (around 4) for maximizing variance retained. **ii.** Same as i, but with each value normalized*
341     *to its maximum to show maximized values at inputs = 4.*
342

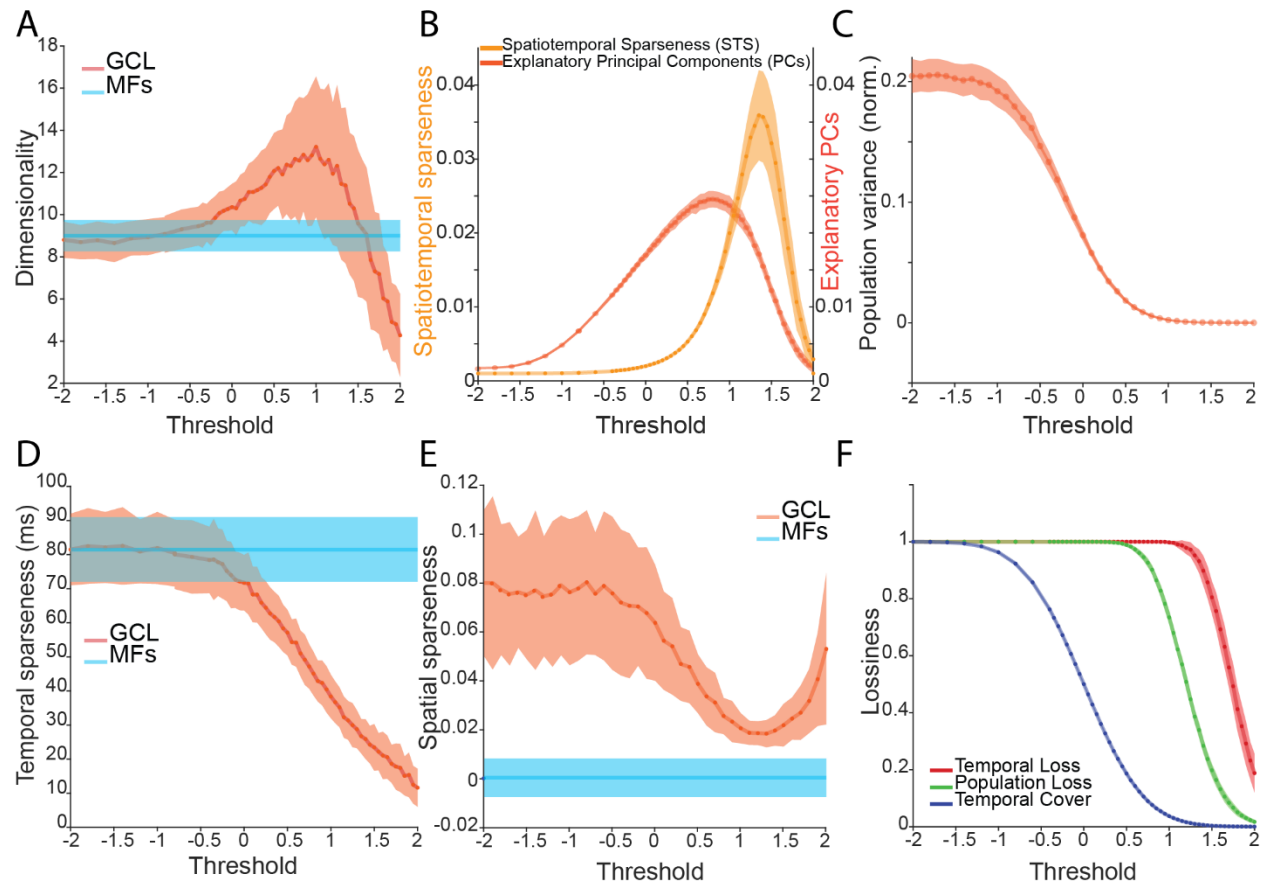343     **General statistical features of GCL population activity**
344     To better relate the present model to previous theoretical studies we looked at a variety of population
345     metrics to help explain how signal filtering by the GCL improves cerebellar learning and why it
346     ultimately fails as the GC threshold is increased.
347

348     The first set of metrics related to pattern separation: dimensionality (Dim), the number of explanatory
349     principal components (PCs), spatiotemporal sparseness (STS), and population variability (See methods
350     for details). (Although STS is a measure of sparseness, it represents an idealized form of separability
351     where GCs represent unique temporal epochs that do not overlap, providing a perfect basis set when
352     maximized, thus is grouped with pattern separation metrics). Dim, PCs, and STS showed non-monotonic
353     relationships with threshold and peaked at thresholds ranging between 0.5 and 1.5 (Fig. 5 A, B), while
354     population variability decreased with increasing thresholds (Fig. 5C). Intuitively, this relationship
355     captures the effect of low thresholds allowing GC activity to relay the mean input, with no pattern
356     separation occurring, and thus minimizing pattern separation metrics. With increasing threshold, GC
357     activity is driven by coincidence detection leading to high dimensional population output. At high
358     thresholds, inputs rarely summate to threshold, leading to lost representation that drives a roll-off in
359     pattern separation within the population. Notably, Dim, PCs, and STS peaked at thresholds greater than
360     peak learning performance, which was optimized at threshold zero, thus none of these three pattern
361     separation metrics alone account for learning performance. Population variability (i.e. GCL variance per
362     unit) is thought to aid classification and separability of GCL output (Cayco-Gajic et al., 2017). This
363     metric's decrease with increasing threshold was likely due to the decrease in overall representation by
364     each unit due to sparsening and diminishing the dynamic range of GC rates due to threshold subtraction
365     (Fig. 1, Fig. 5C).
366

367     The second set of metrics are related to sparse representations: temporal sparseness and spatial
368     sparseness. Temporal sparseness – defined by the exponential decay of GC autocovariance, where smaller
369     values typify signals that change quickly with time -- decreased as a function of threshold because of
370     sparsened representation at higher thresholds (Fig. 5D). The mean pairwise GC correlation, (Fig. 5E) i.e.
371     spatial sparseness, shared a drop-off after a threshold of 0, but increased again at high thresholds because
372     only a few MF signals were retained at high threshold and thus were highly correlated. By experimental
373     design, decorrelation was already maximized in OU inputs. Similar to the pattern separation metrics,
374     these sparseness metrics did not show an obvious relationship to the U-shaped learning performance seen
375     in Fig. 2A.
376

377     Finally, we examined three metrics of lossiness defined to quantify (1) the fraction of the total epoch with
378     no activity in any GC unit (e.g. with "temporal lossiness" of 0.1, 10% of the total epoch has no activity in
379     any GCs) (2) the proportion of granule cells with any activity over the entire epoch ("population
380     lossiness") (3) the mean fraction of the epoch in which each granule cell is active ("temporal cover"). Not
381     surprisingly, each lossiness metric increased with high thresholds (Fig. 5F). However, despite diminishing
382     activity in individual GCs with increasing threshold, (the blue curve Fig. 5F), each GC was resistant to
383     becoming completely silent (green curve drop, Fig. 5F), owing to a few dominant inputs.
384

385    Notably, none of these metrics alone obviously tracked the U-shaped learning performance (Fig. 2A).
386    However, collectively, these descriptive statistics of model GCL population activity set the stage for
387    analyzing how information preprocessing by the basic GCL architecture relates to learning time series,
388    explored below.
389



390
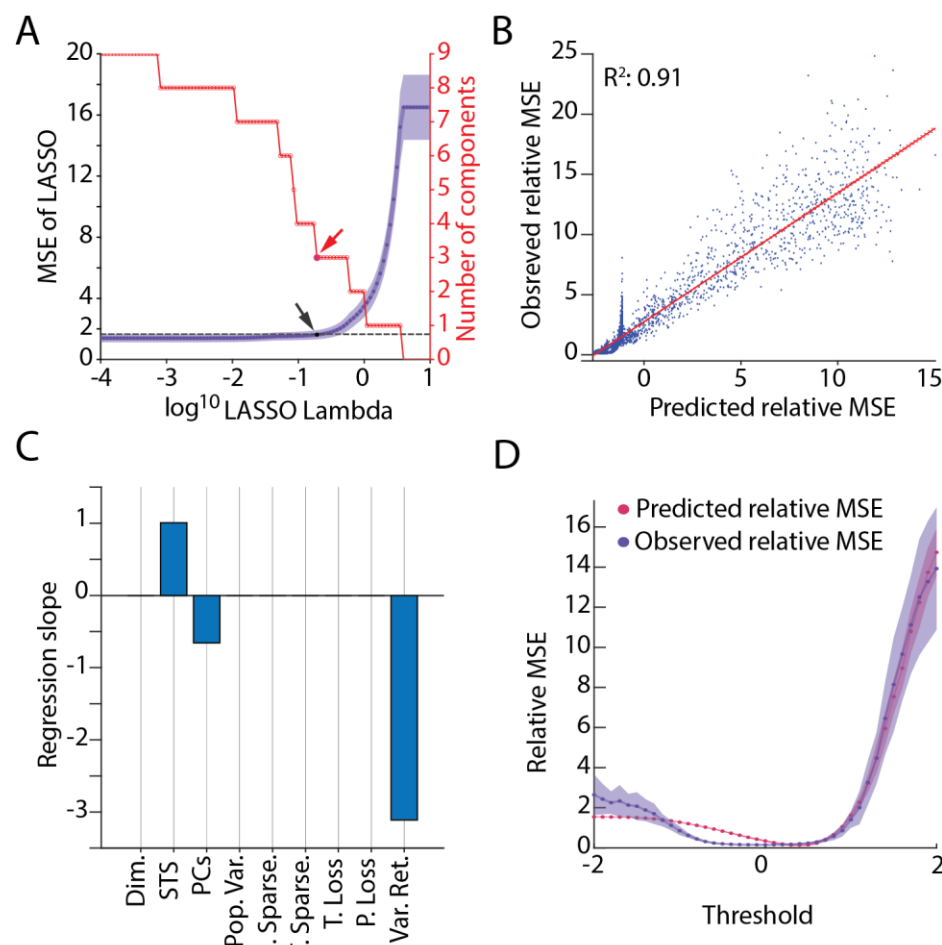391    *Figure 5: Statistical features of GCL output.*
392    *A. GCL dimensionality (red) and MF dimensionality (blue) as a function of threshold. Note peak near a*
393    *threshold of 1 for the GCL. B. Two metrics of pattern separation in GCL output -- STS (light orange) and*
394    *PCs (dark orange) -- as a function of threshold. Note peaks near 1.5 and 0.5, respectively. C. The sum of*
395    *GCL variance produced by the model as a function of threshold. Note monotonic decrease with threshold.*
396    *D. Temporal sparseness as a function of thresholding. Note monotonic decrease in GCL with*
397    *thresholding. E. Mean pairwise correlation of the population plotted as a function of threshold. Note*
398    *trough near 1. F: Three forms of lossiness in GCL output as a function of threshold. Each metric had*
399    *differential sensitivity to thresholding but note that all decrease with increasing threshold. Across*
400    *metrics, function maxima and minima ranged widely and were not obviously related to thresholds of*
401    *optimized learning.*
402
403

**Optimization of learning through GCL transformations**

405    With the knowledge that thresholding drives changes both in learning time series (Fig. 2, 3) and GC
406    population metrics that are theorized to modulate learning (Fig. 4, 5), we next directly investigated the
407    relationships of these metrics to learning performance. To test this, we used a LASSO regression method
408    to identify learning performance-driving variables taken from the metrics described in Figures 4 and 5
409    (Fig. 6A, C).  Using the output of the LASSO model, we found that a three-term model using the most
410    explanatory variables -- STS, the number of explanatory PCs and variance retained (Fig. 6B, C, D) --

411  accounted for 91% of learning variance. The three-term model performance is plotted against the
412  observed MSE over a range of thresholds in Fig. 6D, showing strong similarity.
413
414  These results were somewhat surprising given prior studies showing benefits of population sparseness or
415  decorrelation to learning. To interrogate this seeming disparity, we introduced fictive GCL population
416  activity that had specific statistical features as inputs to P-cells. Consistent with previous reports,
417  decorrelation and temporal sparseness improved learning accuracy, with complete decorrelation and
418  temporally sparse supporting the best performance (Fig. 6 - figure supplement 1; Cayco-Gaijic et al.,
419  2018). Thus, on their own, population, temporal and idealized spatiotemporal sparseness do modulate
420  learning when their contribution is independent, but these features nevertheless do not emerge as features
421  in the naturalistic GCL model as statistical properties that drive performance of time series. This property
422  is a consequence of temporal sparseness and decorrelation covarying with lossiness (captured by the
423  variance retained metric), which drives down performance. Rather, the statistical features produced by the
424  model GCL with naturalistic inputs that best explain learning are the number of explanatory PCs, STS,
425  and the amount of input variance retained -- metrics that may align well with recently described GC
426  population activity during locomotion (Lanore et al., 2021).
427



428
429  *Figure 6. Relationship between sparseness metrics and MSE.*
430  *A. LASSO regression model selection as a function of progression of the Lambda parameter (penalty*
431  *applied to regressor selection). The removal of regressors with increasing Lambda (red steps) selected*
432  *from the following potential regressors: dimensionality (Dim.), spatiotemporal sparseness (STS),*
433  *explanatory principal components of the GC population (PCs), population variability (Pop. Var.), spatial*
434  *sparseness (S. Sparse.), temporal sparseness (T. Sparse.), temporal lossiness (T. Loss.), population*

435 *lossiness (P. Loss), and input variance retained (Var. Ret; Figure 4). Arrow shows selection point of*
436 *LASSO regression MSE using "1SE" (1 standard error) method (see Methods, purple lines, black dot and*
437 *arrow indicating the selected model, with red arrow showing selection point in the parameter reduction*
438 *plot, red).* ***B.*** *Relationship between LASSO model (predicted relative MSE) against the observed relative*
439 *MSE (ratio of GC MSE to MF alone MSE) with fit line and variance explained by regression ($R^2 = 0.91$)*
440 ***C.*** *Regression slopes of the selected LASSO model from A, showing that STS, PCs, and Input Variance*
441 *Retained are the selected regressors, with Var. Ret. being the largest contributing factor. All factors*
442 *normalized to a normal distribution for comparison.* ***D.*** *The output of the selected model and the observed*
443 *MSE plotted against threshold for a comparison of fits, demonstrating high accuracy in the 0-2 range, but*
444 *less accuracy in the -2-0 range.*

445



446
447 ***Figure 6 figure supplement 1: GC population statistics regulate learning accuracy when independently***
448 ***controlled.*** *Fictive population activity with structured statistics were introduced to P-cells to explicitly*
449 *test the roles of population decorrelation and structured spatiotemporal sparseness on learning.* ***A-B****:*
450 *Learning performance (MSE) as a function of temporal sparseness (i.e. autocovariance tau) or spatial*
451 *sparseness (i.e. population correlation). Red dots on A and B indicate values used for input model to GCL*
452 *in Figs 2, 5, and 6.* ***C****: Matrix of effects on MSE when modulating temporal spareness via tau, and spatial*
453 *sparseness via population correlation. Lower values for both (cooler colors) indicate the best learning*
454 *accuracy.* ***D****: The results of these analyses support the idea that GCL filtering benefits learning through*
455 *transformation of statistical structure fed to the P-cell. A remaining caveat was that the number of*
456 *granule cells far exceeded the number mossy fibers, raising the question of whether the learning*
457 *advantage conferred by the GCL is merely a consequence of this difference. To test this, we fed MFs*
458 *directly to the P-cell units and varied their numbers between ranges of 2 to 3000. While learning*
459 *accuracy improved with more MFs, asymptotic MSE values were lower than the GCL, indicating that the*
460 *filtering properties of GCL are indeed important for this learning task. Figure plots the MSE as a*
461 *function of the number of inputs to Purkinje cells, showing that too few MFs are insufficient to produce*
462 *accurate learning, but having a large number makes little difference beyond $10^{1.5} \sim= 31$ MFs.* ***E****: To test*
463 *how the uniqueness of individual unit activity across time contributes to learning we selected population*

*activity that varied in STS from a bank of simulations. Two distribution structures were tested. The first maximized granule cell uniqueness in time and temporal organization – e.g. each granule cell is active only once during the epoch and only one granule cell is active at a given time, such that the population histograms resemble a 'staircase' ("ideal STS basis"). Overlap of active granule cells drives decreases in computed STS, or wider steps in the staircase ("temporally linked overlap"). The second class of STS maximized uniqueness without requiring temporal organization – e.g. any slice of time is unique, but an individual granule cell can occupy an arbitrary number of time bins ("stochastic overlap"). STS drops when a given granule cell activity occupies more time bins, reducing the uniqueness of the granule cells contribution to the population. Figure shows schematic diagram of these different types of spatiotemporal sparseness, with structured overlap "temporal overlap" and "stochastic overlap" illustrating different ways populations could differ. **F:** Effect of STS on MSE, where overlap between units is always local to a particular time point, so that units are only active at a particular continuous temporal range, showing a monotonic decrease in error as STS approaches 1. **G:** Same as F, but the temporal location of overlap between units is random, showing best learning accuracy at STS = 1, and good but less accurate learning at STS = 0. When overlap was decoupled from time in the stochastic overlap case, error was reduced at both maximal and minimal STS simulations with the highest error occurring at intermediate STS values. This may be because the gradient descent algorithm is able to use dense, variable signals, like those seen in very low STS value GCL outputs, to learn essentially as well as the high STS values which have strong isolation in individual unit representation and are guaranteed to be good for learning.*


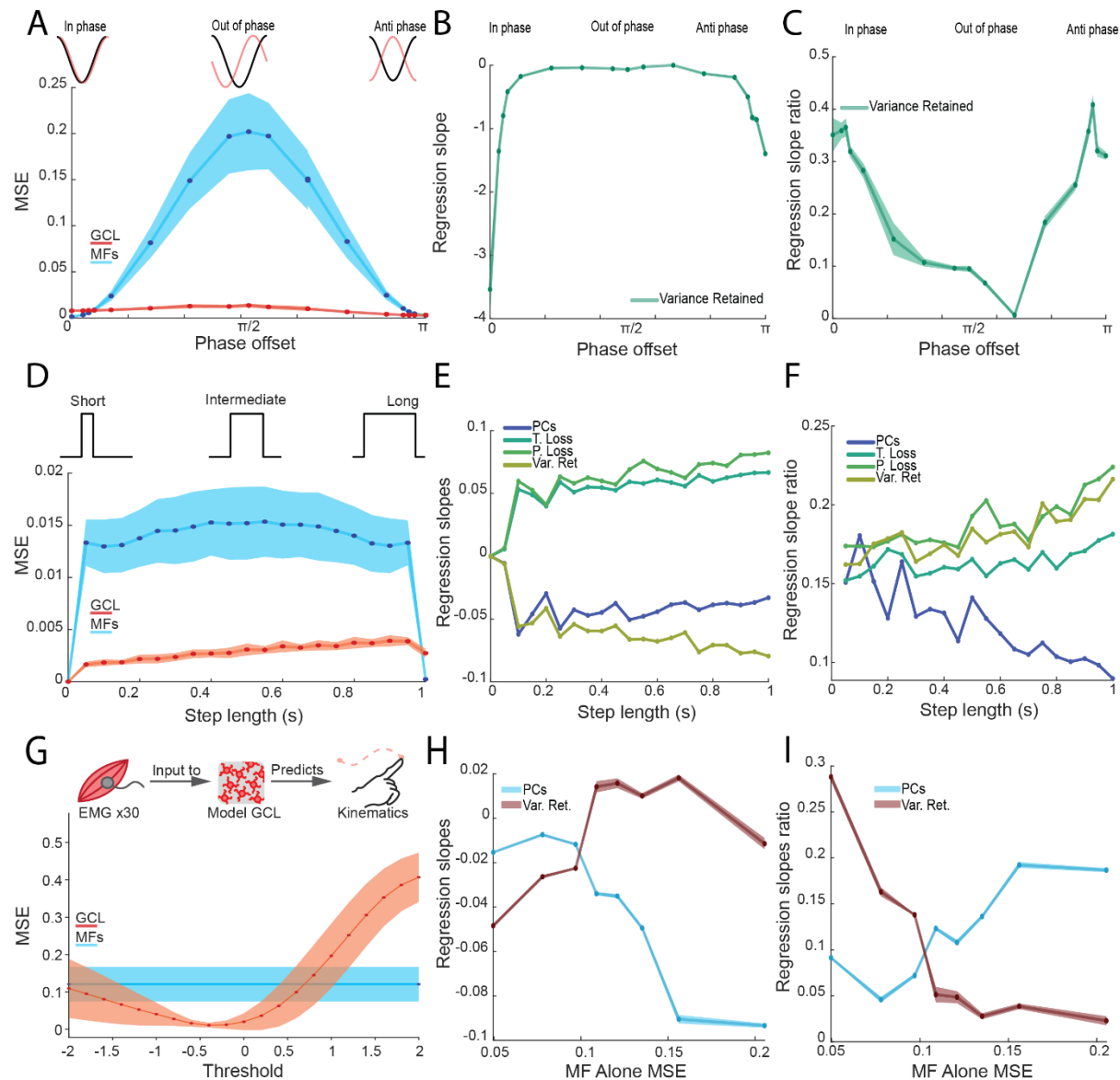**GCL properties that enhance learning in naturalistic tasks**
Together, these models suggest that the GCL can reformat random inputs suitable to support rapid and accurate learning of time-series. The real cerebellum is topographically organized along multiple parasagittal output modules (Apps and Garwicz, 2005; De Zeeuw, 2020). This organization suggests segregated afferents with specific statistical structure could refine specific behaviors. To examine whether different population statistical features might support distinct learning tasks, we utilized the model to perform a series of naturalistic cerebellar tasks: vestibulo-ocular reflex (VOR) phase adaptation (Ito et al. 1974), temporal interval learning (Narain et al., 2018) and kinematic encoding (Herzfeld et al., 2015).

We speculated that the nature of these tasks might influence the contribution of components of the model to learning accuracy. For example, when VOR is kept in phase, it makes intuitive sense that retention of vestibular input, inherently in-phase with the motor output, would be valuable, with reweighting of GC representations of inputs giving rise to amplitude learning as in VOR gain adaptation. However, if the phase is offset, the relationship between vestibular input and ocular output requires complex mapping (Fig. 7A, top middle inset) and selection of GCs representing sparsened OU processes may be selected instead to allow for reconstruction from high-dimensional outputs. The GCL model supported learning of VOR at all phases, but MFs showed especially poor performance in pi/2 phase shifts (Fig. 7A, 'out-of-phase'). As a result of this reliance on GCL reformatting, we predicted that the contribution of 'variance retained' to learning should decrease depending on the phase shift. In other words, the extent to which the input was inherently related to the output would be of scalable importance. We tested the relationship of input variance retention and phase offset using RIDGE regression (which preserves even small contributions of regressor variables to the model in comparison to LASSO) and found that for in-phase and anti-phase learning input variance retention accounted for most of learning, reflected in large slope coefficients, whereas input retention decreased as an important variable in out-of-phase learning, with shallow slope coefficients (Fig. 7B). Furthermore, the relative magnitude of the slope magnitude of variance retrained is reduced in out-of-phase conditions compared to in-phase and anti-phase (Fig. 7C). This suggests that the learning rule can utilize information preserved by the GCL, as in in-phase learning, but, if necessary, it can learn using information that is so highly reformatted that it no longer retains the original vestibular information.

515     Cerebellar timing tasks, such as delay eyeblink conditioning, involve time estimation over intervals
516     spanning 100-500 *ms*. Models of delay eyelid conditioning suggested that the cerebellum represents the
517     time interval through decomposition of an invariant signal into many signals that tile across time. This
518     hypothesis provides an interesting test of whether lossiness differentially affects behavioral outcomes
519     depending on whether short or long intervals are being estimated, defined here as the duration of an 'on'
520     signal. If we assume that an ideal temporal basis (akin to the 'staircase' representation in Fig 6 - figure
521     supplement 1E) represents different points in time of the stimulus, one might speculate that lossiness in
522     populations representing long intervals would be more detrimental than in populations representing short
523     intervals -- given that only the temporally aligned subsection of the input is relevant to the output
524     response and the rest is discarded or ignored. We tested this prediction by systematically altering the
525     length of a step target function to occupy 0% to 100% of the response epoch using OU processes as
526     inputs. The model using a GCL was able to perform this task more accurately than with MF inputs alone
527     (Fig 7D), and the magnitude of slope for lossiness-related metrics increased with interval duration (Fig.
528     7E), suggesting that learning short intervals is less sensitive to lossiness than learning long intervals. The
529     relative contribution of lossiness metrics to the overall regression performance also increased with step
530     duration compared to PCs (Fig. 7F), suggesting that lossiness-related metrics have a more powerful
531     influence on learning outcomes as a function of increasing duration that is not true of pattern separation
532     metrics like PCs.
533

534     We next asked whether naturalistic input statistics, derived from electromyogram (EMG) signaling, could
535     support learning. We used EMG signals from human subjects in a point-to-point reaching task as MF
536     inputs, and tested whether the model could learn associated limb kinematics from this input (Fig. 7G;
537     Delis et al. 2018; Tseng et al. 2007; Miall and Wolpert 1996; Wolpert et al., 1998). The GCL was able to
538     produce more accurate predictions of the kinematics when compared to the EMG as MF inputs alone, and
539     the range of thresholds which produced the best accuracy were comparable to the previous findings (Fig.
540     3A), but were slightly negatively shifted, suggesting retained variance of inputs might be beneficial to
541     learning kinematics from associated muscle activity.
542

543     Finally, since EMGs used as MF inputs to the model had some level of baseline utility in predicting
544     kinematics based on their intrinsic relationships, (reflected in MFs alone MSE varying between 0.04 and
545     0.22), we next asked whether this influenced which features of the GCL output were most related to
546     learning. In keeping with intuition, when MF based learning was excellent (low MSE), the slope of the
547     variance retained metric was highest (Fig. 7H, I, blue). Conversely, when MF based learning was poor
548     (high MSE) variance retained slopes dropped. Interestingly, a few GCL population metrics became more
549     important for learning as MF MSE worsened, such as the number of explanatory PCs (Fig. 7H, I,
550     maroon). Together this suggests that different pattern separation features of GCL reformatting may serve
551     learning under different conditions, with Purkinje cells using diverse 'pattern separation' features
552     depending on the task and input statistics. When intrinsic relationships are valuable, variance retained is
553     an important population statistical feature; when they are more arbitrary, pattern separation features are
554     more valuable for learning relationships between the inputs and output. This shifting landscape was a
555     general feature of our models (Fig. 6 & 7), suggesting that "pattern separation" by the GCL is not one
556     universal transform that has broad utility. This observation raises the possibility that regional circuit
557     specializations within the cerebellar cortex, such as density of unipolar brush cells (Dino et al. 2000),
558     Golgi cells, or neuromodulators could bias GCL information reformatting to be more suitable for learning
559     of different tasks.
560

561

*Figure 7. Task-dependent relationships between granule cell population statistics and learning.*

*A. Task structure of a phase-offset VOR-like task (top) and learning performance as a function of phase-offset for GCL and MFs alone (bottom). Here, the phase between the input function and the target function varies between 0 and pi. GCL (red) or MFs alone (blue) were used as inputs to learn the task. As the difference in phase between inputs and targets approaches pi/2 (out of phase), performance from MF alone degrades while GCL performance remains accurate and stable. B. RIDGE regression slopes of the input variance retained (Var. Ret.) metric as a function of phase offset. Variance retained slope is large when phase offset is in the 'in phase' and 'anti-phase' regions of the task, but is otherwise minimized, suggesting that the utility of this statistical feature varies depending on task. C. Same data as B but normalized to show the relative proportion of all slope magnitudes accounted for by Var. Ret. (slope magnitude of Var. Ret divided by the sum of all slope magnitudes). Var. Ret. is a primary regressor for 'in phase' or 'anti-phase' learning. D. Task structure (top) and learning performance (bottom) of an interval estimation task, where the model is tasked with learning a step function that varies in length. GCs (red) and MFs alone (blue) were used as inputs to the P-cell. As the interval lengthens, learning using MFs alone was generally poorer than using the GCL. E. RIDGE regression slopes of 4 variables (Var. Ret., T. Loss, P. Loss, PCs) as a function of step length, showing that slopes of lossiness-related metrics (P. & T.*

578     *Loss, and Var. Ret.) increase in magnitude as the step length increases, whereas slope magnitude of PCs*
579     *decreases. **F**. Same as E but showing the relative proportion of all slopes accounted for by these 4*
580     *regressors. **G**. Schematic of underlying dataset using recorded EMG as an input to the model GCL to*
581     *predict kinematics (top). Learning performance of model using EMG alone (MFs; blue) or GCL (red)*
582     *across varying thresholds. The GCL outperforms MFs alone at a threshold range similar to that observed*
583     *in Fig. 2. **H**. RIDGE regression slopes of Var. Ret and PCs metrics as a function of the learning*
584     *performance achieved by using MFs (i.e. EMG) alone, showing that Var. Ret. is a stronger driver of*
585     *performance when MFs alone supported accurate learning, but not when MFs alone supported poor*
586     *learning (higher MSE). PCs show the opposite trend, increasing in slope magnitude when MFs alone*
587     *supported poor learning. **I**. Same as H, but showing the relative proportion of slope magnitude accounted*
588     *for by Var. Ret. and PCs.*
589
590

591     **Discussion**
592     Here we asked a simple question: how does the cerebellar granule layer support temporal learning? This
593     question has captivated theorists for decades, leading to a hypothesis of cerebellar learning that posits that
594     the GCL reformats information to best suit associative learning in Purkinje cells. Recent work has called
595     many of these foundational ideas into question, however, including whether GCL activity is sparse; high
596     dimensional; and what properties of 'pattern separation' best support learning (Wagner et al., 2017;
597     Giovannucci et al., 2017; Knogler et al., 2017; Cayco-Gajic et al., 2017; Gilmer and Person 2017). To
598     reconcile empirical observations with theory, we hypothesized that input statistics and task structures
599     influence how the GCL supports learning. Here, we used naturalistic time-varying inputs to a model GCL
600     and identified pattern separation features that supported learning a time series prediction task, with an
601     arbitrary but temporally linked input-output mapping, recapitulating important features of physiological
602     cerebellar learning tasks. This formulation eliminates the possibility of trivial dimensionality changes
603     improving classification performance, thus approaching naturalistic challenges faced by the real circuit.
604     Several important observations stemmed from these simulations: (1) with naturalistic input statistics, the
605     GCL produces temporal basis sets akin to those hypothesized to support learned timing with minimal
606     assumptions; (2) this reformatting is highly beneficial to learning; (3) maximal pattern separation does not
607     support the best learning; (4) rather, tradeoffs between loss of information and reformatting favored best
608     learning at intermediate network thresholds; and finally (5) different "cerebellar" tasks utilized different
609     GCL population statistical features to optimize performance. Together these findings provide insight into
610     the granule cell layer as performing pattern separation of inputs that transform information valuable for
611     gradient descent-like algorithms (akin to Purkinje cell learning rules), but with idiosyncratic population
612     statistics supporting different tasks. This observation makes predictions about the regional specifications
613     that occur across the layer that may specially subserve diverse behaviors.
614

615     *Emergence of spatiotemporal representation and contribution to learning*
616     A perennial question in cerebellar physiology is how the granule cell layer produces temporally varied
617     outputs that could support learned timing (Mauk and Buonomano 2004). While cellular and synaptic
618     properties have been shown to contribute (Chabrol et al., 2015; Duguid et al, 2012; Guo et al., 2021;
619     Crowley et al., 2009; Rudolph et al., 2015; Buonomano and Mauk 1994; Kanichay and Silver 2008;
620     Simat et al., 2007; Mapelli et al., 2009; Rossi et al., 1996; Gall et al., 2005; Armano et al., 2000; Rizwan
621     et al. 2016; Tabuchi et al., 2019; D'Angelo and De Zeeuw 2009), we observed that with naturalistic
622     inputs, temporal basis set formation is a robust emergent property of the threshold-linear input-output
623     function of granule cells receiving multiple independent time-varying inputs (Fig. 1B). But is this
624     reformatting beneficial to learning? We addressed this question by comparing learning of a complex time-
625     series in model Purkinje cells receiving either mossy fibers alone or reformatted output from the GCL.
626     We found that indeed the GCL outperformed MFs alone in all tasks (Figs. 2, 3, 7). Nevertheless, we
627     wondered what features of the population activity accounted for this improved learning. While
628     sparseness, decorrelation, dimensionality and lossless encoding have been put forward as preprocessing

629 steps supporting learning, we found that none of these alone accounted for the goodness of model
630 performance. Rather, disparate pattern separation metrics appear to strike a balance between maximizing
631 sparsenesses without trespassing into lossy encoding space that severely, and necessarily, degrades
632 learning of time-series.

634 Moreover, the value of different population metrics to learning varied with the specific task -- with some
635 tasks relying on input retention for best performance, others relying on absence of lossiness, and some
636 requiring pattern separation to accomplish accurate predictions (Fig. 7). For instance, when input statistics
637 are well suited to learning the specific task, as in in-phase VOR, preservation of input variance drives best
638 performance (Fig. 7A-C). Importantly, although the properties of the GCL selected to improve learning
639 varied across tasks, the underlying architecture of the GCL and thresholding did not. This suggests that
640 the output of the GCL is well structured to support a variety of tasks. Thus, Purkinje cells are able to
641 make use of a spectrum of information formats that, depending on task requirements, are selected to serve
642 best learning.

644 These observations are interesting in light of a long history of work on granule layer function. Marr,
645 Albus, and others proposed that the granule cell layer performs pattern separation useful for classification
646 tasks. In this framework, sparseness is the key driver of performance, and could account for the vast
647 number of granule cells. Nevertheless, large-scale GCL recordings unexpectedly showed high levels of
648 correlation and relatively non-sparse activity (Wagner et al., 2017; Giovannucci et al., 2017; Knogler et
649 al., 2017). Despite methodological caveats, alternate recording methods seem to support the general
650 conclusion that sparseness is not as high as originally thought (Lanore et al. 2021; Kita et al., 2021;
651 Gurgani and Silver 2021). Indeed, subsequent theoretical work showed that sparseness has deleterious
652 properties (Cayco-Gajic et al., 2017; Billings et al., 2014), also observed in the present study, that may
653 explain dense firing patterns seen *in vivo*. Here we found that the best learning occurred when individual
654 granule cell activity occupied around half of the observed epoch (Fig. 5F, blue trace), achieved with
655 intermediate thresholding levels. We also observed temporal organization that is consistent with the firing
656 patterns observed *in vivo*. While these findings seem to suggest that sparseness is not the 'goal' of GCL
657 processing, our findings and others (Litwin-Kumar et al., 2016; Cayco-Gajic et al., 2017) suggest that
658 pattern separation broadly is a positive modulator of GCL support of learning processes.

660 Previous work (Sanger et al., 2020) proposed that time-series prediction was possible with access to a
661 diverse set of geometric functions represented in the GC population. However, that study left open the
662 question of how such a diverse collection of basis functions would emerge. The GCL model used here
663 minimized free parameters by incorporating very few independent circuit elements, suggesting that a
664 single transform is sufficient to produce a basis set which is universally able to learn arbitrary target
665 functions. We used a simple threshold-linear filter with a singular global threshold that relied on sparse-
666 sampling to produce spatiotemporally varied population outputs. This simple function worked to support
667 learning at a broad range of inputs and thresholding values, ultimately allowing the Purkinje cells
668 downstream to associate the spatiotemporally sparser inputs with feedback to learn arbitrary, and often
669 quite difficult, target functions. The emergence of this basis set is remarkable given the very simple
670 assumptions applied, but is also physiologically realistic, given the simple and well characterized
671 anatomical properties of the MF divergence and convergence patterns onto GCs, which are among the
672 simplest neurons in the brain (Jakab and Hamori, 1988; Palay and Chan-Palay, 1974; Palkovits et al.,
673 1971). Although we suggest that the key regulator of thresholding in the system is the feedforward
674 inhibition from Golgi cells, many factors may regulate the transformation between input and GC output in
675 the network, allowing for multiple levels and degrees of control over the tuning of the filter or real
676 mechanism that controls the outcomes of GCL transformations. Golgi cell dynamics may prove critical
677 for enforcing the balance between pattern separation metrics and lossy encoding (Hull 2020) thus are
678 critical players in mean thresholding found here to optimize learning. Additional mechanistic
679 considerations may also play a role, including short-term synaptic plasticity (Chabrol et al. 2015) and

680   network recurrence (Gao et al. 2016; Houck and Person 2014; 2015; Judd et al., 2021), allowing for a
681   more nuanced and dynamic regulatory system than the one shown here.
682
683   *Recapturing input information in the filtered GCL output*
684   Two schools of thought surround what information is relayed to Purkinje cells through GCs. Work in the
685   oculomotor cerebellum and flocculus suggests that Purkinje cells inherit virtually untransformed
686   information encoding eye velocity and visual motion, integrated in P-cells as positional signals (Herzfeld
687   et al., 2020; Krauzlis and Lisberger, 1991). Alternatively, the implication of theories of Marr and Albus
688   suggest that input information is so sparsened that Purkinje cells receive only a small remnant of the
689   sensorimotor information sent to the cerebellum. These divergent views have never been reconciled to our
690   knowledge. We addressed this disconnect by determining the fraction of MF input variance recoverable in
691   GCL output. Interestingly, the GCL population retains sufficient information to recover more than 90%
692   the input variance despite filtering out 50% or more of the original signal (Fig. 4). This information
693   recovery is achieved at the population level and thus requires sufficient numbers of granule cells so that
694   the subset of signals that are subthreshold are also super-threshold in other subsets of GCs through
695   probabilistic integration with other active inputs. While variance recovery is not a true measure of mutual
696   information, it is indicative of the utility that the intersectional filtering performed by the GCL. The
697   expansion of representations in the GCL population achieved by capturing the coincidence of features in
698   the input population creates a flexible representation in the GCL output that has many beneficial
699   properties, including the preservation of information through some degree of preserved mutual
700   information between the GCL and its inputs.
701
702   *Enhanced learning speed*
703   Our model not only improved learning accuracy, but also speed, compared to MFs alone (Fig. 3). Both
704   learning speed and accuracy progressed in tandem: threshold parameter ranges that enhanced overall
705   learning speed also minimized mean squared error, suggesting that speed and accuracy are enhanced by
706   similar features in GCL output. Learning speed was well described by a double exponential function with
707   a slow and fast component. This dual time course in the model with only one learning rule is interesting
708   in light of observations of behavioral adaptation that also follow dual time courses (Herzfeld et al., 2014;
709   Smith et al., 2006). Some behavioral studies have postulated that these time courses suggest multiple
710   underlying learning processes (Yang and Lisberger, 2014). Our model indicates that even with a single
711   learning rule and site of plasticity, multiple time-courses can emerge, presumably because when error
712   becomes low, update rates also slow down.
713
714   Another observation stemming from simulations studying learning speed was that the behavior of the
715   model varied as a function of the learning 'step size' parameter of the gradient descent method (Fig 3 –
716   Fig. Supplement 1). The step size -- ie. the, typically small, scalar regulating change in the weights
717   between GCs and P-cells following an error -- determined the likelihood of catastrophically poor learning:
718   when the step size was too large, it led to extremely poor learning because the total output 'explodes' and
719   fails to converge on a stable output. Nevertheless, the model tolerated large steps and faster learning
720   under some conditions, since the threshold also influenced the likelihood of catastrophic learning.
721   Generally, higher thresholds prevented large weight changes from exploding, suggesting that sparse
722   outputs may have an additional role in speeding learning by supporting larger weight changes in Purkinje
723   cells. Indeed, appreciable changes in simple spike rates occur on a trial-by-trial basis, gated by the
724   theorized update signals that Purkinje cells receive, climbing fiber mediated complex spikes. These
725   plastic changes in rate could reflect large weight updates associated with error. Moreover, graded
726   complex spike amplitudes that alter the size of trial-over-trial simple spike rate changes suggest that
727   update sizes are not fixed (Najafi et al., 2014; Herzfeld et al., 2020; Medina and Raymond 2018). It is
728   possible that the amplitude of synaptic weight changes following a complex spike might be set by tunable
729   circuitry in the molecular layer to optimize learning speed relative to the statistics of the GCL output.
730

731 Together, this study advances our understanding of how the GCL may diversify or isolate components of
732 inputs. A number of behavioral observations might be informed by the present findings. The timecourse
733 of learning for instance varies widely across tasks. Eyeblink conditioning paradigms require hundreds of
734 trials to learn (Millenson 1997; Khilkevich et al., 2016; Lincoln et al., 1982), while saccade adaptation
735 and visuomotor adaptation of reaches, which are also mediated by the cerebellum (Raymond and
736 Lisberger; Martin et al. 1996), requires just tens of trials (Tseng et al., 2007; Shadmehr and Mussa-Ivaldi
737 1994; Ruttle et al., 2021; Calame et al., 2021). This discrepancy in learning rates raised the possibility that
738 the learning algorithm used by the cerebellum is better engaged during naturalistic movements compared
739 to time-invariant cues, such as a conditioning stimulus. Such purely time-invariant cues would be
740 difficult, if not impossible, for our model GCL to reformat and sparsen, as they are incompatible with
741 thresholding-based filtering of input signals used here. Supportive of this view, recent work showed that
742 EBC learning was faster if the animal is locomoting during training (Albergaria et al., 2018). We
743 hypothesize that naturalistic time-variant signals associated with ongoing movements inputted to the
744 cerebellum through MFs support robust temporal pattern separation in the GCL, enhancing learning
745 accuracy and speed, while time invariant associative signals used in typical classical conditioning
746 paradigms result in an impoverished 'basis', making learning more difficult. That this feature is so robust
747 could explain why tasks like eyeblink conditioning are so difficult to learn, sensorimotor tasks can be
748 adapted rapidly. We speculate that the cerebellum is structured to support fast learning in situations where
749 there are physiologically structured inputs, typified by convergent, temporally varying self-generated
750 efference and reafference, within rich sensory and motor environments, as in normal movements during
751 daily life.

752 **Methods**
753 **Model construction**
754 The model presented here incorporated major features of the granule cell layer (GCL) circuit anatomical
755 organization and physiology. The features chosen for the model were the sparse sampling of inputs (GCs
756 have just 4 synaptic input branches in their segregated dendrite complexes on average), which was
757 reflected in the connectivity matrix between the input pool and the GCs, where each GC received 4 inputs
758 with weights of $1/4^{th}$ (i.e. 1 divided by the number of inputs; $1/M$) of the original input strength, summing
759 to a total weight of 1 across all inputs.  The other features were thresholding, representing inhibition from
760 local inhibitory Golgi neurons and intrinsic excitability of the GCs. The degree of inhibition and intrinsic
761 excitability (threshold) was a free parameter of the model, and the dynamics were normalized to the z-
762 score of the summated inputs. This feature reflects the monitoring of inputs by Golgi cells while
763 maintaining simplicity in their mean output to GCs. While this model simplifies many aspects of previous
764 models of the GCL, it recreated many of the important features of those models, suggesting that the
765 sparse sampling and firing are the main components dictating GCL functionality.
766 The model, in total, uses the following formulas to determine GC output:
767
768 $$\text{Eq 1: } GC_i(t) = [(\sum_{k_I}^{k_M} \frac{MF_k(t)}{M}) - \theta]_+$$
769
770 where k is a random selection of M MFs from the MF population. The inputs are summed and divided by
771 the total number of MF inputs to the GC, M, so that their total weight is equal to 1. Unless noted as a
772 variable, we used M = 4, reflecting the mean connectivity between MFs and GCs, and the optimal ratio
773 for expansion recoding (Litwin-Kumar et al. 2017), and the point of best input variance retention (Fig. 4).
774 This function is then linearly rectified, i.e. $[x]_+ = x$  if $x > 0$ and 0 otherwise so that there are no negative
775 rates present in the GC activity. The $\theta$ function which determines the threshold mimics intrinsic
776 excitability and feedforward inhibition was formulated as:

777

$$\text{Eq 2: } \theta = \overline{MF} + (z * \sigma(MF))$$

778

779

780     Here, a function of the mean and standard deviation of the entire MF population, z is a free parameter in
781     the model representing the number of standard deviations from the mean, setting the minimum value
782     below which granule cell activity is suppressed, which is the threshold value reported within this study as
783     'threshold'. Note that the summated MF inputs are divided by the number of inputs per GC (N) in Eq. 1
784     such that their received activity relative to $\theta$ is proportional to the input size, M.

785

786     **Input construction**
787     To provide a range of inputs with physiological-like temporal properties that could be parameterized, we
788     used a class of randomly generated signals called Ornstein-Uhlenbeck Processes (OU), defined by the
789     following formula:

790

$$\text{Eq 3: } OU(t) = \left( OU(t - \Delta t) * e^{\left(-\frac{\Delta t}{\tau}\right)} \right) + \left( \sigma * \sqrt{1 - e^{-2*\frac{\Delta t}{\tau}}} * R \right)$$

791

792     Here t is the time point being calculated, $\Delta t$ is the time interval (the time base is in *ms* and $\Delta t$ is 1 *ms*). $\sigma$ is
793     the predetermined standard deviation of the signal, and R is a vector of normally distributed random
794     numbers. This process balances a decay term, the exponential with e raised to $-\Delta t/\tau$, and an additive term
795     which introduces random fluctuations. Without the additive term, this function decays to zero as time
796     progresses. After the complete function has been calculated, the desired mean is added to the timeseries to
797     set the mean to a predetermined value.

798

799     The vector R can also be drawn from a matrix of correlated numbers, as was the case in Fig. 6 – figure
800     supplement 1 B & C. These numbers were produced with the MATLAB functions randn() for normal
801     random numbers, and mvnrnd() for matrices with a predetermined covariance matrix supplied to the
802     function. The covariance matrix used for these experiments was always a 1-diagonal with a constant,
803     predetermined, covariance value on the off-diagonal coordinates.

804

805     **Learning accuracy and speed assay**
806     In order to understand how the GCL contributed to learning, we constructed an artificial Purkinje cell (P-
807     cell) layer. The P-cell unit learned to predict a target function through a gradient descent mechanism, such
808     that the change in weight for each step was:

809

810

$$\text{Eq 4: } Err(t) = |P(t) - TF(t)|$$

811

812

$$\text{Eq 5: } \Delta W_i = W_i - (Err(t) * GG_i(t) * \eta)$$

813

814     Where P(t) is the output of the P-cell at time t, TF(t) is the target function at time t, $W_i$ is the weight
815     between the Purkinje cell and the $i^{th}$ GC, and $\eta$ is a small scalar termed the 'step size'. $\eta$ was 1E-3 for
816     GCs, and 1E-5 for MF alone in simulations shown in this study where the step size was held fixed, which
817     was chosen to maximize learning accuracy and stability of learning for both populations. The learning
818     process in Eq. 4 and 5 was repeated for T trials at every time point in the desired signal. The number of
819     trials was chosen so that learning reached asymptotic change across subsequent trials. Typically, 1000

820  trials were more than sufficient to reach asymptote, so that value was used for the experiments in this
821  study.
822
823  The overall accuracy of this process was determined by calculating the mean squared error between the
824  predicted and desired function:
825
826  $$\text{Eq 6: } MSE \ = \frac{1}{T}\sum_{t=1}^{T}(P(t) - TF(t))^2$$
827
828  The learning speed was determined by fitting an exponential decay function to the MSE across every trial
829  and taking the tau of the decay (See methods: Model output metrics, Time decay).
830
831  **Model output metrics**
832  To assay the properties of the GCL output that influence learning, we measured the features of GCL
833  output across a spectrum of metrics that have theoretically been associated with GCL functions like
834  pattern separation or expansion, as well as optimization or cost-related metrics developed for this paper.
835  These included: dimensionality, spatiotemporal sparseness, contributing principal components, spatial
836  sparseness (mean population pairwise correlation), temporal sparseness (mean unit autocovariance
837  exponential decay), population variance, temporal lossiness, population lossiness, and temporal cover.
838
839  We considered three forms of lossiness here, two related to the dimensions of sparseness considered
840  above, time and space, and one that is a measure of sparseness on the individual GC level. Temporal
841  lossiness is a measure of the percentage of time points that are not encoded by any members of the GCL
842  population, essentially removing the ability of P-cells to learn at that time point and producing no output
843  at that time in the final estimation of the target function. Increases in the value are guaranteed to degrade
844  prediction accuracy for any target function that does not already contain a zero value at the lossy time
845  point.
846  $$\text{Eq 7:}$$

847  $$Temp.\,Lossiness = \frac{1}{T}\sum_{t=1}^{T}x_t \ where \ x_t \begin{cases} (\sum_{i=1}^{N} GC_i(t)) \leq 0 = 1 \\ else = 0 \end{cases}$$

848
849  Here, T is the total number of points in the encoding epoch, the bracketed portion of the formula is a
850  summation of inputs from all GCs (N = population size) at that timepoint. When all GCs are silent, the
851  sum is 0, and the temporal lossiness is calculated as 1, and when all time points are covered by at least
852  one GC, total temporal lossiness is 0.
853
854  Spatial lossiness, or population lossiness, is the proportion of GCs in the population that are silent for the
855  entirety of the measured epoch. This is thought to reduce total encoding space and deprive downstream P-
856  cells of potential information channels and could potentially impact learning efficacy. It is defined as:
857
858  $$\text{Eq 8:}$$

859

$$Pop.\,Lossiness = \frac{1}{N}\sum_{i=1}^{N} x_i \ where \ x_i \begin{cases} (\sum_{t=1}^{T} GC_t) \leq 0 = 1 \\ \\ else = 0 \end{cases}$$

860

861   Here, N is the total population size of the GCL, and the bracketed portion of the formula is a sum of the
862   activity of GCs across all timepoints, such that if a GC is silent across all timepoints $x_i$ is calculated as 1,
863   indicating the 'loss' of that GC unit's contribution. When all GCs are silent, population lossiness is 1, and
864   when all GCs are active for at least one time point, population lossiness is 0.

865

866   Additionally, we looked at the mean sparseness of activity across the population by measuring the
867   'coverage' or proportion of time points each GC was active during, defined as:
868                                                    Eq 9:

869   $$Coverage = \frac{1}{N}\sum_{i=1}^{N}(\frac{1}{T}\sum_{t=1}^{T} x_i \ where \ x_i \begin{cases} GC_i(t) > 0 = 1 \\ else = 0 \end{cases})$$

870

871   As before, N is the number of cells in the population and T is the total length of the epoch. The bracketed
872   function counts the number of time points where $GC_i$ is active, and divides that by the total time period
873   length to get the proportion of time active. This value is summed across all GCs and divided by N to
874   calculate the average coverage in the population. This value has strong synonymy with population
875   variance, so it was not used for fitting assays in later experiments (Fig. 6), but reflects the effect of
876   thresholding on average activity in the GCL population.

877

878   Dimensionality is a measure of the number of independent dimensions needed to describe a set of signals,
879   similar in concept to the principal components of a set of signals. This measure is primarily influenced by
880   covariance between signals, and when dimensionality approaches the number of signals included in the
881   calculation (n), the signals become progressively independent. The GCL has previously been shown to
882   enhance the dimensionality of input sets and does so in the model presented here too. Dimensionality is
883   calculated with:

884                          Eq 10: $Dim = (\sum_{i=1}^{n} \lambda_i)^2 / (\sum_{i=1}^{n} \lambda_i^2)$

885

886   Provided by Litwin-Kumar, et al, 2016. This is the ratio of the squared sum of the eigenvalues to the sum
887   of the squared eigenvalues of the covariance matrix of the signals.

888

889   Spatiotemporal Sparseness (STS) was a calculated cost function meant to measure the divergence of GC
890   population encoding from a 'perfect' diagonal function where each GC represents one point in time and
891   does not overlap in representation with other units. This form of representation is guaranteed to produce
892   perfect learning, and transformations between the diagonal and any target function can be achieved in a
893   single learning step, making this form of representation an intriguing form of GCL representation, if it is
894   indeed feasible. We calculated the cost as:

895

896                          Eq 11: $STS = (1 - L_t) * (\frac{1}{T}) * (\frac{W}{GC_w})$

897

898     Where $(1 - L_t)$ is the cost of temporal lossiness, defined above (Eq. 7), and T is the total length of the
899     epoch. W is the number of unique combinations (termed 'words', akin to a barcode of activity across the
900     population), of GCs across the epoch at each point of discrete time, and $GC_w$ is the average number of
901     words each GC is active at all within the time-bins chosen (e.g. a binary representation of GC activity).
902     The intuition used here is that when there is no temporal lossiness, all points in time are represented,
903     leading the $1 - L_t$ term to have no effect on the STS equation, and when W, the number of unique
904     combinations of GC activities is equal to T, then each point in time has a unique 'word' associated with it.
905     Finally, when $GC_w$ is 1, $W/GC_w$ is equal to W, which only occurs when each GC contributes to a single
906     word. When these conditions are met, STS = 1, otherwise when GCs contribute to more than one word,
907     $GC_w$ increases and W is divided by a number larger than 1, decreasing STS. Alternately, when there are
908     not many unique combinations, such as when every GC has the exact same output, $W/GC_w$ is equal to
909     (1/T), decreasing STS. Finally, because lossiness causes the occurrence of a 'special', but non-associable,
910     word, we multiplied the above calculations by $(1 - L_t)$ to account for the effect of the unique non-encoding
911     word (i.e. all GCs inactive) on distance from the ideal diagonal matrix.
912
913     Mean temporal decay, i.e. temporal sparseness, is a measure of variance across time for individual
914     signals, where a low value would indicate that the signals coherence across time is weak, meaning that the
915     signal varies quickly, whereas a high value would mean that trends in the signal persist for long periods of
916     time. This value is extracted by fitting an exponential decay function to the autocovariance of each unit's
917     signal and measuring the tau of decay in the function:
918

$$\text{Eq 12: } y = a * e^{(-x/\tau)}$$

920
921     This is converted to the ms form by taking the ratio of $1000/\tau$. $y$ here $\tau$ is a description of the
922     autocovariance of the activity of a MF or GC signal, so when the descriptor $\tau$ is a large number, the decay
923     in autocovariance is longer, or slower, when $\tau$ is a small number, the autocovariance across time decays
924     more quickly, making the change in activity faster.
925
926     While dimensionality and STS are metrics rooted in a principled understanding of potentially desirable
927     properties of population encoding, the gradient descent algorithm can extract utility from population
928     statistics that are much noisier and correlated than the ideal populations that dimensionality and STS
929     account for. To measure a more general pattern separation feature in GCL output that could still be
930     associated with the complex target function, we turned to principal component analysis (PCA) with the
931     intuition that components which explain variance in the GCL output could be utilized by the downstream
932     Purkinje cell units to extract useful features from the input they receive (Lanore et al., 2021). We
933     parameterized the utility of this measure by taking the proportion of the PCs derived from the GCL output
934     which explained variance (of the GCL output) in that population by more than or equal to 1/N, where N is
935     the number of GCs, suggesting that they explain more variance than would be expected from chance.
936
937     Population correlation, was measured by taking the mean correlation between all pairwise combinations
938     of GCs using the corr() function in MATLAB and excluding the diagonal and top half of the resultant
939     matrix.
940

941 Population aggregate variance is a measure related to the expansion or collapse of total space covered by
942 the encoding done by a population, and higher or expanded values in this metric are thought to assist in
943 pattern separation and classification learning.
944
945 $$\text{Eq 13: Pop. Var} = \sum_{n=1}^{N} (x_n - \mu)^2$$
946
947 As shown in Cayco-Gajic et al. (2017). Here x is the activity of one of n cells across a measured epoch,
948 and μ is the mean of that activity. This value is reported relative to the number of GC units, such that Pop.
949 Var reported in Fig. 5 is normalized to Pop. Var / N.
950
951

**Variance retained assay**

953 To test the recovery of inputs by a feedforward network with a granule cell layer (GCL), we used
954 explained variance, $R^2$ , to quantify the quality of recovery of a sequence of normal random variables
955 (Fig. 2) across $N_w = 1000$ numerical experiments. To distinguish this metric from the MSE and $R^2$
956 metrics to evaluate other models in the study, we rename this 'variance retained'. Within each numerical
957 experiment $i$, at each time point, a vector of inputs $x_t$ of length $M$ (representing the mossy fiber, MF,
958 inputs) was drawn from an $M$-dimensional normal distribution with no correlations, $x_t \sim \mathcal{N}(0, I_M)$. This
959 vector is then left-multiplied by a random binary matrix $W$ with $N$ rows and $M$ columns with $n$ 1's per
960 row and the rest zeros, followed by a threshold linearization to obtain the GCL output, $y_t = [W x_t - z]_+$
961 with threshold. This process is then repeated $T = 1000$ times and a downstream linear readout was fit to
962 optimally recover $x_t$ from $y_t$. It can be shown multivariate linear regression (MATLAB's regress()
963 function, employing least squares to minimize mean squared error) solves this problem, identifying for
964 each MF input stream $x_{1:T}^j$, the optimal weighting $B_{1:T}$ from the GCL to estimate $\hat{x}_{1:T}^j = B_{j,1:N} y_{1:T}$.
965 Across time $t = 1:T$, we then computed the squared error across the vector, $MSE_i = \sum_{t=1}^{T} \sum_{j=1}^{M} (\hat{x}_t^j -$
966 $x_t^j)^2$, as well as the summed variance of the actual input, $Var_i = \frac{1}{MT} \sum_{j=1}^{M} \sum_{t=1}^{T} (x_t^j - \bar{x}^j)^2$, where $\bar{x}^j =$
967 $\frac{1}{T} \sum_{t=1}^{T} x_t^j$ is the mean of the $j$th MF input stream. Lastly, to compute variance explained, we take $R^2 =$
968 $1 - \frac{\sum_{i=1}^{N_w} MSE_i}{\sum_{i=1}^{N_w} Var_i}$, so the higher the relative mean squared error is, the lower the variance explained will be.
969 To generate the panels in Fig. 4, we always kept the number of timepoints and experiments the same, but
970 varied (Fig. 4B) the threshold along the axis and the number of inputs $n$ per GC output; (Fig. 4C) the total
971 number of GC outputs $N$ and input per output $n$; (Fig. 4D) number of inputs $M$ and outputs $N$; and finally
972 (Fig. 4E) the number of inputs per GC output $n$ along with the total number of outputs $N$.
973

**Independent measures generation**

975 To determine if the sparseness measures had inherent benefits for learning, we supplemented the GCL
976 output with OU processes with known temporal and correlational properties to examine their effect on
977 learning accuracy (Figure 6 figure supplement 1). We varied the temporal properties by systematically
978 varying the tau value in the exponential decay function. To vary population correlation, the random draw
979 function in the OU process was replaced with a MATLAB function, mvnrnd(), which allowed for preset
980 covariance values to direct the overall covariance between random samples. We used a square matrix with
981 1s on the diagonal and the desired covariance on all off-diagonal locations for this process and varied the

982    covariance to alter the correlation between signals. The OU outputs from this controlled process were
983    then fed into model P cells with randomized OU targets, as per the normal learning condition described
984    above. To vary the effect of the input population size, the size of the supplemented population varied
985    from 10 to 3000 using a tau of 10 and drawing from normal random numbers.

986

987    To measure the effects of STS on learning, a diagonal matrix was used at the input to a Purkinje unit,
988    which represented population activity with an STS of 1 (see Eq 11 in Model output metrics). To degrade
989    the STS metric, additional overlapping activity was injected either by expanding temporal representation
990    or at random, for example, adding an additional point of activity causes inherent overlap in the diagonal
991    matrix, increasing the $GC_w$ denominator of Eq 11 to $(1 + 2/N)$ because the overlapping and overlapped
992    units now each contribute to 1 additional neural word.  This process was varied by increasing the amount
993    of overlap to sample STS from 0 to 1.

994

995    **GCL output metrics fits to learning**
996    To estimate the properties of GCL output that contribute to enhanced learning of time series, we used
997    multiple linear regression to find the fit between measures of GCL population activity and observed MSE
998    in learning. Because there are large inherent correlations between the metrics used (dimensionality,
999    spatiotemporal sparseness, explanatory principal components of the GC population, population
1000    variability, mean pairwise GC correlation, temporal sparseness, temporal lossiness, population lossiness,
1001    and input variance retained) we used two linear regression normalization techniques: LASSO and RIDGE
1002    regression. For Figure 6, LASSO was used to isolate the 'top' regressors, while RIDGE was used in
1003    Figure 7 to preserve small contributions from regressors. The RIDGE regression method was then used to
1004    compare resultant regression slopes (beta coefficients) to changes in task parameters (see Methods on
1005    Simulation of cerebellar tasks).

1006

1007    Regressions were performed using the fitrlinear() function in MATLAB, with LASSO selected by using
1008    the 'SpaRSA' (Sparse Reconstruction by Separable Approximation; Wright et al., 2009) solver, and
1009    RIDGE selected with the 'lbfgs' (Limited-memory BFGS; Nodecal and Wright 2006) solver techniques.
1010    The potential spread of MSE in the models was determined using a K-fold validation technique, with 10
1011    'folds' used, as well as for determining the range of slopes shown in Figures 7, B, C, E, F, H, and I, of
1012    which the mean and standard deviation of cross-validation trials are plotted with solid lines and shaded
1013    polygons, respectively. Models were selected by choosing the model with the least complex fitting
1014    parameters (i.e. the model with the highest Lamba) while still falling within the bounds of the model with
1015    the minimized MSE plus the standard error (a standard '1SE' method).

1016

1017    To convey the overall contribution of regressors to the above models of MSE, both the slope (e.g. 'Beta')
1018    (Fig. 7: B, E, H), and the slope relative to the magnitude of all slopes were used as plotted metrics (Fig. 7:
1019    C, F, I).

1020

1021    **Simulation of cerebellar tasks**
1022    To simulate the input and output relationship observed in cerebellar and cerebellar-related tasks like
1023    vestibulo-ocular reflex adaptation (VOR), interval estimation, and motor-kinematic transformations, we
1024    adjusted the inputs and target functions in the task used above to mimic these scenarios. For the VOR-like
1025    task (Fig. 7 A-C), the inputs were 10% cosines with a fixed period and amplitude (10Hz, Amplitude range

1026     [0, 2]) and the rest were OU processes with taus of 100 and means and standard deviations of 0.5, and 0.2.
1027     The target functions were also cosines whose periods and amplitudes were identical to the inputs, but
1028     which had phase offsets between 0 and pi to mimic phase-offset VOR tasks.
1029

1030     The interval estimation tasks (Fig. 7 D-F) had standard OU inputs with target functions that were step
1031     functions with amplitude ranges from 0 to 1 and intervals that ranged from 0 to 1000 ms, which was the
1032     maximal extent of the epoch.
1033

1034     Finally, to simulate the transformation between motor commands and kinematic predictions, we used
1035     human EMG as a proxy for a motor command-like input signal to the GCL. 30 muscles from 15 bilateral
1036     target muscles were used (Delis et al., 2018; Hilt et al., 2018). The target function was a kinematic
1037     trajectory recorded simultaneously with the recordings of EMG used for the study. Although many body
1038     parts and coordinate dimensions were recorded of the kinematics, we opted to use the kinematic signal
1039     with the largest variance to simplify the experiment to a single target function.
1040

1048

1049     **References**
1050     Achen CH (1982) Interpreting and Using Regression. Sage University Paper Series on Quantitative
1051     Applications in the Social Sciences, vol. 29.
1052

1053     Albergaria C, Silva NT, Pritchett DL, Carey MR (2018) Locomotor activity modulates associative
1054     learning in mouse cerebellum. Nature Neuroscience 21, 725–735. doi:10.1038/s41593-018-0129-x
1055

1056     Albus JS (1971) A theory of cerebellar function. Mathematical Biosciences 10, 25–61. doi:10.1016/0025-
1057     5564(71)90051-4
1058

1059     Albus JS (1975) Data Storage in the Cerebellar Model Articulation Controller (CMAC). ASME. J. Dyn.
1060     Sys., Meas., Control. 97(3): 228–233.
1061

1062     Armano S, Rossi P, Taglietti V, D'Angelo E (2000) Long-term potentiation of intrinsic excitability at the
1063     mossy fiber–granule cell synapse of rat cerebellum. J Neurosci 20:5208–5216, pmid:10884304.
1064

1065     Apps R, Garwicz M (2005) Anatomical and physiological foundations of cerebellar information
1066     processing. Nat Rev Neurosci 6:297-311
1067

1068     Bengtsson F, Jorntell H (2009) Sensory transmission in cerebellar granule cells relies on similarly coded
1069     mossy fiber inputs. Proceedings of the National Academy of Sciences 106, 2389–2394.
1070     doi:10.1073/pnas.0808428106
1071

1072    Billings G, Piasini E, Lőrincz A, Nusser Z, Silver RA (2014) Network Structure within the Cerebellar
1073    Input Layer Enables Lossless Sparse Encoding. Neuron 83, 960–974. doi:10.1016/j.neuron.2014.07.020
1074
1075    Buonomano DV, Mauk MD (1994) Neural Network Model of the Cerebellum: Temporal Discrimination
1076    and the Timing of Motor Responses. Neural Computation 6, 38–55. doi:10.1162/neco.1994.6.1.38
1077
1078    Calame DJ, Becker MI, Person AL (2021) Associative learning underlies skilled reach adaptation.
1079    Biorxiv 2021.12.17.473247
1080
1081    Cayco-Gajic NA, Clopath C, Silver RA (2017) Sparse synaptic connectivity is required for decorrelation
1082    and pattern separation in feedforward networks. Nature Communications 8. doi:10.1038/s41467-017-
1083    01109-y
1084
1085    Cayco-Gajic NA, Silver RA (2019) Re-evaluating Circuit Mechanisms Underlying Pattern Separation.
1086    Neuron 101, 584–602. doi:10.1016/j.neuron.2019.01.044
1087
1088    Chabrol FP, Arenz A, Wiechert MT, Margrie TW, Digregorio DA (2015) Synaptic diversity enables
1089    temporal coding of coincident multisensory inputs in single neurons. Nature Neuroscience 18, 718–727.
1090    doi:10.1038/nn.3974
1091
1092    Crowley JJ, Fioravante D, Regehr WG (2009) Dynamics of fast and slow inhibition from cerebellar Golgi
1093    cells allow flexible control of synaptic integration. Neuron 63:843–853.
1094    doi:10.1016/j.neuron.2009.09.004 pmid:19778512
1095
1096    D'Angelo E, De Zeeuw CI (2009) Timing and plasticity in the cerebellum: focus on the granular layer.
1097    Trends Neurosci 32:30–40. doi:10.1016/j.tins.2008.09.007 pmid:18977038
1098
1099    De Zeeuw CI, Simpson JI, Hoogenraad CC, Galjart N, Koekkoek SK, Ruigrok TJ (1998) Microcircuitry
1100    and function of the inferior olive. Trends Neurosci. 21: 391–400
1101
1102    De Zeeuw CI (2020) Bidirectional learning in upbound and downbound microzones of the cerebellum.
1103    Nat Rev Neurosci 22: 92-110
1104
1105    Dean P, Porrill J (2008) Adaptive-filter Models of the Cerebellum: Computational Analysis. The
1106    Cerebellum 7, 567–571. doi:10.1007/s12311-008-0067-3
1107
1108    Delis I, Hilt PM, Pozzo T, Panzeri S, Berret B (2018) Deciphering the functional role of spatial and
1109    temporal muscle synergies in whole-body movements. Scientific Reports 8. doi:10.1038/s41598-018-
1110    26780-z
1111
1112    Hilt PM, Delis I, Pozzo T, Berret B (2018) Space-by-Time Modular Decomposition Effectively Describes
1113    Whole-Body Muscle Activity During Upright Reaching in Various Directions. Front Comput Neurosci.
1114    2018;12:20. doi:10.3389/fncom.2018.00020
1115
1116    Dino MR, Schuerger RJ, Liu Y, Slater NT, Mugnaini E (2000) Unipolar brush cell: a potential
1117    feedforward excitatory interneuron of the cerebellum. Neuroscience 98:625–636.
1118
1119    Duguid I, Branco T, London M, Chadderton P, Hausser M (2012) Tonic Inhibition Enhances Fidelity of
1120    Sensory Information Transmission in the Cerebellar Cortex. The Journal of Neuroscience 32, 11132–
1121    11143. doi:10.1523/jneurosci.0460-12.2012
1122

Eccles JC, Ito M, Szentágothai J (1967) The Cerebellum as a Neuronal Machine, Springer, New York (1967)

Eriksson, JL, Robert A (1999) The representation of pure tones and noise in a model of cochlear nucleus neurons. The Journal of the Acoustical Society of America 106, 1865–1879. doi:10.1121/1.427936

Fujita M (1982) Adaptive filter model of the cerebellum. Biological Cybernetics 45, 195–206. doi:10.1007/bf00336192

Gall D, Prestori F, Sola E, D'Errico A, Roussel C, Forti L, Rossi P, D'Angelo E (2005) Intracellular calcium regulation by burst discharge determines bidirectional long-term synaptic plasticity at the cerebellum input stage. J Neurosci 25:4813–4822, doi:10.1523/JNEUROSCI.0410-05.2005, pmid:15888657.

Gao Z, Proietti-Onori M, Lin Z, Ten Brinke MM, Boele HJ, Potters JW, Ruigrok TJ, Hoebeek FE, De Zeeuw CI (2016) Excitatory cerebellar nucleocortical circuit provides internal amplification during associative conditioning. Neuron 89:645–657. 10.1016/j.neuron.2016.01.008

Gilmer JI, Person AL (2017) Morphological Constraints on Cerebellar Granule Cell Combinatorial Diversity. The Journal of Neuroscience 37, 12153–12166. doi:10.1523/jneurosci.0588-17.2017

Gilmer JI, Person AL (2018) Theoretically Sparse, Empirically Dense: New Views on Cerebellar Granule Cells. Trends in Neurosciences 41, 874–877. doi:10.1016/j.tins.2018.09.013

Giovannucci A, Badura A, Deverett B, Najafi F, Pereira TD, Gao Z, Ozden I, Kloth AD, Pnevmatikakis E, Paninski L, De Zeeuw CI, Medina JF, Wang SS-H (2017) Cerebellar granule cells acquire a widespread predictive feedback signal during motor learning. Nature Neuroscience 20, 727–734. doi:10.1038/nn.4531

Guo C, Huson V, Macosko EZ, Regehr WG (2021) Graded heterogeneity of metabotropic signaling underlies a continuum of cell-intrinsic temporal responses in unipolar brush cells. Nat Comm 12:5491
Gurnani H, Silver RA (2021) Multidimensional population activity in an electrically coupled inhibitory circuit in the cerebellar cortex. Neuron 109, 1739–1753.e8. doi:10.1016/j.neuron.2021.03.027

Herculano-Houzel S (2010) Coordinated scaling of cortical and cerebellar numbers of neurons. Front Neuroanat 4:12. doi:10.3389/fnana.2010.00012 pmid:20300467

Herzfeld DJ, Hall NJ, Tringides M, Lisberger, SG (2020) Principles of operation of a cerebellar learning circuit. eLife 9. doi:10.7554/elife.55217

Herzfeld DJ, Kojima Y, Soetedjo R, Shadmehr R (2015) Encoding of action by the Purkinje cells of the cerebellum. Nature 526, 439–442. doi:10.1038/nature15693

Houck BD, Person AL (2014) Cerebellar loops: a review of the nucleocortical pathway. Cerebellum 13:378–385. 10.1007/s12311-013-0543-2

Houck BD, Person AL (2015) Cerebellar premotor output neurons collateralize to innervate the cerebellar cortex. J Comp Neurol 523:2254–2271. 10.1002/cne.23787

Huang CC, Sugino K, Shima Y, Guo C, Bai S, Mensh BD, Nelson SB, Hantman AW (2013) Convergence of pontine and proprioceptive streams onto multimodal cerebellar granule cells. Elife 2:e00400. doi:10.7554/eLife.00400 pmid:23467508

1174
1175  Hull C (2020) Prediction signals in the cerebellum: beyond supervised motor learning. Elife. 9:e54073.
1176  doi: 10.7554/eLife.54073. PMID: 32223891; PMCID: PMC7105376.
1177
1178  Ishikawa T, Shimuta M, Häusser, M (2015) Multimodal sensory integration in single cerebellar granule
1179  cells in vivo. eLife 4. doi:10.7554/elife.12916
1180
1181  Ito M, Shiida T, Yagi N, Yamamoto M (1974) Visual influence on rabbit horizontal vestibulo-ocular
1182  reflex presumably effected via the cerebellar flocculus. Brain Res 65:170 –174.
1183
1184  Izawa J, Criscimagna-Hemminger SE, Shadmehr R (2012) Cerebellar Contributions to Reach Adaptation
1185  and Learning Sensory Consequences of Action. The Journal of Neuroscience 32, 4230–4239.
1186  doi:10.1523/jneurosci.6353-11.2012
1187
1188  Jakab RL, Hamori J (1988) Quantitative morphology and synaptology of cerebellar glomeruli in the rat.
1189  Anatomy and Embryology 179, 81–88. doi:10.1007/bf00305102
1190
1191  Judd EN, Lewis SM, Person AL (2021) Diverse inhibitory projections from the cerebellar interposed
1192  nucleus. Elife. 10:e66231. doi: 10.7554/eLife.66231. PMID: 34542410; PMCID: PMC8483738.
1193
1194  Kalmbach BE, Voicu H, Ohyama T, Mauk MD (2011) A Subtraction Mechanism of Temporal Coding in
1195  Cerebellar Cortex. The Journal of Neuroscience 31, 2025–2034. doi:10.1523/jneurosci.4212-10.2011
1196
1197  Kanichay RT, Silver RA (2008) Synaptic and Cellular Properties of the Feedforward Inhibitory Circuit
1198  within the Input Layer of the Cerebellar Cortex. The Journal of Neuroscience 28, 8955–8967.
1199  doi:10.1523/jneurosci.5469-07.2008
1200
1201  Kennedy A, Wayne G, Kaifosh P, Alviña K, Abbott LF, Sawtell NB (2014) A temporal basis for
1202  predicting the sensory consequences of motor commands in an electric fish. Nature Neuroscience 17,
1203  416–422. doi:10.1038/nn.3650
1204
1205  Khilkevich A, Halverson HE, Canton-Josh JE, Mauk MD (2016) Links Between Single-Trial Changes
1206  and Learning Rate in Eyelid Conditioning. The Cerebellum 15, 112–121. doi:10.1007/s12311-015-0690-8
1207
1208  Kita K, Albergaria C, Machado AS, Carey MR, Müller M, Delvendahl I (2021) GluA4 facilitates
1209  cerebellar expansion coding and enables associative memory formation. eLife 10. doi:10.7554/elife.65152
1210
1211  Knogler LD, Markov DA, Dragomir EI, Štih V, Portugues R (2017) Sensorimotor Representations in
1212  Cerebellar Granule Cells in Larval Zebrafish Are Dense, Spatially Organized, and Non-temporally
1213  Patterned. Current Biology 27, 1288–1302. doi:10.1016/j.cub.2017.03.029
1214
1215  Krauzlis RJ, Lisberger SG (1991) Visual motion commands for pursuit eye movements in the cerebellum.
1216  Science 253:568-71
1217
1218  Lanore F, Cayco-Gajic NA, Gurnani H, Coyle D, Silver, RA (2021) Cerebellar granule cell axons support
1219  high-dimensional representations. Nature Neuroscience 24, 1142–1150. doi:10.1038/s41593-021-00873-x
1220
1221  Lincoln JS, Mccormick DA, Thompson RF (1982) Ipsilateral cerebellar lesions prevent learning of the
1222  classically conditioned nictitating membrane/eyelid response. Brain Research 242, 190–193.
1223  doi:10.1016/0006-8993(82)90510-8
1224

1225  Litwin-Kumar A, Harris KD, Axel R, Sompolinsky H, Abbott LF (2017) Optimal Degrees of Synaptic
1226  Connectivity. Neuron 93, 1153–1164.e7. doi:10.1016/j.neuron.2017.01.030
1227
1228  Liu, Y, Tiganj, Z, Hasselmo, ME, Howard, MW (2019) A neural microcircuit model for a scalable scale-
1229  invariant representation of time. Hippocampus. 29: 260– 274.
1230
1231  Mapelli L, Rossi P, Nieus T, D'Angelo E (2009) Tonic activation of GABAB receptors reduces release
1232  probability at inhibitory connections in the cerebellar glomerulus. J Neurophysiol 101:3089–3099.
1233  doi:10.1152/jn.91190.2008 pmid:19339456
1234
1235  Marr D (1969) A theory of cerebellar cortex. The Journal of Physiology 202, 437–470.
1236  doi:10.1113/jphysiol.1969.sp008820
1237
1238  Martin TA, Keating JG, Goodkin HP, Bastian AJ, Thach WT (1996) Throwing while looking through
1239  prisms: I. Focal olivocerebellar lesions impair adaptation. Brain 119, 1183–1198.
1240  doi:10.1093/brain/119.4.1183
1241
1242  Mauk MD, Buonomano DV (2004) THE NEURAL BASIS OF TEMPORAL PROCESSING. Annual
1243  Review of Neuroscience 27, 307–340. doi:10.1146/annurev.neuro.27.070203.144247
1244
1245  Mauk MD, Steinmetz JE, Thompson RF (1986) Classical conditioning using stimulation of the inferior
1246  olive as the unconditioned stimulus. Proc Natl Acad Sci USA, 83 pp. 5349-5353
1247
1248  Mccormick DA, Clark GA, Lavond DG, Thompson RF (1982) Initial localization of the memory trace for
1249  a basic form of learning. Proceedings of the National Academy of Sciences 79, 2731–2735.
1250  doi:10.1073/pnas.79.8.2731
1251  Medina J (2000) Mechanisms of cerebellar learning suggested by eyelid conditioning. Current Opinion in
1252  Neurobiology 10, 717–724. doi:10.1016/s0959-4388(00)00154-9
1253
1254  Miall RC, Wolpert DM (1996) Forward Models for Physiological Motor Control. Neural Netw. 9:1265–
1255  1279.
1256
1257  Millenson JR, Kehoe EJ, Gormezano I (1977) Classical conditioning of the rabbit's nictitating membrane
1258  response under fixed and mixed CS–US intervals. Learn Motiv, 8 pp. 351-366
1259
1260  Najafi F, Giovannucci A, Wang SS, Medina JF (2014) Coding of stimulus strength via analog calcium
1261  signals in Purkinje cell dendrites of awake mice. Elife. 2014;3:e03663. doi:10.7554/eLife.03663
1262
1263  Narain D, Remington ED, De Zeeuw CI (2018) A cerebellar mechanism for learning prior distributions of
1264  time intervals. Nat Commun. 9, 469.
1265
1266  Nocedal J,  Wright SJ (2006) Numerical Optimization, 2nd ed., New York: Springer.
1267
1268  Palay S, Chan-Palay V (1974) Cerebellar cortex: Cytology and Organization, pp 100–132. Berlin-
1269  Heidelberg: Springer-Verlag.
1270
1271  Palkovits M, Magyar P, Szentágothai J (1971) Quantitative histological analysis of the cerebellar cortex
1272  in the cat. Brain Research 32, 15–30. doi:10.1016/0006-8993(71)90152-1
1273
1274

1275   Rancz EA, Ishikawa T, Duguid I, Chadderton P, Mahon S, Häusser M (2007) High-fidelity transmission
1276   of sensory information by single cerebellar mossy fibre boutons. Nature 450, 1245–1248.
1277   doi:10.1038/nature05995
1278
1279   Raymond JL, Lisberger SG (1998) Neural Learning Rules for the Vestibulo-Ocular Reflex. The Journal
1280   of Neuroscience 18, 9112–9129. doi:10.1523/jneurosci.18-21-09112.1998
1281
1282   Raymond JL, Medina JF (2018) Computational Principles of Supervised Learning in the Cerebellum.
1283   Annu Rev Neurosci. 41:233-253. doi: 10.1146/annurev-neuro-080317-061948. PMID: 29986160;
1284   PMCID: PMC6056176.
1285
1286   Rizwan AP, Zhan X, Zamponi GW, Turner RW (2016) Long-Term Potentiation at the Mossy Fiber–
1287   Granule Cell Relay Invokes Postsynaptic Second-Messenger Regulation of Kv4 Channels. The Journal of
1288   Neuroscience 36, 11196–11207. doi:10.1523/jneurosci.2051-16.2016
1289
1290   Rossi P, D'Angelo E, Taglietti V (1996) Differential long-lasting potentiation of the NMDA and non-
1291   NMDA synaptic currents induced by metabotropic and NMDA receptor coactivation in cerebellar granule
1292   cells. Eur J Neurosci 8:1182–1189, doi:10.1111/j.1460-9568.1996.tb01286.x, pmid:8752588.
1293
1294   Rudolph S, Hull C, Regehr WG (2015) Active dendrites and differential distribution of calcium channels
1295   enable functional compartmentalization of Golgi cells. J Neurosci 35:15492–15504.
1296   doi:10.1523/JNEUROSCI.3132-15.2015 pmid:26609148
1297   Ruttle JE, Marius 't Hart B, Henriques DYP (2021) Implicit motor learning within three trials. Sci Rep.
1298   11(1):1627. doi: 10.1038/s41598-021-81031-y.
1299
1300   Sanger TD, Yamashita O, Kawato M (2020) Expansion coding and computation in the cerebellum: 50
1301   years after the Marr–Albus codon theory. The Journal of Physiology 598, 913–928. doi:10.1113/jp278745
1302
1303   Saviane C, Silver RA (2006) Fast vesicle reloading and a large pool sustain high bandwidth transmission
1304   at a central synapse. Nature 439:983–987. doi:10.1038/nature04509 pmid:16496000
1305
1306   Shadmehr R, Mussa-Ivaldi F (1994) Adaptive representation of dynamics during learning of a motor task.
1307   The Journal of Neuroscience 14, 3208–3224. doi:10.1523/jneurosci.14-05-03208.1994
1308
1309   Simat M, Parpan F, Fritschy JM (2007) Heterogeneity of glycinergic and gabaergic interneurons in the
1310   granule cell layer of mouse cerebellum. J Comp Neurol 500:71–83. doi:10.1002/cne.21142
1311   pmid:17099896
1312
1313   Smith MA, Ghazizadeh A, Shadmehr R (2006) Interacting adaptive processes with different timescales
1314   underlie short-term motor learning. PLoS Biol 4: e179, 2006. doi:10.1371/journal.pbio.0040179
1315
1316   Solinas S, Nieus T, D'Angelo E (2010) A realistic large-scale model of the cerebellum granular layer
1317   predicts circuit spatio-temporal filtering properties. Front Cell Neurosci. 2010 May 14;4:12. doi:
1318   10.3389/fncel.2010.00012. PMID: 20508743; PMCID: PMC2876868.
1319
1320   Tabuchi S, Gilmer JI, Purba K, Person AL (2019) Pathway-Specific Drive of Cerebellar Golgi Cells
1321   Reveals Integrative Rules of Cortical Inhibition. The Journal of Neuroscience 39, 1169–1181.
1322   doi:10.1523/jneurosci.1448-18.2018
1323

1324 Tseng YW, Diedrichsen J, Krakauer JW, Shadmehr R, Bastian AJ. (2007) Sensory prediction errors drive
1325 cerebellum-dependent adaptation of reaching. J Neurophysiol. 2007 Jul;98(1):54-62. doi:
1326 10.1152/jn.00266.2007. Epub 2007 May 16. PMID: 17507504.
1327
1328 Tyrrell T, Willshaw D (1992) Cerebellar cortex: its simulation and the relevance of Marr's theory. Philos
1329 Trans R Soc Lond B Biol Sci. 29;336(1277):239-57. doi: 10.1098/rstb.1992.0059. PMID: 1353267.
1330
1331 Wagner MJ, Kim TH, Savall J, Schnitzer MJ, Luo L (2017) Cerebellar granule cells encode the
1332 expectation of reward. Nature 544, 96–100. doi:10.1038/nature21726
1333
1334 Wolpert DM, Miall RC, Kawato, M (1998) Internal models in the cerebellum. Trends in Cognitive
1335 Sciences 2, 338–347. doi:10.1016/s1364-6613(98)01221-2
1336
1337 Wright SJ, Nowak RD, Figueiredo MAT (2009) Sparse Reconstruction by Separable Approximation.
1338 Trans. Sig. Proc., Vol. 57, No 7: 2479–2493.
1339
1340 Yang Y, Lisberger SG (2014) Role of Plasticity at Different Sites across the Time Course of Cerebellar
1341 Motor Learning. The Journal of Neuroscience 34, 7077–7090. doi:10.1523/jneurosci.0017-14.2014
1342
1343 Zhou S, Masmanidis SC, Buonomano DV (2020) Neural Sequences as an Optimal Dynamical Regime for
1344 the Readout of Time. Neuron. 108(4):651-658.e5. doi: 10.1016/j.neuron.2020.08.020. Epub 2020 Sep 17.
1345 PMID: 32946745; PMCID: PMC7825362.