

1 Demystifying unsupervised learning: how it helps and hurts

2 Franziska Bröker^{1,2,3,4,*}, Lori L. Holt⁵, Brett D. Roads⁶, Peter Dayan^{1,7,†}, Bradley C. Love^{6,†}

3 ¹Department of Computational Neuroscience, Max Planck Institute for Biological Cybernetics, Tübingen, Germany

4 ²Gatsby Computational Neuroscience Unit, University College London, London, UK

5 ³Department of Psychology, Carnegie Mellon University, Pittsburgh, US

6 ⁴Neuroscience Institute, Carnegie Mellon University, Pittsburgh, US

7 ⁵Department of Psychology, University of Texas at Austin, Austin, US

8 ⁶Department of Experimental Psychology, University College London, London, UK

9 ⁷University of Tübingen, Tübingen, Germany

10 † equal contribution

11 * Correspondence: franziska.broeker.15@ucl.ac.uk (F. Bröker)

12

13 Published paper: <https://doi.org/10.1016/j.tics.2024.09.005>

14



15 *This work is licensed under CC BY-NC-ND 4.0. To view a copy of this license, visit*

16 <https://creativecommons.org/licenses/by-nc-nd/4.0/>

17 Highlights

18 Humans are not guaranteed to benefit from unsupervised experiences (and neither are machines).

19 Instead, given unsupervised experience, humans self-reinforce their predictions. This can help
20 performance when the predictions are accurate; it can hurt or have no effect when the predictions are
21 inaccurate.

22 Predictions depend on the internal representations of learners which are shaped by prior experiences.
23 Thus, prediction accuracy depends on how well internal representations align with the task. Only by
24 assessing these representations can researchers understand whether and why unsupervised learning
25 helps or hurts in a specific task and in a specific person.

26 Literatures on self-reinforcement and unsupervised learning in humans have largely operated in isolation
27 but would benefit from more crosstalk.

28 Insights also have broad implications for lifelong learning and the design of instruction.

29

30 *“There was, Carter thought, a downside to experience. ‘Experience is making the same mistake*
31 *over and over again, only with greater confidence,’ he said. The line wasn’t his, but he liked it.”*
32 *— Michael Lewis, The Premonition: A Pandemic Story*

33

34 **Abstract**

35 Humans and machines rarely have access to explicit external feedback, or supervision, yet
36 manage to learn. Most modern machine learning systems succeed because they benefit from
37 unsupervised data. Humans are also expected to benefit and yet, mysteriously, empirical results
38 are mixed. Does unsupervised learning help humans or not? We argue that the mixed results
39 are not conflicting answers to this question, but reflect that humans self-reinforce their
40 predictions in the absence of supervision, which can help or hurt depending on whether
41 predictions and task align. We use this framework to synthesize empirical results across various
42 domains to clarify when unsupervised learning will help or hurt. This provides new insights into
43 fundamentals of learning with implications for instruction and lifelong learning.

44

45

46 **Keywords:** unsupervised learning; semi-supervised learning; self-reinforcement; mental
47 representation; representation-to-task alignment

48 **Supervised and unsupervised learning**

49 We live and learn in an environment that rarely provides us with **supervision** (see Glossary) in
50 the form of explicit external **feedback**. For example, we have learned to call some animals
51 “sheep” and others “goats”. Many of us acquired this distinction at a young age when we spent
52 much time around our caretakers. Like an external teacher, they provided us explicitly with the
53 correct labels by naming animals in our field of view. Getting older, we still encounter sheep and
54 goats, as well as animals we have never seen before, but we now rarely have a teacher in tow.
55 Thus, our learning about the world could be helped if we also made use of the information
56 contained in all these unsupervised experiences (Fig. 1).

57 Machine learning faces a conspicuously similar problem. Typically, an abundance of
58 unsupervised data is available for learning (e.g., images of sheep and goats), but supervision
59 (e.g., human-annotated sheep / goat labels for each image) is rare and expensive. This has led
60 to extensive research aiming to harness the information contained in unsupervised data. As a
61 result, we now have powerful **learning algorithms** able to extract statistical information and
62 features from unsupervised data [1] which can be further fine-tuned to specific tasks [2] or used
63 to boost **supervised learning** [3]. Ultimately, the tremendous success of machine learning
64 methods stems from their ability to learn in the absence of supervision.

65 **The mystery of unsupervised learning in humans**

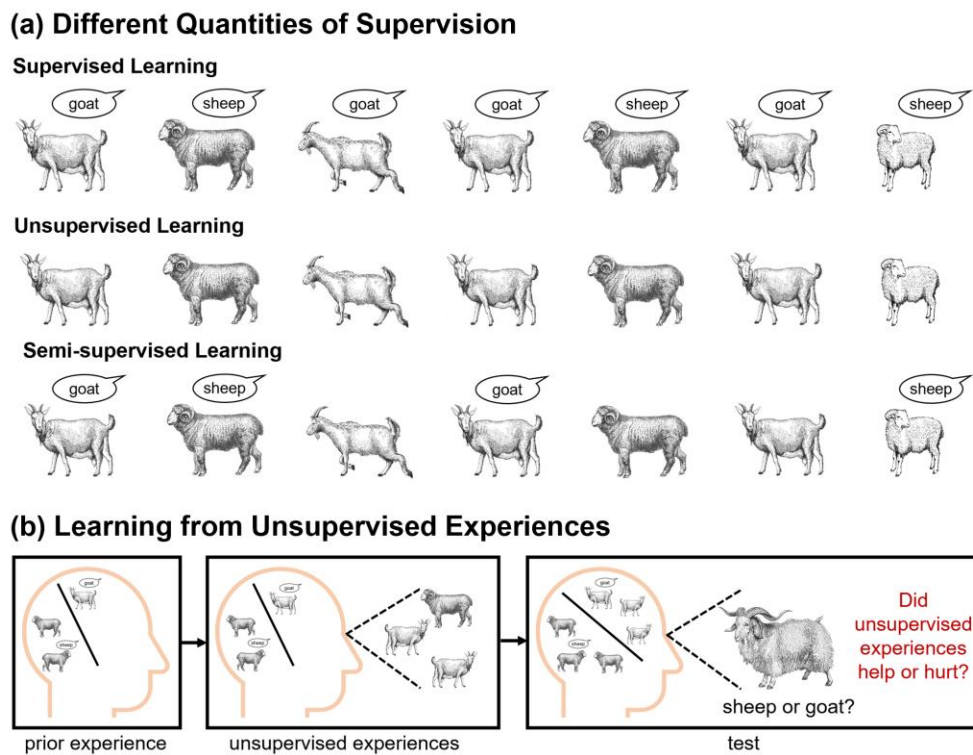
66 It seems clear that both humans and machines benefit from leveraging unsupervised
67 experiences. There has thus been a surge in empirical and computational work over the past
68 decades proposing that humans perform **unsupervised learning** by applying information
69 processing capabilities they share with machine learning algorithms [4–6]. A simple and
70 intuitive prediction results from this: If humans share unsupervised information processing
71 capabilities with machines, and machines show benefits leveraging unsupervised data, then
72 humans should benefit from their unsupervised experience in the same way. That is, humans
73 should be able to recover statistical information from their unsupervised experiences and they
74 should be able to combine it with their rare, supervised experiences.

75 Paradoxically, this is not supported by the scientific literature. In the most basic learning
76 experiments, humans are not guaranteed to extract the statistical information in their
77 unsupervised experiences [7–10] or to boost their supervised learning [11–13]. In fact,
78 unsupervised experiences can reduce performance in category learning [14], language learning
79 [15,16], motor learning [17] and stereotyping [18,19]. So instead of supporting the view that
80 unsupervised experiences help humans in their learning, the literature on lab studies is riddled
81 with equivocal results about their benefit. In one experiment people may need feedback to
82 learn how to distinguish between different visual inputs; in another, they do not [7,20].

83 These results stem from highly influential experimental designs that have shaped our
84 understanding of how humans extract statistical information. Unsupervised studies often use a
85 simple stimulus-response or passive exposure paradigm. These well controlled designs are
86 popular because they parallel supervised designs, allowing comparisons. In unsupervised
87 studies, learners predict task-appropriate responses from stimuli without feedback. The
88 statistics in the stimuli are the only information available for learning. Supervised studies are
89 close analogues which provide additional corrective feedback or correct labels, giving learners
90 more information.

91 Outside the lab, human learning operates on a larger scale in terms of data and time. For
92 example, an abundance of additional information can inform learning about sheep and goats,
93 like separate housing. Learning also serves long-term performance in the world rather than on
94 one specific task. Similarly, modern machine learning solves increasingly large-scale **learning**
95 **problems**. Because machine learning algorithms can be flexibly chosen for specific problems,
96 supervised algorithms now solve unsupervised problems by adapting the objective of the
97 **learning task**, as in **self-supervised learning**. Another example is large language models which
98 learn not by getting feedback on text they generate, but from predicting words in a sequence.
99 This then serves as a foundation model for further supervised fine-tuning on how to engage in
100 friendly chat with users. These developments increase the complexity in technical approaches
101 and terminology that has yet to be reconciled with human learning inside and outside the lab.
102 While an analogy between human and machine unsupervised learning is compelling and often
103 assumed, the devil appears to be in the details.

104 Here, we mainly focus on in-lab studies that test unsupervised or **semi-supervised learning**
 105 using the well-controlled, influential designs. Other unsupervised paradigms exist but are rarer
 106 [21]. Our narrower focus ensures results across various learning contexts are informative about
 107 the same learning principles. We refer to momentary learning from unsupervised experiences in
 108 experimental tasks as simply “unsupervised learning” to differentiate it from momentary
 109 learning with supervisory signals. While focusing on in-lab studies, we also present evidence
 110 suggesting unsupervised learning to be limited more generally, as it can worsen performance in
 111 machines [3] and human learning outside the lab [22]. In fact, as it happens, telling sheep apart
 112 from goats is a task on which many people fail despite recurring exposure (Fig. 1b).
 113



114
 115 **Figure 1. Learning with and without supervision.**
 116 (a) Illustration of supervised, unsupervised, and semi-supervised learning problems. (b) Empirical
 117 results conflict as to whether unsupervised experiences improve human performance in
 118 unsupervised and semi-supervised learning tasks. We refer to the momentary learning from
 119 unsupervised experiences as simply “unsupervised learning” throughout the manuscript. The

120 *reader is encouraged to guess whether the test animal is a sheep or goat. The answer is*
121 *provided in the footnote.¹*

122

123 **The unsupervised snowball effect**

124 How can we explain the mysterious results? When does unsupervised learning help and when
125 does it not? We think that the answer lies in the way that unsupervised learning is affected by
126 the relationship between the experimenter-defined task and the representations subjects have
127 acquired from prior experience (**representation-to-task alignment**, [14]). Concretely, we
128 propose the unsupervised learning mechanism to be **self-reinforcement** by which humans learn
129 from their own predictions so that pre-existing associations between experiences and
130 appropriate responses are strengthened (Fig. 2 b, Key Figure) and decision confidence increases.
131 For example, when seeing the woolly goat in Fig. 1b, readers who categorize by woolliness
132 would incorrectly self-reinforce their predictions that it is a sheep, whereas readers who know
133 to attend to the tail would correctly self-reinforce their prediction that it is a goat. Since
134 strengthening predictions snowballs existing learning without changing its course, self-
135 reinforcement can help or hurt depending on how accurate the predictions are for the task at
136 hand (Fig. 2 a). Self-reinforcing predictions that are largely correct will improve performance in
137 the task. But predictions will only be largely correct if prior experiences shaped the learner’s
138 representations in a way that new experiences elicit appropriate predictions. If this is the case,
139 representations and task are aligned, the task feels “easy”, and supervision is superfluous. By
140 contrast, self-reinforcing predictions that are largely incorrect will have a detrimental, or at best
141 no, effect on performance. Predictions will be largely incorrect if prior experiences have shaped
142 the learner’s representations to be misaligned with the task. In this case, the task feels “hard”,
143 and supervision is necessary to adjust the unhelpful representations and predictions of the
144 learner. That self-reinforcing existing representations results in these kind of learning dynamics
145 has previously been described in the specific context of unsupervised Hebbian (correlational)

¹ The test animal in Fig. 1b is a goat. While most non-experts make their sheep/goat predictions based on unreliable features, such as woolliness, the easiest way to tell them apart is by their tails: goats point their tails upward while sheep cannot lift their tails.

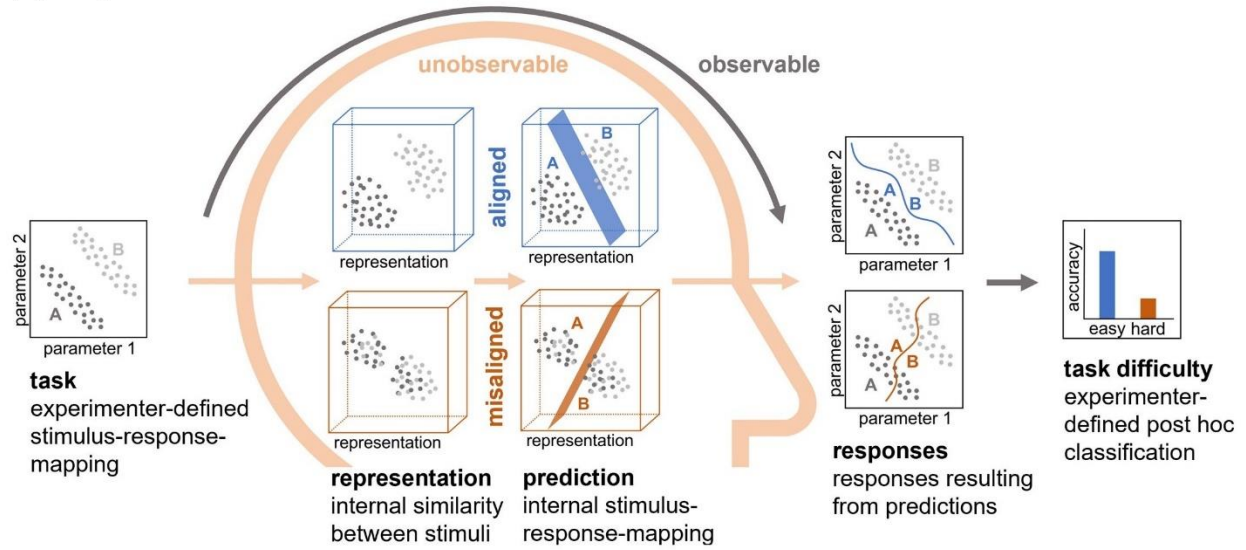
146 learning [23,24]. Our framing of unsupervised learning in terms of representation-to-task
147 alignment and self-reinforcement is more general in that it does not assume specific
148 representations nor a specific computational model of learning.

149 This type of self-reinforcing snowball effect can also be seen when trying to master a new skill,
150 such as playing the violin. This requires practicing with the correct technique because a faulty
151 technique engrains mistakes if left uncorrected. Thus, from our perspective, the equivocal
152 results in the literature about the benefit of unsupervised experiences do not reflect a conflict
153 but are in fact expected from representation-to-task alignment and its interaction with
154 unsupervised self-reinforcement. Our argument not only follows an intuitive logic but is also
155 supported by the theoretical principles that allow machine learning algorithms to leverage
156 unsupervised data on many, but not all, occasions (Box 1).

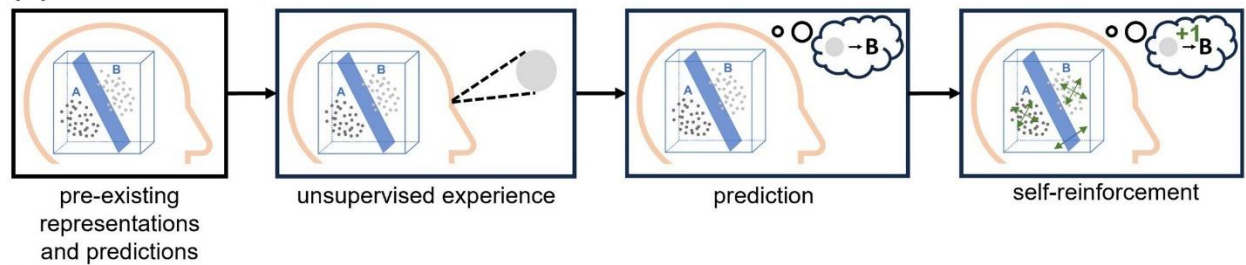
157 In the following, we provide support for this perspective by synthesizing various cognitive
158 science literatures that have long investigated the questions about how feedback influences
159 human learning. The fact that related research fields have largely developed in isolation allows
160 us to test our predictions against their extensive independent evidence. First, we show that
161 representation-to-task alignment correlates with the efficacy of unsupervised learning, as
162 predicted by our hypothesized unsupervised snowball effect. The evidence we consider for the
163 effect of alignment is often somewhat indirect because learners' representations, let alone their
164 alignment with the task, are not typically assessed. We thus leverage the equivalences between
165 representation-to-task alignment, predictions and task difficulty as described in Fig. 2a to
166 contextualize the results. Second, we show that unsupervised self-reinforcement has been
167 reported repeatedly across diverse learning settings. We conclude by discussing the implications
168 of our analysis and promising future avenues.

169

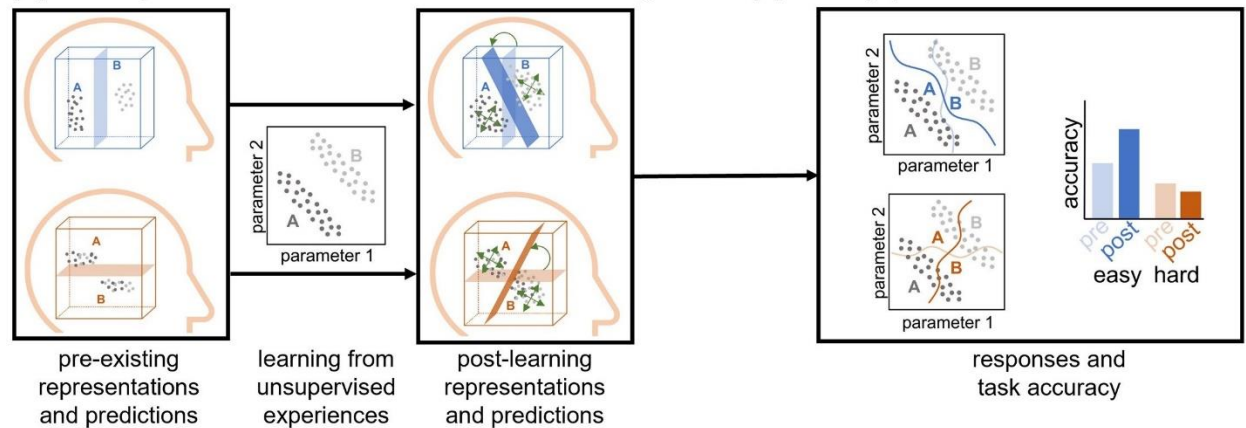
(a) Representation-to-task alignment



(b) Self-reinforcement



(c) Unsupervised snowball effect resulting from (a) and (b)



171

172

173 **Figure 2. The unsupervised snowball effect.**

174 *Two key factors affect unsupervised learning: representation-to-task alignment and self-*
 175 *reinforcement resulting in the unsupervised snowball effect as illustrated for the example of a*
 176 *category learning task. (a) Relationship between experimenter-defined task, its internal*
 177 *representation and the resulting predictions, responses, and accuracy. Factors including prior*

178 *experience, context or attention transform observed stimuli and warp their similarities into an*
179 *internal representational space that might or might not recover experimenter-defined task*
180 *statistics. If learners have a task-aligned representation, stimuli from different categories are*
181 *sufficiently separated in the learner’s representational space such that it supports accurate*
182 *predictions. The task will seem easy and performance will be high. If learners have a task-*
183 *misaligned representation, items from different categories are not well separated in the*
184 *learner’s representational space such that they make incorrect predictions based on whichever*
185 *task-irrelevant statistics their representations reflect. The task will seem hard and performance*
186 *will be low. We can thus assume an equivalence between alignment in representations, accuracy*
187 *of predictions and task difficulty. (b) Self-reinforcement of predictions (adapted from [19]). When*
188 *a stimulus is observed without supervision, an appropriate response is predicted and*
189 *subsequently self-reinforced. This results in changes in the representations and predictions. (c) If*
190 *prior representations and predictions are sufficiently aligned with the task, self-reinforcement*
191 *leads to performance improvement. In the case of misalignment, self-reinforcement has*
192 *detrimental or no effect on performance. This results in a snowball effect, the course of which*
193 *can only be changed if supervision is provided to correct mistakes and align representations with*
194 *the task.*

195

196 **Box 1: Theoretical principles predict unsupervised snowball effect**

197 We propose that human predictions self-reinforce in the absence of supervision. Since self-
198 reinforcement simply snowballs prior learning, it can help or hurt performance depending on
199 whether predictions and their underlying representations align with the task. Unsupervised
200 learning only succeeds in tasks aligned with the learner’s representations.

201 This intuitive reasoning is supported by the theoretical and computational principles that allow
202 unsupervised and semi-supervised machine learning algorithms to be successful. Inevitably,
203 unsupervised learning can only recover ground-truth structure in the data if this structure is
204 reflected in salient data statistics. For example, for clustering to work, similar points must
205 belong to the same cluster and dissimilar points must belong to different clusters (this is known
206 as the cluster assumption, [95]). In other words, clusters need to be sufficiently easy to tell apart
207 to be accurately recovered. In the same way, successful semi-supervised learning requires the
208 to-be-learned input classes to be sufficiently distinctive to work effectively [96–98,3]. Because
209 this is not always guaranteed, and in practice is often difficult to validate, unsupervised data is

210 not guaranteed to boost an algorithm’s supervised performance. In fact, much of the success of
211 semi-supervised machine learning could be due to standard-practice data curation that removes
212 difficult datapoints from unsupervised training with the effect that input classes become more
213 distinct [99]. Thus, while learning from unlabeled data has led to the much-reported
214 performance boosts in machine learning, they can also lead to degradation. In fact, reports of
215 performance degradation following the addition of unsupervised data exist and are likely under-
216 reported [3].

217 Returning to empirical studies, sufficient cluster “distinctiveness” may appear to be a theoretical
218 prerequisite that is easy enough to control experimentally to assess successful, rather than
219 detrimental, unsupervised learning. However, there is a subtle, yet crucial, twist: while
220 experimental tasks may appear to comply with the prerequisite in the experimenter-defined
221 input space, they can simultaneously violate it in the space relevant for learning which is not
222 routinely assessed -- the learner’s internal representations of the input space (Fig. 2a). When
223 overlooked, equivocal results about the benefit of unsupervised experiences can appear
224 conflicting when, in fact, they are predictable. To understand whether results conflict or are
225 simply evidence for the varied directions unsupervised self-reinforcement can take, the
226 alignment between internal representations and experimenter-defined task needs to be
227 considered.

228

229 **Representation-to-task alignment determines efficacy of unsupervised learning**

230 Representation-to-task alignment is a theoretical concept capturing how well a learner’s
231 representations set them up for learning in a new task. Alignment is sufficient when task-
232 relevant statistics are prominent in the representations (e.g., well-separated clusters), when
233 only adaptation of existing representations is needed (e.g., repositioning cluster centers), or
234 both when a beneficial learning sequence builds on prominent representations and
235 subsequently adapts them (e.g., an easy-to-hard curriculum). In these cases, performance is
236 high, and tasks are easy (Fig. 2a). Because representation-to-task alignment is independent of
237 any specific type of representation or task, we can expect to observe its effects on the efficacy

238 of unsupervised learning across all types of learning. Here, we will test this prediction against
239 the evidence from different, independent literatures.

240 ***Perceptual and category learning***

241 Perceptual and category learning experiments share many methodological commonalities.
242 Perceptual learning investigates how perception is changed because of experience with sensory
243 inputs, like the ability to distinguish different line lengths. This fundamental form of learning is
244 often studied by manipulating simple, physical stimulus dimensions like line length. Category
245 learning investigates the process of assigning labels (or other distinct responses) to groups of
246 inputs, such as assigning either “sheep” or “goat” to each input. This is often studied by
247 manipulating stimulus distributions and boundaries defining categories within them. Stimuli can
248 range from simple shapes or sounds, akin to those used in perceptual learning, to complex,
249 high-dimensional artificial objects. In both paradigms, learners are usually presented with
250 stimuli on a trial-by-trial basis and respond by guessing category membership, or in the case of
251 perceptual learning, making a same-different judgment between two stimuli.

252 The perceptual learning literature has extensively studied the effect of different forms of
253 supervision [25,26] and thus serves as a superb source of evidence on the effectiveness of
254 unsupervised learning. Results can be summarized simply: unsupervised perceptual training can
255 help in some, but not all, tasks. It does this in a way that correlates with task difficulty, as
256 predicted by our representation-to-task alignment view that requires sufficient class separation
257 or convenient presentation order. Concretely, unsupervised learning helps if the task is easy and
258 training accuracy is high, as predicted for aligned tasks, [27] or if high-accuracy, easy trials
259 precede or are interleaved with low-accuracy, difficult trials [28–30]. By contrast, feedback
260 appears necessary for learning when task difficulty is high and initial performance is low, as
261 predicted for misaligned tasks [31,27].

262 Unsupervised and semi-supervised categorization studies in adults echo results from perceptual
263 learning: unsupervised experiences facilitate learning in easier tasks, but not in more difficult
264 ones [9,10]. Learning to separate low-variability categories is easy (aligned task) and equally
265 effective with or without feedback, whereas learning to separate high-variability categories is

266 hard (misaligned task) and requires feedback [32]. Extending this finding, category learning is
267 known to be influenced by the degree of within-category variance [8], with unsupervised
268 learning being most effective and robust when categories are statistically dense and category
269 separation is large [33–35]. This further indicates that sufficient class separation is necessary for
270 successful unsupervised learning (Box 1).

271 Moreover, unsupervised experiences can have both beneficial and detrimental effects in the
272 exact same task, depending on the alignment of a learner’s representations [14]. This pattern is
273 also reflected across tasks. In simple category structures, where stimuli vary along a single
274 dimension, learners can recover categories [20] or shift previously supervised category
275 boundaries without feedback [36–40]. By contrast, in two-dimensional tasks, subjects appear
276 unable to recover categories without feedback [7] and the addition of unsupervised experiences
277 does not boost supervised performance [11,12] except under limiting conditions [41–43]. While
278 experimenter-defined task dimensionality does not imply task difficulty per se, in these
279 experiments, representations required to succeed in the two-dimensional tasks were
280 unmistakably less obvious compared to those required for the one-dimensional tasks. In line
281 with these results, prior knowledge relevant to the task can enhance unsupervised learning
282 [44].

283 This pattern of results is echoed in language acquisition. When learning nonnative phonetic
284 contrasts, unsupervised exposure has been shown to be unsuccessful unless it is complemented
285 by sufficient supervised learning [45] or if it only involves shifting boundaries of existing
286 phonetic contrasts [46] or if phonetic contrasts are made distinctive [47,48]. We can rephrase
287 these results within our perspective: Learning new phonetic contrasts is challenging due to their
288 misalignment with the native speech sound space. To make unsupervised exposure succeed, the
289 task needs to be simplified either by providing feedback that fosters the formation of more
290 aligned representations, by changing the task to only involve modulation of existing, sufficiently
291 aligned representations, or by amplifying the to-be-learned contrast as a form of class
292 separation. Similarly, unsupervised exposure to an artificial language leads to simple word
293 learning whereas learning its complex syntactic regularities requires feedback [49]. Further,
294 research on infants’ capacity to integrate labeled and unlabeled exposure to new categories

295 indicates that learning is successful only when labels are introduced initially, but not when they
296 are presented at the end or omitted entirely [50,51]. This lends credence to our prediction that
297 supervision is required to transition from a misaligned to an aligned representational space
298 before unsupervised experiences can improve performance. A study investigating children's
299 acquisition of linguistic category labels revealed that unsupervised exposure to structured,
300 straightforward labels (regular plural nouns) impaired performance on unstructured, difficult
301 labels (irregular plural nouns) among younger, error-prone children who have not yet mastered
302 the regularities and irregularities. Conversely, it boosted performance among older, more
303 proficient children capable of making adequate predictions [15,16]. This underscores that the
304 outcomes of unsupervised training can vary within the same task, contingent on the learners'
305 representations.

306 Pre-exposure studies assess the impact of initial unsupervised exposure on later supervised
307 learning and have received independent attention. The effects of pre-exposure vary with
308 category structures [13] with improvements seen for statistically dense categories [52] and
309 exposure to easy stimuli [53,54]. This is in line with our perspective: unsupervised pre-exposure
310 helps in easy tasks but does not affect, or even hinders, difficult ones. Interestingly, rat studies
311 show the opposite (Box 2). This discrepancy is likely due to humans' ability to reason about
312 tasks [55].

313 ***Selective Feedback***

314 Real-world feedback is selective and action-dependent which can lead to learning traps due to
315 unchallenged false predictions [56]. For example, a negative first impression may deter future
316 interactions, preventing the revision of potentially false initial impressions [57]. Similarly,
317 stereotyping can be perpetuated by initial negative experiences with a group, leading to future
318 avoidance. This selective information sampling prevents updating of false predictions about
319 group members and the likelihood of future avoidance increases when predictions are made
320 without feedback [18]. Consequently, stereotyping intensifies over time, with untested
321 predictions often misremembered as validated [19]. In this way, the selective-feedback
322 literature highlights the detrimental effects of unsupervised learning when predictions are
323 misaligned with reality, as seen in stereotyping.

324 **Expertise**

325 So far, we have seen that unsupervised learning effects vary in controlled lab studies. To gauge
326 whether this generalizes to real-world learning, we can assess uncontrolled, long-term learning.
327 Expertise is the product of extensive learning from varying quantity and quality of supervisory
328 signals outside the lab. For instance, radiologists initially receive supervised training but later
329 get less feedback, often not knowing if their diagnoses were correct. If unsupervised
330 experiences had only beneficial effects, we would expect performance to improve over time,
331 leading to expertise even without supervision. However, this prediction has received substantial
332 opposition [58–62] and has even led critics to claim that “At best, experience is an uncertain
333 predictor of degree of expertise. At worst, experience reflects seniority – and little more.” [60].
334 Biases, a form of prior expectations, can distort learning and hinder steady improvement
335 through experience. For instance, confirmation bias gives more weight to information that
336 aligns with learners’ expectations, skewing learning away from actual evidence [63,64]. In other
337 circumstances, learners may attribute their failure to external factors instead of modifying their
338 erroneous behavior so that performance deteriorates [22].

339 Irrespective of how expert performance is reached, the expertise literature supports, on a more
340 general level, the claim that unsupervised experiences alone do not guarantee improvement.
341 Instead, reliable improvement seems to require rapid and regular feedback on decisions [62].
342 Because acquiring expertise is not easy, but involves learning new skills beyond prior
343 knowledge, these results fit well with our representation-to-task alignment perspective. This is
344 further supported by work showing that initial feedback and guidance are crucial for skill
345 learning [65]. For instance, an in-lab study shows that withdrawing feedback early in motor skill
346 learning, when errors are high (inaccurate predictions), causes performance to deteriorate,
347 whereas doing so later, when errors are low (accurate predictions), enables the skill to be
348 maintained or improved [17].

349

350 **Box 2: Results requiring further attention**

351 **Pre-exposure in rodents.** Interestingly, the effects of unsupervised pre-exposure in rodents are
352 found to be the opposite of those observed in humans. Rodent studies show that unsupervised
353 learning benefits are greater when stimuli are perceptually similar and thus hard to discriminate
354 [100]. Conversely, rodent learning can be hindered when the stimuli are perceptually distinct
355 and thus easier to discriminate [101,102]. This effect is attributed to a combination of two
356 learning principles: unsupervised differentiation, which refines representations over time, and
357 latent inhibition, which reduces the associability between inputs and a response [102]. In this
358 context, latent inhibition could explain the slower learning seen after exposure to stimuli that
359 are easily distinguishable.

360 The opposing effects observed in animals and humans could be due to humans' awareness of
361 their participation in an experiment, leading to heightened attention to stimuli and potential
362 weakening of latent inhibition [55,103]. This is supported by the reversal of pre-exposure effects
363 in rats when using hedonic stimuli which are believed to stimulate attention [104]. Moreover,
364 interleaving unsupervised and supervised trials in mice appears more effective than
365 unsupervised pre-exposure [105], potentially also modulated by attentional factors.

366 **Blocked testing effects.** Understanding learning is important, but it is also important to examine
367 how learning could be helped. Across domains, research on optimal training schedules shows
368 that interleaving supervised training with blocks of unsupervised testing consistently improves
369 human learning compared to no testing or restudying of materials. It helps learning and
370 retention of materials preceding or following testing [106,107] and even replacing interim active
371 testing with passive exposure improves performance [108,109]. While individual studies
372 highlight the benefits of supervised testing, particularly its ability to correct inaccuracies and
373 confirm low-confidence predictions [110], a meta-analysis reveals unsupervised testing benefits
374 are comparable [111]. Taken together, these results appear to suggest that unsupervised testing
375 is exclusively beneficial, a finding that would contradict our unsupervised snowballing theory.
376 However, occasional evidence of performance interactions with learner proficiency and
377 confidence suggest representation-to-task alignment effects may be at play and could simply
378 have gone underreported.

379

380 **Self-reinforcement underlies unsupervised learning**

381 While representation-to-task alignment can predict the effectiveness of unsupervised learning,
382 it does not provide a mechanism. A number of specific learning procedures have been explored
383 in this context, all of which have self-reinforcement at their core, where learning uses the
384 system's own predictions in lieu of ground-truth supervision, snowballing existing learning
385 without altering its direction.

386 *Perceptual learning, category learning and expertise*

387 The perceptual learning literature not only supports representation-to-task alignment, but also
388 offers strong evidence for unsupervised self-reinforcement, formalized by Hebbian learning
389 models. Unsupervised Hebbian learning can improve or degrade performance depending on
390 how well representations serve learning a task [23,24]. A Hebbian model that learns from both
391 unsupervised and supervised experiences by adapting representations and their associations
392 with responses [66,67] is successful in accounting for a broad range of results [27]. While trial-
393 by-trial category learning is only rarely modelled, self-reinforcement models have demonstrated
394 their ability to account for semi-supervised categorization [68,14] and can also predict
395 unsupervised learning trajectories in children acquiring linguistic labels [16]. In expertise
396 studies, computational work is limited. However, theories of closed-loop motor skill learning
397 suggest internal estimates guide learning in the absence of feedback leading to either
398 performance gains or decrements [69].

399 *Selective Feedback*

400 As described earlier, false predictions that remain unchallenged can, for example, lead to the
401 perpetuation of stereotypes. This can be accounted for by models employing unsupervised self-
402 reinforcement [18,19]. Predictions also remain unchallenged when some actions are never
403 followed by feedback (i.e., unsupervised actions). Here, the same self-reinforcement can be
404 observed: humans learn from their own predictions as if they received validation for it
405 (constructivist coding hypothesis, [70–72]) which can be modeled by a self-reinforcement
406 mechanism [71].

407 *Internal feedback signals*

408 Self-reinforcement requires internal learning signals independent of external supervision. While
409 the neural mechanisms involved in external supervision (or at least rewards and punishments)
410 are fairly well understood [73], knowledge of the brain's self-generated feedback signals is
411 limited. Recent studies indicate that brain areas active during external feedback processing are
412 also active when feedback is inferred [74–76]. Moreover, choice consistency and subjective
413 confidence increase in the absence of feedback reflecting self-reinforcement [77] which is in line
414 with evidence that chosen actions carry more internal weight than unchosen ones [78].
415 Subjective rewards can also self-reinforce choices [79]. Large-scale, real-world studies indicate
416 that this can cause people to fall into a learning trap, ceasing exploration and exploiting even
417 when better options exist [80], which an error-driven learning model can account for by aligning
418 subjective preferences with past choices [81]. Neuroimaging also shows that preferences are
419 updated online and only for remembered choices [82]. Moreover, replay, another active
420 research area, involves a form of self-reinforcement in which the brain rehearsed past
421 experiences through offline neural reactivation [83,84]. Overall, research supports the brain's
422 use of unsupervised self-reinforcement mechanisms, with internal signals like confidence
423 playing a key role when feedback is absent.

424

425 **Concluding remarks**

426 In summary, studies across different literatures and learning domains support our perspective:
427 Humans self-reinforce their predictions in the absence of supervision, which can either help or
428 hurt performance depending on the alignment between the learner's representations and the
429 task. While we focused on studies testing unsupervised learning under controlled conditions,
430 the expertise literature suggests that these considerations are also relevant to naturalistic
431 settings. This shift in perspective resolves the paradox to predict learning successes and failures
432 in the lab, and fundamentally alters what we expect from unsupervised learning. Unsupervised
433 learning may not be the knight that battles to save us when we lack supervision; instead, it
434 appears to wield a double-edged sword. This raises new questions and lays the foundation for

435 future research on the role of supervision in learning that will have implications for the design
436 of instruction and learning over the lifespan (see Outstanding Questions).

437 A key implication of this perspective is that a deeper understanding of unsupervised learning
438 requires consideration of the alignment between mental representation and task. This is
439 challenging because alignment depends on specific stimuli, task structures, and learners'
440 representations. Efficiently assessing and modeling alignment to account for individual tasks
441 and learners is an important future direction that can build on recent advances [85–88]. In fact,
442 assessing alignment is also important for predicting supervised learning [89,90], memory [91]
443 and perception [92] which suggests that it also applies to naturalistic, large-scale unsupervised
444 learning. Future models need to make explicit the concrete relationship between alignment and
445 learning and be constrained by neural evidence on biologically supported mechanisms [93].

446 Our efforts to understand when unsupervised learning succeeds and fails have illuminated the
447 rich interconnections between historically separate research areas that can be leveraged in
448 future studies. Beyond the topics discussed here, relevant research also encompasses areas like
449 attention [94] and training schedules (Box 2). Linking results across these domains promotes a
450 more rigorous examination of learning principles.

451 Future research should also go beyond the traditional approach of studying unsupervised
452 learning in isolation. To understand why humans manage to learn despite all difficulties, we
453 need to explore how supervised and unsupervised learning mechanisms interact and relate to
454 feedback sources more akin to reinforcement, self-supervised, or sequential learning that are
455 blended in modern machine learning systems. Crucially, future work should explore how
456 unsupervised self-reinforcement and learning from (self-)supervisory signals coexist in humans,
457 who may use one general-purpose mechanism instead of different special-purpose algorithms
458 like machines. This crosstalk could lead to a more holistic theory of human learning, which is
459 important for understanding real-world learning, like the acquisition of expertise.

460 In conclusion, we advocate for an interdisciplinary approach to studying the mechanisms of
461 unsupervised learning and the broader role of supervision which should integrate
462 representational and neural constraints. This new direction contributes to our understanding of

463 learning fundamentals and can improve the design of instructional systems that better support
464 learning across the lifespan to prevent us from mistaking goats for sheep with ever greater
465 confidence.

466 **Outstanding questions**

467 What exactly is the quantitative relationship between representation-to-task alignment and
468 learning? How does this relate to different sources of the problem, e.g. poor extraction of
469 relevant features versus good feature extraction, but poor cluster separation? How does this
470 relate to different timescales, e.g. short-term learning to direct attention versus long-term
471 representational change?

472 How much representation-to-task alignment is needed for unsupervised learning to help?

473 How can we measure representation-to-task alignment? How can we incorporate
474 representation-to-task alignment into computational models of learning?

475 Does representation-to-task alignment affect supervised and unsupervised learning differently?

476 How is self-reinforcement implemented by the brain? Which role does meta-cognition play in
477 this? Is it affected by brain development?

478 Does self-reinforcement affect supervised learning too?

479 How do supervised and unsupervised learning interact? Are they fundamentally different or can
480 they be unified?

481 How does learning from other feedback signals, like reward, compare with supervised and
482 unsupervised learning?

483 How does unsupervised learning compare in humans and animals? Are there differences
484 between implicit / subconscious and deliberate / conscious unsupervised learning?

485 Which other factors related to the presence and absence of supervision, like motivation, affect
486 learning?

487 How does the sequential order (e.g., blocked supervised and unsupervised exposure) affect
488 unsupervised learning?

489

490

491 **Glossary**

492 **Learning algorithm:** The specific algorithm used to maximize a task objective which can be supervised or
493 unsupervised. A supervised algorithm (e.g., a standard neural network) learns an input / stimulus to
494 output / response mapping and uses supervision to improve its predictions. Apart from solving
495 supervised learning problems, supervised algorithms can also be used to tackle unsupervised learning
496 problems by adapting the task objective (e.g., self-supervised learning). An unsupervised algorithm (e.g.,
497 a standard Bayesian Graphical Model) extracts information from the inputs / stimuli without accessing
498 ground-truth supervision. Unsupervised algorithms are designed to solve unsupervised problems but
499 can also be adapted to tackle supervised learning problems.

500 **Learning problem:** The type of learning problem that is partially defined by the data, especially whether
501 supervision is available or not (i.e., supervised or unsupervised learning problem).

502 **Learning task:** The specific task (or task objective) that is defined within a learning problem and which
503 can be supervised or unsupervised. In experimental studies, the task objective derives from the
504 experimenter-defined stimulus-response-mapping.

505 **Representation-to-task alignment:** The degree to which the internal representations of a learning
506 system create a similarity space that suggests an input / stimulus to output / response mapping that is in
507 agreement with the objective mapping defined by the task.

508 **Self-reinforcement:** A mechanism by which a system learns from its own predictions in lieu of ground-
509 truth supervision. This has the effect that existing predictions from inputs / stimuli to outputs /
510 responses are strengthened. This mechanism can, in principle, be implemented both by supervised and
511 unsupervised algorithms. This mechanism is popular in semi-supervised machine learning and called
512 self-training or pseudo-labelling.

513 **Self-supervised learning:** A machine learning approach that solves an unsupervised learning problem by
514 turning it into a supervised task so that a supervised algorithm can be applied. Since no external
515 supervision is available, supervision is created directly from the unsupervised data.

516 **Semi-supervised learning:** Learning in a problem / task that offers a mixture of supervised and
517 unsupervised inputs / stimuli.

518 **Supervised learning:** Learning in a problem / task that requires the learning of an input / stimulus to
519 output / response mapping and in which ground-truth supervision is available.

520 **Supervision / Feedback:** In machine learning, supervision is defined as the delivery of ground-truth
521 outputs (e.g., labels) following some inputs (e.g., images). In human learning studies, supervision more
522 often refers to the delivery of corrective feedback (e.g., correct / incorrect response) on their response
523 to some preceding stimulus.

524 **Unsupervised learning:** Learning in a problem / task without supervision, simply through extraction of
525 information from the observation of inputs / stimuli.

526

527 **Acknowledgements**

528 This work was supported by the Gatsby Charitable Foundation (FB); Neuroscience Institute of
529 Carnegie Mellon University (FB); the Max Planck Society (FB, PD); the Alexander von Humboldt
530 Foundation (PD); ESRC (ES/W007347/1) and a Royal Society Wolfson Fellowship (18302) to BCL;
531 NSF BCS2420979 and NSF BCS2346989 to LLH.

532 Parts of this work have also been described in the PhD thesis of the first author (F. Bröker, PhD
533 thesis, University College London, 2022).

534 During the preparation of this work the authors used Microsoft Copilot in order to shorten the
535 text and improve the writing. After using this tool/service, the authors reviewed and edited the
536 content as needed and take full responsibility for the content of the publication.

537

538 **References**

- 539 1. Tian, K. *et al.* (2017) Deepcluster: A general clustering framework based on deep learning. in
540 *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2017,*
541 *Skopje, Macedonia, September 18–22, 2017, Proceedings, Part II 17*, pp. 809–825
- 542 2. Devlin, J. *et al.* (2019) BERT: Pre-training of Deep Bidirectional Transformers for Language
543 Understanding. in *Proceedings of the 2019 Conference of the North American Chapter of the*
544 *Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and*
545 *Short Papers)*, Minneapolis, Minnesota, pp. 4171–4186
- 546 3. Van Engelen, J.E. and Hoos, H.H. (2020) A survey on semi-supervised learning. *Mach. Learn.* 109,
547 373–440
- 548 4. Lee, T.S. and Mumford, D. (2003) Hierarchical Bayesian inference in the visual cortex. *JOSA A* 20,
549 1434–1448
- 550 5. Knill, D.C. and Pouget, A. (2004) The Bayesian brain: the role of uncertainty in neural coding and
551 computation. *TRENDS Neurosci.* 27, 712–719
- 552 6. Tenenbaum, J.B. *et al.* (2011) How to grow a mind: Statistics, structure, and abstraction. *Science*
553 331, 1279–1285
- 554 7. Ashby, F.G. *et al.* (1999) On the dominance of unidimensional rules in unsupervised categorization.
555 *Percept. Psychophys.* 61, 1178–1199
- 556 8. Ell, S.W. *et al.* (2012) Unsupervised category learning with integral-dimension stimuli. *Q. J. Exp.*
557 *Psychol.* 65, 1537–1562
- 558 9. Wade, T. and Holt, L.L. (2005) Incidental categorization of spectrally complex non-invariant auditory
559 stimuli in a computer game task. *J. Acoust. Soc. Am.* 118, 2618–2633
- 560 10. Emberson, L.L. *et al.* (2013) Is statistical learning constrained by lower level perceptual
561 organization? *Cognition* 128, 82–102
- 562 11. Vandist, K. *et al.* (2009) Semisupervised category learning: The impact of feedback in learning the
563 information-integration task. *Atten. Percept. Psychophys.* 71, 328–341
- 564 12. McDonnell, J.V. *et al.* (2012) Sparse category labels obstruct generalization of category
565 membership. in *Proceedings of the Annual Meeting of the Cognitive Science Society*, 34
- 566 13. Wills, A.J. *et al.* (2004) The Role of Category Structure in Determining the Effects of Stimulus
567 Preexposure on Categorization Accuracy. *Q. J. Exp. Psychol. Sect. B* 57, 79–88
- 568 14. Bröker, F. *et al.* (2022) When unsupervised training benefits category learning. *Cognition* 221,
569 104984
- 570 15. Ramscar, M. and Yarlett, D. (2007) Linguistic Self-Correction in the Absence of Feedback: A New
571 Approach to the Logical Problem of Language Acquisition. *Cogn. Sci.* 31, 927–960
- 572 16. Ramscar, M. *et al.* (2013) Error and expectation in language learning: The curious absence of
573 mouses in adult speech. *Language* 89, 760–793
- 574 17. Newell, K.M. (1974) Knowledge of results and motor learning. *J. Mot. Behav.* 6, 235–244
- 575 18. Allidina, S. and Cunningham, W.A. (2021) Avoidance begets avoidance: A computational account of
576 negative stereotype persistence. *J. Exp. Psychol. Gen.* 150, 2078–2099
- 577 19. Cox, W.T.L. *et al.* (2022) Untested assumptions perpetuate stereotyping: Learning in the absence of
578 evidence. *J. Exp. Soc. Psychol.* 102, 104380
- 579 20. Rosenthal, O. *et al.* (2001) Forming classes by stimulus frequency: Behavior and theory. *Proc. Natl.*
580 *Acad. Sci.* 98, 4265–4270
- 581 21. Love, B.C. (2002) Comparing supervised and unsupervised category learning. *Psychon. Bull. Rev.* 9,
582 829–835
- 583 22. Kc, D. *et al.* (2013) Learning from my success and from others' failure: Evidence from minimally
584 invasive cardiac surgery. *Manag. Sci.* 59, 2435–2449

- 585 23. McClelland, J.L. (2001) Failures to learn and their remediation: A Hebbian account. In *Mechanisms*
586 *of cognitive development*, pp. 109–134, Psychology Press
- 587 24. McClelland, J.L. (2006) How far can you go with Hebbian learning, and when does it lead you
588 astray. *Process. Change Brain Cogn. Dev. Atten. Perform. Xxi* 21, 33–69
- 589 25. Doshier, B.A. and Lu, Z.-L. (2009) Hebbian reweighting on stable representations in perceptual
590 learning. *Learn. Percept.* 1, 37–58
- 591 26. Doshier, B.A. and Lu, Z.-L. (2017) Visual Perceptual Learning and Models. *Annu. Rev. Vis. Sci.* 3, 343–
592 363
- 593 27. Liu, J. *et al.* (2010) Augmented Hebbian reweighting: Interactions between feedback and training
594 accuracy in perceptual learning. *J. Vis.* 10, 29–29
- 595 28. Ahissar, M. and Hochstein, S. (1997) Task difficulty and the specificity of perceptual learning.
596 *Nature* 387, 401–406
- 597 29. Liu, J. *et al.* (2012) Mixed training at high and low accuracy levels leads to perceptual learning
598 without feedback. *Vision Res.* 61, 15–24
- 599 30. Asher, J.M. and Hibbard, P.B. (2020) No effect of feedback, level of processing or stimulus
600 presentation protocol on perceptual learning when easy and difficult trials are interleaved. *Vision*
601 *Res.* 176, 100–117
- 602 31. Shiu, L.-P. and Pashler, H. (1992) Improvement in line orientation discrimination is retinally local but
603 dependent on cognitive set. *Percept. Psychophys.* 52, 582–588
- 604 32. Homa, D. and Cultice, J.C. (1984) Role of feedback, category size, and stimulus distortion on the
605 acquisition and utilization of ill-defined categories. *J. Exp. Psychol. Learn. Mem. Cogn.* 10, 83
- 606 33. Kloos, H. and Sloutsky, V.M. (2008) What’s behind different kinds of kinds: effects of statistical
607 density on learning and representation of categories. *J. Exp. Psychol. Gen.* 137, 52
- 608 34. Pothos, E.M. *et al.* (2011) Measuring category intuitiveness in unconstrained categorization tasks.
609 *Cognition* 121, 83–100
- 610 35. Vong, W.K. *et al.* (2016) The helpfulness of category labels in semi-supervised learning depends on
611 category structure. *Psychon. Bull. Rev.* 23, 230–238
- 612 36. Zhu, X. *et al.* (2007) Humans perform semi-supervised classification too. in *AAAI*, 2007, pp. 864–
613 870
- 614 37. Lake, B. and McClelland, J. (2011) Estimating the strength of unlabeled information during semi-
615 supervised learning. in *Proceedings of the Annual Meeting of the Cognitive Science Society*, 33
- 616 38. Kalish, C.W. *et al.* (2011) Can semi-supervised learning explain incorrect beliefs about categories?
617 *Cognition* 120, 106–118
- 618 39. Kalish, C.W. *et al.* (2015) Drift in children’s categories: when experienced distributions conflict with
619 prior learning. *Dev. Sci.* 18, 940–956
- 620 40. Gibson, B.R. *et al.* (2015) What causes category-shifting in human semi-supervised learning? in
621 *Proceedings of the Annual Meeting of the Cognitive Science Society*
- 622 41. Rogers, T. *et al.* (2010) Semi-supervised learning is observed in a speeded but not an unspeeded 2D
623 categorization task. in *Proceedings of the Annual Meeting of the Cognitive Science Society*, 32
- 624 42. Rogers, T. *et al.* (2010) Humans Learn Using Manifolds, Reluctantly. in *Advances in Neural*
625 *Information Processing Systems*, 23
- 626 43. Vandist, K. *et al.* (2019) Semisupervised category learning facilitates the development of
627 automaticity. *Atten. Percept. Psychophys.* 81, 137–157
- 628 44. Clapper, J.P. (2007) Prior knowledge and correlational structure in unsupervised learning. *Can. J.*
629 *Exp. Psychol. Rev. Can. Psychol. Expérimentale* 61, 109–127
- 630 45. Wright, B.A. *et al.* (2019) Semi-supervised learning of a nonnative phonetic contrast: How much
631 feedback is enough? *Atten. Percept. Psychophys.* 81, 927–934

- 632 46. Chládková, K. *et al.* (2022) Unattended distributional training can shift phoneme boundaries. *Biling.*
633 *Lang. Cogn.* 25, 827–840
- 634 47. McCandliss, B.D. *et al.* (2002) Success and failure in teaching the [r]-[l] contrast to Japanese adults:
635 Tests of a Hebbian model of plasticity and stabilization in spoken language perception. *Cogn. Affect.*
636 *Behav. Neurosci.* 2, 89–108
- 637 48. Escudero, P. *et al.* (2011) Enhanced bimodal distributions facilitate the learning of second language
638 vowels. *J. Acoust. Soc. Am.* 130, EL206–EL212
- 639 49. Frinsel, F.F. *et al.* (2024) The Role of Feedback in the Statistical Learning of Language-Like
640 Regularities. *Cogn. Sci.* 48, e13419
- 641 50. LaTourrette, A. and Waxman, S.R. (2019) A little labeling goes a long way: Semi-supervised learning
642 in infancy. *Dev. Sci.* 22
- 643 51. LaTourrette, A. and Waxman, S.R. (2022) Sparse labels, no problems: Infant categorization under
644 challenging conditions. *Child Dev.* 93, 1903–1911
- 645 52. Unger, L. and Sloutsky, V.M. (2022) Ready to Learn: Incidental Exposure Fosters Category Learning.
646 *Psychol. Sci.* 33, 999–1019
- 647 53. Milton, F. *et al.* (2014) The effect of pre-exposure on family resemblance categorization for stimuli
648 of varying levels of perceptual difficulty. in *Proceedings of the Annual Meeting of the Cognitive*
649 *Science Society*, 36
- 650 54. Milton, F. *et al.* (2020) The effect of preexposure on overall similarity categorization. *J. Exp. Psychol.*
651 *Anim. Learn. Cogn.* 46, 65–82
- 652 55. Angulo, R. *et al.* (2019) Stimulus comparison: Effects of the pre-exposure schedule and instructions
653 for perceptual learning and attention. *Learn. Motiv.* 65, 20–32
- 654 56. Rich, A.S. and Gureckis, T.M. (2018) The limits of learning: Exploration, generalization, and the
655 development of learning traps. *J. Exp. Psychol. Gen.* 147, 1553
- 656 57. Denrell, J. (2005) Why most people disapprove of me: experience sampling in impression
657 formation. *Psychol. Rev.* 112, 951
- 658 58. Brehmer, B. (1980) In one word: Not from experience. *Acta Psychol. (Amst.)* 45, 223–241
- 659 59. Garb, H.N. (1989) Clinical judgment, clinical training, and professional experience. *Psychol. Bull.*
660 105, 387
- 661 60. Shanteau, J. *et al.* (2002) Performance-based assessment of expertise: How to decide if someone is
662 an expert or not. *Eur. J. Oper. Res.* 136, 253–263
- 663 61. Ericsson, K.A. (2004) Deliberate practice and the acquisition and maintenance of expert
664 performance in medicine and related domains. *Acad. Med.* 79, S70–S81
- 665 62. Kahneman, D. and Klein, G. (2009) Conditions for intuitive expertise: a failure to disagree. *Am.*
666 *Psychol.* 64, 515
- 667 63. Wason, P.C. (1960) On the failure to eliminate hypotheses in a conceptual task. *Q. J. Exp. Psychol.*
668 12, 129–140
- 669 64. Nickerson, R.S. (1998) Confirmation bias: A ubiquitous phenomenon in many guises. *Rev. Gen.*
670 *Psychol.* 2, 175–220
- 671 65. Dunphy, B.C. and Williamson, S.L. (2004) In pursuit of expertise. Toward an educational model for
672 expertise development. *Adv. Health Sci. Educ.* 9, 107–127
- 673 66. Petrov, A.A. *et al.* (2005) The dynamics of perceptual learning: an incremental reweighting model.
674 *Psychol. Rev.* 112, 715
- 675 67. Petrov, A.A. *et al.* (2006) Perceptual learning without feedback in non-stationary contexts: Data and
676 model. *Vision Res.* 46, 3177–3197
- 677 68. Gibson, B.R. *et al.* (2013) Human Semi-Supervised Learning. *Top. Cogn. Sci.* 5, 132–172
- 678 69. Adams, J.A. (1971) A closed-loop theory of motor learning. *J. Mot. Behav.* 3, 111–150

- 679 70. Elwin, E. *et al.* (2007) Constructivist Coding: Learning From Selective Feedback. *Psychol. Sci.* 18,
680 105–110
- 681 71. Henriksson, M.P. *et al.* (2010) What is coded into memory in the absence of outcome feedback? *J.*
682 *Exp. Psychol. Learn. Mem. Cogn.* 36, 1–16
- 683 72. Elwin, E. (2013) Living and Learning: Reproducing Beliefs in Selective Experience: Living and
684 Learning. *J. Behav. Decis. Mak.* 26, 327–337
- 685 73. Schultz, W. (2007) Behavioral dopamine signals. *Trends Neurosci.* 30, 203–210
- 686 74. Daniel, R. and Pollmann, S. (2012) Striatal activations signal prediction errors on confidence in the
687 absence of external feedback. *NeuroImage* 59, 3457–3467
- 688 75. Guggenmos, M. *et al.* (2016) Mesolimbic confidence signals guide perceptual learning in the
689 absence of external feedback. *eLife* 5, e13388
- 690 76. Rouault, M. *et al.* (2019) Forming global estimates of self-performance from local confidence. *Nat.*
691 *Commun.* 10, 1141
- 692 77. Ptasczynski, L.E. *et al.* (2022) The value of confidence: Confidence prediction errors drive value-
693 based learning in the absence of external feedback. *PLOS Comput. Biol.* 18, e1010580
- 694 78. Sakamoto, Y. and Miyoshi, K. (2024) A confidence framing effect: Flexible use of evidence in
695 metacognitive monitoring. *Conscious. Cogn.* 118, 103636
- 696 79. Vinckier, F. *et al.* (2019) Sour grapes and sweet victories: How actions shape preferences. *PLOS*
697 *Comput. Biol.* 15, e1006499
- 698 80. Riefer, P.S. *et al.* (2017) Coherency-maximizing exploration in the supermarket. *Nat. Hum. Behav.* 1,
699 0017
- 700 81. Hornsby, A.N. and Love, B.C. (2020) How decisions and the desire for coherency shape subjective
701 preferences over time. *Cognition* 200, 104244
- 702 82. Voigt, K. *et al.* (2019) Hard decisions shape the neural coding of preferences. *J. Neurosci.* 39, 718–
703 726
- 704 83. McClelland, J.L. *et al.* (1995) Why there are complementary learning systems in the hippocampus
705 and neocortex: insights from the successes and failures of connectionist models of learning and
706 memory. *Psychol. Rev.* 102, 419
- 707 84. Barry, D.N. and Love, B.C. (2022) A neural network account of memory replay and knowledge
708 consolidation. *Cereb. Cortex* 33, 83–95
- 709 85. Houlshby, N.M. *et al.* (2013) Cognitive tomography reveals complex, task-independent mental
710 representations. *Curr. Biol.* 23, 2169–2175
- 711 86. Hebart, M.N. *et al.* (2020) Revealing the multidimensional mental representations of natural
712 objects underlying human similarity judgements. *Nat. Hum. Behav.* 4, 1173–1185
- 713 87. Ma, W.J. and Peters, B. (2020) A neural network walks into a lab: towards using deep nets as
714 models for human behavior. *arXiv* at <<https://doi.org/10.48550/arXiv.2005.02181>>
- 715 88. Roads, B.D. and Love, B.C. (2021) Enriching imagenet with human similarity judgments and
716 psychological embeddings. in *Proceedings of the IEEE/CVF conference on computer vision and*
717 *pattern recognition*, pp. 3547–3557
- 718 89. Aho, K. *et al.* (2022) System alignment supports cross-domain learning and zero-shot
719 generalisation. *Cognition* 227, 105200
- 720 90. Roark, C.L. *et al.* (2022) A neural network model of the effect of prior experience with regularities
721 on subsequent category learning. *Cognition* 222, 104997
- 722 91. Schurgin, M.W. *et al.* (2020) Psychophysical scaling reveals a unified theory of visual memory
723 strength. *Nat. Hum. Behav.* 4, 1156–1172
- 724 92. Zaman, J. *et al.* (2021) Perceptual variability: Implications for learning and generalization. *Psychon.*
725 *Bull. Rev.* 28, 1–19
- 726 93. Golub, M.D. *et al.* (2018) Learning by neural reassociation. *Nat. Neurosci.* 21, 607–616

- 727 94. Hammer, R. *et al.* (2015) Feature saliency and feedback information interactively impact visual
728 category learning. *Front. Psychol.* 6
- 729 95. Chapelle, O. *et al.* (2006) *Semi-supervised learning*, The MIT Press
- 730 96. Singh, A. *et al.* (2008) Unlabeled data: Now it helps, now it doesn't. in *Advances in Neural*
731 *Information Processing Systems*, 21
- 732 97. Zhu, X. and Goldberg, A.B. (2009) *Introduction to Semi-Supervised Learning*, Springer International
733 Publishing
- 734 98. Oymak, S. and Gulcu, T.C. (2020) Statistical and Algorithmic Insights for Semi-supervised Learning
735 with Self-training. *arXiv* at <<http://arxiv.org/abs/2006.11006>>
- 736 99. Ganev, S. and Aitchison, L. (2021) Semi-supervised learning objectives as log-likelihoods in a
737 generative model of data curation. *arXiv* at <<https://doi.org/10.48550/arXiv.2008.05913>>
- 738 100. Oswalt, R.M. (1972) Relationship between level of visual pattern difficulty during rearing and
739 subsequent discrimination in rats. *J. Comp. Physiol. Psychol.* 81, 122
- 740 101. Chamizo, V.D. and Mackintosh, N. (1989) Latent learning and latent inhibition in maze
741 discriminations. *Q. J. Exp. Psychol.* 41, 21–31
- 742 102. Saksida, L.M. (1999) Effects of similarity and experience on discrimination learning: a
743 nonassociative connectionist model of perceptual learning. *J. Exp. Psychol. Anim. Behav. Process.*
744 25, 308
- 745 103. Graham, S. and McLaren, I. (1998) Retardation in human discrimination learning as a consequence
746 of pre-exposure: Latent inhibition or negative priming? *Q. J. Exp. Psychol. Sect. B* 51, 155–172
- 747 104. Sanjuán, M. del C. *et al.* (2014) An easy-to-hard effect after nonreinforced preexposure in a
748 sweetness discrimination. *Learn. Behav.* 42, 209–214
- 749 105. Schmid, C. *et al.* (2024) Passive exposure to task-relevant stimuli enhances categorization learning.
750 *eLife* 12, RP88406
- 751 106. Lee, H.S. and Ahn, D. (2018) Testing prepares students to learn better: The forward effect of testing
752 in category learning. *J. Educ. Psychol.* 110, 203–217
- 753 107. Yang, C. and Shanks, D.R. (2018) The forward testing effect: Interim testing enhances inductive
754 learning. *J. Exp. Psychol. Learn. Mem. Cogn.* 44, 485–492
- 755 108. Wright, B.A. *et al.* (2010) Enhancing Perceptual Learning by Combining Practice with Periods of
756 Additional Sensory Stimulation. *J. Neurosci.* 30, 12868–12877
- 757 109. Wright, B.A. *et al.* (2015) Enhancing speech learning by combining task practice with periods of
758 stimulus exposure without practice. *J. Acoust. Soc. Am.* 138, 928–937
- 759 110. Wang, L. and Yang, J. (2021) Effect of feedback type on enhancing subsequent memory: Interaction
760 with initial correctness and confidence level. *Psych J.* 10, 751–766
- 761 111. Adesope, O.O. *et al.* (2017) Rethinking the use of tests: A meta-analysis of practice testing. *Rev.*
762 *Educ. Res.* 87, 659–701
- 763