

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Computational evidence for hierarchically structured reinforcement learning in humans

### Permalink

<https://escholarship.org/uc/item/04h4d08z>

### Journal

Proceedings of the National Academy of Sciences of the United States of America,  
117(47)

### ISSN

0027-8424

### Authors

Eckstein, Maria K  
Collins, Anne GE

### Publication Date

2020-11-24

### DOI

10.1073/pnas.1912330117

Peer reviewed



# Computational evidence for hierarchically structured reinforcement learning in humans

Maria K. Eckstein<sup>a</sup>  and Anne G. E. Collins<sup>a,1</sup> 

<sup>a</sup>Department of Psychology, University of California, Berkeley, CA 94704

Edited by Danielle S. Bassett, University of Pennsylvania, Philadelphia, PA, and accepted by Editorial Board Member Dale Purves April 15, 2020 (received for review August 2, 2019)

**Humans have the fascinating ability to achieve goals in a complex and constantly changing world, still surpassing modern machine-learning algorithms in terms of flexibility and learning speed. It is generally accepted that a crucial factor for this ability is the use of abstract, hierarchical representations, which employ structure in the environment to guide learning and decision making. Nevertheless, how we create and use these hierarchical representations is poorly understood. This study presents evidence that human behavior can be characterized as hierarchical reinforcement learning (RL). We designed an experiment to test specific predictions of hierarchical RL using a series of subtasks in the realm of context-based learning and observed several behavioral markers of hierarchical RL, such as asymmetric switch costs between changes in higher-level versus lower-level features, faster learning in higher-valued compared to lower-valued contexts, and preference for higher-valued compared to lower-valued contexts. We replicated these results across three independent samples. We simulated three models—a classic RL, a hierarchical RL, and a hierarchical Bayesian model—and compared their behavior to human results. While the flat RL model captured some aspects of participants' sensitivity to outcome values, and the hierarchical Bayesian model captured some markers of transfer, only hierarchical RL accounted for all patterns observed in human behavior. This work shows that hierarchical RL, a biologically inspired and computationally simple algorithm, can capture human behavior in complex, hierarchical environments and opens the avenue for future research in this field.**

computational modeling | reinforcement learning | hierarchy | structure learning | task-sets

**R**esearch in the cognitive sciences has long highlighted the importance of hierarchical representations for intelligent behavior, in domains including perception (1), learning and decision making (2, 3), planning and problem solving (4), cognitive control (5), and creativity (6), among many others (7, 8). The common thread across all these domains is the insight that hierarchical representations (i.e., the simultaneous representation of information at different levels of abstraction) allow humans to behave adaptively and flexibly in complex, high-dimensional, and ever-changing environments. Exhaustive nonhierarchical (“flat”) representations, in contrast, are insufficient to achieve human-like behaviors.

To illustrate, consider the following situation. Mornings in your office, your colleagues are working silently, or quietly discussing work-related topics. After work, they are laughing and chatting loudly at their favorite bar. In this example, a context change induced a drastic change in behavior, despite the same interaction partners (i.e., “stimuli”). Hierarchical theories of cognition capture this behavior by positing that we learn strategies hierarchically, activating different behavioral strategies (or “task-sets”) in different contexts. Although hierarchical representations can incur additional cognitive cost (9), they provide a range of advantages compared to exhaustive flat representations: once a task-set has been selected (e.g., office), attention can be focused on a subset of environmen-

tal features (e.g., just the interaction partner) (10–13). When new contexts are encountered (e.g., new workplace, new bar), entire task-sets can be reused, allowing for generalization (6, 14, 15). Old skills are not catastrophically forgotten (16). In addition, hierarchical representations deal elegantly with incomplete information, for example, when contexts are unobservable (6, 17). All of these advantages are evident in the current study.

Although we know that hierarchical representations are essential for flexible behavior, how humans create these representations and how they learn to use them is still poorly understood. Here, we hypothesize that learning and using hierarchical representations can be explained under a hierarchical reinforcement learning (RL) framework, in which simple RL computations are combined to simultaneously operate at different levels of abstraction. RL theory (18) formalizes how to adjust behavior based on feedback in order to maximize rewards. Standard RL algorithms estimate how much reward to expect when selecting actions in response to stimuli and use these “action-value” estimates to select actions. Old action-values are updated in proportion to the “reward-prediction error,” the discrepancy between action-values and received reward, to produce increasingly accurate estimates. Such “flat” RL algorithms operate over unstructured, exhaustive representations (*SI Appendix, Fig. 3A*), converge to optimal behavior, are computationally inexpensive, and have led to recent breakthroughs in artificial intelligence (AI; ref. 18).

Broad evidence suggests that the brain implements computations similar to RL: dopamine neurons generate reward-prediction errors (19, 20), and a widespread network of frontal cortical regions (21) and basal ganglia (22, 23) represents action-values. Specific brain circuits thereby form “RL loops” (17, 24), in which learning is implemented through the continuous updating of action-values (11, 25). In this sense, estimating

This paper results from the Arthur M. Sackler Colloquium of the National Academy of Sciences, “Brain Produces Mind by Modeling,” held May 1–3, 2019, at the Arnold and Mabel Beckman Center of the National Academies of Sciences and Engineering in Irvine, CA. NAS colloquia began in 1991 and have been published in PNAS since 1995. From February 2001 through May 2019, colloquia were supported by a generous gift from The Dame Jillian and Dr. Arthur M. Sackler Foundation for the Arts, Sciences, & Humanities, in memory of Dame Sackler’s husband, Arthur M. Sackler. The complete program and video recordings of most presentations are available on the NAS website at <http://www.nasonline.org/brain-produces-mind-by>.

Author contributions: M.K.E. and A.G.E.C. designed research; M.K.E. performed research; M.K.E. analyzed data; and M.K.E. and A.G.E.C. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission. D.S.B. is a guest editor invited by the Editorial Board.

Published under the [PNAS license](#).

Data deposition: The data have been deposited in the National Institute of Mental Health Data Archive at <https://dx.doi.org/10.15154/1518660>. Analysis and modeling code is available on GitHub (<https://github.com/MariaEckstein/TaskSets>).

<sup>1</sup> To whom correspondence may be addressed. Email: [annecollins@berkeley.edu](mailto:annecollins@berkeley.edu).

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1912330117/-DCSupplemental>.

First published November 23, 2020.

action-values via RL is an algorithm of special interest to cognition: there is strong evidence that the brain implements a simple mechanism to perform the necessary computations. Nevertheless, RL algorithms have important shortcomings: they suffer from the curse of dimensionality (an exponential drop in learning speed with increasing problem complexity); they lack flexibility for behavioral change; and they cannot easily generalize or transfer old knowledge to new situations. Hierarchical RL (26) mitigates these shortcomings by nesting RL processes at different levels of temporal (27–29) or state abstraction (12, 30).

Recent research has provided support for a plausible implementation of hierarchical RL in the brain: the neural circuit that implements RL is multiplexed, such that distinct RL loops operate at different levels of abstraction along the rostrocaudal axis (10, 24, 31–37). Consistent with this architecture, recent studies have shown signatures of RL values and reward-prediction errors at different levels of abstraction in the human brain (29, 38). However, previous studies did not provide evidence that neural signatures of hierarchical value support learning and generalizing hierarchically structured behavior. Thus, it remains unknown whether hierarchical RL indeed supports hierarchical behavior. The goal of this study is to fill this gap. We investigate hierarchical RL in a paradigm that promotes the creation and reuse of hierarchical structure. We provide a fully fledged computational model that accounts for behavior across a variety of relevant situations: context-dependent learning, context switches, generalization to new contexts, partially observable problems, and choices at different levels of abstraction. This study tests all predictions of hierarchical RL in a single paradigm. Because hierarchical RL makes specific behavioral predictions in each situation, we are able to test the model qualitatively against human behavior (39). We then compare our hierarchical RL model quantitatively to the two most relevant competing models, flat RL and hierarchical Bayes. The former employs RL but without hierarchical structure. The latter assumes that high-level decisions are based on Bayesian inference of task-set reliability, rather than RL using task-set values (14).

In the following, we first introduce our hierarchical RL model and experimental paradigm. We then test whether humans show qualitative behaviors that are predicted by the hierarchical RL model, as well as the two competing models. We first show evidence for hierarchical representations in humans, as predicted by both hierarchical RL and hierarchical Bayes, but not flat RL. We employ multiple independent analyses, including switch-cost measures and positive and negative transfer. We then provide evidence for human hierarchical value learning, which is only consistent with the hierarchical RL model. We next provide quantitative support for these qualitative results and show that

model comparison supports the hierarchical RL model over flat RL and hierarchical Bayes. The majority of results replicates across three independent participant samples.

## Results

**Computational Models.** Our hierarchical RL model is composed of two hierarchically structured RL processes. The high-level process manages behavior at the abstract level by acquiring a “policy over policies”: it learns which task-set to choose in each context, using “task-set values” (the estimated expected reward of selecting a task-set in a given context). The low-level process acquires these task-sets: it learns which actions to choose in response to each stimulus by estimating “action-values” (the estimated expected reward of selecting an action for a given stimulus, within a specific task-set; Fig. 1A).

At the beginning of the task, task-sets and actions are picked randomly, but over time, trial-and-error learning leads to the formation of meaningful task-sets, which represent policies that are specialized for particular contexts. Trial-and-error learning also underlies the policy over task-sets that determines which task-set is selected in each context. Thus, our hierarchical RL model is based on two nested processes, which create an interplay between learning stimulus-action associations (low level) and context-task-set associations (high level). *SI Appendix, Fig. 4* shows a step-by-step visualization of this model.

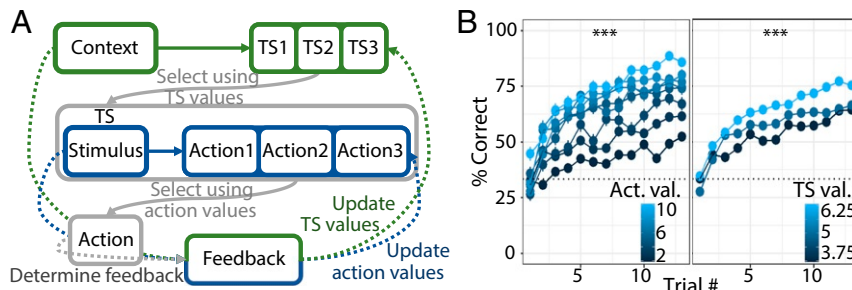
Formally, to select an action  $a$  in response to stimulus  $s$  in context  $c$ , hierarchical RL goes through a two-step process: 1) it selects a task-set  $TS$  based task-set values in the current context,  $Q(TS|c)$ , using  $p(TS|c) = \frac{\exp(Q(TS|c))}{\sum_{TS_i} \exp(\beta_{TS} Q(TS_i|c))}$ . The inverse temperature  $\beta_{TS}$  captures task-set choice stochasticity (Fig. 1A). The chosen task-set  $TS$  provides a set of action-values  $Q(a|s, TS)$ , which are used to 2) select an action  $a$ , according to  $p(a|s, TS) = \frac{\exp(Q(a|s, TS))}{\sum_{a_i} \exp(\beta_a Q(a_i|s, TS))}$ , where  $\beta_a$  captures action choice stochasticity (Fig. 1A; for trial-by-trial behavior, see *SI Appendix, Fig. 4B*). After executing action  $a$  on trial  $t$ , feedback  $r_t$  reflects the continuous amount of reward received, which guides learning at both levels of abstraction (i.e., to update the values of the selected task-set and action):

$$Q_{t+1}(TS|c) = Q_t(TS|c) + \alpha_{TS} (r_t - Q_t(TS|c))$$

$$Q_{t+1}(a|s, TS) = Q_t(a|s, TS) + \alpha_a (r_t - Q_t(a|s, TS)).$$

$\alpha_{TS}$  and  $\alpha_a$  are learning rates at the levels of task-sets and actions (Fig. 1A and *SI Appendix, Fig. 4C*).

The flat RL model uses the same mechanism for value learning and action selection but lacks hierarchical structure: it treats



**Fig. 1.** (A) Schematic of the hierarchical RL model. A high-level RL loop (green) selects a task-set  $TS$  in response to the observed context, using  $TS$  values. The chosen task-set provides action-values, based on which the low-level RL loop (blue) selects an action in response to the observed stimulus. Task-set and action-values are both learned based on action feedback. (B) Human learning curves during the initial learning phase, averaged over blocks. Colors denote underlying action-values (Left) and task-set values (Right), respectively. Stars show that both affect performance (*Learning Curves and Effects of Reward*), consistent with hierarchical RL.  $***P < 0.001$ .

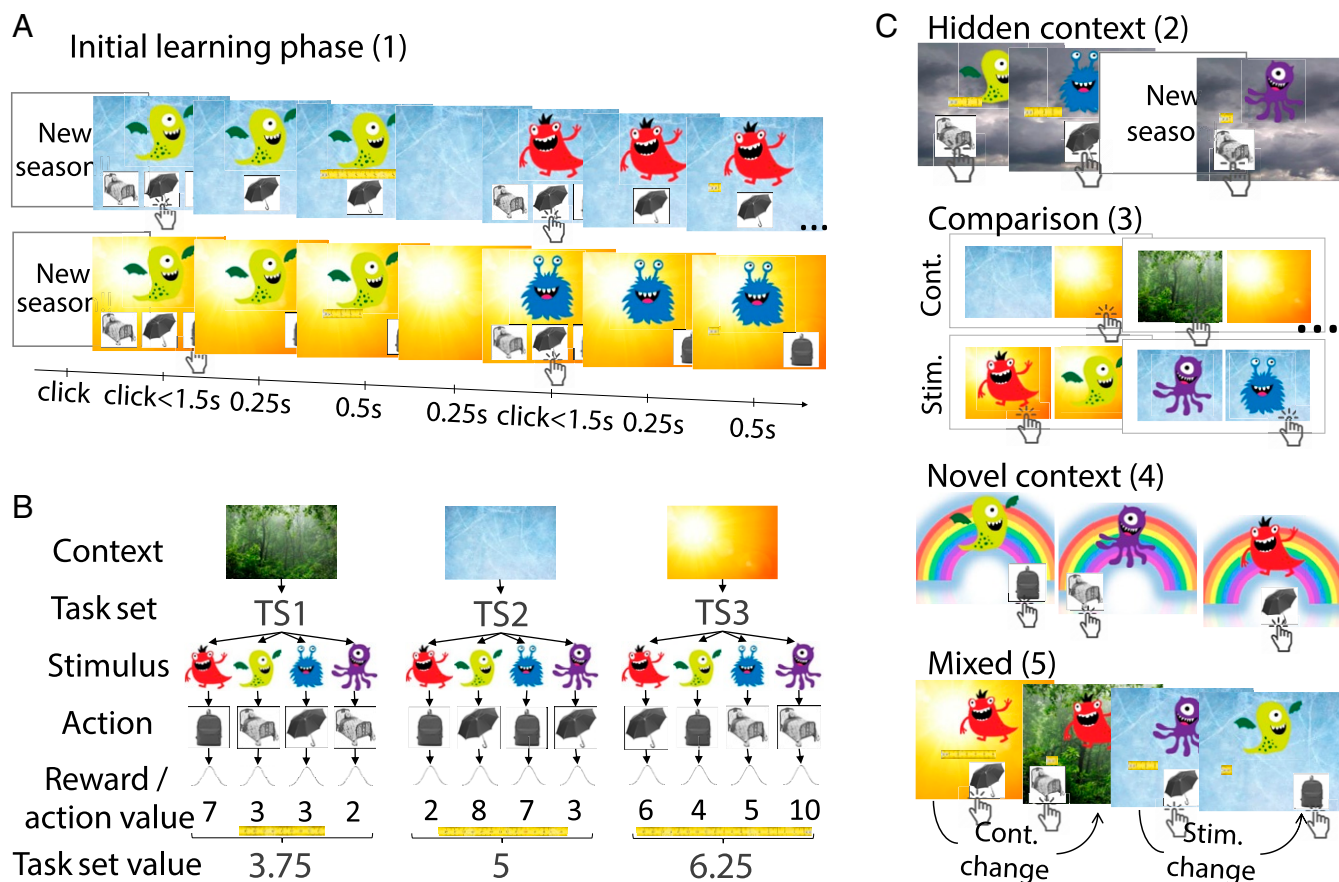
each combination of context and stimulus as a unique state (*Methods*). The hierarchical Bayesian model creates a task-set structure like hierarchical RL but selects task-sets according to their inferred reliability, rather than task-set values (*Methods*).

**Task Design.** We designed a task in which participants learned to select the correct actions for different stimuli (Fig. 2A). The mapping between stimuli and actions varied across three contexts, creating three distinct task-sets (Fig. 2B). Each context appeared in three blocks of 52 trials, for a total of 9 blocks. Contexts differed in average rewards, allowing us to test for RL values at the level of task-sets. After an initial-learning phase of this task (Fig. 2A), participants completed four test phases (Fig. 2C) to hone in on specific predictions of hierarchical RL. Detailed information about the task is provided in Fig. 2, *Methods*, and *SI Appendix*.

**Learning Curves and Effects of Reward.** As expected, participants' performance increased within a block, showing adaptation to context changes (Fig. 1B). We also verified that participants were sensitive to continuous differences in reward magnitudes (tape length). RL predicts better performance for larger rewards

because these lead to larger action-values, which make correct actions more distinguishable from incorrect ones (*SI Appendix*, Fig. 4B for details). Participants indeed showed better performance for high-reward stimuli (Fig. 1B, *Left*). This effect was predicted by both hierarchical and flat RL. Hierarchical RL additionally predicts better performance for high-valued contexts: larger rewards create larger reward-prediction errors at the task-set level, which allow for better discrimination between correct and incorrect task-sets and lead to better task-set selection and performance (see *SI Appendix*, Fig. 4A for details). As predicted, participants also showed an effect of task-set values on performance (Fig. 1B, *Right*).

To quantify both effects, we conducted a mixed-effects logistic regression model predicting trialwise accuracy from action-values, task-set values, and their interaction (fixed effects), specifying participants, trial, and block as random effects. We approximated action-values as average stimulus–action rewards and task-set values as average context–task-set rewards, as shown in Fig. 2B. The model revealed significant effects of both action-values ( $\beta = 0.38$ ,  $P < 0.001$ ) and task-set values ( $\beta = 0.20$ ,  $P < 0.001$ ) on performance (for complete statistics and results in other samples, see *SI Appendix*, Table 1). This provides initial



**Fig. 2.** Task design. (A) In the initial-learning phase, participants saw one of four stimuli (aliens) in one of three contexts (seasons) and had to find the correct action (item) through trial and error. Each context had a different mapping between stimuli and correct actions, and contexts were presented blockwise. Feedback indicated correctness deterministically, but different context–stimulus–action combinations led to different rewards (with Gaussian noise). (B) Example mapping between stimuli and actions for each context, defining three task-set TS1 to TS3. Average rewards (task-set values) differed between contexts. All actions and stimuli had equal average rewards. (C) Additional test phases. The hidden-context phase, presented after initial learning, was identical except that contexts were unobservable (season hidden by clouds). This allowed us to test whether participants reactivated previously learned task-sets. In the comparison phase, participants saw either two contexts (“Cont.”) or two stimuli (“Stim.”) on each trial and selected their preferred one. We used subjective preferences to assess task-set values (contexts) and action-values (stimuli). The novel-context phase was similar to initial learning but had a new context and no feedback, to test how participants generalized previous knowledge to new situations. The final mixed phase was similar to initial learning but not blocked, i.e., both stimuli and contexts could change on every trial, to test for asymmetric switch costs. All test phases were separated by “refresher blocks” similar to initial learning, to alleviate carry-over effects and forgetting.

evidence that human choices were sensitive to RL values at two levels of abstraction—actions and task-sets—as predicted by hierarchical RL.

**Hierarchical Representation.** We tested participants’ abstractions in more detail using three independent analyses: switch costs in the mixed phase of the task, reactivation of task-sets in the hidden-context phase, and task-set selection errors during initial learning.

**Asymmetric Switch Costs.** Asymmetric switch costs can be evidence for hierarchical representations because changes across trials are more challenging at higher than lower levels of abstraction (40, 41). For example, switching contexts is more cognitively costly than switching stimuli within a context. To test for such asymmetries in our paradigm, we compared trials on which a different stimulus was presented from that on the previous trial (but the same context) with those on which a different context was presented (but the same stimulus), using the mixed phase (Fig. 2C). As expected, participants responded significantly slower after context switches than after stimulus switches ( $t(25) = 3.47, P = 0.002$ ). This was not due to participants’ initial surprise about the interleaved presentation of contexts in the mixed phase, as the result held throughout the phase (SI Appendix). Asymmetric switch costs therefore suggest that participants created hierarchical representations, nesting stimuli within contexts, as predicted by hierarchical RL and hierarchical Bayes.

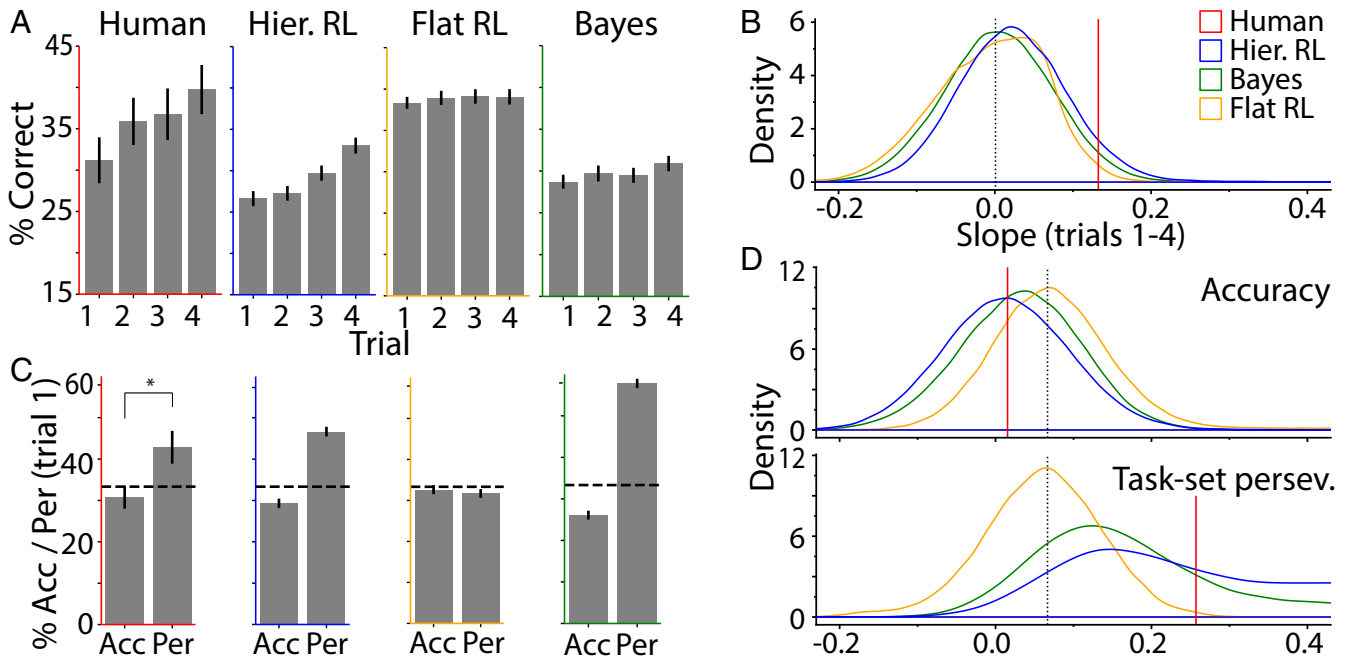
**Reactivating Task-Sets.** Did representing the task hierarchically benefit performance (e.g., did it support positive transfer)? In the hidden-context phase of our task, contexts were not observable, such that participants could either relearn old stimulus–action mappings from scratch (no transfer) or reactivate previous

task-sets, with the correct mappings already in place (transfer). By enabling reactivation of old task-sets, hierarchy has been shown previously to enable better performance and faster learning (6, 14, 34).

If participants reactivated task-sets, we expect a specific pattern of performance in the hidden-context phase, specifically on the first few trials after a context switch, before any stimulus is repeated: because every trial provides feedback about the appropriateness of the chosen task-set, task-set selection should become more accurate on each trial, and, consequently, accuracy should improve. If, on the other hand, participants did not use task-sets and instead relearned stimulus–response associations from scratch, as predicted by flat RL, performance can only increase after a stimulus is repeated. Because no stimuli are repeated until the fifth trial in our task, the first four trials provide the perfect testing ground to pitch these two predictions against each other, as illustrated in Fig. 3A: hierarchical RL simulations show increasing performance, whereas flat RL simulations show no change (simulation details in *Methods* and *Modeling Behavioral Patterns Jointly*).

Human behavior qualitatively matched the predictions of hierarchical RL: performance increased steadily over the first four trials after a context switch (Fig. 3A), evident in the significant correlation between item position (1 to 4) and performance ( $r = 0.19, P = 0.048$ ). This shows that participants recalled previously learned stimulus–action mappings rather than relearning them, a signature of task-set transfer.

We next assessed quantitatively which of our three candidate models captured this behavior best. We compared the models using Bayes Factors (BF), which we estimated using a method related to Approximate Bayesian Computation (*Methods*, SI Appendix, and ref. 42). Our method involved simulating



**Fig. 3.** (A and B) Participants reactivated task-sets in the hidden-context phase. (A) Human performance (Left, red) increased over the first four trials following a context switch, even though different stimuli were presented on each trial. The “best” (*Methods*) simulation based on the hierarchical RL model showed qualitatively similar behavior (blue). The effect was absent in the flat RL model (orange) and present but weaker in the Bayesian hierarchical model (green). (B) Slopes of the performance increase in A, as densities over 50,000 simulations per model, with parameters sampled uniformly at random. These densities approximate marginal model likelihoods for the calculation of Bayes factors. The densities of hierarchical RL and hierarchical Bayes were shifted toward larger slopes, making human-like performance slopes more likely. Dotted line indicates chance. (C and D) Task-set perseveration errors in the initial-learning phase. (C) Percent correct trials (“Acc”) and percent task-set perseveration errors (“Per”) on the first trial after a context switch. Humans (Left, red): star denotes significance in repeated-measures *t* test. Models: hierarchical RL and hierarchical Bayes, but not flat RL, qualitatively reproduced human behavior. (D) Accuracy and task-set perseveration errors for all simulations, as densities.

synthetic data from each model and estimating the likelihood of human behavior under the simulated data, as illustrated in Fig. 3B. Hierarchical RL surpassed both flat RL ( $BF = 5.12$ ) and also hierarchical Bayes ( $BF = 1.96$ ) in model comparison (*SI Appendix, Table 3*). This confirms the qualitative result, showing that human performance in the hidden-context phase was better captured by hierarchical than flat models.

**Task-Set Perseveration Errors.** We showed that hierarchy allowed for positive transfer, enabling participants to reactivate old task-sets. However, hierarchy can also lead to negative transfer: when participants select the wrong task-set, the “correct” action according to that task-set is likely to be incorrect in the current context. We call such errors “task-set selection errors,” and focus on a specific subtype, “task-set perseveration errors.” Here, actions are chosen that would have been correct in the previous context but are incorrect in the current one. Contrary to flat RL, hierarchical models predict task-set perseveration (*Methods* and example in *SI Appendix, Fig. 4A*), reflected in high proportions of task-set perseveration errors and low initial accuracy (Fig. 3C and D).

We tested this prediction on the first trial after each context switch during initial learning, and found that participants were more likely to make task-set perseveration errors than to select correct actions [ $t(25) = 2.1$ ,  $P = 0.046$ ], in accordance with hierarchical model simulations (Fig. 3D). Task-set perseveration persisted several trials into the new block, as evident in a logistic regression predicting task-set perseveration errors from trial index ( $\beta = -6.83\%$ ,  $z = -9.31$ ,  $P < 0.001$ ), task-set values ( $\beta = -2.43\%$ ,  $z = -1.00$ ,  $P < 0.001$ ), and action-values ( $\beta = -14.03\%$ ,  $z = -8.45$ ,  $P < 0.001$ ), controlling for block and specifying random effects of participants.

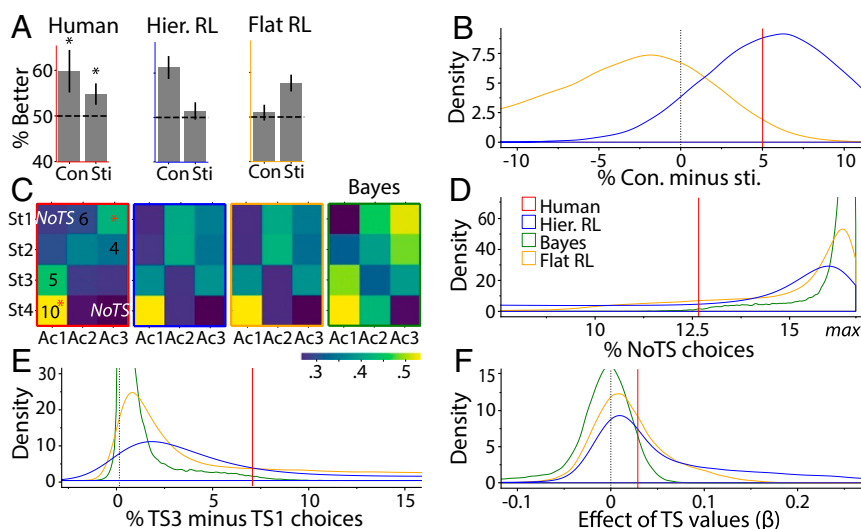
In summary, the presence of task-set perseveration errors in humans is qualitative evidence for hierarchical processing. Quantitative model comparison supports this conclusion, showing that hierarchical models fit human error patterns better than flat RL (hierarchical vs. flat RL:  $BF = 14.99$ ; hierarchical Bayes vs. flat RL:  $BF = 10.32$ ; hierarchical RL vs. Bayes:  $BF = 1.40$ ).

**RL Values at Different Levels of Abstraction.** Our results so far focused on hierarchical representations in general, showing that participants created, reactivated, and transferred task-sets. We now test predictions that are unique to hierarchical RL, assessing whether participants acquired RL values at the level of task-sets as well as actions.

**Task-Set Values Affect Subjective Preference.** A classic approach to assess RL values in humans is to investigate subjective preferences (43). To investigate whether participants acquired values at both levels, we thus used a comparison phase, where participants selected their preferred out of two items on each trial. Items were either two contexts or two stimuli—testing task-set and action-values, respectively (Fig. 2C).

The hierarchical RL model selected contexts based on the task-set values acquired during initial learning and showed a strong preference for high-valued over low-valued contexts (*SI Appendix, Fig. 7A*). The flat RL model selected contexts based on average action-values in this context and showed a much weaker preference (*SI Appendix, Fig. 7A*). The hierarchical Bayesian model did not track values over contexts and was thus not simulated in this phase. As predicted by hierarchical RL, participants preferred high-valued over low-valued contexts [ $t(25) = 2.56$ ,  $P = 0.017$ ], indicating RL values at the level of contexts. Quantitative model comparison (Fig. 4B) strongly favored hierarchical over flat RL ( $BF = 1,171.65$ ). For completeness, we also confirmed participants’ RL values at the level of stimuli, as predicted by both flat and hierarchical RL and evident in the preference for high-valued over low-valued stimuli [ $t(25) = 2.11$ ,  $P = 0.045$ ]. In conclusion, participants’ preferences were best accounted for by the hierarchical RL model.

We next investigated a different model prediction in the comparison phase: the hierarchical RL model takes two steps to retrieve action-values but only one to retrieve task-set values. This suggests stimulus selection should be slower and noisier than context selection. Flat RL, on the contrary, takes one step to retrieve action-values but multiple steps to calculate context-values, suggesting the inverse pattern. Humans showed



**Fig. 4.** Effects of task-set values on behavior. (A and B) Comparison phase. (A) Humans (red) performed better for contexts (“Con”) than stimuli (“Sti”; “% Better”: percentage choosing higher-valued alternative). Stars indicate significant difference from chance (dashed line). The hierarchical RL simulation showed the same qualitative pattern, whereas flat RL showed the opposite. (B) Difference between context and stimulus condition, as simulation-based densities. (C and E) Novel-context phase. (C) Raw action frequencies. “Ac1-3”, actions; “St1-3”, stimuli. Humans (red frame): overlaid numbers show action-values in TS3, the highest-valued task-set, which was chosen frequently. Red stars indicate actions that were correct in multiple task-sets, also selected frequently. “NoTS” indicates actions that were incorrect in all task-sets, selected rarely. Models: hierarchical (blue frame) and flat RL (orange frame) were qualitatively similar to humans; hierarchical Bayes (green frame) made different predictions. (D) NoTS choices in all simulations. (E) Difference between percentage of actions consistent with TS3 and TS1. (F) Initial-learning phase. Regression weights predicting performance from task-set values, showing that values at both levels affected performance more in the RL models.

the patterns predicted by hierarchical RL: response times (RTs) were numerically slower and performance was significantly worse for contexts than for stimuli [mixed-effects regression, RTs:  $\beta = 148.21$ ,  $t(25) = 1.63$ ,  $P = 0.12$ ; Accuracy:  $\beta = 0.28$ ,  $z = 2.0$ ,  $P = 0.048$ ; Fig. 4B]. Although the effect on RTs did not reach significance here, it was strongly significant in the replication (*SI Appendix, Table 1*). Quantitative model comparison strongly favored hierarchical over flat RL in terms of accuracy ( $BF = 39.64$ ).

**Task-Set Values Affect Performance.** As explained above, human initial learning was affected by both action-values and task-set values (Fig. 1B), in accordance with hierarchical RL. To compare our models in this regard, we calculated the effects of task-set values on performance, using a simplified regression model (*SI Appendix*). Supporting our qualitative findings, the hierarchical RL model provided a better fit than valueless hierarchical Bayes ( $BF = 6.62$ ) and, crucially, than flat RL ( $BF = 1.49$ ; Fig. 4F).

**Task-Set Values Affect Generalization.** We showed above that participants preferred high-valued over low-valued contexts (*SI Appendix, Fig. 7A*). We now test whether participants showed similar task-set preferences in the novel-context phase, that is, when generalizing old knowledge to a new context. For simulations, our hierarchical RL model applied its highest-valued task-set throughout the novel-context phase. The hierarchical Bayes model applied its most reliable task-set. The flat RL model chose actions based on average values (*Methods*).

We labeled each action in the novel-context phase as one of the following: correct in task-set TS3, TS2, TS1, both TS3 and TS1, both TS2 and TS1, or not correct in any task-set (NoTS). Despite the lack of feedback, human participants showed consistent preferences for certain stimulus-action combinations over others (Fig. 4C; see *SI Appendix, Fig. 2* for heat maps of task-set values). They chose NoTS actions less often than other actions, controlling for the frequency of each category [ $t(25) = 2.24$ ,  $P = 0.034$ ]. Mappings shared between multiple task-sets (TS2 and TS1; TS3 and TS1) were more frequent than mappings that only occurred in one task-set (TS1, TS2, TS3), controlling for chance level [ $t(25) = 2.83$ ,  $P = 0.0091$ ]. This confirms that participants reused old task-sets for new contexts, in accordance with our findings in the hidden-context phase and prior literature (17). Quantitative model comparison confirmed that the number of NoTS choices was captured better by hierarchical RL than by flat RL ( $BF = 1.78$ ) or hierarchical Bayes ( $BF = 45.60$ ).

Highlighting the role of task-set values, hierarchical RL predicted more actions from the highest-valued TS3 than from the lowest-valued TS1 and a greater difference between the two than flat RL or hierarchical Bayes (Fig. 4E). Humans showed the same pattern, selecting more TS3 than TS1 actions [ $t(25) = 2.58$ ,  $P = 0.016$ ]. Bayes Factors confirmed that this difference was captured better by hierarchical RL than flat RL ( $BF = 1.59$ ) or hierarchical Bayes ( $BF = 32.01$ ). Taken together, our hierarchical RL model captured both the reuse of old task-sets in new contexts and the preference for high-valued over low-valued task-sets.

**Modeling Behavioral Patterns Jointly.** Human behavior followed predictions of hierarchical RL qualitatively, and Bayes Factors confirmed quantitatively that this model fit better than the competing ones. However, we treated each behavioral measure independently. We next sought to confirm that it was possible to obtain all behavioral results simultaneously based on a single set of parameters. To this end, we chose one “best” set of parameters for each model (*Methods*) and showed the behavior of this simulation side-by-side with humans, for each behavioral measure. As expected, neither flat RL nor hierarchical Bayes replicated all qualitative patterns in Figs. 3 A and C and 4 A and C. How-

ever, importantly, a single set of parameters could capture all qualitative patterns in the hierarchical model. Note that because parameters were not obtained through model fitting, behavior can deviate quantitatively from human data.

## Discussion

The goal of the current study was to assess whether human flexible behavior could be explained by hierarchical RL (i.e., the concurrent use of RL at different levels of abstraction) (3, 44). We proposed a hierarchical RL model that acquires low-level strategies—or “task-sets”—using RL and that also learns to choose between these task-sets using RL. We contrasted this model with a flat RL model, to highlight the unique contribution of hierarchy, and to a hierarchical Bayesian model, to highlight the contribution of a hierarchical value representation.

Our hierarchical RL model predicted unique patterns of behavior in a variety of situations. To assess whether humans employed hierarchical RL, we designed a context-based learning task in which multiple subtasks tested these predictions. Indeed, participants’ behavior followed the predictions in all subtasks. The first prediction was that participants would create hierarchical representations. Several independent results supported this claim, including asymmetric switch costs, task-set perseveration errors, and task-set reactivation. These results could not be accounted for by the flat RL model but were also compatible with the hierarchical Bayesian model.

To address the unique predictions of hierarchical RL, we sought evidence of hierarchical values. Hierarchical RL predicts value-based 1) context preferences, 2) performance differences between contexts, and 3) generalization in new contexts. Human behavior showed the predicted patterns: 1) when asked to pick their preferred contexts, participants selected higher-valued ones more often. This suggests that they had formed abstract task-set values, in addition to low-level action-values. Participants also performed better when choosing between high-level contexts than low-level stimuli, in accordance with the “blessing of abstraction” (45, 46). 2) Task-set values affected performance, with better performance of higher-valued task-sets. This shows that hierarchical representation can explain performance differences between contexts. 3) When faced with a new context, participants reused previous task-sets, preferring higher-valued over lower-valued ones. This suggests that task-set values guided generalization of old knowledge to new situations.

In summary, human behavior showed all qualitative patterns predicted by hierarchical RL. To quantify the differences with hierarchical Bayes and flat RL, we conducted formal model comparison using Bayes Factors. Because marginal model likelihoods were intractable, we approximated them using simulations, similar to refs. 42, 47, and 48. Bayes Factors instantiate an implicit Occam’s razor that accounts for differences in model complexity, such as the larger number of parameters in the hierarchical models compared to flat RL, differences in the functional form of each model, and differences in parameter spaces (47, 49). In this way, Bayes Factors implement a more comprehensive tradeoff between parsimony and goodness-of-fit than traditional methods.

In our paradigm, Bayes Factors showed that hierarchical RL and hierarchical Bayes captured behavioral aspects of hierarchy better than flat RL (e.g., task-set reactivation, task-set perseveration), whereas flat RL and hierarchical RL captured value-based aspects better (e.g., value-based generalization, effects of values on performance). Furthermore, hierarchical RL uniquely captured the influence of two sets of values on behavior. Overall, Bayes Factors favored the hierarchical RL model over flat RL and hierarchical Bayes. Based on this quantitative confirmation, we next asked whether all results could be jointly observed when simulating the hierarchical RL model with a single set of parameters, to confirm that different parameters

were not responsible for different behaviors. We used simulation summary statistics to identify a “best” set of parameters for each model. Only the hierarchical RL simulation qualitatively replicated all human behaviors but not flat RL or hierarchical Bayes. This shows that seemingly different behaviors, including trial-and-error learning (initial-learning phase), “inference” of missing information (hidden-context phase), subjective preferences (comparison phase), and generalization (novel-context phase), can all be explained in the same overarching hierarchical RL framework.

Note that we have not explored the full space of possible models. In particular, it would be possible to construct a hierarchical Bayesian model that tracks task-set and action-values rather than their reliability but uses Bayesian inference rather than RL to perform updates. This model might capture the behavioral patterns we observed here. Indeed, our results show evidence for humans’ ability to track values at multiple levels of hierarchy in support of generalizable behavior but do not speak directly to the exact update process. However, we favor the hierarchical RL formulation of such updates because it is inspired by a rich literature on brain circuits that makes its implementation plausible, and because it is algorithmically simple, with the ability to account for complex cognitive processes.

Many computational models have addressed cognitive hierarchy. How are they related to our model? One important class of hierarchical models is purely Bayesian (7, 50, 51). These models aim to explain, on a computational level of analysis (52), the fundamental purpose of hierarchy for cognitive agents. Our model, on the other hand, is algorithmic, like many pure-RL models: it aims to describe dynamically which cognitive steps humans take when they make decisions in complex environments. Our model is also inspired by the structure of human neural learning circuits (24, 32, 35), thereby extending to the implementational level of analysis.

Some models of hierarchical cognition are method hybrids: some combine Bayesian inference at the abstract level with RL at the lower level (6, 10). Other, resource-rational models, combine Bayesian principles of rationality with cognitive constraints (53). Frank and Badre (10, 31) proposed a hybrid model that uses Bayesian inference to arbitrate between multiple types of hierarchy and flat RL. In general, hybrid models assume a role for Bayesian inference at higher levels of hierarchy, contrary to our hierarchical RL model. This is an important difference: hierarchical RL mimics a form of inference (for example, identifying the latent task-set at the beginning of a block; *SI Appendix, Results 2.1*) but cannot do it optimally. It is an important direction for future research to identify whether human behavior is suboptimal in the same way.

Computational models at different levels of analysis (52) are not mutually exclusive. Bayesian inference offers a perspective based on optimality, but it is often intractable and approximations are computationally expensive. RL, on the other hand, uses values to approximate expectations instead of calculating them exactly. Because of its relative computational simplicity, and because it is biologically well supported, RL has often been used as an algorithmic and implementational model. Recent research showed that a neural network implementing hierarchical RL approximated the results of Bayesian inference (17). In other words, hierarchical RL might allow for optimal behavior using simpler computations.

Hierarchical RL was initially proposed in AI (26, 54). A number of AI algorithms has recently been used to model human cognition as well (28, 29, 55, 56), showcasing how intertwined the two fields have become (18, 57, 58). Nevertheless, most hierarchical RL algorithms in AI focus on hierarchy over the time scale of choices (“temporal abstraction,” e.g., breaking up long-term goals into subgoals). Our hierarchical model, in contrast, focuses on “choice abstraction” (i.e., allowing choice at the level

of task-sets and motor actions), a still rare approach in AI (but see ref. 59).

To conclude, classic RL has been a powerful model for simple decision making in animals and humans, but it cannot explain hallmarks of intelligence like flexible behavioral change, continual learning, generalization, and inference of missing information. Recent advances in AI have proposed hierarchical RL as a solution to a number of such shortcomings, and we found that human behavior showed many signs of hierarchical RL, which were captured better by our hierarchical RL model than competing ones.

There is no debate that achieving goals and receiving punishment are some of the most fundamental motivators that shape our learning and decision making. Nevertheless, almost all decisions humans face pose more complex problems than what can be achieved by flat RL. Structured hierarchical representations have long been proposed as a solution to this problem, and our hierarchical RL model uses only simple RL computations, known to be implemented in our brains, to solve complex problems that have traditionally been tackled with intractable Bayesian inference. This research aims to model complex behaviors using neurally plausible algorithms and provides a step toward modeling human-level, everyday-life intelligence.

## Methods

**Participants.** We tested three independent groups of participants, with approval from University of California, Berkeley’s institutional review board. All were university students, gave written informed consent and received course credit for participation.

The pilot sample had 51 participants (26 women; mean age  $\pm$  SD:  $22.1 \pm 1.5$  y), 3 of whom were excluded due to past or present psychological or neurological disorders. Due to a technical error, data were not recorded in the comparison phase for this sample. The second and main sample had 31 participants (22 women; mean age  $\pm$  SD:  $20.9 \pm 2.1$  y), 4 of whom were excluded due to disorders, and 1 of whom was excluded because average performance in the initial-learning phase was below 35% (chance is 33%). We added the mixed testing phase for this sample. The third sample had 32 participants (15 women; mean age  $\pm$  SD:  $=20.8 \pm 5.0$  y), 2 of whom were excluded due to disorders. Five participants did not complete the experiment and were excluded when data were missing. The task was minimally adapted for electroencephalography. All statistical tests were conducted in all samples (*SI Appendix, Table 1 and Fig. 1*), and *SI Appendix* discusses sample differences in detail.

**Task Design.** Participants first received instructions and underwent the initial-learning phase of the task. The purpose of initial learning was for participants to acquire distinct task-sets (i.e., specific stimulus-action mappings for each context). We also used the initial-learning phase to test for the effects of action-values and task-set values on performance and to assess errors types predicted by hierarchical RL.

In the beginning, participants were instructed to “feed aliens to help them grow as much as possible.” A tutorial with instructed trials followed, and then participants practiced a simplified task without contexts: on each trial, participants saw one of four stimuli and selected one of three actions by pressing J, K, or L on the keyboard (Fig. 2A). Feedback was given in form of a measuring tape whose length indicated the amount of reward. Correct actions produced consistent long (mean: 5.0) and incorrect actions short tapes (mean: 1.0; Fig. 2). When no action was selected, participants were reminded to respond faster next time, and the trial was counted as missed. Participants received 10 training trials per stimulus (40 total), with a maximum response time of 3,000 ms. Order was pseudorandomized such that each stimulus appeared once in four trials, and the same stimulus never appeared twice in a row.

The initial-learning phase had the same structure as training, but stimuli were presented in one of three contexts, each with a unique mapping between stimuli and actions (Fig. 2B). The context remained the same for a block of 52 trials. At the end of a block, a context change was explicitly signaled, before the next block began with a new context. Participants went through 9 blocks (3 per context) for a total of 468 trials. Participants needed to respond within 1.5 s and then received reward. Rewards varied between 2 to 10 for correct actions (Fig. 2B); rewards for incorrect actions remained 1.



We chose these numbers to maximize differences between contexts, while controlling for differences between stimuli and actions. The hidden-context phase was identical to initial learning, and participants knew they would encounter the same contexts as before, but, this time, they were “hidden” (Fig. 2C). There were 9 blocks with 10 trials per stimulus per block (360 total). Context switches were signaled. The purpose of the comparison phase was to assess participants’ subjective preferences for contexts and stimuli, as estimates of their task-set and action-values. Participants were shown two contexts (context condition), or two stimuli in the same context (stimulus condition), and selected their preferred one (Fig. 2C). Participants saw each of 3 pairs of contexts 5 times and each of 18 pairs of stimuli 3 times, for a total of  $15 + 198 = 213$  trials. Participants had 3 s to respond.

The purpose of the novel-context phase was to probe generalization, specifically the reuse of old task-sets in a new context. This phase was identical to the initial-learning phase, except that it introduced a new context in extinction (i.e., without feedback) (Fig. 2C). Participants received 3 trials per stimulus (12 total). The purpose of the final mixed phase was to probe switch costs, assessing whether switching contexts was more costly than switching stimuli, indicating hierarchical representation. The mixed phase was identical to the initial-learning phase, except that contexts as well as stimuli could change on every trial. Participants received 3 blocks of 84 trials (252 total), each with 7 repetitions per stimulus–context combination. To alleviate carry-over effects and forgetting between test phases, we interleaved them with refresher blocks, shorter 120-trial versions of the initial-learning phase. More details on task design are provided in *SI Appendix*.

**Computational Models.** We will address in turn how each model behaves in each phase. During initial learning, the flat RL model implemented classic model-free (“delta-rule”) RL (18): it treated every combination of a context and a stimulus as a unique state and learned one RL value for each state and action, as visualized in *SI Appendix, Fig. 3A*. Using the notations introduced in *Results*, values were updated based on  $Q_{t+1}(a|s, c) = Q_t(a|s, c) + \alpha (r - Q_t(a|s, c))$ , and actions were selected based on  $p(a|s, c) = \frac{\exp(Q(a|s, c))}{\sum_{a'} \exp(\beta Q(a'|s, c))}$ .

The flat RL model acquired 36 action-values, based on 3 parameters ( $\alpha$ ,  $\beta$ , and  $f$ ), whereas the hierarchical RL model acquired 9 task-set values and 36 action-values (45 total), with 6 free parameters ( $\alpha_a$ ,  $\alpha_{TS}$ ,  $\beta_a$ ,  $\beta_{TS}$ ,  $f_a$ , and  $f_{TS}$ ; equations in *Results*). *SI Appendix Fig. 3* visualizes the difference between both models, and *SI Appendix, Fig. 4* explains hierarchical RL behavior trial-by-trial. The forgetting parameters  $f \in [f_a, f_{TS}]$  captured value decay in both models:  $Q_{t+1} = (1 - f) Q_t + f Q_{init}$ .

The hierarchical Bayes model also learned task-sets but acquired their action-values based on correct–incorrect rather than continuous feedback:  $Q_{t+1}(a|s, TS) = Q_t(a|s, TS) + \alpha (\text{correct} - Q_t(a|s, TS))$ . The main difference to hierarchical RL was the selection of task-sets: the Bayesian model chose task-sets based on estimated reliability rather than task-set values, using Bayes theorem to obtain task-set reliabilities:  $p_{t+1}(TS|c) = \frac{p(r|s, TS, a) p_t(TS|c)}{p(r|s, a)}$ , with  $p(r|s, TS, a) = Q(a|s, TS)$ . Another difference was that hierarchical RL updated  $Q(TS|c)$  only for the chosen task-set, whereas hierarchical Bayes kept  $p(TS|c)$  up-to-date at all times for all task-sets (6, 14).

$Q$  values for both models were initialized at the expected reward of chance performance,  $Q_{init} = 1.67$ . The subsequent testing phases started from the  $Q$  values obtained at the end of initial learning.

In the hidden-context phase, contexts were not shown, such that models could not directly reuse acquired values that depended on contexts [flat RL:  $Q(a|c, s)$ ; hierarchical RL:  $Q(TS|c)$ ; Bayes:  $p(TS|c)$ ]. All models instead initialized these values at  $Q_{init}$  after each context switch and then relearned them using the same update equations as before. For flat RL, this resulted in learning an entire new policy  $Q(a|c, s)$ . For hierarchical models, only high-

level information [ $Q(TS|c)$  for RL,  $p(TS|c)$  for Bayes] had to be relearned but not action-values  $Q(a|s, TS)$ . This ability to transfer learned values is one of the main advantages of hierarchy.

For the comparison phase, we only simulated RL models because the Bayesian model does not provide values at the level of contexts. To select between two stimuli, RL models first computed the “state value” (18) of each, based on action-values:  $V(c, s) = \max_a Q(a|c, s)$  (flat RL) and  $V(c, s) = \max_a Q(a|s, TS) p(TS|c)$ , where  $p(TS|c) = \text{softmax}(Q(TS|c))$  (hierarchical RL). Models then selected one stimulus based on a softmax over the two state values. To select between contexts, the hierarchical model repeated the same computation for task-set values:  $V(c) = \max_{TS} Q(TS|c)$ . The flat model, lacking task-set values, used averages over action-values to estimate context preferences on-the-fly:  $V(c) = \text{mean}_s V(c, s)$ .

In the novel-context phase, models were faced with a context for which they had not learned values. Flat RL used averages over previous action-values to choose  $Q(a|c_{new}, s) = \text{mean}_c Q(a|c, s)$ . Hierarchical RL [Bayes] applied the previously highest-valued [most reliable] task-set:  $Q(TS|c_{new}) = \max_c Q(TS|c)$  [ $p(TS|c_{new}) = \max_c p(TS|c)$ ].

**Model Comparison.** The Bayes Factor  $BF$  quantifies the support for one model  $M_1$  over another model  $M_2$  by assessing the ratio between their marginal likelihoods,  $BF = \frac{p(\text{data}|M_1)}{p(\text{data}|M_2)}$ .  $BF > 1$  provides evidence for  $M_1$ . Marginal model likelihoods represent the probability of the data under the model, marginalizing over model parameters  $\theta$ :  $p(\text{data}|M) = \int p(\text{data}|M, \theta) p(\theta) d\theta$ .

For each model, we simulated datasets by drawing model parameters  $\theta$  uniformly at random. Due to uniform sampling,  $p(\theta)$  is equal for all  $\theta$ , such that the empirical distribution over simulations approximates the marginal likelihood. To obtain Bayes Factors, we computed the same summary statistics  $s_m$  as for humans for each individual simulation (e.g., performance slope in hidden-context phase). We estimated model densities  $\hat{s}_m$  based on a large number of simulations. We obtained marginal model likelihoods as the probability of the human summary statistic  $s_h$  under the model,  $p(s_h|\hat{s}_m)$ . Bayes Factors are given by  $BF = \frac{p(s_h|\hat{s}_{m1})}{p(s_h|\hat{s}_{m2})}$ .

We drew parameters uniformly at random in a range allowing as broad coverage of possible behavior as possible:  $0 < \alpha_a, \alpha_{TS}, f_a, f_{TS} < 1$  and  $1 < \beta_a, \beta_{TS} < 20$ . Each synthetic dataset consisted of 26 agents simulated on the exact same inputs received by the 26 participants, such that the noise in the synthetic statistics was identical to the one in the human dataset. We simulated 50,000 datasets for each model to ensure convergence of the density estimates.

We presented one example datasets for each model in the bar graphs of Figs. 3 A and C and 4A. These datasets were obtained by first selecting all of the 50,000 model simulations that fell within a certain range of human behavior for all summary statistics (50 to 150% for flat and hierarchical RL; 10 to 190% for hierarchical Bayes). We then simulated one new dataset per model based on the median parameter values of the selected models. *SI Appendix* provide a detailed discussion of our model comparison method and selection of the example datasets.

**Data Availability.** All data for this study have been made available for only researchers through the National Institute of Mental Health Data Archive (60). Analysis and modeling code is available on GitHub (<https://github.com/MariaEckstein/TaskSets>).

**ACKNOWLEDGMENTS.** We thank Lucy Whitmore and Sarah Master for their contributions. This project was supported by NIH Grant 1R01MH119383-01 (to A.G.E.C.).

1. T. S. Lee, D. Mumford, Hierarchical Bayesian inference in the visual cortex. *J. Opt. Soc. Am. A Opt. Image Sci. Vis.* **20**, 1434–1448 (2003).
2. M. Botvinick, A. Weinstein, A. Solway, A. Barto, Reinforcement learning, efficient coding, and the statistics of natural tasks. *Curr. Opin. Behav. Sci.* **5**, 71–77 (2015).
3. M. Botvinick, Y. Niv, A. C. Barto, Hierarchically organized behavior and its neural foundations: A reinforcement learning perspective. *Cognition* **113**, 262–280 (2009).
4. W. G. Chase, H. A. Simon, Perception in chess. *Cognit. Psychol.* **4**, 55–81 (1973).
5. E. K. Miller, J. D. Cohen, An integrative theory of prefrontal cortex function. *Annu. Rev. Neurosci.* **24**, 167–202 (2001).
6. A. G. E. Collins, E. Koehlin, Reasoning, learning, and creativity: Frontal lobe function and human decision-making. *PLoS Biol.* **10**, e1001293 (2012).
7. J. B. Tenenbaum, C. Kemp, T. L. Griffiths, N. D. Goodman, How to grow a mind: Statistics, structure, and abstraction. *Science* **331**, 1279–1285 (2011).
8. T. L. Griffiths et al., Doing more with less: Meta-reasoning and meta-learning in humans and machines. *Curr. Opin. Behav. Sci.* **29**, 24–30 (2019).
9. A. G. E. Collins, The cost of structure learning. *J. Cognit. Neurosci.* **29**, 1646–1655 (2017).
10. M. J. Frank, D. Badre, Mechanisms of hierarchical reinforcement learning in cortico-striatal circuits 1: Computational analysis. *Cerebr. Cortex* **22**, 509–526 (2012).
11. Y. Niv et al., Reinforcement learning in multidimensional environments relies on attention mechanisms. *J. Neurosci.* **35**, 8145–8157 (2015).
12. Y. C. Leong, A. Radulescu, R. Daniel, V. DeWoskin, Y. Niv, Dynamic interaction between reinforcement learning and attention in multidimensional environments. *Neuron* **93**, 451–463 (2017).
13. R. C. Wilson, Y. Niv, Inferring relevance in a changing world. *Front. Hum. Neurosci.* **5**, 189 (2012).
14. M. Donoso, A. G. E. Collins, E. Koehlin, Foundations of human reasoning in the prefrontal cortex. *Science* **344**, 1481–1486 (2014).
15. N. A. Taatgen, The nature and transfer of cognitive skills. *Psychol. Rev.* **120**, 439–471 (2013).
16. T. Flesch, J. Balaguer, R. Dekker, H. Nili, C. Summerfield, Comparing continual task learning in minds and machines. *Proc. Natl. Acad. Sci. U.S.A.* **115**, E10313–E10322 (2018).

17. A. G. E. Collins, M. J. Frank, Cognitive control over learning: Creating, clustering, and generalizing task-set structure. *Psychol. Rev.* **120**, 190–229 (2013).
18. R. S. Sutton, A. G. Barto, *Reinforcement Learning: An Introduction* (MIT Press, Cambridge, MA; London, UK, ed. 2, 2017).
19. W. Schultz, P. Dayan, P. Read Montague, A neural substrate of prediction and reward. *Science* **275**, 1593–1599 (1997).
20. H. M. Bayer, P. W. Glimcher, Midbrain dopamine neurons encode a quantitative reward prediction error signal. *Neuron* **47**, 129–141 (2005).
21. D. Lee, H. Seo, M. W. Jung, Neural basis of reinforcement learning and decision making. *Annu. Rev. Neurosci.* **35**, 287–308 (2012).
22. B. Abler, H. Walter, S. Erk, H. Kammerer, M. Spitzer, Prediction error as a linear function of reward probability is coded in human nucleus accumbens. *Neuroimage* **31**, 790–795 (2006).
23. L.-H. Tai, A. Moses Lee, N. Benavidez, A. Bonci, L. Wilbrecht, Transient stimulation of distinct subpopulations of striatal neurons mimics changes in action value. *Nat. Neurosci.* **15**, 1281–1289 (2012).
24. G. E. Alexander, M. R. DeLong, P. L. Strick, Parallel organization of functionally segregated circuits linking basal ganglia and cortex. *Annu. Rev. Neurosci.* **9**, 357–381 (1986).
25. W. Schultz, Updating dopamine reward signals. *Curr. Opin. Neurobiol.* **23**, 229–238 (2013).
26. G. Konradis, On the necessity of abstraction. *Curr. Opin. Behav. Sci.* **29**, 1–7 (2019).
27. M. Botvinick, Hierarchical reinforcement learning and decision making. *Curr. Opin. Neurobiol.* **22**, 956–962 (2012).
28. I. Momennejad *et al.*, The successor representation in human reinforcement learning. *Nat. Hum. Behav.* **1**, 680–692 (2017).
29. J. Ribas Fernandes *et al.*, A neural signature of hierarchical reinforcement learning. *Neuron* **71**, 370–379 (2011).
30. S. Farashahi, K. Rowe, Z. Aslami, D. Lee, A. Soltani, Feature-based learning improves adaptability without compromising precision. *Nat. Commun.* **8**, 1768 (2017).
31. D. Badre, M. J. Frank, Mechanisms of hierarchical reinforcement learning in cortico-striatal circuits 2: Evidence from fMRI. *Cerebr. Cortex* **22**, 527–536 (2012).
32. W. H. Alexander, J. W. Brown, Hierarchical error representation: A computational model of anterior cingulate and dorsolateral prefrontal cortex. *Neural Comput.* **27**, 2354–2410 (2015).
33. M. Haruno, M. Kawato, Heterarchical reinforcement-learning model for integration of multiple cortico-striatal loops: fMRI examination in stimulus-action-reward association learning. *Neural Network.* **19**, 1242–1254 (2006).
34. E. Koehlin, Prefrontal executive function and adaptive behavior in complex environments. *Curr. Opin. Neurobiol.* **37**, 1–6 (2016).
35. D. Badre, Cognitive control, hierarchy, and the rostro-caudal organization of the frontal lobes. *Trends Cognit. Sci.* **12**, 193–200 (2008).
36. D. Badre, M. D'Esposito, Is the rostro-caudal axis of the frontal lobe hierarchical?. *Nat. Rev. Neurosci.* **10**, 659–669 (2009).
37. B. W. Balleine, A. Dezfouli, M. Ito, K. Doya, Hierarchical control of goal-directed action in the cortical-basal ganglia network. *Curr. Opin. Behav. Sci.* **5**, 1–7 (2015).
38. C. Diuk, K. Tsai, J. Wallis, M. Botvinick, Y. Niv, Hierarchical learning induces two simultaneous, but separable, prediction errors in human basal ganglia. *J. Neurosci.* **33**, 5797–5805 (2013).
39. S. Palminteri, V. Wyart, E. Koehlin, The importance of falsification in computational cognitive modeling. *Trends Cognit. Sci.* **21**, 425–433 (2017).
40. S. Monsell, Task switching. *Trends Cognit. Sci.* **7**, 134–140 (2003).
41. A. G. E. Collins, J. F. Cavanagh, M. J. Frank, Human EEG uncovers latent generalizable rule structure during learning. *J. Neurosci.* **34**, 4677–4685 (2014).
42. M. Sunnaker *et al.*, Approximate bayesian computation. *PLoS Comput. Biol.* **9**, e1002803 (2013).
43. G. Jocham, T. A. Klein, M. Ullsperger, Dopamine-mediated reinforcement learning signals in the striatum and ventromedial prefrontal cortex underlie value-based choices. *J. Neurosci.* **31**, 1606–1613 (2011).
44. C. Diuk *et al.*, “Divide and conquer: Hierarchical reinforcement learning and task decomposition in humans” in *Computational and Robotic Models of the Hierarchical Organization of Behavior* (Springer, Berlin, Heidelberg, Germany, 2013), pp. 271–291.
45. J. S. Gershman, On the blessing of abstraction. *Q. J. Exp. Psychol.* **70**, 361–365 (2017).
46. C. Kemp, A. Perfors, J. B. Tenenbaum, Learning overhypotheses with hierarchical Bayesian models. *Dev. Sci.* **10**, 307–321 (2007).
47. H. Steingrover, R. Wetzels, E.-J. Wagenmakers, Bayes factors for reinforcement-learning models of the Iowa gambling task. *Decision* **3**, 115–131 (2016).
48. M. D. Lee, How cognitive modeling can benefit from hierarchical Bayesian models. *J. Math. Psychol.* **55**, 1–7 (2011).
49. D. J. C. MacKay, Bayesian interpolation. *Neural Comput.* **4**, 415–447 (1992).
50. M. S. Tomov, S. Yagati, A. Kumar, W. Yang, S. J. Gershman, Discovery of hierarchical representations for efficient planning. *PLoS Comput. Biol.* **14**, e1007594 (2019).
51. A. Solway *et al.*, Optimal behavioral hierarchy. *PLoS Comput. Biol.* **10**, e1003779 (2014).
52. D. Marr, *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information* (Henry Holt and Co., Inc., New York, NY, 1982).
53. F. Lieder, T. L. Griffiths, Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behav. Brain Sci.* **43**, e1 (2019).
54. A. S. Vezhnevets *et al.*, “FeUdal networks for hierarchical reinforcement learning” in *Proceedings of the 34th International Conference on Machine Learning*, D. Precup, Y. W. Teh, Eds. (PMLR, 2017), vol. 70, pp. 3540–3549.
55. R. S. Sutton, D. Precup, S. Singh, Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artif. Intell.* **112**, 181–211 (1999).
56. J. X. Wang *et al.*, Prefrontal cortex as a meta-reinforcement learning system. *Nat. Neurosci.* **21**, 860–868 (2018).
57. B. M. Lake, T. D. Ullman, J. B. Tenenbaum, S. J. Gershman, Building machines that learn and think like people. *Behav. Brain Sci.* **40**, e253 (2017).
58. A. G. E. Collins, Reinforcement learning: Bringing together computation and cognition. *Curr. Opin. Behav. Sci.* **29**, 63–68 (2019).
59. A. S. Vezhnevets, Y. Wu, R. Leblond, J. Z. Leibo, Options as responses: Grounding behavioural hierarchies in multi-agent RL. arXiv:1906.01470 (6 June 2019).
60. M. K. Eckstein, A. G. E. Collins, Data for computational evidence for hierarchical reinforcement learning in humans. NIMH NDA. <https://dx.doi.org/10.15154/1518660>. Deposited 15 January 2020.