

Neural and behavioral signatures of policy compression in cognitive control

Shuze Liu^{1,*}, Atsushi Kikumoto^{2,3}, David Badre^{2,4}, and Samuel J. Gershman^{5,6}

¹Program in Neuroscience, Harvard University, 52 Oxford Street, 02138, Cambridge, MA, United States

²Department of Cognitive and Psychological Sciences, Brown University, 190 Thayer Street, 02912, Providence, RI, United States

³RIKEN Center for Brain Science, 2-1 Hirosawa, Saitama 351-0106, Wako, Japan

⁴Carney Institute for Brain Science, Brown University, 164 Angell Street, 02906, Providence, RI, United States

⁵Department of Psychology, Harvard University, 52 Oxford Street, 02138, Cambridge, MA, United States

⁶Center for Brain Science, Harvard University, 52 Oxford Street, 02138, Cambridge, MA, United States

*Corresponding author: Email: shuzeliu@fas.harvard.edu

Making context-dependent decisions incurs cognitive costs. Cognitive control studies have investigated the nature of such costs from both computational and neural perspectives. In this paper, we offer an information-theoretic account of the costs associated with context-dependent decisions. According to this account, the brain's limited capacity to store context-dependent policies necessitates “compression” of policies into internal representations with an upper bound on codelength, quantified by an information-theoretic measure (policy complexity). These representations are decoded into actions by sequentially inspecting each bit, such that longer codes take more time to decode. When a response deadline is imposed, the account predicts that policy complexity should increase with the deadline. Higher policy complexity is associated with several behavioral signatures: (i) higher accuracy; (ii) lower variability; and (iii) lower perseveration. Analyzing electroencephalography data from a rule-based action selection task, we found evidence supporting all of these predictions. We further hypothesized that complex policies require higher neural dimensionality (which constrains the code space). Consistent with this hypothesis, we found that policy complexity correlates with a measure of neural dimensionality in a rule-based decision task. This finding brings us a step closer to understanding the neural implementation of policy compression and its implications for cognitive control.

Keywords: action selection; cognitive control; policy compression; rate-distortion theory; representational dimensionality.

Introduction

A prominent hallmark of human cognition is our ability to select actions appropriate for the current context. The study of cognitive control seeks to understand this core capability (Fan 2014; Shenhav 2017; Badre 2024). A growing body of research has highlighted the cognitive costs incurred by such control processes. These costs are evident in human tendencies to avoid tasks that demand greater cognitive control (Kool et al. 2010; Shenhav 2017; Sayali et al. 2023), and increase control under greater reward incentives (Umemoto and Holroyd 2015; Krebs and Woldorff 2017; Shenhav 2017).

The ubiquity of cognitive costs has prompted researchers to probe their underlying mechanisms. Computational models have proposed a range of explanations, from metaphorical force fields to theories of limited mental resources and reward-based accounts of effort allocation (Kruglanski 2012; Botvinick and Braver 2015; Shenhav 2017). In parallel to the computational perspectives above, neuroimaging studies have localized control-related signals in the brain, including but not limited to the anterior cingulate cortex and lateral prefrontal cortex (Koechlin and Summerfield 2007; Badre 2008; McGuire and Botvinick 2010). While these perspectives are in principle reconcilable (Botvinick and Cohen 2014), their integration has been limited by the

absence of a unifying computational resource formulation—one that can both reconcile these diverse theories and be empirically tested.

In this paper, we take a step toward such a unifying perspective by refining an information-theoretic framework of cognitive control. Building on previous theoretical and empirical work, we explore how this perspective can bridge behavioral and neural accounts of cognitive effort. The idea of connecting information theory to cognitive control is not new: early cognitive control studies on the PFC proposed a hierarchical architecture along a rostro-caudal axis, with anterior regions supporting higher-order, more abstract-level control. These ideas were formalized using information-theoretic measures including mutual information and conditional entropy, interpreted as measures of remaining uncertainty within the hierarchy (Koechlin and Summerfield 2007; Badre 2008). While these information-theoretic measures have the benefit of being domain-general and thus generalizable across tasks, the above line of work did not specify how such information-theoretic constructs may reconcile the various cognitive cost formulations at the computational level, or be tested using behavioral data.

More recently, the importation of rate-distortion theory (Cover 1999) into cognitive science has filled this gap. By leveraging

Received: May 7, 2025. Revised: July 9, 2025. Accepted: July 19, 2025

© The Author(s) 2025. Published by Oxford University Press. All rights reserved. For commercial re-use, please contact reprints@oup.com for reprints and translation rights for reprints. All other permissions can be obtained through our RightsLink service via the Permissions link on the article page on our site—for further information please contact journals.permissions@oup.com.

its formalism of constrained optimization—balancing reward against channel information rate (Tishby and Polani 2010; Sims 2016; Zenon et al. 2019; Lai and Gershman 2021)—this line of work has offered normative models of human decision-making that naturally connect the resource and reward-based views of cognitive control (Lai and Gershman 2021). Additionally, the identified solution bears similarities with Kullback–Leibler regularized control and the way it penalizes deviations from a default policy (Todorov 2009), enabling closer connections to the force field perspective. The resulting *policy compression* framework has been empirically supported in contextual bandit tasks (Gershman 2020; Lai and Gershman 2024; Liu et al. 2024; Liu and Gershman 2025), and has begun to find traction in neuroscience (Gershman and Lak 2025). While some connections to cognitive control—such as task switching and the overriding of habits—have been established (Zenon et al. 2019), few empirical studies have quantitatively assessed the framework’s predictions for cognitive control.

In this paper, we introduce the policy compression framework as a normative theory of context-dependent decision-making, articulating its implications for three major issues in cognitive control: the formulation of cognitive cost as the mutual information between context and action (policy complexity), the emergence of perseverative behavior toward default actions, and the link between control and response times. We also propose an algorithmic-level implementation of the framework based on entropy coding, connecting the framework to recent findings on neural representational dimensionality and their role in cognitive control (Rigotti 2013; Bernardi 2020; Kikumoto et al. 2024a). In particular, the framework predicts that representational dimensionality should scale with policy complexity, under the assumption that more complex policies consume more representational resources. To evaluate these predictions, we analyze the behavior and neural activity of human participants performing a rule-based action selection task under varying response deadlines. Similar to prior work in contextual multi-armed bandit tasks (Gershman 2020; Lai and Gershman 2024; Liu and Gershman 2025), the current dataset features within-participant manipulations of response deadlines—specifically, variations in stimulus-onset asynchrony (SOA). This design allows us to examine spontaneous, within-subject fluctuations in policy complexity as well as their behavioral and neural consequences.

The policy compression framework

Here we describe the proposed policy compression framework in detail. We first outline its foundation in rate-distortion theory, which formalizes context-dependent action selection as a reward optimization problem under information-theoretic constraints. We then elaborate on the framework’s implications for cognitive control, highlighting how it accounts for perseverative behavior during action selection and elucidates the links between state-dependent policies, response times, and representational dimensionality.

The nervous system operates under numerous constraints (Shenhav 2017). These constraints have been formalized at multiple levels of analysis, ranging from computational-level accounts of sampling costs and computational complexity (Ma et al. 2014; Vul et al. 2014; Bossaerts et al. 2019; Zhou et al. 2024), to physiological-level costs of interference and neural metabolism (Gailliot and Baumeister 2007; Musslick 2016). Here, we specifically focus on the influence of channel capacity, the maximum information that can be transmitted across a noisy channel (Shannon 1948), on decision-making processes.

The framework we propose is an application of rate-distortion theory to action selection. Rate-distortion theory prescribes how to construct an optimal channel that minimizes some notion of error (the distortion), or maximizes reward, subject to a constraint on the information transmission rate (ie an information bottleneck) (Cover 1999). The utility of rate-distortion theory lies in its task-general formulation of cognitive constraints. Beyond action selection (Lai and Gershman 2021, 2024; Liu et al. 2024; Liu and Gershman 2025), it has been successfully applied to various cognitive processes including visual working memory (Sims et al. 2012; Sims 2015; Jakob and Gershman 2023), perception (Gershman and Burke 2023), intertemporal decision-making (Gershman and Bhui 2020), and cognitive abstraction formation (Genewein et al. 2015). Previous works have linked information theory to various facets of cognitive control, from task-switching costs (Zenon et al. 2019) to the hierarchical organization of executive function (Koechlin and Summerfield 2007; Badre 2008). These studies support a *modulatory* view of control, in which top-down processing of higher-level contextual cues enhances or inhibits lower-level stimulus-response associations to guide behavior (Aron 2007; Goghari and MacDonald 2009; Badre et al. 2021). Such connections enable us to apply rate-distortion theory to cognitive control through the lens of policy compression, which we will elaborate on below.

We assume a contextual decision-making setup, in which the agent selects actions based on some environmental state to maximize reward (Koechlin and Summerfield 2007; Badre 2008; Zenon et al. 2019; Kikumoto et al. 2022, 2024a). In this paper we will use “context” and “state” interchangeably. The agent’s policy $\pi(a|s)$ is a probabilistic mapping from states $s \in \mathcal{S}$ to actions $a \in \mathcal{A}$. Here we make the simplifying assumption that all contextual information is encapsulated into a (potentially high-dimensional) state s that informs action selection (Fig. 1A).

It is well known that context-dependent action selection incurs cognitive costs that affect human behavior (Shenhav 2017; Zenon et al. 2019; Sayali et al. 2023). For a resource-rational agent, we formalize the cognitive cost as the mutual information between states and actions, which we call the *policy complexity*:

$$I^{\pi}(\mathcal{S}; \mathcal{A}) = \sum_s P(s) \sum_a \pi(a|s) \log \frac{\pi(a|s)}{P(a)} \quad (1)$$

where $P(s)$ is the state probability distribution and $P(a) = \sum_s P(s) \pi(a|s)$ is the marginal probability of choosing action a under the agent’s policy $\pi(a|s)$. Intuitively, high-complexity policies preserve state information (eg deterministic mappings from states to actions) whereas low-complexity policies discard state information (eg random actions). One can additionally decompose policy complexity into individual policy cost terms for each state-action pair, $\log \frac{\pi(a|s)}{P(a)}$, whose trial-wise average defines the policy complexity.

If the agent has infinite cognitive resources, it would be optimal to map each state to the most rewarding action. However, we assume that policies are subject to a capacity constraint C (a concept frequently invoked in the cognitive control literature; see (Botvinick and Cohen 2014)), which, in our formulation, acts as an upper bound on policy complexity. The value of C can be voluntarily set by the agent (Liu et al. 2024), or, under cognitive load or response deadlines, limited by the cognitive resources available (Lai and Gershman 2024; Liu and Gershman 2025). In this formulation, for a given task with state probability distribution $P(s)$ and state-specific action rewards $Q(s, a)$, agents would maximize trial-averaged reward $V^{\pi} = \sum_s P(s) \sum_a \pi(a|s) Q(s, a)$ subject to their constraint $I^{\pi}(\mathcal{S}; \mathcal{A}) \leq C$. We can express this constrained

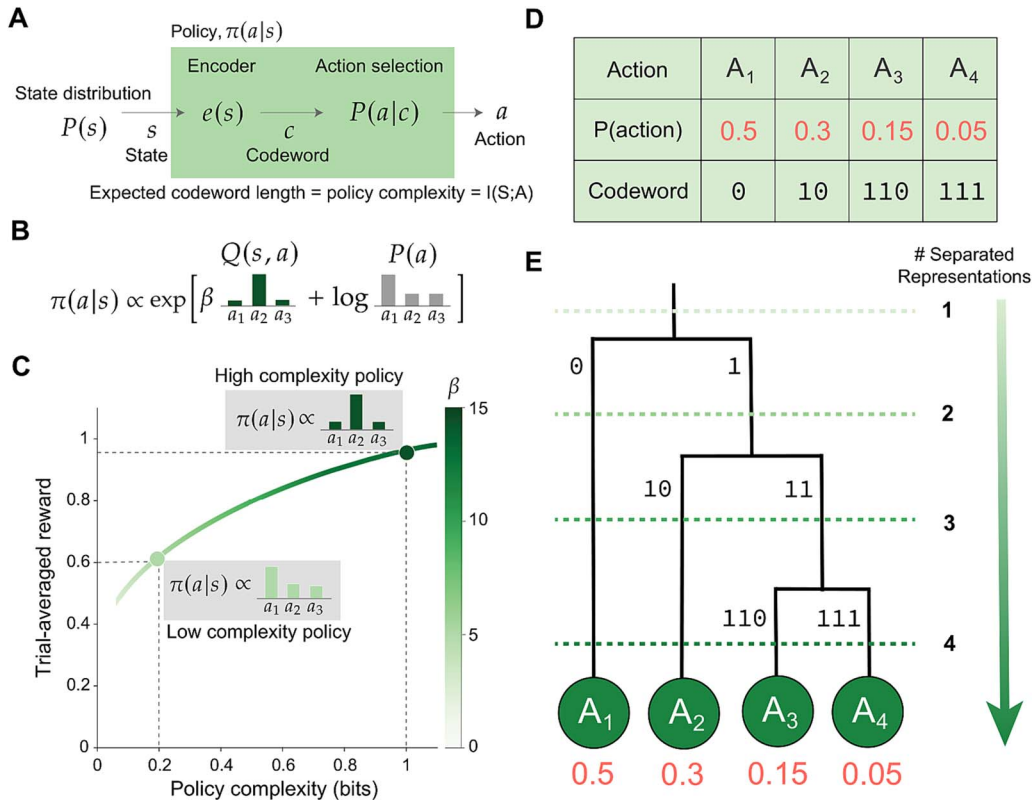


Fig. 1. The policy compression framework and optimal entropy codes. A) Context-dependent action selection as an information channel. The context/state is encoded as a codeword and subsequently decoded to reveal its policy-assigned action. The information channel is capacity-limited, which upper bounds the policy complexity, defined as the mutual information between states and actions. This is equivalent to a bound on the average codeword length. B) The optimal solution prescribed by rate-distortion theory, featuring state-dependent $Q(s, a)$ and state-independent $P(a)$ components. The parameter β adjusts the relative contributions of the two components. C) The reward-complexity Frontier for an example task, derived by varying the β parameter in B) and finding the corresponding optimal policies. The optimal policy at high complexity is strongly influenced by state-dependent rewards $Q(s, a)$ (dark green distribution shared with B)). In contrast, the optimal policy at low complexity is further compressed (less state-dependent) and more strongly influenced by its own marginal action distribution $P(a)$. A–C) are adapted from (Lai and Gershman 2024). D) An example Huffman code with optimal codeword assignments that minimizes average code length. The code lengths are assigned based on the probability of decoding different symbols (leaf nodes). In policy compression, the codetree would be tailored to decode actions. E) The same Huffman code visualized as a codetree. On the process level, we hypothesize that each bifurcation in the codetree manifests as an increase in neural representation dimensionality, which is required to separate the two possible branches downstream. Under trial-by-trial response deadlines, codetree traversal would be cut off at earlier bifurcations (colored dotted lines), thus capping the number of representation separations and downstream policy complexity.

optimization problem in Lagrange form:

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \beta V^\pi - I^\pi(S; A) + \sum_s \lambda(s) \left(\sum_a \pi(a|s) - 1 \right) \quad (2)$$

where $\beta \geq 0, \lambda(s) \geq 0 \forall s \in S$ are Lagrange multipliers to ensure $I^\pi(S; A) \leq C$ and proper policy normalization: $\sum_a \pi(a|s) = 1$. Solving the Lagrange form yields the following optimal policy:

$$\pi^*(a|s) \propto \exp(\beta Q(s, a) + \log P^*(a)) \quad (3)$$

where $P^*(a) = \sum_s P(s) \pi^*(a|s)$ is the optimal marginal action distribution, and β is a Lagrange multiplier whose value depends on C in a task-specific manner. Despite the recursive nature of Equation (3), one can numerically find the optimal policy $\pi^*(a|s)$ by iterating updates of Equation (3) and the defining equation of $P^*(a)$ for different β values. This numerical process is known as the Blahut-Arimoto algorithm, which is guaranteed to converge (Arimoto 1972; Blahut 1972; Tishby and Polani 2010).

Intuitively, at high policy complexity (corresponding to large C), the value of β is large and the optimal policy is dominated by Q -values, which renders it state-dependent. At low policy

complexity (small C), the value of β is close to 0 and Q -values have minimal impact on the optimal policy. Moreover, low-complexity policies are dominated by the $\log P^*(a)$ term, which manifests as perseveration toward more frequently chosen actions in a state-independent manner (Fig. 1B). In general, high-complexity policies enable more trial-averaged reward than low-complexity policies due to their state-dependence. By varying β and calculating the optimal policy, we can trace out the reward-complexity Frontier, which delimits the maximal trial-averaged reward obtainable for a given policy complexity (Fig. 1C).

Connections to habits and perseveration

The formation of default actions or habits—and the cognitive cost required to override them—has been a central topic in cognitive control (Miller and Cohen 2001; Koehlin and Summerfield 2007; Shenhav 2017; Zenon et al. 2019). The policy compression framework offers a resource-rational account of such habitual behavior: when cognitive resources are limited, biasing action selection toward frequently chosen actions is part of the optimal strategy. This arises because the optimal policy depends on the marginal action distribution $P^*(a)$. When $P^*(a)$ is uniform, the optimal policy simplifies to the standard softmax choice rule,

$\pi(a|s) \propto \exp[\beta Q(s, a)]$ (Sutton and Barto 2018). However, when $P^*(a)$ is biased—reflecting that some actions have been chosen more frequently than others—the optimal policy will favor these actions, making them appear as default choices controlling for differences in $Q(s, a)$. As cognitive resources become more limited (low policy complexity and small β), this bias is amplified: $P^*(a)$ increasingly dominates Equation 2, leading to stronger perseveration toward default actions. Furthermore, assuming that participants build and update their $P(a)$ on a trial-by-trial basis (eg via an iterative update process, as modeled in Lai and Gershman (2021); Bari et al. (2024); Lai and Gershman (2024); Liu et al. (2024)), they should demonstrate a greater tendency to repeat actions at low complexity. Such perseverative tendencies—either toward more frequently-taken actions or manifested in the process-level implementation of action repetition—are not predicted by traditional softmax models with variable noise levels, where the policy converges to randomness as noise increases. In recent work, we have identified perseverative signatures in human decision-making that are consistent with policy compression but not predicted by the softmax model (Lai and Gershman 2024; Liu et al. 2024; Liu and Gershman 2025), supporting the framework's relevance. Unlike previous accounts that attribute habits to model-free learning (Krueger and Griffiths 2018; Pauli et al. 2018), policy compression explains habits as value-independent tendencies shaped purely by past action frequencies. This perspective aligns with and provides a normative justification for recent findings suggesting that habits can form and exert influence independently of value representations (Miller et al. 2019; Nebe et al. 2024; Zhang et al. 2024).

Connections to representational dimensionality

Recent studies of cognitive control have proposed a novel account of action selection, informed by the representational dimensionality of neural population codes. Representational dimensionality is defined as the minimum number of dimensions required to capture the variability of neural activity across task states (Fusi et al. 2016; Ahlheim and Love 2018; Badre et al. 2021; Jazayeri and Ostojic 2021). In this view, high dimensionality allows separating representations of different states, which enables state-dependent action selection. Conversely, low-dimensional representations reduce such separability, limiting the ability to tailor action selection to the current state (Badre et al. 2021; Badre 2024). Specific to rule-based decision-making, the benefit of high dimensionality manifests as the formation of conjunctive state representations integrating stimuli and rules. Consistent with this view, empirical work in cognitive control has observed transient increases in neural representational dimensionality during action selection, coinciding with higher-quality rule-based decisions (Kikumoto et al. 2022, 2024a).

The above representational account has rarely been linked to the information-theoretic view of cognitive control, which typically emphasizes top-down modulatory architectures instead (Koechlin and Summerfield 2007; Fan 2014; Badre et al. 2021). However, a theoretical bridge emerges when we consider algorithmic-level entropy codes inspired by information theory. To illustrate the connection, let us consider representing states as binary codewords generated through entropy coding, with the Huffman code (Huffman 1952) as a canonical example (Fig. 1D). In this scheme, each state codeword maps to a specific action by traversing a binary tree structure. State-dependent action selection requires that codewords corresponding to different states be sufficiently distinct to reach different leaf nodes. This necessitates longer codewords for readout, which in turn requires

traversing more bifurcations in the codetree (Fig. 1E). In the regime of errorless transmission, the policy's complexity (in bits) corresponds to the average codelength, under an optimal entropy coding (eg Huffman coding) that minimizes this quantity (Cover 1999). By adopting this entropy coding view of policy compression, we hypothesize that the codetree bifurcations during readout map onto transient increases in neural representational dimensionality. Consequently, behavioral policy complexity should reflect the representational dimensionality required to support the policy.

Optimal entropy codes, such as Huffman codetrees, apply primarily to errorless transmission. However, when placed under time or memory constraints (as in the cognitive control dataset analyzed), agents may be unable to fully traverse the (optimal) codetree and thus operate in a lossy compression regime (Lai and Gershman 2021). In our dataset, where response deadlines vary across trials, participants would be forcefully cut off at earlier bifurcations. Consequently, action selection would become less sensitive to the specific state, leading to reduced policy complexity and possible suboptimality (compared to a policy optimized for such resource constraints along the reward-complexity Frontier).

Connections to response times

It is well established that exerting cognitive control leads to longer RTs (Matsumoto and Tanaka 2004; Kool et al. 2017). Through the lens of optimal entropy codes, policy compression rationalizes longer RTs under high-complexity policies. Specifically, executing a high-complexity policy entails reading longer state codewords on average, which requires additional time to traverse the codetree (Lai and Gershman 2021). This viewpoint closely mirrors neuroscience studies showing that transient increases in representational dimensionality are temporally extended during decision tasks, and that forceful cutoffs at an earlier timepoint diminishes rule-based performance (Kikumoto et al. 2024a). Similar information-theoretic explanations for RTs have long been applied to decision tasks that vary the number of available actions, encapsulated in the Hick-Hyman Law. This law formalizes the empirical observation that RT increases logarithmically with the number of possible actions, or equivalently, linearly with the amount of information transmitted (Hick 1952; Hyman 1953). Whether policy complexity similarly predicts RTs in cognitive control tasks remains an open empirical question.

Materials and methods

To directly test the predictions of policy compression for cognitive control tasks, we reanalyzed data from a previous electroencephalography (EEG) study (Kikumoto et al. 2024a). We briefly summarize the study methods here, and refer readers to the original paper for more details.

Participants

Forty-two participants (27 female, mean age 22 yr) were recruited. The recruitment followed procedures approved by the Human Subjects Committee at the RIKEN (approval number: RIKEN-W1-2022-030), and all participants gave informed consent. The sex and gender of participants were determined based on self-report. They all had normal or corrected-to-normal vision and had no history of neurological or psychiatric disorders. No statistical method was used to predetermine the sample size. After preprocessing the EEG data, one participant was removed due to excessive amounts of artifacts (ie more than 25% of trials; see EEG recordings and preprocessing for details).

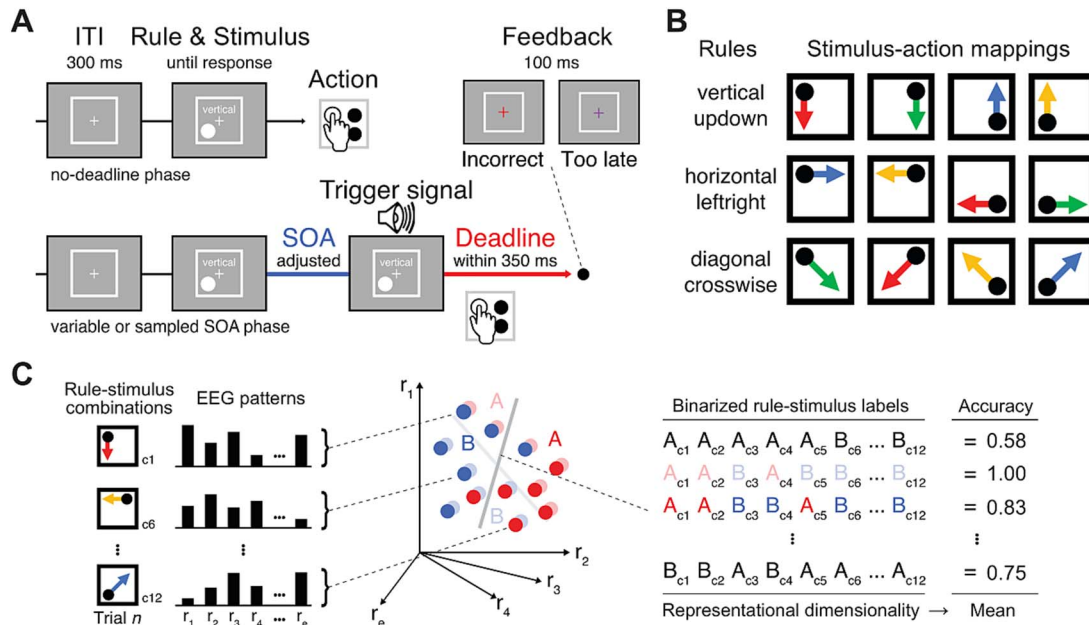


Fig. 2. Behavioral task and EEG shattering dimensionality estimation. A) Task pipeline and SOA conditions. On every trial, participants observe one randomly sampled rule and stimulus, and respond by choosing an action on the number pad. In the variable and sampled SOA phase, participants receive an auditory SOA signal some time interval after trial onset, and they must respond within 0 to 350 ms after the signal's presentation. In the variable SOA phase, the SOA intervals are determined for each subject separately to achieve three predefined accuracy levels. The SOA condition for each trial is sampled randomly, without block-level regularity that allows participants to adjust their strategy beforehand. B) The stimulus-action mappings are determined by high-level rules. There are 3 possible rules, 4 possible stimuli (hence $3 \times 4 = 12$ rule-stimulus combinations, randomly sampled for each trial), and 4 possible actions. C) Shattering dimensionality estimation. EEG signals (31 electrodes \times 5 frequency bands = 155 dimensions) are fed into a linear decoder, to decode each of the $(2^{12} - 2) = 4,096$ binary partitions of true rule-stimulus identity for each trial. The resulting trained decoder accuracy rates are used to estimate dimensionality of the underlying EEG representation. Adapted from Kikumoto et al. (2024a).

Behavioral task

Participants completed a rule-based action selection task. On each trial, the participant simultaneously sees one of three rules ("horizontal," "vertical," and "diagonal") and one of four stimuli (situated in a 2-by-2 matrix), randomly sampled with equal probability for each trial, yielding a total of 12 rule-stimulus combinations (ie "states"). Participants must choose one of four actions, arranged in a 2-by-2 matrix on a number pad (Fig. 2A). The rule determines the correct stimulus-action mappings. For example, if the current rule is "horizontal" and the stimulus is on the top-left, the correct action is to press the top-right action key on the number pad, as this position is horizontally adjacent to the stimulus (Fig. 2B).

Trials were organized into experimental blocks, each lasting 18 s. For each block, participants were instructed to complete as many trials as possible. Trials that were initiated within the 18 s block duration but extended beyond it were allowed to finish.

Blocks were further organized into task phases. Participants completed three task phases in order: the no-deadline phase, the variable SOA phase, and the sampled SOA phase. The no-deadline phase contained 25 blocks, where participant RTs are used to adjust the SOA interval for later phases.

In the variable SOA phase containing 35 blocks, for every trial, participants received an auditory signal some time after the trial onset. Upon this SOA onset, they were required to select an action within 350 ms after hearing the auditory signal. If participants made a response before SOA onset or after 350 ms after SOA onset, the trial response is considered invalid, and participants receive no reward. The time from trial onset to auditory signal presentation is called the SOA interval, which could take on 12 values determined for each participant. In each trial, the SOA interval was randomly sampled.

In the sampled SOA phase containing 185 blocks, SOAs were still present after trial onset, but each participant only received three unique SOA intervals, determined in a participant-specific fashion to induce accuracy rates of 50%, 70%, and 90% across trials (both invalid and incorrect responses are counted as inaccurate). The SOA interval was randomly sampled for every trial and thus interleaved. These three SOA intervals constitute task conditions and are labeled "short," "medium," and "long SOA" respectively. EEG data were collected in this phase.

Estimation of SOA functions

To separately analyze the three accuracy levels during the sampled SOA phase, we preprocessed each participant's RTs during the no-deadline and variable SOA phases. For the no-deadline phase, we computed each participant's RT distribution, excluding trials with RTs slower than 5 SDs. We model the RT distribution using the ex-Gaussian distribution using the exgauss toolbox in Matlab 2019B.

In the variable SOA phase, each participant receives 12 possible SOA interval values. The values were determined based on deviations from the participant's mean RT in the no-deadline phase (estimated as the mean of the fitted ex-Gaussian distribution). The deviations range from -450 ms to $+200$ ms with 50 ms increments. Negative RTs were dropped during estimation. This allowed us to construct SOA functions that map SOA interval values to accuracy rates.

We modeled SOA functions using an exponential function: $p(\text{correct}) = \lambda(1 - \exp(-\beta(t - \Omega)))$ if $t > \Omega$, and $p(\text{correct}) = 0$ otherwise. The parameters (λ, β, Ω) were estimated by a nonlinear least square solver via the lsqcurvefit function in Matlab. The fitted SOA functions allowed us to identify SOA interval values

that would produce particular accuracy levels (see (Kikumoto et al. 2024a) for visualizations of fitted SOA functions).

EEG recordings and wavelet analysis

EEG recordings were collected during the variable SOA phase. The EEG signals were recorded using a Brain Products actiCHamp recording system (Brain Products GmbH), featuring 31 electrodes from a broad set of scalp sites. The scalp EEG and EOG were amplified with an SA Instrumentation amplifier with a band-pass of 0.01 to 45 Hz, and signals were downsampled at 250 Hz using EEGLab93. We aligned trial-wise EEG recordings, spanning -800 ms to $+200$ ms relative to action onset (0 ms).

After obtaining trial-wise EEG recordings, we decomposed them into 5 frequency bands via complex wavelet analysis (2 to 3 Hz for the delta-band, 4 to 7 Hz for the theta-band, 8 to 12 Hz for the alpha-band, 13 to 30 Hz for the beta-band, and 31 to 35 Hz for the gamma-band), yielding a power measure for each timestep and frequency band. This produced 155 EEG features (31 electrodes \times 5 frequency bands) for every participant, trial, and timestep (each lasting $1\text{ s}/250\text{ Hz} = 4\text{ ms}$). For more details on EEG recordings and wavelet analysis, please see Kikumoto et al. (2024a).

Estimation of EEG representational dimensionality

To characterize neural representational dimensionality, we computed a shattering dimensionality estimate informed by the number of binary partitions of rule-stimulus combinations that are linearly separable based on the underlying neural representation (Rigotti 2013; Bernardi 2020; Courellis 2024; Kikumoto et al. 2024a).

The EEG instantaneous power vectors (31 electrodes \times 5 frequency bands) were further processed before decoder training. First, trials where responses occurred before SOA onset and trials where responses were completely omitted were excluded. Further, all trials in the first block (of the sampled SOA phase) were excluded. The remaining instantaneous power data was further averaged into 20 ms time bins. For each participant and frequency band, the input vector entries were z-transformed across electrodes to remove effects that scaled all electrodes uniformly.

Based on how shattering dimensionality is typically estimated, we trained linear decoders on the EEG instantaneous power vectors to recover information on the trial's presented rule-stimulus combination. Specifically, for each participant and 20 ms time bin, we label each underlying trial's presented rule-stimuli combination. We then bi-partition the 12 possible combinations, leading to $(2^{12} - 2 = 4,096)$ possible binary partitions (each of the 12 combinations may be included in or excluded from Group 0; minus the two binary partitions leading to all 0-labels or 1-labels). Thus, for each binary partition, we obtain a dataset over trials where the input consists of the EEG instantaneous power vector (5 frequency-bands \times 31 electrodes = 155 dimensions), and the output are 0/1 trial-specific binary labels (Fig. 2C left). For each of the 4096 binary partitions, we train a linear decoder to recover the binary labels from the input vector (Fig. 2C center). This leads to a decoding accuracy value for the specific participant, time bin, and binary partition. The above process is repeated under repeated five-fold cross-validation, where the folds themselves were repeatedly partitioned 5 times through a random process. The resulting decoding accuracies are averaged.

After obtaining repeated, cross-validated decoding accuracies for each of the 4,096 binary partitions, we further averaged over the decoding accuracies over binary partitions. This *partition-averaged decoding accuracy* measures the linear separability of

neural representations across all possible binary partitions. Consistent with past work (Bernardi 2020; Courellis 2024), we use these partition-averaged decoding accuracies to construct a proxy of the EEG signals' neural representational dimensionality (Fig. 2C right). Specifically, the above partition-averaging process was done separately for each SOA task condition, leading to condition-specific representational dimensionalities for each participant and timepoint. We then aggregated the partition-averaged decoding accuracies over timepoints to compute the mean decoding accuracy across time bins for a trial, and additionally the trial-averaged mean of the maximum decoding accuracy over each trial's time bins. The latter measure is our proxy for representational dimensionality.

Note that unlike some previous studies (Bernardi 2020; Courellis 2024), the rule-stimulus combinations used in this task did not feature strict dichotomies. Consequently, we decided to use a slightly different procedure in computing decoding accuracies, incorporating all $(2^{12} - 2)$ possible binary partitions. While many of these binary partitions induce classification imbalance (eg all but one rule-stimulus combination being assigned the label 0), they do not pose a significant problem for our subsequent analysis. This is because the imbalance is present for all participants, SOA conditions, trials, and timesteps, and we only focus on comparing relative differences across partition-averaged decoding accuracies.

Estimation of policy complexity

We defined policy complexity as the mutual information between the observed states and chosen actions. Following prior work (Gershman 2020; Lai and Gershman 2024; Liu et al. 2024), we estimated the policy complexity of each participant in each SOA condition using the Hutter estimator (Hutter 2001). Specifically, for each of the 12 states (rule-stimulus combinations), we assume a symmetric Dirichlet prior with $\alpha = 0.01$ for all actions chosen, and use the empirical action counts to reach a posterior Dirichlet distribution over action probabilities. We then estimate policy complexity as the mutual information of the posterior mean policy. The above procedure is informed by previous literature, reporting that the resulting estimates exhibit reasonably good performance when the joint distribution is sparse (Archer et al. 2014). The choice of $\alpha = 0.01$ is informed by rate-distortion theory, stating that empirical trial-averaged reward values cannot be above the reward-complexity Frontier. We have chosen $\alpha = 0.01$ empirically to obey this property.

We also computed policy cost values $\log \frac{\pi(a|s)}{P(a)}$ for individual trials. Given the blocked structure of the experiment, as well as the dynamically evolving nature of $P(a)$, we decided to track $\pi(a|s)$ and $P(a)$ in an online manner for each block separately.

Statistical analysis and modeling

We assessed the behavioral predictions of policy compression using two-sided paired t-tests across the two most extreme SOA conditions (low versus high). Regarding trial-averaged RT, we quantitatively connect it to policy compression through the fitting of linear mixed-effects models (LME), using each participant-condition's policy complexity (both fixed and participant-specific random effects) to predict corresponding trial-averaged RTs. The resulting LME model was compared to null models with only intercepts via the Bayesian Information Criterion (BIC).

For the EEG data, we connected policy compression to the shattering representational dimensionality measures derived according to previous sections. We only analyze responsive trials with available EEG data, focusing on the decoding accuracies of time

bins from -800 ms to 0 ms aligned to action onset (ie discarding post-action decoding accuracies from 0 ms to $+200$ ms). The dimensionality measure—trial-averaged mean of the maximum decoding accuracy over each trial's time bins—was of particular interest, due to its connection with policy complexity—the average readout codelength of the policy's optimal entropy coding. On the single-trial level, we studied the time-averaged decoding accuracy and policy costs. The mean is taken over time bins instead of the maximum, due to the noisiness of trial-wise EEG data. We again fitted LME models that used behavioral policy compression or policy cost measures (fixed and participant-specific random effects) to predict dimensionality or trial-wise mean decoding accuracy. The fitted LMEs were compared against null models with only intercepts or order effects (eg trial number) using BIC.

Perseveration analysis

Given the policy compression framework's rationalization of perseveration behavior, we computed probabilities of repeating the previous trial's action, reapplying the previous trial's rule to the current trial's stimulus, and choosing actions based on the previous trial's stimulus and the current trial's rule. The perseveration analyses is done in a block-wise manner, so as to filter out the influence of resting time between blocks. Before computing the above probabilities, we excluded all omitted trials (in which participants did not respond), all trials preceded by an omitted trial, and the first trials of every block. Given the rule application analysis, we also excluded trials in which the participant had chosen an action incorrect under every possible rule (ie choosing the action that has the same number-key location as the stimulus). The application of these criteria excluded 9,113 (8.33%) responsive trials across participants.

Results

Experimental predictions

The policy compression framework makes the following predictions: longer SOA conditions should be associated with 1) higher policy complexity; 2) higher trial-averaged RT; 3) higher accuracy. Furthermore, 4) the change in RT should be largely explained by changes in policy complexity, such that an LME model using participant-specific policy complexity levels to predict RT should identify positive fixed effects and win model comparison against a null LME model. Predictions 1, 2, and 4 arise due to the RT implications of policy complexity as discussed in the framework's introduction.

Prediction 3 derives from the reward-complexity Frontier associated with the task, coupled with the normative assumption that participants should achieve maximally attainable trial-averaged rewards at their chosen policy complexity levels. Policy compression also predicts that 5) the action entropy $H(A|S)$ —capturing choice variability under the same state—would increase, as informed by the optimal policies at different β levels. Regarding perseveration patterns, the framework predicts that 6) the probability of repeating the previous trial's action, should increase. Given the hierarchical rule-stimulus structure of the task, we further postulate 7) increases in the probability of applying the previous trial's rule on the current trial's stimulus. These predictions arise from the $P^*(a)$ term in Equation (3), and the fact that this marginal probability distribution must be dynamically updated across trials.

Regarding representational dimensionality, the policy compression framework makes the following predictions: on the

single trial level, 8) an LME model using policy costs to predict single-trial mean decoding accuracy should perform better than a null model, which either features only fixed and random intercepts, or additionally non-policy-compression predictors including block-order effects, trial-order effects, and SOA condition. On the trial-aggregate level, we should see 9) higher trial-averaged maximum decoding accuracy over timesteps (proxy for representational dimensionality) for longer SOA conditions, and 10) an LME using participant-condition-specific policy complexity values should predict representational dimensionality better than a null model with only intercept effects. These predictions derive from the connection of optimal entropy codelengths to representational dimensionality necessary for policy execution, as discussed in the earlier policy compression section.

Behavioral results

Participants completed an average of $2,518 \pm 24.28$ trials each (mean \pm SEM), resulting in a total of 103,256 responsive trials for subsequent analysis.

Participant behavior was closely aligned with the task's reward-complexity Frontier (Fig. 3A). This suggests that participants have achieved near maximal trial-averaged reward as allowed by their policy complexity level, supporting the relevance of the policy compression framework for cognitive control tasks. However, we also observed systematic deviations from the Frontier at low levels of policy complexity. These deviations suggest an inefficient use of cognitive resources, caused by suboptimal behavioral patterns that we will elaborate on below.

The behavioral predictions of policy compression were predominately supported. In longer SOA conditions, participants adopted higher policy complexity ($t(40) = -18.7, p < 10^{-20}$; Fig. 3B) while incurring longer RTs ($t(40) = -15.1, p < 10^{-17}$; Fig. 3C), higher accuracy rates ($t(40) = -15.2, p < 10^{-17}$; Fig. 3D), and lower action entropy ($t(40) = 19.1, p < 10^{-20}$; Fig. 3E). To determine whether the increase in RTs was indeed associated with increased policy complexity, we fitted an LME model with fixed and random effects for both the intercept and policy complexity predictors to predict trial-averaged RTs. As expected, the fitted LME identified positive fixed effects for policy complexity ($0.236 \pm 0.0188, t(121) = 12.5, p < 10^{-22}$, random effects $SD = 0.0875$; Fig. 3F), and outperforms a null LME model with only fixed and random effects for the intercept ($\Delta BIC = -122$).

We next examined perseveration patterns as predicted by the framework. At the level of low-level actions, there was no evidence of increased perseveration: the probability of repeating the previous trial's action did not significantly change across SOA conditions ($t(40) = 0.289, P = 0.774$; Fig. 4A). In contrast, we observed prominent perseveration tendencies at the higher level. Under shorter SOAs, participants were more likely to reapply the rule they behaviorally followed on the previous trial to the current trial's stimulus, suggesting rule-based perseveration ($t(40) = 4.77, p < 10^{-4}$; Fig. 4B). Additionally, they were also more likely to apply the current trial's rule to the previous trial's stimulus, which could be summarized as stimulus-based perseveration ($t(40) = 2.75, P = 0.00889$; Fig. 4C).

To better understand the empirical suboptimalities observed in Fig. 3A and the lack of significant differences in Fig. 4A, we conducted follow-up analyses stratified by action identities. The observed suboptimality at low policy complexity likely resulted from a biased marginal action distributions $P(a)$. Specifically, in the short SOA condition, participants disproportionately favored two specific actions— a_1 and a_2 , corresponding to the top two keys

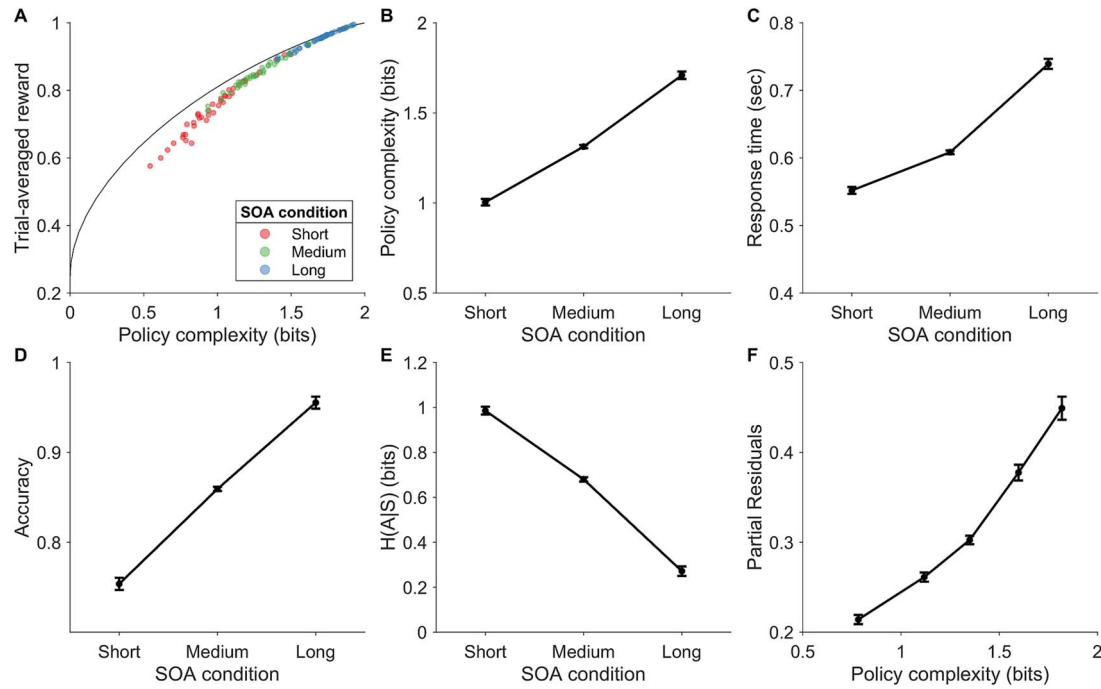


Fig. 3. Behavioral results. Color denotes SOA condition. A) Reward-complexity Frontier and human data, where each of the $3 \times 4 = 12$ rule-stimulus combinations is treated as a separate state. Trial-averaged reward is binarized according to correct/incorrect responses provided, regardless of whether the response is valid (ie provided after the SOA onset but before another 350 ms passed by). B) Policy complexity, C) RT, D) Accuracy, E) Action entropy $H(A|S)$, F) Partial residual plot for the fixed effects of the RT LME model, visualizing the effect of policy complexity on trial-averaged RT. Observations were binned into quantiles for visualization. For A–E), error bars denote mean \pm standard error of the mean (SEM) across participants (Cousineau 2005); for F), error bars denote mean \pm SEM across binned observations.

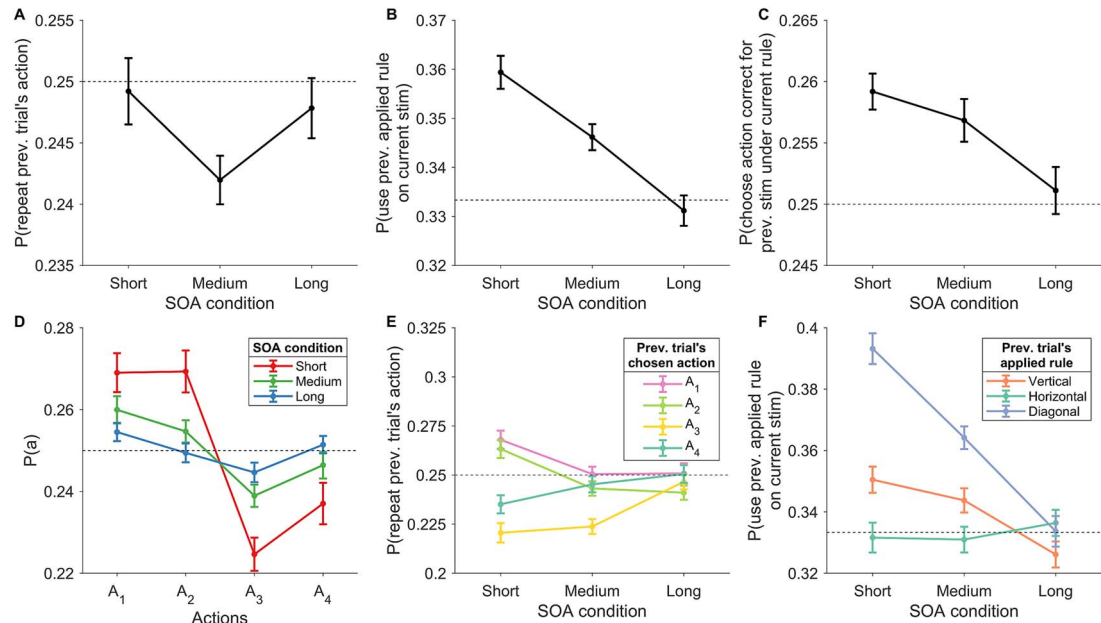


Fig. 4. Perseveration. Color denotes different actions/rules. A) Probability of repeating the previous trial's action; B) Probability of reapplying "the rule used by the agent on the previous trial" to the current trial's stimulus. C) Probability of choosing the previous trial's stimulus's correct action under the current trial's rule. D) $P(a)$ across SOA conditions. E) A) stratified by action identity. B) stratified by rule identity. The dotted lines in each panel denote chance levels of action selection or rule application. Error bars denote mean \pm SEM across participants.

on the number pad (mean $P(a)$ for a_1 and a_2 versus that for a_3 and a_4 : $t(40) = 5.78, p < 10^{-6}$; Fig. 4D). This preference deviates from the framework's normative predictions, which prescribe optimal policies with equiprobable marginal distributions $P(a)$ across β values.

We conducted similar follow-up analyses of perseveration patterns, stratifying by either action or rule identity. We again observed the empirical preference for a_1 and a_2 in the short SOA condition (Fig. 4E). Additionally, participants increasingly favored the diagonal rule under shorter SOAs (Fig. 4F).

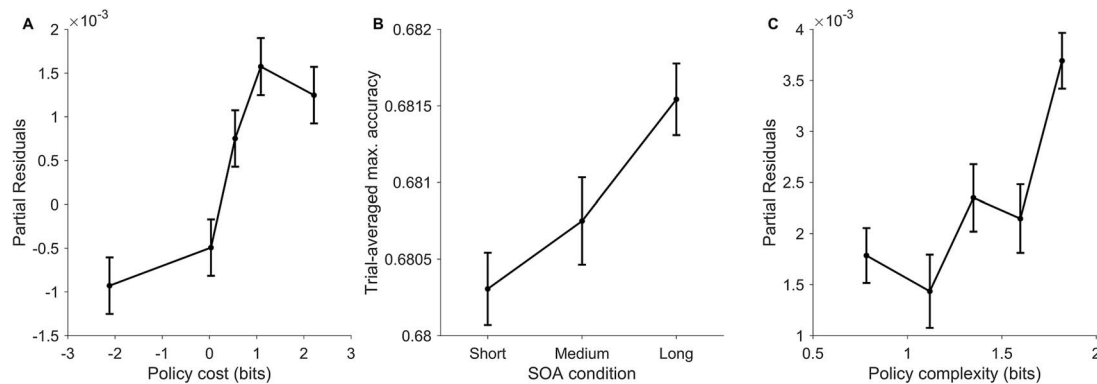


Fig. 5. Neural dimensionality. A) Partial residual plot for the fixed effects of the trial-by-trial LME, visualizing the effect of policy cost on decoding accuracy. Error bars reflect mean \pm SEM across LME observations. B) Mean \pm SEM of participant maximum dimensionality across SOA conditions. For each participant, the maximum decoding accuracy for each trial is computed over time, and then averaged across trials to reach the dimensionality measure. C) Partial residual plot for the fixed effects of the trial-aggregated LME, visualizing the effect of policy complexity on the maximum dimensionality as computed in (B). Error bars denote mean \pm SEM across LME observations. For A) and C), LME observations were binned into quantiles for visualization.

EEG dimensionality results

On the single-trial level, we assessed whether policy cost—whose average across trials defines the policy complexity—predicts trial-wise mean EEG decoding accuracy. The corresponding LME yielded positive fixed effects for policy cost (0.000887 ± 0.000177 , $t(96630) = 5.01$, $p < 10^{-6}$, random effects SD = 0.000500; Fig. 5A), while also outperforming a null model containing only fixed and random effects on the intercept ($\Delta\text{BIC} = -16.5874$). To rule out potential confounds, we constructed alternative null models that included block order, trial order, and SOA condition as predictors, either individually or in combination. Model comparison consistently preferred the policy cost LME ($\Delta\text{BIC} \in [-138.1, -42.03]$ across all pairwise comparisons).

At the trial-aggregated level, we examined whether policy complexity predicts representational dimensionality, measured by the trial-averaged maximum decoding accuracy across timesteps. As predicted, representational dimensionality was higher for longer SOA conditions ($t(40) = -3.31$, $P = 0.00198$; Fig. 5B). Additionally, an LME model predicting representational dimensionality from policy complexity revealed positive fixed effects for the latter (0.00168 ± 0.000567 , $t(121) = 2.97$, $P < 0.00362$, random effects SD = 0.000327; Fig. 5C). However, in contrast to the single-trial results, the above model lost model comparison against a null model with only fixed and random effects for the intercept ($\Delta\text{BIC} = 5.71$), indicating ambiguous evidence for the relationship at the trial-aggregate level.

Discussion

In this paper, we introduced policy compression as a normative framework for understanding cognitive control. Extending its previous applications to multiarmed bandit tasks, we conceptualized the core challenge of cognitive control—dynamically selecting actions based on context-specific goals (Fan 2014; Badre 2024)—as a constrained optimization problem: agents must maximize external rewards under cognitive resource constraints. We formalized these resource constraints using policy complexity, defined as the mutual information between environmental states and the policy-assigned actions. This formulation captures the informational cost of context-sensitive decision-making, as postulated by previous accounts (Fan 2014; Shenhav 2017; Zenon et al. 2019).

We demonstrated the theoretical utility of policy compression by linking it to several phenomena central to cognitive control. This includes the emergence of default actions and habits, the effortful overriding of habitual responses, the neural representational dimensionality required to support control, and the prolonged response times associated with increased context-sensitivity. In doing so, policy compression offers a unifying explanation for these diverse observations. Notably, our approach advances traditional information-theoretic views of cognitive control (Koechlin and Summerfield 2007; Fan 2014) by adopting a resource-rational perspective that explicitly incorporates cognitive costs (Shenhav 2017), and by grounding the theory in recent neuroscientific findings on representational geometry and dimensionality (Rigotti 2013; Bernardi 2020; Kikumoto et al. 2024a).

To evaluate the behavioral and neural predictions of policy compression for cognitive control tasks, we analyzed a previously collected dataset featuring rule-dependent action selection and EEG recordings. Our behavioral results support the framework's core predictions: across SOA conditions, behavior was close to the optimal reward-complexity Frontier, exhibiting varying levels of policy complexity that explain corresponding changes in response times and action stochasticity. In contrast, the framework's predictions on perseveration are partially supported. While action-level perseveration did not significantly change across conditions (which could result from a suboptimal over-reliance on two of the four actions), participants demonstrated a marked tendency to repeat previously applied rules or stimuli, suggesting compression at higher levels of abstraction.

Combined with previous studies on multiarmed bandits, the above behavioral findings contribute to a growing body of work suggesting that policy compression provides a domain-general account of context-sensitive behavior. A limitation of this prior work is the requirement that participants learn their policy from feedback; the policy compression framework is fundamentally about limitations on action selection, not learning, and thus learning dynamics obscure the framework's predictions. Cognitive control tasks—where learning is deliberately minimized by the task design—offer a cleaner test of the theoretical predictions. Our results suggest that human behavior, in the absence of reward learning, aligns even more closely with the framework's optimality predictions. These promising findings encourage future work connecting policy compression to neighboring topics in cognitive control, such as task-switching costs

(Zenon et al. 2019) and aversion to hierarchical task abstraction (Sayali et al. 2023), both of which can be interpreted as consequences of limited information-processing capacity.

More behavioral work is needed to further assess policy compression's predictions for cognitive control tasks. First, the $P^*(a)$ term in Equation (3) could be operationalized as repeating previous actions as we hypothesized, but it is fundamentally a bias toward more frequently taken actions. The current dataset did not feature manipulations over the reward structure $Q(s, a)$, which, when present, could induce changes in the optimal marginal action distribution $P^*(a)$ by biasing it toward certain low-level actions or high-level rules (Lai and Gershman 2024; Liu et al. 2024; Liu and Gershman 2025). Consequently, we were only able to test the framework's predictions of (optimal) flat marginal distributions and heightened repetition in rule/action application under low policy complexity, but not the normative bias toward actions that ought to be more frequently taken. Novel task designs that incorporate such $P^*(a)$ manipulations would allow us to test whether and how human adaptively compress their policies in response to environmental regularities.

Second, our current policy formulation is limited in that it groups both rules and stimuli into a joint state representation. However, participants demonstrate different perseveration patterns at action and rule levels, featuring significant inter-conditional differences in one but not the other. One possible explanation is that different levels of the control hierarchy are subject to distinct information-processing bottlenecks (Sayali et al. 2023)—an idea that merits further investigation. To address this limitation, we plan to develop a hierarchical extension of the policy compression framework that models rule application and action selection as a two-stage process with separate information bottlenecks.

Third, the usage of SOA conditions may delay matured responses, as participants are discouraged from responding before SOA onset. However, it is unlikely that the current RT results stem solely from this delaying effect, as we have observed significant changes in policy complexity, accuracy, and action entropy across SOA conditions. The above phenomenon could not be explained by the delaying account, which would predict similar behavioral statistics despite differences in RT. To further mitigate delay effects, we have repeated the RT regression analysis for the short SOA condition only, which is least prone to relevant issues. As predicted by the framework, the regression model featuring policy complexity again yielded positive effects (0.149 ± 0.0353 , $t(39) = 4.20$, $P < 0.001$), and won against a null model containing only the intercept ($\Delta\text{BIC} = -11.0$). Future empirical studies could remove this confound by replacing SOA deadlines with RT deadlines common to policy compression studies (Lai and Gershman 2024; Liu and Gershman 2025), which would similarly compel responses but never delay them.

While previous studies have already identified neural correlates of policy cost in striatal dopamine (Gershman and Lak 2025) and hypothalamic hypocretin/orexin (Tesmer 2025), the initial site of policy cost computation, and how it is used by downstream areas, remains unclear. In addition, previous studies have focused on assessing the existence of policy cost representations in the brain, without providing a process-level account of how such costs were derived in the first place.

Our contribution here is introducing neural representational dimensionality as a cortical correlate of policy compression, thereby connecting it to the information-theoretic view of cognitive control. We found that greater decoding accuracy over

binary state label partitions—a trial-level proxy for increased neural dimensionality—was correlated with higher policy cost. This finding aligns with the entropy coding perspective, where high-dimensional, fine-grained representations are required for supporting highly state-dependent policies. Beyond single-trial correlations, we also observed significant differences in neural representational dimensionality across SOA conditions. Although the relationship was relatively weak, the observed trends encourage future experiments specifically designed to test the neural-level predictions of policy compression.

The neural correlates of policy complexity remains a large topic that warrants additional investigation. One prominent future direction entails developing biologically plausible neural circuits for implementing near-optimal entropy codes. Combined with causal intervention studies, such models could further inform how the brain may implement policy compression. Another potential direction involves examining the development and plasticity of these high-dimensional neural representations. While policy compression itself does not prescribe a specific learning algorithm, studying how neural representations gradually attain optimality through practice (Kikumoto et al. 2024b) may help explain human deviations from the normative predictions of policy compression. Lastly, recent studies have highlighted the role of oscillations within EEG frequency bands in supporting goal-directed action (Beste et al. 2023; Ericson et al. 2025). Although our study focused on the representational dimensionality of EEG signals, future work could explore potential links between our cross-frequency dimensionality measures and oscillatory activity within individual bands.

Acknowledgments

We thank Jacob Zavatone-Veth for helpful comments and suggestions. We also thank Apoorva Bhandari and Kazuhisa Shibata for their contribution to a previous paper introducing the analyzed dataset.

Author contributions

Shuze Liu (Conceptualization, Formal analysis, Investigation, Methodology, Project administration, Visualization, Writing—original draft, Project conception), Atsushi Kikumoto (Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Software, Writing—review & editing), David Badre (Funding acquisition, Writing—review & editing), and Samuel J. Gershman (Conceptualization, Methodology, Supervision, Writing—review & editing).

Funding

This work was supported by funding from a Schmidt Sciences Polymath Award to S.G., a fellowship from the Kempner Institute of the Study of Natural and Artificial Intelligence to S.L., the National Institute of Mental Health (R01 MH125497), the National Institute of Neurological Disorders and Stroke (R21 NS108380), a Multidisciplinary University Research Initiative award from the Office of Naval Research (N00014-16-1-2832) to D.B. and from JSPS KAKENHI Grant Number 19H01041, 20H05715, JST Moonshot R&D JPMJMS2013 to Kazuhisa Shibata and the JSPS Overseas Research Fellowships to A.K.

Conflict of interest statement: None declared.

Data and code availability

The dataset is available at <https://neurodata.riken.jp/id/20240831-001>. The analysis code is available at https://github.com/LSZ2001/policycompression_neuraldim.

References

- Ahlheim C, Love BC. 2018. Estimating the functional dimensionality of neural representations. *NeuroImage*. 179:51–62. <https://doi.org/10.1016/j.neuroimage.2018.06.015>.
- Archer E, Park IM, Pillow JW. 2014. Bayesian entropy estimation for countable discrete distributions *J Mach Learn Res*. 15:2833–2868.
- Arimoto S. 1972. An algorithm for computing the capacity of arbitrary discrete memoryless channels. *IEEE Trans Inf Theory*. 18:14–20. <https://doi.org/10.1109/TIT.1972.1054753>.
- Aron AR. 2007. The neural basis of inhibition in cognitive control. *Neuroscientist*. 13:214–228. <https://doi.org/10.1177/1073858407299288>.
- Badre D. 2008. Cognitive control, hierarchy, and the rostro-caudal organization of the frontal lobes. *Trends Cogn Sci*. 12:193–200. <https://doi.org/10.1016/j.tics.2008.02.004>.
- Badre D. 2024. Cognitive control. *Annu Rev Psychol*. 76:167–195.
- Badre D, Bhandari A, Keglovits H, Kikumoto A. 2021. The dimensionality of neural representations for control. *Curr Opin Behav Sci*. 38:20–28. <https://doi.org/10.1016/j.cobeha.2020.07.002>.
- Bari BA, Krystal AD, Pizzagalli DA, Gershman SJ. 2024. Computationally-informed insights into anhedonia and treatment by kappa-opioid receptor antagonism. *medRxiv*. <https://doi.org/10.1101/2024.04.09.24304873>.
- Bernardi S et al. 2020. The geometry of abstraction in the hippocampus and prefrontal cortex. *Cell*. 183:954–967.e21. <https://doi.org/10.1016/j.cell.2020.09.031>.
- Beste C, Munchau A, Frings C. 2023. Towards a systematization of brain oscillatory activity in actions. *Commun Biol*. 6:137. <https://doi.org/10.1038/s42003-023-04531-9>.
- Blahut R. 1972. Computation of channel capacity and rate-distortion functions. *IEEE Trans Inf Theory*. 18:460–473. <https://doi.org/10.1109/TIT.1972.1054855>.
- Bossaerts P, Yadav N, Murawski C. 2019. Uncertainty and computational complexity. *Philos Trans R Soc B*. 374:20180138. <https://doi.org/10.1098/rstb.2018.0138>.
- Botvinick M, Braver T. 2015. Motivation and cognitive control: from behavior to neural mechanism. *Annu Rev Psychol*. 66:83–113. <https://doi.org/10.1146/annurev-psych-010814-015044>.
- Botvinick MM, Cohen JD. 2014. The computational and neural basis of cognitive control: charted territory and new Frontiers. *Cogn Sci*. 38:1249–1285. <https://doi.org/10.1111/cogs.12126>.
- Courellis HS et al. 2024. Abstract representations emerge in human hippocampal neurons during inference. *Nature*. 632:841–849. <https://doi.org/10.1038/s41586-024-07799-x>.
- Cousineau D. 2005. Confidence intervals in within-subject designs: a simpler solution to Loftus and Masson's method. *Tutor Quant Methods Psychol*. 1:42–45. <https://doi.org/10.20982/tqmp.01.1.p042>.
- Cover TM. 1999. *Elements of Information Theory*. John Wiley & Sons.
- Ericson J, Ruiz Ibáñez N, Lundqvist M, Klingberg T. 2025. Low frequency oscillations—neural correlates of stability and flexibility in cognition. *Nat Commun*. 16:1–14. <https://doi.org/10.1038/s41467-025-60821-2>.
- Fan J. 2014. An information theory account of cognitive control. *Front Hum Neurosci*. 8:680.
- Fusi S, Miller EK, Rigotti M. 2016. Why neurons mix: high dimensionality for higher cognition. *Curr Opin Neurobiol*. 37:66–74. <https://doi.org/10.1016/j.conb.2016.01.010>.
- Gailliot MT, Baumeister RF. 2007. The physiology of willpower: linking blood glucose to self-control. *Personal Soc Psychol Rev*. 11:303–327. <https://doi.org/10.1177/1088868307303030>.
- Genewein T, Leibfried F, Grau-Moya J, Braun DA. 2015. Bounded rationality, abstraction, and hierarchical decision-making: an information-theoretic optimality principle. *Front Robot AI*. 2:27.
- Gershman SJ. 2020. Origin of perseveration in the trade-off between reward and complexity. *Cognition*. 204:104394. <https://doi.org/10.1016/j.cognition.2020.104394>.
- Gershman SJ, Bhui R. 2020. Rationally inattentive intertemporal choice. *Nat Commun*. 11:3365. <https://doi.org/10.1038/s41467-020-16852-y>.
- Gershman SJ, Burke T. 2023. Mental control of uncertainty. *Cogn Affect Behav Neurosci*. 23:465–475. <https://doi.org/10.3758/s13415-022-01034-8>.
- Gershman SJ, Lak A. 2025. Policy complexity suppresses dopamine responses. *J Neurosci*. 45:e1756242024. <https://doi.org/10.1523/JNEUROSCI.1756-24.2024>.
- Goghari VM, MacDonald AWIII. 2009. The neural basis of cognitive control: response selection and inhibition. *Brain Cogn*. 71:72–83. <https://doi.org/10.1016/j.bandc.2009.04.004>.
- Hick WE. 1952. On the rate of gain of information. *Q J Exp Psychol*. 4:11–26. <https://doi.org/10.1080/17470215208416600>.
- Huffman DA. 1952. A method for the construction of minimum-redundancy codes. *Proc IRE*. 40:1098–1101. <https://doi.org/10.1109/JRPROC.1952.273898>.
- Hutter M. 2001. Distribution of mutual information. *Adv Neural Inf Proces Syst*. 14.
- Hyman R. 1953. Stimulus information as a determinant of reaction time. *J Exp Psychol*. 45:188–196. <https://doi.org/10.1037/h0056940>.
- Jakob AM, Gershman SJ. 2023. Rate-distortion theory of neural coding and its implications for working memory. *Elife*. 12:e79450. <https://doi.org/10.7554/eLife.79450>.
- Jazayeri M, Ostojic S. 2021. Interpreting neural computations by examining intrinsic and embedding dimensionality of neural activity. *Curr Opin Neurobiol*. 70:113–120. <https://doi.org/10.1016/j.conb.2021.08.002>.
- Kikumoto A, Mayr U, Badre D. 2022. The role of conjunctive representations in prioritizing and selecting planned actions. *Elife*. 11:e80153. <https://doi.org/10.7554/eLife.80153>.
- Kikumoto A, Bhandari A, Shibata K, Badre D. 2024a. A transient high-dimensional geometry affords stable conjunctive subspaces for efficient action selection. *Nat Commun*. 15:8513. <https://doi.org/10.1038/s41467-024-52777-6>.
- Kikumoto A, Shibata K, Nishio T, Badre D. 2024b. Practice reshapes the geometry and dynamics of task-tailored representations bioRxiv.
- Koechlin E, Summerfield C. 2007. An information theoretical approach to prefrontal executive function. *Trends Cogn Sci*. 11:229–235. <https://doi.org/10.1016/j.tics.2007.04.005>.
- Kool W, McGuire JT, Rosen ZB, Botvinick MM. 2010. Decision making and the avoidance of cognitive demand. *J Exp Psychol Gen*. 139:665–682. <https://doi.org/10.1037/a0020198>.
- Kool W, Shenhav A, Botvinick MM. 2017. Cognitive control as cost-benefit decision making. In *The Wiley Handbook of Cognitive Control*. Wiley, pp 167–189.
- Krebs RM, Woldorff MG. 2017. Cognitive control and reward. In *The Wiley Handbook of Cognitive Control*. Wiley, pp 422–439.

- Krueger PM, Griffiths TL. 2018. Shaping model-free habits with model-based goals. In: *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume. 40.
- Kruglanski AW et al. 2012. The energetics of motivated cognition: a force-field analysis. *Psychol Rev*. 119:1–20. <https://doi.org/10.1037/a0025488>.
- Lai L, Gershman SJ. 2021. Policy compression: an information bottleneck in action selection. In: Federmeier K.(ed) *Psychology of Learning and Motivation*, Vol. 74. Elsevier, pp 195–232.
- Lai L, Gershman SJ. 2024. Human decision making balances reward maximization and policy compression. *PLoS Comput Biol*. 20:1–32. <https://doi.org/10.1371/journal.pcbi.1012057>.
- Liu S, Gershman S. 2025. Action subsampling supports policy compression in large action spaces. *PsyArXiv*. <https://doi.org/10.31234/osf.io/bjt7d>.
- Liu S, Lai L, Gershman SJ, Bari BA. 2024. Time and memory costs jointly determine a speed-accuracy trade-off and set-size effects. *PsyArXiv*.
- Ma WJ, Husain M, Bays PM. 2014. Changing concepts of working memory. *Nat Neurosci*. 17:347–356. <https://doi.org/10.1038/nn.3655>.
- Matsumoto K, Tanaka K. 2004. Conflict and cognitive control. *Science*. 303:969–970. <https://doi.org/10.1126/science.1094733>.
- McGuire JT, Botvinick MM. 2010. Prefrontal cortex, cognitive control, and the registration of decision costs. *Proc Natl Acad Sci*. 107: 7922–7926. <https://doi.org/10.1073/pnas.0910662107>.
- Miller EK, Cohen JD. 2001. An integrative theory of prefrontal cortex function. *Annu Rev Neurosci*. 24:167–202. <https://doi.org/10.1146/annurev.neuro.24.1.167>.
- Miller KJ, Shenhav A, Ludvig EA. 2019. Habits without values. *Psychol Rev*. 126:292–311. <https://doi.org/10.1037/rev0000120>.
- Musslick S et al. 2016. Parallel processing capability versus efficiency of representation in neural networks. *Network*. 8.
- Nebe S, Kretschmar A, Brandt MC, Tobler PN. 2024. Characterizing human habits in the lab. *Collabra Psychol*. 10:1–45.
- Pauli WM, Cockburn J, Pool ER, Perez OD, O'Doherty JP. 2018. Computational approaches to habits in a model-free world. *Curr Opin Behav Sci*. 20:104–109. <https://doi.org/10.1016/j.cobeha.2017.12.001>.
- Rigotti M et al. 2013. The importance of mixed selectivity in complex cognitive tasks. *Nature*. 497:585–590. <https://doi.org/10.1038/nature12160>.
- Sayali C, Rubin-McGregor J, Badre D. 2023. Policy abstraction as a predictor of cognitive effort avoidance. *J Exp Psychol Gen*. 152: 3440–3458. <https://doi.org/10.1037/xge0001449>.
- Shannon CE. 1948. A mathematical theory of communication. *Bell Syst Tech J*. 27:379–423. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>.
- Shenhav A et al. 2017. Toward a rational and mechanistic account of mental effort. *Annu Rev Neurosci*. 40:99–124. <https://doi.org/10.1146/annurev-neuro-072116-031526>.
- Sims CR. 2015. The cost of misremembering: inferring the loss function in visual working memory. *J Vis*. 15:2–2. <https://doi.org/10.1167/15.3.2>.
- Sims CR. 2016. Rate-distortion theory and human perception. *Cognition*. 152:181–198. <https://doi.org/10.1016/j.cognition.2016.03.020>.
- Sims CR, Jacobs RA, Knill DC. 2012. An ideal observer analysis of visual working memory. *Psychol Rev*. 119:807–830. <https://doi.org/10.1037/a0029856>.
- Sutton RS, Barto AG. 2018. *Reinforcement Learning: an Introduction*. MIT press.
- Tesmer AL et al. 2025. Neurometabolic signaling and control of policy complexity bioRxiv.
- Tishby N, Polani D. 2010. Information theory of decisions and actions. In: Cutsuridis V, Hussain A, Taylor J. (eds) *Perception-Action Cycle: Models, Architectures, and Hardware*. Springer, New York, NY, pp 601–636.
- Todorov E. 2009. Efficient computation of optimal actions. *Proc Natl Acad Sci*. 106:11478–11483. <https://doi.org/10.1073/pnas.0710743106>.
- Umemoto A, Holroyd CB. 2015. Task-specific effects of reward on task switching. *Psychol Res*. 79:698–707. <https://doi.org/10.1007/s00426-014-0595-z>.
- Vul E, Goodman N, Griffiths TL, Tenenbaum JB. 2014. One and done? Optimal decisions from very few samples. *Cogn Sci*. 38:599–637. <https://doi.org/10.1111/cogs.12101>.
- Zenon A, Solopchuk O, Pezzulo G. 2019. An information-theoretic perspective on the costs of cognition. *Neuropsychologia*. 123:5–18. <https://doi.org/10.1016/j.neuropsychologia.2018.09.013>.
- Zhang C, van Wissen A, Dotsch R, Lakens D, W. A. IJsselstein. 2024. A sequential sampling approach to the integration of habits and goals. *Comput Brain Behav*. 7:480–501.
- Zhou H, Nagy DG, Wu CM. 2024. Harmonizing program induction with rate-distortion theory arXiv preprint arXiv:2405.05294.