

# Cognitive graphs: Representational substrates for planning

Jungsun Yoo<sup>a, \*</sup>, Elizabeth R. Chrastil<sup>a, b, c</sup>, Aaron M. Bornstein<sup>a, b</sup>

<sup>a</sup>Department of Cognitive Sciences, University of California, Irvine, Irvine, CA, USA, 92697

<sup>b</sup>Center for the Neurobiology of Learning and Memory, University of California, Irvine, CA, USA, 92697

<sup>c</sup>Department of Neurobiology & Behavior, University of California, Irvine, Irvine, CA, USA, 92697

\*To whom correspondence should be addressed: jungsuy@uci.edu

**Abstract.** Making plans for upcoming actions is a computationally demanding process. To mitigate these demands, individuals can build extensive internal models of their environment – states, actions, and their sequential relationships – that allow for plans to be developed with minimal computational costs. Initially, these models reflect elaborate networks of learned associative relationships, which can be used to generate plans for reward through more iterative computations such as trajectory sampling. After sufficient experience, compressed forms of these models can efficiently capture long-range sequential structure, allowing them to be used for rapid planning even in pursuit of novel or changing rewards. Here, we review recent work on the multitude of representations that can support different forms of planning. We discuss how *cognitive graphs*, a framework with roots in both cognitive psychology and computer science, can provide a unifying view of these representations and their relationships to one another. Conceptualizing internal models as forms of graphs situates them on a spectrum where different kinds of structured sequences can be queried to support both planning and the formation of iteratively more compressed predictive representations. We discuss how each of these kinds of cognitive graphs are created during learning, and used to transfer and generalize knowledge across environments. Taken together, this review highlights the significant impact that the various associative structures of memory have on planning.

**Keywords:** cognitive graph, planning, reinforcement learning.

## 1 Introduction

Planning is a common, and complex, form of decision-making. It requires both representing actions, along with their precedents and consequences, and sequencing them appropriately. This process consists of *offline* stages, during which predictive representations of the environment are formed and refined, and *online* stages, where extant representations of relevant past experience are interrogated and their predictions for the outcomes of planned choices are arbitrated.

The term ‘representation’ used here refers to how elements of a given decision problem are encoded in memory and associated with each other. Choices are strongly influenced by the format in which decision elements are arranged when presented to an individual – for instance, risk attitudes often vary considerably when options are presented as explicit frequencies, rather than summary probabilities (Kahneman and Tversky, 1979). More recently, researchers have begun to systematically explore how choices depend on the ways in which decision-relevant information, such as state spaces, are represented *internally* by the decision-maker (Doya et al., 2002; Wang, Feng, and Bornstein, 2022). This work shows that the choice of internal representations can have a similarly dramatic influence on the outcome of a decision. For example, individuals who remember their local environment as a series of routes they have taken (“egocentric” representation) may be unlikely to try a novel route in face of a detour, unlike individuals who have integrated their experiences to form a map-like (“allocentric”) summary of the environment (Chrastil and Warren,

2014). This example highlights how internal, unlike external, representations, can be transformed from one format into another given sufficient experience and/or time; critically, this process can yield many intermediate formats, where some information is retained and other information is lost.

In the case of planning, representation format is critical in part because much planning occurs ahead of time, by constructing a semi-flexible *policy* that establishes the rules by which sequences of actions are to be taken. In these cases, the selection of which *kind* of internal model provides the state space over which the policy is defined, and thus critically determines the actions ultimately taken (Ho et al., 2022). The importance of representational format in planning is further underscored by its role in *transfer learning*, which requires first identifying similar situations from the past and subsequently selecting the relevant aspects of that previously learned structure. When experience in related environments is extensive, allowing the agent to infer common latent structure, one could apply compact, “map-like” representations that allow for efficient planning with minimal error (Geerts et al., 2022; Whittington et al., 2020). However, as the overlap between well-learned settings and the current environment decreases, one must rely on approximations to identify relevant instances of previous experiences with the current or similar environments (Zhao, Richie, and Bhatia, 2022). Internal simulations informed by these sorts of instance samples can be used for iterative, vicarious evaluation of decision problems that not only informs the decision at hand, but allows the agent to accelerate the inference of more general latent structure (George et al., 2021).

An implication of this *representation-centric* view of planning is that a key problem for agents to solve is how to summarize the available experience in a way that best supports efficient and effective planning and transfer learning. The type of summary representation best suited to each situation thus depends on the complexity of the environment, the amount of experience the agent has in it, and the time and computational resources available to evaluate candidate policies; these quantities are often dynamic or not known ahead of time, thus licensing the agent to maintain multiple representations that can be leveraged to different degrees in different settings (Doya et al., 2002; Wang, Feng, and Bornstein, 2022).

We propose that these many distinct forms of internal representations – associative relationships – can be fruitfully understood as types of graphs (Butts, 2009). Here, environmental states are represented as nodes and the transitions between them are shown as various types of edges (Schapiro et al., 2013; Lynn et al., 2020), depending on the information available (Chrastil and Warren, 2014). The edges could be either unidirectional when describing causality or irreversible transitions, or could be bidirectional when these conditions are not assumed. For instance, a decision tree is a specific example of a graph that encodes sequential, or unidirectional, relationships between states (Bertsekas, 2012). Formalizing these structures as graphs can allow researchers to formally connect seemingly disparate types of planning, to reason about their related algorithmic and implementational properties (Zhang, Yang, and Stadie, 2021), and to determine how and

which information is transferred (“consolidated”) from one format to another, e.g. during sleep (Feld et al., 2022).

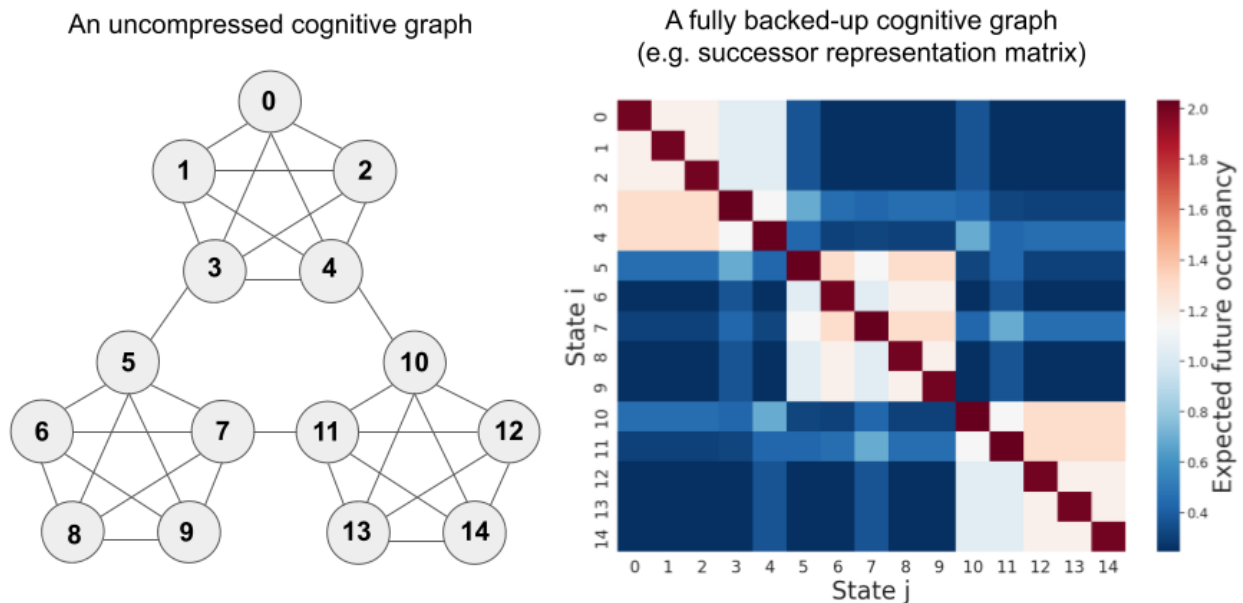


Figure 1: A graphical illustration of the two extremes of representations as a function of compression: an uncompressed cognitive graph (left) and a fully backed-up cognitive graph (successor representation used as an example here) (right). **left** Numbers indicate labels of a node, or a discrete state, in a cognitive graph. Edges between two nodes depict the transition between two states. **right** A backed-up version of a cognitive graph that fully captures future trajectories from a given node (state). Row identification numbers (ID;  $i$ ) indicate the current state, and column IDs ( $j$ ) indicate the successor state. The values in the matrix represent the expected general future occupancy of  $j$  from  $i$ , and are color-coded for visualization. Note that while future occupancy statistics preserve the coarse community structure, route information is diminished (e.g. the adjoining gateway nodes, such as 3 and 11 for the cluster of states 5-9, are only slightly distinguished.)

Graph-like forms useful in planning range between extremes – at one end, sets of instances of individual pairwise associations; at the other, compact, long-range multi-step contingencies – with many points along the spectrum between these (Chrastil and Warren, 2014, Figure 1). Recent work supports the simultaneous creation and updating of multiple graph-like knowledge structures in support of planning. These internal models are distinguished by their content, format, and also in what they entail for the dynamics of their learning and use in deliberative decision-making (Doya et al., 2002; Bornstein and Daw, 2012, 2013; Smith and Graybiel, 2013; Tambini et al., 2023). Below, we review findings that suggest that they influence behavior in accordance to their suitability to the task at hand, and that the apparent shift in behavioral control from one form to others is characterized by the transformation of information between representational formats, with attendant trade-offs in function and fidelity.

Theoretical (Weber and Johnson, 2006) and empirical (Otto et al., 2022; Palminteri et al.,

2015; Wu, Schulz, and Gershman, 2021) work supports the idea that the evaluation of an option depends in part on how that option is remembered — for instance, if it is remembered as part of a set of related options, with ranked preferences within that set (e.g. a favored restaurant among those of similar cuisine), or if its visible features are associated with other latent features (e.g. a food attribute linked to allergic reactions). Foundational work has demonstrated that successive memory retrievals are related to the underlying associative structure of memory (Howard and Kahana, 2002), supporting a form of trajectory sampling (Gershman and Daw, 2017; Wang, Feng, and Bornstein, 2022), and that the content of extended memory retrieval at the time of choice has meaningful influence on preferences (Bornstein and Norman, 2017). Taken together, this work supports a critical influence on choice of the associative structure of memory. Therefore, it is important to understand the different forms this structure can take, and to identify commonalities and points of divergence relevant to choice behavior.

## 2 Cognitive graphs

These associative structures can be understood as forms of cognitive graphs, that range from “uncompressed” to “compressed” (Figure 1). The most uncompressed form, in which states are encoded as experienced sequences with minimal latent structure inference, conceptually aligns with previous articulations of “cognitive graphs” (Muller, Stead, and Pach, 1996; Chrastil and Warren, 2014; George et al., 2021), and that is proposed to support types of model-based reinforcement learning (Daw, Niv, and Dayan, 2005; Lengyel and Dayan, 2007; Gershman and Daw, 2017). A cognitive graph can be characterized as a directed graph (Muller, Stead, and Pach, 1996), with nodes representing states and edges indicating state transitions. These edges may be labeled, augmenting the topology with local metric information (Chrastil and Warren, 2014; Warren, 2019). They may also be weighted, reflecting the transition probability between states (Natarajan and Kolobov, 2022; George et al., 2021; Sutton and Barto, 2018). A cognitive graph is formed through learning how different sequences of state transitions connect at intersections (Stiso et al., 2022), enabling agents to flexibly navigate conceptual and spatial networks by recombining the segments in novel ways (Warren, 2019; Mark et al., 2020; Peer et al., 2021). Additionally, their abstract nature supports counterfactual simulations and generalizations to novel environments, thereby accelerating the learning process (Zhu et al., 2020). Though the entire continuum of representations are graph-structured, we will for clarity refer to the most uncompressed extreme form as “full” or “flexible” graphs, and the most compact representations as “backed-up” or “compressed” predictive representations.

At the other end of the spectrum, backed-up, predictive representations contain information that is fully predictive of the  $N$ -step consequence of taking a given action  $a$  in the current state  $s$  (Figure 1, right). To elaborate, a standard model of choice describes preferences between options as formed after a unitary expected value is computed by combining the reward distributions

implied by each options’ features (Rangel, Camerer, and Montague, 2008). These values – both unitary and the components – can be represented in different ways, each of which has different implications for the preference construction process. Backed-up representations enable fast, cheap evaluation of N-step plan outcomes, using an operation akin to matrix multiplication (though the neural instantiation of this process has yet to be fully described (Gershman, 2018) and may be approximated by sampling (Gershman et al., 2012)). For example, model-free reinforcement learning of action values (Sutton and Barto, 2018) captures this unitary value as a recency-weighted average of the discounted total reward obtained in past episodes where the agent took the given action in the given state. Here, the outcome values of multi-step actions are mediated by a discount factor,  $\delta$ , applied at each update operation. An alternative approach to constructing unitary values is to use a backed-up representation of the discounted N-step state occupancy alone, irrespective of reward obtained, which allows decoupling the environmental state dynamics, which may be more stable, from reward contingencies that may fluctuate more often or be entirely trial-unique. Such *successor representations* (or their mirror, *predecessor representations*; Jeong et al., 2022) compress occupancy of sequences following or preceding a given state (Dayan, 1993), which can be used to derive biological cell response types matching those observed in subfields of the hippocampal formation (Stachenfeld, Botvinick, and Gershman, 2017). There are several related formats that differ in what information is included in the backed-up representation, such as successor features (Barreto et al., 2017) – which generalize the state-space learning approach to a space over option dimensions (e.g., desirability for food) – and first-occupancy representations (Moskovitz, Wilson, and Sahani, 2021) that only consider the first-time visits to each state. Inspired by the need to bridge the gap between behavioral economics and reinforcement learning,  $\lambda$ R incorporates the concept of diminishing marginal utility (Moss, 1984) into reinforcement learning by discounting multiple visits to a state, thereby providing an intermediate representation between successor representations and first-occupancy representations (Moskovitz et al., 2023).

Between the extremes of compressed versus flexible-model representations, cognitive graphs with intermediate modes of approximation can also be identified. We described above how the discount factor allows the successor representation to be parametrically distinguished from the outcomes of Monte Carlo trajectory sampling from a full model. Another axis along which these representations can vary in their approximation of the full environment dynamics is the degree to which they reflect hierarchical structure. For example, agents may cluster or abstract related states as intermediate “sub-goals” that exist in multiple levels hierarchically to plan efficiently (Tomov et al., 2020; Noh et al., 2023). Compression can also occur by compressing *actions or policies* into higher-level actions, referred to as option or skill discovery (Sutton, Precup, and Singh, 1999). Automated discovery of options at multiple levels has facilitated learning in artificial agents (Fox et al., 2017). Likewise, humans appear to adopt policy compression to balance cognitive costs and maximizing reward (Lai, Huang, and Gershman, 2022; Lai and Gershman, 2021). Similar to this,

extracting a causal relationship between events at various levels of granularity could be seen as an abstraction or compression of the environment (Kinney and Lombrozo, 2023a,b).

In this paper, we initially delve into the differences between planning predicated on the flexible recombination of action sequences and planning employing compressed representations. Subsequently, we propose cognitive graphs as a potential unifying framework supporting both decisions based on sampling potential future sequences and decisions based on full long-run occupancy statistics. We discuss how these functions require key mechanisms – in particular, merging disjoint sequences and splitting aliased states – offered by some implementations of cognitive graphs. We conclude with a discussion of further research directions, in particular understanding how the spectrum of forms of cognitive graph may support distinct control strategies. This, in turn, could potentially clarify the differential use of types of control in different stages of learning.

### 3 Planning as a function of representational compactness

#### 3.1 Planning in a Markov Decision Process

For simplicity, planning is often conceptualized within the context of a Markov Decision Process (MDP). In a classical MDP, the environment in which an agent plans is characterized as a tuple of  $\langle S, A, T, R, \pi, \gamma \rangle$  where  $S$  is a finite and discrete state space that is comprised of states, and  $A$  is a set of actions that can be executed in each state  $s \in S$ .  $\gamma$  refers to the discount factor that represents how future rewards are valued in comparison to immediate rewards. The models consist of two functions, where  $T(s, a, s')$  is the transition function for each  $s \in S$  and  $a \in A$ , and  $R(s, a, s')$  is a reward function that provides the immediate reward or value obtained after taking action  $a$  in state  $s$  and transitioning to state  $s'$ .  $\pi$  refers to the policy, or the probability distribution of the actions at a certain state. We assume that an agent starts from an initial state  $s_0$  and executes a sequence of possible actions in the successor states ( $s'$ ) up to a terminal or goal state  $s \in S_G$ . The agent’s goal in planning is to learn and execute actions based on an optimal solution, or policy, that maximizes the cumulative value from an initial state to a goal state.

One of the most crucial components for successful planning is having an accurate internal model of the environment, because the model is used for simulating or predicting behavior; inaccurate models could entail incorrect predictions and thereby result in a chain of sub-optimal actions (Talvitie, 2017). It is also important to adopt the most suitable models for each specific context, given that optimal type of model to use may vary depending on the relationship of the model to the environment – for instance, whether the model is known with certainty to correspond exactly to the environment (Jiang et al., 2015). Below, we delve into the kind of model utilization that may be optimal in scenarios where agents are still in the preliminary phases of environment interaction (*Section 3.2.1 Planning with uncompressed representations: sampling instances*), or in circumstances where they possess sufficient experience for compression of representations to occur (*Section 3.2.2 Planning with backed-up state/action sequences*).



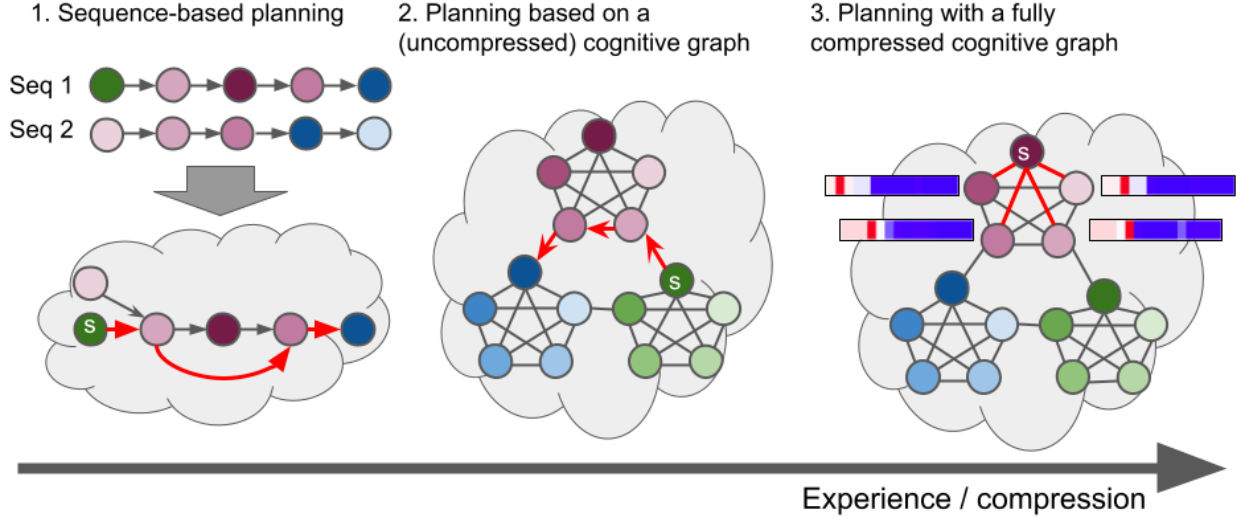


Figure 2: A graphical illustration of planning based on the suggested spectrum of representations – specifically, the degree to which they are pre-compiled – as a function of experience. Nodes indicated with ‘s’ represent starting points. Arrows, or directed edges, describe possible plans for the agent in the starting node to execute. The grey clouds represent the form of model the agent is using to plan in a given phase. **1. Sequence-based planning** This phase represents the early learning phase where an agent has not yet constructed an integrated model of the environmental dynamics. Here, agents are assumed to plan based on sampling instances previously experienced trajectories. First, two sample trajectories are shown here, labeled ‘seq 1’ and ‘seq 2.’ From these two trajectories, an agent is able to create a combined representation and plan efficiently with it (e.g., taking the shortcut as seen in the red arrows). **2. Planning based on an uncompressed cognitive graph** After a few experiences, an agent is able to build cognitive graphs by conjoining past trajectories. Agents are thought to iteratively sample next actions based on the cognitive graph. As agents gain more experience, simultaneously, a ‘diffusion-like’ process is thought to take place such that information about neighboring nodes are integrated into each node. **3. Planning with a fully compressed cognitive graph** After sufficient experience, a compressed trajectory from a given node to each other node is available in a summary format. Rows of such a *successor representation* are shown here. The availability of information in this representational format allows agents to plan for novel or changing rewards in a statistically efficient manner.

### 3.2 Planning with learned cognitive graphs: uncompressed vs. backed-up

Models in planning capture statistical regularities of the environment, and could be either given *a priori* or learned from experience. If an agent has full information about the transition structure of the environment, then the agent is able to plan even without experience. This is conceptually relevant to a classical control problem or search algorithm (Korf, 1987): for example, for the game Tic-Tac-Toe, an agent can be endowed with a complete model (or a human can be verbally instructed about the rules of the game). Given this starting point, the player can construct a tree-like graph of possible future states and actions, and perform search to find the optimal decision (Sriram et al., 2009). However, in more naturalistic contexts, the dynamics of the environment are

unknown to us initially and our internal models develop and change with our experience with the environment (Schrittwieser et al., 2020; Lengyel and Dayan, 2007). We confine further discussion to these latter, learned models of the environment.

Cognitive graphs at different levels of compression could serve as models that can support different forms of planning (Figure 2). Raw, uncompressed cognitive graphs support planning via iterative sampling of subsequent states or actions from a given state, or node. Here, individual instances or nodes have minimal information about other nodes, thus making it crucial to traverse graphs based on the relationship between nodes, or edges. Thus, in this form of planning, the sampling algorithm is critical.

At the other extreme, actions and states in a cognitive graph are fully backed-up – for instance, successor or predecessor representations. In the successor representation, each node-state contains the expected future state occupancy given a current state and according to a given policy (Dayan, 1993); these can be thought of as integrated trajectories sampled from the current state. Conversely, predecessor representations can be thought of as fully bootstrapped versions of eligibility traces, a memory-like mechanism that assigns credit to past states and actions from a given state (Bailey and Mattar, 2022; Sutton and Barto, 2018). Predecessor representations could be seen as a hindsight version of successor representations in that it bootstraps the possible trajectories *that could have lead to* a current state. Whether directed forward or backward, once these bootstrapped representations converge, the expected cumulative reward or credit can be efficiently computed for planning, just by taking the product of the representation and a separate reward function. With fully backed-up cognitive graphs, the need for edge-based sequences, or sampling successor states, becomes negligible.

### 3.2.1 Planning with uncompressed representations: sampling instances

During early stages of learning an environment, planning could be facilitated by instance-based methods instead of relying on an explicit model, or rule-based methods. Take, for example, the task of choosing a restaurant to dine in an area that one has just moved to and thus has limited experience with. It might be more effective to decide based on a few recent dining experiences rather than attempting to decide based on a general summary of what little experience one has. A model of decision-making under uncertainty captures this intuition as *case-based* decision theory (Gilboa and Schmeidler, 1995), which suggests that to make decisions under uncertainty, people rely on memory of similar cases that had worked well in the past.

This non-parametric, kernel-based method offers several cognitive advantages that could potentially bolster decision efficiency, and provides a better account of human decision making compared to rule-based methods. For instance, a small group of samples reduces memory load (Barron and Erev, 2003), simplifies the decision rule (Fiedler, 2000; Hertwig and Pleskac, 2010), facilitates generalization to unseen observations (Wimmer and Shohamy, 2012; Barron, Dolan, and Behrens,



2013), and reduces time (Hertwig et al., 2004; Fox and Hadar, 2006). Also, instance sampling has been shown to be a superior explanation of decision behavior in several laboratory tasks (Bhatia, 2014; Hotelling and Kellen, 2022; Zhao, Richie, and Bhatia, 2022; Bornstein et al., 2017; Wang, Feng, and Bornstein, 2022). For example, in a repeated decision task, individuals appear to rely on small numbers of samples of recent experiences. When intermittent reminder probes were added to the task, manipulating the apparent recency of past experiences, these probes had a significant effect on subsequent choice (Bornstein et al., 2017). Critically, the choice of instances is sensitive to current task demands: Recent experiences may be more likely to be sampled in an environment that does not have an explicit structure, but the introduction of periodic structure can lead to more adaptive sampling of relevant events (Plonsky, Teodorescu, and Erev, 2015).

An example of case-based decision theory applied to reinforcement learning is episodic control (Dayan, 2008; Lengyel and Dayan, 2007). Episodic control enables agents to make an informed guess about the value of unseen states by averaging the values of the stored past instances that are most similar to the current state. The kernel-based nonparametric approach that underlies case-based decision theory improves sample efficiency since the same amount of observations could be used to inform estimates about a greater number of states, compared to classic reinforcement learning, as well as providing a method for generalization, which is particularly important in continuous state spaces (Gershman and Daw, 2017; Bhui, 2018). Simulation results show that this advantage renders episodic control superior to model-based or model-free control during initial learning stages, as it accelerates the learning process under a low-data limit compared to other control methods (Lengyel and Dayan, 2007; Blundell et al., 2016). One drawback of episodic control is that the search process could be inefficient as the number of episodes stored increases. For scalability, neural episodic control (NEC) uses deep learning methods to embed the keys of each state into a fixed-length vector (Pritzel et al., 2017). Embedded inputs are then fed into a differentiable neural dictionary, or a learnable episodic memory system that maps keys to values. The final value of an observation is obtained by the weighted sum of the values in the differentiable neural dictionary, where the weights are computed by the similarity between the current key and the keys of states in the memory system.

Sampling-based accounts of human multi-step planning have also provided descriptive value in at least two aspects. First, an extension of decision field theory toward the realm of planning – named decision field theory-planning (DFT-P; Hotelling, 2020) – can explain human planning behavior better than backward induction, at least in situations where multi-step plans contend with high payoff variability. Here, the unreliability of experience may be a critical factor favoring this instance-based approach. In large, continuous, and highly uncertain environments an agent would require unrealistically extensive experience to develop stable, compressed, and predictive representations. Silver and Veness (2010) show that in these environments, asymptotically optimal plans can be constructed using Monte Carlo trajectory sampling over an iteratively updated internal

model. A second advantage of representing the full, uncompressed, model of the state space, with all its intermediate states is that it supports effective exploration strategies: in particular, one can perform ‘far’ jumps across state spaces to distal, weakly connected nodes (Zhu, Sanborn, and Chater, 2018); the resulting “Lévy flight” behavior matches observations of biological agents exploring novel environments (Hunt et al., 2021) and performing memory search (Rhodes and Turvey, 2007).

Research has shown that in an environment based on graph-like reward structure (e.g., subway maps), people leverage learned graph structure to guide sampling-based decisions (Wu, Schulz, and Gershman, 2019, 2021). Nevertheless, whether people are able to *spontaneously* construct cognitive graphs from sequential experiences in graph-like structures and still leverage this to guide decision has not yet been directly investigated.

### 3.2.2 *Planning with backed-up state/action sequences*

Earlier, we introduced the concept of backed-up representations as a way of incrementally learning compact summaries of multi-step contingencies. Successor representations have been devised to balance the possible computational intractability of fully model-based methods and inflexibility of computationally cheap model-free methods, providing a robust solution to this problem (Dayan, 1993). These compressed, predictive representations summarize expected future occupancy of successor states from a current state given a policy. Using successor representations compresses the multi-step planning process into a single-step process, since long-range outcomes of all possible future trajectories are considered at once (Dayan, 1993). This not only reduces computational complexity, but it also facilitates generalization and learning when adapting to variable reward contingencies. Empirical evidence from studies conducted on humans (Momennejad et al., 2017) and artificial agents (Barreto et al., 2017) suggests that using transition dynamics compressed in a successor representation lead to faster adaptation to value-function changes, because only the reward function requires re-learning, thus significantly enhancing learning efficiency.

Another example of compressing sequences of observation, or states, is seen in robust predictable control (Eysenbach, Salakhutdinov, and Levine, 2021). This algorithm is explicitly encouraged to find a compressed policy by penalizing complexity, which is operationally defined as the amount of information needed from observations for a policy to make decisions. The intuition behind this is that agents will rely less on gathering information from observations as they become better at predicting the future accurately. Agents trained on compressed policies are less susceptible to unknown or missing observations (i.e., perturbations), since compressed policies have been trained to use fewer bits of information per observation. This leads to improved *open-loop control* – producing a plan of action sequences at the beginning and executing it without checking the progress along the way.

In sum, compressed representations lower the cost of planning by reducing complexity at the representational level. This kind of representation also fosters open-loop planning by enabling the execution of action sequences as a single operation (Eysenbach, Salakhutdinov, and Levine, 2021). This could be efficient in environments where transition dynamics are relatively well-known and unchanging. On the other hand, when models of the environment have not been fully developed yet, instance-based control can be useful. In particular, sampling trajectories of instances to preserve the sequential nature of experiences provides a method with less complexity and greater scalability, while still maintaining high performance. In the following section, we discuss how graph-structured representations can improve trajectory-based planning.

#### **4 Ways in which uncompressed cognitive graphs could facilitate planning**

Recent approaches transform planning into a graph-search problem (Savinov, Dosovitskiy, and Koltun, 2018; Liu et al., 2020). One study leveraged graph-based representations to identify landmarks or subgoals in latent graphs, and then performed graph search on the nodes (Zhang, Yang, and Stadie, 2021). Here, edges between the nodes are weighted with “reachability” between nodes, making it as a form of a labeled graph. In the domain of spatial navigation, algorithms construct graphs based on subgoals and then plan based on the constructed graphs for efficiency (Bagaria, Senthil, and Konidaris, 2021).

It has also been found that people spontaneously construct graph-like representations when observing a sequence of events, where these latent graphs could be either correlational (Rmus et al., 2022; Solomon et al., 2019; Kahn et al., 2018; Schapiro et al., 2013, undirected graphs) or causal (Gopnik and Schulz, 2004; Gopnik et al., 2004; Sommerville and Woodward, 2005a,b, directed graphs). Furthermore, people have been shown to be able to capture the topological structure of an underlying graph (i.e., identifying bottleneck states; Schapiro et al., 2013; Solway et al., 2014), even after passive observation of trajectories through the graph space. Intriguingly, the general tendency to use plans over model-free approaches appears to be correlated with the ability to infer latent graph-based structure from jumbled sequences of experiences (Rmus et al., 2022), potentially underscoring the utility of learning graph-structured representations in planning.

##### *4.1 Mechanisms by which cognitive graphs could facilitate planning*

States and observations or instances may not be mapped onto each other in a one-to-one fashion. This phenomenon, referred to as perceptual aliasing, could potentially destabilize control in reinforcement learning (Whitehead and Ballard, 1991). To overcome this, agents must employ an accurate and parsimonious representation of experience that is able to split identical observations into different underlying states or merge seemingly different observations into a single state for generalization, depending on the context (Niv, 2019). In other words, correctly identifying the underlying latent state associated with an observation is crucial. Latent state inference thus plays an

important role in constructing cognitive graphs – especially during early learning, when the small amount of experience can lead to highly uncertain estimates of the state structure, which adaptive decision-makers must account for (Jiang et al., 2015; Harhen and Bornstein, 2023).

When observations are aliased with respect to the underlying latent states, inferring the generative structure requires interpreting each observation relative to the others; formally, conditioning inference on some subset of the history of observations, rather than just the one presently available sensory input (Whittington et al., 2022). Hidden Markov Models (HMM) provide a computational solution to latent state inference by decoupling the transition structure of the latent states (transition matrix) and the probability that a given observation maps onto latent states (emission matrix; George et al., 2021; Mark et al., 2020); this dichotomy is also dubbed as “stimulus-stimulus” associations and “stimulus-context latent” associations of content representations, respectively (see Wang, Feng, and Bornstein, 2022).

A clone-structured cognitive graph is a version of the HMM that conditions the transition of latent spaces on actions (George et al., 2021). To elaborate, a given observation is explained in terms of two components: a transition tensor which accounts for the action-conditioned transitions between latent states, and an emission matrix that assigns probabilities to the latent states given an observation. Within the transition tensor, each latent state in a sequence is identified in relation to its previous latent state and action, and whenever a new context – or a new combination of previous latent state and action – is encountered, a new clone is created. Clone-structured cognitive graphs have been able to capture phenomena thought to be important to structure learning in both spatial (George et al., 2021) and non-spatial (Swaminathan et al., 2023) domains: *splitting*, the ability to recover the ground-truth space from aliased observations, as well as *merging*, the ability to stitch overlapping latent states together from two disjoint observations. Thus, the clone-structured cognitive graph is an exciting proposal for how an agent can simultaneously learn both the structure (i.e., nodes) of the environment as well as its transition dynamics (i.e., edges). Within the framework we discuss here, the resulting representation is considered *uncompressed*, as it is attempting to capture the full, flexible environment model. Backed-up representations can be built by querying the resulting graph, as it stabilizes with sufficient experience (Wittkuhn, Krippner, and Schuck, 2022).

Another variant of the HMM-based cognitive graph explicitly assumes the idea of predefined schemas for identifying the transition structure. Here, it is postulated that the transition dynamics emerge from predefined structural forms such as hexagonal grids or community structures (Mark et al., 2020), which could be grounded in the wider notion of inherent basis sets (Kemp and Tenenbaum, 2008; Tenenbaum et al., 2011; Luettgau et al., 2023) or generative grammar of sequences (Dragoi, 2023). The idea that cognitive graphs are constructed using the prior knowledge of structures could be empirically supported by results that human transfer learning is best explained by these models (Mark et al., 2020; Luettgau et al., 2023). The hippocampal-entorhinal system has

been proposed to underlie decoupling, or factorizing, structure and sensory observations (Whittington et al., 2018). Here, the medial entorhinal cortex contains grid cells (Hafting et al., 2005) that provide a basis set along which transition structure is defined, and the lateral entorhinal cortex supports sensory representations. The conjunctive code of the transition structure and “emission” is hypothesized to be reflected in the hippocampus (Whittington et al., 2018). These distinct representational forms each play a critical role in the use of hippocampal replay to infer compositional structure across environments, permitting the construction of more compressed representations that can support efficient planning in novel environments (Kurth-Nelson et al., 2023). An area for future research is whether endowing artificial agents with this representational decomposition and algorithmic approach to replaying and recombining structure elements can allow them to perform an efficient approximation to graph compression.

Below, we describe how cognitive graphs support both merging and splitting in sequences of observations, and what specific mechanisms an unfolded graph could provide to facilitate early stage learning.

#### *4.1.1 Merging: fast generalization by extrapolating trajectories*

Associative memory could be seen as the building block of cognitive graphs. One such instantiation is transitive inference, which is an example of leveraging relational information of instances for faster generalization, observed in humans and animals (Bryant and Trabasso, 1971; Gillan, 1981; Davis, 1992). When an agent experiences  $A > B$  and  $B > C$ , the unobserved relationship between  $A > C$  can be inferred without direct experience (Eichenbaum et al., 1999). This can be achieved through forming supraordinate representations, comparable to cognitive graphs, such that  $A > B > C$ , which has been found to be supported by the hippocampus (Greene et al., 2006; Dusek and Eichenbaum, 1997; Zalesak and Heckers, 2009). Similarly, disparate fragments of event trajectories can be fused together, creating graph-like formations by leveraging the intersections of these trajectories (Eichenbaum and Cohen, 2014; Rmus et al., 2022). From these graphs, inferences can be made between instances that were not directly experienced together, supporting flexible recombination and fast generalization (Eichenbaum, 2004). After learning sequences of objects that are generated based on graphs that are either hexagonal or community-structured, humans are able to infer unobserved links using the transition structure of the latent graphs (Mark et al., 2020). This study provides direct evidence that people are able to extract long-run transition structure from sequences of events and also are able to transfer it for generalization.

Implementing this associative-memory-based cognitive graph leads to efficient planning algorithms. For example, an episodic reinforcement learning algorithm called Episodic Reinforcement Learning with Associative Memory (ERLAM) augmented with associative memory showed increased sample efficiency compared to benchmarks (Zhu et al., 2020). In ERLAM, experienced trajectories are reorganized into graphs, which speeds the propagation of value learned from one in-

stance to other related instances, thereby enhancing sample efficiency. In addition, clone-structured cognitive graphs introduced earlier have been shown to be capable of performing transitive inferences (George et al., 2021). This capability was demonstrated in the spatial domain, in agents navigating a larger environment divided into discrete “rooms.” Here, two separate rooms are stitched together to form an overlapping region. Agents navigate each room separately and are tested on whether they can travel from a non-overlapping region of one room to a region exclusive to the other room. Results show that agents are able to construct a latent map by stitching sequential observations from two disjoint episodes; overlapping observations from different trajectories are correctly assigned to the same hidden state.

In addition to conjoining separate sequences, associative memory binds seemingly independent choice options together into a temporal context, so that learning the value of a chosen option also influences the value of unchosen options (Biderman and Shohamy, 2021). This is referred to as counterfactual reasoning, another example of associative memory accelerating learning since information about an instance can be propagated to related experiences. Counterfactual reasoning is observed in reinforcement learning: humans not only deploy “factual” information through direct trial-and-error, but also incorporate counterfactual learning (Boorman, Behrens, and Rushworth, 2011; Fischer and Ullsperger, 2013). Interestingly, counterfactual learning engages cognitive graphs for *both* model-based and model-free learning (Moran, Dayan, and Dolan, 2021). In this process, the model-free values of options are positively reinforced by direct rewards and negatively influenced by the value of counterfactual options. Associative memory strength between options in reinforcement learning being correlated with how much learning about one option influences other unchosen options suggests that counterfactual learning operates on a cognitive graph where edge weights are defined by associative memory strength between items (Biderman, Gershman, and Shohamy, 2023). An open question is whether factual and counterfactual learning are performed on the same cognitive graph. Some evidence points to a single representation supporting both kinds of reasoning (Boorman et al., 2009; Fischer and Ullsperger, 2013), whereas other evidence supports these forms of learning update distinct representations (Lohrenz et al., 2007; Li and Daw, 2011; Kishida et al., 2016). A common finding is that individuals are generally biased towards reinforcing their own choices (“confirmation bias”, or the tendency to collect information partially according to the preexisting belief or action (Nickerson, 1998)), in a way that they incorporate more information when the chosen option is more rewarded (i.e., greater learning rate for positive prediction errors of factual options) and when the *unchosen* option turns out to be *less* rewarding (i.e., greater learning rate for negative prediction errors of counterfactual options; Palminteri et al., 2017). Asymmetric updating in the other direction (more negative than positive) has been observed in individuals diagnosed with psychiatric disorders (e.g. depression; Rouhani and Niv, 2019); though this pattern has itself been shown to arise from individual differences in representational precision (Harhen and Bornstein, 2024).



ERLAM provides an example of leveraging counterfactual combinatorial trajectories to facilitate learning of artificial agents (Zhu et al., 2020). In this algorithm, trajectories are reorganized into graphs by merging the common elements (nodes) of the two trajectories. This agent would have an advantage over an agent who uses pure episodic memory in cases such as right after experiencing two intersecting trajectories that each lead to reward (e.g.,  $A > B > C > \text{reward}$ ) and no-reward (e.g.,  $D > B > E > \text{no-reward}$ ); while the ERLAM agent will be able to leverage the graph to plan an unexperienced route (e.g.,  $D > B > C > \text{reward}$ ), an agent that only relies on episodic reinforcement learning would associate  $D$  with reward only after direct experience. Recently, *expected eligibility traces* have been introduced as a form of leveraging counterfactual trajectories to accelerate learning (Hasselt et al., 2021). Eligibility trace is a mechanism in reinforcement learning that provides a hindsight credit assignment with regard to the current state by keeping a trace of past experiences weighted by their recency (Singh and Sutton, 1996; Sutton and Barto, 2018). *Expected* eligibility traces improves the limitation of eligibility traces – that only one directly experienced trace is updated each time – by considering multiple counterfactual sequences that could have preceded a current state. Mirroring the relationship between the full forward model and successor representations, the predecessor representation is the fully backed-up version of the state tree supporting expected eligibility traces (Bailey and Mattar, 2022).

#### 4.1.2 Splitting: Recovering latent structure from aliased sequences

It is possible that two different states are “aliased”, or mapped onto overlapping observations. In this situation, as opposed to the example above that agents should be able to create a graph that *merges* two sequences –  $A > D > C$  and  $B > D > E$  –, an agent should be able to *split*  $D$  into two different nodes according to their contexts. Clone-structured cognitive graphs are able to accurately reconstruct correct latent graphs from sequences of aliased sensory observations by making clones of observations (George et al., 2021). Impressively, clone-structured cognitive graphs are not only able to both split aliased observation into latent states, but also able to merge the reconstructed graphs as in transitive inference.

Indeed, as implied above, these are exactly the sorts of environments in which clone-structured cognitive graphs have an advantage over backed-up representations. Specifically, when presenting a clone-structured cognitive graph agent with a sequence of aliased observations from a graph with community structure (e.g. Figure 1 left), the agent is able to recover the modular structure. However, an agent that used a successor representation was not able to recover this structure (George et al., 2021). This suggests that environments in which modular structure is important to the task at hand benefit from having available less-compressed representations of experience. This idea aligns with the finding that sequences of observations generated by a modular vs. lattice graph – where the two graphs only differ in terms of their higher-order structure – lead to more robust latent representations (Kahn et al., 2023).

## 5 Discussion

Multi-step planning is a critical ability for autonomous organisms. Extensive research has identified multiple kinds of planning, each with their own benefits and appropriate to specific situations. These distinct approaches rely on different representational substrates and have different algorithmic commitments. Which kind of planning an individual performs in a given setting can dramatically change the outcome of their decisions. Therefore, it can be valuable to judgment and decision-making research to understand how to characterize the commonalities, differences, and appropriate uses of each form of planning. Here, we suggest that these seemingly distinct representational forms that support planning can be described as varying types of *cognitive graph*, where these various manifestations of graphs exist along a spectrum of compression.

At one end of the spectrum, fully uncompressed representations capture every element of an associative network in detail. This full-featured model of the environment allows for flexible trajectory sampling at the time of decision, and supports plans that are robust to changes in contingency and reward structure. In addition, this kind of uncompressed representation is a necessary first step for building more compressed representations, because the latent structure (edges) in compressed representations requires inferring across multiple experiences (Wittkuhn, Krippner, and Schuck, 2022; Lynn and Bassett, 2020). Since observations are often aliased, uniquely characterizing their latent state (nodes) and structure (edges) requires them to be placed in a sequence (Whittington et al., 2022). Network model simulations show that uncompressed sequences of events are necessary for building latent graphs that enable complex functions such as rapid value propagation (Zhu et al., 2020) or extracting higher-order structures (George et al., 2021).

At the other extreme, a fully compressed graph – such as the *successor* and *predecessor* representations – captures summary-level statistical structure. These graphs are formed by “bootstrapping” – repeated sampling of the full model to identify the long-run relationships between each pair of nodes in the network. These representations allow for fast, cheap multi-step planning as they cache previous trajectory samples into a compact matrix format. Their *factorized* form, separating transition (edge) information from reward values, allows for replanning in the face of changing reward outcomes. However, the kind of planning they support is not robust to changes in contingency structure – these must be re-learned, slowing planning until stable estimates can be obtained again. This is because backed-up representations like the successor representation are conditional on the specific policy that generated the compressed graphs; in other words, if the goal changes, the optimal action in each state should be re-learned, thereby not transferable (Lehnert, Tellex, and Littman, 2017). Linear reinforcement learning, which incentivizes learning a “default” policy distributed uniformly across possible successor states, is a framework that addresses this limitation and explains flexible replanning in humans (Piray and Daw, 2021).

These different kinds of representations are learned simultaneously, which allows the agent to arbitrate between the most reliable representation at a given moment (Wang, Feng, and Bornstein,

2022). In situations of high uncertainty, such as during the early phase of reinforcement learning where there is not enough data to construct a reliable model (Lengyel and Dayan, 2007) or in volatile environments (Nicholas, Daw, and Shohamy, 2022), consulting on a subset of episodic samples provides a more reliable approximation of the value of observations. Arbitration between different representations has been proposed to be reflected in discontinuous “jumps” of subjective evidence (“jump-diffusions”) in evidence accumulation models, where these sudden jumps during the sequential evidence sampling could indicate alternations to other sets of representations (Wang, Feng, and Bornstein, 2022). The constantly changing ensemble of representations that lead to these jumps is hypothesized to occur in a bottom-up manner, akin to product-of-experts in machine learning (Wang, Feng, and Bornstein, 2022)

By which mechanism does compression happen, such that more experience gradually leads to more compression? One possible mechanism could be the diffusion of information between nodes through replay of events. The transition from the uncompressed graph to the fully backed-up form occurs via repeatedly sampling and aggregating features from neighboring nodes, analogous to message passing algorithms (Hamilton, Ying, and Leskovec, 2017; Parr et al., 2019; George et al., 2021). At the beginning of the learning process, the cognitive graph resembles an undiffused graph where a node, or a given state, holds limited information about others, thus requiring the agent in a state to explicitly traverse edges to infer about other states. At the same time, uncompressed graphs provide full representations of the contingency structure between states and actions, which allow for flexibility at the cost of greater computation time and behavioral variability. With more experience, the cognitive graphs undergo a transformation into a bootstrapped representation where information about future states is aggregated into each adjacent state, making explicit edge-based inferences between states less important. Caching these distal outcomes subserves rapid planning, while still retaining sensitivity to changes in reward availability. However, without additional mechanisms, it also confers a relative insensitivity to contingency changes that may be undesirable in novel or volatile environments. Replay of events could be a biological instantiation of message passing, given that the construction of backed-up representations is mediated by on-task replay in humans (Wittkuhn, Krippner, and Schuck, 2022). A possible future direction for research would be to investigate whether replay contributes to maintaining and arbitrating between multiple kind of representations.

One interesting direction to expand this concept of representational spectrum would be to test whether different modes of control arise as a function of the degree of compression (Moskovitz et al., 2022). Recall that agents using compressed representations should be adept at open-loop control because they can in principle select action sequences into a single operation. This eliminates the need for intermittent re-planning during action execution (Eysenbach, Salakhutdinov, and Levine, 2021). However, if agents plan by sequentially sampling next actions using uncompressed cognitive graphs, taking small steps could be more efficient than open-loop control, since

the model has not been compiled yet to provide reliable future trajectories from a state. This edge-based planning is conceptually more similar to closed-loop planning, where an agent stops at each transition to re-plan. This leads us to the overarching question of whether the utility of using closed-loop vs. open-loop planning aligns with the degree of compression in cognitive graphs. This alignment would be similar to the evolution of an episodic control system to model-based, and then finally to model-free systems (Lengyel and Dayan, 2007). Based on an interpretation that the seemingly model-free behaviors could actually be action sequences (Dezfouli and Balleine, 2012, 2013), an interesting hypothesis is that the model-free system at the end of the spectrum could be in fact representing action sequences formed by open-loop control, likely a result of using highly compressed models.

## 6 Conclusions

To conclude, we highlight the kind of representations that could be used to support instance-based planning at early stages of learning – uncompressed cognitive graphs – and suggest that they could be in a spectrum, rather than discrete concepts, with backed-up representations at the other end. Further research may investigate whether this of spectrum of *representations* directly induces a continuum of planning *algorithms*, such as closed- versus open-loop control.

## 7 Acknowledgements

The authors would like to thank Dr. Bruce McNaughton and Dr. Jeffrey L. Krichmar for insightful discussions. Funding was provided by NIA R21AG072673 (to AMB), and NINDS R01NS119468 (to ERC).

## 8 Statements and Declarations

The authors declare no competing interest.

## References

- Bagaria, Akhil, Jason K Senthil, and George Konidaris (2021). “Skill discovery for exploration and planning using deep skill graphs”. In: *International Conference on Machine Learning*. PMLR, pp. 521–531.
- Bailey, Duncan and Marcelo Mattar (2022). “Predecessor Features”. In: *arXiv preprint arXiv:2206.00303*.
- Barreto, André et al. (2017). “Successor features for transfer in reinforcement learning”. In: *Advances in neural information processing systems* 30.
- Barron, Greg and Ido Erev (2003). “Small feedback-based decisions and their limited correspondence to description-based decisions”. In: *Journal of behavioral decision making* 16.3, pp. 215–233.

- Barron, Helen C, Raymond J Dolan, and Timothy EJ Behrens (2013). “Online evaluation of novel choices by simultaneous representation of multiple memories”. In: *Nature neuroscience* 16.10, pp. 1492–1498.
- Bertsekas, Dimitri (2012). *Dynamic programming and optimal control: Volume I*. Vol. 4. Athena scientific.
- Bhatia, Sudeep (2014). “Sequential sampling and paradoxes of risky choice”. In: *Psychonomic Bulletin & Review* 21.5, pp. 1095–1111.
- Bhui, Rahul (2018). “Case-based decision neuroscience: Economic judgment by similarity”. In: *Goal-directed decision making*. Elsevier, pp. 67–103.
- Biderman, Natalie, Samuel J Gershman, and Daphna Shohamy (2023). “The role of memory in counterfactual valuation.” In: *Journal of Experimental Psychology: General* 152.6, p. 1754.
- Biderman, Natalie and Daphna Shohamy (2021). “Memory and decision making interact to shape the value of unchosen options”. In: *Nature communications* 12.1, p. 4648.
- Blundell, Charles et al. (2016). “Model-free episodic control”. In: *arXiv preprint arXiv:1606.04460*.
- Boorman, Erie D, Timothy E Behrens, and Matthew F Rushworth (2011). “Counterfactual choice and learning in a neural network centered on human lateral frontopolar cortex”. In: *PLoS biology* 9.6, e1001093.
- Boorman, Erie D et al. (2009). “How green is the grass on the other side? Frontopolar cortex and the evidence in favor of alternative courses of action”. In: *Neuron* 62.5, pp. 733–743.
- Bornstein, Aaron M and Nathaniel D Daw (2012). “Dissociating hippocampal and striatal contributions to sequential prediction learning”. In: *European Journal of Neuroscience* 35.7, pp. 1011–1023.
- (2013). “Cortical and hippocampal correlates of deliberation during model-based decisions for rewards in humans”. In: *PLoS computational biology* 9.12, e1003387.
- Bornstein, Aaron M and Kenneth A Norman (2017). “Reinstated episodic context guides sampling-based decisions for reward”. In: *Nature neuroscience* 20.7, pp. 997–1003.
- Bornstein, Aaron M et al. (2017). “Reminders of past choices bias decisions for reward in humans”. In: *Nature Communications* 8.1, pp. 1–9.
- Bryant, Peter E and Thomas Trabasso (1971). “Transitive inferences and memory in young children”. In: *Nature* 232, pp. 456–458.
- Butts, Carter T (2009). “Revisiting the foundations of network analysis”. In: *science* 325.5939, pp. 414–416.
- Chrastil, Elizabeth R and William H Warren (2014). “From cognitive maps to cognitive graphs”. In: *PloS one* 9.11, e112544.
- Davis, Hank (1992). “Transitive inference in rats (*Rattus norvegicus*).” In: *Journal of Comparative Psychology* 106.4, p. 342.
- Daw, Nathaniel D, Yael Niv, and Peter Dayan (2005). “Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control”. In: *Nature neuroscience* 8.12, pp. 1704–1711.
- Dayan, Peter (1993). “Improving generalization for temporal difference learning: The successor representation”. In: *Neural computation* 5.4, pp. 613–624.
- (2008). “The role of value systems in decision making.” In.
- Dezfouli, Amir and Bernard W Balleine (2012). “Habits, action sequences and reinforcement learning”. In: *European Journal of Neuroscience* 35.7, pp. 1036–1051.

- Dezfouli, Amir and Bernard W Balleine (2013). “Actions, action sequences and habits: evidence that goal-directed and habitual action control are hierarchically organized”. In: *PLoS computational biology* 9.12, e1003364.
- Doya, Kenji et al. (2002). “Multiple model-based reinforcement learning”. In: *Neural computation* 14.6, pp. 1347–1369.
- Dragoi, George (2023). “The generative grammar of the brain: a critique of internally generated representations”. In: *Nature Reviews Neuroscience*, pp. 1–16.
- Dusek, Jeffery A and Howard Eichenbaum (1997). “The hippocampus and memory for orderly stimulus relations”. In: *Proceedings of the National Academy of Sciences* 94.13, pp. 7109–7114.
- Eichenbaum, Howard (2004). “Hippocampus: cognitive processes and neural representations that underlie declarative memory”. In: *Neuron* 44.1, pp. 109–120.
- Eichenbaum, Howard and Neal J Cohen (2014). “Can we reconcile the declarative memory and spatial navigation views on hippocampal function?” In: *Neuron* 83.4, pp. 764–770.
- Eichenbaum, Howard et al. (1999). “The hippocampus, memory, and place cells: is it spatial memory or a memory space?” In: *Neuron* 23.2, pp. 209–226.
- Eysenbach, Ben, Russ R Salakhutdinov, and Sergey Levine (2021). “Robust predictable control”. In: *Advances in Neural Information Processing Systems* 34, pp. 27813–27825.
- Feld, GB et al. (2022). “Sleep targets highly connected global and local nodes to aid consolidation of learned graph networks”. In: *Scientific Reports* 12.1, p. 15086.
- Fiedler, Klaus (2000). “Beware of samples! A cognitive-ecological sampling approach to judgment biases.” In: *Psychological review* 107.4, p. 659.
- Fischer, Adrian G and Markus Ullsperger (2013). “Real and fictive outcomes are processed differently but converge on a common adaptive mechanism”. In: *Neuron* 79.6, pp. 1243–1255.
- Fox, Craig R and Liat Hadar (2006). ““Decisions from experience”= sampling error+ prospect theory: Reconsidering Hertwig, Barron, Weber & Erev (2004)”. In: *Judgment and Decision Making* 1.2, pp. 159–161.
- Fox, Roy et al. (2017). “Multi-level discovery of deep options”. In: *arXiv preprint arXiv:1703.08294*.
- Geerts, Jesse P et al. (2022). “A probabilistic successor representation for context-dependent prediction”. In: *bioRxiv*.
- George, Dileep et al. (2021). “Clone-structured graph representations enable flexible learning and vicarious evaluation of cognitive maps”. In: *Nature communications* 12.1, p. 2392.
- Gershman, Samuel J (2018). “The successor representation: its computational logic and neural substrates”. In: *Journal of Neuroscience* 38.33, pp. 7193–7200.
- Gershman, Samuel J and Nathaniel D Daw (2017). “Reinforcement learning and episodic memory in humans and animals: an integrative framework”. In: *Annual review of psychology* 68, p. 101.
- Gershman, Samuel J et al. (2012). “The successor representation and temporal context”. In: *Neural Computation* 24.6, pp. 1553–1568.
- Gilboa, Itzhak and David Schmeidler (1995). “Case-based decision theory”. In: *The quarterly Journal of economics* 110.3, pp. 605–639.
- Gillan, Douglas J (1981). “Reasoning in the chimpanzee: II. Transitive inference.” In: *Journal of Experimental Psychology: Animal Behavior Processes* 7.2, p. 150.
- Gopnik, Alison and Laura Schulz (2004). “Mechanisms of theory formation in young children”. In: *Trends in cognitive sciences* 8.8, pp. 371–377.



- Gopnik, Alison et al. (2004). “A theory of causal learning in children: causal maps and Bayes nets.” In: *Psychological review* 111.1, p. 3.
- Greene, Anthony J et al. (2006). “An fMRI analysis of the human hippocampus: inference, context, and task awareness”. In: *Journal of cognitive neuroscience* 18.7, pp. 1156–1173.
- Hafting, Torkel et al. (2005). “Microstructure of a spatial map in the entorhinal cortex”. In: *Nature* 436.7052, pp. 801–806.
- Hamilton, Will, Zhitao Ying, and Jure Leskovec (2017). “Inductive representation learning on large graphs”. In: *Advances in neural information processing systems* 30.
- Harhen, Nora C and Aaron M Bornstein (2023). “Overharvesting in human patch foraging reflects rational structure learning and adaptive planning”. In: *Proceedings of the National Academy of Sciences* 120.13, e2216524120.
- (2024). “Interval timing as a computational pathway from early life adversity to affective disorders”. In: *Topics in Cognitive Science* 16.1, pp. 92–112.
- Hasselt, Hado van et al. (2021). “Expected eligibility traces”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 35. 11, pp. 9997–10005.
- Hertwig, Ralph and Timothy J Pleskac (2010). “Decisions from experience: Why small samples?” In: *Cognition* 115.2, pp. 225–237.
- Hertwig, Ralph et al. (2004). “Decisions from experience and the effect of rare events in risky choice”. In: *Psychological science* 15.8, pp. 534–539.
- Ho, Mark K et al. (2022). “People construct simplified mental representations to plan”. In: *Nature* 606.7912, pp. 129–136.
- Hotaling, Jared M (2020). “Decision field theory-planning: A cognitive model of planning on the fly in multistage decision making.” In: *Decision* 7.1, p. 20.
- Hotaling, Jared M and David Kellen (2022). “Dynamic decision making: Empirical and theoretical directions”. In: *Psychology of Learning and Motivation*. Vol. 76. Elsevier, pp. 207–238.
- Howard, Marc W and Michael J Kahana (2002). “A distributed representation of temporal context”. In: *Journal of mathematical psychology* 46.3, pp. 269–299.
- Hunt, LT et al. (2021). “Formalizing planning and information search in naturalistic decision-making”. In: *Nature neuroscience* 24.8, pp. 1051–1064.
- Jeong, Huijeong et al. (2022). “Mesolimbic dopamine release conveys causal associations”. In: *Science* 378.6626, eabq6740.
- Jiang, Nan et al. (2015). “The dependence of effective planning horizon on model accuracy”. In: *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pp. 1181–1189.
- Kahn, Ari E et al. (2018). “Network constraints on learnability of probabilistic motor sequences”. In: *Nature human behaviour* 2.12, pp. 936–947.
- Kahn, Ari E et al. (2023). “Network structure influences the strength of learned neural representations”. In: *bioRxiv*, pp. 2023–01.
- Kahneman, Daniel and Amos Tversky (1979). “Prospect Theory: An Analysis of Decision under Risk”. In: *Econometrica* 47.2, pp. 263–291. ISSN: 00129682, 14680262.
- Kemp, Charles and Joshua B Tenenbaum (2008). “The discovery of structural form”. In: *Proceedings of the National Academy of Sciences* 105.31, pp. 10687–10692.
- Kinney, David and Tania Lombrozo (2023a). “Building Compressed Causal Models of the World”. In.

- Kinney, David and Tania Lombrozo (2023b). “Lossy Compression and the Granularity of Causal Representation”. In: *NeurIPS 2023 workshop: Information-Theoretic Principles in Cognitive Systems*.
- Kishida, Kenneth T et al. (2016). “Subsecond dopamine fluctuations in human striatum encode superposed error signals about actual and counterfactual reward”. In: *Proceedings of the National Academy of Sciences* 113.1, pp. 200–205.
- Korf, Richard E (1987). “Planning as search: A quantitative approach”. In: *Artificial intelligence* 33.1, pp. 65–88.
- Kurth-Nelson, Zeb et al. (2023). “Replay and compositional computation”. In: *Neuron* 111.4, pp. 454–469.
- Lai, Lucy and Samuel J Gershman (2021). “Policy compression: An information bottleneck in action selection”. In: *Psychology of Learning and Motivation*. Vol. 74. Elsevier, pp. 195–232.
- Lai, Lucy, Ann Zixiang Huang, and Samuel J Gershman (2022). “Action chunking as policy compression”. In: .
- Lehnert, Lucas, Stefanie Tellex, and Michael L Littman (2017). “Advantages and limitations of using successor features for transfer in reinforcement learning”. In: *arXiv preprint arXiv:1708.00102*.
- Lengyel, Máté and Peter Dayan (2007). “Hippocampal contributions to control: the third way”. In: *Advances in neural information processing systems* 20, pp. 889–896.
- Li, Jian and Nathaniel D Daw (2011). “Signals in human striatum are appropriate for policy update rather than value prediction”. In: *Journal of Neuroscience* 31.14, pp. 5504–5511.
- Liu, Kara et al. (2020). “Hallucinative topological memory for zero-shot visual planning”. In: *International Conference on Machine Learning*. PMLR, pp. 6259–6270.
- Lohrenz, Terry et al. (2007). “Neural signature of fictive learning signals in a sequential investment task”. In: *Proceedings of the National Academy of Sciences* 104.22, pp. 9493–9498.
- Luettgau, Lennart et al. (2023). “Decomposing dynamical subprocesses for compositional generalization”. In: .
- Lynn, Christopher W and Danielle S Bassett (2020). “How humans learn and represent networks”. In: *Proceedings of the National Academy of Sciences* 117.47, pp. 29407–29415.
- Lynn, Christopher W et al. (2020). “Abstract representations of events arise from mental errors in learning and memory”. In: *Nature communications* 11.1, p. 2313.
- Mark, Shirley et al. (2020). “Transferring structural knowledge across cognitive maps in humans and models”. In: *Nature communications* 11.1, p. 4783.
- Momennejad, Ida et al. (2017). “The successor representation in human reinforcement learning”. In: *Nature human behaviour* 1.9, pp. 680–692.
- Moran, Rani, Peter Dayan, and Raymond J Dolan (2021). “Human subjects exploit a cognitive map for credit assignment”. In: *Proceedings of the National Academy of Sciences* 118.4, e2016884118.
- Moskovitz, Ted, Spencer R Wilson, and Maneesh Sahani (2021). “A First-Occupancy representation for reinforcement learning”. In: *arXiv preprint arXiv:2109.13863*.
- Moskovitz, Ted et al. (2022). “A unified theory of dual-process control”. In: *arXiv preprint arXiv:2211.07036*.
- Moskovitz, Ted et al. (2023). “A State Representation for Diminishing Rewards”. In: *arXiv preprint arXiv:2309.03710*.
- Moss, Laurence S (1984). “The Laws of Human Relations and the Rules of Human Action Derived Therefrom. By Hermann Heinrich Gossen. Translated by Rudolph C. Blitz with an introductory

- essay by Nicholas Georgescu-Roegen. Cambridge: MIT Press, 1983. Pp. 460. 47.50.”. In: *The Journal of Economic History* 44.4, pp. 1130–1132.
- Muller, Robert U, Matt Stead, and Janos Pach (1996). “The hippocampus as a cognitive graph.” In: *The Journal of general physiology* 107.6, pp. 663–694.
- Natarajan, Mausam and Andrey Kolobov (2022). *Planning with Markov decision processes: An AI perspective*. Springer Nature.
- Nicholas, Jonathan, Nathaniel D Daw, and Daphna Shohamy (2022). “Uncertainty alters the balance between incremental learning and episodic memory”. In: *Elife* 11, e81679.
- Nickerson, Raymond S (1998). “Confirmation bias: A ubiquitous phenomenon in many guises”. In: *Review of general psychology* 2.2, pp. 175–220.
- Niv, Yael (2019). “Learning task-state representations”. In: *Nature neuroscience* 22.10, pp. 1544–1553.
- Noh, Sharon Mina et al. (2023). “Memory precision and age differentially predict the use of decision-making strategies across the lifespan.” In.
- Otto, A Ross et al. (2022). “Context-dependent choice and evaluation in real-world consumer behavior”. In: *Scientific reports* 12.1, p. 17744.
- Palminteri, Stefano et al. (2015). “Contextual modulation of value signals in reward and punishment learning”. In: *Nature communications* 6.1, p. 8096.
- Palminteri, Stefano et al. (2017). “Confirmation bias in human reinforcement learning: Evidence from counterfactual feedback processing”. In: *PLoS computational biology* 13.8, e1005684.
- Parr, Thomas et al. (2019). “Neuronal message passing using Mean-field, Bethe, and Marginal approximations”. In: *Scientific reports* 9.1, p. 1889.
- Peer, Michael et al. (2021). “Structuring knowledge with cognitive maps and cognitive graphs”. In: *Trends in cognitive sciences* 25.1, pp. 37–54.
- Piray, Payam and Nathaniel D Daw (2021). “Linear reinforcement learning in planning, grid fields, and cognitive control”. In: *Nature communications* 12.1, p. 4942.
- Plonsky, Ori, Kinneret Teodorescu, and Ido Erev (2015). “Reliance on small samples, the wavy recency effect, and similarity-based learning.” In: *Psychological review* 122.4, p. 621.
- Pritzel, Alexander et al. (2017). “Neural episodic control”. In: *International Conference on Machine Learning*. PMLR, pp. 2827–2836.
- Rangel, Antonio, Colin Camerer, and P Read Montague (2008). “A framework for studying the neurobiology of value-based decision making”. In: *Nature reviews neuroscience* 9.7, pp. 545–556.
- Rhodes, Theo and Michael T Turvey (2007). “Human memory retrieval as Lévy foraging”. In: *Physica A: Statistical Mechanics and its Applications* 385.1, pp. 255–260.
- Rmus, Milena et al. (2022). “Humans can navigate complex graph structures acquired during latent learning”. In: *Cognition* 225, p. 105103.
- Rouhani, Nina and Yael Niv (2019). “Depressive symptoms bias the prediction-error enhancement of memory towards negative events in reinforcement learning”. In: *Psychopharmacology* 236, pp. 2425–2435.
- Savinov, Nikolay, Alexey Dosovitskiy, and Vladlen Koltun (2018). “Semi-parametric topological memory for navigation”. In: *arXiv preprint arXiv:1803.00653*.
- Schapiro, Anna C et al. (2013). “Neural representations of events arise from temporal community structure”. In: *Nature neuroscience* 16.4, pp. 486–492.

- Schrittwieser, Julian et al. (2020). “Mastering atari, go, chess and shogi by planning with a learned model”. In: *Nature* 588.7839, pp. 604–609.
- Silver, David and Joel Veness (2010). “Monte-Carlo planning in large POMDPs”. In: *Advances in neural information processing systems* 23.
- Singh, Satinder P and Richard S Sutton (1996). “Reinforcement learning with replacing eligibility traces”. In: *Machine learning* 22.1-3, pp. 123–158.
- Smith, Kyle S and Ann M Graybiel (2013). “A dual operator view of habitual behavior reflecting cortical and striatal dynamics”. In: *Neuron* 79.2, pp. 361–374.
- Solomon, Ethan A et al. (2019). “Hippocampal theta codes for distances in semantic and temporal spaces”. In: *Proceedings of the National Academy of Sciences* 116.48, pp. 24343–24352.
- Solway, Alec et al. (2014). “Optimal behavioral hierarchy”. In: *PLoS computational biology* 10.8, e1003779.
- Sommerville, Jessica A and Amanda L Woodward (2005a). “Infants’ sensitivity to the causal features of means-end support sequences in action and perception”. In: *Infancy* 8.2, pp. 119–145.
- (2005b). “Pulling out the intentional structure of action: the relation between action processing and action production in infancy”. In: *Cognition* 95.1, pp. 1–30.
- Sriram, Sivaraman et al. (2009). “Implementing a no-loss state in the game of tic-tac-toe using a customized decision tree algorithm”. In: *2009 International Conference on Information and Automation*. IEEE, pp. 1211–1216.
- Stachenfeld, Kimberly L, Matthew M Botvinick, and Samuel J Gershman (2017). “The hippocampus as a predictive map”. In: *Nature neuroscience* 20.11, pp. 1643–1653.
- Stiso, Jennifer et al. (2022). “Neurophysiological evidence for cognitive map formation during sequence learning”. In: *Eneuro* 9.2.
- Sutton, Richard S and Andrew G Barto (2018). *Reinforcement learning: An introduction*. MIT press.
- Sutton, Richard S, Doina Precup, and Satinder Singh (1999). “Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning”. In: *Artificial intelligence* 112.1-2, pp. 181–211.
- Swaminathan, Sivaramakrishnan et al. (2023). “Schema-learning and rebinding as mechanisms of in-context learning and emergence”. In: *arXiv preprint arXiv:2307.01201*.
- Talvitie, Erik (2017). “Self-correcting models for model-based reinforcement learning”. In: *Thirty-First AAAI Conference on Artificial Intelligence*.
- Tambini, Arielle et al. (2023). “Structured memory representations develop at multiple time scales in hippocampal-cortical networks”. In: *bioRxiv*, pp. 2023–04.
- Tenenbaum, Joshua B et al. (2011). “How to grow a mind: Statistics, structure, and abstraction”. In: *science* 331.6022, pp. 1279–1285.
- Tomov, Momchil S et al. (2020). “Discovery of hierarchical representations for efficient planning”. In: *PLoS computational biology* 16.4, e1007594.
- Wang, Shaoming, Samuel F Feng, and Aaron M Bornstein (2022). “Mixing memory and desire: How memory reactivation supports deliberative decision-making”. In: *Wiley Interdisciplinary Reviews: Cognitive Science* 13.2, e1581.
- Warren, William H (2019). “Non-euclidean navigation”. In: *Journal of Experimental Biology* 222.Suppl.1, jeb187971.

- Weber, Elke U and Eric J Johnson (2006). “Constructing preferences from memory”. In: *The Construction of Preference*, Lichtenstein, S. & Slovic, P.,(eds.), pp. 397–410.
- Whitehead, Steven D and Dana H Ballard (1991). “Learning to perceive and act by trial and error”. In: *Machine Learning* 7, pp. 45–83.
- Whittington, James et al. (2018). “Generalisation of structural knowledge in the hippocampal-entorhinal system”. In: *Advances in neural information processing systems* 31.
- Whittington, James CR et al. (2020). “The Tolman-Eichenbaum machine: unifying space and relational memory through generalization in the hippocampal formation”. In: *Cell* 183.5, pp. 1249–1263.
- Whittington, James CR et al. (2022). “How to build a cognitive map: insights from models of the hippocampal formation”. In: *arXiv preprint arXiv:2202.01682*.
- Wimmer, G Elliott and Daphna Shohamy (2012). “Preference by association: how memory mechanisms in the hippocampus bias decisions”. In: *Science* 338.6104, pp. 270–273.
- Wittkuhn, Lennart, Lena M Krippner, and Nicolas W Schuck (2022). “Statistical learning of successor representations is related to on-task replay”. In: *bioRxiv*, pp. 2022–02.
- Wu, Charley M, Eric Schulz, and Samuel J Gershman (2019). “Generalization as diffusion: human function learning on graphs”. In: *BioRxiv*, p. 538934.
- (2021). “Inference and search on graph-structured spaces”. In: *Computational Brain & Behavior* 4.2, pp. 125–147.
- Zalesak, Martin and Stephan Heckers (2009). “The role of the hippocampus in transitive inference”. In: *Psychiatry Research: Neuroimaging* 172.1, pp. 24–30.
- Zhang, Lunjun, Ge Yang, and Bradley C Stadie (2021). “World model as a graph: Learning latent landmarks for planning”. In: *International Conference on Machine Learning*. PMLR, pp. 12611–12620.
- Zhao, Wenjia Joyce, Russell Richie, and Sudeep Bhatia (2022). “Process and content in decisions from memory.” In: *Psychological Review* 129.1, p. 73.
- Zhu, Guangxiang et al. (2020). “Episodic reinforcement learning with associative memory”. In.
- Zhu, Jianqiao, Adam Sanborn, and Nick Chater (2018). “Mental sampling in multimodal representations”. In: *Advances in neural information processing systems* 31.