

# A mechanistic theory of planning in prefrontal cortex

Kristopher T. Jensen<sup>@1,3</sup>, Peter Doohan<sup>2</sup>, Mathias Sablé-Meyer<sup>1</sup>, Sandra Reinert<sup>1</sup>, Alon Baram<sup>3</sup>, Thomas Akam<sup>1,2</sup>, and Timothy E. J. Behrens<sup>1,3</sup>

<sup>1</sup> Sainsbury Wellcome Centre, University College London, London, UK

<sup>2</sup> Department of Experimental Psychology, University of Oxford, Oxford, UK

<sup>3</sup> Oxford Centre for Integrative Neuroimaging, University of Oxford, Oxford, UK

<sup>@</sup> Corresponding author (Kris.Torp.Jensen@gmail.com)

## Abstract

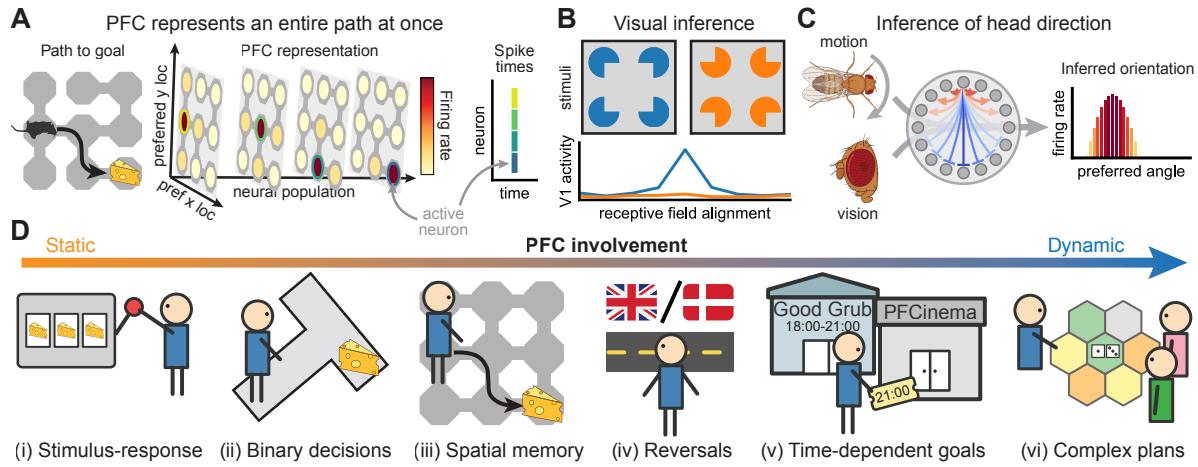
Planning is critical for adaptive behaviour in a changing world, because it lets us anticipate the future and adjust our actions accordingly. While prefrontal cortex is crucial for this process, it remains unknown how planning is implemented in neural circuits. Prefrontal representations were recently discovered in simpler sequence memory tasks, where different populations of neurons represent different future time points. We demonstrate that combining such representations with the ubiquitous principle of neural attractor dynamics allows circuits to solve much richer problems including planning. This is achieved by embedding the environment structure directly in synaptic connections to implement an attractor network that infers desirable futures. The resulting ‘spacetime attractor’ excels at planning in challenging tasks known to depend on prefrontal cortex. Recurrent neural networks trained by gradient descent on such tasks learn a solution that precisely recapitulates the spacetime attractor – in representation, in dynamics, and in connectivity. Analyses of networks trained across different environment structures reveal a generalisation mechanism that rapidly reconfigures the world model used for planning, without the need for synaptic plasticity. The spacetime attractor is a testable mechanistic theory of planning. If true, it would provide a path towards detailed mechanistic understanding of how prefrontal cortex structures adaptive behaviour.

## 1 Introduction

2 While different cortical areas support different functions, common computational principles are shared across many areas. For functions as disparate as sensory processing, spatial reasoning, and language comprehension, features of the environment are inferred from partial information. To do so, structural knowledge about the world must be embedded in synaptic connections (Ko et al., 2011; Burak and Fiete, 2009; Iacaruso et al., 2017; Turner-Evans et al., 2020). This constrains neural circuits to represent meaningful interpretations of the environment, and inputs select between these interpretations. It is not known whether similar principles generalise to complex prospective behaviours, such as planning extended action sequences to achieve a distant goal. Recent recordings from prefrontal cortex give an important clue. When mice or monkeys must execute a sequence of actions, neurons represent the entire sequence concurrently (El-Gaby et al., 2023; Xie et al., 2022). Different neuronal populations encode different steps of the future behaviour. If such a representation could be inferred from inputs indicating goals and constraints, the network could plan the future. Excitingly, this solution would use algorithmic principles similar to those known to infer features of the present in other cortical areas.

28 This paper has four overlapping aims. (i) To develop a detailed circuit model that infers explicit  
29 representations of the future from partial inputs.  
30 (ii) To discern the principles of synaptic connectivity  
31 that enable such a model to solve complex problems  
32 including planning. (iii) To understand when and  
33 why this algorithm succeeds while simpler circuit  
34 models fail. (iv) To explore how it relates to pre-  
35 frontal representations, connectivity, and function.  
36

37 **Sequence representations in PFC** Recent  
38 work has uncovered PFC representations underlying  
39 sequence working memory (Xie et al., 2022; El-Gaby  
40 et al., 2023; Whittington et al., 2023; Botvinick  
41 and Plaut, 2006). Separate neural ‘subspaces’, or  
42 groups of neurons, represent the expected or de-  
43 sired state of the world at different steps along the  
44 sequence of future behaviour. In other words, some  
45 neural subspace ‘A’ represents the present, while  
46 subspace ‘B’ represents the immediate future, and  
47 subspace ‘C’ the more distant future (Figure 1A).  
48 Critically, these subspaces are active simultaneously.  
49 Together, they instantaneously represent the entire  
50 behavioural sequence. Importantly, El-Gaby et al.  
51 (2023) showed that the PFC subspaces representing  
52 different steps of the future are not independent.  
53 Neuronal correlations reflect the structure of the  
54 task being performed, even during rest.



**Figure 1: Background.** **(A)** Recent work has identified explicit sequence representations in prefrontal cortex during working memory (Xie et al., 2022; El-Gaby et al., 2023). When an animal has to execute a behavioural sequence (left), individual neurons represent conjunctions of location and sequence element (centre). Separate populations (planes) therefore represent the expected location at different times in the future. The entire sequence is represented concurrently by the simultaneous firing of different neurons that encode the expected location at each time in the future (right). **(B)** An example V1 cell fires when its receptive field is aligned with an inferred line (blue), but not for a control stimulus with no inference (orange). Such visual inference is mediated by structural priors embedded in the circuit connectivity, where neurons representing consistent visual features excite each other (Iacaruso et al., 2017; Shin et al., 2023). Figure adapted from Lee and Nguyen (2001). **(C)** Structural knowledge is embedded in the synapses of the head direction system (centre; Turner-Evans et al., 2020), which constrains the network to represent single angles (Kim et al., 2017). Visual and proprioceptive inputs (left) determine which angle should be represented (right). We suggest that a mechanism like (B) and (C) infers the representation in (A). **(D)** Prefrontal cortex is particularly important in dynamic environments. Stimulus-response associations and repeated choices are robust to prefrontal lesions (i; ii), and acortical mice can solve spatial memory tasks (iii; Zheng et al., 2024). However, PFC is needed for reversal learning (iv; Walton et al., 2010) and when different goals are important, or ‘rewarding’, at different times (v; Shallice and Burgess, 1991). In multiplayer board games (vi), different resources are valuable at different stages, and opponents can dynamically change their strategy.

55 **Planning by inferring the future** The correlation structure observed in PFC resembles other  
 56 attractor circuits known to infer the current state of  
 57 the world, such as the instantaneous visual input or  
 58 head direction of an animal (Figure 1B-C; Ko et al.,  
 59 2011; Chaudhuri et al., 2019). If a similar inference  
 60 process acted on the future representation in PFC, it would extend its function beyond sequence  
 61 memory to situations where entire action sequences  
 62 are inferred from partial cues. Planning could be  
 63 solved by inferring sequences of desired actions from  
 64 a set of goals in a network implementation of ‘planning  
 65 as inference’ (Botvinick and Toussaint, 2012;  
 66 Levine, 2018). The resulting algorithm would nat-  
 67 urally cope with dynamic environments, because it  
 68 represents each time in the future separately. This  
 69 is intriguing because PFC is particularly important  
 70 for problems that require the correct behaviour to  
 71 be expressed at an appropriate time (Figure 1D;  
 72 Shallice and Burgess, 1991; Volle et al., 2011). Plan-  
 73 ning via attractor dynamics is also consistent with

76 winner-take-all dynamics identified in frontal cortex  
 77 during non-sequential behaviours (Ruff et al., 2025;  
 78 Inagaki et al., 2019). Such planning as inference  
 79 differs from most planning algorithms studied in  
 80 cognitive science and machine learning, which often  
 81 rely on sequential search (Callaway et al., 2022;  
 82 Schrittwieser et al., 2020). Search is easily adapted  
 83 to new environments, but it is slow at decision time.  
 84 In familiar environments where the structure is em-  
 85 bedded in cortical connections, attractor dynamics  
 86 provide a complementary mechanism for rapid eval-  
 87 uation of many possible futures in parallel.

88 **A mechanistic theory** In this paper, we show  
 89 that known features of PFC representations and  
 90 connectivity are sufficient to implement a powerful  
 91 planning algorithm with minimal additional assump-  
 92 tions. The resulting ‘spacetime attractor’ (STA)  
 93 instantiates an explicit world model in the synaptic  
 94 connections between neurons, which allows it to  
 95 plan by inferring optimal future trajectories. This

algorithm resembles other cortical circuits known to infer features of the present, thereby unifying our understanding of PFC with the rest of cortex. The spacetime attractor excels at dynamic problems with changing reward and transition structures, which PFC is critical for and existing mechanistic models cannot solve. RNNs trained on dynamic tasks implement a spacetime attractor in their dynamics, suggesting that it is an efficient solution. Our findings provide a precise mechanistic theory of adaptive behaviour that reconciles prior work on PFC representations with other known cortical computations and deficits from lesions.

## Results

### An attractor network in space and time

We now introduce the spacetime attractor in more detail. We first review ring and grid attractors, which illustrate how simple circuit motifs can guide inference from partial information. We then explain how the spacetime attractor uses similar principles to infer future behaviour. Neurons are assumed to encode single environment features, but similar ideas apply when individual neurons represent combinations of variables (Clark et al., 2025).

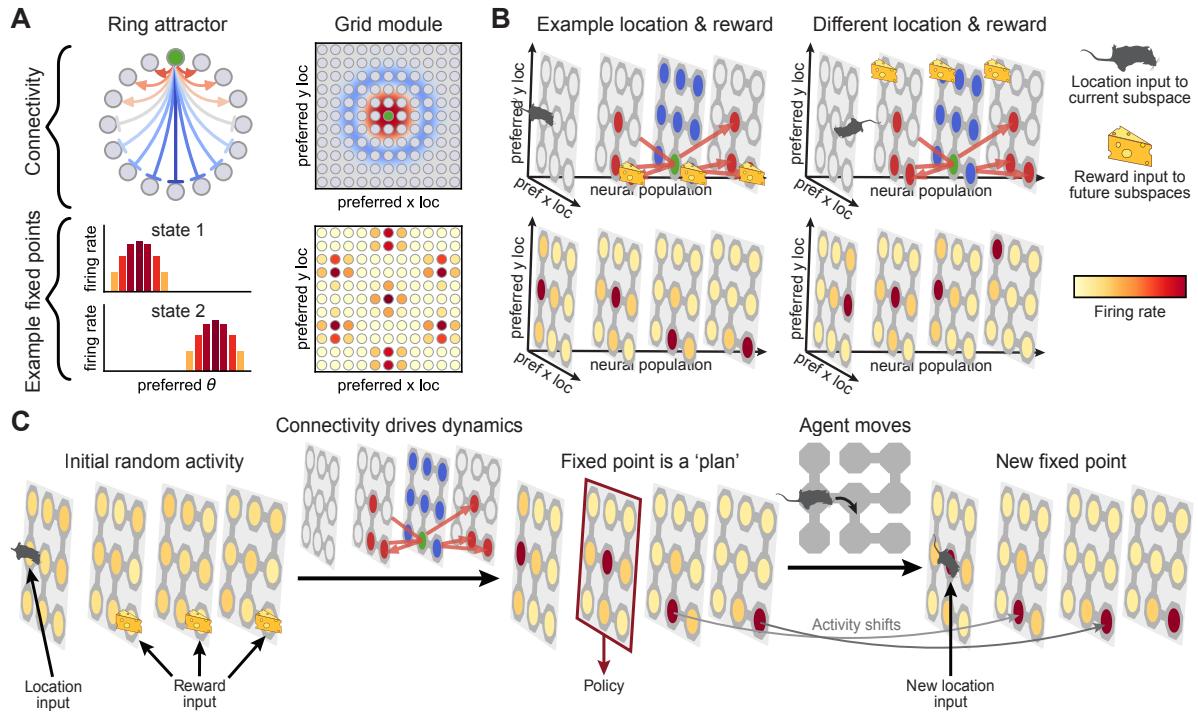
**Ring attractors** Ring attractors are neural networks that can only encode circular variables such as angles (Zhang, 1996; Ben-Yishai et al., 1995). Different neurons have different ‘preferred orientations’, and the network dynamics have a set of ‘fixed points’ (stable activity patterns) that encode a particular angle. At these fixed points, neurons with a preferred angle near the encoded angle are active, and the remaining neurons are silent (Figure 2A, bottom left). This is achieved by a network where neurons are mutually excitatory if they have similar preferred orientations, and inhibitory otherwise (Figure 2A, top left). The network then *infers* which angle is most compatible with noisy inputs such as vision and proprioception. This property underlies the ability of ring attractors to integrate angular velocity in the head direction system (Turner-Evans et al., 2017; Skaggs et al., 1994).

**Grid attractors** Grid attractors instead encode position in two-dimensional space. This is achieved by neurons that are mutually excitatory if they represent nearby locations in space and inhibitory if they represent more distant locations (Figure 2A, top right; Fuhs and Touretzky, 2006; Burak and Fiete, 2009). The fixed points of the network are hexagonal patterns of activity that resemble canonical grid cells (Figure 2A, bottom right; Hafting

et al., 2005), and the inputs determine which location is represented at a given moment in time. These properties underlie the ability of grid cells to perform path integration (Burak and Fiete, 2009).

**Spacetime attractors** Inspired by these known attractor networks in the brain, we propose the existence of a *spacetime attractor* in prefrontal cortex. The STA is a network that infers explicit representations of the future, and its fixed points therefore have to be entire paths through space and time (Figure 2B, bottom). Given such fixed points, any input would automatically be converted into a representation of future behaviour – an inferred plan. Similar to the ring and grid attractors, this is achieved by connecting STA neurons such that *consistent* states excite each other and *inconsistent* states inhibit each other. In the spacetime attractor, each neuron has both a preferred location and a preferred delay  $\delta$ , which determines how far into the future its tuning curve is defined. If a neuron has  $\delta = 0$ , it will fire when the agent is currently at the preferred location of that neuron. If a neuron has  $\delta = 3$ , it will fire when the agent expects to be at the preferred location after three actions. In such a network, consistent states are those that can be part of a single trajectory, while inconsistent states cannot be part of the same trajectory. The connectivity between neurons in subspaces with preferred delays  $\delta$  and  $\delta+1$  should therefore correspond to the structure – more specifically the adjacency matrix – of the environment (Figure 2B, top).

**Reward input selects the future** By embedding a world model in its connections, the spacetime attractor creates fixed points that are future trajectories through space and time. It can then use reward information from the environment to bias the representation towards trajectories that represent desirable futures. This resembles how visual and proprioceptive inputs to a ring attractor bias its representation towards particular orientations (Kim et al., 2019). Reward information must similarly be an input to the STA to enable fast adaptation to changing rewards without rewiring the synaptic connections. Importantly, the STA can accommodate time-varying reward structures known to engage PFC (Figure 1D; Shallice and Burgess, 1991; Volle et al., 2011; Carlesimo et al., 2014). This is possible because different neural populations that represent different times in the future can receive different inputs. Given a restaurant booking in 2 hours and a cinema ticket in 4 hours, the  $\delta = 2$  neurons would receive reward input at the restaurant, and the  $\delta = 4$  neurons at the cinema. The network dynam-



**Figure 2: The spacetime attractor.** **(A)** Left; in a one-dimensional ring attractor, neurons representing similar directions excite each other (top, red), and different directions inhibit each other (blue). Connections are shown for the green example cell. Stable network states (fixed points) have neurons at a particular encoded angle being most active (bottom). Right; grid cells (bottom) emerge from a two-dimensional attractor network (top), where cells with similar preferred location excite each other (red) and intermediate distances inhibit each other (blue). **(B)** The spacetime attractor is a three-dimensional generalisation, where neurons have both a preferred location (planes) and delay (horizontal axis). This resembles the PFC representation in [Figure 1A](#). Cells that represent adjacent locations in both space and time excite each other (top, red), and different locations at the same time inhibit each other (blue). Given inputs indicating the current location (mouse) and future reward (cheese), the fixed points represent reward-maximising paths through space and time (bottom). While this example has a single static goal, reward inputs can also differ between subspaces if the reward is dynamic. **(C)** Dynamics of the spacetime attractor. Initial activity is diffuse (left), followed by convergence to a stable representation of a plan (centre). At convergence, a policy is read out from the subspace that represents the next location along the desired trajectory. When the agent moves to the next state, neural activity updates to represent the remaining plan-to-go (right; [El-Gaby et al., 2023](#)).

ics would then infer a future where the agent is at each location at the appropriate time.

To summarise, the spacetime attractor has four main components: (i) different neural populations represent different times in the future; (ii) the neurons are connected according to the structure of the environment; (iii) the current location is an input to the ‘present’ subspace; and (iv) the reward structure is an input to all future subspaces. The resulting fixed points ([Figure 2B](#), bottom) represent trajectories that maximise cumulative reward. The network dynamics implement a gradual relaxation process, where reward inputs bias the representation in each subspace towards high-reward locations, while inputs from neighbouring subspaces constrain their representations to be consistent. Together, these

two processes converge to a stable representation of a coherent future plan ([Figure 2C](#)).

**The STA guides behaviour** After computing an explicit representation of a plan, it can be used to inform behaviour. In particular, an agent should simply take whichever action leads it to the next location on the inferred trajectory ([Figure 2C](#), middle). The entire representation of the future then needs to move by one ‘action’, so the state that was previously represented in subspace  $\delta$  is now in subspace  $\delta - 1$ . The resulting dynamics resemble a ‘conveyor belt’ that allows the STA to execute entire trajectories without recomputing the policy after every action. Such conveyor belt dynamics have been observed experimentally during sequence working memory ([El-Gaby et al., 2023](#)). Impor-

tantly, the representation remains a fixed point of the spacetime attractor dynamics, because the location and reward inputs are updated due to (i) the movement of the agent, and (ii) the passing of time.

**World models for planning** Planning requires access to an internal world model that predicts the consequences of different actions. However, the spacetime attractor uses such a world model differently from many other planning algorithms. Most algorithms apply a world model sequentially to simulate actions one by one. The spacetime attractor instead instantiates multiple copies of the world model explicitly in the synapses between different neural subspaces. This allows the network to simulate many possible futures in parallel. Entorhinal grid cells are also thought to embed a world model in their connectivity (McNaughton et al., 2006), but they do not represent the distant future explicitly (Ouchi and Fujisawa, 2024). Instead, they infer a single location at a time. The spacetime attractor therefore suggests that circuit principles in prefrontal cortex resemble other cortical areas that use structural knowledge to infer features of the world. We propose the major difference to be that PFC instantaneously represents the world at many points in time, which extends known circuit principles to complex planning problems. This also requires knowledge of the reward available in different states, which could be estimated separately by other neural circuits. In this work, we simply assume access to ground-truth rewards.

### STAs are flexible planners

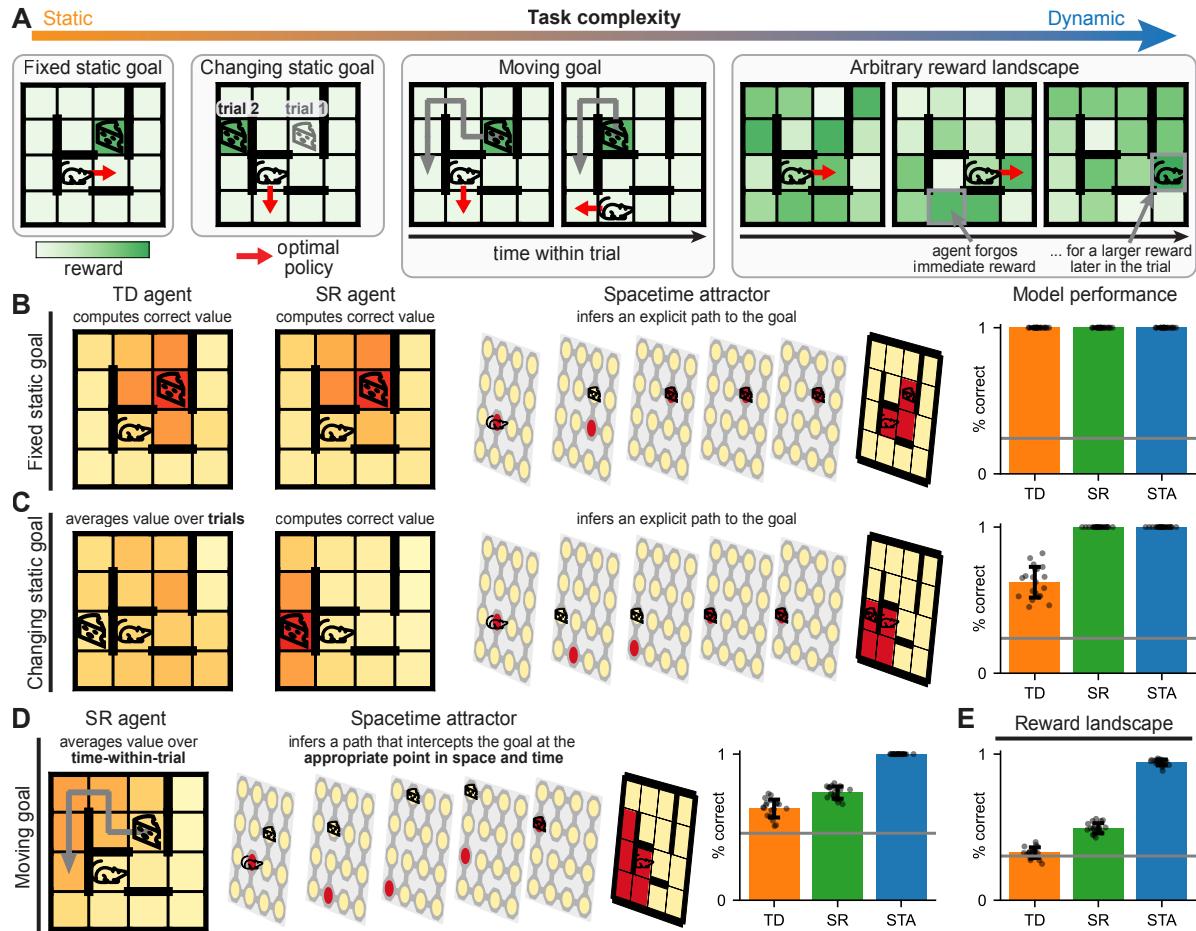
Prefrontal cortex is particularly important for tasks that require flexible behaviour in changing environments (Figure 1D; Burgess and Wu, 2013). To understand whether the spacetime attractor is a good model of planning in PFC, we therefore need to (i) study its behaviour and performance in such dynamic tasks, and (ii) characterise when and how it differs from existing models. We will see that the spacetime attractor is well-suited to dynamic ‘PFC-like’ problems, which other mechanistic models struggle to solve.

**From static to dynamic tasks** We designed a set of four tasks that vary in how much the reward changes in space and time (Figure 3A). The tasks are all embedded in Euclidean space for simplicity of exposition, but the underlying principles generalise to any environment with known structure. Task 1 is simple navigation towards a static goal that remains constant across trials. In task 2, the static goal changes between trials. Task 3 is

navigation to a goal that also moves within each trial, where each location is only rewarded when the goal is there. Task 4 generalises the idea of time-dependent goals to a non-binary reward landscape. Reward magnitudes are sampled independently between -1 and +1 for each location at each time point in each trial. The objective is to maximise cumulative reward over the trial, which requires balancing immediate reward with the potential for future reward (Figure 3A, right). This is reminiscent of the example from Figure 1D, where the restaurant and cinema are desirable at different times. To better understand how the STA solves these tasks, we compare it to two models commonly studied in systems neuroscience. The first is temporal difference (TD) learning, which has neural correlates in striatum (Sutton, 1988; Schultz et al., 1997). The second is the successor representation (SR), which has neural correlates in hippocampus (Dayan, 1993; Stachenfeld et al., 2017). These models assume fixed rewards across trials (TD) or within a trial (SR). Unlike the spacetime attractor, they do not generalise well when the reward changes rapidly.

**Simpler models can solve static tasks** TD learners gradually propagate value from rewarded locations to all other locations, and optimal policies are only learned when rewards are constant across both time and trials (Figure 3B-C, orange). The SR agent uses the environment structure and trial-specific reward function to compute values, and it can solve tasks where the goal changes between trials (Figure 3B-C, green). The STA can also solve these tasks, but it does so by inferring an entire spatial trajectory to the goal instead of computing a value function (Figure 3B-C, blue). Because the reward is constant throughout a trial, each subspace receives the same reward input. However, the inferred representation differs across subspaces because it is also constrained to satisfy the current agent location and the transition structure of the world. As we will see, the inference of an explicit spacetime representation by the STA increases its flexibility compared to the amortised state values computed by TD and SR agents.

**Only the STA solves dynamic tasks** While the SR agent can adapt to rewards that change across trials, the reward structure still needs to remain constant for the duration of the planning horizon. This is because the SR computes time-averaged occupancy and lacks fine-grained information about when the agent is where. When the reward changes within a trial, the SR is limited to computing average values across time-within-trial, and it fails



**Figure 3: Model comparisons in static and dynamic tasks.** (A) Example tasks with different degrees of dynamic reward. Colours indicate reward at each location (white to green), and red arrows indicate the optimal policy. In the ‘moving goal’ task, the agent intercepts a target that moves along a different trajectory in each trial (arrow). In the ‘reward landscape’ task, the reward is sampled independently in space and time. The agent has to maximise cumulative reward, and it can be optimal to forgo immediate reward for a later payoff. (B) Representations and performance in the static goal task with a fixed goal. Left: value functions learned by TD and SR agents. Centre: STA representation at convergence, which encodes a path through space and time. The final plane is a max projection that summarises the path. Right: performance of each model in the task (grey: random baseline). (C) As in (B), now for the task where the static goal changes between trials. (D) Representations and performance in the moving goal task. Left: the SR computes a value function that averages reward across time-within-trial. Middle: the STA takes into account the moving goal and computes a path that intercepts it. Right: performance of each agent. (E) Performance in the reward landscape task. All error bars indicate 1 standard deviation across 20 agents and environments. Dots indicate individual agent performance.

337 to intercept the goal at the correct point in space  
 338 and time (Figure 3D). In contrast to the TD and  
 339 SR agents, the STA can solve the moving goal task  
 340 because the input to each subspace is the reward at  
 341 that specific moment in time. The representation  
 342 is therefore biased towards coinciding with the goal  
 343 in both space and time, and the network dynamics  
 344 relax to a future trajectory that correctly intercepts  
 345 the goal (Figure 3D). The advantage of the STA is  
 346 exacerbated in the more dynamic reward landscape  
 347 task, where very different locations can be rewarded

348 at different times (Figure 3E). These results suggest  
 349 that different algorithms could contribute differently  
 350 to decision making. TD learning could drive rapid  
 351 decisions in environments with stable rewards; the  
 352 SR would be more flexible when rewards change  
 353 on intermediate timescales; and the STA could en-  
 354 able adaptive behaviour in environments with rapid  
 355 changes. Finally, sequential search would facilitate  
 356 slower planning in novel environments and could be  
 357 guided by partial plans from a spacetime attractor.

358 **STAs are efficient planners**

359 We have seen that a handcrafted spacetime attractor  
360 can solve dynamic problems that simpler algorithms  
361 cannot. However, this requires many neurons to  
362 explicitly represent the future. We therefore investi-  
363 gated whether superior solutions exist to these  
364 ‘PFC-like’ problems by training a recurrent neural  
365 network (RNN) to solve them efficiently (Mante  
366 et al., 2013; Stroud et al., 2023). We will see that  
367 regularised RNNs solve dynamic planning problems  
368 with an STA-like algorithm, suggesting that it is ef-  
369 ficient. This result also demonstrates that the STA  
370 is consistent with a prominent theory that PFC  
371 resembles a recurrent meta-learner (Wang et al.,  
372 2018). In this view, PFC learns connections over  
373 long timescales that implement adaptive behaviour  
374 on short timescales through recurrent dynamics.  
375 This is exactly how planning happens in the STA,  
376 which embeds task structure in the weights and  
377 uses recurrent dynamics for rapid decision making  
378 with different rewards. We extend the meta-  
379 learning theory of PFC by showing that an STA  
380 is the mechanism implemented in the dynamics of  
381 RNNs meta-trained on challenging planning tasks.

382 The three defining features that allow a spacetime  
383 attractor to plan via inference are: (i) it has an  
384 explicit representation of the future; (ii) the connec-  
385 tivity embeds a model of the world; and (iii) attractor  
386 dynamics infer the future representation. We  
387 will see that RNNs trained on the reward landscape  
388 task exhibit all of these properties. The RNNs were  
389 trained with supervised learning to generate optimal  
390 actions from inputs that indicated the current loca-  
391 tion and trial-specific reward structure (Figure S1;  
392 Methods; see Supplementary Note for a discussion  
393 of different modelling choices). The reward input  
394 was only provided during an initial ‘planning phase’  
395 prior to the ‘execution phase’ of the task, and the  
396 loss function also penalised the magnitude of neural  
397 firing rates and network weights to encourage an  
398 efficient solution.

399 **RNNs learn explicit future representations**

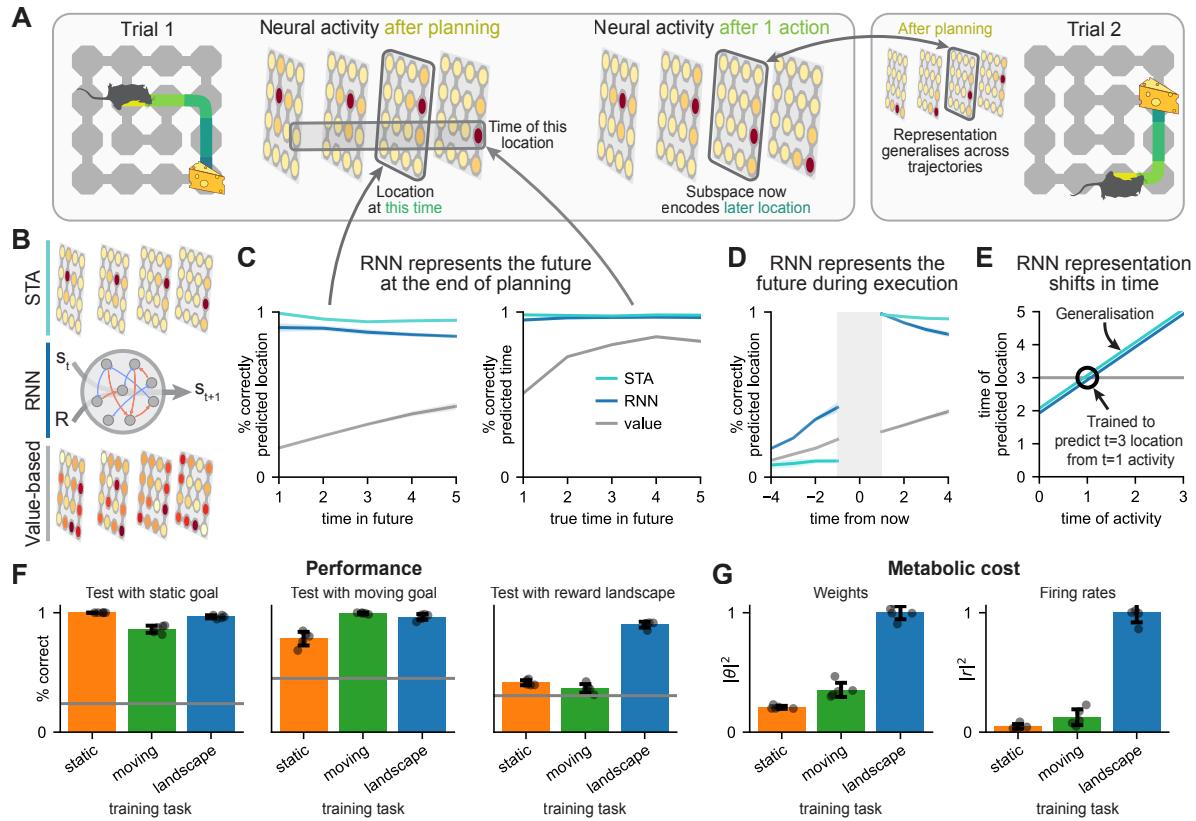
400 Spacetime attractors compute explicit representa-  
401 tions of the future (Figure 4A-B). Critically, these  
402 representations generalise over trajectories. In other  
403 words, the representation in subspace  $\delta$  depends  
404 only on the expected location in  $\delta$  actions, and  
405 not on the rest of the trajectory or the history of  
406 the agent (Figure 4A, right). To confirm that the  
407 RNN learned such a representation, we trained lin-  
408 ear decoders to predict the future from its hidden  
409 state. To ensure generalisation across trajectories,  
410 we trained the decoders while holding out each ‘cur-

411 rent location’ and computed the test accuracy only  
412 from those held out locations. The RNN had an  
413 explicit representation of the future during both  
414 planning and execution (Figure 4C-D; Figure S2).  
415 In a spacetime attractor, the same subspace always  
416 encodes location  $\delta$  actions into the future. The fu-  
417 ture representation of the RNN also exhibited such  
418 ‘conveyor belt’ dynamics (Figure 4E). The neural  
419 representation of the RNN is thus consistent with  
420 a spacetime attractor.

421 Intriguingly, RNNs trained on the simpler static and  
422 moving goal tasks did not reliably learn explicit fu-  
423 ture representations (Figure S3). Additionally, the  
424 RNN trained on the reward landscape task could  
425 solve these simpler tasks, while RNNs trained on  
426 the simpler tasks failed in the reward landscape  
427 task (Figure 4F). Alternative solutions therefore  
428 exist in the simpler tasks, and these solutions are  
429 favored by a pressure to be energetically efficient  
430 (Figure 4G). We identify a tradeoff between gen-  
431 erality and efficiency, where the PFC-like solution  
432 is general but needs more neurons and synapses,  
433 while simpler algorithms can solve simpler tasks  
434 more efficiently. Unlike other theories of PFC such  
435 as value coding, this suggests a reason why PFC  
436 occupies such a large cortical territory.

437 **RNNs learn a world model** In a spacetime at-  
438 tractor, subspaces that represent expected future lo-  
439 cations are connected according to the environment  
440 structure. Testing this prediction in the RNN re-  
441 quires interpretation of the network weights, which  
442 is challenging because the functions of individual  
443 neurons are unknown. We therefore project the  
444 weights into a coordinate system, where the axes  
445 are orthogonal directions in neural state space that  
446 predict each point in spacetime (Methods). While  
447 this does not perfectly recover the true subspaces  
448 in the handcrafted STA, it is a good estimate. The  
449 projected RNN weights can therefore be interpreted  
450 as interactions between representations of different  
451 points in spacetime (Figure 5A).

452 Remarkably, the recurrent weights between adjacent  
453 subspaces closely resembled the adjacency matrix  
454 of the environment (Figure 5B; Figure S4; mean  
455  $\pm$  std correlation of  $0.91 \pm 0.07$  vs.  $0.72 \pm 0.06$   
456 for control environments). The RNN thus learned  
457 a world model explicitly in its recurrent weights.  
458 Consistent with a spacetime attractor, the  $\delta = 0$   
459 subspace received input indicating the current agent  
460 location, and subspace  $\delta$  received input indicating  
461 the reward in  $\delta$  actions (Figure 5C). The RNN  
462 also learned weaker connections between more dis-

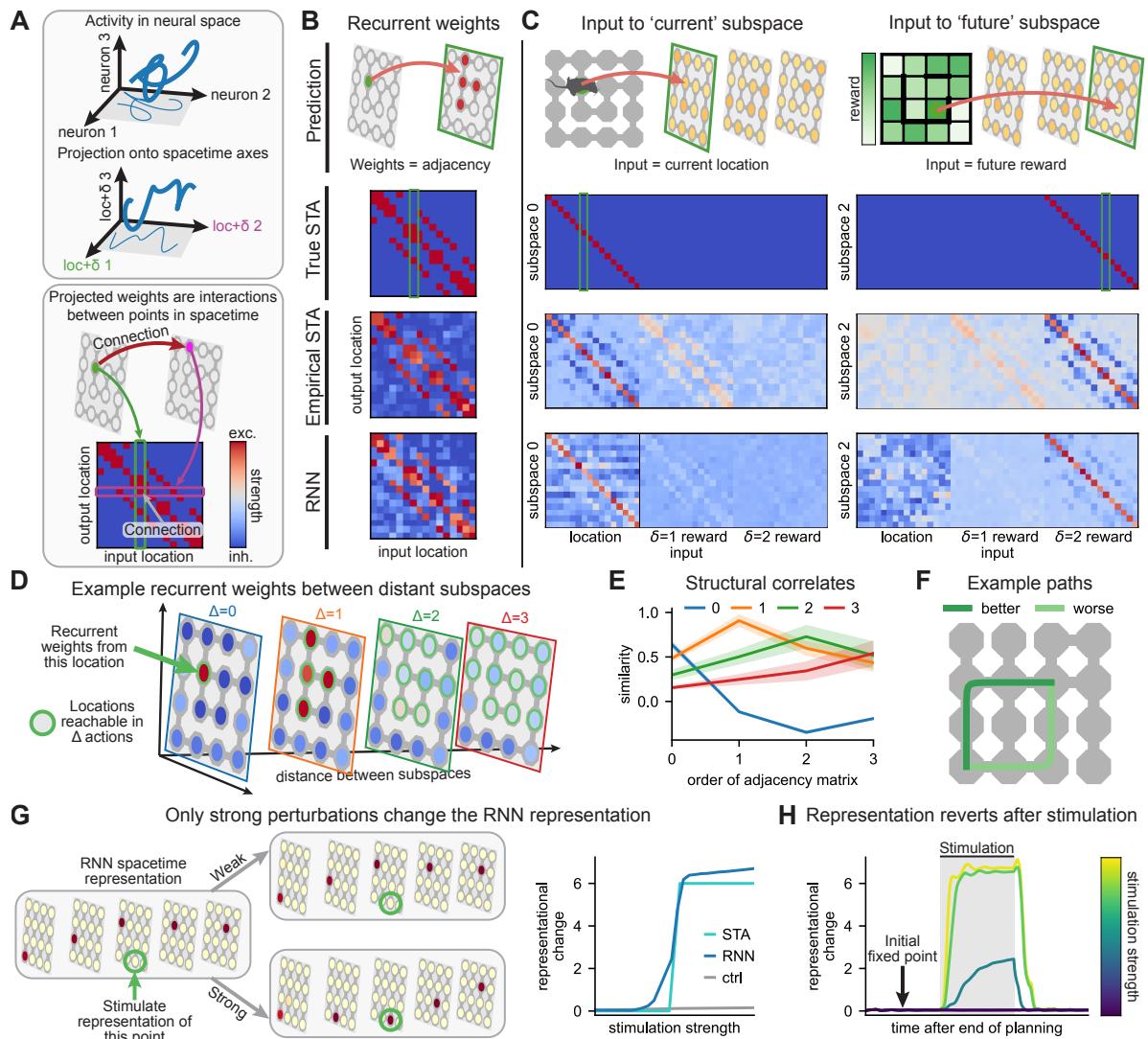


**Figure 4: RNNs learn spacetime representations in the reward landscape task.** (A) An STA infers the entire future during planning and represents the ‘future-to-go’ during execution (left). The representation in each subspace generalises across all trajectories passing through that location in spacetime (right). (B) We compare a spacetime attractor; an RNN trained on the reward landscape task; and an agent that computes an exact value function in space and time. The value-based agent computes an optimal policy from ‘neural activity’ containing (i) the value function, (ii) the agent location, and (iii) the time-within-trial (Methods). (C) Decoding accuracy at the end of planning for: left; agent location at each time in the future. Right; the time at which the agent will be at a given location, plotted as a function of the actual time the location was visited. Decoders were trained in crossvalidation across the current agent location (Methods). (D) Decoding accuracy during execution of location at each time in the past or future. (E) We trained a single decoder to predict location at time 3 from neural activity at time 1 (black circle). The same decoder predicted location at time  $t + 2$  from neural activity at any other  $t$ , demonstrating ‘conveyor belt’ dynamics. (F) Performance in the static goal (left), moving goal (centre), and reward landscape (right) tasks for RNNs trained on either task (x-labels; colours). (G) Normalised parameter magnitudes of the three RNNs (left) and average firing rates in the static goal task (right). All error bars indicate 1 standard deviation across 5 RNNs (dots).

tant subspaces, and these connections reflected the environment structure (Figure 5D-E; Supplementary Note). In summary, the RNN learned all of the structural components that allow a spacetime attractor to infer future actions by integrating expected reward across time using a world model.

**RNNs learn attractor dynamics** In a spacetime attractor, the structured connections give rise to attractor dynamics with fixed points corresponding to explicit representations of desirable futures. To demonstrate similar attractor dynamics in the trained RNN, we artificially perturbed its neural activity at the end of planning. This resembles how

attractor dynamics have been demonstrated in biological circuits (Inagaki et al., 2019; Kim et al., 2017; Vinograd et al., 2024). An attractor network should be robust to small perturbations, and the representation should be more sensitive to perturbations towards other attractor states than perturbations along random directions in neural state space. We analysed the RNN in a setting with two ‘good’ paths that were close in value (Figure 5F). Like the hand-crafted STA, the RNN initially converged to a stable representation of the better path, which persisted in the presence of weak perturbations. Stronger perturbation of a single location on the alternate



**Figure 5: RNNs learn a world model.** (A) The weights of RNNs trained on the reward landscape task are projected into an orthonormal coordinate system with axes that predict different points in spacetime (top). The projected weights are interactions between points in spacetime (bottom). (B) The average recurrent weights between subspaces separated by a single action resemble the environment adjacency matrix (bottom, ‘RNN’). ‘Empirical STA’ is the same analysis performed on approximate subspaces estimated from neural activity in the handcrafted STA. ‘True STA’ indicates weights between the ground truth future-coding subspaces in the handcrafted STA. Green box in ‘True STA’ indicates weights between (i) the single location denoted by a green circle in ‘Prediction’, and (ii) all locations in the subspace indicated by a green square. (C) Input weights to the ‘current’ ( $\delta = 0$ ; left) and a ‘future’ ( $\delta = 2$ ; right) subspace. Consistent with an STA, the current subspace of the RNN receives location input, and future subspaces receive reward information. (D) Average recurrent weights between an example location (green arrow) and other locations at increasing time differences ( $\Delta$ ; planes). Locations in subspaces  $\Delta$  apart are more strongly connected if they can be reached in  $\Delta$  actions. Light green circles indicate these  $\Delta^{\text{th}}$  order adjacency matrices ( $\Delta = 0$  corresponds to the identity). (E) Correlation between (i) the average connectivity between subspaces separated by  $\Delta$  actions (lines; legend), and (ii) different order adjacency matrices (x-axis). Shading indicates standard deviation across 5 RNNs. (F) Example environment with a high-value path and a lower-value path. (G) Weak stimulation does not affect the spacetime representation of the RNN, but strong stimulation switches it to the lower-value path. (H) Magnitude of representational change over time across stimulation strengths (Methods).

path caused a discrete switch to a representation of this entire path, including in subspaces not directly perturbed (Figure 5G). Perturbations of similar

magnitude in random control directions had little effect (Figure 5G, grey). Unlike the handcrafted STA, the RNN representation relaxed back to the

495 original fixed point after removing the perturbation  
496 (**Figure 5H; Figure S5**), which suggests that the  
497 RNN learned an attractor landscape that is less  
498 susceptible to local minima.

499 Together, our analyses show that RNNs trained on  
500 a dynamic planning task learn to closely approx-  
501 imate a spacetime attractor. This was also true  
502 across variations in model architecture (**Methods**;  
503 **Figure S6; Figure S7**). Finally, RNNs with too few  
504 hidden units to learn a spacetime attractor failed  
505 to solve the task (**Figure S8**), suggesting that other  
506 solutions are not readily learned by gradient de-  
507 scent. The spacetime attractor is thus an efficient  
508 solution to dynamic planning problems, and it is  
509 consistent with a prominent theory that PFC can  
510 be understood as a recurrent meta-learner.

### 511 STAs can generalise across environments

512 Hardwiring the transition structure of the environ-  
513 ment into the synapses of a spacetime attractor  
514 seems prohibitive, because it might prevent the net-  
515 work from generalising across different structures  
516 (**Figure 6A**). Here, we show that is not the case.  
517 RNNs trained in environments with changing struc-  
518 ture still learn spacetime attractors. The network  
519 dynamics are adapted to each environment through  
520 input-mediated gating of a generic scaffold. This is  
521 a simple mechanism for rapid adaptation that pre-  
522 serves the ability to plan under changing rewards.

### 523 Representations adapt to the environment

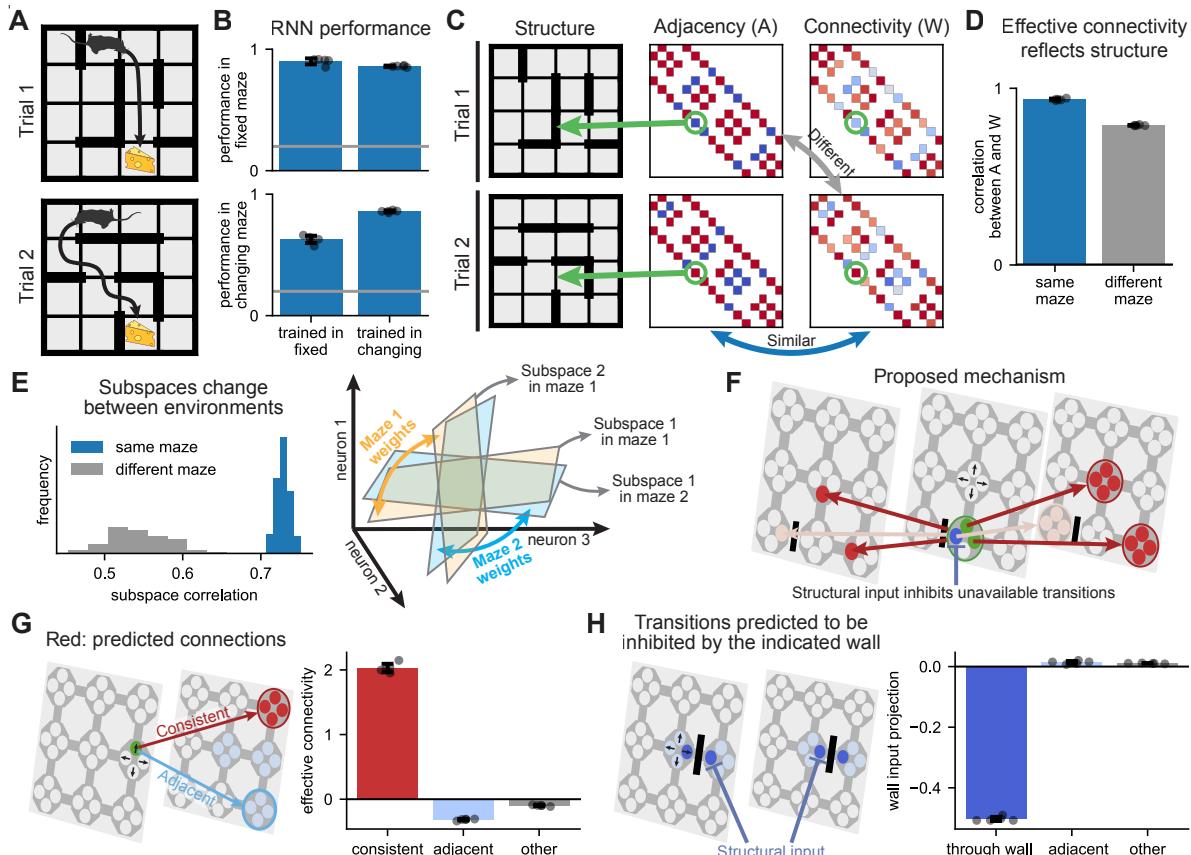
524 To study planning in environments with changing  
525 structure, we trained RNNs on a variant of the re-  
526 ward landscape task, where the transition structure  
527 was different on every trial (**Methods**). The struc-  
528 ture was provided to the RNN as a binary input that  
529 indicated whether pairs of otherwise adjacent states  
530 were separated by a wall. Importantly, the network  
531 had to adapt to each new environment through its  
532 dynamics while keeping the weights fixed across en-  
533 vironments (Wang et al., 2018; Jensen et al., 2024).  
534 The generalising RNNs performed nearly as well  
535 as networks trained in a single environment (**Fig-**  
536 **ure 6B**), raising the question of how generalisation  
537 is achieved. The neural representation resembled  
538 a spacetime attractor in any single environment  
539 (**Figure S9B**). However, the effective connectivity  
540 between subspaces changed between environments.  
541 Remarkably, the connectivity reflected the transi-  
542 tion structure of whichever environment the agent  
543 was currently in (**Figure 6C-D**). This indicates a  
544 general solution that allows recurrent networks to  
545 use attractor dynamics for planning in any similar  
546 environment, without requiring synaptic plasticity.

547 **A model of generalisation** We can understand  
548 how the network adapts to the environment at both  
549 the ‘population level’ and at a mechanistic level.  
550 At a population level, the directions in neural state  
551 space that represent each point in spacetime were  
552 more similar in different trials from the same envi-  
553 ronment than across trials from different environ-  
554 ments (**Figure 6E**, left). This change allows the  
555 network to align its representation with the ap-  
556 propriate components of the connectivity matrix  
557 (**Figure 6E**, right). How is this achieved mechani-  
558 cally? We hypothesised that the network explicitly  
559 represents future *transitions* rather than future loca-  
560 tions. In other words, there are directions in neural  
561 state space that represent ‘*in  $\delta$  actions from now,*  
562 *move from state  $s_i$  to state  $s_j$* ’, which we denote  $\tau_{\delta}^{ij}$ .  
563 This is different from the vanilla spacetime attrac-  
564 tor, which has a representation of future locations  
565 that is invariant to the subsequent state.

566 While it requires more neurons, explicitly rep-  
567 resenting future transitions allows the network to use  
568 input-mediated inhibition to adaptively gate off  
569 transitions that are not available in any particular  
570 environment. The learned structure of the remain-  
571 ing transitions is then used for planning (**Figure 6F**).  
572 Importantly, the scaffold only needs to include transi-  
573 tions that are possible in at least one environment.  
574 This model makes three predictions that we verified:  
575 (i) the network explicitly represents future transi-  
576 tions (**Figure S9C**), (ii) the representations of future  
577 transitions are connected by a world model (**Fig-**  
578 **ure 6G**), and (iii) input indicating a wall between  
579  $s_i$  and  $s_j$  specifically inhibits representations of fu-  
580 ture transitions between those states (**Figure 6H**).  
581 These analyses show that the spacetime attractor  
582 is an efficient planning algorithm in structurally  
583 changing environments, and they suggest a general  
584 mechanism for adaptive behaviour.

### 585 Discussion

586 We have taken inspiration from known cortical attrac-  
587 tor dynamics and recent findings on prefrontal  
588 sequence representations (**Figure 1**) to develop a  
589 new theory of planning. This spacetime attractor  
590 model uses different neural subspaces to represent  
591 the expected location of an agent at different times  
592 in the future. The subspaces are connected ac-  
593 cording to the structure of the environment, which  
594 allows the network to infer optimal future trajec-  
595 tories (**Figure 2**). The spacetime attractor can solve  
596 problems with dynamic rewards, which PFC is im-  
597 portant for and simpler algorithms struggle with  
598 (**Figure 3**). It solves these problems using inter-



**Figure 6: Spacetime attractors can adapt to changing structure.** (A) Changing structure requires adaptation. (B) Performance of RNNs trained on the reward landscape task in a single maze or with a different structure on every trial, when evaluated in either a single maze (top) or across changing mazes (bottom). (C) Two example mazes (left) and their corresponding adjacency matrices (centre). The effective connectivity between adjacent subspaces in a single RNN (right) reflects the structure of each maze. (D) Average correlation between subspace connectivity in a maze and the structure of either the same (blue) or a random (grey) maze. (E) The spacetime representation is more similar across distinct trials from the same maze than trials from different mazes (left). By using slightly different subspaces in different environments, the network can match their connectivity to the environment structure (right; orange vs. blue). (F) Putative mechanism for structural generalisation. Instead of representing each future *location*, neurons encode expected future *transitions* (black arrows in example state). Each transition  $\tau_{\delta}^{ij}$  (green) connects to transitions that can follow ( $\tau_{\delta+1}^{jk}$ ) or precede ( $\tau_{\delta-1}^{li}$ ) it (red arrows). Structural input to PFC inhibits transitions that are not available in a given environment (blue), preventing planning between states that would otherwise be connected (light red arrows). (G) Effective connectivity between directions in neural state space that encode ‘consistent’ consecutive future transitions ( $\tau_{\delta}^{ij}$  and  $\tau_{\delta+1}^{jk}$ ; green to red), ‘adjacent’ transitions (green to blue), or any other transitions (green to grey). (H) Projection of the input from wall  $w_{ij}$  onto representations of future transitions *through* the wall ( $\tau^{ij}$ ; dark blue), transitions to other adjacent states ( $\tau^{ik}$ ; light blue), or any other transitions (grey). All error bars indicate 1 standard deviation across 5 RNNs (dots).

599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
nal representations that closely resemble prefrontal  
representations during sequence working memory,  
thereby unifying working memory and planning  
in PFC. RNNs trained to solve dynamic planning  
tasks learn to implement the spacetime attractor  
algorithm in their internal dynamics, suggesting  
that it is an efficient solution (Figure 4; Figure 5).  
Finally, spacetime attractors can generalise across  
environments with different transition structures

through rapid gating of a general scaffold to reflect  
each particular environment (Figure 6).

**Experimental predictions** The spacetime attractor is inspired by data and makes precise predictions that can be tested in future experiments:

- After planning a behavioural sequence, different subspaces of PFC activity should represent distinct steps of the plan.

616 • Optogenetic activation of neurons in a future sub-  
617 space should bias representations and behaviour  
618 towards the stimulated state at a delay corre-  
619 sponding to the stimulated subspace.

620 • The effective connectivity between subspaces  
621 should reflect the structure of the environment.  
622 This could be tested in noise correlations across  
623 neurons or by explicit holographic stimulation.

624 • Different patterns of neurons should represent  
625 the future in environments with different transi-  
626 tion structures.

627 • The neurons active in each environment should  
628 be connected according to its structure. Inputs  
629 that mediate this gating could come through  
630 sensory cortex in observable environments or  
631 hippocampus when the structure is learned.

632 The first two predictions would also be true of inde-  
633 pendent subspaces used to store sequence memories  
634 ([Xie et al., 2022](#); [El-Gaby et al., 2023](#)), while the  
635 last three are unique predictions of a structured  
636 network that plans behavioural sequences.

637 **Decision making across the brain** Our com-  
638 parisons of different models across different tasks  
639 suggest complementary contributions of different  
640 brain regions to decision making. Striatum is  
641 thought to implement temporal difference learn-  
642 ing ([Schultz et al., 1997](#)), which facilitates rapid  
643 responses in stable environments. Hippocampus em-  
644 beds the structure of the world, but it is thought to  
645 either average representations over future timesteps  
646 in a successor representation ([Stachenfeld et al.,  
647 2017](#)) or represent one state at a time in a sequence  
648 ([George et al., 2021](#); [Whittington et al., 2020](#); [Jensen  
649 et al., 2024](#)). This provides an efficient solution  
650 to problems with consistent spatial structure and  
651 rewards that change on intermediate timescales.  
652 A spacetime attractor in frontal cortex facilitates  
653 adaptive behaviour in dynamic environments with  
654 a familiar structural scaffold that has been embed-  
655 ded in PFC. Finally, explicit search would facilitate  
656 slower planning in novel environments and could be  
657 guided by partial plans from a spacetime attractor.  
658 While spacetime attractors are particularly useful  
659 in dynamic environments, they can also solve sim-  
660 pler problems. It is therefore possible that animals  
661 have developed a spacetime attractor because they  
662 sometimes need it for planning, and then reuse it  
663 in simpler tasks. However, it is also plausible that  
664 many laboratory behaviours instead engage simpler  
665 algorithms in the basal ganglia or hippocampus.  
666 If that is the case, richer spacetime problems may  
667 be needed to query the representations used for  
668 planning in prefrontal cortex.

## 669 **Interactions with other planning algorithms**

670 Planning via inference in a spacetime attractor dif-  
671 fers from most models of planning in cognitive sci-  
672 ence. In particular, the spacetime attractor posits  
673 that planning can happen via *recognition*, where  
674 the sequence of steps needed to reach some de-  
675 sired state is directly inferred. In contrast, many  
676 studies of human planning focus on explicit search,  
677 where different paths are simulated and evaluated  
678 sequentially. We suggest that these two processes  
679 coexist, with spacetime attractor dynamics facil-  
680 itating rapid planning as inference in familiar envi-  
681 ronments, where the structure has been embedded  
682 in prefrontal connections. The STA could also help  
683 focus explicit search towards putative high-value  
684 paths and evaluate the utility of different paths.  
685 This could happen through interactions with hip-  
686 pocampal activity sequences that have been pro-  
687 posed to facilitate planning (either replays or theta  
688 sequences; [Foster, 2017](#); [Widloski and Foster, 2022](#);  
689 [Jensen et al., 2024](#)). In particular, [Jensen et al.  
690 \(2024\)](#) suggest that PFC biases which sequences are  
691 replayed in hippocampus. PFC would then update  
692 its representation to make high-reward sequences  
693 more likely in an iterative policy improvement pro-  
694 cess. Hippocampal replay has also been proposed  
695 to implement a ‘DYNA’ algorithm, where past ex-  
696 perience is used to learn a model-free policy during  
697 rest or sleep ([Mattar and Daw, 2018](#)). This is com-  
698 plementary to both the spacetime attractor and  
699 explicit search. DYNA facilitates rapid decision  
700 making when the optimal policy remains stable  
701 across time, while decision-time planning allows  
702 adaptation to changing environments.

703 **Learning a spacetime attractor** In the hand-  
704 crafted spacetime attractor, we embedded copies of  
705 the environment adjacency matrix in the connec-  
706 tions between every pair of consecutive subspaces.  
707 The RNN analyses show that such structure can be  
708 learned from repeated experience. However, learning  
709 copies of the same parameters independently  
710 for each pair of subspaces is inefficient. It would  
711 be more efficient to store a cache of experienced  
712 trajectories that can be used to learn all of the  
713 parameters. Hippocampal replay has been proposed  
714 to build cognitive maps via such interactions with  
715 prefrontal cortex ([Bakermans et al., 2023](#); [Ou et al.,  
716 2025](#)). In particular, [Ou et al. \(2025\)](#) suggest that  
717 hippocampal replay consolidates structural infor-  
718 mation from hippocampus into cortex. Replay from  
719 hippocampus to PFC during sleep or rest could  
720 therefore provide a mechanism for learning all of  
721 the STA parameters from the same data.

722 **The state space of planning** We have assumed  
723 that planning happens at the level of individual  
724 locations in the environment. However, humans  
725 often plan in more abstract spaces, which improves  
726 efficiency by reducing the required planning depth.  
727 We also plan hierarchically by first computing an  
728 abstract plan that can be refined in increasing lev-  
729 els of detail (Eckstein and Collins, 2020). Multiple  
730 spacetime attractors operating at different levels  
731 of abstraction can be coupled to implement such  
732 a hierarchical planner. One ‘abstract’ STA would  
733 compute a high-level plan. A second STA would  
734 treat the next abstract state as a goal and compute  
735 a detailed plan to get there. Stacking spacetime at-  
736 tractors in this way allows inference of plans that are  
737 exponentially long in the depth of the hierarchy, and  
738 therefore the number of neurons, in contrast to the  
739 linear scaling of a non-hierarchical STA. However,  
740 learning appropriate abstractions and embedding  
741 them in an STA remains an unsolved challenge.

742 **Bidirectional environment interactions** We  
743 have focused on environments that evolve indepen-  
744 dently of the agent. The rewards and structure  
745 of the environment can therefore be estimated up  
746 front and provided as inputs to a spacetime at-  
747 tractor in PFC. However, many problems involve  
748 structure that depends on the behaviour of an agent.  
749 One example is ‘key-door’ problems, where moving  
750 through a locked door becomes possible only after  
751 picking up the key. We posit that such problems can  
752 be solved by structural gating of a learned scaffold,  
753 similar to the STA that generalises across environ-  
754 ment structures (Figure 6). The structural gating  
755 would no longer be an external input, but instead a  
756 function of the spacetime representation itself. For  
757 example, activating a cell that represents ‘being at  
758 the key’ could disinhibit ‘transitioning through the  
759 door’ at a later time.

760 Another important example of bidirectional interac-  
761 tions is social behaviours, where agents can change  
762 their plans in response to each other. A system of  
763 two spacetime attractors could simulate the joint  
764 dynamics of two such agents. The first STA ‘A’  
765 infers future behaviour for agent A. The reward  
766 input would be the output of a different STA ‘B’  
767 that predicts the behaviour of another agent. The  
768 reward for STA B would itself be an output of STA  
769 A, which couples the predicted behaviour of the  
770 two agents. The combined system relaxes to a fixed  
771 point where the behaviour of A is optimal given  
772 the predicted behaviour of B and vice versa – a  
773 putative neural implementation of ‘theory of mind’.

774 **Outlook** We have developed a new theoretical  
775 framework for planning in prefrontal cortex. It  
776 builds on recently discovered prefrontal working  
777 memory representations and known attractor dy-  
778 namics in other neural circuits. The spacetime  
779 attractor extends these principles to prospective  
780 behaviours in complex and changing environments –  
781 a setting that has previously eluded mechanistic cir-  
782 cuit models. Existing data does not allow us to test  
783 these ideas explicitly. Instead, we have provided  
784 a series of concrete predictions for future exper-  
785 iments. We hope this will inspire new work in both  
786 experimental and computational neuroscience.

## 787 Author contributions

788 KTJ and TEJB developed the conceptual frame-  
789 work with input from TA, PD, MSM, SR, and AB.  
790 KTJ ran the simulations and performed the anal-  
791 yses. KTJ made the figures with help from MSM,  
792 SR, and AB. KTJ and TEJB wrote the manuscript  
793 with input from all authors.

## 794 Code availability

795 Code is available at [github.com/KrisJensen/pysta](https://github.com/KrisJensen/pysta).  
796 This includes code for running example models and  
797 reproducing all results from the paper.

## 798 Acknowledgments

799 We are grateful to Diksha Gupta, Mohamady El-  
800 Gaby, Will Dorrell, and Tom Mrsic-Flogel for  
801 helpful feedback on the manuscript. This work  
802 was supported by a Wellcome Principal Research  
803 Fellowship (219525/Z/19/Z; TEJB, KTJ, AB);  
804 the Gatsby Initiative for Brain Development and  
805 Psychiatry (GAT3955; TEJB); the Jean Francois  
806 and Marie-Laure de Clermont Tonerre Foundation  
807 (TEJB); a Wellcome Trust Career Development  
808 Award (225926/Z/22/Z; TA); the Oxford Claren-  
809 don Fund (PD); the Human Frontier Science Pro-  
810 gram (LT0040/2024-L; SR); EMBO (ALTF 651-  
811 2023; SR); and a FYSSEN postdoctoral study grant  
812 (MSM). KTJ, TEJB, MSM, and SR were also sup-  
813 ported by the Sainsbury Wellcome Centre’s core  
814 provided by Wellcome (219627/Z/19/Z) and the  
815 Gatsby Charitable Foundation (GAT3755).

## 816 References

- 817 Bakermans, J. J., Warren, J., Whittington, J. C.,  
818 and Behrens, T. E. (2023). Constructing future  
819 behaviour in the hippocampal formation through  
820 composition and replay. *bioRxiv*, pages 2023–04.  
821 Ben-Yishai, R., Bar-Or, R. L., and Sompolinsky,

- 822 H. (1995). Theory of orientation tuning in visual cortex. *Proceedings of the National Academy of Sciences*, 92(9):3844–3848.
- 823
- 824
- 825 Blanco-Pozo, M., Akam, T., and Walton, M. E. (2024). Dopamine-independent effect of rewards on choices through hidden-state inference. *Nature Neuroscience*, 27(2):286–297.
- 826
- 827
- 828
- 829 Botvinick, M. and Toussaint, M. (2012). Planning as inference. *Trends in cognitive sciences*, 16(10):485–488.
- 830
- 831
- 832 Botvinick, M. M. and Plaut, D. C. (2006). Short-term memory for serial order: a recurrent neural network model. *Psychological review*, 113(2):201.
- 833
- 834
- 835 Burak, Y. and Fiete, I. R. (2009). Accurate path integration in continuous attractor network models of grid cells. *PLoS computational biology*, 5(2):e1000291.
- 836
- 837
- 838
- 839 Burgess, P. W. and Wu, H. (2013). Rostral prefrontal cortex (Brodmann area 10). *Principles of frontal lobe function*, pages 524–544.
- 840
- 841
- 842 Callaway, F., van Opheusden, B., Gul, S., Das, P., Krueger, P. M., Griffiths, T. L., and Lieder, F. (2022). Rational use of cognitive resources in human planning. *Nature Human Behaviour*, 6(8):1112–1125.
- 843
- 844
- 845
- 846
- 847 Carlesimo, G. A., di Paola, M., Fadda, L., Caltagirone, C., and Costa, A. (2014). Prospective memory impairment and executive dysfunction in prefrontal lobe damaged patients: is there a causal relationship? *Behavioural neurology*, 2014(1):168496.
- 848
- 849
- 850
- 851
- 852
- 853 Chaudhuri, R., Gercek, B., Pandey, B., Peyrache, A., and Fiete, I. (2019). The intrinsic attractor manifold and population dynamics of a canonical cognitive circuit across waking and sleep. *Nature neuroscience*, 22(9):1512–1520.
- 854
- 855
- 856
- 857
- 858 Chen, J., Zhang, C., Hu, P., Min, B., and Wang, L. (2024). Flexible control of sequence working memory in the macaque frontal cortex. *Neuron*.
- 859
- 860
- 861 Clark, D. G., Abbott, L., and Sompolinsky, H. (2025). Symmetries and continuous attractors in disordered neural circuits. *bioRxiv*, pages 2025–01.
- 862
- 863
- 864
- 865 Dayan, P. (1993). Improving generalization for temporal difference learning: The successor representation. *Neural computation*, 5(4):613–624.
- 866
- 867
- 868 Dorrell, W., Hsu, K., Hollingsworth, L., Lee, J. H., Wu, J., Finn, C., Latham, P. E., Behrens, T. E., and Whittington, J. C. (2024). Don't cut corners: Exact conditions for modularity in biologically inspired representations. *arXiv preprint arXiv:2410.06232*.
- 869
- 870
- 871
- 872
- 873
- 874 Eckstein, M. K. and Collins, A. G. (2020). Computational evidence for hierarchically structured reinforcement learning in humans. *Proceedings of the National Academy of Sciences*, 117(47):29381–29389.
- 875
- 876
- 877
- 878
- 879 El-Gaby, M., Harris, A. L., Whittington, J. C., Dorrell, W., Bhomick, A., Walton, M. W., Akam, T., and Behrens, T. E. (2023). A cellular basis for mapping behavioural structure. *bioRxiv*, pages 2023–11.
- 880
- 881
- 882
- 883
- 884 Foster, D. J. (2017). Replay comes of age. *Annu. Rev. Neurosci.*, 40(581-602):9.
- 885
- 886 Fuhs, M. C. and Touretzky, D. S. (2006). A spin glass model of path integration in rat medial entorhinal cortex. *Journal of Neuroscience*, 26(16):4266–4276.
- 887
- 888
- 889
- 890 George, D., Rikhye, R. V., Gothoskar, N., Guntupalli, J. S., Dedieu, A., and Lázaro-Gredilla, M. (2021). Clone-structured graph representations enable flexible learning and vicarious evaluation of cognitive maps. *Nature communications*, 12(1):2392.
- 891
- 892
- 893
- 894
- 895
- 896 Hafting, T., Fyhn, M., Molden, S., Moser, M.-B., and Moser, E. I. (2005). Microstructure of a spatial map in the entorhinal cortex. *Nature*, 436(7052):801–806.
- 897
- 898
- 899
- 900 Iacaruso, M. F., Gasler, I. T., and Hofer, S. B. (2017). Synaptic organization of visual space in primary visual cortex. *Nature*, 547(7664):449–452.
- 901
- 902
- 903
- 904 Inagaki, H. K., Fontolan, L., Romani, S., and Svoboda, K. (2019). Discrete attractor dynamics underlies persistent activity in the frontal cortex. *Nature*, 566(7743):212–217.
- 905
- 906
- 907
- 908 Jensen, K. T. (2023). An introduction to reinforcement learning for neuroscience. *arXiv preprint arXiv:2311.07315*.
- 909
- 910
- 911 Jensen, K. T., Hennequin, G., and Mattar, M. G. (2024). A recurrent network model of planning explains hippocampal replay and human behavior. *Nature neuroscience*, 27(7):1340–1348.
- 912
- 913
- 914

- 915 Kim, S. S., Hermundstad, A. M., Romani, S., Abbott, L., and Jayaraman, V. (2019). Generation 962  
916 of stable heading representations in diverse visual 963  
917 scenes. *Nature*, 576(7785):126–131. 964  
918
- 919 Kim, S. S., Rouault, H., Druckmann, S., and Jayaraman, V. (2017). Ring attractor dynamics in the 965  
920 drosophila central brain. *Science*, 356(6340):849– 966  
921 853. 967  
922
- 923 Kingma, D. P. and Ba, J. (2015). Adam: A method 968  
924 for stochastic optimization. In Bengio, Y. and 969  
925 LeCun, Y., editors, *3rd International Conference 970  
926 on Learning Representations, ICLR 2015, San 971 Diego, CA, USA, May 7-9, 2015, Conference 972  
927 Track Proceedings*. 973
- 928 Ko, H., Hofer, S. B., Pichler, B., Buchanan, 974  
929 K. A., Sjöström, P. J., and Mrsic-Flogel, T. D. 975  
930 (2011). Functional specificity of local synaptic 976  
931 connections in neocortical networks. *Nature*, 977  
932 473(7345):87–91. 978  
933
- 934 Lee, T. S. and Nguyen, M. (2001). Dynamics of 979  
935 subjective contour formation in the early visual 980  
936 cortex. *Proceedings of the National Academy of 981  
937 Sciences*, 98(4):1907–1911. 982  
983
- 938 Levine, S. (2018). Reinforcement learning and 984  
939 control as probabilistic inference: Tutorial and 985  
940 review. *arXiv preprint arXiv:1805.00909*. 986
- 941 Mante, V., Sussillo, D., Shenoy, K. V., and New- 987  
942 some, W. T. (2013). Context-dependent compu- 988  
943 tation by recurrent dynamics in prefrontal cortex. 989  
944 *nature*, 503(7474):78–84. 990  
991
- 945 Mattar, M. G. and Daw, N. D. (2018). Prioritized 992  
946 memory access explains planning and hippocam- 993  
947 pal replay. *Nature neuroscience*, 21(11):1609– 994  
948 1617. 995
- 949 McNaughton, B. L., Battaglia, F. P., Jensen, O., 996  
950 Moser, E. I., and Moser, M.-B. (2006). Path inte- 997  
951 gration and the neural basis of the ‘cognitive map’. 998  
952 *Nature Reviews Neuroscience*, 7(8):663–678. 999  
1000
- 953 Ou, J., Qu, Y., Xu, Y., Xiao, Z., Behrens, T., and 1001  
954 Liu, Y. (2025). Replay builds an efficient cog- 1002  
955 nitive map offline to avoid computation online. 1003  
956 *bioRxiv*, pages 2025–01. 1004
- 957 Ouchi, A. and Fujisawa, S. (2024). Predictive grid 1005  
958 coding in the medial entorhinal cortex. *Science*, 1006  
959 385(6710):776–784. 1007
- 960 Pedregosa, F., Varoquaux, G., Gramfort, A., 1008  
961 Michel, V., Thirion, B., Grisel, O., Blondel, 1009
- M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830.
- Ruff, D. A., Markman, S. K., Kim, J. Z., and Cohen, M. R. (2025). Linking neural population formatting to function. *bioRxiv*.
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., et al. (2020). Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609.
- Schultz, W., Dayan, P., and Montague, P. R. (1997). A neural substrate of prediction and reward. *Science*, 275(5306):1593–1599.
- Shallice, T. and Burgess, P. W. (1991). Deficits in strategy application following frontal lobe damage in man. *Brain*, 114(2):727–741.
- Shin, H., Ogando, M. B., Abdeladim, L., Durand, S., Belski, H., Cabasco, H., Loefler, H., Bawany, A., Hardcastle, B., Wilkes, J., et al. (2023). Recurrent pattern completion drives the neocortical representation of sensory inference. *bioRxiv*.
- Skaggs, W., Knierim, J., Kudrimoti, H., and McNaughton, B. (1994). A model of the neural basis of the rat’s sense of direction. *Advances in neural information processing systems*, 7.
- Stachenfeld, K. L., Botvinick, M. M., and Gershman, S. J. (2017). The hippocampus as a predictive map. *Nature neuroscience*, 20(11):1643–1653.
- Stroud, J. P., Watanabe, K., Suzuki, T., Stokes, M. G., and Lengyel, M. (2023). Optimal information loading into working memory explains dynamic coding in the prefrontal cortex. *Proceedings of the National Academy of Sciences*, 120(48):e2307991120.
- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine learning*, 3:9–44.
- Tian, Z., Chen, J., Zhang, C., Min, B., Xu, B., and Wang, L. (2024). Mental programming of spatial sequences in working memory in the macaque frontal cortex. *Science*, 385(6716):eadp6091.
- Turner-Evans, D., Wegener, S., Rouault, H., Franconville, R., Wolff, T., Seelig, J. D., Druckmann, S., and Jayaraman, V. (2017). Angular velocity integration in a fly heading circuit. *Elife*, 6:e23496.

- 1010 Turner-Evans, D. B., Jensen, K. T., Ali, S., Patterson, T., Sheridan, A., Ray, R. P., Wolff, T., Lauritzen, J. S., Rubin, G. M., Bock, D. D., et al. (2020). The neuroanatomical ultrastructure and function of a biological ring attractor. *Neuron*, 108(1):145–163. 1041  
1042  
1043  
1044  
1045  
1046
- 1016 Vinograd, A., Nair, A., Kim, J. H., Linderman, S. W., and Anderson, D. J. (2024). Causal evidence of a line attractor encoding an affective state. *Nature*, 634(8035):910–918. 1047  
1048  
1049  
1050
- 1020 Volle, E., Gonen-Yaacovi, G., de Lacy Costello, A., Gilbert, S. J., and Burgess, P. W. (2011). The role of rostral prefrontal cortex in prospective memory: a voxel-based lesion study. *Neuropsychologia*, 49(8):2185–2198. 1051  
1052  
1053  
1054  
1055
- 1025 Walton, M. E., Behrens, T. E., Buckley, M. J., Rudebeck, P. H., and Rushworth, M. F. (2010). Separable learning systems in the macaque brain and the role of orbitofrontal cortex in contingent learning. *Neuron*, 65(6):927–939. 1056  
1057  
1058  
1059
- 1030 Wang, J. X., Kurth-Nelson, Z., Kumaran, D., Tirumala, D., Soyer, H., Leibo, J. Z., Hassabis, D., and Botvinick, M. (2018). Prefrontal cortex as a meta-reinforcement learning system. *Nature neuroscience*, 21(6):860–868. 1060  
1061  
1062  
1063  
1064
- 1035 Whittington, J. C., Dorrell, W., Behrens, T. E., Ganguli, S., and El-Gaby, M. (2023). On prefrontal working memory and hippocampal episodic memory: Unifying memories stored in weights and activation slots. *bioRxiv*, pages 2023–11. 1065  
1066  
1067  
1068  
1069
- 1040 Whittington, J. C., Muller, T. H., Mark, S., Chen, G., Barry, C., Burgess, N., and Behrens, T. E. (2020). The Tolman-Eichenbaum machine: Unifying space and relational memory through generalization in the hippocampal formation. *Cell*, 183(5):1249–1263. 1041  
1042  
1043  
1044  
1045  
1046
- 1016 Widloski, J. and Foster, D. J. (2022). Flexible rerouting of hippocampal replay sequences around changing barriers in the absence of global place field remapping. *Neuron*, 110(9):1547–1558. 1047  
1048  
1049  
1050
- 1020 Xie, Y., Hu, P., Li, J., Chen, J., Song, W., Wang, X.-J., Yang, T., Dehaene, S., Tang, S., Min, B., et al. (2022). Geometry of sequence working memory in macaque prefrontal cortex. *Science*, 375(6581):632–639. 1051  
1052  
1053  
1054  
1055
- 1025 Zhang, K. (1996). Representation of spatial orientation by the intrinsic dynamics of the head-direction cell ensemble: a theory. *Journal of Neuroscience*, 16(6):2112–2126. 1056  
1057  
1058  
1059
- 1030 Zheng, J., Guimaraes, R., Hu, J. Y., Perona, P., and Meister, M. (2024). Mice in the manhattan maze: Rapid learning, flexible routing and generalization, with and without cortex. *Cognitive Computational Neuroscience*. 1060  
1061  
1062  
1063  
1064
- 1035 Zintgraf, L., Shiarlis, K., Igl, M., Schulze, S., Gal, Y., Hofmann, K., and Whiteson, S. (2019). VariBAD: A very good method for Bayes-adaptive deep RL via meta-learning. *arXiv preprint arXiv:1910.08348*. 1065  
1066  
1067  
1068  
1069

## 1070 Methods

### 1071 Tasks

1072 We used three different classes of tasks to train and compare models in this work – the ‘static goal’ task, the  
1073 ‘moving goal’ task, and the ‘reward landscape’ task. All tasks had agents navigate mazes on a 4x4 grid, with  
1074 walls preventing transitions between some pairs of otherwise adjacent states. The wall configurations defining  
1075 the mazes were sampled as described by [Jensen et al. \(2024\)](#). For most analyses, the wall configuration  
1076 remained fixed across all trials, and only the reward function changed. For the analyses in [Figure 6](#), the wall  
1077 configuration changed across trials, and the agent had to adapt to the new transition structure through its  
1078 recurrent network dynamics. In all tasks, the agent location was sampled randomly at the beginning of each  
1079 trial. At the beginning of each trial, the agent location and reward function were ‘frozen’ for 5-7 (randomly  
1080 sampled on each trial) iterations of the environment. This constituted an initial ‘planning’ phase that was  
1081 followed by an ‘execution’ phase, where (i) the agent took actions that changed its location, and (ii) the  
1082 reward function updated as described in more detail below. For all analyses, the full reward function was  
1083 provided to the agent during the initial planning period. For the handcrafted STA and the RNN in [Figure S6](#),  
1084 the input provided reward information that was relative in time during execution – it indicated what the  
1085 reward would be at a given location *in*  $\delta$  *steps* rather than *at time t*. The reward input was zero for all  $\delta$  that  
1086 corresponded to time points beyond the end of the trial. For all other RNN analyses, the reward input was  
1087 set to zero during execution.

### 1088 Moving goal task

1089 In this task, a full goal ‘trajectory’ was sampled on each trial. The goal trajectory was sampled as a random  
1090 walk that could only turn around if it reached a dead end. The start location of the agent was restricted  
1091 to not coincide with the start location of the goal. The trial terminated when the agent was at the same  
1092 location as the goal at a given moment in time, or after a maximum of six actions. The reward input to the  
1093 agent,  $R \in \mathbb{R}^{T \times N}$ , was a matrix indicating the location of the goal at every point in the future. The reward  
1094 input consisted of ‘+0.6’ for the goal location and ‘-0.6’ for all other locations at each moment in time. In  
1095 other words,  $R_{\delta i}$  was ‘+0.6’ if the goal would be at location  $s_i$   $\delta$  steps into the future, and –0.6 otherwise.

### 1096 Static goal task

1097 This task was identical to the moving goal task, except that the goal remained stationary for the duration  
1098 of each trial. The static goal task was implemented in two different variants – one where the goal remained  
1099 fixed across all trials, and one where the goal was resampled at random in every new trial. The handcrafted  
1100 models were evaluated in both of these tasks, while the RNNs were only trained with a goal that changed  
1101 between trials.

### 1102 Reward landscape task

1103 In this task, every element  $R(t, s)$  of the reward function  $\mathcal{R}$  was sampled uniformly and independently between  
1104 -1 and +1. The complete future reward structure was provided as an input to the agent as described above.  
1105 Every trial finished when the agent had taken six actions.

### 1106 Quantification of performance

1107 For most performance comparisons, we computed the probability of choosing the optimal first action in  
1108 each trial. We define the optimal first action as the first action of the trajectory that maximises cumulative  
1109 reward over the entire trial. We use this metric rather than the actual cumulative reward for two reasons. (i)  
1110 We are interested in the process of planning, whereby an agent balances immediate and long-term reward.  
1111 This is most challenging for early actions, while the greedy policy is optimal for the last action. (ii) The  
1112 probability of choosing an optimal action is readily interpretable as a number between 0 and 1. All results  
1113 were qualitatively similar if we instead used the probability of choosing an optimal action averaged over the  
1114 entire trial as a performance metric, or if we used the average empirical reward.

### 1115 Handcrafted models

1116 Here we provide an overview of the handcrafted STA, TD, and SR agents used in the paper.

1117 **Spacetime attractor model**

1118 The spacetime attractor is a recurrent neural network with an exponential nonlinearity and within-subspace  
 1119 normalisation. We define  $z_{\delta i}$  as the ‘potential’ and  $r_{\delta i}$  as the ‘firing rate’ of a neuron that represents ‘being  
 1120 at location  $s_i$  in  $\delta$  actions’. The spacetime attractor then has the following network dynamics:

$$\tau \dot{z}_{\delta i} = -z_{\delta i} + \hat{R}_{\delta i} + \max \left( \epsilon, \log \sum_j A_{ij} r_{\delta-1,j} \right) + \max \left( \epsilon, \log \sum_k A_{ki} r_{\delta+1,k} \right) + \eta \quad (1)$$

$$z_{\delta i} = \max \left( \epsilon, z_{\delta i} - \log \sum_j r_{\delta j} \right) \quad (2)$$

$$r_{\delta i} = e^{z_{\delta i}}. \quad (3)$$

1121 Here,  $\tau = 50$  iterations is the time constant of the dynamics,  $\epsilon = -100$  is a small constant that thresholds  
 1122 the maximum inhibition between and within subspaces, and  $\eta \sim \mathcal{N}(0, \sigma = 0.1)$  is white noise added to the  
 1123 network dynamics.  $A_{ij}$  are weights corresponding to the adjacency matrix, which indicates whether  $s_i$  can be  
 1124 reached from  $s_j$  in one action.  $\hat{R}_{\delta i}$  is an input that reflects the normalised reward available at location  $s_i$  in  $\delta$   
 1125 actions. In particular,  $e^{\hat{R}_{\delta i}} \propto e^{\beta R_{\delta i}}$ , where  $R_{\delta i}$  is the trial-specific reward described above, and we set  $\beta = 9$ .  
 1126 To provide an input indicating the current agent location, we set  $R_{0i} = 20$  if  $s_i$  is the current location and 0  
 1127 otherwise. Since the activity is explicitly normalised within each subspace, the resulting  $r_{\delta i}$  can be interpreted  
 1128 as a distribution over desired locations in  $\delta$  actions. When simulating the behaviour of the STA, we took  
 1129 the policy of the agent to be greedy in the representation in subspace  $\delta = 1$  of locations accessible from the  
 1130 current state,  $a = \operatorname{argmax}_{i \in \mathcal{N}(s)} [r_{1i}]$ . We also added a small amount of noise  $\eta_{ij} \sim \mathcal{U}(-0.025, -0.015)$  to  
 1131 each element of the weight matrix instead of using the exact adjacency matrix. Noise was added to ensure  
 1132 robustness, and it was restricted to be negative by adding a bias term that prevents representations from  
 1133 ‘teleporting’ between locations (Supplementary Note). An additional ‘feedforward component’ in the form of  
 1134 the identity matrix was added to the connectivity between subspace  $\delta$  and  $\delta - 1$  during the first 100 network  
 1135 iterations (2 time constants) after every action. This is inspired by the ‘update neurons’ of the fruit fly  
 1136 head direction circuit (Turner-Evans et al., 2017) and stabilises the ‘conveyor belt’ dynamics, but it is not  
 1137 necessary for any of the main results in the paper (Supplementary Note). We ran the STA dynamics for 400  
 1138 network iterations before each action to ensure convergence to a stable representation.

1139 **Baselines**

1140 Here we provide details of the temporal difference and successor representation baselines that we compare the  
 1141 spacetime attractor to. For further details, we refer to Jensen (2023). Common to both of these frameworks  
 1142 is that they explicitly estimate the ‘value function’ under some policy  $\pi$ :

$$V^\pi(s) = \mathbb{E}_{\tau \sim p_\pi(\tau)} \left[ \sum_{t' \geq t} \gamma^{t' - t} r_{t'} | s_t = s \right]. \quad (4)$$

1143 Here,  $\mathbb{E}_{\tau \sim p_\pi(\tau)}[\cdot]$  indicates an expectation taken over trajectories  $\tau$  resulting from the agent following  $\pi$ . The  
 1144 ‘TD’ and ‘SR’ heatmaps in Figure 3 show the computed value functions. In both cases, we take the policy of  
 1145 the agent to be greedy in the value function evaluated at all locations accessible from the current state.

1146 **Temporal difference learning** We implement vanilla temporal difference learning, which computes a  
 1147 value function by iteratively applying the update

$$\Delta V(s_t) = \alpha(-V(s_t) + r_t + \gamma V(s_{t+1})). \quad (5)$$

1148 We set the temporal discount factor to  $\gamma = 1$ , since we are interested in maximising the non-discounted  
 1149 cumulative reward. We allowed the TD agent to interact with the environment for 4,000 trials with a learning  
 1150 rate of  $\alpha = 0.05$  before analysing its performance.

1151 **Successor representation** In the successor representation formalism, the value function is decomposed as

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0} \gamma^t r_t | s_0 = s \right] \quad (6)$$

$$= \sum_{t=0} \gamma^t \sum_{s'} p_\pi(s_t = s' | s_0 = s) r(s') \quad (7)$$

$$= \mathbf{r}^T \mathbf{m}_s^\pi. \quad (8)$$

1152 Here,  $\mathbf{r}$  is a vector of the average reward associated with each state, and  $\mathbf{m}_s^\pi$  is a vector of the expected  
1153 discounted future occupancy of state  $s'$  if the agent starts in state  $s$  and follows policy  $\pi$ :

$$M_{ss'}^\pi = \sum_{t=0} \gamma^t p_\pi(s_t = s' | s_0 = s). \quad (9)$$

1154 The full matrix  $\mathbf{M}^\pi$ , constructed from stacking the  $\mathbf{m}_s^\pi$  corresponding to all states  $s$ , is denoted the ‘successor  
1155 matrix’, and it allows us to write down a vector of expected rewards from all states as

$$\mathbf{v}^\pi = \mathbf{M}^\pi \mathbf{r}. \quad (10)$$

1156  $\mathbf{M}^\pi$  can be learned using a TD-like algorithm as above. However, since the STA algorithm also has explicit  
1157 access to the transition structure of the environment, we simply computed the exact successor matrix as  
1158 the geometric series  $\mathbf{M}^\pi = \mathbf{I} + \gamma \mathbf{T}^\pi + \gamma^2 (\mathbf{T}^\pi)^2 + \dots = (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1}$  with  $\gamma = 0.95$ . When computing the  
1159 successor matrix, we take  $\pi$  to be the diffusion policy, and  $\mathbf{T}^\pi$  is therefore the diffusion matrix. This is similar  
1160 to how the adjacency matrix serves as a global structural prior in the STA.

1161 **Spacetime value agent** In Figure 4, we compare representations to an agent that computes a value  
1162 function in a state space consisting of space *and* time. This agent uses dynamic programming to compute the  
1163 value of being at location  $s$  at time  $t$  under an optimal policy for every combination of  $s$  and  $t$ :

$$V(T, s) = R(T, s) \quad (11)$$

$$V(t, s) = R(t, s) + \max_{s' \in \mathcal{N}(s)} V(t+1, s'). \quad (12)$$

1164 At time  $t$ , the optimal policy is then greedy in the value of accessible locations at time  $t+1$ . For the decoding  
1165 analyses in Figure 4, we took the ‘neural representation’ to be the concatenation of (i) the flattened value  
1166 function  $\mathbf{v}_f \in \mathbb{R}^{TN}$ , (ii) a one-hot representation of the current location, and (iii) a one-hot representation of  
1167 the current time-within-trial. These three quantities are sufficient to compute an optimal policy.

## 1168 Recurrent neural networks

1169 In this section, we provide details of how recurrent neural networks were trained and analysed. All networks  
1170 were fully connected with  $N_{\text{rec}} = 800$  hidden units (except Figure S8) and ReLU nonlinearities. The network  
1171 dynamics were given by

$$\tau \dot{\mathbf{z}} = -\mathbf{z} + \mathbf{W}_{\text{in}} \mathbf{x} + \mathbf{W}_{\text{rec}} \mathbf{r} + \mathbf{b}_{\text{rec}} + \eta, \quad (13)$$

$$\mathbf{r} = [\mathbf{z}]_+ \quad (14)$$

$$\mathbf{y} = \mathbf{W}_{\text{out}} \mathbf{r} + \mathbf{b}_{\text{out}}. \quad (15)$$

1172 Here,  $\mathbf{z}$  is the network potential,  $\mathbf{r}$  is the ‘firing rate’,  $\mathbf{x}$  is the input,  $\theta = \{\mathbf{W}_{\text{in}}, \mathbf{W}_{\text{rec}}, \mathbf{W}_{\text{out}}, \mathbf{b}_{\text{rec}}, \mathbf{z}_0, \mathbf{b}_{\text{out}}\}$  are  
1173 the network parameters, and  $\eta \sim \mathcal{N}(0, \sigma = 10^{-3})$  is Gaussian noise. All simulations used a time constant of  
1174  $\tau = 5$  iterations. The output policy was defined in global allocentric coordinates as the desired next state,  
1175  $\pi \propto e^{\mathbf{y}}$ . In Figure S7, we also analyse a network that produced a policy in ‘local’ coordinates that indicated  
1176 the desired movement direction. The input consisted of (i) a one-hot representation of the current location  
1177 in the environment; (ii) a binary representation of the location of all walls (Jensen et al., 2024); and (iii)  
1178 the reward  $R_{\delta i}$  as described for the tasks above. Gaussian noise with a standard deviation of  $\sigma = 10^{-3}$  was  
1179 added to the inputs before passing them to the RNN.

1180 For all analyses in the main text, the RNN performed 10 network iterations per environment iteration (action).  
 1181 Reward input was only provided during the planning phase and set to 0 during the execution phase. This  
 1182 was done to avoid biasing the execution-time representation towards a ‘relative’ subspace representation  
 1183 by providing reward input in that format ([Supplementary Note](#)). For comparison with this RNN, we also  
 1184 trained an RNN with reward input during both the planning and execution phase, which learned qualitatively  
 1185 similar representations and dynamics ([Figure S6](#)). For this network, the reward input was given relative to  
 1186 the current time as in the handcrafted STA, and the number of network iterations was randomly sampled  
 1187 between 9 and 11 before each environment iteration.

1188 We optimised all network parameters  $\theta$  to minimise the loss function

$$\mathcal{L}(\theta) = \lambda_\theta \mathcal{L}_{\text{params}}(\theta) + \mathbb{E}_{\tau \sim \pi_{\theta, \text{opt}}} \left[ \sum_t \mathcal{L}_{\text{acc}}^t + \lambda_r \mathcal{L}_{\text{rate}}^t + \lambda_e \mathcal{L}_{\text{ent}}^t \right] \quad (16)$$

$$\mathcal{L}_{\text{params}}(\theta) = |\theta|_2^2 \quad (17)$$

$$\mathcal{L}_{\text{acc}}^t = \sum_{a \in \mathcal{A}_{\text{opt}}^t} -\pi_t(a) \quad (18)$$

$$\mathcal{L}_{\text{rate}}^t = |\mathbf{r}_t|_2^2 \quad (19)$$

$$\mathcal{L}_{\text{ent}}^t = \sum_a \pi_t(a) \log \pi_t(a). \quad (20)$$

1189  $\mathcal{A}_{\text{opt}}^t$  is the set of optimal actions at time  $t$ , and  $\mathbb{E}_{\tau \sim \pi_{\theta, \text{opt}}}[\cdot]$  indicates an expectation over trajectories induced  
 1190 by the policy of the agent, renormalised over optimal actions. That is, we consider an imitation learning  
 1191 setting where ties between equally optimal actions are broken according to the actual policy of the agent.  
 1192 We used  $\lambda_\theta = 2 \times 10^{-7}$ ,  $\lambda_r = 10^{-5}$ , and  $\lambda_e = 10^{-4}$  for all analyses. The RNNs were trained using Adam  
 1193 ([Kingma and Ba, 2015](#)) with a learning rate of  $3 \times 10^{-4}$  for 200,000 batches of 200 trials (250,000 batches for  
 1194 the RNNs trained in changing mazes).

1195 All results are reported as mean and standard deviation across 5 separate RNNs that were trained from  
 1196 different random seeds and with different environment transition structures.

## 1197 Analyses

1198 In this section, we describe the analyses used to compare different handcrafted models and trained RNNs.

### 1199 Decoding analyses

1200 To decode future locations from neural activity ([Figure 4C-D](#)), we used scikit-learn ([Pedregosa et al., 2011](#))  
 1201 to train L2-regularised logistic regression models that predicted location at time  $t_L$  from neural activity at  
 1202 time  $t_N$  with an inverse regularisation strength of  $C = 1.0$ . We performed this analysis in cross-validation  
 1203 across locations at time  $t_N$ . In other words, we trained a decoder on data where the agent was in any state  
 1204  $s_{t_N} \neq s$  to predict all locations at time  $t_L$ , and we then tested this decoder in trials where the agent was in  
 1205 state  $s$  at time  $t_N$ . We repeated this analysis across all test locations  $s$  and averaged performance across the  
 1206 resulting 16 folds. We did this to test whether the RNN had a generalisable representation of future location,  
 1207 rather than an encoding of e.g. current location and neighboring values.

1208 In [Figure 4C](#) (left), we plot the performance of decoders trained on  $t_N = -1$  to predict location at any  
 1209  $t_L \geq 1$ . In [Figure 4D](#), we trained decoders on every pair of  $t_N \in [0, 5]$  and  $t_L \in [0, 5]$ . We then averaged the  
 1210 performance across every ‘delay’  $t_L - t_N$ .

1211 To investigate the generalisation of decoders in [Figure 4E](#), we trained a single decoder using neural activity  
 1212 at time  $t_N = 1$  to predict location at time  $t_L = 3$ . We then applied the same decoder to neural activity at all  
 1213 times  $t'_N$  and quantified how well it predicted location at all times  $t'_L$ . This analysis was again performed in  
 1214 cross-validation. A separate decoder was trained while holding out each location  $s_1 = s$ , and then tested only  
 1215 on trials where  $t'_N = s$ . Performance was averaged across all held-out locations and across 5 independently  
 1216 trained RNNs. [Figure 4E](#) shows the time at which the average predictive performance was highest for each  
 1217  $t'_N$ . [Figure S2](#) shows the full generalisation behaviour of the decoder.

1218 To predict the time at which the agent would be at a particular location in [Figure 4C](#) (right), we analysed

every state  $s$  separately, and then averaged over all  $s$ . For each  $s$ , we identified trials where the RNN passed through  $s$  exactly once. We trained a decoder to predict the time at which  $s$  was visited and computed the test performance as a function of the true time at which  $s$  was visited. As before, all decoders were trained in crossvalidation across current agent location.

## Comparisons of RNN performance and efficiency

In Figure 4F-G, we compare three classes of RNNs, which were trained on either the reward landscape task, the moving goal task, or the static goal task with a goal that changed between trials. The performance of each RNN was quantified in each of the three tasks as the probability of choosing the optimal first action (see above). We also computed the average parameter magnitudes of all the networks,  $\|\theta\|_2^2$ . Finally, we computed the average firing rate of each network in the static goal task as  $\mathbb{E}_{\tau \sim \pi} [\sum_t |\mathbf{r}_t|_2^2]$ . All other RNN analyses apart from Figure S3 used networks trained on the reward landscape task.

## Subspace identification

Prior work has investigated the extent to which working memory subspaces are orthogonal (Xie et al., 2022; Dorrell et al., 2024). However, we are primarily interested in the *dynamics* between subspaces, which are easier to interpret in an orthonormal coordinate system. We therefore asked whether there exists a set of orthogonal subspaces that predict the location of the agent at every time in the future. To do so, we first simulated 6,000,000 trials and collected pairs of neural activity at time  $t$  and location at  $t'(\delta)$  separately for each  $\delta$ . We did this in two different ways. To estimate ‘planning’ subspaces, subspace  $\delta$  was defined as a decoder that predicts location at time  $t'(\delta) = \delta$  from neural activity at times  $t \in \{-2, -1\}$ . To estimate ‘execution’ subspaces, subspace  $\delta$  was defined as a decoder that predicts location at  $t'(\delta) = t + \delta$  from neural activity at any  $t \geq 0$ .

We then defined a predictive distribution parametrised by  $\phi_\delta$  for each  $\delta$ :

$$p_{\phi_\delta}(s_{t'} = s_i) \propto \exp(\mathbf{c}_{\delta,i} \mathbf{r}_t + \mathbf{b}_\delta). \quad (21)$$

Finally, we minimised an objective function that combines the accuracy across  $\delta$ s and the overlap between subspaces:

$$\mathcal{L}(\theta) = \sum_{\delta} \left[ \mathbb{E}_{\mathbf{r}_t, s_{t'(\delta)} \sim \mathcal{D}} [-\log p_{\phi_\delta}(s_{t'(\delta)})] + \alpha_1 |\phi_\delta|_1 + \alpha_2 |\phi_\delta|_2^2 \right] + \alpha_{\text{orth}} \mathcal{L}_{\text{orth}}, \quad (22)$$

$$\mathcal{L}_{\text{orth}} := \sum_{\delta_k, \delta_l} \sum_{i,j} \hat{\mathbf{c}}_{\delta_k, i}^T \hat{\mathbf{c}}_{\delta_l, j}. \quad (23)$$

$\hat{\mathbf{c}}_{\delta_k, i}$  indicates the normalised parameter vector that predicts being at location  $s_i$  at a delay of  $\delta_k$ . The parameters were optimised using ADAM with a learning rate of  $5 \times 10^{-3}$  until convergence or for a maximum of 2000 iterations. We used  $\alpha_1 = 10^{-4}$ ,  $\alpha_2 = 10^{-3}$ , and  $\alpha_{\text{orth}}$  was annealed from 0 to  $2 \times 10^{-3}$  over 500 iterations. These hyperparameters were chosen because they resulted in a good approximation to the ‘true’ subspaces in the handcrafted STA.

## Estimating effective connectivity

To compute the effective connectivity between representations of different points in spacetime, we projected the learned network parameters into a coordinate system defined by the parameter vectors  $\mathbf{C} \in \mathbb{R}^{NT \times N_{\text{rec}}}$ . Each row of  $\mathbf{C}$  is a normalised vector  $\hat{\mathbf{c}}_{\delta, i}$  that predicts a particular point in spacetime. In this coordinate system, the input weights are given by  $\mathbf{W}_{\text{in}}^{\text{eff}} = \mathbf{C} \mathbf{W}_{\text{in}}$ ; the output weights by  $\mathbf{W}_{\text{out}}^{\text{eff}} = \mathbf{W}_{\text{out}} \mathbf{C}^T$ ; and the recurrent weights by  $\mathbf{W}_{\text{rec}}^{\text{eff}} = \mathbf{C} \mathbf{W}_{\text{rec}} \mathbf{C}^T$ . To analyse weights between ‘adjacent’ subspaces, we averaged the blocks of  $\mathbf{W}_{\text{rec}}^{\text{eff}}$  that corresponded to weights from any subspace  $\delta$  to  $\delta + 1$  and from any subspace  $\delta$  to  $\delta - 1$ . To avoid our analyses being biased by the fact that the networks were trained with supervised learning to predict optimal locations that could only be adjacent, we did not include the weights between subspaces 0 and 1 in this average.

To quantify the similarity between the recurrent weights in this spacetime coordinate system and different order adjacency matrices for the environment, we computed point-biserial correlations. The  $\Delta^{\text{th}}$  order adjacency matrix  $\mathbf{A}_\Delta \in \mathbb{R}^{N \times N}$  was defined as a binary matrix with elements equal to 1 for pairs of states

1261 that can be reached from one another in  $\Delta$  actions, and 0 for all other pairs of states. The 0<sup>th</sup> order adjacency  
1262 matrix was defined as the identity matrix.

### 1263 Perturbation analyses

1264 For the perturbation analyses in **Figure 5F-H**, we constructed an environment where the reward function  $R(t, s)$   
1265 was (i) 1.0 for  $(t, s) \in \{(0, 0), (1, 1), (2, 2), (3, 6), (4, 10), (5, 10), (6, 10)\}$ , (ii) 0.7 for  $(t, s) \in \{(1, 4), (2, 8), (3, 9)\}$ ,  
1266 and -1 for all other points in spacetime. The RNN reliably converged to a representation of the optimal path.  
1267 We first ran the network dynamics for 10 environment iterations without perturbation after the end of the  
1268 normal planning period. We then continued to run the network dynamics for 10 environment iterations with a  
1269 bias term defined by  $b_{\text{rec}}^{\text{stim}} = b_{\text{rec}} + \alpha \hat{c}_{2,8}$ , where  $\alpha$  is the stimulation strength. The perturbed representation  
1270 was defined as the representation at the end of this perturbation period. The two example representations in  
1271 **Figure 5G** used  $\alpha = 0.3$  ('weak') and  $\alpha = 10$  ('strong'). The quantification in **Figure 5G** used a range of  $\alpha$   
1272 from 0 to 10 for the RNN, and from 0 to 500 for the handcrafted STA. This is because the attractor wells are  
1273 deeper in the handcrafted STA, and a stronger perturbation is therefore required to push the network out of  
1274 an attractor state. The control analyses (grey lines in **Figure 5G** and **Figure S5C**) were performed by running  
1275 the same analysis on the RNNs, but with stimulation of the same magnitude in a random direction in neural  
1276 state space. For the analysis in **Figure 5H**, we also removed the perturbation and ran the network dynamics  
1277 for 10 environment iterations after the end of the perturbation period. All analyses in the main text focused  
1278 on representations in the space of implied future trajectories,  $p_{\phi_\delta}(s_{t'} = s_i)$ . The 'representational change'  
1279 was quantified as the L1 norm of the difference from the spacetime representation at the end of the normal  
1280 planning period. See **Figure S5** for an analysis of the raw firing rates.

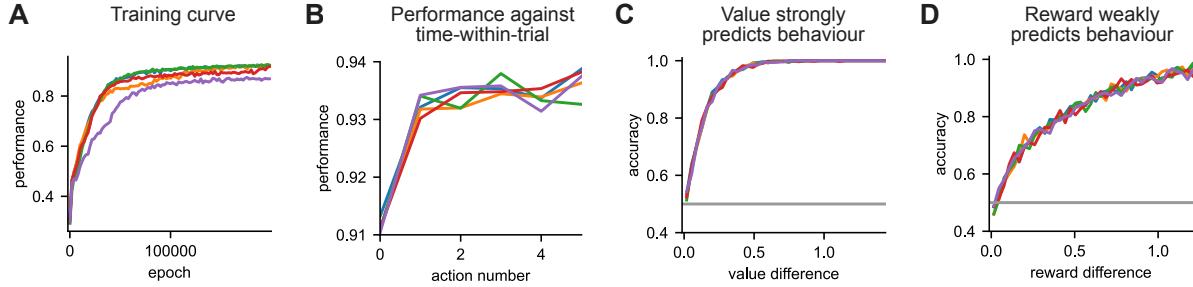
### 1281 RNNs trained in changing environments

1282 For the analyses in **Figure 6**, we trained another set of RNNs in a version of the reward landscape task  
1283 where the transition structure changed between trials. The locations of all walls in the environment were  
1284 provided as a binary input to the agent (see [Jensen et al., 2024](#) for details). For the performance comparison  
1285 in **Figure 6B** (top), we evaluated the performance of these networks in the environments that the 'single  
1286 structure' networks had been trained on.

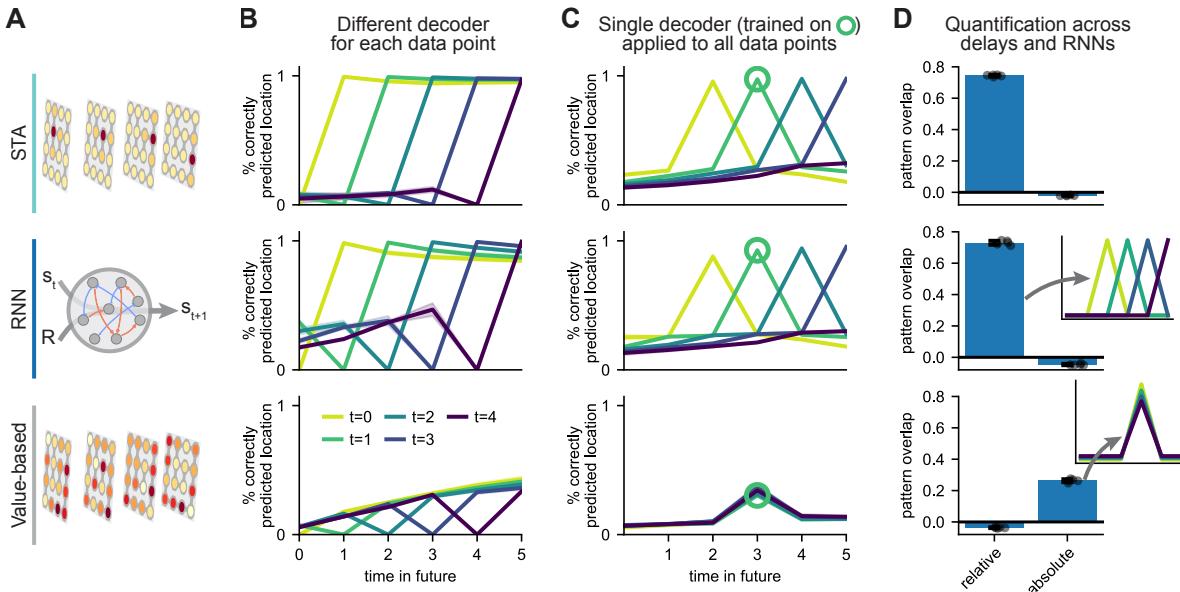
1287 The effective connectivity in **Figure 6C-D** was computed as in the RNNs trained on a single structure. For  
1288 these analyses, we identified the subspaces from 1,000,000 trials in a single environment, and repeated this  
1289 analysis across 30 different environments. We computed similarities between (i) the effective connectivity  
1290 estimated in an environment and the adjacency matrix of the same environment, and (ii) the effective  
1291 connectivity and the adjacency matrix from a different control environment. In **Figure 6E**, the 'subspace  
1292 similarity' was defined as the correlation between the set of vectors that define the subspaces, averaged  
1293 over all points in spacetime. We computed the similarity between (i) subspaces identified from two sets of  
1294 independent trials from the same maze, and (ii) the same number of trials from two different mazes. Recall  
1295 that the effective recurrent weights between subspaces are given by  $\mathbf{W}_{\text{rec}}^{\text{eff}} = \mathbf{C} \mathbf{W}_{\text{rec}} \mathbf{C}^T$ . By using different  
1296 subspaces in different environments, the RNN changes  $\mathbf{C}$  between environments, which changes the effective  
1297 connectivity between pairs of future subspaces (**Figure 6E**).

1298 For the analyses in **Figure 6G-H**, we repeated the subspace identification procedure, but with two notable  
1299 differences. First, we used trials across many environments to find generalised directions in neural state space  
1300 that predict the future in *any* environment. Second, we defined the objective function in terms of future  
1301 transitions  $\tau_\delta^{ij}$  instead of locations. **Figure 6G** quantifies the effective connectivity between future transitions  
1302 in consecutive subspaces ( $\tau_\delta^{ij}$  and  $\tau_{\delta+1}^{kl}$ ) that are either 'consistent' ( $k = j$ ), 'adjacent' ( $j \neq k$  but  $s_i$  and  $s_k$   
1303 are adjacent in an environment with no walls), or any other transitions. In **Figure 6H**, we first projected the  
1304 input corresponding to a given wall location  $w_{ij}$  onto each subspace and normalised the projection within  
1305 each subspace. We then calculated the dot product between this projection and the normalised directions in  
1306 neural state space that predict either (i) transitions between  $s_i$  and  $s_j$ , (ii) transitions from  $s_i$  or  $s_j$  to some  
1307 other state  $s_k$ , or (iii) transitions that do not originate at  $s_i$  or  $s_j$ . Projection magnitudes were averaged  
1308 over transitions within each of these three groups, then across subspaces, and finally the mean and standard  
1309 deviation were computed across 5 RNNs.

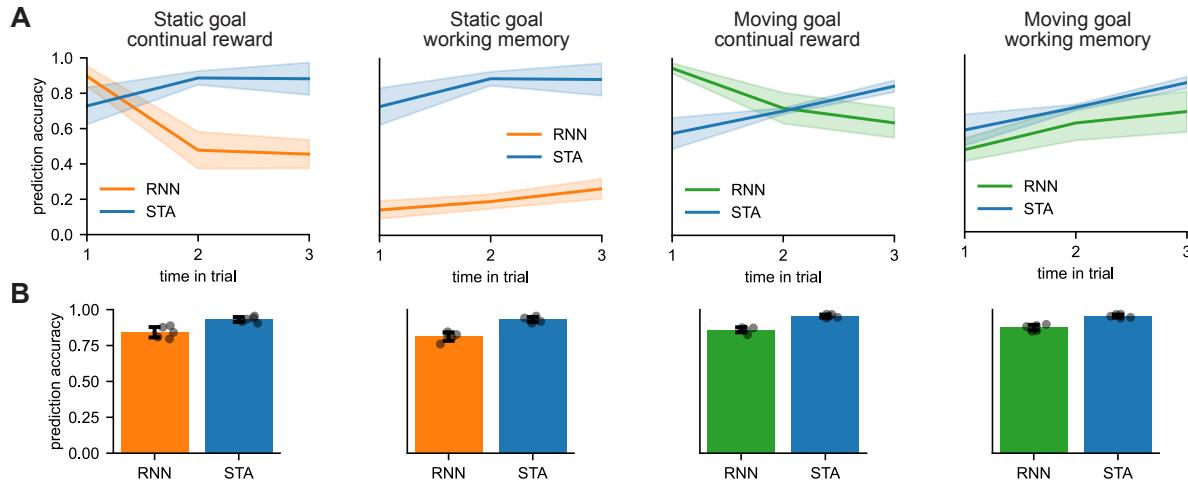
1310 **Supplementary figures**



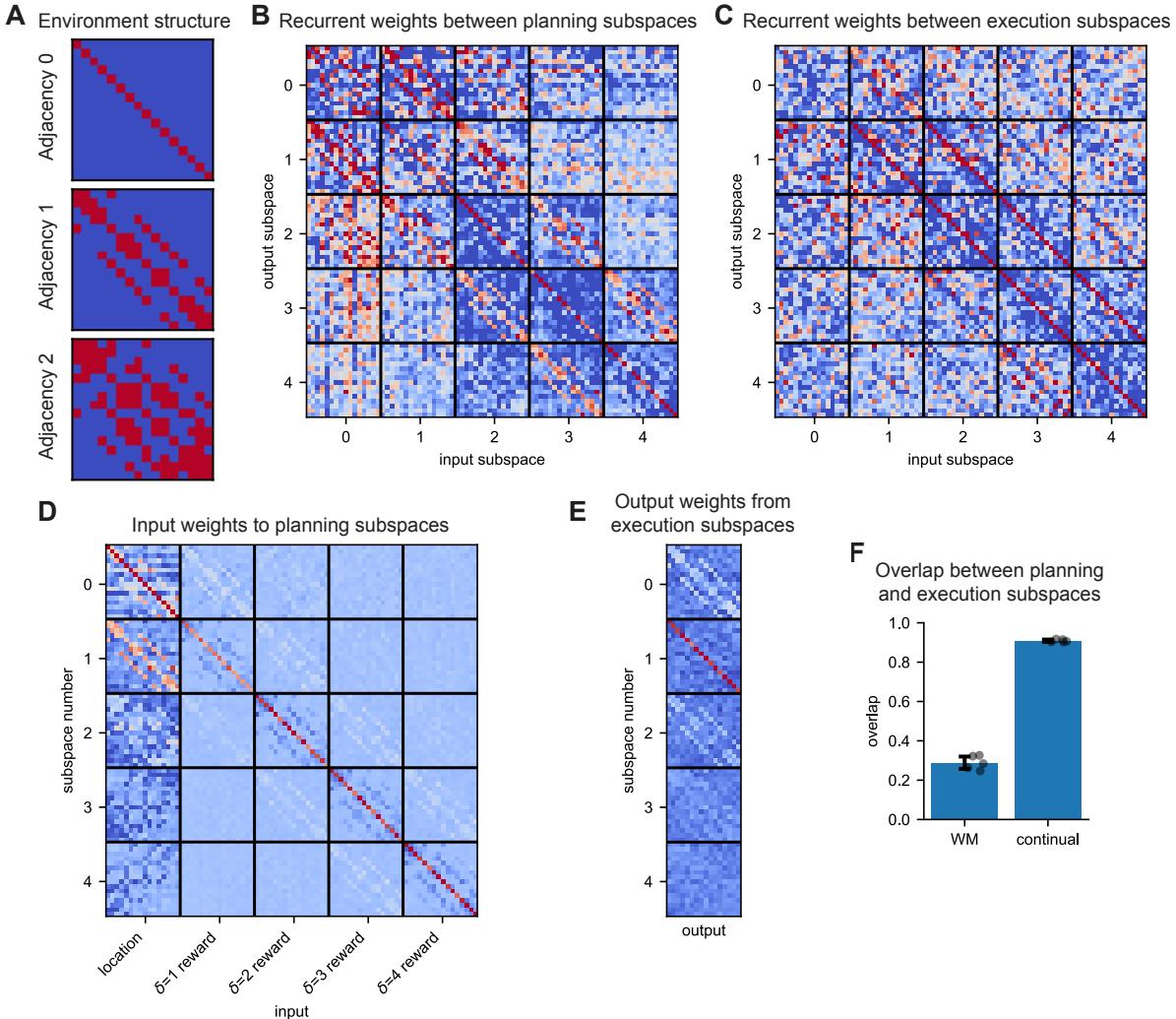
**Figure S1: RNN performance during and after training.** Each line in this figure corresponds to one of the five RNNs that were used for analyses in the main text. **(A)** Performance over the course of training, averaged over all actions within each trial. **(B)** Performance at the end of training as a function of the action number within the trial. When assessing performance at time  $t$ , trials were only included that had optimal choices up to time  $t - 1$ . **(C)** Probability of choosing the action with highest value as a function of the value difference between the two actions with highest value. This analysis shows that errors are only made then the optimal action is close in value to the second best action. **(D)** Probability of choosing the action with highest reward as a function of the reward difference between the two actions with highest reward. Reward is less predictive of behaviour than value, confirming that the RNNs compute long-term value rather than relying on greedy reward.



**Figure S2: Additional analyses of learned RNN representations.** **(A)** We compare a spacetime attractor; an RNN trained on the reward landscape task; and an agent that analytically computes a full spacetime value function. The value-based agent computes an optimal policy from ‘neural activity’ containing (i) the value function, (ii) the current location, and (iii) the time-within-trial (Methods). **(B)** Decoding accuracy of agent location at different times (x-axis) from neural activity at every other time (lines; legend). All decoders were trained in crossvalidation across the current agent location (Methods). This is why the accuracy is zero when decoding location from activity at the same time. **(C)** We trained a single decoder to predict location at time 3 from neural activity at time 1 (green circle). The same decoder predicts location at time  $t + 2$  (x-axis) from neural activity at any other  $t$  (lines). **(D)** Similarity of decoding patterns to idealised representations of future location in ‘relative’ or ‘absolute’ time (schematics).



**Figure S3: RNNs trained on simpler tasks do not learn spacetime representations.** In this figure, we analyse the representations of four RNNs trained on all combinations of (i) the static goal or the moving goal task, and (ii) reward input throughout the task ('continual') or reward input only during the planning phase ('working memory'). For all analyses in this figure, we only included trials where an optimal agent would intercept the goal in 3 to 6 actions. **(A)** Decoding accuracy for agent location at different times (x-axis) from neural activity at the end of the planning period. Decoders were trained in crossvalidation across the current agent location. Only the RNN trained on the moving goal task in a working memory setting seems to learn a generalisable representation of the future. This network is also unlikely to have learned a full STA, since it fails catastrophically on the reward landscape task (Figure 4F). Note that the decoding accuracy generally *increases* slightly for the true STA as a function of time-within-trial. This is because the navigation tasks have stronger correlations between consecutive positions, which leads to some degree of overfitting on the training data. This overfitting is less prominent later in a trial, where the space of possible locations conditioned on the current location is larger. **(B)** In this analysis, we trained a decoder to predict whether the agent would be at a particular location at *any* time in the trial from neural activity at the end of planning. Binary decoders were trained for each possible future location in crossvalidation across the current agent location. Bars indicate the average predictive accuracy across all binary decoders and current locations. The simpler networks seem to learn a representation of whether they will be at a given location at some point in the future.



**Figure S4: Parameters learned by the reward landscape RNN.** Network weights are projected into an orthonormal coordinate system with axes that maximally predict future locations. All weight matrices are for a single example RNN, since the environment differs between networks, and the connectivity is therefore slightly different. **(A)** Structure of the environment that the RNN was trained in, illustrated as the 0<sup>th</sup> order adjacency matrix (the identity matrix), the 1<sup>st</sup> order adjacency matrix, and the 2<sup>nd</sup> order adjacency matrix **(B)** Recurrent weight matrix estimated during the planning period, which shows structure resembling the environment adjacency matrix in the off-diagonal blocks. **(C)** Recurrent weight matrix estimated during the execution period, which shows an additional ‘feedforward’ component that copies information from later to earlier subspaces. We posit that this component of the connectivity matrix helps implement the conveyor belt dynamics identified in [Figure 4E \(Supplementary Note\)](#). **(D)** Input weight matrix estimated during the planning period. The ‘current’ subspace receives location input, and future subspaces receive reward corresponding to the appropriate time in the future. **(E)** Output weight matrix estimated during the execution period. The policy is read out from the ‘immediate future’ subspace as expected in a spacetime attractor. **(F)** Overlap between subspaces estimated during the planning and execution periods. This analysis was performed both for the standard RNN (‘WM’), and for a network trained with continual reward input throughout the trial instead of only during the planning phase (‘continual’; [Figure S6](#)). The WM RNN uses separate subspaces for computation of the plan and subsequent execution, consistent with the different connectivity patterns in (B) and (C). The continual RNN can use the same subspaces for planning and execution since it always receives the same type of input.

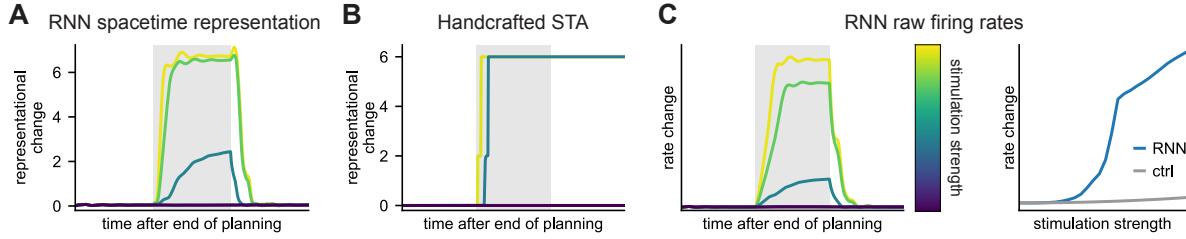


Figure S5: **Additional analyses of attractor dynamics.** **(A)** Change in implied spacetime representation over time in the trained RNN for different perturbation strengths (reproduced from Figure 5H). **(B)** Change in spacetime representation over time in the handcrafted spacetime attractor for different perturbation strengths. In contrast to the trained RNN, the ‘low value’ path is a fixed point of the perturbation-free network dynamics in the handcrafted network. At the end of a sufficiently strong perturbation, the representation can therefore stay in this new fixed point. **(C)** Change in RNN firing rates for different perturbation strengths. While the change in implied spacetime representation completely saturates with perturbation strength, the change in firing rates continues to increase with perturbation strength. This is expected because the network has a non-saturating ReLU nonlinearity. Small external perturbations are still quenched when quantifying the change in representation using the raw firing rates instead of the implied spacetime representation.

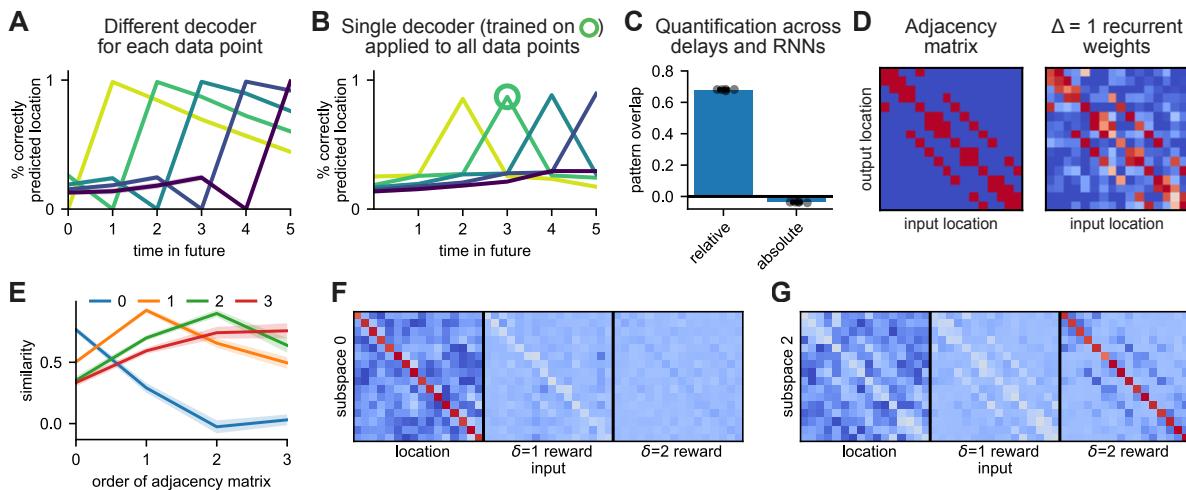
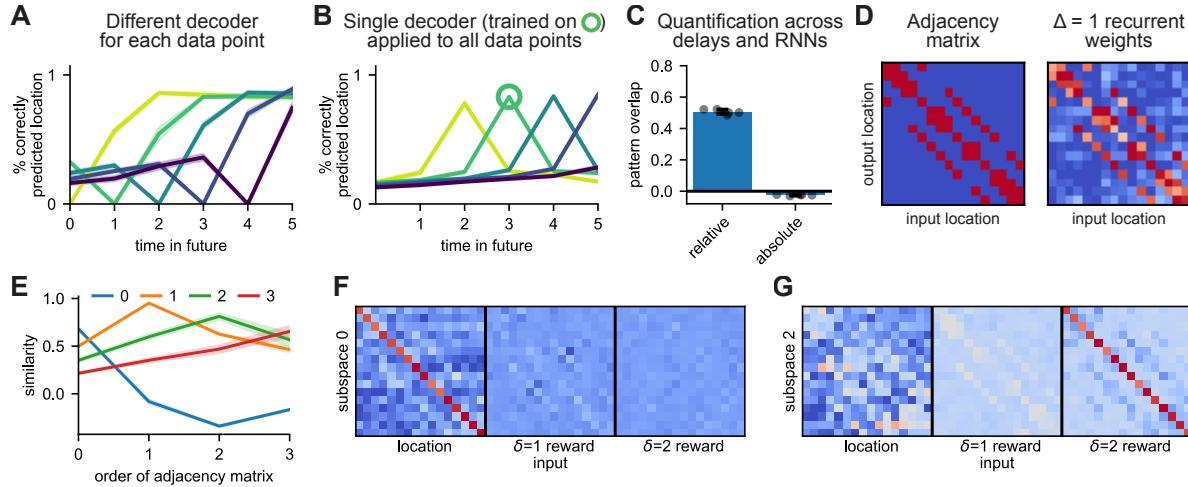
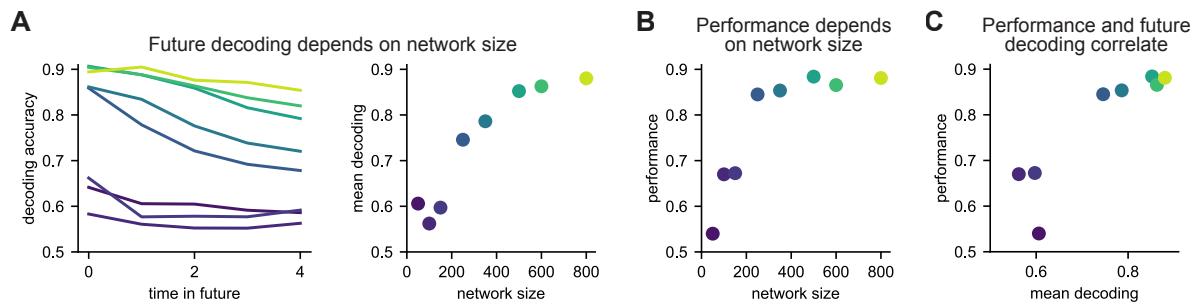


Figure S6: **RNNs trained with continual reward input also learn spacetime attractors.** In the main text, we focused on an RNN that was trained with reward input provided during an initial ‘planning phase’, while no information was given about the reward function during subsequent ‘execution’. In this figure, we perform some of the same analyses on an RNN that receives reward input throughout the entire task. In this setting, the task could in theory be solved using a ‘feedforward’ strategy that does not rely on recurrent dynamics at all. However, the RNNs still learn a spacetime attractor-like solution. **(A)** Decoding accuracy of agent location at different times (x-axis) from neural activity at every other time. Each line corresponds to predictions from neural activity at a different time in the trial from  $t = 0$  (yellow) to  $t = 4$  (blue). Decoders were trained in crossvalidation across the current agent location. **(B)** We trained a single decoder to predict location at time 3 from neural activity at time 1 (green circle). The same decoder predicts location at time  $t + 2$  (x-axis) from neural activity at any other  $t$  (lines). **(C)** Similarity of decoding patterns to idealised representations of future location in ‘relative’ or ‘absolute’ time (Figure S2). **(D)** The average recurrent weights between subspaces separated by a single action resemble the adjacency matrix of the environment. **(E)** Correlation between (i) the average connectivity between subspaces separated by  $\Delta$  actions (lines; legend), and (ii) different order adjacency matrices (x-axis). **(F)** Input weights to the ‘current’ subspace ( $\delta = 0$ ). **(G)** Input weights to a ‘future’ subspace ( $\delta = 2$ ).



**Figure S7: RNNs with a local action space also learn spacetime attractors.** In the main text, we analysed an RNN that generated a global allocentric policy, consisting of a probability distribution over all locations that was renormalised over ‘adjacent’ locations before sampling an action. In this figure, we perform some of the same analyses on a network that outputs a ‘local’ policy in an action space consisting of ‘north’, ‘south’, ‘east’, and ‘west’. This RNN also learns a spacetime attractor. **(A)** Decoding accuracy of agent location at different times (x-axis) from neural activity at every other time. Each line corresponds to predictions from neural activity at a different time in the trial from  $t = 0$  (yellow) to  $t = 4$  (blue). Decoders were trained in crossvalidation across the current agent location. **(B)** We trained a single decoder to predict location at time 3 from neural activity at time 1 (green circle). The same decoder predicts location at time  $t + 2$  (x-axis) from neural activity at any other  $t$  (lines). **(C)** Similarity of decoding patterns to idealised representations of future location in ‘relative’ or ‘absolute’ time (Figure S2). **(D)** The average recurrent weights between subspaces separated by a single action resemble the adjacency matrix of the environment. **(E)** Correlation between (i) the average connectivity between subspaces separated by  $\Delta$  actions (lines; legend), and (ii) different order adjacency matrices (x-axis). **(F)** Input weights to the ‘current’ subspace ( $\delta = 0$ ). **(G)** Input weights to a ‘future’ subspace ( $\delta = 2$ ).



**Figure S8: RNN representations and performance across network sizes.** We trained a series of RNNs with different network sizes ranging from 50 (dark blue) to 800 (yellow) hidden units. **(A)** Networks with approximately 300 or more units learned a spacetime representation, and the future could be decoded from the hidden state of the network at the end of the planning period. **(B)** Task performance saturated as a function of network size at approximately 300 hidden units. **(C)** Task performance increased with the ability of the network to represent the entire future explicitly. These results mirror the findings of Whittington et al. (2023) that RNNs trained on working memory tasks learn a similar ‘slot-like’ solution only if the network is large enough.

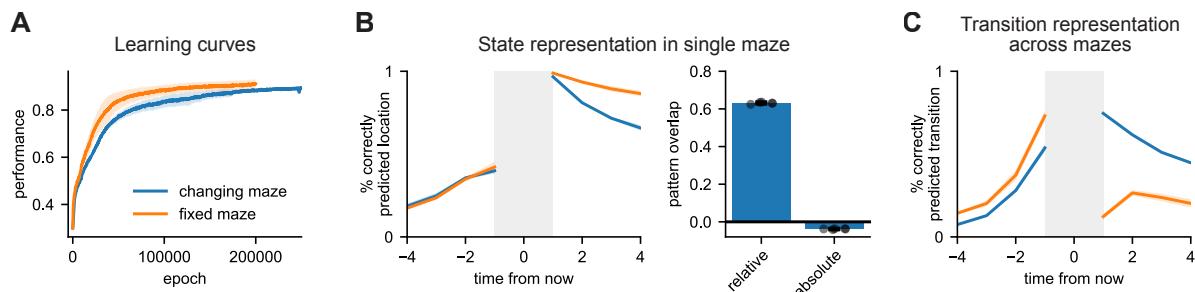


Figure S9: **Additional analyses of RNNs trained with changing environment structure.** (A) Learning curves of RNNs trained on the reward landscape task in a single maze ('fixed maze'; orange) or with a different structure on every trial ('changing maze'; blue). (B) We took the RNNs trained across many mazes and evaluated them in a single maze (blue). Future states could be decoded almost as well as in the RNNs trained in a single maze (orange). Additionally, the representations in each maze were more similar to idealised representations of future location in 'relative' than 'absolute' time (right). (C) When considering data across many mazes, future *transitions* could be decoded in a way that generalised across current transition and maze structure (blue). Such a representation did not exist in RNNs trained in a single maze (orange).

## 1311 Supplementary discussion of architecture and modelling choices

1312 In this section, we discuss some of the many architectural and modelling choices that went into our work.  
1313 As is the case for much work in modern computational neuroscience, the space of models was vast – and  
1314 larger than we could feasibly explore fully in a single paper. In what follows, we aim to further motivate the  
1315 choices that were made in the main paper and to provide additional intuition for the importance and effect of  
1316 various architectural choices and hyperparameters in our work. This note is also not exhaustive, but we hope  
1317 that it will be useful both for the reader looking to gain a deeper understanding of our work, and for those  
1318 who want to draw inspiration from it in their own research.

### 1319 Reward input format

1320 When training RNNs with a reward function that changes in time, there are multiple different ways this  
1321 could be provided as an input to the agent. The three most natural choices are:

- 1322 1. A relative encoding, where a given input channel always reflects the reward in  $\delta$  actions.
- 1323 2. An absolute encoding, where a given input channel always reflects reward at time  $t$ .
- 1324 3. A ‘working memory’ setting, where reward input is only provided during a planning period prior to the  
1325 first action. At this time, the relative and absolute representations are identical, and no choice has to  
1326 be made between the two.

1327 One might expect a relative reward input to bias RNNs towards relative future representations, and ‘conveyor  
1328 belt’ dynamics, and an absolute reward input to bias RNNs towards an absolute code. While RNNs trained  
1329 with relative reward input did show strong relative coding of the future (Figure S6), RNNs trained with an  
1330 absolute input learned somewhat mixed representations that were harder to interpret.

1331 To avoid having to choose between these two, most of the analyses in the paper were done on an RNN trained  
1332 in the ‘working memory’ setting. We also saw slightly stronger ‘spacetime’ representations in the working  
1333 memory setting. In fact, Figure S6 shows that the effective planning horizon of the ‘relative’ model is only 4-5  
1334 steps, and predictive performance of future states drops off beyond that. This is presumably because a shorter  
1335 planning horizon is sufficient for near-optimal performance in the task, and the regularisation encourages  
1336 networks to not encode more information than necessary. The RNN therefore effectively re-plans at every  
1337 step using a spacetime attractor with a depth of approximately 4.

1338 The handcrafted spacetime attractor always received a ‘relative’ reward input. Otherwise, it would need to  
1339 include either (i) an additional ‘memory’ component, (ii) a transformation from absolute to relative inputs,  
1340 or (iii) a mechanism for changing the readout between subspaces. We expect that if biological agents use  
1341 spacetime attractor-like dynamics, there may be settings where the code is relative and settings where it is  
1342 absolute. This might depend on whether the plan is shaped around constraints in relative time (e.g. having  
1343 to meet someone in 10 minutes) or absolute time (e.g. having to meet someone at 1 pm).

### 1344 Planning period

1345 When training the RNNs, we included a planning period prior to the first action. During this period, the  
1346 output of the network had no effect on the environment. We did this (i) to separate the ‘planning’ period  
1347 from the ‘execution’ period in the working memory setting, and (ii) because it slightly improved performance.  
1348 In the handcrafted STA, there is no distinction between planning and execution, since the inputs are the  
1349 same in both cases. However, convergence is slower before the initial action because the network state starts  
1350 further from a fixed point. This means that the STA can be run for fewer iterations after the first action,  
1351 which effectively corresponds to a longer ‘planning phase’ followed by faster execution.

### 1352 Network iterations per action

1353 For all RNNs, we used a variable number of planning iterations. This is because we were interested in stable  
1354 representations of future behaviour, rather than networks learning to time their dynamics to initiation. In the  
1355 RNNs with a relative reward input, we also varied the number of network iterations between each action for  
1356 the same reason. However, this is not possible in the working memory RNN, since it would not know when it

had taken an action. There are two possible solutions to this problem. (i) we can provide a specific ‘action’ input that tells the network when the environment has updated. (ii) we can use a fixed number of network iterations between every action and allow the network to learn the speed of the environment. We opted for the second option to keep the inputs as simple as possible. We suspect that this could lead to changes in the subspaces used at different ‘phases’ of the dynamics, but we did not test this explicitly. Instead, we focused on the planning period when analysing the connectivity of the working memory RNNs to circumvent this potential complication. When analysing the activity during ‘execution’ in Figure 4D, we used only neural activity right before each action.

### Execution-time dynamics

The fixed points of the STA dynamics are input-dependent, and the inputs to the handcrafted model change when it takes an action, because the location and/or future reward will be different. For this reason, the representation can automatically update to reflect the ‘path-to-go’ after each action. However, this involves recomputing the future to some extent, which is ‘wasteful’ since it has already been computed once. Additionally, the representation can get stuck in local minima because the barrier to switching between representations is non-zero. We therefore included an additional component in the STA weight matrix, which was 1 for all off-diagonal elements that implement a transfer of information to earlier subspaces. This transfer component was only active for 100 network iterations (two time constants) after each action to update the representation before settling into a new fixed point. This is similar to how the fruit fly head direction circuit uses populations of ‘shift’ (PEN-1) neurons to update the internal heading, and these PEN-1 neurons are gated by angular velocity input (Turner-Evans et al., 2017, 2020). Including an explicit shift component in the STA is not necessary for any of the main results in the paper, but it stabilises the dynamics to some extent. We posit that the RNNs learn something similar, as suggested by the feedforward component of the ‘execution period’ parameters in Figure S4C.

### Noise sensitivity

The handcrafted STA is more or less susceptible to different types of noise. The model is very robust to noise added to the neural potential (‘ $z$ '; Methods), which can be interpreted as a spacetime distribution in log probability space. This is because non-desirable locations in spacetime can have very negative values, which are not very noise sensitive. The STA is more sensitive to noise added to the firing rates (‘ $r$ '; Methods), which can be interpreted as a spacetime distribution in actual probability space. This is because addition of positive noise to some location  $s_i$  in subspace  $\delta$  will propagate through the adjacency matrix to all neighboring locations at time  $\delta + 1$ . This can lead to representations that ‘teleport’ between distant locations if  $s_i$  is near the reward location and therefore receives strong input from the reward function or future subspaces.

The STA is sensitive to structural noise for a similar reason. Adding a small positive value to weights corresponding to elements of the adjacency matrix that are meant to be zero leads to dynamics that include a finite probability of teleporting between these distant locations. In this work, we mitigate the structural noise sensitivity by using a ‘pessimistic’ estimate of the adjacency matrix as the base weights before adding noise (Methods). In other words, we subtract a small constant from all weights to ensure that distant locations are connected with weights that are zero or negative, even though they are noisy. This may be less of a problem in biological networks, since Dale’s law ensures that synapses are either excitatory, inhibitory, or absent.

The sensitivity to both rate noise and structural noise is higher when the input strength is larger ( $\beta$ ; Methods), which biases the representation more strongly towards rewarded locations. Conversely, the converged representation is more diffuse if the input strength is weaker, because there is a smaller bias towards rewarded locations. This is particularly true when planning towards distant rewards. The strength of the reward input therefore has to balance the planning depth with the susceptibility to teleporting representations. If there is no noise in the system, the input strength can safely be very large, which leads to robust performance for long planning horizons. If there is more noise in the system, the input strength should be smaller, which increases noise robustness at the expense of a shorter effective planning horizon. When training RNNs across tasks, they will naturally learn to balance the robustness of the representation with the required planning depth. Indeed, the analyses in Figure S5 suggest that RNNs learn to do so better than our handcrafted models.

1406 **Choice of RNN learning algorithm**

1407 We used supervised learning to train all RNNs in this paper. In other words, the RNNs were trained to  
1408 predict the behaviour of an optimal agent. An alternative would have been to train the networks end-to-end  
1409 using reinforcement learning. We opted for the supervised setting for two reasons. Firstly, supervised learning  
1410 is more stable, which leads to more robust results that are less sensitive to hyperparameters and use less  
1411 compute. These features make it much simpler for others to build on our work. Secondly, we do not think  
1412 cortical representations are learned from scratch via reinforcement learning. Instead, we are of the opinion  
1413 that cortical representations are likely learned via predictive or ‘semi-supervised’ learning. The basal ganglia  
1414 can then use reinforcement learning to map these representations onto actions (Blanco-Pozo et al., 2024;  
1415 Zintgraf et al., 2019). We expect that training the RNNs with reinforcement learning would yield similar  
1416 results, but we have not tested this explicitly.

1417 **Convergence of RNN training**

1418 Convergence was in general fairly good, but it did vary with hyperparameters to some extent. For some com-  
1419 binations of regularisation strengths, the networks sometimes got trapped in local minima that corresponded  
1420 to incomplete structural learning. These networks would partially learn the structure of the environment,  
1421 but they would fail to learn connections between some pairs of states that were actually connected, which  
1422 impaired performance. Similarly, some RNNs trained on the changing maze task erroneously ‘hard coded’  
1423 some transitions instead of having them flexibly modulated by the inputs. In general, the accuracy of the  
1424 learned structure was correlated with model performance for both the networks trained in a single maze and  
1425 networks trained in changing mazes. Additionally, the overall loss was higher for the networks that failed to  
1426 learn the full task structure, suggesting that it is an issue of convergence rather than regularisation making  
1427 the partially learned solution optimal.

1428 **Long distance parameters in the RNN**

1429 In Figure 5 and Figure S4, we saw that the trained RNNs learn some long-range connections between  
1430 subspaces separated by more than one action. This differs somewhat from the handcrafted STA, which only  
1431 has connections between adjacent subspaces. However, the presence of long-range connections in the RNN is  
1432 not too surprising. In particular, we expect this additional structure to stabilise fixed points corresponding to  
1433 possible trajectories, since it inhibits any trajectory that includes impossible n-step transitions. We suspect  
1434 that stabilising the dynamics through weaker connectivity between all subspaces is cheaper than strong  
1435 connectivity exclusively between adjacent subspaces. This may be a consequence of the L2 regularisation used  
1436 to train the networks, which favors many small parameters over few large parameters. In future work, it could  
1437 be interesting to explore whether the degree of long range connectivity is lower when using L1 regularisation  
1438 instead. The result of such an analysis may also depend on how disentangled the spacetime representations  
1439 are, which we did not explore in the present paper.

1440 **Discrete space and time**

1441 We have discretised space and time throughout this work. This makes the models and analyses simpler,  
1442 because the action space is discrete and one action always leads to a step change in the environment. However,  
1443 we expect that the basic ideas extend to continuous space and time as well. In this case, neurons would still  
1444 represent particular points in spacetime. Pairs of neurons would be connected as a function of their difference  
1445 in preferred location in a way that reflects which locations can be reached in a unit time. This is similar to  
1446 the mechanism used for angular velocity integration in ring attractors and path integration in grid attractors.  
1447 If the speed of the agent can vary, it might be necessary to represent different speeds in different connections  
1448 or neurons, and the desired speed at every point in time could be inferred together with the trajectory.

1449 **Stochastic environments**

1450 We have worked with deterministic environments throughout this paper. This choice greatly simplifies the  
1451 spacetime attractor, since the deterministic adjacency matrix can be built into the network connectivity. When  
1452 working in a stochastic environment, we would intuitively want the connections to represent  $\max_a p(s_{t+1}|s_t, a)$ ,  
1453 which is a generalisation of the adjacency matrix for deterministic environments. We have not explored this  
1454 explicitly but consider it an important extension for future work. One potential challenge in the stochastic  
1455 case is that the spacetime attractor as formulated here effectively performs planning as inference under the

1456 assumption that the posterior distribution over locations factorises across time. This assumption may have  
1457 more severe consequences in stochastic environments, where correlations are more important. In particular,  
1458 this assumption generally leads to a collapse to a ‘modal’ representation of a single trajectory rather than  
1459 representations of entire distributions of trajectories. In stochastic environments, that means it is possible to  
1460 converge to a single representation of a trajectory that is unlikely to happen in reality, even with the correct  
1461 choice of actions.

#### 1462 Relationship to sequence working memory

1463 The spacetime attractor is strongly inspired by representations identified for sequence working memory  
1464 ([El-Gaby et al., 2023](#); [Xie et al., 2022](#); [Chen et al., 2024](#); [Tian et al., 2024](#); [Whittington et al., 2023](#)). These  
1465 sequence memory tasks have notable similarities to the reward landscape task studied in this work. In  
1466 particular, sequence working memory often involves presenting an animal with a sequence of 2-4 options  
1467 (e.g. ‘a’, ‘b’, and ‘c’) that are sampled from a set of N possible states or actions ([Xie et al., 2022](#); [El-Gaby](#)  
1468 [et al., 2023](#)). The animal is then rewarded for repeating the sequence at response time. The task therefore  
1469 has a reward function that is categorical (only one element is rewarded at a given moment in time) and  
1470 changes in time (first ‘a’ is rewarded, then ‘b’...). This is very similar to the changing reward function in the  
1471 reward landscape task. A notable difference is that the sequence memory tasks have no constraints on the  
1472 possible transitions – any element can be produced before or after any other element. For this reason, each  
1473 sequence element can be treated independently, and there is no need to pass reward or value information  
1474 between subspaces. The independent sequence working memory task can therefore be seen as a special case  
1475 of the reward landscape task, where (i) only one state is rewarded at each point in time, and (ii) any state  
1476 can be reached from any other state, so the adjacency matrix is uniform. Interestingly, theoretical work  
1477 shows that explicit representations of the future preferentially emerge for sequence working memory when  
1478 correlations (or more precisely, ‘range dependence’) between subsequent elements are weak, and the space of  
1479 possible sequences is therefore large ([Dorrell et al., 2024](#)). This is reminiscent of our finding that spacetime  
1480 representations for planning preferentially emerge in RNNs trained on the reward landscape task, where the  
1481 space of possible optimal trajectories is larger than in the static and moving goal tasks.

#### 1482 Comparisons with TD and SR agents

1483 In [Figure 3](#), we compare the representations and performance of the STA to temporal difference learners and  
1484 successor representation agents. We show that the STA solves ‘dynamic’ problems that these algorithms  
1485 struggle with. This is in some sense a property of the state space rather than the decision making algorithm.  
1486 It would be possible to construct TD and SR agents in a ‘space-by-time’ state space, which would allow them  
1487 to solve these dynamic problems. Our message is not that this is not possible. Instead we are highlighting  
1488 that the way these algorithms are usually implemented involves representations that are ‘flat’ across time,  
1489 and we use them as a comparison to show why spacetime representations can be useful. If a TD learner was  
1490 implemented with a space-by-time state space, it would be able to solve tasks where the reward changes  
1491 within a trial but the same pattern is seen across all trials. The SR agent with a space-by-time state space  
1492 could solve the general reward landscape task. In the simplest implementation, this would require inversion  
1493 of a matrix  $\mathbf{T} \in \mathbb{R}^{NT \times NT}$ , which is computationally expensive. However, it is possible that the structure of  
1494  $\mathbf{T}$  could be exploited to invert it more efficiently, which would be an interesting avenue for future research.

#### 1495 Summary

1496 As is evident from this discussion, many choices went into this work that could have been different. We  
1497 do not claim to have explored the full space of models and tasks, and we are not trying to argue that the  
1498 spacetime attractor is a silver bullet for planning and decision making. Instead, we have tried to argue that  
1499 STA-like models are interesting solutions to a range of problems that are relevant to prefrontal cortex and not  
1500 widely studied in systems neuroscience. However, many open questions remain, some of which we have briefly  
1501 motivated here. We therefore hope that this paper will inspire others to further explore these questions both  
1502 experimentally and computationally.