## B. SIGNIFICANCE

Uncertainty about the state of the world arises due to ambiguous sensory information, which prevents the brain from inferring the external events that caused each pattern of sensory activity with certainty. Probabilities are therefore the most suitable framework for formalizing this process (Jaynes, 2003) and for developing hypotheses about its brain mechanisms. In this project, we propose to use normative modeling to generate predictions about behavior and neural mechanisms underlying causal inference of motion and depth during animals' self-motion in the world, which we will test in neural population recordings of Projects B and C.

Our models will build on the seminal work of Körding et al. (2007), who used Bayesian decision theory to explain human behavior in a simple audio-visual integration task, laying the foundation for research in causal inference. Previous work showed how vestibular signals and visual optic flow are combined to infer self-motion (Acerbi, Dokka, Angelaki, & Ma, 2018), and how humans infer object motion in the presence of vestibular signals (Dokka, Park, Jansen, DeAngelis, & Angelaki, 2019). Yet no study has modeled the causal inference underlying segmentation of a visual scene into self-motion versus object motion in the world, nor how inference about an object's depth from binocular cues relates to inference about its motion.

We propose to extend current theory in multiple ways. In Aim 1, we will develop a probabilistic model of the causal inference process that segments a visual scene into elements that are stationary in the world and those that are moving. In Aim 2, we will focus on dynamic causal inference and control, namely how we would expect the latent variables underlying causal inference to evolve over time, given a stream of sensory evidence, and how these latent variables in turn might support sequences of goal-directed actions. In both aims, the parameters underlying these models will reflect the animals' assumptions about the world, and we will tune them to best match the observed behaviors. Furthermore, we will compare their predictions to those of simplified models, to see if such simplifications better account for behavior. The latent variables predicted by the best-fitting models will provide a basis for data analysis to identify their neural correlates (Brunton, Botvinick, & Brody, 2013; Gold & Shadlen, 2007), and thus the neural circuits involved in causal inference.

**B.1. Significance of Aim 1**. Optimal motion processing involves at least three key elements: (1) computing object motion while simultaneously inferring self-motion, (2) transforming between different coordinate systems (e.g., retinal vs. world-centered), and (3) inferring joint posterior beliefs that respect the geometric relationships between motion and depth. Some of these computations have been studied psychophysically in isolation (Rushton & Warren, 2005; Warren & Rushton, 2007, 2009) but their relationships and the neural basis of these computations are mostly unknown. In Aim 1, we will develop a unified, normative causal inference model that can make principled predictions for the neural responses representing these computations. Our model will link three variables—object motion, self motion, and depth—that may be computed and represented in different neural populations and even brain areas, allowing us to determine how the brain's beliefs about these variables influence each other and propagate between different parts of the brain.

**B.2. Significance of Aim 2**. Like most phenomena in perception, causal inference is a dynamic process that evolves over time, but it has been studied mainly in a simplified trial-by-trial structure in which this temporal dynamic can be ignored (e.g., Dokka et al., 2019; Körding et al., 2007; Shams & Beierholm, 2010). Previous work on the dynamics of causal inference has been limited. The neural time-course of human causal inference (Rohe, Ehlis, & Noppeney, 2019), for example, was studied with a static, potentially confounding model. To link causal inference to neural circuit dynamics and investigate its role in the dynamic control tasks of Project C, we must explicitly take time into account. In Aim 2, we will design models that describe such causal inference in continuous time, and how its outcome in turn affects a continuous-time control policy to achieve a desired goal. In contrast to other heuristic dynamical models (Daemi, Harris, & Crawford, 2016), we will rely on the same successful normative basis as Körding et al. (2007) and extend it to operate in real-time.

The benefits of such a real-time extension are manifold, as real-world tasks are always dynamic. When applied to tasks with fixed-duration stimuli, our models can be reduced to previous models (e.g., Acerbi et al., 2018; Dokka et al., 2019), with their predictions as special cases. However, we can additionally ask questions like how the belief about object motion is expected to evolve over time (see **Fig. 4**) to support inference in dynamic experiments, such as those in Project C. Lastly, we can evaluate temporally predictive models in the brain. For example, some of our dynamic experiments require planning a path toward an obscured object that is potentially moving. Our dynamic model of inference and control will allow us to determine whether the brain encodes only the expected point of interception or dynamically updates the predicted object location in real time. All of these properties make our dynamic model an indispensable foundation for moving the study of causal inference toward dynamic real-world scenarios.

## C. INNOVATION
- We will develop unified probabilistic inference models that can explain psychophysical data and generate hypotheses for the underlying neural representations and computations in cortical motion areas.
- Our models will perform joint causal inference over self-motion, object motion, and object depth.
- Our static model makes strong predictions about how neural responses depend on stimulus context in the tasks of Project B and how responses of neurons representing object motion, self-motion, and depth should be correlated, enabling principled and systematic tests of neural coding in visual cortical areas.
- Our dynamic models will strictly generalize normative causal inference models to continuous-time, dynamic inference, as required in real-world tasks, letting us study the neural circuit dynamics of causal inference.
- Our dynamic models will embed causal inference within real-time control to identify the neural correlates of control-related continuous predictions based on sensory inputs.

## D. APPROACH
### D.1. Aim 1: Develop a causal inference model for motion and depth perception during self-motion

Procedure. *Generative model of sensory inputs*

In preliminary work, we developed a causal inference model (**Fig. 1**) for motion processing that is applicable to the task in Project B, Aim 1. We describe the generative model for the subject's sensory input under the simplifying assumption that the visual input consists of a discrete number of elements (here, dots). Inference in this generative model yields a posterior belief over all unobserved variables (object motion, self-motion, and depth), given an observation of the sensory input (stimulus). The self-motion vector ($v_{self}$) produces the subject's vestibular input ($o_{vestibular}$, non-visual, received from the vestibular system and quantified in terms of equivalent self-motion). Together with the distance to each dot ($d^{dot;k}$), it also determines the retinal motion of each dot (indexed by k) if the dot were stationary ($v^{dot;k;stationary}_{retina}$). The actual retinal motion ($v^{dot;k}_{retina}$) of dot k is then given by combining what it would be purely due to self-motion with a relative velocity ($v^{dot;k}_{relative}$) corresponding to the dot's velocity in world coordinates projected onto the



**Figure 1. Generative model of sensory inputs.** Gray boxes indicate variables potentially representing subjective percepts.

retina (i.e., without self-motion). The square box around the variables (**Fig. 1**) indicates multiple plates, representing all the dots in the display. We model the conditional dependencies (arrows) on sensory observations as Gaussian distributions defined by mean and standard deviation representing sensory noise. Standard geometry relates distance (d) to $v^{stationary}$.

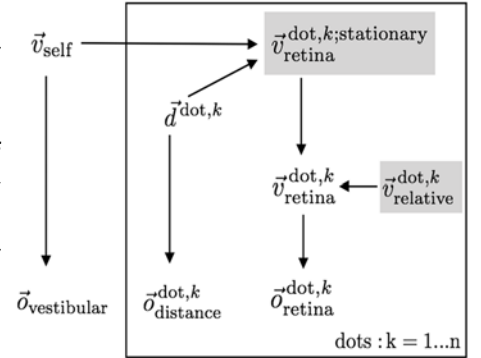To make this a causal inference model after Körding et al. (2007), the key is that the prior distributions over $v_{self}$ and $v_{relative}$ are mixtures of a delta-distribution at zero and a Gaussian distribution around zero (marginalizing out the causal variable "C" in their model). As a result, inference in this model amounts to causal inference of whether the subject is stationary (with a Gaussian slow speed prior if not, as in Stocker & Simoncelli (2006)), or whether an individual dot is stationary in the world (with a Gaussian slow speed prior if not). These inferences are mathematically equivalent to inferring whether the difference between visual and auditory location was zero (indicating a shared cause). Both priors in our model are motivated by the natural statistics of these variables: most objects, including observers, are stationary, and the others are more likely to be moving slowly than fast.

Without vestibular signals and with many dots, this model infers a belief about self-motion such that the overall world motion of all dots is zero. Self-motion trades off with dot motion probabilistically and, when a vestibular signal is present, infers a self-motion belief that combines all cues and accounts for their uncertainty.

Assuming that motion relative to the world ($v^{dot}_{relative}$) underlies the subject's percept, this model can explain data from classic flow-parsing experiments, in which the perceived direction of an object's motion is strongly and predictably influenced by background optic flow. When a subset of dots makes up the optic flow and a single dot (or subset of dots) represents motion at some location on the retina, this dot's motion is not perceived as $v^{dot}_{retina}$ but approximately as $v^{dot}_{retina} - v^{dot,stationary}_{retina}$. This subtraction of local motion implied by self-motion from local retinal motion is central to the flow-parsing hypothesis (Rushton & Warren, 2005). However, because $v_{self}$ is generally smaller than the self-motion implied by the optic flow field in isolation (primarily due to a vestibular signal implying zero self-motion in most experiments), $v^{dot,stationary}_{retina}$ is scaled down below what would be expected based on optic flow alone. Exactly these results were found in classic experiments in which the subtraction appears incomplete (Warren & Rushton, 2007, 2009). Furthermore, our model predicts that the bias due to optic flow decreases with the number of dots in the optic flow, also as observed by us (data

not shown) and others (Warren & Rushton, 2009). Finally, as our model's perceptual bias is mediated by self-motion rather than retinal motion directly, this bias generalizes across the visual field even when optic flow is restricted to parts of it, as has also been shown (Warren & Rushton, 2009).

　　　Our model goes beyond existing flow-parsing studies by using causal inference to infer whether the movement of a visual element should be flow-parsed in the first place. When the retinal motion of a visual element deviates only slightly from what would be expected if it were stationary, then the brain should combine its observation with the expectation from stationarity—a process that can be understood as cue integration of local and global motion. The behavioral signature of this effect is increased uncertainty for a range of motion that is equally well explained by object motion or stationarity corrupted by sensory noise. We observed this increased variability in reports for patch velocity angles of 5-20 degrees with respect to the optic flow (**Fig. 2**).

Relation to experiments. Aim 1 of Project B will investigate the neural basis of this hallmark of causal inference: transition from cue integration to segregation (flow parsing). Our model will make quantitative predictions for how posterior beliefs change in mean and shape (e.g., variance), allowing us to search for their correlates in neural responses (see Overall Strategy for details). For instance, based on fits to behavioral data, our model will quantitatively predict the strength of the anticorrelation between decoded self-motion and object motion.

Aim 2 of Project B will test a specific dependency between two key variables believed to be represented by MT neurons: depth and object motion. Due to their geometric relationship imposed by the 3D structure of the world (captured by the model in the arrows between $d^{dot;k}$ and $v^{dot;k,stationary}_{retina}$, and between $v_{self}$ and $v^{dot;k,stationary}_{retina}$), shifts in the posterior belief over one variable should align with shifts in the posterior belief over the other variable. For example, this dependency predicts that depth and object motion as decoded from measured neural responses should also be correlated, along with the individual neural responses (as a function of each neuron's coupling coefficient with each variable). In preliminary human behavioral data, we



**Figure 2. Human psychophysical data support causal inference modulation of flow parsing.** The subject reports motion direction of a dot with (red) and without (blue) surrounding optic flow. Individual direction estimates often lie between the predictions for flow-parsing (green) and forced integration (orange). Black dotted line: unity-slope diagonal. Note increased report variability for intermediate directions (5-20 deg), as expected by the causal inference model.

observed this dependency (Fig. 7, Overall Strategy, supporting our model and justifying our top-down, normative approach.

*Decision rule.* Additional assumptions are necessary to relate the posteriors in this generative model to decision behavior. We will consider optimal decision-making (reporting the most probably correct choice), probability matching, and schemes intermediate to both extremes (Shivkumar, DeAngelis, & Haefner, 2019).

**Aim 1.1:** *Convert the dot-based model to a texture-based model*
While our preliminary model describes a world of dots, neurons in dorsal cortex are thought to represent particular areas of the visual input (receptive fields) that are approximately stationary in space. To more directly link the variables in our model to neural data, we will convert our model to a texture-based one, in which individual variables represent the motion energy within receptive fields that tile the visual field. We will follow Gershman et al. (2016) in modeling the motion texture in each receptive field as a Gaussian process.

**Aim 1.2:** *Infer model parameters from psychophysical data*
We will apply methods (Bayesian adaptive directive search and variational Bayesian Monte Carlo) that we have used with other causal inference models to fit our generative model to monkey behavioral data (Acerbi, 2018; Acerbi & Ma, 2017). These methods yield posteriors over model parameters, which will allow us to confirm that monkeys perform causal inference as expected from our preliminary human data (**Fig. 2**).

**Aim 1.3:** *Generate quantitative neural hypotheses*
We will compute approximate joint posteriors over all latent variables in our model given the sensory inputs on each trial for the tasks in Project B using standard Markov Chain Monte Carlo sampling methods. These posteriors, in particular their means and (co-)variances, will provide single-trial predictions for the latent states that may be represented in neural recordings. Our use of inference by sampling does not assume that the brain uses sampling, as it will also allow us to extract any parameters of the resulting distribution (e.g., variances or
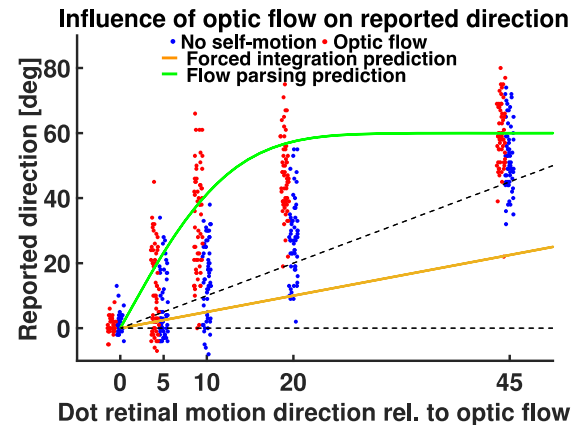
covariances) to test whether they can be decoded from neural responses, and whether they improve the predictability of neural responses beyond the posterior means.

Expected results. We will approximate the internal model of each monkey based on fits to its own behavioral data in the motion tasks proposed in Project B to produce trial-by-trial joint posteriors over the variables in this internal model, and hence hypotheses for the latent states in the neural activity from cortical motion areas. We will predict sampling-based (posterior variability should relate to neural variability) and parametric (posterior variability should relate to linear projections of mean neural activity) representations. These predictions are important not just to reveal the neural basis of causal inference in this context but also to test whether parametric or sampling-based models better describe neural representations.

Potential challenges and alternatives. Psychophysical data might not sufficiently constrain the generative model parameters, leading to large uncertainties about the hypotheses. We believe this outcome to be unlikely due to the large number of trials we will obtain from each monkey. Also, the resulting model fits might not describe the data well (e.g., in variance). Such a finding would be interesting because it would suggest a deviation from optimality in the task. In that case, we would look for minimal modifications of our model to fit the empirical data to quantify those deviations and generate hypotheses for neural data. Finally, we might have unanticipated problems converting the dot-based generative model to a texture-based one. In that case, we would use the dot-based model, which we found to explain behavior well in preliminary data. Normatively, the dot-model is equivalent to a texture one in which there is at most one dot in a neuron's receptive field.

## D.2. Aim 2: Develop a model of dynamical causal inference for perception and control

Rationale. Standard causal inference models make the simplifying assumption that all sensory information is available instantaneously. However, in a brief time interval, the stream of sensory input provides only limited information, which needs to be accumulated over time to form reliable estimates (Drugowitsch, Deangelis, Klier, Angelaki, & Pouget, 2014; Drugowitsch, Moreno-Bote, Churchland, Shadlen, & Pouget, 2012; Gold & Shadlen, 2007). Using causal inference to determine if an object is moving requires an estimate of our motion relative to this object. Certainty for some variables depends on how long the visual scene is examined, such as when self-motion direction and speed are estimated from visual flow. Furthermore, self- or object-motion might change with time. If these variables evolve over time, so will their neural correlates. Therefore, we propose to construct a dynamic model for perception and control to identify the neural population dynamics underlying continuous behavior (**Fig. 3**).

**Figure 3. Dynamic causal inference model.** The model simultaneously tracks an estimate of self- and object motion under the assumption of a moving object (blue shading), and an estimate of self-motion under the assumption of a stationary object (red shading), along with the likelihood of each alternative. It updates these estimates dynamically over time, given noisy momentary sensory evidence, and supports changes of object state between moving and stationary (dashed arrows).

We will develop this model in two stages. In Aim 2.1, we will describe the dynamic computations underlying the estimate of self- and object-motion from uncertain perceptual evidence. To do so, we will leverage our tight experimental control over the animal's percepts, whether the animal was a passive observer (Project B) or controlled its motion trajectories (Project C). In Aim 2.2, we will model goal-based rational control, which aims to identify a continuous stream of actions that minimize the cost of control while maximizing overall reward. This model (**Fig. 4**) will close the perception-action loop by determining rational control strategies that predict how behavior depends on the current percept. From observed behavior, we will reverse-engineer the percepts that led to the behavior. The resulting control strategies will predict the temporal evolution of the underlying variables to support identification of their neural correlates.
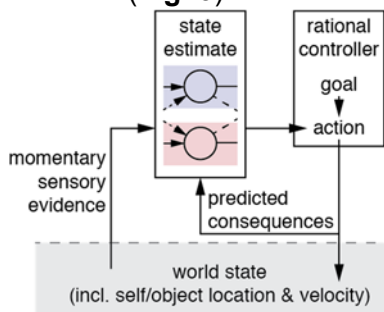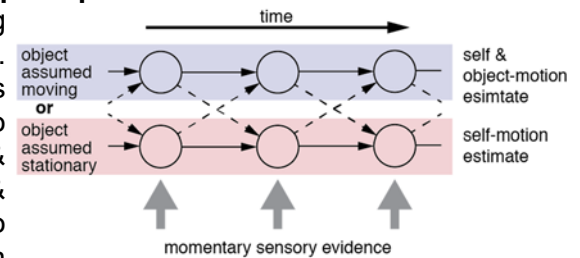
**Figure 4. A rational controller closes the perception-action loop**. This controller finds a policy that turns current state estimates into actions. Actions and natural dynamics update the world state, and the resulting momentary sensory evidence is used to update the state estimate. Real-time state estimation (**Fig. 3**) is required for rational control.

*Aim 2.1: Model dynamic causal inference for perception*
Procedure. We will develop our model for dynamic causal inference by combining continuous-state state-space models with causal inference. Such models are used to infer some latent state that evolves over time in response to noisy inputs. In our case, this latent state consists of self- and object-motion velocities. The

inputs are sensory percepts that, at each point in time, yield noisy measurements of these velocities. State-space models provide the best current velocity estimates given past sensory evidence, but commonly do not differentiate between stationary and moving objects. To support this distinction, we will merge them with causal inference models. The resulting model tracks two latent states—one assuming a stationary and the other a moving object—together with the likelihood of each state given sensory evidence so far (**Fig. 3**) to support real-time causal inference from noisy momentary information.

In most proposed tasks, neither self- nor object-motion changes during passive viewing within a single trial. That is, if an object is stationary or moving at the beginning of the trial, it remains stationary or moves at the same velocity throughout the trial. In the real world, in contrast, self- or object motion may not be constant over time. Furthermore, the causal inference machinery innate to mice and monkeys might take this possibility into account even during tasks with immutable object states. Thus, we will include possible state changes in our state-space model. Specifically, at each point in time, the model will assume that self- and object-motion velocities might change slightly, according to a physically realistic process, and that the object might start moving, or stop if sufficiently slow. The dynamics of these events will be described by parameters that can be tuned and reflect the animals' assumptions about the world. The assumption of an immutable object state will correspond to specific settings of these parameters. This model extension makes exact inference intractable, so we will use projection filtering (a continuous-time variant of assumed density filtering; (Brigo, Hanzon, Gland, & Gland, 1999; Murphy, 2012)) to approximate inference. Overall, this model extension will allow us to estimate animals' belief states even if there is a mismatch between their world model and the conditions imposed by the tasks.
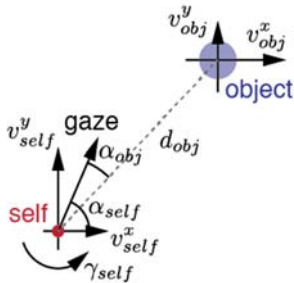


**Figure 5. Three trials of dynamic causal inference with different object velocities.** The observer estimates self- and object-motion and belief that the object is moving from a stream of noisy sensory evidence of self-motion velocity and object-motion velocity relative to self. A slow-moving object leads to extended uncertainty about object motion and estimated self-motion velocity.



**Figure 6. Sensory inputs and variables tracked by full model.** This model estimates self-motion velocities, gaze direction, object-motion velocities, and distance, based on noisy sensory inputs of self-motion velocities, yaw (visual flow field), and relative object direction and distance.

We demonstrated dynamic causal inference in a simplified proof-of-concept model (**Fig. 5**) with two noisy one-dimensional sensory inputs: self-motion velocity and object-motion velocity relative to oneself. The model estimates one-dimensional self- and object-motion velocity from these inputs, and whether the object is stationary or moving. The full dynamical causal inference model will be more complex (**Fig. 6**), tracking self-motion velocity, gaze direction, object distance, direction, and motion velocity in two planar directions, along with whether the object is moving. It will do so based on visual flow field information about self-motion in two planar directions, yaw, and object distance and current direction relative to self. We do not need to assume a percept of object velocity, as a change of the object's direction over time implicitly provides this information. In the simplified model, all variables of interest are linearly related. In the full model, we will handle non-linear variable relationships by projection filtering. For tasks that only provide some of this sensory information (e.g., task without yaw), we will make the absent sensory modalities uninformative. We will also design simplified models that replace some tracked uncertainties with model parameters, and will yield alternative hypotheses. The resulting dynamic causal models will be able to describe real-time perceptual causal inference in all of the proposed tasks.

Expected results. This aim will produce models for dynamic causal inference that take a stream of sensory observations as input, and at each point in time provide a posterior distribution (or point estimates) over one's own state and that of the object relative to self. This posterior will include estimates of self- and object-motion, as well as the belief about whether the object is stationary or moving. The models thus will predict at all times all the variables that the animals' brains need to track to perform the various proposed tasks. They will allow us to predict goal-directed behavior that relies, by task design, on these variables. The exact predictions will depend on the parameters that describe the animal's assumptions about the dynamics of the world, which we will tune to best explain the animal's behavior. Different model simplifications will yield different predictions, which we will distinguish by model comparison.

Relation to experiments. To model the latent state dynamics in Aims 1 and 2 of Project B, we will use reduced models with only lateral self- and object motion, as well as object distance. Furthermore, we will make the object
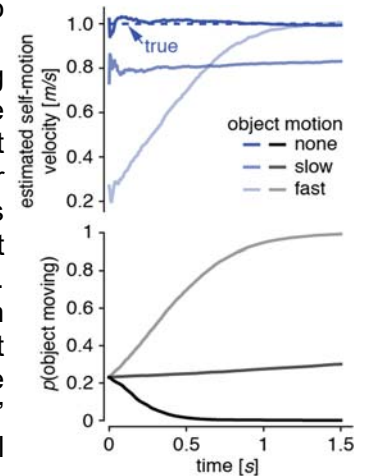
state immutable (i.e., within a trial, the object continues moving or remains stationary). These models will then predict how latent state estimates evolve over time and the animal's final choice at the end of each trial. We will adjust model parameters to best explain the animal's choices.

For Project C, we will use the full model with all degrees of freedom, and will consider both the immutable and mutable object models, to see which one explains the behavior better. Given the stimulus sequence provided to the animal, and the observed actions, the model will predict all the estimates the animals need to track to perform the tasks. These include one's own lateral and forward velocity and yaw, object's distance relative to self, and object velocity. The object's distance will be tracked, or predicted when it is occluded. Relating model estimates to behavior, and adjusting the model's parameters, will require the rational control framework developed in Aim 2.2.

Potential pitfalls and solutions. If our initial assumption of immutable object states cannot explain behavior, we will move to mutable object states and tune the model parameters that reflect the animal's assumptions about the world. If this approach also produces unsatisfactory behavioral fits, we will analyze prediction errors and re-design the model assumptions to correct them. The proposed inference by projection filtering is approximate. If this approximation results in significant errors, we will switch to particle filters, which can be made arbitrarily accurate by increasing the number of particles, at the cost of increased computation time.

*Aim 2.2: Develop a model of dynamic causal inference for rational control*
Procedure. Building on Daptardar et al. (2019), we will use an artificial neural network controller to identify a family of rational controllers that generalize across our tasks. The complexity of our state space precludes the use of simpler controllers, like linear-quadratic-Gaussian control. Thus, we will leverage progress in optimal control for complex problems (e.g., the Atari-game playing Deep Q-Learning; Mnih et al. (2015)), and use the deep deterministic policy gradient algorithm (Lillicrap et al., 2016) to identify rational control policies. These policies will be used to reverse-engineer the animal's percepts, by asking which sequence of percepts and resulting state estimates are most compatible with its actions (akin to Kalman smoothing). This procedure effectively generalizes choice correlations (Britten, Newsome, Shadlen, Celebrini, & Movshon, 1996), which also infer the animal's noisy percept from observed choices, to continuous time and actions.

We will focus on control policies that assume animals act rationally for a certain model of the world, but do not assume that they maintain the correct model. For example, animals might misjudge the size of the goal area that leads to reward, or have a strong prior belief about the distance they need to travel, leading to biased actions. This approach implies that seemingly suboptimal actions might nonetheless be rational. To capture such suboptimalities, we will parameterize a family of world models and subjective reward functions, and identify the parameters that best predict the animal's actions.

Expected results. This aim will produce a rational control policy that describes the best action (e.g., joystick deflection or running direction) for an animal's current state estimate, as inferred by the models of Aim 2.1. This action will depend on assumed costs and rewards that are ideally determined by the task. However, we will support parametric deviations from this ideal, to take into account the possibility that animals may have suboptimal world models. These parameters will be tuned to best match the animal's observed behavior. We will perform model comparison to distinguish among the different dynamic causal inference models of Aim 2.1

Relation to experiments. Rational control becomes essential to model continuous actions of the tasks in Project C. Latent state estimates for this task are described in Aim 2.1. The rational controller will predict for each latent state the best action to achieve the task's aim (i.e., steer toward the goal area) for a particular assumption about rewards and costs. As described above, we will simultaneously tune the model's assumptions about rewards and costs (thus influencing the rational policy linking state estimates to predicted actions), as well as the model's assumptions about the world dynamics (linking momentary sensory evidence to state estimates, Aim 2.1) to best capture the animal's behavior. By this process, we will predict the sequence of latent variables that the animal is most likely to maintain throughout the task, which we will use to identify their neural correlates in Projects B and C (see Overall Strategy).

Potential pitfalls and solutions. Identifying the control parameters on this task that best explain the animal's actions has precedent and is thus feasible (Daptardar et al., 2019). Nonetheless, if we find systematic deviations from predicted and observed actions, we will revise the assumed goals and control costs to correct them. Rational control with artificial neural networks requires identifying a suitable network structure. We will start with a radial basis function architecture, and move to deeper, more complex architectures if our model recovery procedure (see Data Science Core) identifies errors that are too large.