

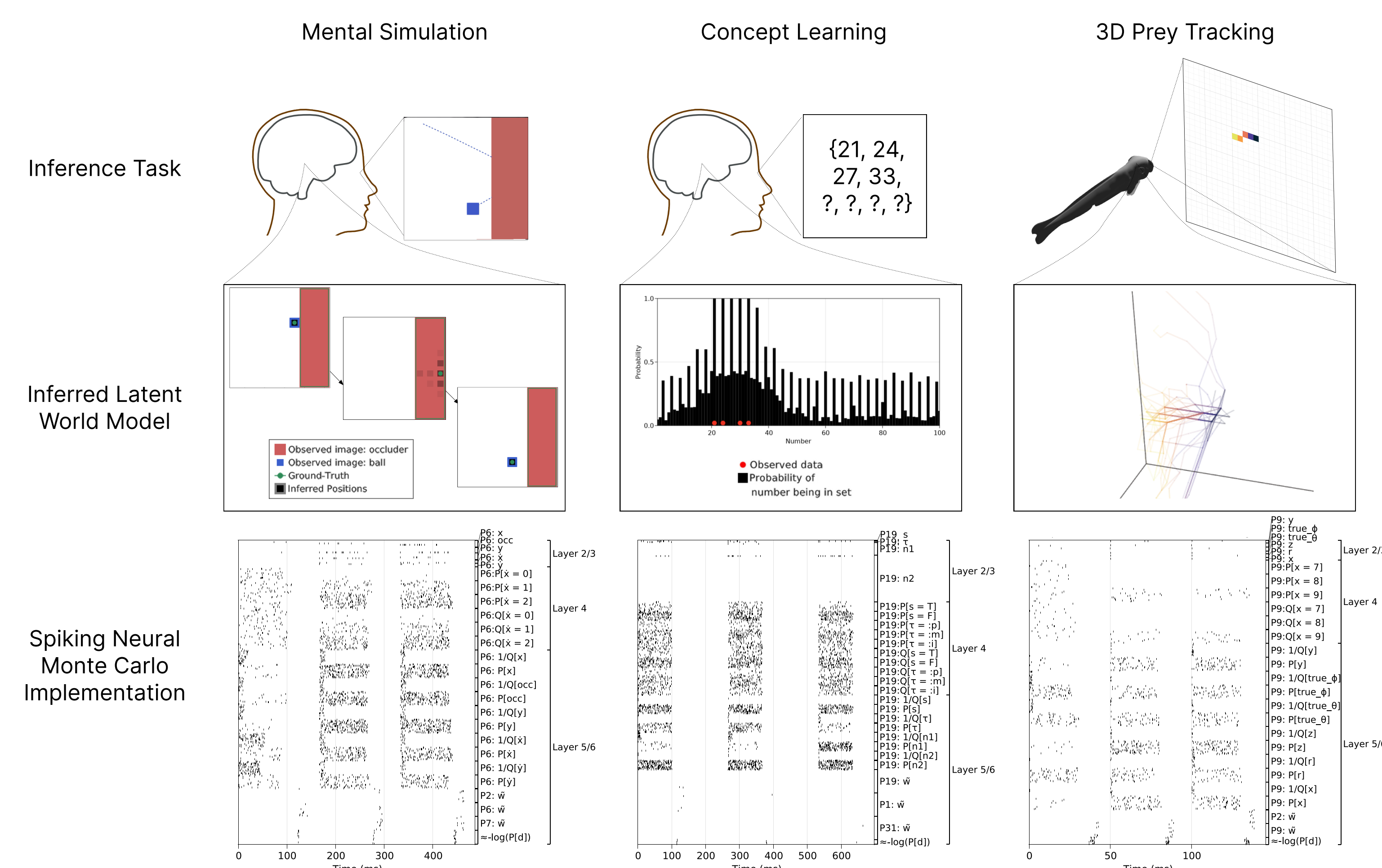
Brain computation as fast spiking neural Monte Carlo inference in probabilistic programs

George Matheos*, Andrew D Bolton*, McCoy Becker, Cameron Freer, Vikash Mansinghka

Introduction

How can slow, spiking neurons implement the fast probabilistic inferences needed to explain perception and cognition? Biological neurons are millions of times slower than electronic computers, yet they can somehow approximate probabilistic inferences in complex probabilistic programs with many latent variables in real-time. Here we show how neurons could perform probabilistic inference, using massively parallel spiking assemblies to implement a novel neural coding scheme, called a dynamically weighted Monte Carlo spiking code. We prove that these assemblies generate approximate samples and make unbiased estimates of probabilities and importance weights, enabling sound approximate inference. Sampled latent variables are sparsely coded, but probabilities and weights are densely coded, and can be read via time-varying, divisively normalized decoding of the dense spiking from specific sub-populations. We show how to implement data-driven artificial neural networks for making fast, bottom-up proposals that are scored and corrected using a structured generative model, yielding new hybrids of distributed and localized neural representations. These spiking neural Monte Carlo architectures scale exponentially better than probabilistic population codes, and are neurally mappable, but unlike deep learning models, they also enable sound implementations of state-of-the-art model-driven AI architectures and inference processes from Bayesian cognitive science. We demonstrate generality by providing spiking circuits for probabilistic program models of visual prey tracking by larval zebrafish, mental physics simulation by primates, and human concept learning. We also present empirical support for this theory, confirming predictions for neural connectivity, coding, and dynamics using data from multiple brain regions and model organisms.

SNMC scales to real time perception and cognition



	Latent Variables	Observed Variables	Size of Spiking Neural Representation	Weighted Monte Carlo (this paper)	ENS Codes, Standard PPCs
1D object tracking	$\{x_t, \hat{x}_t\}_t$	$\{d_t^x\}_t$	Sparse: 27 Dense: 5	Sparse: 27 Dense: 5	Dense: 140
2D object tracking	$\{x_t, y_t, \hat{x}_t, \hat{y}_t\}_t$	$\{d_t^x, d_t^y\}_t$	Sparse: 30 Dense: 10	Sparse: 30 Dense: 10	Dense: 2500
Mental Physics Simulation	$\{x_t, y_t, \hat{x}_t, \hat{y}_t, \alpha_t\}_t$	$\{((d_t^{x,y})_{i=1}^{10})_{i=1}^2\}_t$	Sparse: 38 Dense: 110	Sparse: 38 Dense: 110	Dense: 2500
3D object tracking from 2D observations	$\{x_t, y_t, \hat{x}_t, \hat{y}_t, z_t, r_t, \phi_t, \theta_t\}_t$	$\{d_t^x, d_t^y\}_t$	Sparse: 160 Dense: 20	Sparse: 160 Dense: 20	Dense: 23,180,062,500
Recursive Concept Learning (Sizes are for $D=2, M=10$)	$\bigcup_{k=1}^D \bigcup_{a=1}^M \{s_k^{(a,b)}, r_k^{(a,b)}, \tau_k^{(a,b)}\}$	$\{d_t\}_{t=1}^M$	Sparse: 180 Dense: 40	Sparse: 180 Dense: 40	Dense: 5.92704×10^{11}

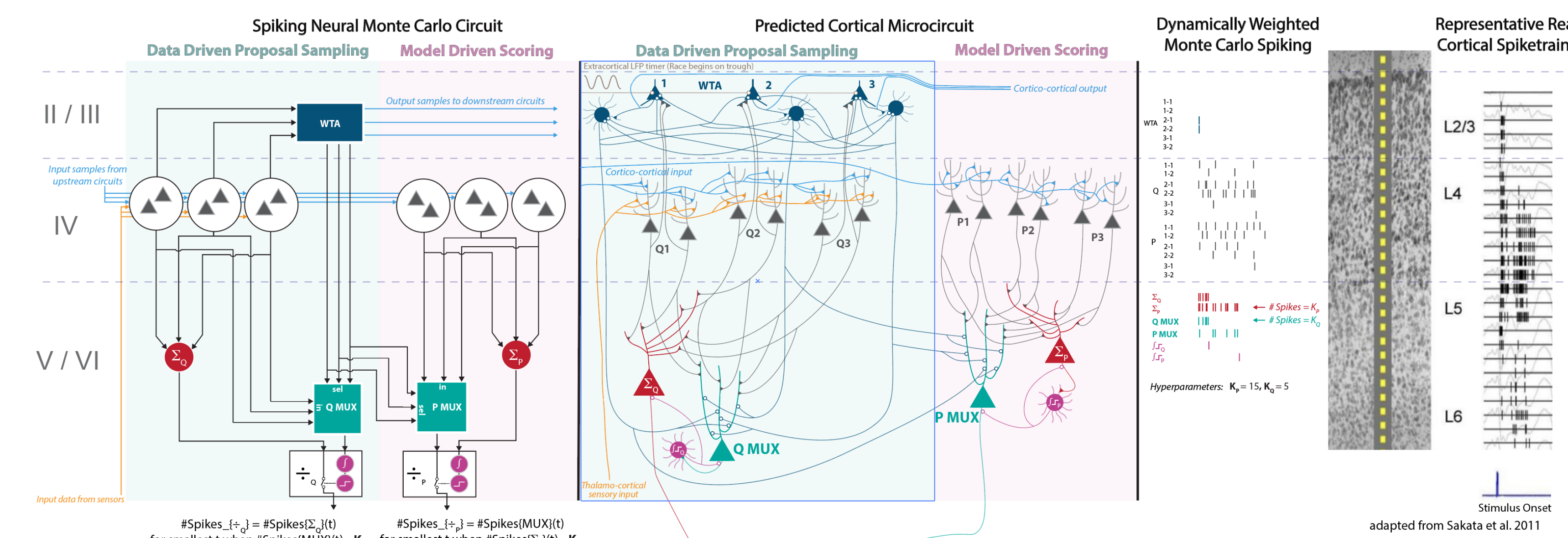
Spiking Neural Monte Carlo requires exponentially fewer neurons than standard probabilistic population codes and ENS spiking codes. For low-dimensional probabilistic programs that only make a small number of latent choices, the difference can be modest in absolute terms. As the number of latent variables in the probabilistic program grows, the cost of the neural representation for previously proposed schemes grows exponentially, rendering them impractical for the majority of perceptual and cognitive inferences. (Ma, Beck, et al 2006, Legenstein and Maass 2014).

Spiking Neural Monte Carlo circuits

Micro-scale Predictions of SNMC

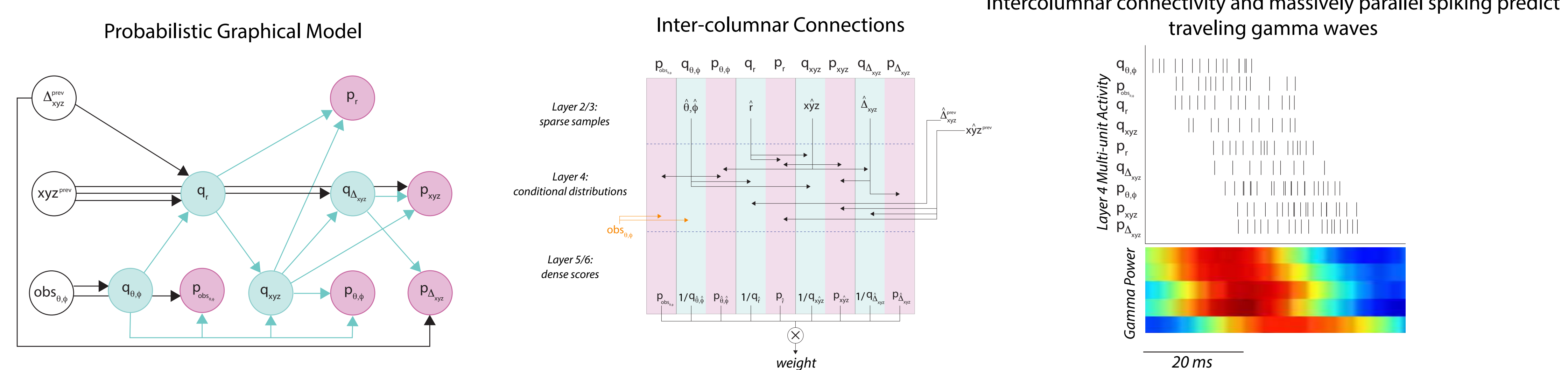
SNMC importance sampler circuit topology, sorted given known sensory (thalamo-cortical) L4 input according to known intra-layer cortical connectivity, predicts multiple findings in hodology, synaptic physiology, and extracellular spike & field electrophysiology:

1. Observed L2/3 sparse and L4/5 dense spiking and relative timing
2. The existence of thalamo-cortical L4 inputs
3. The existence of specific cell types & synaptic physiology in L4/5



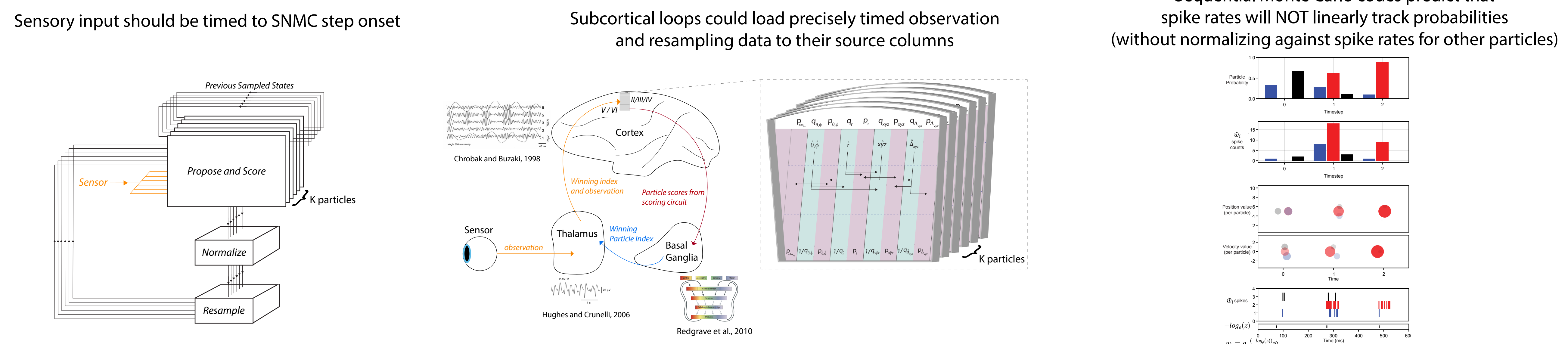
Meso-scale Predictions of SNMC

Successive sample-score epochs, whose order is defined by dependency structure, produce phase shifted bursts of activity in neighboring columns

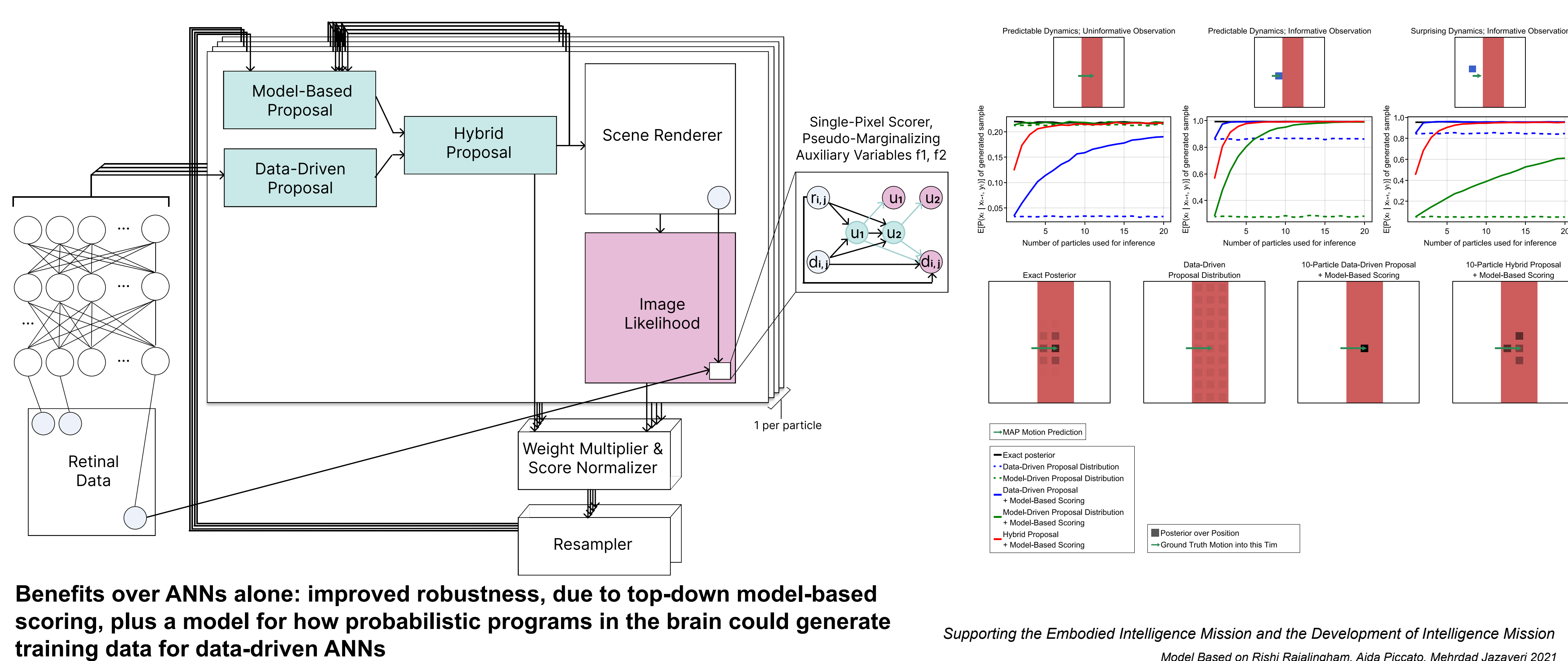


Macro-scale Predictions of SNMC

Topology and timing of sequential Monte Carlo - loading data, proposing variables, scoring variables, and resampling particles - predicts parallel cortico-subcortical loops for each particle, synchronized within particles (across columns) and across columns, via cortico-thalamic theta rhythms



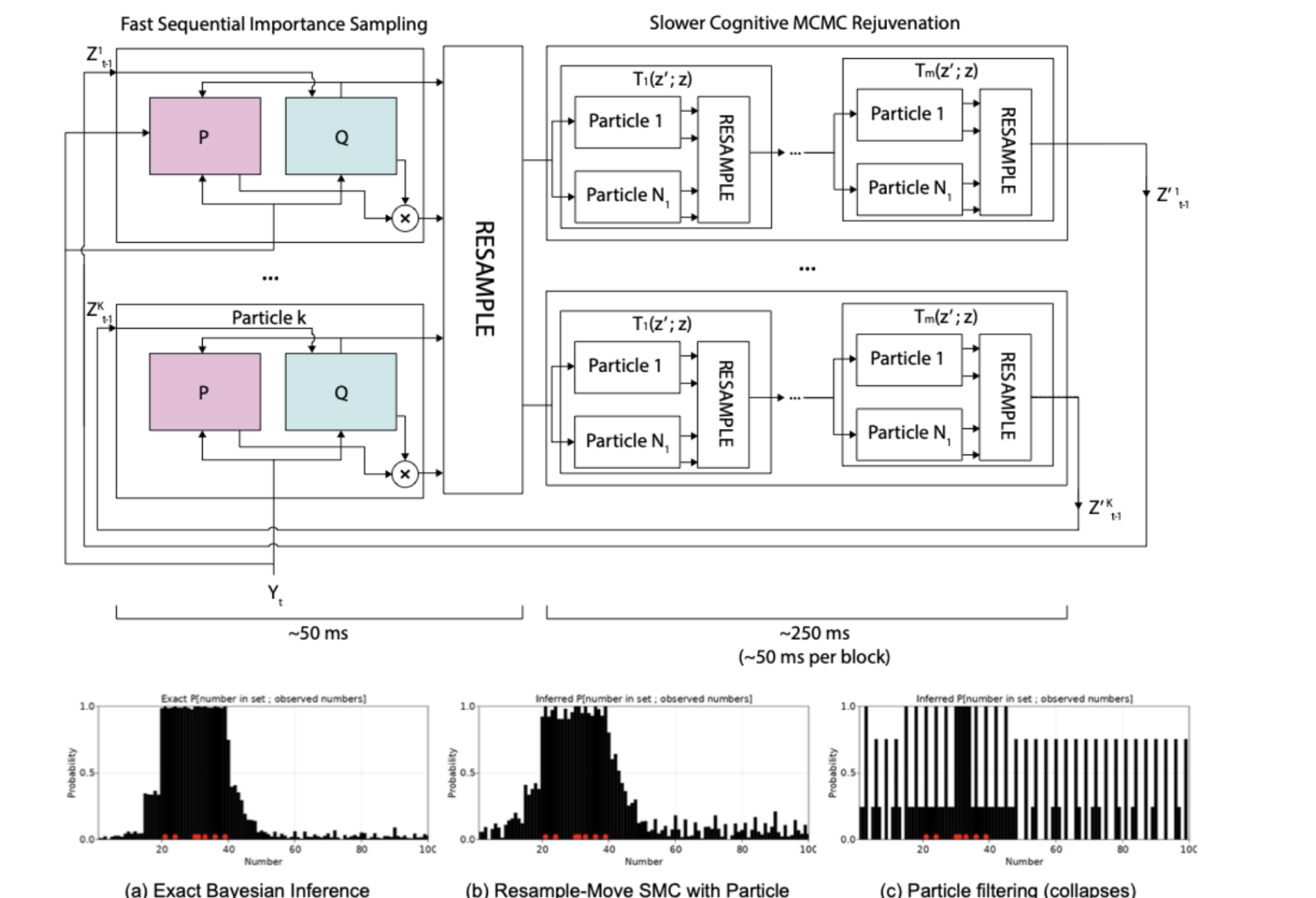
SNMC model of primate physical scene understanding, including data-driven ANN proposals



Benefits over ANNs alone: improved robustness, due to top-down model-based scoring, plus a model for how probabilistic programs in the brain could generate training data for data-driven ANNs

Supporting the Embodied Intelligence Mission and the Development of Intelligence Mission

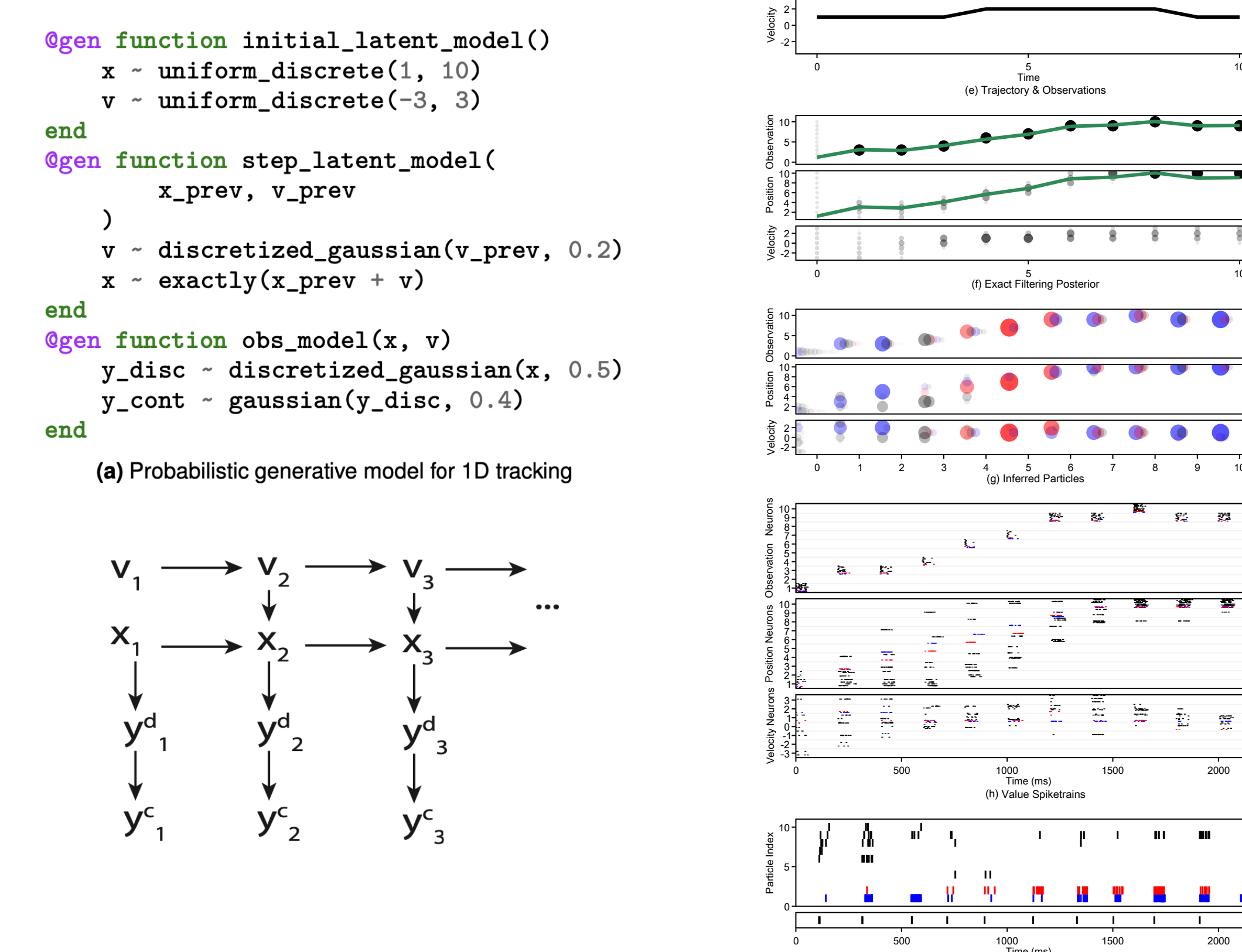
SNMC model of human concept learning via data-driven cognitive MCMC



Benefits over vanilla SMC / particle filtering: robust convergence even for high dimensional problems & unlikely data, due to data-driven proposals and MCMC updates (Gills and Beruzini 2001)

Supporting the Supporting the Development of Intelligence Mission

Constructing SNMC circuits from Gen probabilistic programs



Math

1 Equations for Spike Counts
1.1 P scoring
 For random variable X takes values in $\{1, 2, \dots, |X|\}$. To represent the distribution $P(X=x)$, there are $|X|$ assemblies. Given the value of X is present variable i , the i th assembly spikes as a Poisson process with rate $\lambda_i = \lambda \cdot P(X=i)$.
 After n seconds, the number of spikes from the i th assembly is $N_i(n) \sim \text{Poisson}(\lambda_i n)$.
 For sampled value $x \in \{1, 2, \dots, |X|\}$, we estimate $P(X=x)$ using $\hat{p} = \frac{N_x(n)}{\sum_{i=1}^{|X|} N_i(n)}$.
 For any integer hyperparameter K_x (e.g. $K_x = 15$), if we take $n = \tau_x$ where $\tau_x = \frac{K_x}{\lambda_i}$ is the smallest time τ s.t. $\sum_{i=1}^{|X|} N_i(\tau) = K_x$.
 Then $\hat{p} = P(X=x)$.

1.2 Q sampling and scoring
 To represent the distribution $Q(X=x)$, there are $|X|$ assemblies. Given the value of X , the i th assembly spikes as a Poisson process with rate $\lambda_i = \lambda \cdot Q(X=i)$.
 For $j = 1, 2, \dots$, the index of the assembly which emitted the j th spike, A_j , is uniformly sampled from Q .
 In particular, $A_1 \sim Q(X)$.
 After n seconds, the number of spikes from the i th assembly is $N_i(n) \sim \text{Poisson}(\lambda_i n)$.
 For sampled value $x \in \{1, 2, \dots, |X|\}$, we estimate $q(x) = \frac{N_x(n)}{n}$.
 For any integer hyperparameter K_q (e.g. $K_q = 5$), if we take $n = \tau_q$ where $\tau_q = \frac{K_q}{\lambda_i}$ is the smallest time τ s.t. $N_x(\tau) = K_q$.
 Then $\hat{q} = \frac{N_x(\tau_q)}{\tau_q} = q(x)$.

2 Pseudo-Marginal Importance Sampling for One Variable
 A joint probability distribution $P(X, Y)$ is encoded in the spiking neural network parameters as a proposal latent $Q(X, Y)$.
2.1 Importance Sampling
 Given observed value y , the importance sampling size N copies of a neural circuit to sample, for each $i \in \{1, 2, \dots, N\}$, $x_i^{(i)}$.
 Each sampling circuit is paired with a scoring circuit which mechanically computes values w_i such that $E[w_i] = \frac{P(x_i, y)}{Q(x_i, y)}$.
 Let w_i^* be the normalized weight $w_i^* = \frac{w_i}{\sum_{i=1}^N w_i}$.

2.2 Convergence Guarantees
 Under random i.i.d. condition on Q , the weighted particle cloud $\{w_i^*, x_i^{(i)}\}_{i=1}^N$ converges to the posterior distribution P .
 For any finite N , the normalized particle weights form an unbiased estimator of P : $E[\sum_{i=1}^N w_i^* \delta_{x_i^{(i)}}] = P(x)$.
 These guarantees are obtained in importance sampling. SNMC differs from standard IS by using stochastic, low-probability estimates of $\frac{P(x, y)}{Q(x, y)}$. The spiking Monte Carlo analogs between them [1], we prove that these formal guarantees still hold when using low-probability neural estimates of importance weights.

2.3 Unbiased Weighting in SNMC
 In Spiking Neural Monte Carlo, we use $w_i^* = \frac{P(x_i, y)}{Q(x_i, y)}$.
 In spiking neural networks, w_i^* and $q(x_i)$ are mechanically computed by a neural circuit in each a way that $E[w_i^*] = P(x_i, y)$, $E[q(x_i)] = \frac{1}{Q(x_i, y)}$.

3 Using Artificial Neural Networks within Spiking Neural Monte Carlo
 For an animal to efficiently track an object moving in 2D, using images on a flat retina, we need an efficient proposal distribution $Q(x_t, y_t, \hat{x}_t, \hat{y}_t, z_t, r_t, \phi_t, \theta_t)$ where $Q(x_t, y_t, \hat{x}_t, \hat{y}_t, z_t, r_t, \phi_t, \theta_t)$.

To encode such a proposal distribution efficiently, we use an artificial neural network: $Q(x_t, y_t, \hat{x}_t, \hat{y}_t, z_t, r_t, \phi_t, \theta_t) = \sigma(\mathbf{w}^T \mathbf{x}_t + b)$ where σ denotes the space of one-hot encodings of variables V , \mathbf{w} accepts a rich encoding of $x_{t-1}, y_{t-1}, \hat{x}_{t-1}, \hat{y}_{t-1}, z_{t-1}, r_{t-1}, \phi_{t-1}, \theta_{t-1}$, and outputs distribution over x_t, y_t, z_t .
 This is implemented by copying the output codes from the last layer of the neural network to the assembly in the proposal distribution in Spiking Neural Monte Carlo circuit for sampling and scoring. $Q(x_t, y_t, \hat{x}_t, \hat{y}_t, z_t, r_t, \phi_t, \theta_t)$ is implemented in the spiking network as a fully-connected network of neurons with projections if needed to ensure dense and consistent connectivity, using outer-product approximation of neural values. Our work does not modify the input of this approximation as ANN performance predicts, and how the brain could implement backpropagation.

4 Resample-Move SMC with Particle Gibbs Rejuvenation
4.1 Resample-Move Sequential Monte Carlo
 Dynamic Probabilistic Programs: Resample-Move Sequential Monte Carlo inference is a dynamic probabilistic program $P(x_{1:T}, y_{1:T}) = \prod_{t=1}^T P(x_t | x_{1:t-1}, y_{1:t-1})$.
 Each x_t and y_t may be a collection of values for many different variables (e.g. in a grid giving the intensity of color at every point on an animal's retina).
 Resample-Move SMC: Resample-Move SMC interleaves Importance Sampling steps with Resampling steps and Particle Gibbs MCMC "rejuvenation" steps. Importance sampling uses two data-driven proposal distributions $Q_t(x_t | x_{1:t-1}, y_{1:t-1})$ and $Q_t(y_t | x_{1:t-1}, y_{1:t-1})$.
 Initial importance sampling: At $t=0$, N initial particles are generated and weighted. Size $N \in \{1, 2, \dots, N\}$.
 Resampling: At each timestep $t > 0$, particles from the previous step are resampled, by selecting N "ancestors" indices $a_i \in \{1, 2, \dots, N\}$.
 Here, w_i^* is the normalized weight $w_i^* = \frac{w_i}{\sum_{i=1}^N w_i}$.
 MCMC Rejuvenation: Each resampled particle is optionally rejuvenated by running a MCMC step on it.
 Particle extension: Finally, each particle is extended to the next timestep, and weighted: $w_i^* = \frac{w_i^* P(x_t | x_{1:t-1}, y_{1:t-1})}{Q_t(x_t | x_{1:t-1}, y_{1:t-1})}$.
4.2 Particle Gibbs Rejuvenation
 The family of MCMC kernels in Particle Gibbs: Spiking Neural Monte Carlo can be used to implement any Particle Gibbs algorithm. Here we present the simplest case of Particle Gibbs rejuvenation: "resample-move".
 Proposing modifications to the latent state: Given current latent state $x_{1:t-1}$ and evidence $y_{1:t-1}$, Particle Gibbs resamples x_t by proposed modifications to the value using a proposal distribution $Q_t(x_t | x_{1:t-1}, y_{1:t-1})$. We take $x_t^* = x_t$ and for $i \in \{1, \dots, N\}$:
 There are restrictions on the family of valid proposal distributions $Q_t(x_t | x_{1:t-1}, y_{1:t-1})$. One by restriction is that $Q_t(x_t | x_{1:t-1}, y_{1:t-1})$ may not condition on the value of a variable in $x_{1:t-1}$ if $Q_t(x_t | x_{1:t-1}, y_{1:t-1})$ is a proposal modification to the value of a variable in $x_{1:t-1}$.
 Selecting the next state: Importance weights are computed: $w_i^* = \frac{w_i^* P(x_t | x_{1:t-1}, y_{1:t-1})}{Q_t(x_t | x_{1:t-1}, y_{1:t-1})}$.
 Add a particle index is sampled: $i^* \sim \text{Categorical}(w_1^*, \dots, w_N^*)$.
 when $w_i^* = \frac{w_i^* P(x_t | x_{1:t-1}, y_{1:t-1})}{Q_t(x_t | x_{1:t-1}, y_{1:t-1})}$.
 Particle-Gibbs returns x_t^* as the next MCMC state $x_t^* = x_{1:t-1} \cup \{x_t^*\}$.
 Stationarity: Particle-Gibbs is stationary for $P(x_t)$, meaning that $E[P(x_t)] = P(x_t)$.
 When $N_{t-1} = 2$, this class of Particle-Gibbs reduces to a variant of Metropolis-Hastings.