

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

George Matheos^{a,b,*}, Andrew D Bolton^{b,c,*}, McCoy Becker^b, Cameron Freer^b, and Vikash Mansinghka^b

^aUC Berkeley; ^bMIT Probabilistic Computing Project; ^cHarvard; *Equal Contribution

This manuscript was compiled on November 30, 2022

1 **How can slow, spiking neurons implement the fast probabilistic in-**
2 **ferences needed to explain perception and cognition? Biological**
3 **neurons are millions of times slower than electronic computers, yet**
4 **they appear to robustly approximate probabilistic inferences in com-**
5 **plex probabilistic programs with many latent variables in real-time.**
6 **Here we show how biologically realistic, massively parallel assem-**
7 **blies of spiking neurons can perform real-time probabilistic infer-**
8 **ence. Our approach, based on novel weighted Monte Carlo spiking**
9 **codes that leverage spike rates as well as coarse spike timing, re-**
10 **quires exponentially fewer neurons than standard probabilistic pop-**
11 **ulation codes. It also scales to real-time inference via massively par-**
12 **allel hybrids of model-based Monte Carlo and data-driven neural net-**
13 **works, and works for high-dimensional probabilistic programs that**
14 **previous spiking neural inference architectures do not handle. We il-**
15 **lustrate generality by providing neurally mappable implementations**
16 **of resource-rational variants of Bayesian cognitive models for pri-**
17 **mate mental physical simulation, human learning of numerical con-**
18 **cepts, and 3D prey tracking by larval zebrafish. We also confirm**
19 **predictions of the spiking neural Monte Carlo theory using empiri-**
20 **cal data drawn from the hodology, functional neuroanatomy, synap-**
21 **tic physiology, and extracellular spike and field electrophysiology of**
22 **multiple brain regions and model organisms.**

Brain computation | probabilistic programming | Monte Carlo | deep learning | spiking neural networks | probabilistic inference | visual perception | mental simulation | concept learning

1 **T**his paper addresses two questions: In theory, how can
2 slow, spiking neurons possibly implement the fast ap-
3 proximate probabilistic inferences needed for perception and
4 cognition? And can such a theory predict empirical data
5 from studies of hodology, functional neuroanatomy, synaptic
6 physiology, and extracellular spike and field electrophysiology?
7 This paper introduces new, massively parallel architectures for
8 spiking neural Monte Carlo inference in probabilistic programs
9 with many latent variables, overcoming scaling limitations of
10 previous work on spiking neural inference. It also confirms mul-
11 tiple predictions about fundamental biophysical mechanisms,
12 micro-scale circuits, meso-scale networks, and macro-scale ar-
13 chitectures and dynamics using empirical data from multiple
14 brain regions and model organisms.

15 The idea that everyday perception and cognition relies on
16 probabilistic inference in rich, flexible generative models can
17 be traced back at least as far as Helmholtz (1) and Laplace
18 (2). Probabilistic inference in structured probabilistic models
19 played a central role in multiple generations of artificial intelli-
20 gence systems (3–5) and computational models of cognition (6).
21 Probabilistic inference, and especially sampling-based, Monte
22 Carlo approximate inference approaches, are also central to
23 reverse-engineering approaches in computational cognitive sci-
24 ence, especially the traditions of “resource-rational” analysis of

cognitive inference processes (7), and in the “Bayesian brain”
(8) or “sampling hypothesis” frameworks (9, 10). Unfortu-
nately, despite the conceptual appeal of this perspective, it
has proved difficult to bridge the gap between computational
theories of inference and neural representation (11, 12).

Probabilistic programming (13–17) provides a computa-
tional formalism for generalizing and scaling implementations
of inference in generative models. Probabilistic programs with
many latent variables are increasingly central to state-of-the-
art architectures for real-time 3D computer vision (17, 18) and
theory of mind via inverse planning (19, 20) and also to com-
putational cognitive science (21). Probabilistic programs offer
new possibilities for solving problems central to embodied in-
telligence by integrating data-driven and model-driven modes
of inference (22, 23), and support state-of-the-art hybrids of
sequential Monte Carlo (16, 24) with variational inference
(25, 26). Probabilistic programs can even encode risk-sensitive
action selection and decision-theoretic planning (13, 27, 28).
Unfortunately, thus far, there have been no spiking neural
architectures that can scale to perform real-time, high-quality
approximate probabilistic inference in probabilistic programs
with many latent variables. There is thus a fundamental gap
between computational models of intelligence and biologically
realistic models of brain computation.

Significance Statement

Cognitive science, neuroscience, and artificial intelligence have not yielded an integrative theory of how probabilistic inference is implemented in the mind and brain, leaving fundamental gaps between phenomenological, causal, and computational accounts of intelligence. Spiking neural Monte Carlo narrows these gaps, offering a theory for reverse-engineering brain computation that is more computationally general, cognitively realistic, and biologically grounded than artificial neural networks on their own. It gives a unifying explanation of micro-scale, meso-scale, and macro-scale features of neural connectivity, coding, and dynamics. It shows how to automatically construct implementations of a broad class of probabilistic programs that encode Bayesian models, and test their predictions against both behavioral and neural data. Finally, it exposes massive micro-scale, meso-scale, and macro-scale parallelism inherent in probabilistic programming, yielding a new brain-like scaling route for engineering intelligent machines.

G.M. and A.B. and V.M. performed research; C.F. and M.B. assisted with research; V.M. designed and oversaw research; and V.M., A.B., and G.M. wrote the paper.

Please declare any competing interests here.

^{1,2}George Matheos and Andrew Bolton contributed equally to this work.

^{1, 2, 5}To whom correspondence should be addressed. E-mails: georgematheos@berkeley.edu, andrewbolton@fas.harvard.edu, and vkmm@mit.edu

49 Consider that inference in probabilistic programs is ordi- 110
 50 narily implemented using electronic computers that perform 111
 51 hundreds of millions of instructions per second. These com- 112
 52 puters are in turn implemented via logic gates that transition 113
 53 billions of times per second. In contrast, biological neurons 114
 54 spike millions of times slower, yet many perceptual inferences 115
 55 require just hundreds of milliseconds, and many cognitive in- 116
 56 ferences require just seconds. This in turn means that the 117
 57 brain must somehow approximate probabilistic inference using 118
 58 massively parallel circuits that integrate new data without 119
 59 the long sequential chains of operations that can be used 120
 60 in software. Influential feedforward models of visual object 121
 61 recognition have just 5 layers (29). 122

62 Artificial neural network (ANN) models can be trained 123
 63 to provide low-latency approximate probabilistic inferences. 124
 64 They have been used to build neurally mappable models of 125
 65 primate vision (30–32) and mental simulation (33), as well 126
 66 as larger-scale models (34) simulation. However, fundamen- 127
 67 tal limitations of ANN models, both as AI technology and 128
 68 as models of visual perception, are also increasingly widely 129
 69 recognized (35–38) 130

70 Even proponents of ANN models see fundamental open 131
 71 problems, such as how to account for the role of top-down 132
 72 connections in visual cortex, and therefore the computational 133
 73 interactions between bottom-up, data-driven processing and 134
 74 top-down, model-based feedback (39) This limitation appears 135
 75 related to failures of ANN models in practice. Consider that 136
 76 even state-of-the-art extensions to CNNs and RNNs for visual 137
 77 perception, trained on massive datasets via algorithms that 138
 78 lack biologically plausible implementations, exhibit striking 139
 79 failures that biological vision systems do not (37). These 140
 80 include adversarial examples (40), and also other failures of 141
 81 common sense, such as falsely positing everyday 3D objects 142
 82 floating inexplicably in unoccupied space (18), and failing 143
 83 to miss visually salient objects such as pedestrians, trucks, 144
 84 and emergency vehicles (41, 42) Also, standard ANN models 145
 85 do not reflect many fundamental characteristics of biological 146
 86 neural networks, ranging from the laminar structure of cortex 147
 87 to combinations of dense and sparse coding to widespread 148
 88 oscillations and synchrony. 149

89 Inspired in part by these challenges, there is a rich liter- 150
 90 ature on spiking neural architectures and other massively 151
 91 parallel circuit formalisms for probabilistic inference. Promi- 152
 92 nent examples include probabilistic population codes (43, 44); 153
 93 spiking neural Gibbs samplers suitable for inference in discrete 154
 94 Bayesian network models with sufficiently sparse connectivity 155
 95 (45, 46); spiking neural Bayesian filters that extend probabilis- 156
 96 tic population codes for real-time tracking (47, 48). However, 157
 97 these previous proposals for spiking probabilistic inference 158
 98 cannot implement state-of-the-art schemes for real-time se- 159
 99 quential Monte Carlo inference in complex probabilistic pro- 160
 100 grams. There have also been proposals for stochastic digital 161
 101 circuits for massively parallel, low precision, real-time Monte 162
 102 Carlo (49–51) inference in probabilistic graphical models and 163
 103 non-parametric Bayesian models with tens of thousands of 164
 104 variables; and other, more specialized neural inference schemes. 165
 105 But these stochastic digital circuits do not explain how to per- 166
 106 form robust, real-time probabilistic inference using components 167
 107 that are as slow as biological neurons. 168

108 This paper introduces spiking neural Monte Carlo circuits, 169
 109 including new weighted Monte Carlo spiking codes and mas-

sively parallel spiking neural assemblies. It shows that these 110
 architectures enable hybrids of data-driven and model-driven 111
 Monte Carlo inference that suffice for real-time probabilistic 112
 inference in probabilistic programs. It includes architectures 113
 for generating approximate samples for latent variables and for 114
 unbiased estimation of probability densities and importance 115
 weights. Crucially, these architectures can also be used recur- 116
 sively. They apply at the scale of individual latent variables, 117
 and to larger collections of latent variables, arising in both 118
 target models and in proposal distributions. This approach 119
 thus enables complex Monte Carlo inference architectures, 120
 with proposals defined by data-driven probabilistic programs, 121
 including artificial neural networks, whose outputs are re- 122
 weighted and corrected via generative model-driven Monte 123
 Carlo inference. Expressiveness is illustrated via three exam- 124
 ples: visual prey tracking by larval zebrafish; mental physics 125
 simulation by both humans and non-human primates; and 126
 recursive concept learning by human adults. Figure 1 shows 127
 three inference tasks, each previously studied in Bayesian cog- 128
 nitive science, for which our approach provides the first spiking 129
 neural implementations. 130

1. Spiking neural Monte Carlo 131

A. Weighted Monte Carlo spiking codes. A dynamic proba- 132
 bilistic program defines a joint density over a sequence $\mathbf{z}_{1:T}$ of 133
 latent states, and a sequence of observed data, $\mathbf{d}_{1:T}$: 134

$$P(\mathbf{z}_{1:T}, \mathbf{d}_{1:T}) = P(\mathbf{z}_1) \prod_{t=2}^T P(\mathbf{z}_t | \mathbf{z}_{t-1}) \prod_{t=1}^T P(\mathbf{d}_t | \mathbf{z}_t) \quad [1] \quad 135$$

Inference in a dynamic probabilistic program consists of 136
 estimating the conditional distribution $P(\mathbf{z}_{1:t} | \mathbf{d}_{1:t})$ for each 137
 t . This can be done sequentially, using the inference about 138
 $P(\mathbf{z}_{1:t-1} | \mathbf{d}_{1:t-1})$ to form inferences about $P(\mathbf{z}_{1:t} | \mathbf{d}_{1:t})$, using 139
 the following recursion: 140

$$P(\mathbf{z}_{1:T} | \mathbf{d}_{1:T}) = P(\mathbf{z}_{1:T-1} | \mathbf{d}_{1:T-1}) \times Q(\mathbf{z}_T; \mathbf{z}_{T-1}, \mathbf{d}_T) \frac{P(\mathbf{z}_T, \mathbf{d}_T | \mathbf{z}_{T-1})}{Q(\mathbf{z}_T; \mathbf{z}_{T-1}, \mathbf{d}_T) P(\mathbf{d}_T | \mathbf{d}_{T-1})} \quad [2] \quad 141$$

The sequential Monte Carlo algorithm implements this re- 142
 cursion for an approximation each distribution $P(\mathbf{z}_{1:t} | \mathbf{d}_{1:t})$ re- 143
 presented as a set $\{(\mathbf{z}_{1:t}^i, w_t^i)\}_{i=1}^N$ of “weighted particles” meant 144
 to approximate the distribution. At each time t , the existing 145
 particles of the form $\mathbf{z}_{1:t-1}^i$ are extended, and the weights are 146
 updated, according to

$$\mathbf{z}_t^i \sim Q(\cdot; \mathbf{d}_t, \mathbf{z}_{t-1}^i), \quad w_t^i = w_{t-1}^i \frac{P(\mathbf{z}_t^i | \mathbf{z}_{t-1}^i) P(\mathbf{d}_t | \mathbf{z}_t^i)}{Q(\mathbf{z}_t^i; \mathbf{d}_t, \mathbf{z}_{t-1}^i)} \quad [3] \quad 147$$

**B. Massively parallel micro-scale spiking assemblies and mi- 136
 cro-circuits for individual latent variables.** A conditional proba- 137
 bility distribution $P(z | \text{par}(z))$ is implemented using a col- 138
 lection of neural assemblies, one for each value i that z can 139
 take. The i th assembly has $a_P^{z=i}$ neurons in it, each spiking 140
 at rate $r_P^{z=i}$, given a particular value of $\text{par}(z)$. The overall 141
 rate of the assembly is therefore $\lambda_P^{z=i} = a_P^{z=i} \times r_P^{z=i}$. The 142
 assemblies correctly implement $P(z | \text{par}(z))$ if the rate of the 143
 each assembly is equal to the probability of the corresponding 144
 value of z , up to some proportionality constant γ_P : 145

$$\lambda_P^{z=i} = \gamma_P P(z = i | \text{par}(z)) \quad [4] \quad 146$$

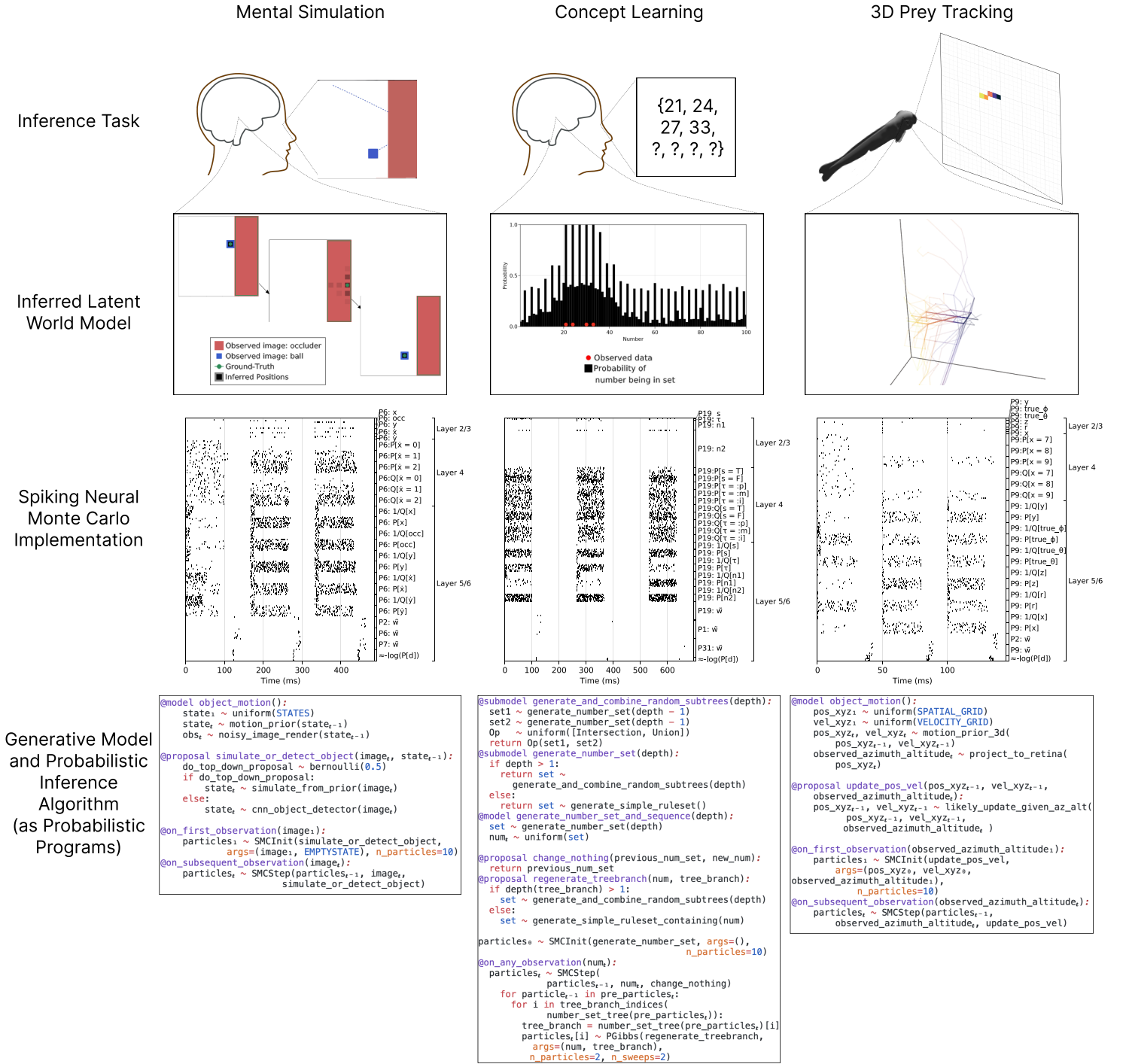


Fig. 1. Real-time spiking neural Monte Carlo models for diverse probabilistic inferences in perception and cognition. Tasks are primate mental simulation (left column), human concept learning (middle column), and 3D prey tracking (right column). Each task (top row) requires the model organism to infer latent world models (second row). This is achieved by a spiking neural Monte Carlo implementation (third row) of probabilistic programs (bottom row) that encode a generative model and a sequential Monte Carlo inference algorithm, implementing hybrids of data-driven and model-driven inference. Spike rasters show weighted Monte Carlo spiking representations distributed across model neurons from superficial (L2/3), middle (L4), and deeper layers of cortex (L5/6) including sparse codes, dense codes, and power oscillations in the gamma and theta bands. The theory in this paper shows how to automatically construct spiking models such as these that implement biologically realistic, massively parallel, real-time inference for a broad class of probabilistic programs.

147 Likewise, a proposal distribution $Q(z; d)$ is implemented
148 using a collection of assemblies, the i th having $a_Q^{z=i}$ neurons
149 at rate $r_Q^{z=i}$ given d to achieve total rate $\lambda_Q^{z=i} = a_Q^{z=i} \times r_Q^{z=i}$,

where

$$\lambda_Q^{z=i} = \gamma_Q Q(z = i; d) \quad [5] \quad 151$$

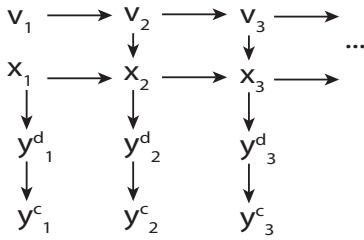
For $k > 0$ let $s_P^{i,k}$ be the time when the k th spike is emitted 152

```

@gen function initial_latent_model()
  x ~ uniform_discrete(1, 10)
  v ~ uniform_discrete(-3, 3)
end
@gen function step_latent_model(
  x_prev, v_prev
)
  v ~ discretized_gaussian(v_prev, 0.2)
  x ~ exactly(x_prev + v)
end
@gen function obs_model(x, v)
  y_disc ~ discretized_gaussian(x, 0.5)
  y_cont ~ gaussian(y_disc, 0.4)
end

```

(a) Probabilistic generative model for 1D tracking



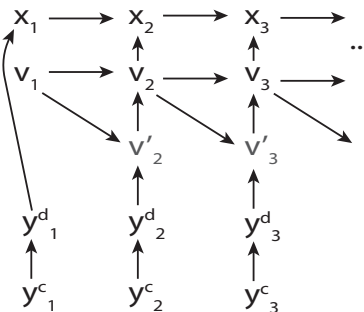
(b) Graphical model corresponding to (a).

```

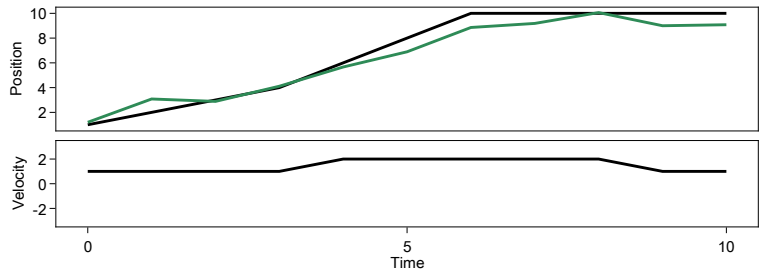
# activation of each tuning-curve, given y_cont
tuning_curves(y_cont) = normalize([
  pdf(gaussian, (mean, 0.4), y_cont)
  for mean in 1:10
])
@gen function initial_proposal(y_cont)
  y_disc ~ cat(tuning_curves(y_cont))
  v ~ uniform_discrete(-3, 3)
  x ~ discretized_gaussian(y_disc, 0.5)
end
@gen function step_proposal(x_prev, v_prev, y)
  y_disc ~ cat(tuning_curves(y))
  v' = y_disc - x_prev # apparent velocity
  v ~ discretized_gaussian((v_prev + v')/2, 1)
  x ~ exactly(x_prev + v)
end

```

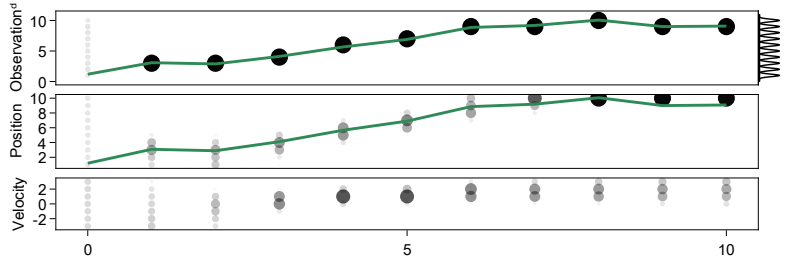
(c) Data-driven inference proposals



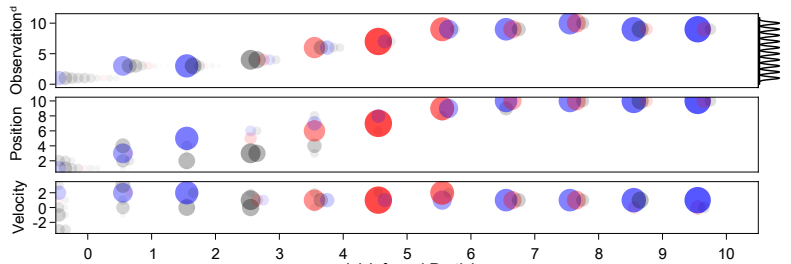
(d) Graphical model corresponding to (c).



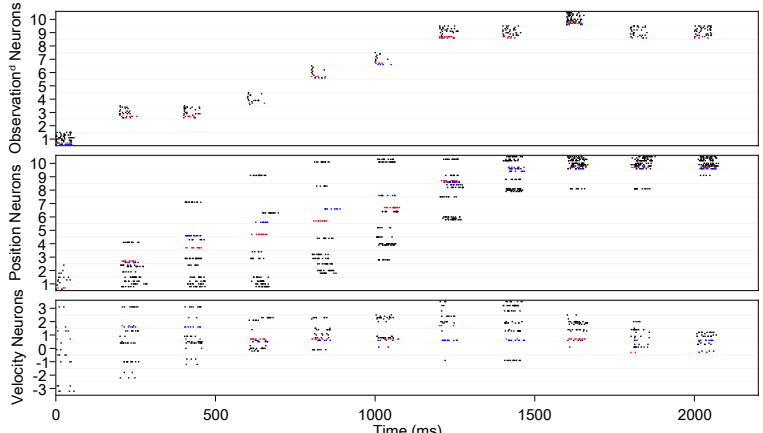
(e) Trajectory & Observations



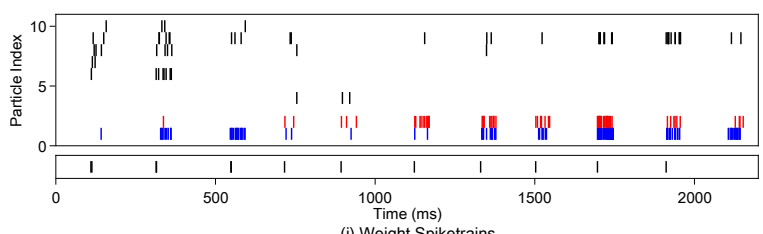
(f) Exact Filtering Posterior



(g) Inferred Particles



(h) Value Spiketrains



(i) Weight Spiketrains

Fig. 2. A weighted Monte Carlo spiking code for a 1 dimensional visual prey tracking problem. (a) shows the target probabilistic generative model, implemented in Gen, and (b) shows its top-down dependency structure as a directed graphical model. (c) shows bottom-up, data-driven inference proposals, along with (d) their dependence structure. (e) shows the inference problem being solved: tracking a prey moving in 1 dimension. The ground truth prey trajectory is shown in black, and the observed position data is shown in green. (f) shows an exact Bayesian filtering solution to this problem. Note that observed position is discretized via Gaussian tuning curves before conditioning the Bayes filter. (g) shows a particle filtering approximation, with particle importance weights encoded by circle size. (h) shows sparse codes for sampled values representing the particles from (g), and (i) shows dense codes for estimated importance weights from (g).

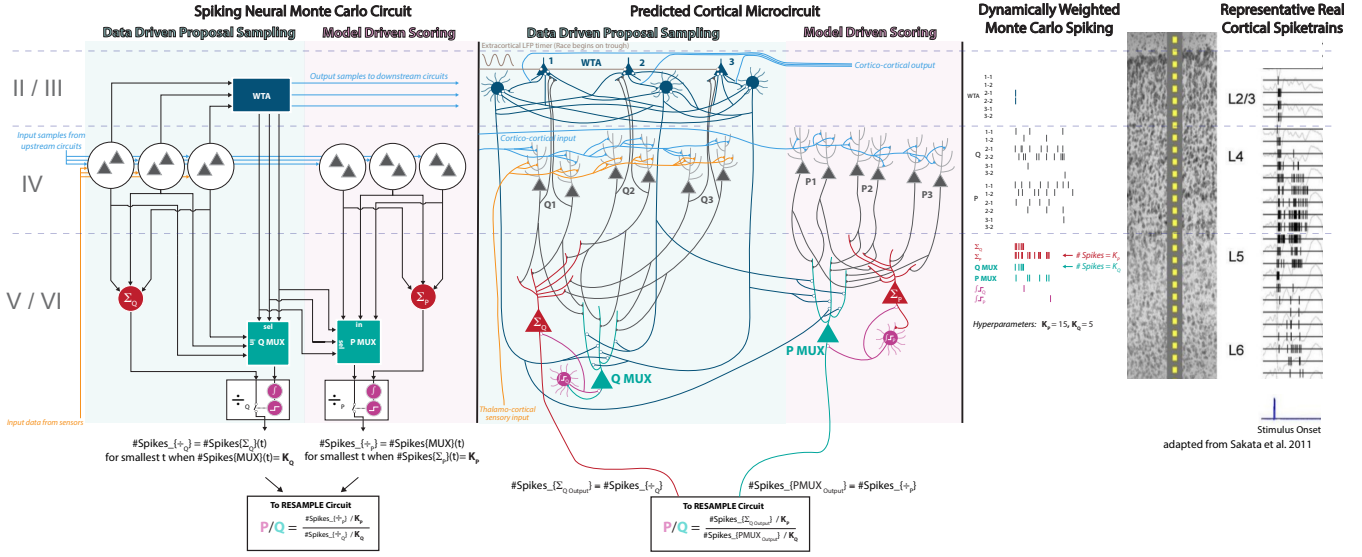


Fig. 3. Micro-scale spiking assemblies and spiking neural Monte Carlo micro-circuits predict connectivity, coding, and dynamics of real cortical micro-circuits. From left to right, this figure shows the architecture for proposing and scoring a single latent variable; the predicted connectivity and synaptic characteristics of a real cortical micro-circuit implementing this design; representative spiking at each layer; and real cortical spiking from a depth electrode recording. Each spiking assembly represents a possible value for a random variable. The sampled proposal value is generated via a winner-take all race, which feeds into MUXs to generate spike counts. These spike counts yield provably unbiased estimates of proposal and target probabilities, and can be combined into importance weights. This architecture turns out to predict multiple features of the synaptic physiology, connectivity, and spiking dynamics of cortical micro-circuits.

153 from the P -assembly corresponding to outcome $z = i$, and let
 154 $s_P^{i,0} = 0$. We model the overall assembly as a Poisson Process,
 155 meaning that for each i and each k , the time between the k
 156 and $k + 1$ th spike is exponentially distributed:

$$157 \quad s_P^{i,k} - s_P^{i,k-1} \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(\lambda_P^{z=i}) \quad [6]$$

158 $S_P^0 = \{s_P^{i,k}\}_{i,k}$ is the set of the times at which any P
 159 assembly spikes. Let S_P^k be the set of all spike-times except
 160 the first k , and let t_P^k be the time at which the k th spike is
 161 emitted from any assembly:

$$162 \quad t_P^k = \inf S_P^{k-1}, \quad S_P^k = S_P^{k-1} \setminus \{t_P^k\} \quad [7]$$

163 Let z_P^k be the value i of z corresponding to the assembly which
 164 emitted the k th spike among all the spikes:

$$165 \quad z_P^l = i : t_P^l \in \{s_P^{i,k}\}_k \quad [8]$$

166 Likewise we let $s_Q^{i,k}$ denote the time of the k th spike from
 167 the i th Q -assembly, and we define S_Q^k , t_Q^k , and z_Q^k analogously:

$$168 \quad \forall k > 0, s_Q^{i,k} - s_Q^{i,k-1} \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(\lambda_Q^{z=i}) \quad [9]$$

$$169 \quad \forall k > 0, t_Q^k = \inf S_Q^{k-1}, \quad S_Q^k = S_Q^{k-1} \setminus \{t_Q^k\} \quad [10]$$

$$170 \quad \forall k > 0, z_Q^k = i : t_Q^k \in \{s_Q^{i,k}\}_k \quad [11]$$

173 where $S_Q^0 = \{s_Q^{i,k}\}_{i,k}$ and $s_Q^{i,0} := 0$.

174 It turns out that since the assemblies are Poisson Processes
 175 (Eqns. 6, 9) with appropriately set rates (Eqns. 4, 5), the
 176 values z_P^k and z_Q^k are fair samples:

$$177 \quad z_P^k \stackrel{\text{i.i.d.}}{\sim} P(z|\text{par}(z)), \quad z_Q^k \stackrel{\text{i.i.d.}}{\sim} Q(z;d) \quad [12]$$

178 From this it is evident how to draw a sample from the Q
 179 proposal distribution: simply use one of the z_Q^k values! In our

180 proposed circuits, the identity $i^* = z_Q^1$ of the first assembly to
 181 spike is selected by a Winner-Take-All circuit to be used as
 182 the sampled value, and output in the sparse code.

183 Since each spike index is a fair sample from P or Q , we can
 184 obtain approximations of $P(z = i^*|\text{par}(z))$ and $Q(z = i^*; d)$
 185 using simple Monte Carlo estimates. To do this for a P -score,
 186 we use a circuit that considers what fraction of the first c_P
 187 spikes came from the i^* th assembly. Let $N_P^{z=i}$ be the number
 188 of spikes from the i th assembly, out of the first c_P spikes

$$189 \quad N_P^{z=i} = \sum_{k=1}^{c_P} 1_{z_P^k=i}, \quad [13]$$

190 and let N_P^z denote the count for the sampled value, $N_P^z =$
 191 $N_P^{z=i^*}$. Let \hat{p} be the fraction of the first c_P spikes to come
 192 from the i^* th assembly:

$$193 \quad \hat{p} = \frac{1}{c_P} N_P^z \quad [14]$$

194 Then \hat{p} is an unbiased estimate of $p := P(z = i^*|\text{par}(z))$:
 195 $\mathbb{E}[\hat{p}] = p$.

196 For the Q distribution we obtain a probability-value esti-
 197 mate slightly differently, because (1) we wish to obtain an
 198 unbiased estimate of $q^{-1} = 1/Q(z = i^*; d)$, rather than of q ,
 199 and (2) the first spike z_P^1 equals the sampled value i^* (since
 200 this is how the circuit samples i^*) and therefore this first spike
 201 must be ignored. Due to (1), rather than waiting for a fixed
 202 number of spikes to occur from any assembly, the $1/Q$ -scoring
 203 circuit waits until c_Q spikes occur *in the i^* th assembly*. Let
 204 $N_Q^{z=i}$ be the number of spikes from the i th assembly by the
 205 time that the $c_Q + 1$ th spike is emitted by the i^* th assembly,

$$206 \quad N_Q^{z=i} = \sum_k 1_{s^{i,k} \leq s^{i^*,c_Q+1}}, \quad [15]$$

207 and let N_Q^z be the number of spikes from *all* the assemblies
 208 in this time, excluding the first spike: $N_Q^z = \sum_i N_Q^{z=i} - 1$.

209 Then let \hat{q}^{-1} be the ratio of the total number of spikes in any
 210 assembly to the number from the i^* th assembly in this time:

$$211 \quad \hat{q}^{-1} = \frac{1}{c_Q} N_Q^z \quad [16]$$

212 Then \hat{q}^{-1} is an unbiased estimate of q : $\mathbb{E}[\hat{q}^{-1}] = q$.

213 Let $\tau_{Q,z}^{\text{sample}}$ be the amount of time to generate a sample
 214 of z from Q , and let $\tau_{Q,z}^{\text{score}}$ be the amount of time needed to
 215 Q -score the sample after it has been drawn,

$$216 \quad \tau_{Q,z}^{\text{sample}} = t_Q^1, \quad \tau_{Q,z}^{\text{score}} = s_Q^{i^*, c_Q} \quad [17]$$

SO

$$\tau_{Q,z}^{\text{sample}} \sim \text{Exponential}\left(\sum_i \lambda_Q^i\right),$$

$$\tau_{Q,z}^{\text{score}} - \tau_{Q,z}^{\text{sample}} \sim \text{Erlang}(c_Q + 1, \lambda_Q^{i^*})$$

217 The amount of time $\tau_{P,z}^{\text{score}}$ needed to P -score z is defined
 218 analogously.

219 **C. Massively parallel meso-scale spiking networks for high-**
 220 **dimensional probabilistic programs.** For simplicity of present-
 221 ation we assume that the latent state and data vectors \mathbf{z}_t and
 222 \mathbf{d}_t have fixed length regardless of t , and we use the indices
 223 $\{1, \dots, |\mathbf{z}_t|\}$ to refer to the variables in \mathbf{z}_t and the indices
 224 $\{|\mathbf{z}_t| + 1, \dots, |\mathbf{z}_t| + |\mathbf{d}_t|\}$ to refer to the variables in \mathbf{d}_t .

225 Since the step model is a probabilistic program, the density
 226 $P(\mathbf{z}_t | \mathbf{z}_{t-1})$ decomposes into the product

$$227 \quad P(\mathbf{z}_t | \mathbf{z}_{t-1}) = \prod_{i=1}^{|\mathbf{z}_t|} P(\mathbf{z}_t^i | \{\mathbf{z}_t^j\}_{j \in \text{par}_P^t(i)}, \{\mathbf{z}_{t-1}^j\}_{j \in \text{par}_P^{t-1}(i)}) \quad [18]$$

228 where $\text{par}_P^t(i)$ and par_P^{t-1} are the indices of the parent variables
 229 of variable i in \mathbf{z}_t and \mathbf{z}_{t-1} respectively.

230 Likewise, the Q proposal probabilistic program decomposes
 231 into a product. Since values sampled from Q proposal distri-
 232 butions must be sampled in the topological order induced by
 233 the probabilistic program, we emphasize that the variables
 234 proposed by Q are organized into Depth_Q layers L_Q^1 through
 235 $L_Q^{\text{Depth}_Q}$, which together form a partition of the set of variable
 236 indices $\{1, \dots, |\mathbf{z}_t|\}$. $Q(\mathbf{z}_t; \mathbf{z}_{t-1}, \mathbf{d}_t)$ decomposes into

$$Q(\mathbf{z}_t; \mathbf{z}_{t-1}, \mathbf{d}_t) = \prod_{k=1}^{\text{Depth}_Q} \prod_{i \in L_Q^k} Q(\mathbf{z}_t^i; \{\mathbf{z}_t^j\}_{j \in \text{par}_Q^t(i)}, \{\mathbf{z}_{t-1}^j\}_{j \in \text{par}_Q^{t-1}(i)}, \{\mathbf{d}_t^j\}_{j \in \text{par}_Q^d(i)}) \quad [19]$$

237 par_Q^t , par_Q^{t-1} , and par_Q^d are the indices of the parent variables
 238 of variable i in \mathbf{z}_t , \mathbf{z}_{t-1} , and \mathbf{d}_t respectively:

239 Let $\tau_{Q,L_Q^k}^{\text{sample}}$ be the maximum time needed to sample any
 240 variable in L_Q^k , let τ_{Q}^{score} denote the maximum time needed to
 241 score any variable in \mathbf{z}_t , and let τ_P^{score} denote the maximum
 242 time needed to score any variable in \mathbf{z}_t or \mathbf{d}_t .

$$243 \quad \tau_{Q,L_Q^k}^{\text{sample}} = \max_{i \in L_Q^k} \tau_{Q,i}^{\text{sample}} \quad [20]$$

$$245 \quad \tau_Q^{\text{score}} = \max_{i=1}^{|\mathbf{z}_t|} \tau_{Q,i}^{\text{score}}, \quad \tau_P^{\text{score}} = \max_{i=1}^{|\mathbf{z}_t| + |\mathbf{d}_t|} \tau_{P,i}^{\text{score}} \quad [21]$$

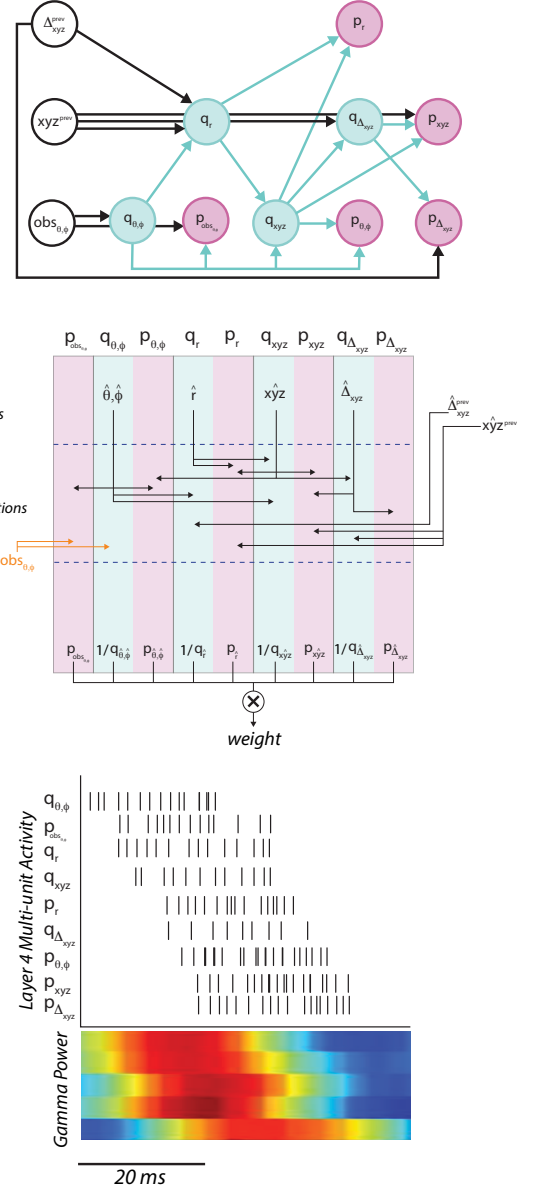


Fig. 4. Meso-scale spiking monte Carlo networks for sampling and scoring multiple variables. (top) the dependence structure of one time slice in the 3D prey tracking model, showing how proposed values and scoring can be interleaved and parallelized. (middle) Each variable's data-driven proposal sampler and model-based scoring circuit is located in its own micro-circuit. The graph structure above is implemented via inter-micro-circuit connections at the appropriate layers. (bottom) This architecture predicts traveling spiking cascades, spreading across dependent columns at the speed of gamma oscillations, that have been confirmed in multiple model organisms. Note that latency is low — only long enough to get a single sampled value from the layer II/III WTAs — because all scoring can be done via massive parallelism. Also note that all variables in the target model can be scored at the same time, regardless of its size.

	Latent Variables	Observed Variables	Size of Spiking Neural Representation Weighted Monte Carlo (this paper)	ENS Codes, Standard PPCs
1D object tracking	$\{x_t, \dot{x}_t\}_t$	$\{d_t^x\}_t$	Sparse: 27 Dense: 5	Dense: 140
2D object tracking	$\{x_t, y_t, \dot{x}_t, \dot{y}_t\}_t$	$\{d_t^x, d_t^y\}_t$	Sparse: 30 Dense: 10	Dense: 2500
Mental Physics Simulation	$\{x_t, y_t, \dot{x}_t, \dot{y}_t, o_t\}_t$	$\{((d_t^{(i,j)})_{i=1}^{10})_{j=1}^{10}\}_t$	Sparse: 38 Dense: 110	Dense: 2500
3D object tracking from 2D observations	$\{x_t, y_t, z_t, \dot{x}_t, \dot{y}_t, \dot{z}_t, r_t, \phi_t, \theta_t\}_t$	$\{d_t^\phi, d_t^\theta\}_t$	Sparse: 160 Dense: 20	Dense: 23,180,062,500
Recursive Concept Learning (Sizes are for $D = 2, M = 10$)	$\bigcup_{h=1}^D \bigcup_{b=1}^{2^{h-1}} \{s^{(h,b)}, \tau_r^{(h,b)}, \tau_c^{(h,b)}, n_1^{(h,b)}, n_2^{(h,b)}\}$	$\{d_t\}_{t=1}^M$	Sparse: 180 Dense: 40	Dense: 5.92704×10^{11}

Table 1. Weighted Monte Carlo spiking 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.

Then the overall latency needed at each timestep of SMC to sample \mathbf{z}_t and estimate the importance weight update w_t/w_{t-1} , τ , is bounded by

$$\tau \leq \sum_{k=1}^{\text{Depth}_Q} \tau_{Q,L^k}^{\text{sample}} + \tau_Q^{\text{score}} + \tau_P^{\text{score}} \quad [22]$$

Observe that while the sampling time grows linearly in the *depth* of Q , the scoring time does not depend on the depth of Q or P , and given a fixed Q -depth, adding more variables to the model also does not increase latency. That is, arbitrarily large models P can be used to provide top-down feedback without linearly* increasing the latency of the circuit, since all the variables in the model can be scored in parallel.

D. Massively parallel macro-scale spiking circuits for real-time sequential Monte Carlo. The cortical columns for particle i corresponding to the Q sampler, Q scorer, and P scorer for each variable in \mathbf{z}_t^i , and the columns for P -scoring each variable in \mathbf{d}_t^i , output a vector \mathbf{z}_t^i of sampled variable values represented in the sparse code, and output two collections of spike counts, $\{N_P^{i,j}\}_{j=1}^{|\mathbf{z}_t^i|+|\mathbf{d}_t^i|}$ and $\{N_Q^{i,j}\}_{j=1}^{|\mathbf{z}_t^i|}$ (these are the counts N_P^z and N_Q^z used in Eqns. 14, and 16 where z is the j th variable in the i th particle). These spike counts encode the set of probability estimates $\{\hat{p}_{i,j}\}_{j=1}^{|\mathbf{z}_t^i|+|\mathbf{d}_t^i|}$ and $\{(\hat{q}_{i,j})^{-1}\}_{j=1}^{|\mathbf{z}_t^i|}$ via the relations

$$\hat{p}_{i,j} = N_P^{i,j}/c_P, \quad \hat{q}_{i,j}^{-1} = N_Q^{i,j}/c_Q \quad [23]$$

While each of these score terms can be approximated reasonably well using a spike count from a single neuron or assembly that is proportional to the value \hat{p} or \hat{q} by a fixed constant c_P or c_Q , the overall importance weight estimate $\hat{w}_i = \prod_{j=1}^{|\mathbf{z}_t^i|+|\mathbf{d}_t^i|} \hat{p}_{i,j} \prod_{j=1}^{|\mathbf{z}_t^i|} \hat{q}_{i,j}^{-1}$ can have enormous dynamic range, and so cannot be represented with a spike count proportional to the value with a fixed constant of proportionality.

*In our model of assemblies as Poisson-Processes, where the $\tau_{P,i}^{\text{score}}$ values are Erlang-distributed, the expected latency will *slightly* increase as more variables are added to the model, because this increases the probability that one of the variables happens to take unusually long to score. In other models of neural scoring assemblies – e.g. as units which either spike or do not during each fixed time window (like (47)) – increasing the number of variables would not at all increase the expected scoring latency.

† To alleviate this issue, when multiplying \hat{p} and \hat{q} terms to compute the importance weight estimates \hat{w}_i , we use a circuit that dynamically chooses the constant of proportionality used in the spiking representation of \hat{w} , and represents this value using a count c_{\log} on a logarithmic scale (so that large ranges of values can be represented). The constant of proportionality is chosen to be the same for each particle, so that the different importance weights $(\hat{w}_i)_i$ can be compared directly to one another. This is implemented using the **MultAutonorm** circuit, which has $N + 1$ output assemblies, one to output the value c_{\log} , and N to output a spiking rate λ_i for each particle:

$$(c_{\log}, (\lambda_i)_{i=1}^N) \sim \text{MultAutonorm}(\{N_P^{i,j}\}_{j=1}^{|\mathbf{z}_t^i|+|\mathbf{d}_t^i|}, \{N_Q^{i,j}\}_{j=1}^{|\mathbf{z}_t^i|})_{i=1}^N \quad [24]$$

These outputs convey the importance weight values in a representation called *neural floating point*, where one line is on a logarithmic scale and the other is on a direct scale (similarly to the way floating-point numbers are represented in digital computers):

$$\hat{w}^i = b^{c_{\log}} \lambda_i \quad [25]$$

where b is the base of the logarithm. The spike rates λ_i can be read-out into spike counts and sent to other parts of the circuit as estimates of $\hat{w}^i/b^{c_{\log}}$, or these assemblies can themselves be used as a sampler (using the same WTA circuit used in the Q samplers for sampling each variable) to choose a particle with probability proportional to its importance weight. This is the key operation needed to perform resampling. Crucially, the **MultAutonorm** will choose a c_{\log} value so that the sum $\sum_i \lambda_i$ of the output spiking rates has a small dynamic range (roughly, a range of 10-100 Hz), so that it is possible to read-out spikes from some of the assemblies relatively quickly, and none of the assemblies are saturated. For the details of this circuit, see Appendix ??.

† For example, if a neuron's maximum rate is 100Hz and the maximum importance weight value is 1, we need 100Hz neurons to correspond to $\hat{w} = 1$. But then in cases where $\hat{w} = 10^{-10}$, we'd need to read-out a spike count from a neuron with rate 10^{-8} Hz. Reading this out at any reasonable precision would require waiting for years, or using an assembly with millions of neurons!

Characteristics of fundamental building blocks of spiking neural Monte Carlo	Experimental evidence for their existence in biological neural systems
Samples are represented via sparse codes from WTAs, but scores are represented via dense codes from MUXes and assemblies	Sparse & dense codes coexist in multiple brain regions (57)
Samples are generated via first-to-spike races between Poisson processes	EPSCs and IPSCs in all neural systems studied to date follow an exponentially-distributed spacing rule, i.e. the number of spikes in a given time window follows a Poisson process (58, 59). Increased synaptic input yields a change in Poisson rate, i.e. probability of race victory scales directly with input.
Winner-take-all samplers rely on fast inhibition.	Ephaptic coupling enables inhibition at the speed of electrical propagation (60, 61)
Scoring units rapidly and accurately count spikes from MUXes and assemblies.	NMDAR plateau potentials are a recently discovered non-decaying synaptic current (50ms) that can stack linearly with other arriving potentials, providing a mechanism for short timescale counting of presynaptic spikes. (62, 63)

occupancy grids, 3D scene graphs (18), and hierarchical object models that adjust resolution based on perceptual uncertainty (77). In both non-human primates and rodents, these models could be compared to fine-grained neural data using relatively well-established techniques (78), and also simultaneously compared to fine-grained behavioral measurements of reaction time and accuracy. It seems appealing to use spiking neural Monte Carlo to integrate empirical constraints that neither Bayesian cognitive models nor artificial neural networks can precisely account for, such as the number of neurons and connectivity of the dorsal and ventral streams, and quantitative latency/accuracy tradeoffs that are observed both neurally and behaviorally. It remains to be seen whether quantitative predictions can be made precisely enough to motivate direct comparison of interventions on model networks to interventions on biological neurons.

Another approach, grounded in cognitive neuroscience, is to seek spatiotemporally coarser behavioral and neural correlates that are easier to measure via neuroimaging techniques. For example, the timecourse of traveling gamma waves, aligned with connectomic data, constrains the dependence structure of data-driven proposals for 3D perception via inverse graphics (22, 79), and also the dependence structure of top-down generative models. Quantitative similarities between weighted Monte Carlo spiking activity could also potentially be compared to behavioral similarity measures and to similarity between stimulus-induced BOLD activity.

Complementary tests can be obtained via smaller model organisms and also via in vitro studies, leveraging their vastly greater levels of observability and control. For example, the spiking model of 3D visual prey tracking from this paper already gives a more detailed causal account of prey capture than previous phenomenological models grounded in Bayesian cognitive science (80), and suggests an approach to depth estimation that could explain recent data on 3D barrier avoidance (81). It also seems appealing to implement spiking neural Monte Carlo circuits using detailed biophysical simulators (82, 83), and to compare implementations against quantitative data from in vitro experiments.

A. Scaling to richer forms of cognition and learning. Bayesian cognitive scientists can directly apply the theory in this paper to make more fine-grained resource-rational models of

causal reasoning. For example, spiking neural circuits for model-driven particle Gibbs MCMC could potentially be fit to population-level inference latency and accuracy. Unlike standard resource-rational models, spiking neural Monte Carlo model fits could incorporate quantitative assumptions about the number of neurons and the level of parallelism that is recruited by the thought process. Spiking neural Monte Carlo can also be used to implement richer models of thinking processes that leverage inference-based value-of-information estimators (84, 85).

Real-time perceptual learning, real-time inference over dynamic data structures (13), and structure learning of probabilistic programs (86) all present additional challenges. For example, although probabilistic programs can learn models of novel objects from just a handful of images (18), they have not yet been shown to simultaneously learn the structure of new generative models for objects and the structure of new efficient data-driven proposals for recognizing those objects. This paper shows how the training data for self-supervised learning of data-driven neural network proposals could be generated, and how probabilistic losses could be rapidly estimated. It does not reveal how to scalably estimate gradients, even for shallow models. It is unclear if biologically realistic deep learning implementations can be developed, or if it will be more fruitful to pursue alternatives based on shallow learning and geometric modeling (87, 88) or online synthesis of provably near-optimal data-driven proposals (89). It is also unclear which truncated representations of latent data structures (such as 3D scene graphs, the plans of other agents, syntactic parse trees and logical representations of grounded natural language semantics, and even symbolic probabilistic program source code representing learned concepts) will lead to practical spiking neural Monte Carlo circuits for real-time inference.

B. Risk-sensitive control, action selection, and planning via inference. Brain computation requires risk-sensitive action selection, not just uncertain inference about world structure. The theory of spiking neural Monte Carlo can be applied to yield neurally mappable architectures for risk-sensitive model-based predictive sensorimotor control, action selection, and planning. This can be achieved via well-known reductions of those problems to probabilistic inference, such as (90–92).

Characteristics of micro-circuits for single-variable importance sampling	Characteristics of biological cortical micro-circuits
WTA units fire only sparsely, at the beginning of each sample/score cycle, to enforce a single race winner.	Layer II/III fires the most sparsely of the cortical layers. (64, 65)
Assembly neurons receive parents' sampled values from parents' WTA samplers, as well as sensory observations.	Layer IV receives intracortical input from Layer II/III (WTAs) of other microcolumns, plus sensory input from the thalamus (66)
WTA neurons control which assembly's spikes pass through the MUX	Layer II/III sends inhibitory projections to Layer V dendrites (67)
Multiplexer units collect spikes from all assemblies, but only output spikes from the assembly chosen by the WTA	Dendritic segmentation of input channels has recently been discovered, providing a biophysically realistic implementation mechanism (68, 69)
Characteristics of meso-scale multivariate importance sampling	Meso-scale characteristics of biological neural systems
Multivariate importance sampling requires latency sufficient for proposal cascades and synchronized scoring, yielding traveling cascades of layer 4 spiking across micro-circuits, at the latency needed for individual samples	Gamma-band (30-100Hz) oscillations and traveling waves (70) are widely observed
Latency τ for sampling and scoring depends on data and parent values, and is thus variable across brain states and regions	Gamma-band oscillations are predicted to have variable frequency, as is widely observed
Assemblies representing more probable values will spike earlier relative to traveling cascades	Phase precession of spiking with respect to gamma oscillations has been observed (71)
Global posterior probabilities can only be read via normalized weights, not directly via spike rates, thus "beliefs" are only implicit and challenging to extract from weighted Monte Carlo spiking	Direct mappings of posterior probabilities and environmental probabilities onto firing rates are not yet strongly supported by empirical evidence (12)
Characteristics of macro-scale sequential Monte Carlo	Macro-scale connectivity and dynamics of biological neural systems
Resampling draws on weights that span multiple microcircuits and returns new particle indices (for new sources of parent variables) back to source microcircuits	There is a cortically-stratified cortico-subcortical-thalamic loop that sends information from lower cortical layers (where SNMC predicts weights are stored) to the basal ganglia, and then back through the thalamus to source cortical layers (52)
Resampling of particles containing multiple variables takes more time than sampling single variables, and is synchronized across particles.	There exist larger-scale oscillations (e.g. alpha, theta) that embed gamma within them (53, 72), with rhythms generated at/near putative thalamic source of resampling (54)
Resampled particles can be used to make high-quality multivariate proposals (25, 73). Higher-quality proposals (i.e. closer to the local posterior distribution, generating higher posterior probability values for latent variables) will have higher weights, and can therefore be proposed more quickly (relative to the cycle at which particles are resampled).	Traveling waves are observed for slower oscillations, e.g. alpha and theta (74). Also, phase precession is observed relative to slower theta rhythms, e.g. for hippocampal place cells, which spike earlier relative to theta oscillations for place cells representing more probable places (75) and of activity in other regions such as the mPFC (76)

418 Structured, "program-like" policies for selecting actions can
419 also be inferred using Monte Carlo inference in probabilistic
420 programs (93). However, despite the potential engineering
421 appeal of these approaches, their potential reverse-engineering
422 value has yet to be evaluated. It also is possible that the brain
423 leverages specialized neural mechanisms for optimizing risk-
424 sensitive action selection beyond what can easily be achieved
425 via Monte Carlo inference.

C. Fundamental limits on scale and efficiency. This paper has
426 focused on sequential Monte Carlo approximations whose ac-
427 curacy is difficult to analyze. However, it is important to
428 note that there are high-dimensional, non-convex energy land-
429 scapes for which approximate sampling is provably efficient,
430 even when optimization is provably NP-hard (94). But these
431 results do not directly address the unreasonable in-practice
432 effectiveness of sophisticated hybrids of data-driven and model-
433

434 driven Monte Carlo for high-dimensional probabilistic pro- 500
435 grams. There is a widespread need for new theory that can 501
436 guide the design of real-time Monte Carlo approximations, ac- 502
437 counting for design tradeoffs between latency, parallelism, and 503
438 variance. One promising approach could be to try to extend 504
439 recent spiking circuit formalisms from theoretical computer 505
440 science to illuminate representational tradeoffs for parallel 506
441 sampling circuits. Consider that given sufficient hardware, 507
442 inference latency can be driven down to the number of serial 508
443 steps needed to generate proposals, regardless of the number 509
444 of variables being inferred or the complexity of the causal 510
445 dependencies among them. But how do the mind and brain 511
446 automatically constrain the structure of their generative mod- 512
447 els, so that “good enough, shallow enough” inference processes 513
448 can be automatically generated?

449 Despite these open questions, the integrative theory intro-
450 duced in this paper offers a candidate unifying framework for
451 simultaneous reverse-engineering of the mind and brain at the
452 computational, cognitive, and neural levels. It has survived an
453 initial battery of empirical tests. We hope it enables neurosci-
454 entists and cognitive scientists to collaborate more closely with
455 each other, and with artificial intelligence researchers, using a
456 modeling formalism for intelligent computation that simultane-
457 ously addresses phenomenological, causal, and computational
458 considerations. It also invites an intriguing question, in this
459 time of excitement and concern about artificial intelligence:
460 what useful intelligent systems can we build, if this brain-like
461 model of computation is implemented using silicon that spikes
462 millions of times faster than biological neurons?

463 Materials and Methods

464 **Spiking Neural Monte Carlo Emulation in the Gen Prob-**
465 **abilistic Programming System..** We have built an extension
466 to the Gen probabilistic programming language (17) that allows
467 scientists to run Monte Carlo inference algorithms in Gen probabilis-
468 tic models and visualize the spiketrains that (one implementation
469 of) Spiking Neural Monte Carlo may produce when running that
470 inference algorithm on a given dataset. This extension also cor-
471 rupts the probability calculations used to implement Monte Carlo
472 inference algorithms, so that the probability computations are per-
473 formed at low-precision, to perfectly match the scheme for obtaining
474 probability-estimates used by our circuits for P -scoring, Q -scoring,
475 and auto-normalization. To produce spiketrains corresponding to
476 the Q -sampling, $1/Q$ -scoring, and P -scoring for a given variable
477 in a given particle on a given inference run, the emulator logs the
478 noisy probability estimates it uses during inference. When asked to
479 produce a spiketrain (for some subset of the neurons in cortex Layer
480 2-6 need for these operations), the emulator samples spiketrains
481 consistent with these probability estimates and sampled values from
482 the conditional distribution over spiketrains under the distribution
483 over spiketrains implied by the circuits described in (the conditional
484 distribution has a closed form due to the nice properties of Poisson
485 Processes). The emulator can also produce spiketrains correspond-
486 ing to the autonormalize-and-multiply operation. In effect this
487 library will allow scientists to (1) experiment with how the quality
488 of probabilistic inferences under a given model vary as the latency,
489 assembly-sizes, and particle counts are varied to change the level
490 of noise in the neural computations, and (2) produce spiketrains
491 consistent with inference under given models, to compare to real
492 biological activity to help them test the empirical predictions of
493 the SNMC theory (and also to help them test what models and
494 proposal distributions are actually present in the brains of different
495 organisms).

496 **Event-driven simulation of Spiking Neural Monte Carlo**
497 **circuits..** To more thoroughly test the feasibility of Spiking Neural
498 Monte Carlo, we have also implemented a compiler which can
499 compile any model and algorithm in a restricted subset of the Gen

probabilistic programming language into a spiking neural network
which runs SNMC inference under the given model using the given
inference algorithm. We used an event-driven simulator to simulate
these neural networks and verify that they produce reasonable
inferences in our simpler models. These simulations are run by
feeding in spiking input events to the simulator at fixed intervals,
to input the observed data to the neural circuit over time, and
inference results are read out of the circuit in the Weighted Monte
Carlo Spiking Code. The neural networks which are produced
consist entirely of neurons whose spiking behavior is described by
an inhomogenous Poisson Process; all the neurons we use have
biologically realistic rates, except for some of the neurons used to
implement logic-gating circuitry (which we expect is implemented
using sub-neuronal mechanisms in the brain).

ACKNOWLEDGMENTS. The authors would like to thank
Mehrdad Jazayeri and Melissa Kline for helpful comments on
drafts.

1. Hv Helmholtz, *Treatise on Physiological Optics*. (Thoemmes Continuum, Bristol), (1856).
2. PS Laplace, *Philosophical Essay on Probabilities*. (1825).
3. S Russell, P Norvig, *Artificial Intelligence: A Modern Approach*. (Prentice Hall), (1995).
4. S Thrun, W Burgard, D Fox, *Probabilistic Robotics*. (MIT Press), (2005).
5. F Saad, Ph.D. thesis (Massachusetts Institute of Technology) (2022).
6. BM Lake, TD Ullman, JB Tenenbaum, SJ Gershman, Building machines that learn and think like people (2016).
7. F Lieder, TL Griffiths, Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behav. brain sciences* **43** (2020).
8. DC Knill, A Pouget, The bayesian brain: the role of uncertainty in neural coding and computation. *TRENDS Neurosci.* **27**, 712–719 (2004).
9. J Fiser, P Berkes, G Orbán, M Lengyel, Statistically optimal perception and learning: from behavior to neural representations. *Trends cognitive sciences* **14**, 119–130 (2010).
10. P Berkes, G Orbán, M Lengyel, J Fiser, Spontaneous cortical activity reveals hallmarks of an optimal internal model of the environment. *Science* **331**, 83–87 (2011).
11. SW Linderman, SJ Gershman, Using computational theory to constrain statistical models of neural data. *Curr. opinion neurobiology* **46**, 14–24 (2017).
12. WJ Ma, M Jazayeri, et al., Neural coding of uncertainty and probability. *Annu. review neuroscience* **37**, 205–220 (2014).
13. N Goodman, V Mansinghka, DM Roy, K Bonawitz, JB Tenenbaum, Church: a language for generative models. *arXiv preprint arXiv:1206.3255* (2012).
14. B Milch, et al., 1 blog: Probabilistic models with unknown objects. *Stat. relational learning*, 373 (2007).
15. V Mansinghka, D Selsam, Y Perov, Venture: a higher-order probabilistic programming platform with programmable inference. *arXiv preprint arXiv:1404.0099* (2014).
16. VK Mansinghka, et al., Probabilistic programming with programmable inference in *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*. pp. 603–616 (2018).
17. MF Cusumano-Towner, FA Saad, AK Lew, VK Mansinghka, Gen: a general-purpose probabilistic programming system with programmable inference in *Proceedings of the 40th acm sigplan conference on programming language design and implementation*. pp. 221–236 (2019).
18. N Gothoskar, et al., 3dp3: 3d scene perception via probabilistic programming in *NeurIPS*. (2021).
19. T Zhi-Xuan, J Mann, T Silver, J Tenenbaum, V Mansinghka, Online bayesian goal inference for boundedly rational planning agents. *Adv. Neural Inf. Process. Syst.* **33**, 19238–19250 (2020).
20. MF Cusumano-Towner, A Radul, D Wingate, VK Mansinghka, Probabilistic programs for inferring the goals of autonomous agents. *arXiv preprint arXiv:1704.04977* (2017).
21. BM Lake, R Salakhutdinov, JB Tenenbaum, Human-level concept learning through probabilistic program induction. *Science* **350**, 1332–1338 (2015).
22. TD Kulkarni, P Kohli, JB Tenenbaum, V Mansinghka, Picture: A probabilistic programming language for scene perception in *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4390–4399 (2015).
23. N Roy, et al., From machine learning to robotics: Challenges and opportunities for embodied intelligence. *arXiv preprint arXiv:2110.15245* (2021).
24. M Cusumano-Towner, B Bichsel, T Gehr, M Vechev, VK Mansinghka, Incremental inference for probabilistic programs in *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*. pp. 571–585 (2018).
25. AK Lew, M Cusumano-Towner, VK Mansinghka, Recursive monte carlo and variational inference with auxiliary variables. *arXiv preprint arXiv:2203.02836* (2022).
26. R Ranganath, S Gerrish, D Blei, Black box variational inference in *Artificial intelligence and statistics*. (PMLR), pp. 814–822 (2014).
27. D Wingate, ND Goodman, DM Roy, LP Kaelbling, JB Tenenbaum, Bayesian policy search with policy priors in *Twenty-second international joint conference on artificial intelligence*. (2011).
28. T Rainforth, TA Le, JW van de Meent, MA Osborne, F Wood, Bayesian optimization for probabilistic programs. *Adv. Neural Inf. Process. Syst.* **29** (2016).
29. M Riesenhuber, T Poggio, Hierarchical models of object recognition in cortex. *Nat. neuroscience* **2**, 1019–1025 (1999).
30. C Zhuang, et al., Unsupervised neural network models of the ventral visual stream. *Proc. Natl. Acad. Sci.* **118**, e2014196118 (2021).
31. DLK Yamins, et al., Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proc. Natl. Acad. Sci.* **111**, 8619–8624 (2014).

- 580 32. CF Cadieu, et al., Deep neural networks rival the representation of primate it cortex for core
581 visual object recognition. *PLoS computational biology* **10**, e1003963 (2014). 664
- 582 33. R Rajalingham, A Piccato, M Jazayeri, The role of mental simulation in primate physical
583 inference abilities. *bioRxiv* (2021). 665
- 584 34. GR Yang, MR Joglekar, HF Song, WT Newsome, XJ Wang, Task representations in neural
585 networks trained to perform many cognitive tasks. *Nat. neuroscience* **22**, 297–306 (2019). 666
- 586 35. CL Alan L. Yuille, Limitations of deep learning for vision, and how we might fix them (2021). 667
- 587 36. U of California Los Angeles., Can artificial intelligence tell a teapot from a golf ball? (2019). 668
- 588 37. T Serre, Deep learning: The good, the bad, and the ugly. *Annu. review vision science* (2019). 669
- 589 38. JS Bowers, et al., Deep problems with neural network models of human vision (2022). 670
- 590 39. DL Yamins, JJ DiCarlo, Eight open questions in the computational modeling of higher sensory
591 cortex. *Curr. Opin. Neurobiol.* **37**, 114–120 (2016) Neurobiology of cognitive behavior. 671
- 592 40. A Kurakin, I Goodfellow, S Bengio, Adversarial examples in the physical world (2016). 672
- 593 41. BBC, Uber in fatal crash had safety flaws say us investigators (2019). 673
- 594 42. Guardian, Tesla driver dies in first fatal crash while using autopilot mode (2016). 674
- 595 43. A Pouget, K Zhang, S Deneve, PE Latham, Statistically efficient estimation using population
596 coding. *Neural Comput.* **10**, 373–401 (1998). 675
- 597 44. WJ Ma, JM Beck, PE Latham, A Pouget, Bayesian inference with probabilistic population
598 codes. *Nat. neuroscience* **9**, 1432–1438 (2006). 676
- 599 45. D Pecevski, L Buesing, W Maass, Probabilistic inference in general graphical models through
600 sampling in stochastic networks of spiking neurons. *PLoS computational biology* **7**, e1002294
601 (2011). 677
- 602 46. D Pecevski, W Maass, Learning probabilistic inference through spike-timing-dependent plas-
603 ticity. *eneuro* **3** (2016). 678
- 604 47. L Buesing, J Bill, B Nessler, W Maass, Neural dynamics as sampling: a model for stochas-
605 tic computation in recurrent networks of spiking neurons. *PLoS computational biology* **7**,
606 e1002211 (2011). 679
- 607 48. R Legenstein, W Maass, Ensembles of spiking neurons with noise support optimal probabilis-
608 tic inference in a dynamically changing environment. *PLOS Comput. Biol.* **10**, 1–27 (2014). 680
- 609 49. VK Mansinghka, EM Jonas, JB Tenenbaum, Stochastic digital circuits for probabilistic infer-
610 ence. *Massachusetts Inst. Technol. Tech. Rep. MITCSAIL-TR 2069* (2008). 681
- 611 50. EM Jonas, Ph.D. thesis (Massachusetts Institute of Technology) (2014). 682
- 612 51. V Mansinghka, E Jonas, Building fast bayesian computing machines out of intentionally
613 stochastic, digital parts. *arXiv preprint arXiv:1402.4914* (2014). 683
- 614 52. P Redgrave, et al., Goal-directed and habitual control in the basal ganglia: implications for
615 parkinson's disease. *Nat. Rev. Neurosci.* **11**, 760–772 (2010). 684
- 616 53. J Chrobak, G Buzsáki, Gamma oscillations in the entorhinal cortex of the freely behaving rat.
617 *J. Neurosci.* **18**, 388–398 (1998). 685
- 618 54. SW Hughes, V Crunelli, Just a phase they're going through: The complex interaction of
619 intrinsic high-threshold bursting and gap junctions in the generation of thalamic α and θ
620 rhythms. *Int. J. Psychophysiol.* **64**, 3–17 (2007). 686
- 621 55. P Thaker, JB Tenenbaum, SJ Gershman, Online learning of symbolic concepts. *J. Math.*
622 *Psychol.* **77**, 10–20 (2017). 687
- 623 56. J Tenenbaum, Rules and similarity in concept learning. *Adv. neural information processing*
624 *systems* **12** (1999). 688
- 625 57. F Rieke, D Warland, R de Ruyter van Steveninck, W Bialek, *Spikes: Exploring the neural*
626 *code*. (MIT Press), (year?). 689
- 627 58. MA Phillips, et al., A synaptic strategy for consolidation of convergent visuotopic maps. *Neu-*
628 *ron* **71**, 710–724 (2011). 690
- 629 59. CF Stevens, Quantal release of neurotransmitter and long-term potentiation. *Cell* **72**, 55–63
630 (1993). 691
- 631 60. Y Zhang, et al., Asymmetric ephaptic inhibition between compartmentalized olfactory recep-
632 tor neurons. *Nat. communications* **10**, 1–16 (2019). 692
- 633 61. A Blot, B Barbour, Ultra-rapid axon-axon ephaptic inhibition of cerebellar purkinje cells by the
634 pinceau. *Nat. neuroscience* **17**, 289–295 (2014). 693
- 635 62. G Major, ME Larkum, J Schiller, Active properties of neocortical pyramidal neuron dendrites.
636 *Annu. review neuroscience* **36**, 1–24 (2013). 694
- 637 63. G Major, A Polsky, W Denk, J Schiller, DW Tank, Spatiotemporally graded nmda spike/plateau
638 potentials in basal dendrites of neocortical pyramidal neurons. *J. neurophysiology* **99**, 2584–
639 2601 (2008). 695
- 640 64. AL Barth, JF Poulet, Experimental evidence for sparse firing in the neocortex. *Trends neuro-*
641 *sciences* **35**, 345–355 (2012). 696
- 642 65. S Sakata, KD Harris, Lamina structure of spontaneous and sensory-evoked population ac-
643 tivity in auditory cortex. *Neuron* **64**, 404–418 (2009). 697
- 644 66. NMM Amorim Da Costa, K Martin, Whose cortical column would that be? *Front. neu-*
645 *roanatomy*, 16 (2010). 698
- 646 67. A Naka, H Adesnik, Inhibitory circuits in cortical layer 5. *Front. neural circuits* **10**, 35 (2016). 699
- 647 68. XJ Wang, GR Yang, A disinhibitory circuit motif and flexible information routing in the brain.
648 *Curr. opinion neurobiology* **49**, 75–83 (2018). 700
- 649 69. GR Yang, JD Murray, XJ Wang, A dendritic disinhibitory circuit mechanism for pathway-
650 specific gating. *Nat. communications* **7**, 1–14 (2016). 701
- 651 70. A Bahramisharif, et al., Propagating neocortical gamma bursts are coordinated by traveling
652 alpha waves. *J. Neurosci.* **33**, 18849–18854 (2013). 702
- 653 71. M Vinck, et al., Gamma-phase shifting in awake monkey visual cortex. *J. neuroscience* **30**,
654 1250–1257 (2010). 703
- 655 72. J Lisman, O Jensen, The theta-gamma neural code. *Neuron* **77**, 1002–1016 (2013). 704
- 656 73. C Andrieu, GO Roberts, The pseudo-marginal approach for efficient monte carlo computa-
657 tions. *The Annals Stat.* **37**, 697–725 (2009). 705
- 658 74. H Zhang, AJ Watrous, A Patel, J Jacobs, Theta and alpha oscillations are traveling waves in
659 the human neocortex. *Neuron* **98**, 1269–1281 (2018). 706
- 660 75. J O'Keefe, ML Recce, Phase relationship between hippocampal place units and the eeg theta
661 rhythm. *Hippocampus* **3**, 317–330 (1993). 707
- 662 76. MW Jones, MA Wilson, Phase precession of medial prefrontal cortical activity relative to the
663 hippocampal theta rhythm. *Hippocampus* **15**, 867–873 (2005). 708
77. K Ok, K Liu, N Roy, Hierarchical object map estimation for efficient and robust navigation in
2021 *IEEE International Conference on Robotics and Automation (ICRA)*. (IEEE), pp. 1132–
1139 (2021). 709
78. M Schrimpf, et al., Brain-score: Which artificial neural network for object recognition is most
brain-like? *BioRxiv*, 407007 (2020). 710
79. I Yildirim, M Belledonne, W Freiwald, J Tenenbaum, Efficient inverse graphics in biological
face processing. *Sci. advances* **6**, eaax5979 (2020). 711
80. AD Bolton, et al., Elements of a stochastic 3d prediction engine in larval zebrafish prey cap-
ture. *ELife* **8**, e51975 (2019). 712
81. H Zwaka, et al., Covert attention to obstacles biases escapes via the mauthner cell. *bioRxiv*
(2022). 713
82. DF Goodman, R Brette, Brian: a simulator for spiking neural networks in python. *Front.*
neuroinformatics, 5 (2008). 714
83. M Stimberg, R Brette, DF Goodman, Brian 2, an intuitive and efficient neural simulator. *eLife*
8, e47314 (2019). 715
84. F Saad, M Cusumano-Towner, V Mansinghka, Estimators of entropy and information via infer-
ence in probabilistic models in *International Conference on Artificial Intelligence and Statistics*.
(PMLR), pp. 5604–5621 (2022). 716
85. M Wick, A McCallum, Query-aware mcmc. *Adv. Neural Inf. Process. Syst.* **24** (2011). 717
86. FA Saad, MF Cusumano-Towner, U Schaechtle, MC Rinard, VK Mansinghka, Bayesian syn-
thesis of probabilistic programs for automatic data modeling. *Proc. ACM on Program. Lang.*
3, 1–32 (2019). 718
87. MA Fischler, RC Bolles, Random sample consensus: a paradigm for model fitting with appli-
cations to image analysis and automated cartography. *Commun. ACM* **24**, 381–395 (1981). 719
88. MF Cusumano-Towner, VK Mansinghka, Using probabilistic programs as proposals. *arXiv*
preprint arXiv:1801.03612 (2018). 720
89. A Lew, M Agrawal, D Sontag, V Mansinghka, Pclean: Bayesian data cleaning at scale
with domain-specific probabilistic programming in *International Conference on Artificial In-*
telligence and Statistics. (PMLR), pp. 1927–1935 (2021). 721
90. P Dayan, GE Hinton, Using expectation-maximization for reinforcement learning. *Neural*
Comput. **9**, 271–278 (1997). 722
91. M Hoffman, A Doucet, N Freitas, A Jasra, Bayesian policy learning with trans-dimensional
mcmc. *Adv. neural information processing systems* **20** (2007). 723
92. N Kantas, J Maciejowski, A Lecchini-Visintini, Sequential monte carlo for model predictive
control in *Nonlinear model predictive control*. (Springer), pp. 263–273 (2009). 724
93. D Wingate, ND Goodman, DM Roy, LP Kaelbling, JB Tenenbaum, Bayesian policy search
with policy priors in *Twenty-second international joint conference on artificial intelligence*.
(2011). 725
94. YA Ma, Y Chen, C Jin, N Flammarion, MI Jordan, Sampling can be faster than optimization.
Proc. Natl. Acad. Sci. **116**, 20881–20885 (2019). 726

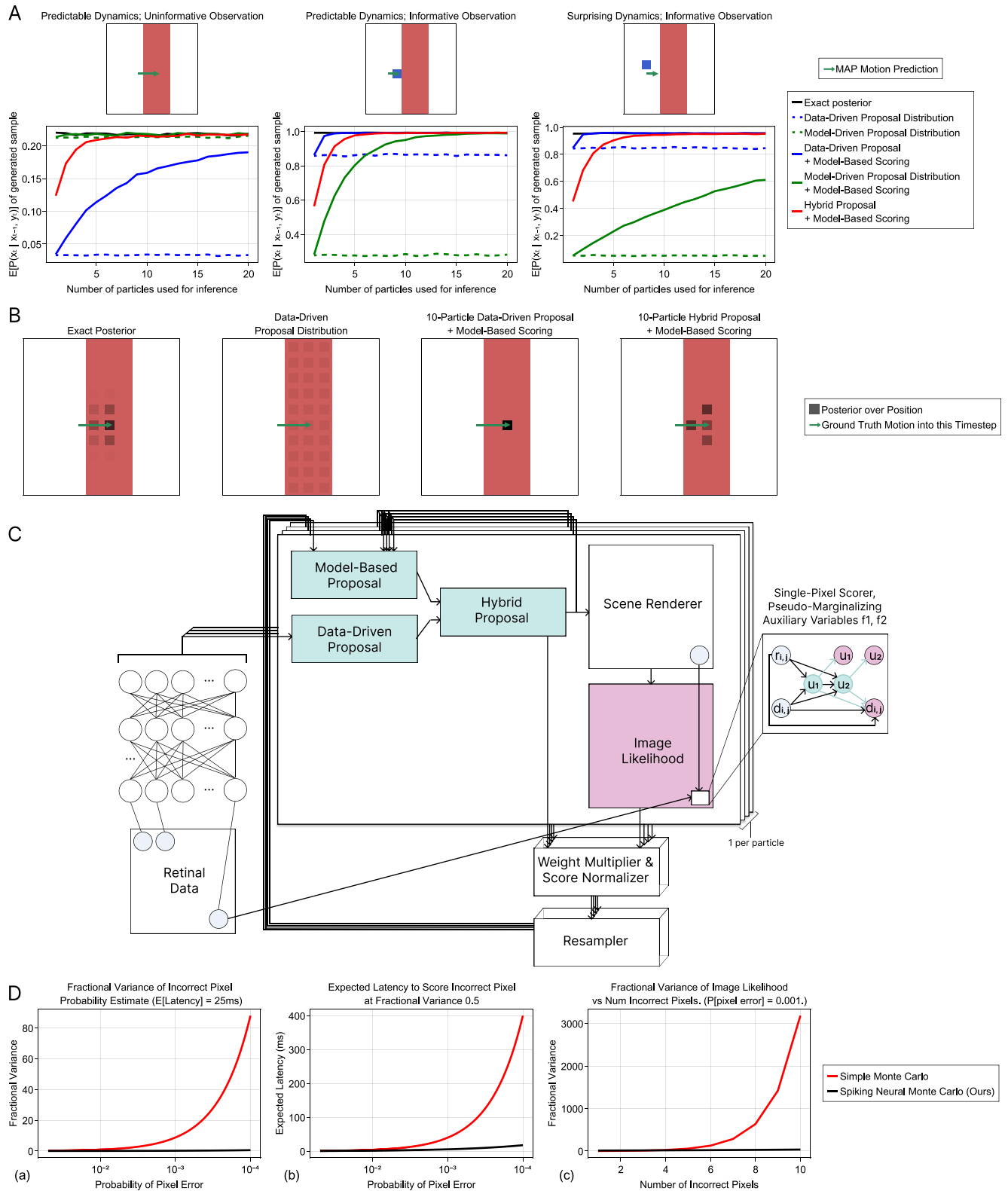


Fig. 7. Scaling up to mental physics simulation and inverse 2D graphics. (A) Hybrids of data-driven and model-driven inference scale better (top, red) than either data-driven or model-driven inference on its own. (B) The hybrid proposal closely approximates the exact Bayes filter, whereas the data-driven proposal struggles when the distinguished object is not visible. (C) The spiking neural Monte Carlo model combines a data-driven neural network (for bottom-up proposals) with top-down model-based inference. Accurate scoring of low-probability data (arising either when data is noisy or the internal model has large errors) is handled by nesting spiking neural Monte Carlo circuits for fast data-driven inference-based scoring (inset on right) within slower spiking neural Monte Carlo circuits for updating scene variables. (D) Inference-based scoring, in which spiking neural Monte Carlo inference over auxiliary variables is used to estimate rare event probabilities, scales to much lower probability data than simple Monte Carlo.

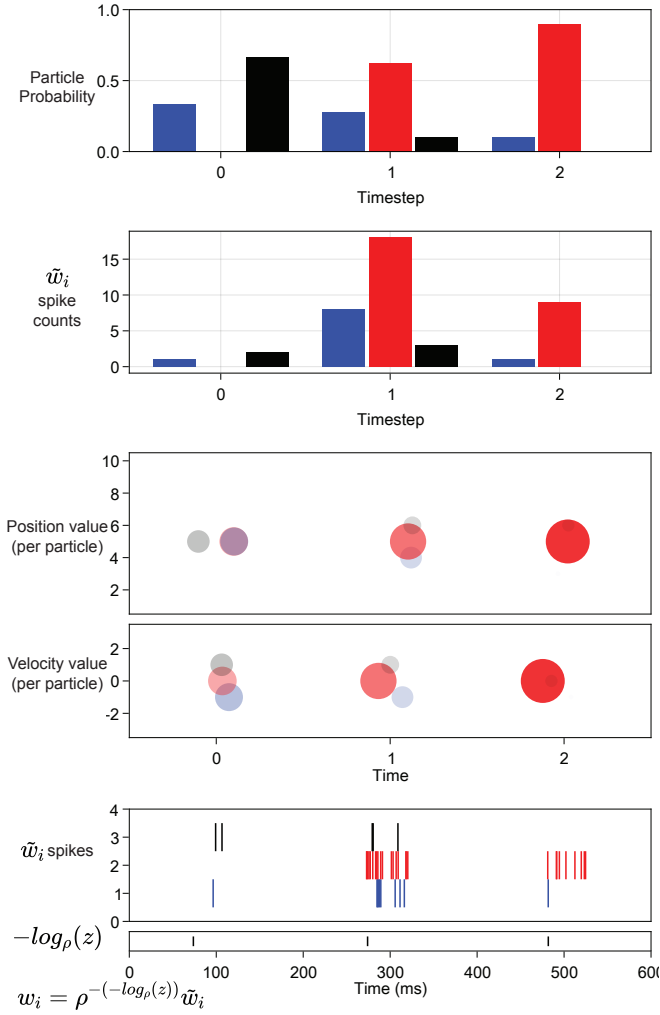


Fig. 8. Exact decoding of weights and probabilities (top) in a dynamically weighted Monte Carlo spiking code (bottom) requires a time-varying non-linear decoder. Spike rate does not always correlate with probability. For example, the probability of the red particle (top) increases significantly from timestep 1 to timestep 2, even when spike count (middle) drops so significantly that it is visible on the spike raster (bottom). Converting importance weights to probabilities requires normalization against the whole set of weights, i.e. $p^t(x_i) = f^t(w_i^t)$ with $f^t(w_i^t) = \frac{w_i^t}{\sum_j w_j^t}$. Thus if all other particles lose nearly all of their weight, the remaining particle's probability will increase, even if its weight drops somewhat significantly.

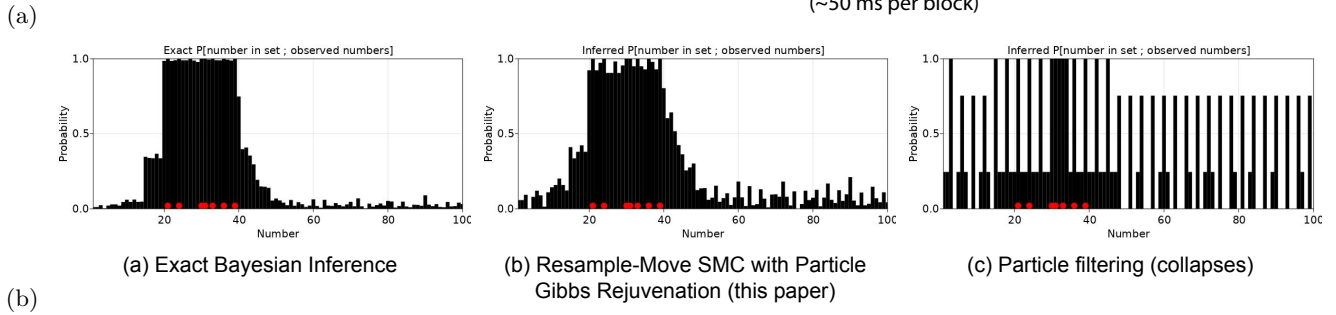
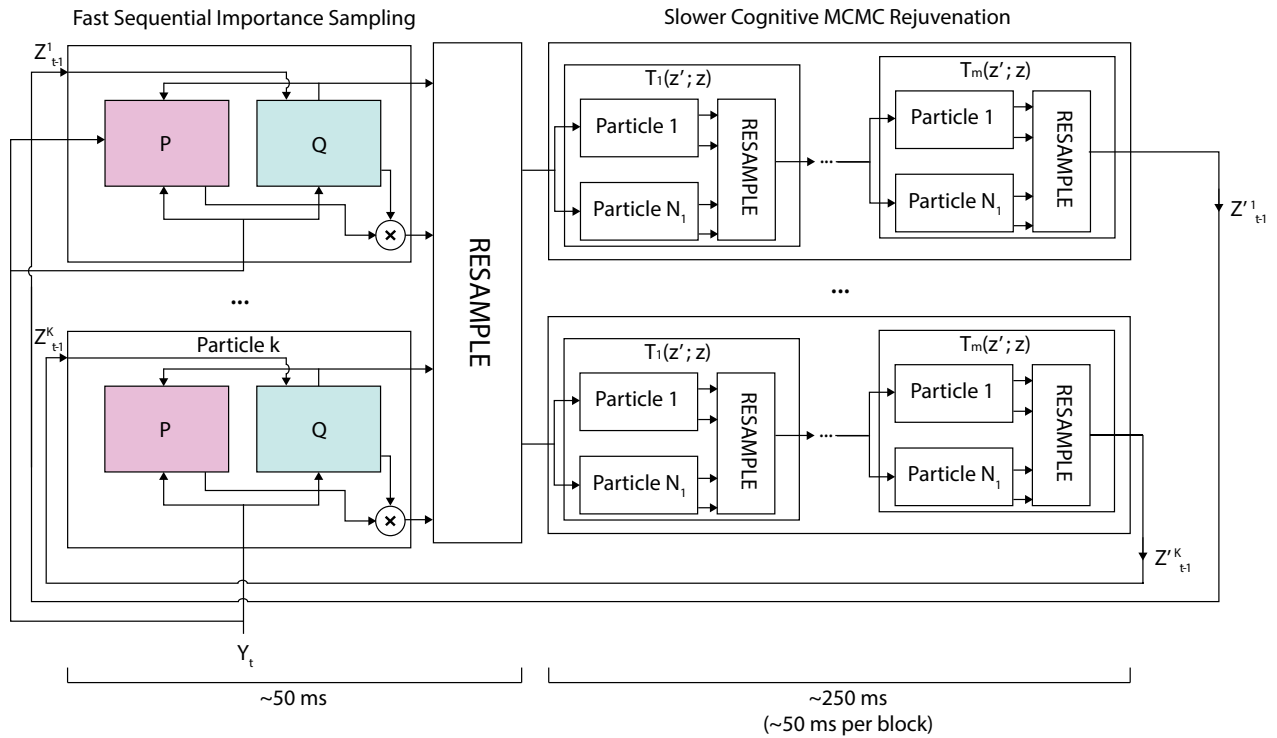


Fig. 9. Online Bayesian concept learning via hybrids of data-driven proposals with iterative, model-based MCMC. (top) Inference relies on fast data-driven proposals followed by model-based MCMC, using proposals that update small subsets of highly-coupled variables via sequential Monte Carlo with multiple particles. (bottom) This architecture makes it possible to explore a broad space of resource-rational, neurally mappable online approximations (55) to exact Bayesian inference (bottom, left) in a classic model of human concept learning (56). The spiking neural Monte Carlo circuit presented here (bottom, middle) better matches both exact Bayesian inference and the behavioral data from (55) than standard particle filtering (bottom, right), which fails to converge in this large hypothesis space.