# Bayesian inference of circular variables in ring attractor networks

Anna Kutschireiter[1]*, Melanie A Basnak[1], Jan Drugowitsch[1]

[1]Harvard Medical School, Department of Neurobiology, 200 Longwood Avenue, Boston, MA-02115, United States.

*Correspondence: anna_kutschireiter@hms.harvard.edu, jan_drugowitsch@hms.harvard.edu

# Summary

Working memories are thought to be held in attractor networks in the brain. Because working memories are often based on uncertain information, memories should ideally come with a representation of this uncertainty for strategic use in behavior. However, the attractor states that hold these memories in attractor networks commonly do not represent such uncertainty. Focusing here on ring attractor networks for encoding head direction, we show that these networks in fact feature all the motifs required to represent uncertainty in head direction estimates. Specifically, they could do so by transiently modulating their overall activity by uncertainty, in line with a circular Kalman filter that performs near-optimal statistical circular estimation. More generally, we show that ring attractors can perform near-optimal Bayesian computation if they can flexibly deviate from their attractor states. Finally, we show that the basic network motifs sufficient for such statistical inference are already known to be present in the brain. Overall, our work demonstrates that ring attractors can in principle implement a dynamic Bayesian inference algorithm in a biologically plausible manner.

**Keywords**
Working memory; Bayesian inference; Ring attractor networks; Head direction neurons; Kalman filtering; Population coding; Drosophila; Central complex;

# Introduction

Many brain functions - including motor control, classification, and pattern completion - have been attributed to attractor networks, and they have proven particularly useful in modeling working memory[1,2]. More specifically, these networks support neural population activity patterns that persist even in the absence of inputs, endowing them with the ability to retain past information across time[3]. A change in the memory's content then corresponds to a change in the network's population activity pattern. At these *attractor states*, the networks only store "point estimates" of these memories, without an associated sense of uncertainty. As this stands in conflict with the observation that memories include a sense of uncertainty (e.g., refs[4,5]), do we need to discard the

1

34  idea of memories being stored in attractor networks? Our work shows that this does not need to
35  be the case.

36  A ring attractor is a special case of an attractor whose set of stable activity profiles forms a ring in
37  neural activity space and thus has the ability to represent circular variables. Head direction (HD)
38  is a classic example of a circular variable that is encoded by a ring attractor network in the brain[6].
39  Many features of mammalian HD neurons are highly suggestive of ring attractors[7–11]. Moreover,
40  recent work has revealed HD cells in the *Drosophila* brain, which not only function as a ring
41  attractor, but also form a topographic map of HD[12–14]. Importantly, the brain often estimates HD
42  under conditions of high uncertainty -- e.g., in unfamiliar environments, or in darkness. Ideally,
43  these HD networks would respond differently to a new piece of information, depending on the
44  current level of uncertainty in the HD estimate. Such an uncertainty-weighted response is a
45  hallmark of Bayesian inference[15]. How exactly the ring attractor networks that track HD could
46  implement Bayesian inference without an explicit notion of uncertainty, however, remains largely
47  unknown.

48  To address potential neural mechanisms for doing so, we took a normative modeling approach,
49  and established how ring attractor networks could maintain and update uncertainty along with the
50  encoded estimate (see Fig. 1 for an overview of our approach). Specifically, we first asked how
51  uncertain HD estimates ought to be updated from unreliable information, irrespective of how these
52  estimates are encoded in the activity of neural populations. We then combined the resulting
53  Bayesian ideal observer model with a neural representation of uncertain HD estimates to arrive
54  at a neural network architecture that can well-approximate the required computations.
55  Interestingly, this network has the general connectivity structure of a ring attractor network.
56  However, its ability to perform near-Bayesian inference depends on its connectivity strengths. A
57  tightly connected, "strict" attractor network performs worse than a weakly connected, "loose"
58  attractor network. This is because strict attractor networks rapidly decay back to their attractor
59  states, while loose attractor networks can persistently deviate from these states. As we show,
60  these deviations are essential to perform the required Bayesian computations. Nonetheless, the
61  networks do not need to be finely tuned to achieve close-to-optimal HD tracking performance.
62  Indeed, a large range of loose networks can adequately combine uncertain HD estimates with
63  unreliable sensory information. Lastly, we showed that model ring attractors can implement
64  dynamic Bayesian inference even after we incorporate constraints from neural connectivity data.
65  In summary, our work provides a principled theoretical foundation for how attractor networks can
66  maintain a sense of uncertainty in their memories, even without an explicit notion of uncertainty.
67  Although we focus on HD encoding as a concrete example, our results are potentially also
68  relevant to other ring attractors in the brain (e.g., the grid cell representation of an animal's path).
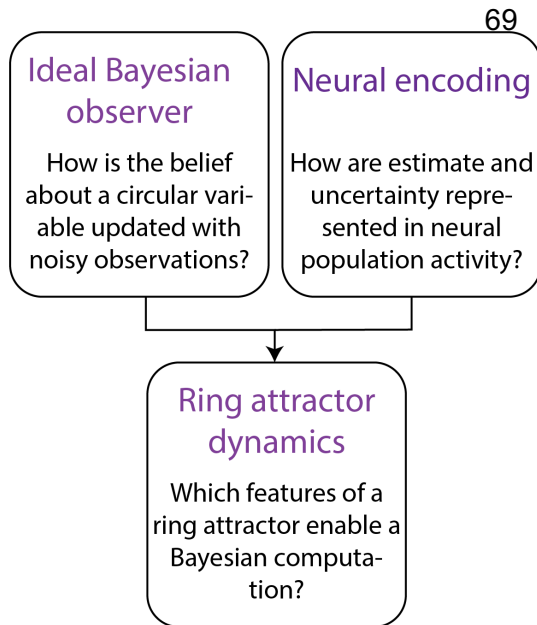
2

69

**Ideal Bayesian observer**

How is the belief about a circular variable updated with noisy observations?

**Neural encoding**

How are estimate and uncertainty represented in neural population activity?

**Ring attractor dynamics**

Which features of a ring attractor enable a Bayesian computation?

**Figure 1.** Our approach combines an ideal Bayesian observer model for circular variables with a 'bump' encoding to derive ring attractor dynamics that perform a Bayesian computation.

83

# Results

## Circular Kalman filtering: a Bayesian ideal observer model for tracking circular variables

We first focus on how uncertain memory ought to be updated from new sensory information irrespective of how this memory is encoded in the activity of a neural network. We do so with the example of HD tracking, by deriving the statistically best HD estimate from a continuous stream of unreliable absolute and relative HD information. This results in a generic algorithm - the circular Kalman filter - that tells us how an estimate of HD, or of any other dynamic circular variable (e.g., visual orientation, time of day, etc.), ought to be updated over time, and the role uncertainty plays in these updates. In the sections that follow we ask how this algorithm can be implemented by neural networks, and analyze the properties of these networks.

HD estimates are informed by two qualitatively different types of sensory inputs (Fig. 2a). Relative HD observations (or **angular velocity observations**), e.g., vestibular or proprioceptive signals, provide information about changes in HD. As they tend to be noisy, integrating them over time results in gradual error accumulation, and a HD estimate that increasingly deviates from the true HD. **Absolute head direction observations**, such as the position of a visual landmark, provide direct HD information that can be used to re-calibrate the HD estimate. Since these observations are also noisy, they should be combined with the internal HD estimate according to their respective reliabilities.

Here, we use dynamic Bayesian inference to properly handle the uncertainties arising from the aforementioned unreliable sensory inputs. We assume access to both angular velocity

3

105 observations $v_t \in R$ and absolute HD observations $z_t \in [-\pi, \pi]$, which provide noisy information
106 about true angular velocity $\dot{\phi}_t \in R$ and HD $\phi_t \in [-\pi, \pi]$, respectively. Specifically, angular velocity
107 observations are corrupted by Gaussian noise that limits the precision of these observations (with
108 precision $\kappa_v$, larger $\kappa_v$ = more reliable), while absolute HD observations are corrupted by von
109 Mises noise with precision $\kappa_z$, the Gaussian equivalent for circular variables. Dynamic Bayesian
110 inference accounts for uncertainties arising from these noisy observations, by forming a **posterior**
111 **belief** of HD $p(\phi_t | z_{0:t}, v_{0:t})$ that is continuously updated in light of new incoming sensory
112 evidence. Importantly, this belief constitutes a whole probability distribution, rather than a single
113 point estimate, which automatically includes a measure of uncertainty around the best HD
114 estimate[15,16].

115 Estimating circular variables, such as HD, precludes the use of standard dynamic Bayesian
116 inference schemes such as the Kalman filter[17,18] to update the posterior belief $p(\phi_t | z_{0:t}, v_{0:t})$ over
117 time. Instead, statistical inference turns out to be analytically intractable[19] and needs to be
118 approximated (see Methods). Here, we approximate this belief at each point in time by a von
119 Mises distribution, $p(\phi_t | z_{0:t}, v_{0:t}) \approx VM(\mu_t, \kappa_t)$, which is fully characterized by its mean $\mu_t$, which
120 is the current best HD estimate, and its precision $\kappa_t$, which measures the estimate's certainty (Fig.
121 2b). As these two posterior, or belief, parameters fully specify the HD belief, updates of the belief
122 in light of sensory evidence simplify to updating these two parameters. We derived the parameter
123 update dynamics by a technique called projection filtering[20,21], resulting in

$$d\mu_t = v_t dt + \frac{\sqrt{2\kappa_z dt}}{\kappa_t} \sin(z_t - \mu_t), \quad (1)$$

$$d\kappa_t = -\frac{f(\kappa_t)}{2\kappa_v} \kappa_t dt + \sqrt{2\kappa_z dt} \cos(z_t - \mu_t). \quad (2)$$

126 Here, $f(\kappa_t)$ is a monotonically increasing nonlinear function that controls the speed of decay of
127 one's certainty $\kappa_t$ (see Eq. (10) in Methods). Equations (1) and (2) together define an algorithm
128 that we call the **circular Kalman filter** (circKF)[21]. The circKF provides a general solution for
129 estimating the evolution of a circular variable over time from noisy data.
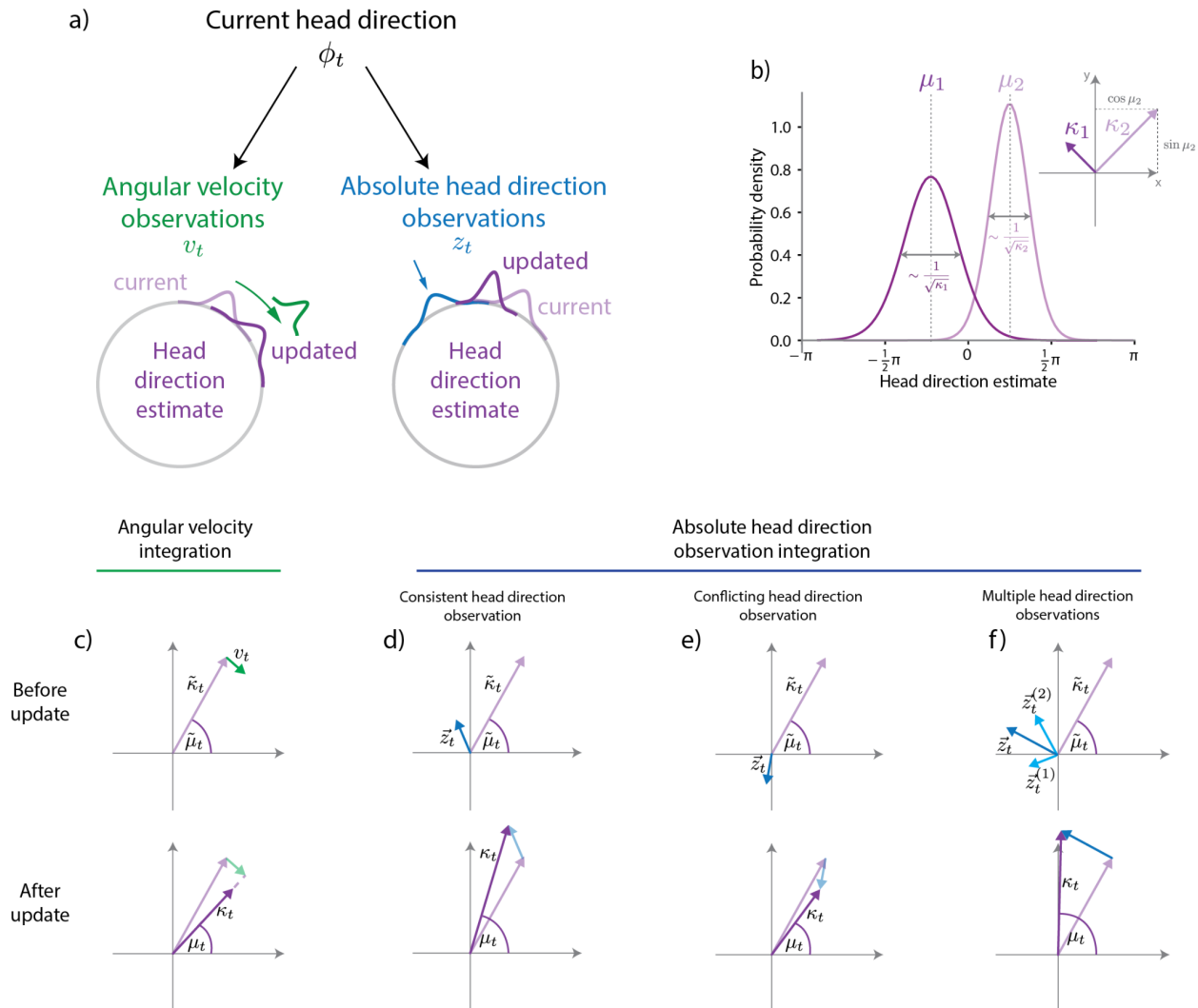
130 To provide intuition for the filter's operation, let us represent the belief parameters in their polar
131 coordinate form as a vector on the 2D plane (Fig. 2b inset). Then, the vector's direction determines
132 the mean HD estimate $\mu_t$, and its length the precision $\kappa_t$. Let us now consider how this vector is
133 updated in light of angular velocity and absolute head direction observations.

134 **Angular velocity observations.** Without absolute head direction observations, i.e., $\kappa_z = 0$, the
135 HD estimate $\mu_t$ is fully determined by integrating angular velocity observations $v_t$ over time, i.e.,
136 angular path integration (Fig. 2c). In our vector representation, angular velocity information (the
137 $v_t dt$ term in Eq. (1)) rotates the brain's HD belief by adding a tangential vector (Fig. 2c, green
138 vector). The increasing error from such angular path integration comes with an associated drop
139 in the belief's certainty $\kappa_t$ ($f(\kappa_t)$-related term in Eq. (2)), which causes the belief vector to shrink
140 (Fig. 2c, bottom). Interestingly, angular velocity observations always decrease certainty. The
141 decrease might be more modest for more precise velocity observations (i.e., $\kappa_v$ large), but
142 nonetheless persists. Thus, if only angular velocity observations are present, the posterior

4

143   certainty $\kappa_t$ will inevitably decay towards zero (uniform posterior distribution, i.e., complete lack
144   of knowledge), with a speed of decay that is determined by the angular velocity observation's
145   "informativeness" $\kappa_v$.

146   **Absolute head direction observations.** Absolute HD information, like observing a visual
147   landmark, directly informs about the current HD, and thus calibrates the internal HD estimate. To
148   weigh the reliability of such information against the current HD estimate's certainty, its impact in
149   Eqs. (1) and (2) is scaled by $\kappa_z$: if the cue's reliability $\kappa_z$ is large, the observation $z_t$ will
150   substantially change the mean $\mu_t$ towards the direction of the cue. Conversely, if the current
151   certainty $\kappa_t$ is large compared to the cue's reliability, an absolute HD observation $z_t$ will hardly
152   update the existing estimate. In vector form, this weighting by reliability corresponds to adding an
153   absolute HD information vector to the current belief vector (Fig. 2d/f; see Methods). The direction
154   and length of this HD information vector are determined by the observation's position $z_t$ and
155   reliability $\kappa_z$, respectively (Fig. 2d, blue vector). Depending on how well the observation is aligned
156   with the current belief (as measured by the cosine in Eq. (2)), the certainty $\kappa_t$ can either increase
157   or, in the case of a strongly conflicting stimulus, even *decrease* (Fig. 2e). This interesting result is
158   a consequence of the circular nature of the inference task[22], and stands in contrast to the Kalman
159   filter where absolute information *always* increases the estimate's certainty[23]. It is thus a key
160   distinction between the Kalman filter and the circKF.

161   In a dynamic setting, both angular velocity and absolute HD observations are available as a
162   continual stream. That is, at every point in time, the belief is updated according to Eqs. (1) and
163   (2). In summary, angular velocity observations rotate the HD estimate and reduce certainty.
164   Absolute HD observations, in contrast, update the HD estimate weighted by their reliability, and
165   either increase certainty (if compatible with the current belief) or reduce certainty (if strongly
166   conflicting with the current belief). These operations are continuously repeated to bring the current
167   belief in line with the latest observations.

**Figure 2. Tracking circular variables with the circular Kalman filter.**

The circular Kalman filter performs dynamic Bayesian inference for circular variables. Its operation is illustrated here for tracking HD.

a) Two different types of observations inform the brain's estimated head direction $\phi_t$: angular velocity observations $v_t$ (green) provide noisy information about the true angular velocity $\dot{\phi}_t$, with precision $\kappa_v$ (larger = more reliable), and absolute HD observations $z_t$ (blue) provide noisy information about the true HD $\phi_t$, with precision $\kappa_z$ (larger = more reliable).

b) At every point in time, the belief $p(\phi_t | v_{0:t}, z_{0:t})$ about HD is approximated by the unimodal von Mises distribution, the Gaussian equivalent for circular variables. It is fully characterized by its mean parameter $\mu_t$, which determines the position of the mode, and its precision parameter $\kappa_t$, which determines our belief's certainty. Interpreted as the polar coordinates in the 2D plane, these parameters provide a convenient vector representation of the belief (inset).

c) Angular velocity observations $v_t$ rotate the current belief vector in the direction of the observations (angular path integration). Error accumulation from angular path integration

6

184 comes with an associated drop in certainty and a corresponding drop in the vector's length
185 (top vs. bottom).
186 d) Integrating absolute HD observations corresponds to adding the absolute HD observation
187 vector (cyan) to the current belief vector (purple).
188 e) Absolute HD observations that are in conflict with the current belief (e.g., >120deg from
189 the current estimate) result in a shortening of the belief vector (top vs. bottom) and an
190 associated reduction of the belief's certainty.
191 f) Integration of multiple absolute HD cues, such as wind and vision, can be considered as
192 a sum of multiple observation vectors.

## Neural encoding of HD estimate and uncertainty

194 To link our ideal observer model to neural networks, we need to specify how the model's belief
195 might be encoded by this activity pattern. In other words, we need to link our "algorithmic model"
196 to a network model. Consider a ring attractor network where the peak of a localized increase in
197 activity, or *bump*, encodes the estimate $\mu_t$ of the circular variable -- here, HD[7,24]. Here we assume
198 that the bump's amplitude scales with the encoded certainty $\kappa_t$. This assumption is supported by
199 some experimental evidence from the head direction system[10,25,26]. In any network where this
200 assumption is correct, the activity of a neuron i with preferred head direction $\phi_i$ can be written as
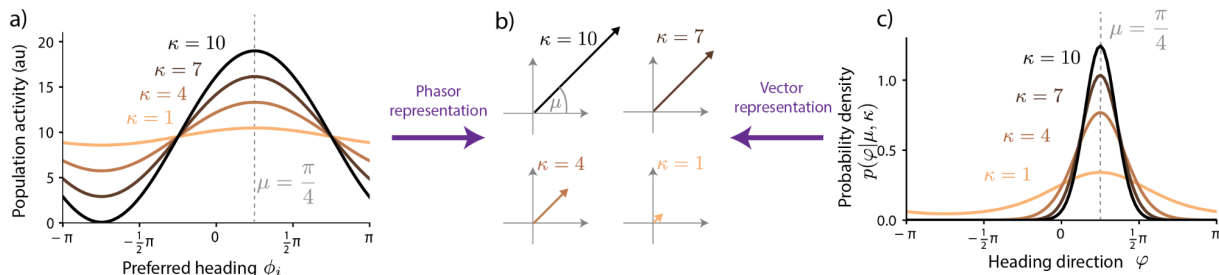201 (Fig. 3a)

$$r_t^{(i)} = \kappa_t \, cos \, (\phi_i - \mu_t) \, + \text{other components} \qquad (3)$$

203 where $\mu_t$ and $\kappa_t$ are the encoded belief's mean and certainty, and the "other components" might
204 be a finite activity baseline or minor contributions of higher-order Fourier components to the
205 activity. Note that Eq. (3) does *not* imply that the tuning curve must be cosine-shaped. Rather, it
206 implies that the cosine component of the tuning curve is modulated by certainty. This is
207 satisfied, for example, by any unimodal bump profile (as the one in Fig. 2a) whose overall gain
208 is governed by certainty. A particularly interesting case that matches Eq. (3) is a linear
209 probabilistic population code[27,28] with von Mises-shaped tuning curves and independent Poisson
210 neural noise (see SI text and Fig. S1).

211 This simple encoding scheme allows the network to encode both mean estimate $\mu$ and associated
212 certainty $\kappa$, as required for implementing the Bayesian update rules (Eqs. (1) and (2)). Moreover,
213 the natural parameters of the von Mises belief, $\theta_1 = \kappa_t \, cos \, (\mu_t)$ and $\theta_2 = \kappa_t \, sin \, (\mu_t)$, can be
214 recovered by taking a weighted sum of the neural population's activity (Methods). This makes
215 these parameters accessible to downstream neurons via simple (linear) neural operations.

216 Interestingly, $\theta_1$ and $\theta_2$ represent the von Mises distribution in terms of Cartesian vector
217 coordinates in the 2D plane, whereas $\mu$ and $\kappa$ are its polar coordinates (cf. Fig. 2b). Such a
218 representation is related to the phasor representation of neural activity[29], which also translates
219 bump position and amplitude to polar coordinates in the 2D plane (Fig. 3b). Since in our model
220 the activity bump is scaled by certainty, the phasor representation of neural activity equals the
221 vector representation of the von Mises distribution (Fig 3b,c).

7

222 Based on our ideal observer model (our "algorithmic model"), we know the vector operations
223 required to implement the circular Kalman filter (Fig 2c-f). Thus, in what follows, we make use of
224 this equality to implement the operations of the ideal observer model through neural dynamics. In
225 other words, we show how the circKF algorithm could be implemented by a neural network.



226

**Figure 3. Encoding the HD belief in neural population activity.**

227
228     **a)** Neural population activity profile (e.g., average firing rate) encoding the HD estimate $\mu = \pi/4$ with different values of certainty $\kappa$. Neurons are sorted by preferred head directions $\phi_i$.
229
230
231     **b)** Vector representation of estimate $\mu = \pi/4$ for different values of certainty $\kappa$. This vector
232 representation can be obtained by linearly decoding the population activity in a) ("phasor
233 representation"). It also corresponds to the vector representation of the von Mises
234 distribution in c), and thus connects neural activities with the probability distributions they
235 encode.
236     **c)** Von Mises probability densities for different values of certainty $\kappa$ and fixed HD estimate
237 $\mu = \pi/4$. Note that, unlike the population activity in a), the density sharpens around the
238 mean with increasing certainty.

## Recurrent neural networks can track Bayesian HD estimates

240 Linking the belief parameters to neural population activity (Eq. (3)) reveals the population activity
241 dynamics required to implement our ideal HD tracking model (Eqs. (1) & (2)). We now ask how
242 these dynamics can be implemented by a recurrent neural network (RNN). We start with an
243 idealized network with a single neural population, similar to many generic ring attractor networks
244 (e.g., refs[3,7]). Later, we will build on this idealized network to construct a more distributed network
245 that satisfies the known constraints of a biological ring attractor that encodes HD.

246 Simple and analytically accessible network dynamics that implement the circular Kalman filter
247 (Eqs. (1) & (2)) are of the form

$$\dot{\mathbf{r}}_t = -\frac{1}{\tau}\mathbf{r}_t - g(\mathbf{r}_t)\mathbf{r}_t + W \cdot \mathbf{r}_t + \mathbf{I}_t^{\text{ext}}$$

248     (4)

249 where $r_t$ denotes a vector of neural activities, with neurons ordered by their preferred head
250 directions $\varphi^{(i)}$, $\tau$ is the network time constant (leak), $W$ is the recurrent connectivity matrix, and
251 $I_t^{ext}$ is a vector of external inputs to the network. The synaptic inhibition nonlinearity $g(r_t)$ is

8

252    closely related to the nonlinearity $f(\kappa_t)$ in Eq. (2): it is tuned such that its output increases with
253    bump amplitude, and thus implements nonlinear global inhibition.

254    The network dynamics in Eq. (4) allow us to attribute specific wiring patterns (or motifs) to the
255    effect they have on the population activity vector, mimicking the transformations required to
256    implement computations in the circKF (Fig. 2c,d). In particular, probabilistic angular path
257    integration is implemented by an interplay between recurrent connectivity ($W$), leak ($1/\tau$), and
258    synaptic inhibition ($g(r_t)$). The matrix of recurrent connectivity $W$ can be divided into symmetric
259    (even) and asymmetric (odd) components (Fig. 4a). The even component holds the bump of
260    activity at its current location in the absence of any other input. Meanwhile, the odd component
261    can push the bump of activity around the ring -- e.g., in response to an angular velocity
262    observation (Fig. 4b). Leak and global inhibition together cause the amplitude of the bump to
263    decay over time (Fig. 4c), corresponding to the progressive decay in certainty in the absence of
264    new HD information. Absolute HD observations enter the network via the external input vector
265    $I_t^{ext}$, in form of a cosine-shaped bump with amplitude modulated by perceptual reliability $\kappa_z$ (Fig.
266    4d). This input activity effectively implements the vector addition required for proper absolute HD
267    observation integration. Then, the external information's weight is determined by the ratio
268    between input amplitude and bump amplitude, in line with the weighting between cue reliability
269    and own certainty required by the circKF. Bump position and amplitude dynamics derived from a
270    network with these basic motifs well-approximate the parameter dynamics of the circKF (Eqs. (1)
271    & (2); Fig. 4e; see Eqs. (13) & (14) in Methods for bump parameter dynamics).
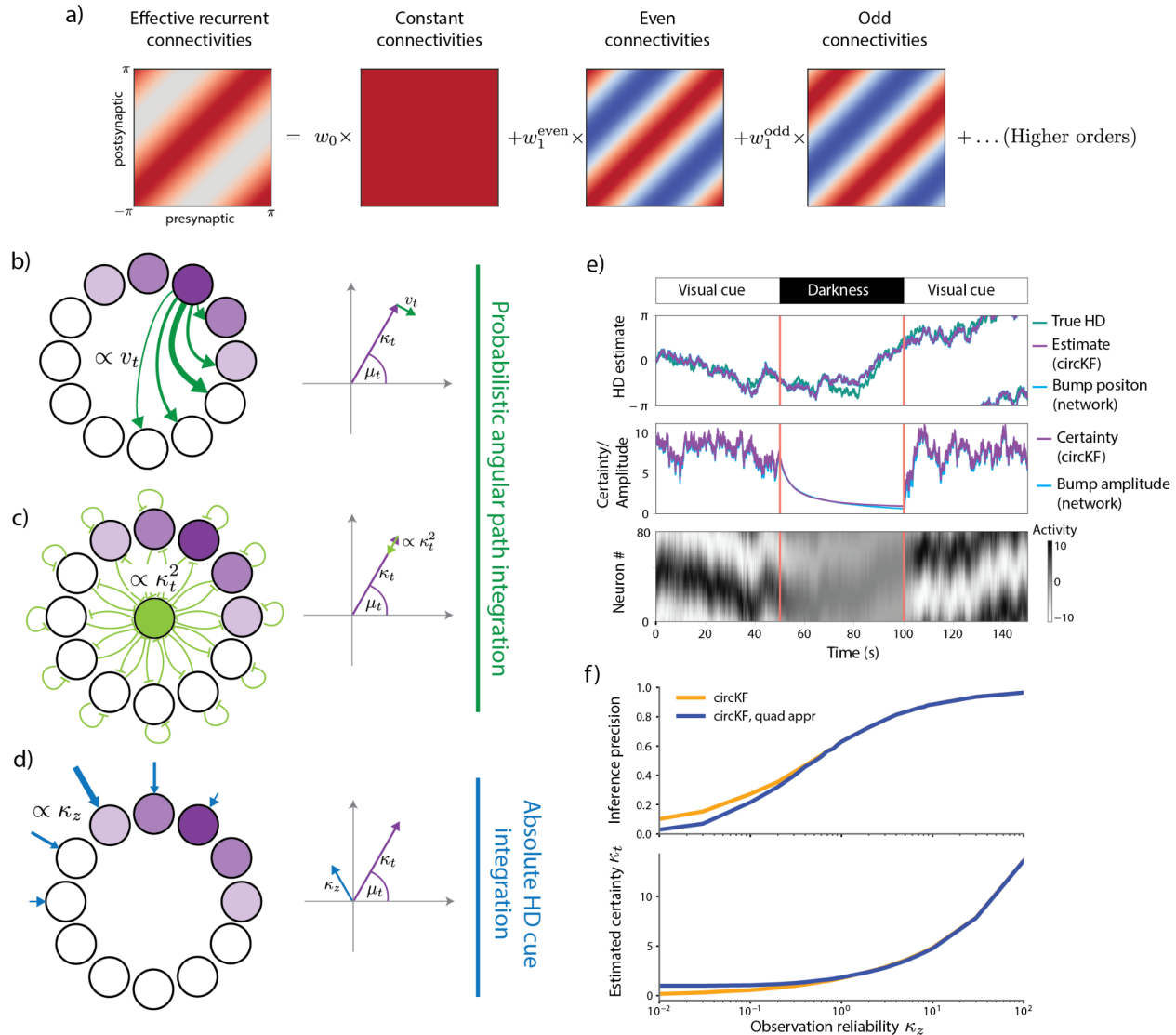
272    In the limit of infinitely many neurons, the network can even be tuned to implement the circKF
273    exactly. Importantly, a network that exactly implements the circular Kalman filter is *not* a ring
274    attractor network: in the absence of external absolute HD input, $I_t^{ext} = 0$, the activity bump decays
275    towards its single attracting state with zero amplitude (Fig. 4e). In contrast, ring attractor networks
276    commonly settle on a constant activity bump with non-zero amplitude (the "attractor state") if input
277    is removed. In our "exact inference" network, activity decay is by design, and reflects the
278    continuously decreasing certainty arising from pure angular path integration in the absence of
279    absolute HD observations.

280    Activity dynamics change qualitatively if we choose $g(r_t)$ such that the second term in Eq. (4)
281    becomes a quadratic function of the bump amplitude (quadratic inhibition). This change in $g(\cdot)$
282    introduces ring attractor states with non-zero network activity, and has the additional advantage
283    of making the network dynamics analytically accessible (see Methods). As a result, we can tune
284    the network parameters such that this network implements **a quadratic approximation** to the
285    circKF. This approximation becomes precise in the limit of large posterior certainties $\kappa_t$. In other
286    words, the bump's amplitude dynamics will correctly reflect the posterior's certainty for large bump
287    amplitudes, but will deviate from it in the small-certainty/small-amplitude limit.

288    Such a Bayesian ring attractor has two operating regimes: a regime close to the attractor state
289    with constant, low bump amplitude, encoding approximately constant certainty, and a dynamic,
290    high amplitude regime away from the attractor state, where the network correctly implements
291    dynamic Bayesian inference. Numerical simulations confirm the existence of these two regimes:
292    the network tracks the HD estimate and its associated certainty just like the circular Kalman filter

293    in the dynamic regime, but features a slightly lower HD tracking precision, and overestimates its
294    confidence, close to the attractor state (Fig. 4f, orange vs. blue; SI Fig. S2 shows that performance
295    is largely independent of ring attractor population size). We will analyze these two regimes further
296    in the next section.

297    In summary, the following three network motifs support the implementation of Bayesian inference
298    in ring attractor networks (Fig. 4b-d): (i) asymmetric (odd) recurrent connectivity with strength
299    modulated by angular velocity observations $v_t$, (ii) global inhibition that is approximately quadratic
300    in bump amplitude, and (iii) a cosine-shaped external input at the position of the absolute HD
301    observation, whose strength is modulated by the reliability $\kappa_z$ of this observation. Motifs (i) and (ii)
302    implement probabilistic angular path integration, whereas motif (iii) updates the network's current
303    HD estimate in light of uncertain absolute HD observations. Interestingly, these motifs are
304    common in many generic ring attractor networks, and have been discussed in terms of their
305    function individually (see e.g. refs[7,30]). Here, we show that, together, they can implement
306    approximate dynamic Bayesian inference for circular variables - inference that becomes more
307    precise in the limit of large amplitudes, away from the attractor state.

**Figure 4. A recurrent neural network implementation of the circular Kalman filter.**

a)  Rotation-symmetric recurrent connectivities (here: neurons are sorted according to their preferred HD) can be decomposed into constant, cosine-shaped (even), sine-shaped (odd) and higher-order frequency components (basis function). Red and blue denote excitatory and inhibitory components, respectively.

b)  Network motifs sufficient to implement the circular Kalman filter (b-d). Rotations of the HD estimate are mediated by sine-shaped (or odd) recurrent connectivities, whose strength is modulated by angular velocity observations.

c)  Decay in amplitude arises from leak and global inhibition.

d)  A cosine-shaped input to the network provides external absolute HD cue input. The strengths of this input is modulated by observation reliability $\kappa_z$.

e)  The dynamics of the network implement the dynamics of the ideal observer's belief, as shown in a simulation of a network with 80 neurons. Here, we assume that vision provides the network with absolute HD information. When a 'visual cue' was present, both absolute

11

323    HD observations and angular velocity observations were available. During 'darkness', only
324    angular velocity observations were available.
325    f)  The network implementation with quadratic leak approximation (circKF, quadratic approx)
326        tracks the HD estimate with the same precision (top; higher = lower circular distance to
327        true HD) as the circular Kalman filter (circKF, Eqs. (1) and (2)) if absolute HD observations
328        are reliable (large $\kappa_z$), but with slightly lower precision once they become less reliable
329        (small $\kappa_z$). This drop co-occurs with an overestimate in the estimate's confidence $\kappa_t$
330        (bottom). Plots are averages over 5'000 simulations (see Methods for simulation details).

## Ring attractors approximate Bayesian inference for HD tracking through amplitude dynamics

333    Our Bayesian ring attractor network qualitatively differs in two ways from classical ring attractor
334    networks for working memory[1,7,31]. First, classical networks are not explicitly designed to represent
335    uncertainty, and therefore assign no interpretation to their bump's amplitude. Second, ring
336    attractors are usually designed to operate close to their attractor states, where the bump
337    amplitude tends to vary little. We now ask how important it is for network activity - including bump
338    amplitude - to deviate from these attractor states to implement Bayesian inference.
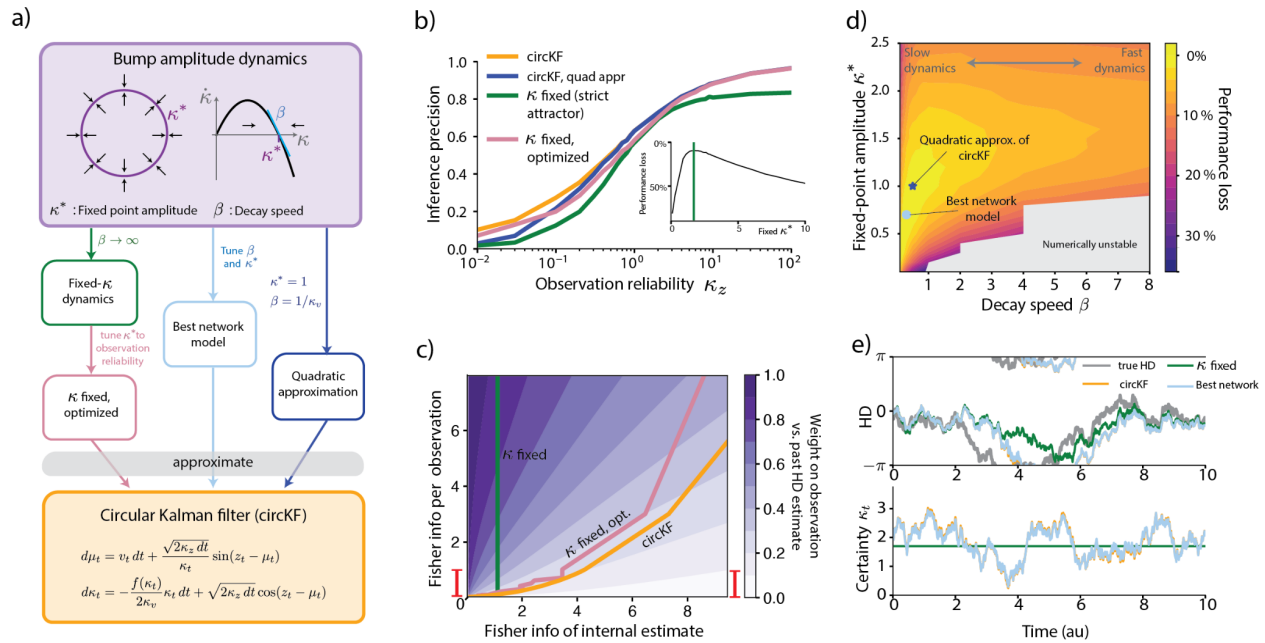
339    Our RNN with quadratic inhibition is a ring attractor network that exhibits attracting states with a
340    finite amplitude. In networks of such structure, bump position changes upon absolute HD inputs
341    in proportion to the ratio between input amplitude (≙ reliability of absolute HD cue) and bump
342    amplitude (≙ own certainty) (see Eq. (13) in Methods). Jointly tuning the network parameters, in
343    particular recurrent weights, network time constant, and inhibitory weights, allows us to change
344    the bump amplitude dynamics (Eq. (15) in Methods) to explore different regimes of network
345    operation. Specifically, we tuned the network parameters to modulate two amplitude
346    characteristics. The first is the attractive amplitude fixed point $\kappa*$ of the population activity bump
347    (specified by the parameter $\kappa*$, which defines both the bump amplitude and the uncertainty it
348    encodes). The second is the effective decay speed $\beta$, which describes how fast the amplitude
349    approaches this fixed point (larger $\beta$ implies faster dynamics, Fig. 5a).

350    In the limit of fast dynamics, $\beta \rightarrow \infty$, the network becomes a "strict attractor" with a bump
351    amplitude that never moves away from its fixed point. As the bump amplitude encodes the HD
352    estimate's uncertainty, such a strict attractor never updates its uncertainty, and consequently
353    lacks proper reliability weighting of absolute HD observations. In general, we expect that such a
354    network is not able to estimate HD as accurately as one that correctly implements Bayesian
355    inference, as it does not properly account for the observation's reliability. Strict attractors with a
356    numerically optimized fixed point amplitude $\kappa*$ (Fig. 5b inset) can still perform HD estimation
357    reasonably well (Fig. 5b, green curve), but perform systematically worse than the tuned network
358    (Fig. 5b, circKF with quadratic approximation, blue curve) or the circKF (Fig. 5b, orange curve).
359    Here, we measure performance by how much the estimate (bump position) deviates on average
360    from the ground-truth HD (circular average distance, see Methods). Adjusting the fixed point
361    amplitude for each level of information reliability *individually* further boosted the network's

362    estimation performance, and effectively re-established proper Bayesian reliability weighting of
363    incoming absolute HD observations, similar to the circKF (Fig. 5c). Even then, strict attractors lack
364    temporal adjustment of their internal certainty estimates, as required by Bayes-optimal evidence
365    integration.

366    We next asked whether we could increase performance by loosening the attractor. Such a
367    relaxation would permit the network to operate farther away from its attracting state, to which it
368    will decay in the absence of absolute HD observations (the dynamics that implement the quadratic
369    approximation of the circKF are a special case of this network, see Fig. 5a). Indeed, the slower
370    dynamics (slower decay speed $\beta$) of such a "loose attractor" boosted overall performance (Fig.
371    5d). In fact, with network parameters tuned numerically to maximize performance, HD estimation
372    performance becomes practically indistinguishable from that of the ideal Bayesian observer (light
373    blue dot in Fig. 5d, light blue line in Fig. 5e). In this regime, the HD estimate and bump amplitude
374    dynamics become almost identical to the dynamics of certainty representations in the circKF. On
375    the other hand, more rigid networks with faster decay speed $\beta$ (such as the strict attractor as an
376    extreme case) clearly deviate from the circKF, despite optimized fixed-point values (Fig. 5e).
377    Interestingly, the optimal network parameters do not necessarily coincide with the quadratic
378    approximation of the circKF, which we found by analytically matching the certainty dynamics in
379    the large-certainty limit rather than by numerical optimization. In fact, a wide range of network
380    parameters lead to a relatively small performance loss (<10%, Fig. 5d). Therefore, accurate
381    parameter tuning might be unnecessary, as long as the network dynamics remain sufficiently
382    slow.

383    Overall, this demonstrates that proper HD estimation relies on weighting absolute HD
384    observations both globally (Figs. 5 b,c), i.e. according to the average level of reliability, and
385    dynamically (Figs. 5d,e), according to the dynamics of one's own certainty. Nonetheless,
386    reasonable performance can be achieved over a wide range of network parameters. This may
387    indicate a "built-in" implicit reliability weighting in attractor networks through their amplitude
388    dynamics. As we just demonstrated, this requires sufficiently slow attractor dynamics around the
389    fixed point and the possibility for deviations from the attractor state. This may explain why ring
390    attractor networks perform evidence integration reasonably well in practice, even though they are
391    unlikely to be precisely tuned to the task.

**Figure 5. Attractor models with slow dynamics approximate Bayesian inference**

**a)** A linear RNN with quadratic inhibition can operate in different regimes. Its bump amplitude dynamics can be characterized by fixed point amplitude $\kappa*$ and decay speed $\beta$. Note that the bump position dynamics is described by the same equation across all compared regimes (Eq. (13) in Methods). However, the position dynamics depend on bump amplitude, whose dynamics differ across regimes (Eq. (15) in Methods). This causes HD tracking behavior to differ across network regimes.

**b)** HD estimation performance as measured by inference precision (as defined by $1 - circVar$, see Methods). Here, the blue curve shows performance of the analytically tuned ring attractor network, implementing the quadratic approximation to the circKF (yellow). For the strict attractor (green curve), we chose $\kappa*$ to numerically maximize performance averaged across all levels of observation reliability, weighted by a prior $p(\kappa_z)$ on this reliability (see Methods). For the optimized, but still strict, network (pink curve), we found the performance-maximizing $\kappa*$ separately for each level of observation reliability.

**c)** The weight with which a single observation contributes to the HD estimate varies with informativeness of both the absolute HD observations and the current HD estimate. We here illustrate this for an absolute HD observation that is orthogonal to the current HD estimate, resulting in the largest possible estimate change ($|z_t - \mu_t| = 90\text{deg}$ in Eq. (1)). The weight itself quantifies how much the observation impacts the HD estimate as a function of how informative this observation is (vertical axis, measured by Fisher information of a 10ms observation) and our certainty in the HD estimate (horizontal axis, also measured by Fisher information) before this observation. A weight of one implies that the observation replaces the previous HD estimate, whereas a weight of zero implies that the observation does not impact this estimate. The close-to-optimal update weight of the circKF (yellow) forms a nonlinear curve through this parameter space. Fisher information per observation is directly related to the observation reliability $\kappa_z$, and the vertical red bar shows the equivalent range of observation reliabilities, $\kappa_z \in [10^{-2}, 10^2]$, shown in panel b.

14

420    Update weights for the tuned network (circKF with quadratic approximation) are not shown
421    as they would be visually indistinguishable from that of the circKF, and only deviate from
422    it for very uninformative observations (see SI Fig. S3).
423  **d)**  Overall inference performance loss (compared to a particle filter; performance measured
424    by avg. inference precision, as in **b**, 0%: same average inference precision as particle
425    filter, 100%: average. inference precision = 0), averaged across all levels of observation
426    reliability (weighted by prior $p(\kappa_z)$, see Methods) as a function of the bump amplitude
427    parameters $\kappa*$ and $\beta$. For too small fixed point amplitudes and too fast dynamics,
428    numerical simulations become unstable (grey area).
429  **e)**  Simulated example trajectories of HD estimate/bump positions of HD estimate/bump
430    positions (top) and certainties/bump amplitudes (bottom). To avoid cluttering, we are not
431    showing the quadratic approximation of the circKF (visually indistinguishable from circKF
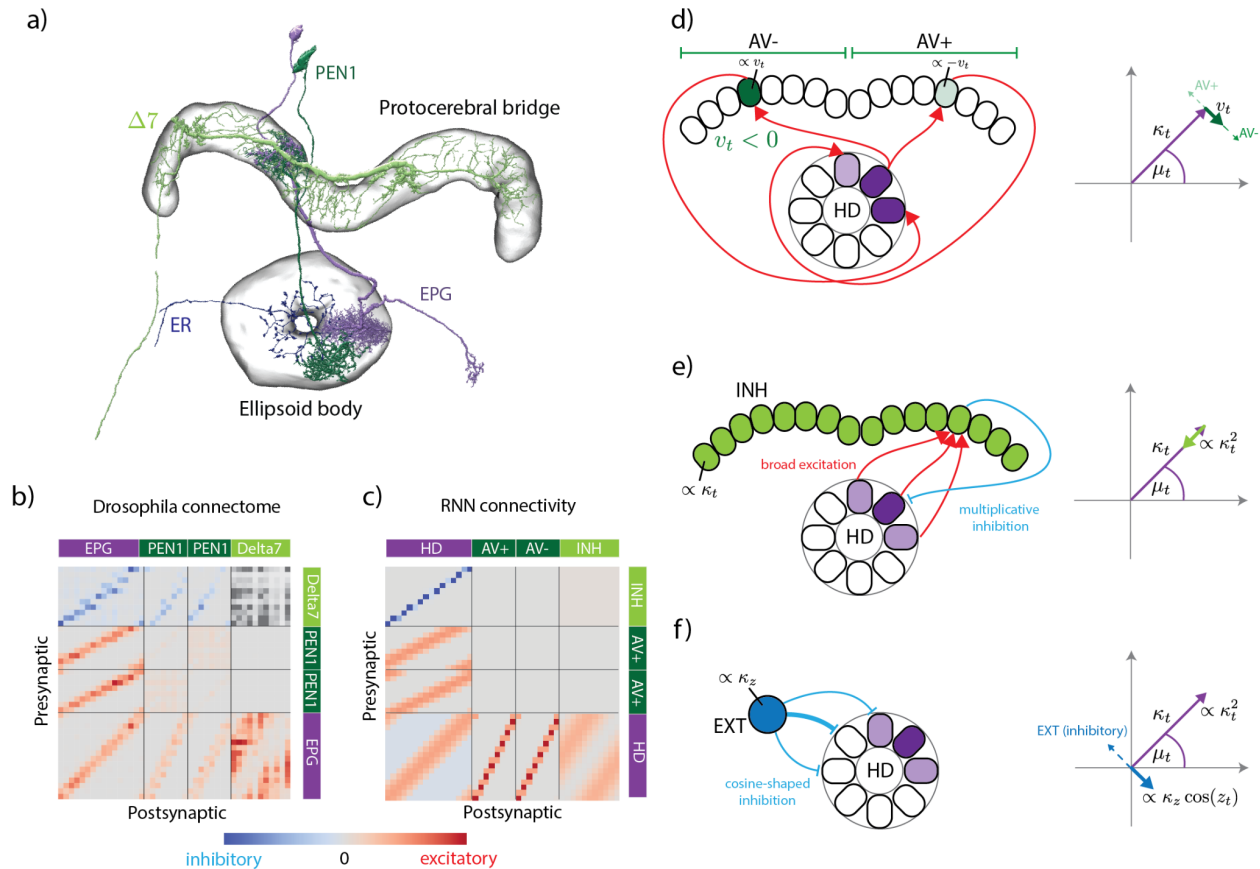432    and best network).

## 433  A biological ring attractor can implement the circular Kalman filter

434  Having established the network motifs sufficient for implementing dynamic Bayesian inference in
435  ring attractor networks, and the network parameter regimes that lead to good HD tracking
436  performance, we finally asked if biological networks are in principle able to implement such
437  inference. A biological implementation is plausible because the critical motifs of our model
438  networks are actually common in many generic ring attractor networks. The most well-studied
439  biological ring attractor network is the HD system of the fruit fly *Drosophila*[13]. Here we show how
440  the motifs of this network -- and, by extension, any biological ring attractor network -- could
441  potentially implement dynamic Bayesian inference.

442  The ring attractor in the *Drosophila* HD system is composed of three core cell types, called EPG,
443  PEN1 and Δ7 neurons[32–34], cf. Fig. 6a-c. Head direction is represented as a bump of neural activity
444  in the EPG population[12]. These neurons are recurrently connected with excitatory PEN1 neurons.
445  When the fly turns, this differentially activates PEN1 neurons in the right and left brain
446  hemispheres, and because PEN1 neurons have asymmetric (shifted) projections back to EPG
447  neurons, they can rotate the bump of EPG activity in accordance with the fly's rotation[14,35]. This
448  motif effectively establishes the velocity-modulated odd recurrent connectivity required to initiate
449  turns in ring attractor networks (Fig. 6d). Moreover, EPG neurons are recurrently connected with
450  inhibitory Δ7 neurons, which establishes broad inhibition (Fig. 6e). Finally, EPG neurons receive
451  inhibitory inputs from so-called ER neurons, which send absolute HD information to EPG
452  neurons[36–38] (Fig. 6f). In summary, the fly's HD system is equipped with the basic motifs to
453  implement a Bayesian ring attractor.

454  To demonstrate that these motifs can in principle implement a Bayesian ring attractor, we
455  analytically tuned the relative connection strength between (rather than within) the populations of
456  our idealized network in Fig. 6c such that the dynamics of the bump parameters in the HD
457  population implement the quadratic approximation to the circKF (see Eqs. (16)-(20) in Methods
458  for network dynamics, SI for derivation). To achieve this, we additionally assumed that the broad
459  inhibition implemented by the inhibitory population (Fig. 6c,e) was achieved by a subtractive signal
460  that resulted from a multiplicative interaction between activities of INH and HD neurons. This

multiplicative interaction achieves the quadratic certainty decay required to approximate the circKF. We found that this network achieves a HD tracking performance indistinguishable to that of our idealized Bayesian ring attractor network (SI Figure S4). Thus, even when we add the constraints dictated by the actual connectivity patterns of neural networks in the brain, the resulting network is still able to implement dynamic Bayesian inference.



**Figure 6. A Drosophila-like network implementing the circular Kalman filter.**

**a)** Cell types in the *Drosophila* brain that could contribute to implementing the circular Kalman filter.

**b)** Connectivity between EPG, *Δ7* and PEN1 neurons, as recovered from the hemibrain:v1.2.1 database[33]. ER neurons were omitted because they only form the inputs to the recurrently connected ring attractor. Here, neurons were grouped according to anatomical region as a proxy for preferred HD, and we used the total number of synaptic connections between two neurons to indicate connection strength. *Δ7* to *Δ7* connectivities are omitted, as the polarity of these connections (inhibitory or excitatory) remains unclear.

**c)** The RNN connectivity profile that implements an approximation of Bayesian inference algorithm is strikingly similar to the connectivity of neurons in the *Drosophila* HD system. To avoid confusion with actual neurons, we refer to the neuronal populations in this idealized RNN as head direction (HD), angular velocity (AV+ and AV-, in reference to the two hemispheres), inhibitory (INH) and external input (EXT) populations.

**d)** Differential activation of AV populations (left/right: high/low) in the two hemispheres as well as a shifted feedback connectivity from AV to HD populations effectively implement

483        the odd (or shifted) connectivity needed to turn the bump position (here: clockwise shift
484        for anti-clockwise turn).

485    **e)**  Broad excitation of the INH population by the HD population, together with a one-to-one
486        multiplicative interaction between INH and HD population, implement the quadratic decay
487        of the bump amplitude needed for the reduction in certainty arising from probabilistic path
488        integration.

489    **f)**  External input is mediated by inhibiting HD neurons with preferred direction *opposite* the
490        location of the absolute HD observation, effectively implementing a vector sum of belief
491        with absolute HD input.

# Discussion

493  We have shown that ring attractor networks - prominent models for working memory of circular
494  variables - can encode and compute with a sense of uncertainty, even when their attractor states
495  are unable to do so. They can achieve this by operating in a dynamic regime away from these
496  attractor states. In this regime, their bump amplitude can vary and thus can encode uncertainty.
497  Such deviations from the attractor state are only possible in loose attractors with sufficiently weak
498  connectivity strengths. Stronger connectivity leads to strict attractors that operate closer to their
499  attractor states and feature worse performance. For a canonical working memory of a circular
500  variable - our sense of head direction - we have shown that network motifs common to ring
501  attractor networks are sufficient to implement the basic computations for dynamic Bayesian
502  inference: (i) angular velocity-modulated odd recurrent connectivities implement incremental
503  changes to the HD estimate, (ii) global inhibition implements the required decay in certainty over
504  time, and (iii) reliability-modulated external input implements reliability-weighted absolute HD
505  integration. We expect these findings to translate to working memories of other circular variables,
506  like those that follow circadian rhythms, or encode memory about visual orientations[5,39]. We
507  further found that close-to-optimal estimation does not require exact tuning of the ring attractor
508  network's connectivities, as long as the networks feature the aforementioned motifs and are
509  flexible enough to deviate from their attractor states. Lastly, we demonstrated that a network with
510  realistic biological constraints still supports the implementation of such a Bayesian ring attractor.
511  Our findings thus suggest that ring attractor models can implement Bayesian computations for
512  working memory.

513  A key element of our approach is the representation of uncertainty as the amplitude of a neural
514  activity bump. This differentiates our work from recent network models that only performed
515  reliability-weighted cue integration at the level of the inputs[40,41], without considering the resulting
516  certainty of the HD estimate. In our framework, this certainty determines the weight with which
517  new external evidence enters the estimate through the bump amplitude. As such, it plays a central
518  computational role for updating the estimate, rather than being a passive measure of
519  precision[25,26]. It predicts that the speed with which the activity bump reacts to changing absolute
520  HD observations should depend on the HD estimation's certainty, and thus bump amplitude: low
521  bump amplitudes (low certainty) should lead to rapid bump shifts, whereas high bump amplitudes
522  (high certainty) show lead to slower ones. Recent experimental evidence[10] suggests that bump

523 amplitude varies in navigating rodents, and this amplitude modulates the speed with which their
524 HD system reacts to changing absolute HD observations - in line with our predictions.

525 By restricting ourselves to an analytically tractable ring attractor network, we were able to almost
526 exactly map the certainty dynamics of the ideal-observer circKF to the bump amplitude dynamics.
527 Having the network implement the circKF rather than a standard Kalman filter fully accounts for
528 the circular symmetry of HD estimation. Thus, unlike previous work[23], our network does not suffer
529 from imprecise inference once absolute HD observations strongly deviate from the current HD
530 estimate. As a result, it yields fundamentally different predictions for strongly conflicting absolute
531 HD direction cues (Fig. 1e). Specifically, since in the circular Kalman filter a conflicting absolute
532 observation (>90 deg from the current estimate) could yield a reduction in certainty, our network
533 dynamics would predict a transient decrease in bump amplitude following a conflicting
534 observation. Further, our network automatically adjusts its cue integration weights (Fig. 5c) to
535 perform close-to-optimal Bayesian inference for absolute HD observations of varying reliability -
536 from highly reliable to very unreliable or even completely absent observations. This stands in
537 contrast to previous approaches[42], that required hand-tuned weights to show that continuous ring
538 attractors can track orientation and compute the running circular average of an absolute HD
539 stimulus. Lastly, our network is to our knowledge the first to fully account for the effect of
540 *probabilistic* angular path integration in a principled way: unlike, e.g., the disc attractor in ref[43], the
541 bump amplitude decay in our network matches the quadratic certainty decay of the ideal Bayesian
542 observer in absence of absolute HD observations. We would expect to observe such a decay in
543 biological ring attractors implementing Bayesian inference once absolute HD observations are
544 removed.

545 Even though our Bayesian HD tracking algorithm requires keeping track of the HD estimate's
546 uncertainty, we have shown that imperfectly tuned ring attractor networks can track head direction
547 reasonably well. In fact, even strict attractor networks with a fixed amplitude, and fixed associated
548 uncertainty, can perform close-to-Bayesian cue integration (Figs. 5b; cf. also ref[44]). This result
549 raises the question of why neurons should encode uncertainty in the first place. First and foremost,
550 for some animals, uncertainty influences their behavior directly to improve their performance (e.g.,
551 refs[45–47]). As a prime example, the homing behavior of the desert ant[48] suggests that the
552 performance gained from tracking one's uncertainty justifies the added complexity for doing so.
553 Further, uncertainty appears to impact the neural encoding of other navigation-related variables.
554 For example, when absolute visual cues are in conflict with path integration cues, grid cells in
555 mouse medial entorhinal cortex are more likely to remap when the visual cues are more reliable[49].
556 Identifying how uncertainty ought to be reflected in their neural activity, as we do here, is required
557 for a comprehensive understanding of the role of uncertainty in the brain's computations.

558 In summary, our work shows how ring attractors could implement dynamic Bayesian inference,
559 even in networks that obey some biological constraints, such as the *Drosophila*'s HD system. We
560 expect similar network motifs to be present in the HD systems of other animals, such as that of
561 mice[9,10], monkeys[50], humans[51], or even in systems that yield three-dimensional HD cells, as those
562 of bats[52]. More generally, we demonstrated how classic network motifs, like those common in ring
563 attractor networks, can perform close-to-optimal Bayesian inference when considered in

564 combination, and expect our results to generalize to other circular variables that are represented
565 in ring attractor networks.

# Acknowledgements

# Author contributions

578 Conceptualization, A.K., M.A.B., J.D.; Methodology, A.K., J.D.; Software, A.K.; Formal analysis,
579 A.K., J.D.; Investigation, A.K, M.A.B., J.D.; Resources, J.D; Writing - Original Draft: A.K., J.D.;
580 Writing - Review & Editing: A.K., M.A.B., J.D.; Visualization, A.K.; Supervision, J.D.; Funding
581 Acquisition, A.K., J.D.

# Declaration of interests

583 The authors declare no competing interests.

## Methods

### Ideal observer model: the circular Kalman filter

Our ideal observer model - the circular Kalman filter (circKF)[21] - performs dynamic Bayesian inference for circular variables. It computes the posterior probability of an unobserved (true) HD $\phi_t \in [-\pi, \pi]$ at each point in time $t$, conditioned on a continuous stream of noisy angular velocity observations $v_{0:t} = \{v_0, v_{dt}, \dots v_t\}$ with $v_\tau \in \mathbb{R}$, and absolute HD observations $z_{0:t} = \{z_0, z_{dt}, \dots z_t\}$ with $z_\tau \in [-\pi, \pi]$. Specifically, we assume that these observations are generated from some true angular velocity $\dot{\phi}_t$ and HD $\phi_t$, whose observations are corrupted by zero-mean noise at each point in time, via

$$v_t | \dot{\phi}_t \sim \mathcal{N}\left(\dot{\phi}_t, \frac{1}{\kappa_v \, dt}\right), \tag{5}$$

$$z_t | \phi_t \sim \mathcal{VM}\left(\phi_t, \sqrt{2\kappa_z \, dt}\right). \tag{6}$$

Here, $\mathcal{N}(\mu, \sigma^2)$ denotes a Gaussian with mean $\mu$ and variance $\sigma^2$, $\mathcal{VM}(\mu, \kappa)$ denotes a von Mises distribution of a circular random variable with mean $\mu$ and precision $\kappa$, and $\kappa_v$ and $\kappa_z$ denote the precision of the angular velocity and absolute HD observations, respectively. Note that as $dt \to 0$, the precision of both angular velocity and absolute HD observations approach 0, in line with the intuition that reducing a time step size $dt$ results in more observations per unit time, which should be accounted for by less precision per observation to avoid "oversampling". More formally, the square-root scaling of the absolute HD observation precision with $\sqrt{2\kappa_z \, dt}$ ensures that the Fisher information of the observations about the true HD scales linearly in time and $\kappa_z$ in the continuum limit $dt \to 0$ (ref[21], Theorem 2). The same applies to the $dt^{-1}$ scaling of the Gaussian variance of the angular velocity observations, again achieving a Fisher information that scales linearly in time.

To support integrating information over time, the model assumes that current HD $\phi_t$ depends on the past HD $\phi_{t-dt}$. Specifically, in absence of further evidence, the model assumes that HD diffuses on a circle,

$$\phi_t | \phi_{t-dt} \sim \mathcal{N}\left(\phi_{t-dt}, \frac{dt}{\kappa_\phi}\right) \mod 2\pi, \tag{7}$$

with a diffusion coefficient that decreases with $\kappa_\phi$. In Results, we assume $\kappa_\phi \to 0$, implying that HD can change arbitrarily across consecutive time steps, which was sufficient to convey intuition into the algorithm's workings. However, when simulating stochastic HD trajectories, we assume they evolve according to Eq. (7) with $\kappa_\phi > 0$, which needs to be accounted for when performing inference. Thus, we here assume a non-zero $\kappa_\phi$ for completeness and reproducibility.

The circKF in Eqs. (1) and (2) assumes that the posterior distribution over HD can be approximated by a von Mises distribution with time-dependent mean $\mu_t$ and certainty $\kappa_t$, i.e. $p(\phi_t | v_{0:t}, z_{0:t}) \approx \mathcal{VM}(\phi_t; \mu_t, \kappa_t)$. Such an approximation is justified if the posterior is sufficiently unimodal, and can, for instance, be compared to a similar approximation employed by extended Kalman filters for non-circular variables.

An alternative parametrization of the von Mises distribution to its mean $\mu_t$ and precision $\kappa_t$, are its natural parameters, $\theta_t = (\kappa_t \cos \mu_t, \kappa_t \sin \mu_t)^T$. Geometrically, the natural parameters can

be interpreted as the Cartesian coordinates of a "probability vector", and $(\mu_t, \kappa_t)$ as its polar co-ordinates (Fig. 2b). As we show in the SI, the natural parameter parametrization makes including absolute HD observations (Eq. (6)) in the circKF straightforward. In fact, it becomes a vector addition. In contrast, including angular velocity observations (Eq. (5)) is mathematically intractable, such that the circKF relies on an approximation method called projection filtering[20] to find closed-form dynamic expressions for posterior mean and certainty (see ref[21] for technical details, and the SI for a more accessible description of the circKF).

Taken together, the circKF for the model specified by Eqs. (5)-(7) reads:

$$d\mu_t = \frac{\kappa_v}{\kappa_\phi + \kappa_v} v_t \, dt + \frac{\sqrt{2\kappa_z \, dt}}{\kappa_t} \sin(z_t - \mu_t), \tag{8}$$

$$d\kappa_t = -\frac{f(\kappa_t)}{2(\kappa_\phi + \kappa_v)} \kappa_t dt + \sqrt{2\kappa_z \, dt} \sin(z_t - \mu_t), \tag{9}$$

where $f(\kappa_t)$ is a monotonically increasing nonlinear function,

$$f(\kappa) = \frac{A(\kappa)}{\kappa_t - A(\kappa) - \kappa A(\kappa)^2}, \quad \text{with } A(\kappa) = \frac{I_1(\kappa)}{I_0(\kappa)}, \tag{10}$$

and $I_0(\cdot)$ and $I_1(\cdot)$ denote the modified Bessel functions of the first kind of order 0 and 1, respectively. Setting $\kappa_\phi \to 0$ yields Eqs. (1) and (2). Importantly, setting $\kappa_\phi \to 0$ does not conceptually change the general vector operations we present in Fig. 2.

For a sufficiently large $\kappa$ (i.e., high certainty), the nonlinearity $f(\kappa)$ approaches the linear function, $f(\kappa) \to 2\kappa - 2$. In our **quadratic approximation**, we thus replace the non-linearity by a quadratic decay:

$$d\kappa_t = -\frac{1}{\kappa_\phi + \kappa_v} \left( \kappa_t^2 - \kappa_t \right) dt + \sqrt{2\kappa_z \, dt} \sin(z_t - \mu_t), \tag{11}$$

which well-approximates the circKF in the high certainty regime.

## Network model

We derived a rate-based network model that implements (approximations of) the circKF, by encoding the von Mises posterior parameters in activity $\mathbf{r}_t \in \mathbb{R}^N$ of a neural population with $N$ neurons. Thereby, we focused on the simplest kind of network model that supports such an approximation, which is of the form:

$$d\mathbf{r}_t = -\frac{1}{\tau} \mathbf{r}_t \, dt - g(\mathbf{r}_t) \mathbf{r}_t \, dt + W \cdot \mathbf{r}_t \, dt + \mathbf{I}_t^{\text{ext}}, \tag{12}$$

where $\tau$ is the network time constant, $g : \mathbb{R}^N \to \mathbb{R}_+$ is a scalar nonlinearity, and the elements of $\mathbf{r}_t$ are assumed to be ordered by the respective neuron's preferred HD, $\phi_1, \ldots, \phi_N$ (see Eq. (3)). We decomposed the recurrent connectivity matrix into $W = \frac{w_0}{2} W^{\text{const}} + w_1^{\text{even}} W^{\text{cos}} + w_1^{\text{odd}} W^{\text{sin}}$, where $W^{\text{const}}$ denotes a matrix with constant entries, and $W^{\text{cos}}$ and $W^{\text{sin}}$ refer to cosine- and sine-shaped connectivity profiles (Fig. 4a). Specifically, due to the network's circular symmetry, the entries of these matrices only depend on the relative distance in preferred HD, and are given by $W_{ij}^{\text{const}} = \frac{2}{N}$,

$W_{ij}^{\text{cos}} = \frac{2}{N} \cos(\phi_i - \phi_j)$, and $W_{ij}^{\text{sin}} = \frac{2}{N} \sin(\phi_i - \phi_j)$. The scaling factor $\frac{2}{N}$ was chosen to facilitate matching our analytical results from the continuum network to the network structure outlined here. We further considered a cosine-shaped external input of the form $I_t^{\text{ext}}(\phi_i) = I_t(dt) \cos(\Phi_t - \phi_i)$ that is peaked around an input location $\Phi_t$. Here, $I_t(dt)$ denotes the maximum input in the infinitesimal time bin $dt$.

As described in Results, we assume the population activity $\mathbf{r}_t$ to encode the HD belief parameters $\mu_t$ and $\kappa_t$ in the phase and amplitude of the activity's first Fourier component. As we show in the SI, the described network dynamics thus lead to the following dynamics of the cosine-profile parameters $\mu_t$ and $\kappa_t$:

$$d\mu_t = w_1^{\text{odd}} \, dt + \frac{I_t}{\kappa_t} \sin(\Phi_t - \mu_t), \tag{13}$$

$$d\kappa_t = \left( w_1^{\text{even}} - \frac{1}{\tau} \right) \kappa_t \, dt - g(\mathbf{r}_t)\kappa_t \, dt + I_t \cos(\Phi_t - \mu_t). \tag{14}$$

To derive these dynamics, we assumed the following:

1. The network is *rate-based*.

2. Our analysis assumes a continuum of neurons, i.e. $N \to \infty$. For numerical simulations, and the network description below, we used a finite-sized network of size $N$ that corresponds to a discretization of the continuous network. SI Fig. S2 demonstrates only a very weak dependence of our results on the exact number of neurons in the network.

3. Our analysis and simulations focused on the first Fourier mode of the bump profile, and is thus independent of the exact shape of the profile (as long as Eq. (3) holds).

**Network parameters for Bayesian inference**

Having identified how the dynamics of the $\mu_t$ and $\kappa_t$ encoded by the network (Eqs. (13) & (14)) depend on the network parameters, we now tuned these parameters to match these dynamics to those of the mean and certainty of the circKF (Eqs. (8) & (9)). Specifically, we find for the network parameters:

- Odd recurrent connectivities are modulated by angular velocity observations, $w_1^{\text{odd}} = \frac{\kappa_v}{\kappa_\phi + \kappa_v} v_t$, which shifts the activity profile without changing its amplitude[7,8].

- Absolute HD observations $z_t$ are represented as the peak position $\Phi_t$ of a cosine-shaped external input whose amplitude is modulated by the reliability of the observation, i.e., $I_t = \sqrt{2\kappa_z \, dt}$. The inputs might contain additional Fourier modes (e.g., a constant baseline), but those do not affect the dynamics in Eqs. (13) and (14).

- The even component of the recurrent excitatory input needs to exactly balance the internal activity decay, i.e., $w_1^{\text{even}} - \frac{1}{\tau} = 0$.

- The decay nonlinearity is modulated by the reliability of the angular velocity observations, and is given by $g(\mathbf{r}_t) = \frac{1}{2(\kappa_\phi + \kappa_v)} f\big(\kappa(\mathbf{r}_t)\big)$, where $f(\cdot)$ equals the nonlinearity that governs

22

675      the certainty decay in the circKF (Eq. (10)). This can be achieved, e.g., through interaction

676      with an inhibitory neuron (or a pool of inhibitory neurons) with activation function $f(\cdot)$ that

677      computes the activity bump's amplitude $f\big(\kappa(\mathbf{r}_t)\big)$.

### Quadratic approximation

679 To gain additional analytical tractability, we further approximated the recurrent inhibition by one that

680 takes the form $g(\mathbf{r}_t)\mathbf{r}_t \to w^{\mathsf{quad}}\left(\pi\sum_{i=1}^{N}[r_t^{(i)}]_+\right)\cdot\mathbf{r}_t$, where $[\cdot]_+$ denotes the rectification nonlinearity.

681 The resulting recurrent inhibition can be shown to be quadratic in the amplitude $\kappa_t$, and has the

682 further benefit of introducing an attractor state at a positive bump aplitude (see below). In the large

683 population limit, $N \to \infty$, this leads to the amplitude dynamics (see SI for derivation)

$$d\kappa_t = \left(w_1^{\mathsf{even}} - \frac{1}{\tau}\right)\kappa_t\,dt - w^{\mathsf{quad}}\kappa_t^2 dt + I_t\cos(\Phi_t - \mu_t). \tag{15}$$

684 The dynamics of the phase $\mu_t$ does not depend on the form of $g(\cdot)$ and thus remains to be given

685 by Eq. (13). If we set the network parameters to $w^{\mathsf{quad}} = \frac{1}{\kappa_\phi + \kappa_v}$ and $w^{\mathsf{even}} - \frac{1}{\tau} = \frac{1}{\kappa_\phi + \kappa_v}$, while

686 sensory input, i.e. angular velocity $v_t$ and absolute HD observations $z_t$, enter in the same way as

687 before, the network implements the quadratic approximation to the circKF (Eqs. (8) & (11)).

### Ring-attractor networks

689 In absence of absolute HD observations ($I_t = 0$), the amplitude dynamics in Eq. (15) has a stable

690 **fixed point** at $\kappa^* = \frac{w^{\mathsf{even}} - 1/\tau}{w^{\mathsf{quad}}}$ and no preferred phase, making it a ring-attractor network. Lin-

691 earizing the $\kappa_t$ dynamics around this fixed points reveals that it is approached with **decay speed**

692 $\beta = w^{\mathsf{even}} - \frac{1}{\tau}$. A large value of $\beta$ denotes faster dynamics and thus indicates more rigid attractor

693 dynamics. In the limit of $\beta \to \infty$ we consider the attractor to be a "strict" attractor that, upon

694 any perturbation, immediately moves back to its attractor state. For the quadratic approximation

695 network, we find $\kappa^* = 1$ and $\beta = \frac{1}{\kappa_\phi + \kappa_v}$. Further, in our simulations in Fig. 5, we explored network

696 dynamics with a range of $\kappa^*$ and $\beta$ values by adjusting network parameters accordingly.

### Multiple population network

698 We extended the single population network dynamics, Eq. (12), to encompass five populations:

699 a HD population, which we designed to track HD estimate and certainty with its bump parameter

700 dynamics, two angular velocity populations (AV+ and AV-), which are tuned to HD and are differ-

701 entially modulated by angular velocity input, an inhibitory population (INH), and a population that

702 mediates external input (EXT), corresponding to absolute HD observations. The resulting network

703 dynamics become (see SI for details):

23

$$\dot{\mathbf{r}}_t^{HD} = -\frac{1}{\tau_{HD}}\mathbf{r}_t^{HD} + W_{HD\leftarrow HD}\cdot\mathbf{r}_t^{HD} + W_{HD\leftarrow AV+}\cdot\mathbf{r}_t^{AV+} + W_{HD\leftarrow AV-}\cdot\mathbf{r}_t^{AV-} \tag{16}$$

$$+ (W_{HD\leftarrow INH}\mathbf{r}_t^{INH}) \circ \mathbf{r}_t^{HD} + \mathbf{I}_t^{ext}, \tag{17}$$

$$\tau_{AV+}\dot{\mathbf{r}}_t^{AV+} = -\mathbf{r}_t^{AV+} + (o^{AV} + v_t)W_{AV+\leftarrow HD}\cdot\mathbf{r}_t^{HD}, \tag{18}$$

$$\tau_{AV-}\dot{\mathbf{r}}_t^{AV-} = -\mathbf{r}_t^{AV-} + (o^{AV} - v_t)W_{AV-\leftarrow HD}\cdot\mathbf{r}_t^{HD}, \tag{19}$$

$$\tau_{INH}\dot{\mathbf{r}}_t^{INH} = -\mathbf{r}_t^{INH} + W_{INH\leftarrow HD}\cdot[\mathbf{r}_t^{HD}]_+ + W_{INH\leftarrow INH}\cdot\mathbf{r}_t^{INH}. \tag{20}$$

Here, the $W_{to\leftarrow from}$ denote connectivities within and between populations, and $o^{AV}$ is a constant activity baseline in the AV populations.

The network parameters were tuned such that the activity profile in the HD population tracks the dynamics of the circKF quadratic approximation, in the same way as for the single-population network, Eq. (12). To limit the degrees of freedom, we further constrained the connectivity structure between HD and AV+/- and INH populations by the known connectome of the *Drosophila* HD system (hemibrain dataset[33]). Specifically, we focused on the connectivities between EPG, PEN1 and $\Delta 7$ neurons (which in our model corresponds to HD, AV+/- and INH neurons), sorted according to anatomical regions within the ellipsoid body and the protocerebral bridge (Fig. 6b). Thereby, we used total number of synaptic connections between two regions as a proxy for connection strength. We further assumed that interactions within AV+/- populations and between AV+/- and INH populations were negligible. The resulting connectivity profile in Fig. 6c was determined by matching the *Drosophila* connectome as closely as possible, while allowing for modulation of the across-population connection strengths $c_0^{HD}, c_1^{HD}$, $c^{AV\pm\leftarrow HD}$, $c^{HD\leftarrow AV\pm}$., $c_0^{INH\leftarrow HD}$, $c_1^{INH\leftarrow HD}$, $c_0^{INH}$, $c_1^{INH}$, and $c^{HD\leftarrow INH}$. We specify the specific analytic functions we used to create the connectivity matrix in Fig. 6c in the SI, where we also compute the connection strengths analytically.

## Simulation details

### Numerical integration

Our simulations in Figs. 4 and 5 used artificial data that matched the assumptions underlying our models. In particular, the 'true' HD $\phi_t$ followed a diffusion on the circle, Eq. (7), and observations were drawn at each point in time from Eqs. (5) and (6). To simulate trajectories and observations, we used the Euler-Maruyama scheme[54], which supports the numerical integration of stochastic differential equations. Specifically, for a chosen discretization time step $\Delta t$, this scheme is equivalent to drawing trajectories and observations from Eqs. (7), (5) and (6) directly while substituting $dt \to \Delta t$. The same time-discretization scheme was used to numerically integrate the SDEs of the circKF, Eqs (8) and (9), its quadratic approximation, Eq. (11), and the network dynamics, Eqs. (12) and (16)-(20).

### Performance measures

To measure performance, in Figs. 4f, 5b and 5d we computed the circular average distance[53] of the estimate $\mu_T$ from the true HD $\phi_T$ at the end of a simulation of length $T = 20$ from $P = 5'000$

734   simulated trajectories by $m_1 = \frac{1}{P} \sum_{k=1}^{P} \exp\left(i\left(\mu_T^{(k)} - \phi_T^{(k)}\right)\right)$. The absolute value of the imaginary-
735   valued circular average, $0 \leq |m_1| \leq 1$ denotes an empirical precision (or 'inference precision'), and
736   thus measures how well the estimate $\mu_T$ matches the true HD $\phi_T$. Here, a value of 1 denotes an
737   exact match. The inference precision is related to the circular variance via $\text{Var}_{circ} = 1 - |m_1|$. In
738   SI Fig. S5, we provide histograms with samples $\mu_T - \phi_T$ with different numerical values of $|m_1|$, to
739   provide some intuition for the spread of estimates for a given value of the performance measure.

740     We estimated performance through such averages for all absolute HD observation reliabilities
741   $\kappa_z$ in Figs. 4f and 5b. For the inset of Fig. 5b, and for Fig. 5d, we additionally performed a grid
742   search over the fixed-point amplitude $\kappa^*$ (inset of Fig. 5b), or both the fixed-point amplitude $\kappa^*$ and
743   of the inverse time constant $\beta$ (Fig. 5d). For each setting of $\kappa^*$ and $\beta$ we assessed the performance
744   by computing an average over this performance for a range of observation reliability $\kappa_z$, weighted
745   by how likely each observation reliability is a-priori assumed to be. The latter was specified by a
746   log-normal prior, $p(\kappa_z) = \text{Lognormal}(\mu_{\kappa_z}, \sigma_{\kappa_z}^2)$, favouring intermediate reliabilitiy levels. We chose
747   $\mu_{\kappa_z} = 0.5$ and $\sigma_{\kappa_z}^2 = 1$ for the prior parameters, but our results did not strongly depend on this
748   parameter choice. The performance loss shown in Fig. 5d also relied on such a weighted average
749   across $\kappa_z$'s for a particle filter benchmark (PF, see SI for details). The loss itself was then defined
750   as $1 - \frac{\text{Performance}}{\text{Performance PF}}$.

### Update weights

752   In Fig. 5c, we computed the weight with which a single observation with $|z_t - \mu_t| = 90°$ changes
753   the HD estimate. We defined this weight as the change in HD estimate, normalized by the value
754   of the maximum possible change, $w = \frac{\Delta\mu_t}{\pi} = \frac{1}{\pi}\tan^{-1}\frac{\alpha(\kappa_z\,dt)}{\kappa_t}$. Here, $\alpha(\kappa_z\,dt)$ denotes a function
755   that ensures a linear scaling of the Fisher information with sampling time step (see ref[21], Theorem
756   2, for details about this function). Thus, by design of the observation model, the Fisher information
757   of a single observation with reliability $\kappa_z$ during a time interval $\Delta t$ is given by $I_{z_t}(\phi_t) = \kappa_z\,\Delta t$.
758   We plot the weight as a function of the Fisher information of a single update (how reliable is the
759   observation?) and the Fisher information of the current HD estimate (how certain is the current
760   estimate?), which is given by

$$I_{\mu_t,\kappa_t}(\phi_t) = \mathbb{E}\left[\left(\frac{\partial}{\partial\phi}\log \mathcal{VM}(\phi, \mu_t, \kappa_t)\right)^2\right] = \kappa_t \frac{I_1(\kappa_t)}{I_0(\kappa_t)}. \tag{21}$$

### Details on numerical simulations

762   In our network simulations, we set the network decay constant $\tau$ to an arbitrary, but non-zero,
763   value. Effectively, this resulted in a cosine-shaped activity profile. Note that by setting higher-order
764   recurrent connectivities accordingly, other profile shapes could be realized, without affecting the
765   validity of our analysis above. From the neural activity vector $\mathbf{r}_t$, we retrieved the natural parame-
766   ters $\boldsymbol{\theta}_t$ with a decoder matrix $A = (\cos(\boldsymbol{\phi}^{(i)}), \sin(\boldsymbol{\phi}^{(i)}))^T$, such that $\boldsymbol{\theta}_t = A \cdot \mathbf{r}_t$, and subsequently
767   computed the position of the bump by $\phi_t = \arctan 2(\theta_2, \theta_1)$, and the encoded certainty (length of
768   the population vector) by $\kappa_t = \sqrt{\theta_1^2 + \theta_2^2}$.

769     In all our simulations, times are measured in units of inverse diffusion time constant $\kappa_\phi$, where
770   we set $\kappa_\phi = 1s$ for convenience. Figures were generated based on simulations with the following

parameters:

- Figure 4e: $\kappa_v = 2$, $\kappa_z = 10$ (during 'Visual cue' bout), $\kappa_z = 0$ (during 'Darkness' bout), $\Delta t = 0.01$.

- Figure 4f, 5b, 5d: $\kappa_v = 1$, $T = 20$, $\Delta t = 0.01$. Results are averages over $P = 5000$ simulation runs.

- Figure 5e: $\kappa_v = 1$, $\kappa_z = 1$, $T = 10$, $\Delta t = 0.01$.

Trajectory simulations and general analyses were performed on a MacBook Pro (Mid 2019) running 2.3 GHz 8-core Intel Core i9. Parameter scans were run on the Harvard Medical School $O_2$ HPC cluster. For all our simulations, we used Python 3.9.1 with NumPy 1.19.2. Jupyter notebooks, Python scripts, and data to reproduce the figures will be made available upon acceptance of the manuscript.

# References

1. Wang, X.-J. Synaptic reverberation underlying mnemonic persistent activity. *Trends Neurosci.* **24**, 455–463 (2001).

2. Compte, A. Computational and in vitro studies of persistent activity: Edging towards cellular and synaptic mechanisms of working memory. *Neuroscience* **139**, 135–151 (2006).

3. Hansel, D. & Sompolinsky, H. Modeling Feature Selectivity in Local Cortical Circuits. *Methods Neuronal Model. Ions Netw.* 69 (1998).

4. Rademaker, R. L., Tredway, C. H. & Tong, F. Introspective judgments predict the precision and likelihood of successful maintenance of visual working memory. *J. Vis.* **12**, 21 (2012).

5. Li, H.-H., Sprague, T. C., Yoo, A. H., Ma, W. J. & Curtis, C. E. Joint representation of working memory and uncertainty in human cortex. *Neuron* **109**, 3699-3712.e6 (2021).

6. Knierim, J. J. & Zhang, K. Attractor Dynamics of Spatially Correlated Neural Activity in the Limbic System. *Annu. Rev. Neurosci.* **35**, 267–285 (2012).

7. Zhang, K. Representation of Spatial Orientation by the Intrinsic Dynamics of the Head-Direction Cell Ensemble: A Theory. *J. Neurosci.* **16**, 2112–2126 (1996).

8. Skaggs, W. E., Knierim, J. J., Kudrimoti, H. S. & McNaughton, B. L. A Model of the Neural Basis of the Rat's Sense of Direction. *Adv. Neural Inf. Process. Syst.* 10 (1995).

9. Peyrache, A., Lacroix, M. M., Petersen, P. C. & Buzsáki, G. Internally organized mechanisms of the head direction sense. *Nat. Neurosci.* **18**, 569–575 (2015).

10. Ajabi, Z., Keinath, A. T., Wei, X.-X. & Brandon, M. P. Population dynamics of the thalamic head direction system during drift and reorientation. 2021.08.30.458266 (2021) doi:10.1101/2021.08.30.458266.

11. Redish, A. D., Elga, A. N. & Touretzky, D. S. A coupled attractor model of the rodent head direction system. *Netw. Comput. Neural Syst.* **7**, 671–685 (1996).

12. Seelig, J. D. & Jayaraman, V. Neural dynamics for landmark orientation and angular path

integration. *Nature* **521**, 186–191 (2015).

13. Kim, S. S., Rouault, H., Druckmann, S. & Jayaraman, V. Ring attractor dynamics in the Drosophila central brain. *Science* **356**, 849–853 (2017).

14. Turner-Evans, D. *et al.* Angular velocity integration in a fly heading circuit. *eLife* **6**, e23496 (2017).

15. Knill, D. C. & Pouget, A. The Bayesian brain: the role of uncertainty in neural coding and computation. *Trends Neurosci.* **27**, 712–719 (2004).

16. Dehaene, G. P., Coen-Cagli, R. & Pouget, A. Investigating the representation of uncertainty in neuronal circuits. *PLOS Comput. Biol.* **17**, e1008138 (2021).

17. Kalman, R. E. A New Approach to Linear Filtering and Prediction Problems. *Trans. ASME J. Basic Eng.* **82**, 35–45 (1960).

18. Kalman, R. E. & Bucy, R. S. New Results in Linear Filtering and Prediction Theory. *J. Basic Eng.* **83**, 95–108 (1961).

19. Kurz, G., Gilitschenski, I. & Hanebeck, U. D. Recursive Bayesian filtering in circular state spaces. *IEEE Aerosp. Electron. Syst. Mag.* **31**, 70–87 (2016).

20. Brigo, D., Hanzon, B. & Le Gland, F. Approximate nonlinear filtering by projection on exponential manifolds of densities. *Bernoulli* **5**, 495–534 (1999).

21. Kutschireiter, A., Rast, L. & Drugowitsch, J. Projection Filtering with Observed State Increments with Applications in Continuous-Time Circular Filtering. *IEEE Trans. Signal Process.* 1–1 (2022) doi:10.1109/TSP.2022.3143471.

22. Murray, R. F. & Morgenstern, Y. Cue combination on the circle and the sphere. *J. Vis.* **10**, 15–15 (2010).

23. Wilson, R. & Finkel, L. A Neural Implementation of the Kalman Filter. *Adv. Neural Inf. Process. Syst.* **22**, 9 (2009).

24. Ben-Yishai, R., Bar-Or, R. L. & Sompolinsky, H. Theory of orientation tuning in visual cortex. *Proc. Natl. Acad. Sci.* **92**, 3844–3848 (1995).

25. Georgopoulos, A., Kettner, R. & Schwartz, A. Primate motor cortex and free arm movements to visual targets in three- dimensional space. II. Coding of the direction of movement by a neuronal population. *J. Neurosci.* **8**, 2928–2937 (1988).

26. Johnson, A., Seeland, K. & Redish, A. D. Reconstruction of the postsubiculum head direction signal from neural ensembles. *Hippocampus* **15**, 86–96 (2005).

27. Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic population codes. *Nat. Neurosci.* **9**, 1432–8 (2006).

28. Beck, J. M., Latham, P. E. & Pouget, A. Marginalization in Neural Circuits with Divisive Normalization. *J. Neurosci.* **31**, 15310–15319 (2011).

29. Lyu, C., Abbott, L. F. & Maimon, G. *A neuronal circuit for vector computation builds an allocentric traveling-direction signal in the* Drosophila *fan-shaped body*. http://biorxiv.org/lookup/doi/10.1101/2020.12.22.423967 (2020) doi:10.1101/2020.12.22.423967.

30. Xie, X., Hahnloser, R. H. R. & Seung, H. S. Double-ring network model of the head-direction system. *Phys. Rev. E - Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.* **66**, 9–9 (2002).

31. Compte, A. Synaptic Mechanisms and Network Dynamics Underlying Spatial Working Memory in a Cortical Network Model. *Cereb. Cortex* **10**, 910–923 (2000).

32. Turner-Evans, D. B. *et al.* The Neuroanatomical Ultrastructure and Function of a Biological Ring Attractor. *Neuron* S0896627320306139 (2020) doi:10.1016/j.neuron.2020.08.006.

33. Scheffer, L. K. *et al.* A connectome and analysis of the adult Drosophila central brain. *eLife* **9**, e57443 (2020).

34. Hulse, B. K. *et al.* A connectome of the Drosophila central complex reveals network motifs suitable for flexible navigation and context-dependent action selection. *eLife* **10**, e66039 (2021).

35. Green, J. *et al.* A neural circuit architecture for angular integration in Drosophila. *Nature* **546**, 101–106 (2017).

36. Omoto, J. J. *et al.* Visual Input to the Drosophila Central Complex by Developmentally and Functionally Distinct Neuronal Populations. *Curr. Biol.* **27**, 1098–1110 (2017).

37. Fisher, Y. E., Lu, J., D'Alessandro, I. & Wilson, R. I. Sensorimotor experience remaps visual input to a heading-direction network. *Nature* (2019) doi:10.1038/s41586-019-1772-4.

38. Kim, S. S., Hermundstad, A. M., Romani, S., Abbott, L. F. & Jayaraman, V. Generation of stable heading representations in diverse visual scenes. *Nature* 1–6 (2019) doi:10.1038/s41586-019-1767-1.

39. van Bergen, R. S., Ji Ma, W., Pratte, M. S. & Jehee, J. F. M. Sensory uncertainty decoded from visual cortex predicts behavior. *Nat. Neurosci.* **18**, 1728–1730 (2015).

40. Sun, X., Mangan, M. & Yue, S. An analysis of a ring attractor model for cue integration. *Lect. Notes Comput. Sci. Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinforma.* **10928 LNAI**, 459–470 (2018).

41. Sun, X., Yue, S. & Mangan, M. A decentralised neural model explaining optimal integration of navigational strategies in insects. *eLife* **9**, e54026 (2020).

42. Esnaola-Acebes, J. M., Roxin, A. & Wimmer, K. *Bump attractor dynamics underlying stimulus integration in perceptual estimation tasks.* http://biorxiv.org/lookup/doi/10.1101/2021.03.15.434192 (2021) doi:10.1101/2021.03.15.434192.

43. Carroll, S., Josić, K. & Kilpatrick, Z. P. Encoding certainty in bump attractors. *J. Comput. Neurosci.* **37**, 29–48 (2014).

44. Takiyama, K. Bayesian estimation inherent in a Mexican-hat-type neural network. *Phys. Rev. E* **93**, 052303 (2016).

45. Ernst, M. O. & Banks, M. S. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature* **415**, 429–433 (2002).

46. Piet, A. T., Hady, A. E., Brody, C. D., El Hady, A. & Brody, C. D. Rats adopt the optimal timescale for evidence integration in a dynamic environment. *Nat. Commun.* **9**, 1–12 (2018).

47. Fetsch, C. R., Turner, A. H., DeAngelis, G. C. & Angelaki, D. E. Dynamic reweighting of visual and vestibular cues during self-motion perception. *J. Neurosci.* **29**, 15601–15612 (2009).

48. Merkle, T. & Wehner, R. Desert ants use foraging distance to adapt the nest search to the uncertainty of the path integrator. *Behav. Ecol.* **21**, 349–355 (2010).

49. Campbell, M. G., Attinger, A., Ocko, S. A., Ganguli, S. & Giocomo, L. M. Distance-tuned neurons drive specialized path integration calculations in medial entorhinal cortex. *Cell Rep.* **36**, 109669 (2021).

50. Robertson, R. G., Rolls, E. T., Georges-François, P. & Panzeri, S. Head direction cells in the primate pre-subiculum. *Hippocampus* **9**, 206–219 (1999).

51. Baumann, O. & Mattingley, J. B. Medial Parietal Cortex Encodes Perceived Heading Direction in Humans. *J. Neurosci.* **30**, 12897–12901 (2010).

52. Finkelstein, A. *et al.* Three-dimensional head-direction coding in the bat brain. *Nature* **517**, 159–164 (2015).

53. Mardia, K. V. & Jupp, P. E. *Directional Statistics.* 3 (John Wiley & Sons, 2000). doi:10.1002/9780470316979.

54. Kloeden, P. E. & Platen, E. *Numerical solution of stochastic differential equations.* (Springer, 2010).