
Latent state dynamics in female *Drosophila* during social interactions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Social interactions across animal species are governed by the interplay between
2 multimodal sensory information and an animal's internal states. Here, we inves-
3 tigate this interplay in female *Drosophila* as she engages with the male during
4 courtship, a highly dynamic social behavior. While male behaviors during courtship,
5 such as song production, have been well characterized moment by moment, the
6 female's actions have not been described with similar temporal precision, despite
7 her central role in determining copulation outcomes. Her behavior displays high
8 variability and is often viewed as volitional, raising questions about its structure
9 and predictability at fine temporal resolution. To address this, we used a state-space
10 model that combines generalized linear models (GLMs) with a Hidden Markov
11 Model (HMM) to uncover latent states that modulate the relationship between male
12 sensory cues and female responses. We find that overall, male cues weakly predict
13 female behavior, but that predictive power varies substantially across inferred states:
14 some states exhibit clear cue-driven structure, while others show reduced sensitivity
15 to external cues and more internally-driven behavior. At short timescales (33ms),
16 female behavior appears only weakly predictable and highly variable, yet, at longer
17 timescales (1min), a rich latent state structure emerges, hinting at internal gating
18 and evaluation of social signals over time. This work provides the first moment-
19 by-moment characterization of female behavior during courtship, taking a crucial
20 step toward closing the loop in social modeling. By capturing how internal states
21 shape responses to multimodal sensory cues, it offers a foundation for identifying
22 underlying circuit pathways for multisensory integration in the brain, via the female
23 whole-brain connectome.

1 Introduction

25 Animal communication unfolds through dynamic exchanges shaped by both external signals and
26 internal states. In *Drosophila melanogaster*, males produce structured songs and pursuit behaviors
27 that have been quantified at millisecond resolution, providing a model system for studying social
28 signals. Male behavior during courtship follows a relatively stereotyped structure, with consistent
29 song motifs and spatial positioning near the female [1, 2, 3]. Female trajectories during courtship
30 vary widely—when grouped by male position and song onset, female trajectories show no consistent
31 directional bias (Figure 1b,c). So, female responses, although decisive in determining courtship
32 outcomes, are often treated as noisy and volitional. This imbalance leaves open a fundamental
33 question: how do internal states structure the female's engagement with male cues during social
34 interaction?

35 The female's velocity time series presents a unique statistical challenge, as she remains largely
36 stationary for over 70% of the time even when he is actively courting her, presumably "assessing" the
37 male [4]. This results in a distribution heavily skewed toward zero values, making it difficult to model

her movement dynamics and to infer the timescales of her decision-making. Because such prolonged stationary-like bouts are common across insects, modeling them explicitly could also provide insight into any variability in responsiveness to social cues. The female fly offers a good candidate for this approach: in a social context where she is expected to be active and responsive to the male, she also spends extended periods in stillness.

Here, we take a first step by first comprehensively modeling the sensory signals received by the female to predict a wide range of her responses. Unlike prior discrete-only approaches, we analyze continuous locomotor dynamics alongside discrete actions, including wing flicking—a rejection behavior that plays a critical role in communication. Using a GLM-HMM, we find that male cues weakly predict female actions overall, but predictive power varies sharply across latent states: some states exhibit clear cue-driven dynamics, while others appear internally dominated. At short timescales, responses are highly variable; over longer windows, a structured latent dynamics emerges.

Recent work has shown that unsupervised latent-state models can uncover “behavioral syllables” that combine into higher-order structures resembling a grammar [5, 6, 7]. These approaches have revealed how individual animals organize locomotor or foraging sequences but have rarely been extended to social contexts. This work provides a temporal account of low-engagement animal during a social dynamic behavior and how latent states shape its behaviors in response to multimodal cues. More broadly, it highlights how AI or ML approaches can reveal latent communicative structures in animal behavior, offering new entry points for comparative studies of communication across species.

2 Methodology

Social behavior quantification. To investigate how female flies respond to male cues during natural courtship, we used high-resolution pose tracking (SLEAP) to extract the trajectories and body keypoints of pairs of male and female flies as they interacted (Figure 1a). Simultaneously, we recorded the male’s courtship song using a 9-mic array, allowing precise segmentation of his song into pulse and sine components (Figure 1a). We extracted detailed behavioral readouts by tracking 13 keypoints on the female’s body (Figure 1d), allowing us to quantify her forward, lateral, and angular velocity, as well as discrete rejection behaviors such as wing flicking (Figure 5b). These outputs capture both locomotor responses and elements of behavioral rejection. To relate female behaviors to male actions, we characterized the set of multimodal sensory inputs or feedback cues the female receives from the male, spanning visual, auditory, and tactile channels (Figure 1f).

A model with hidden states (GLM-HMM) to predict female behavior. We fit a GLM-HMM in which the female transitions between discrete latent states, each with its own linear mapping from male sensory inputs to behavioral outputs. Specifically, male sensory features from the preceding 3 seconds were used to predict female forward, lateral, and angular velocities, as well as wing flicking, with regression weights that depend on the current latent state (Figure 3a; See Appendix). Each latent state z_t defines a separate generalized linear model (GLM) relating the male cue history s_t to the observed female behavior y_t . The female’s movement variables, such as her forward, lateral, and angular velocity, are continuous and modeled as a Gaussian distribution whose parameters depend on her current latent state z_t and cue history s_t :

$$p(y_t | z_t, s_t) = \mathcal{N}(y_t | w_{z_t} s_t + b_{z_t}, \sigma_{z_t}^2)$$

For her wing flicking ($y_t = 1$) which is a discrete behavior, the model uses a logistic function to describe the probability of flicking:

$$p(y_t | z_t, s_t) = \sigma(w_{z_t} s_t + b_{z_t})$$

where $w_k \in \mathbb{R}^M$ denotes the GLM weights for latent state $k \in \{1 \dots, K\}$. The full set of model parameters, $\theta \equiv \{\pi, A, w_k, b_k, \sigma_k^2\}$, is learned using the expectation–maximization (EM) algorithm (Appendix).

3 Results

Five-state GLM-HMM. We used the GLM-HMM to predict female behavior and evaluated its predictive performance on held-out data. To assess model quality, we computed the difference in log-likelihood between the GLM-HMM and the Chance model (more details in Appendix). We fit

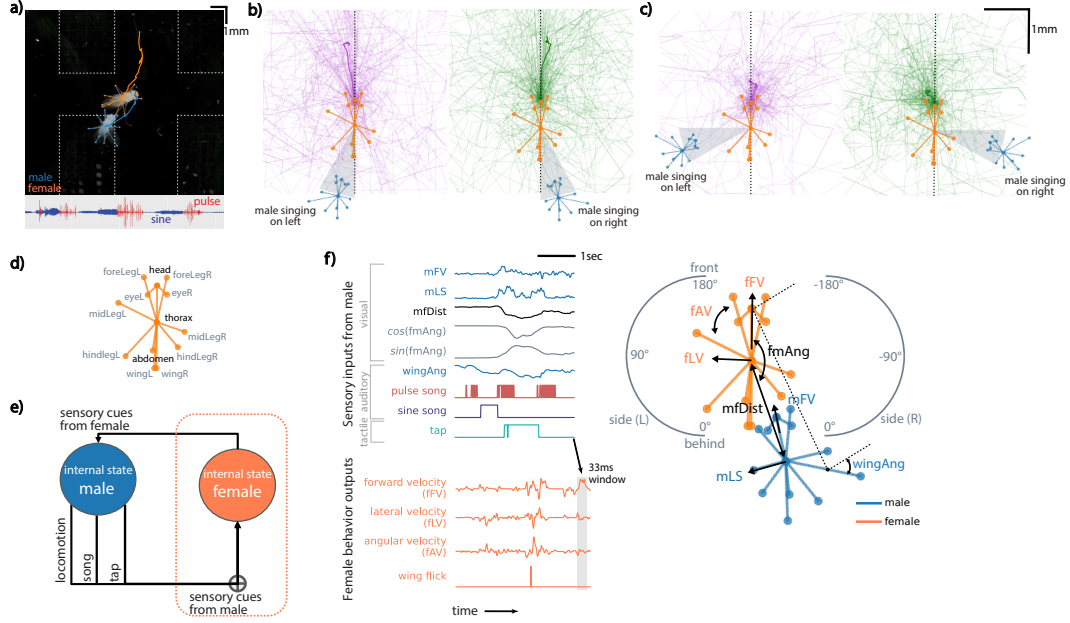


Figure 1: Female behavior during courtship is variable and multimodal. (a) Top: Example frame from a courtship video showing SLEAP-tracked male and female flies. Dashed white lines indicate 9 microphones used to record courtship song. Bottom: recorded courtship song from the male segmented into sine and pulse components. (b, c) Female trajectories following male song onset. Example 1-second trajectories of the female (orange) following male pulse song onset, separated by whether the male was singing from her left or right. Trajectory samples are aligned to the female's position and orientation at song onset ($n=200$). The male is positioned within the gray 30 sector. Trajectories reveal high variability in female responses. Bold indicates the average female trajectory in each condition. (d) Schematic showing the 13 keypoints tracked on the fly body using SLEAP. (e) Schematic of the bidirectional social loop in *Drosophila* courtship. The green box highlights the portion of the loop studied in this work. (f) Left: Example 3-second time series of male sensory cues (top) used to predict female behavior in the next 33ms time bin (below). Sensory cues: visual (male forward velocity - mFV, lateral speed - mLS, angle - fmAng, mfDist - distance), auditory (sine, pulse), and tactile (tap). Outputs: female forward velocity - fFV, lateral velocity - fLV, angular velocity (fAV), and wing flicks. Right: relationship between sensory inputs and female output variables. See Supp Figure 5 for more details.

89 GLM-HMMs with varying numbers of latent states. Note that the one-state GLM-HMM is simply a
90 standard GLM with no internal state. We found that a five-state GLM-HMM achieved a substantial
91 improvement in predictive performance on the heldout fly pairs: 220 bits/s over the Chance model,
92 compared to 115 bits/s for the GLM. Increasing the number of states beyond five did not yield
93 significant additional gains (Figure 3b). Models with 5–7 states consistently outperformed the GLM,
94 suggesting that latent internal dynamics offer predictive value beyond what can be captured by
95 sensory cues alone.

96
97 **Evaluating performance by behavior, by state and overall.** We quantified how well the
98 model captured moment-by-moment variation in female behavior by computing the Pearson
99 correlation between predicted and actual velocities. As shown in Figure 3c, the model achieved
100 consistent predictive performance across animals in both training and held-out datasets. A breakdown
101 across all three velocity components—forward, lateral, and angular—revealed highest correlations
102 for forward motion, with modest correlations for lateral and angular velocities (Figure 3d). Finally,
103 we assessed how well the model captured discrete behaviors such as wing flicking using F1 score
104 (Figure 3d).

105 We next examined the properties of the latent states inferred by the GLM-HMM. Individual flies
106 occupied multiple states during courtship, with most animals spending substantial time in States
107 3–5 (Figure 3f), suggesting that these states capture common, shared behavioral modes. To better

understand how predictive performance varied by state, we computed the correlation between predicted and actual velocities within each state. Some states (e.g., State 2) were more predictable than others (Figure 3g). While correlations decreased slightly on held-out data (Figure 3g, right), the relative pattern was preserved, suggesting that these differences reflect meaningful structure rather than overfitting.

Retrieved latent states are structured across timescales. We next investigated what behavioral features the GLM-HMM states correspond to and when they occur during courtship. State 1 appeared infrequently, typically at the start of sessions, and was marked by high velocities in both animals, large inter-fly distances, and minimal song and tap cues (Figure 2a; Figure 8b). We interpret this as a chamber introduction state at the start of a session (Figure 2e-f), not reflective of courtship, and exclude it from further analysis due to low occupancy and noisy GLM parameters. State 2 reflects an active locomotor state of the female, marked by moderate to fast movement while the male follows at some distance, circling and producing both song and tapping cues. This state occurs throughout interactions but declines over time. State 3 represents moderate engagement, with the female moving slowly in close proximity to the male, who circles behind her with consistent song and tapping. Video inspection suggests the female responds subtly through side-stepping, turning, or small shifts in position, and this state is sustained across the session (Figure 2c). State 4 corresponds to a low-activity period in which the female is largely stationary and often engaged in grooming, while the male remains nearby with minimal movement and sparse sine song or tapping. Finally, State 5 reflects full stillness of the female, with neither locomotion nor grooming. Here, the courtship song and tap cues are minimal, and the male remains positioned closely with little motion. This state often persists for several seconds and becomes increasingly common as courtship progresses (Figure 2e), though it drops sharply just before copulation, when the female typically transitions back into State 3. This prolonged stillness may reflect a critical phase of behavioral assessment or decision-making, but that interpretation remains speculative.

Latent states are defined by distinct mappings between feedback cues and female locomotion behavior. Although States 3, 4, and 5 all have low female velocity (Figure 2a, Figure 8b), they differ markedly in how her behavior couples to male cues. A summary of GLM filters broken down by behavior in Figure 4a confirm this distinction, revealing mostly low cue sensitivity in States 3–4 but tuning to male cues in State 2 and 5. In State 3, the female shows weak but detectable responsiveness, occasionally adjusting her position in response to male circling, song, or tapping—reflected in low but non-zero filter weights (Figure 4). In State 4, interpreted as a grooming state, her motion appears largely self-generated and decoupled from male behavior, with filters near zero regardless of male activity. In contrast, State 5 shows strong sensorimotor coupling: although the female remains stationary, she is poised to respond to male cues and capable of producing large responses, as shown by its high, input-dependent filter weights (Figure 4). These state-dependent sensory filters also depict modality-specific drive of female locomotion. A work in progress is to experiment with female flies with perturbed visual or auditory senses and model their responses and see how they "compensate".

4 Conclusion

Overall, prediction performance remained low across all states (Figure 3h), suggesting that female fly behavior may be harder to predict from male cues alone. These patterns held across both training and test sets, suggesting that while certain aspects of behavior are robustly predicted by the model, others, such as her turning (lateral vel) or side-stepping (angular vel), may depend more heavily on internal state, be influenced by sensory cues not measured here, or simply be noisier and thus harder to predict at this timescale. In particular, states associated with grooming or prolonged stillness also pose additional challenge in modeling animal behavior. Such states make it harder to observe her behavior and contribute to overall low prediction. These states may reflect slower internal processes with behavioral inertia. Together, these findings suggest that fine-timescale models may be insufficient for fully capturing female behavior in these states, and point toward models that incorporate internal dynamics on longer timescales are better suited for her behavior—consistent with prior work where a nonlinear integration of male cues and slow adaptation over scale of minutes were predictive of her walking speed [8].

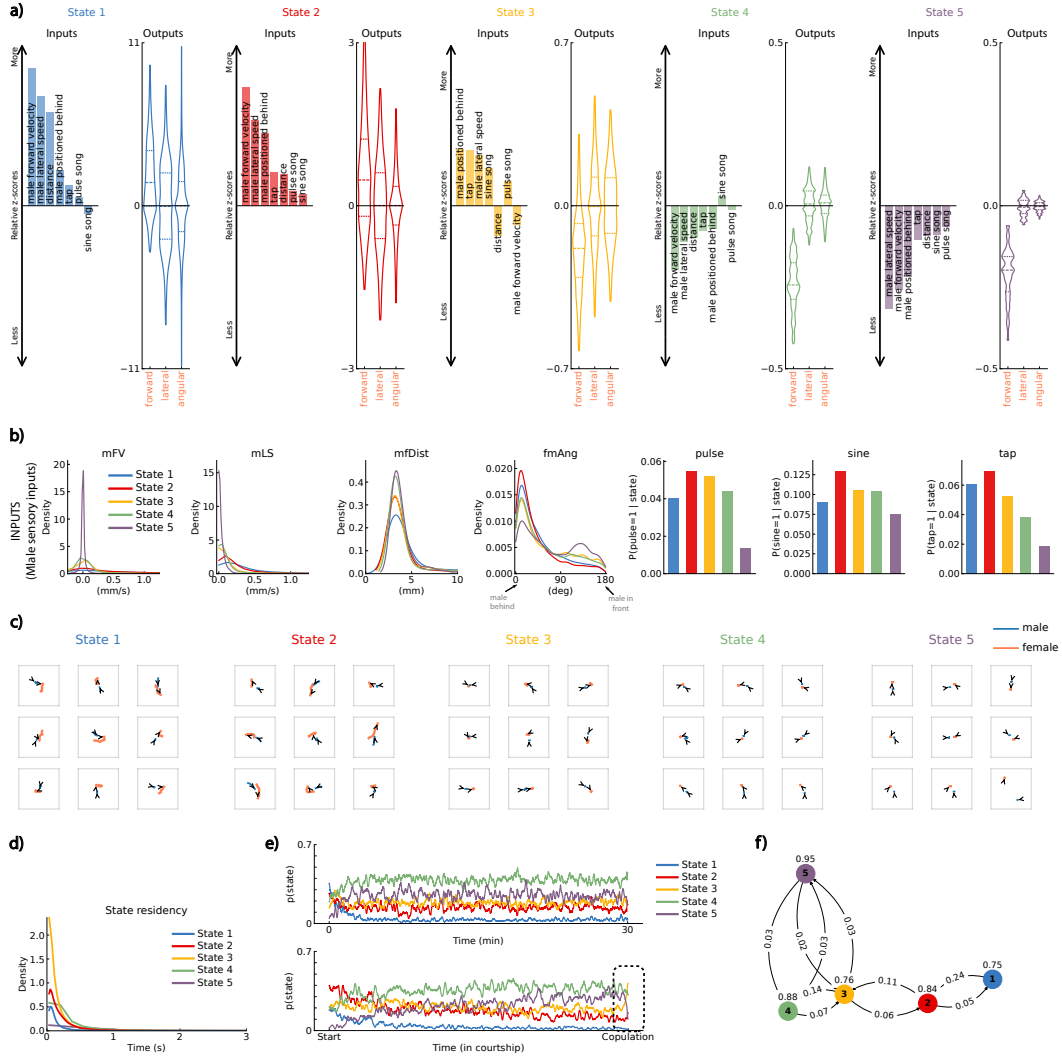


Figure 2: Retrieved latent states uncover structure in female behavior across timescales. (a) Left: For each state, sensory input features are ordered by their relative difference from the across-state mean. Right: Distributions of female forward, lateral, and angular velocities - within each state. Note that the axis scales differ across states. (b) Distributions of values for some of the sensory cues and for each state. Although a state may have features that are larger or smaller than average, the distributions are highly overlapping. (c) Representative traces of male and female movement trajectories in each state. Arrows indicate fly orientation at the end of 300 ms. (d) The dwell times of the five states across all of the data (including both training and test sets). Data from all 75 animals. (e) The mean probability across flies of being in each state fluctuated only slightly over time when aligned to absolute time (top) or the time of copulation (bottom). Immediately before copulation, there was a slight increased probability of being in State 3 (bottom). Data are from all 75 fly pairs. (f) Fitted state transition diagram representing the inferred dynamics of the five-state GLM-HMM. Arrow labels indicate the probability of transitioning from one state to another. Together, panels e and f show that while states $1 \rightarrow 2 \rightarrow 3$ tend to occur in sequence early in a session, states 3, 4, and 5 frequently transition among one another, forming a regime of behaviorally similar dynamics that emerge as courtship progresses.

References

- [1] Herman T Spieth. Courtship behavior in drosophila. *Annual review of entomology*, 19(1): 385–405, 1974.
- [2] Christelle Lasbleiz, Jean-François Ferveur, and Claude Everaerts. Courtship behaviour of drosophila melanogaster revisited. *Animal Behaviour*, 72(5):1001–1012, 2006.
- [3] Philip Coen, Jan Clemens, Andrew J. Weinstein, Diego A. Pacheco, Yi Deng, and Mala Murthy. Dynamic sensory cues shape song structure in Drosophila. *Nature*, 507(7491): 233–237, March 2014. ISSN 0028-0836, 1476-4687. doi: 10.1038/nature13131. URL <http://www.nature.com/articles/nature13131>.
- [4] Kaiyu Wang, Fei Wang, Nora Forknall, Tansy Yang, Christopher Patrick, Ruchi Parekh, and Barry J. Dickson. Neural circuit mechanisms of sexual receptivity in Drosophila females. *Nature*, 589(7843):577–581, January 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-020-2972-7. URL <https://www.nature.com/articles/s41586-020-2972-7>.
- [5] Gordon J Berman, Daniel M Choi, William Bialek, and Joshua W Shaevitz. Mapping the stereotyped behaviour of freely moving fruit flies. *Journal of The Royal Society Interface*, 11(99):20140672, 2014.
- [6] Alexander B. Wiltschko, Matthew J. Johnson, Giuliano Iurilli, Ralph E. Peterson, Jesse M. Katon, Stan L. Pashkovski, Victoria E. Abaira, Ryan P. Adams, and Sandeep Robert Datta. Mapping Sub-Second Structure in Mouse Behavior. *Neuron*, 88(6):1121–1135, December 2015. ISSN 08966273. doi: 10.1016/j.neuron.2015.11.031. URL <https://linkinghub.elsevier.com/retrieve/pii/S0896627315010375>.
- [7] Scott Linderman, Annika Nichols, David Blei, Manuel Zimmer, and Liam Paninski. Hierarchical recurrent state space models reveal discrete and continuous dynamics of neural activity in c. elegans. *BioRxiv*, page 621540, 2019.
- [8] Rich Pang, Christa A. Baker, Mala Murthy, and Jonathan Pillow. Inferring neural population codes for *Drosophila* acoustic communication. *Proceedings of the National Academy of Sciences*, 122(21):e2417733122, May 2025. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.2417733122. URL <https://pnas.org/doi/10.1073/pnas.2417733122>.
- [9] Talmo D. Pereira, Nathaniel Tabris, Arie Matsliah, David M. Turner, Junyu Li, Shruthi Ravindranath, Eleni S. Papadoyannis, Edna Normand, David S. Deutsch, Z. Yan Wang, Grace C. McKenzie-Smith, Catalin C. Mitelut, Marielisa Diez Castro, John D’Uva, Mikhail Kislin, Dan H. Sanes, Sarah D. Kocher, Samuel S.-H. Wang, Annegret L. Falkner, Joshua W. Shaevitz, and Mala Murthy. SLEAP: A deep learning system for multi-animal pose tracking. *Nature Methods*, 19(4):486–495, April 2022. ISSN 1548-7091, 1548-7105. doi: 10.1038/s41592-022-01426-1. URL <https://www.nature.com/articles/s41592-022-01426-1>.
- [10] Adam J. Calhoun, Jonathan W. Pillow, and Mala Murthy. Unsupervised identification of the internal states that shape natural behavior. *Nature Neuroscience*, 22(12):2040–2049, December 2019. ISSN 1097-6256, 1546-1726. doi: 10.1038/s41593-019-0533-x. URL <http://www.nature.com/articles/s41593-019-0533-x>.
- [11] janclemenslab/glm_utils. glm_utils. https://github.com/janclemenslab/glm_utils. Accessed: July 2025.
- [12] Sean Escola. Hidden Markov Models for the Stimulus-Response Relationships of Multistate Neural Systems. 2011.
- [13] Zoe C. Ashwood, Nicholas A. Roy, Iris R. Stone, The International Brain Laboratory, Anne E. Urai, Anne K. Churchland, Alexandre Pouget, and Jonathan W. Pillow. Mice alternate between discrete strategies during perceptual decision-making. *Nature Neuroscience*, 25(2):201–212, February 2022. ISSN 1097-6256, 1546-1726. doi: 10.1038/s41593-021-01007-z. URL <https://www.nature.com/articles/s41593-021-01007-z>.
- [14] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [15] Yoshua Bengio and Paolo Frasconi. An input output hmm architecture. *Advances in neural information processing systems*, 7, 1994.

- 213 [16] Scott W. Linderman, Peter Chang, Giles Harper-Donnelly, Aleya Kara, Xinglong Li, Gerardo
214 Duran-Martin, and Kevin Murphy. Dynamax: A python package for probabilistic state space
215 modeling with jax. *Journal of Open Source Software*, 10(108):7069, 2025. doi: 10.21105/joss.
216 07069. URL <https://doi.org/10.21105/joss.07069>.

A Dataset

We confirm that our research complies with all relevant ethical regulations.

Behavioral chambers. We analyzed behavioral data from 75 male–female pairs of *Drosophila melanogaster* engaged in natural courtship. Behavioral experiments followed the protocol described in the SLEAP software [9]. In brief, flies interacted in custom-fabricated behavioral chambers with a 30mm × 30mm 3D-printed base (Formlabs Form 2, Black V3) and a clear PETG vacuum-molded dome (WidgetWorks Unlimited). Overhead video was captured using a Blackfly S 13YM3-M USB3 camera (FLIR) equipped with an MVL35M23 35mm FL C-mount lens (Thorlabs) and a 25-mm premium 850nm longpass filter (Thorlabs FELH0850). Illumination was provided by 850nm infrared LED strips positioned for side lighting. The arena floor included nine embedded microphone inlets arranged in a 3×3 grid beneath a fine 3D-printed mesh, allowing simultaneous acoustic and behavioral recording. Data acquisition was handled by custom-built workstations with Intel i7-8700K CPUs, 64GB RAM, 4TB Samsung 860 Evo SSDs, and EVGA GeForce GTX 1080 Ti (11GB) GPUs. Videos were recorded from above at 150 frames per second (fps) with a 5ms exposure time and a frame size of 1024×1024 pixels (1 channel), yielding a spatial resolution of 30.3 pixels/mm. Real-time image compression was performed using the Motif recording system and API (Loopbio GmbH), with GPU-accelerated H.264 encoding via the libx264 library (superfast preset). This setup produced nearly lossless videos with independently seekable frames.

Flies. All behavioral experiments were conducted using virgin male and female *Drosophila melanogaster* (wild-type strain NM91), aged 3–5 days post-eclosion, following the protocol in Coen et al., 2014 [3]. Fly bottles were kept at 25C and 60% relative humidity. Experiments were initiated within two hours of incubator lights turning on. Males were single-housed, while females were group-housed prior to experiments. To prevent ceiling walking, the plastic dome of the behavioral chamber was coated with Sigmacote (SL2, Sigma-Aldrich) and allowed to dry under a fume hood for at least 30 minutes before use. Flies were gently introduced into the behavioral chamber using a custom-made aspirator. Recordings were terminated upon copulation or after 30 minutes, whichever occurred first. In total, we recorded 75 NM91 male–female pairs, yielding approximately 22 hours of courtship behavior.

Fly pose estimation and tracking via SLEAP. Fly poses were automatically tracked and manually proofread in all videos using SLEAP. We used the pre-trained ‘flies13’ model published in Pereira et al., 2022 [9], which defines a 13-node skeleton capturing prominent anatomical landmarks: head, thorax, abdomen, left and right wings (wingL, wingR), forelegs (forelegL4, forelegR4), midlegs (midlegL4, midlegR4), hindlegs (hindlegL4, hindlegR4), and eyes (eyeL, eyeR). The skeleton includes 12 edges connecting: thorax to head; thorax to abdomen; each wing to thorax; each leg to thorax; and head to each eye. Male and female identities were tracked using SLEAP’s flow-shift-based identity tracking, followed by manual proofreading and correction of identity switches using the SLEAP GUI. Final joint coordinates and associated confidence scores were exported to .h5 files via the SLEAP API and used for all subsequent analyses.

Song segmentation. Audio was segmented into courtship song using previously described methods [10, 3], with an added pose-based filter to reduce false-positive sine detections. For each audio recording, the segmentation algorithm provided the onset and offset times of pulse bouts and sine trains, as well as the center of each detected pulse. Due to acoustic limitations in the behavioral setup used here, sine detection was prone to noise. Specifically, we retained only those sine bouts where at least one of the male’s wing angles exceeded a threshold of 3deg during the bout.

B Designing sensory inputs

To model female locomotion and wing flicking behavior, we transformed the tracked fly trajectories into a set of f behavioral feedback cues, which served as inputs to the GLM–HMM. For each cue, we extracted a 3 s window of history preceding the current time bin, sampled at 150 Hz, resulting in 450 time points per cue. These temporal windows were projected onto a set of four raised cosine basis functions, yielding four filter coefficients per cue. This produced a $4 \times f$ -dimensional feature vector (f

267 cues \times 4 basis functions). We appended a constant term to this vector to model an intercept, resulting
268 in a final $4 \times f + 1$ -D input vector per time bin.

269 For each fly, we constructed a design matrix of size $T \times (4 \times f + 1)$, where T is the number of time bins
270 after discarding the first 3 s (used for constructing the temporal history). Design matrices from all
271 flies were concatenated to form a population-level dataset, enabling us to fit a single GLM-HMM
272 model across animals.

273 We used a raised cosine basis set to capture temporal structure in the feedback cues. These basis
274 functions are approximately orthonormal and spaced to provide smooth, overlapping temporal filters
275 over the 3-second window, allowing the model to learn coarse-to-fine temporal dependencies while
276 reducing dimensionality. We used the `glm_utils` [11] library to construct raised cosine basis functions
277 and transform or inverse-transform the design matrices.

278 **Output-specific feedback cues.** We fit a single GLM-HMM model with a shared set of latent states
279 across all outputs. However, the set of input features used to predict each behavioral output—forward
280 velocity (fFV), lateral velocity (fLV), angular velocity (fAV) and wing flicking—was distinct, re-
281 flecting the sensory cues most relevant to that dimension of behavior. In particular, the input sets for
282 fLV and fAV included interaction terms (e.g., cue \times side) to capture directional effects that are not
283 applicable to fFV or wing flicking behavior. This was implemented by applying an output-specific
284 mask during the M-step in the EM algorithm, such that only the designated subset of input features
285 contributed to the GLM weights for each output.

286 Inputs used to predict the female’s forward velocity (fFV) included $f=7$ unsigned cues (Figure 1):
287 male forward velocity (mFV), male–female distance (mfDist), male lateral speed (mLS), female
288 heading to male thorax angle (fmAng), binary pulse and sine song, and male tapping. Each of these
289 cues was temporally smoothed using four raised cosine basis functions, resulting in 4 coefficients per
290 cue. This produced a feature vector of size $4 \times 7 + 1 = 29$, where the final element is a bias offset term.

291 For lateral (fLV) and angular velocity (fAV), the same base cues were used, but with directional
292 (signed) information introduced through cue \times side interaction terms. These allowed the model to
293 capture the effects of asymmetric stimuli—such as song played to the left or right of the female—on
294 her turning and lateral movement. The full set of signed cues included: mFV \times side, mLS \times side,
295 mfDist \times side, wing alignment \times side, pulse \times side, sine \times side, and tap \times side. This resulted in a
296 feature vector of size $4 \times 7 + 1 = 29$ each. The sine component, $\sin(\text{fmAng})$, corresponds to the “side”
297 variable here (and as shown in Figure 1), indicating whether the male is to the left or right of the
298 female.

299 For predicting wing flicking behavior, $f=8$ unsigned cues were used: mFV, mLS, mfDist, fmAng,
300 wing alignment, male tapping, and binary pulse and sine song. This resulted in a feature vector of
301 size $4 \times 8 + 1 = 33$.

302 **Encoding male position relative to female.** To represent the relative angular position of the male
303 with respect to the female, we used both the cosine and sine of the angle “fmAng” between their
304 orientations (Figure 1). This circular encoding captures the full 360 directional relationship in a
305 smooth and continuous way, avoiding discontinuities at the angle wraparound (e.g., near 0/180).

306 **Z-scoring sensory inputs.** All input features were z-scored independently for each fly to ensure
307 comparability across individuals and to standardize the scales of different cues. For binary features
308 such as pulse and sine song, we applied “safe” z-scoring: if a feature had near-zero variance ($\text{std} <$
309 $1e-2$, e.g., present in only a few frames), it was set to zero entirely to avoid instability during model
310 fitting. Otherwise, the feature was z-scored normally.

311 **Z-scoring female behavioral outputs.** All behavioral output variables—forward velocity (fFV),
312 lateral velocity (fLV), and angular velocity (fAV)—were standardized independently for each fly by
313 subtracting the mean and dividing by the standard deviation. Binary outputs, such as wing flicking,
314 were left untransformed (0 or 1).

315 **Smoothing.** All male sensory input variables were smoothed using a causal half-Gaussian kernel
316 ($\sigma = 3$ frames (20 ms)), truncated at 4σ (12 frames, 80ms). Female behavioral outputs were also
317 smoothed using the same kernel applied causally to avoid introducing future information. Female

318 outputs were then downsampled to 30 Hz by averaging within non-overlapping 33 ms windows (5
 319 frames at 150 Hz). This preprocessing reduced high-frequency noise while preserving fast behavioral
 320 dynamics relevant to model prediction.

321 C Modeling

322 Chance model

323 We constructed the Chance baseline model using data from the entire courtship dataset, regardless
 324 of state or sensory input. This model estimates the emission distribution as a Gaussian with mean
 325 and covariance computed from all behavioral observations pooled together, for each continuous
 326 emission. For the discrete output (wing flicking), the Chance model uses a Bernoulli distribution
 327 with probability equal to the fraction of wing flicks observed in the dataset (i.e., n/N where n is
 328 the number of time points during the courtship with wing flicking and N is the total number of time
 329 points during the entire courtship).

330 GLM-HMM

331 The relationship between internal state, sensory input, and behavioral output is effectively modeled
 332 by a Generalized Linear Model–Hidden Markov Model (GLM-HMM) [10, 12, 13], which captures
 333 discrete latent states corresponding to distinct mappings from sensory cues to behavior (Figure 3). In
 334 this framework, each latent state z_t defines a separate generalized linear model relating the male cue
 335 history s_t to the observed female behavior y_t , including forward, lateral, and angular velocities, as
 336 well as wing flicking. The emission distribution at time t is:

$$p(y_t | z_t, s_t) = \mathcal{N}(y_t | w_{z_t} s_t + b_{z_t}, \sigma_{z_t}^2) \quad \text{for continuous outputs (velocities)} \quad (1)$$

$$p(y_t | z_t, s_t) = \sigma(w_{z_t} s_t + b_{z_t}) \quad \text{for binary outputs (wing flicking)} \quad (2)$$

337 **Inference of GLM-HMM parameters.** When fitting the GLM-HMM, the goal is to estimate the
 338 parameters that govern both the latent state dynamics and the emission model of female behavior.
 339 These parameters include the initial state distribution $\pi \in \mathbb{R}^K$, the state transition matrix $A \in \mathbb{R}^{K \times K}$,
 340 and the set of emission weights $w_k^{(j)} \in \mathbb{R}^M$, biases $b_k^{(j)}$, and (for continuous emissions) covariances
 341 $\sigma_k^{(j)} \in \mathbb{R}$ for each latent state k . Here j indexes emission variables (e.g., forward, lateral, angular,
 342 wing flicking). We denote the full set of parameters as $\theta = \{\pi, A, w_k^{(j)}, b_k^{(j)}, \sigma_k^{2(j)}\}$.

343 These parameters were fit to the female behavioral data using maximum a posteriori (MAP) estimation,
 344 implemented via the Expectation-Maximization (EM) algorithm. The EM algorithm has previously
 345 been adapted to fit hidden Markov models with external inputs [12, 14, 10, 13, 15]. However,
 346 since several implementation details are application-specific, we include a full description of the
 347 procedure here for completeness. The EM algorithm seeks to maximize the log-posterior of the
 348 model parameters given the female behavior data Y and sensory input features S . The log-posterior
 349 is given by, up to an unknown constant:

$$\begin{aligned} \log p(\theta | Y, S) &= \log p(Y | S, \theta) + \log p(\theta) \\ &= \log \sum_Z p(Y, Z | S, \theta) + \log p(\theta) && \text{(marginalization)} \\ &= \log \sum_{z_{1:T}} p(y_{1:T}, z_{1:T} | s_{1:T}, \theta) + \log p(\theta) && \text{(expanding)} \end{aligned} \quad (3)$$

350 where the sum is taken over all K^T possible latent state sequences $z_{1:T}$. The first term represents the
 351 log-likelihood or the log-posterior of the observed data under the model, and the second term is a
 352 prior on the parameters.

353 **Priors.** The prior distribution over the model parameters θ was assumed to factorize as follows:

$$\begin{aligned}
p(\theta) &= p(\{w_k^{(j)}\}) \cdot p(A) \cdot p(\pi) \\
&= \left[\prod_{k=1}^K \prod_j \mathcal{N}(w_k^{(j)} \mid 0, \lambda_j^{-1}) \right] \left[\prod_{k=1}^K \text{Dirichlet}(A_k \mid \alpha_k) \right] \text{Dirichlet}(\alpha_\pi) \quad (4)
\end{aligned}$$

354 We placed a zero-mean Gaussian prior on each GLM weight vector $w_k^{(j)}$ where λ_j is the inverse
355 variance and controls the strength of regularization for emission variable j . Larger values of λ_j
356 have a shrinking effect on the fitted weights, biasing them toward zero. For the continuous emission
357 variables (forward, lateral, and angular velocity), we set the regularization parameter $\lambda_j = 10^{-6}$. For
358 the discrete wing flicking emission, which uses a Bernoulli emission model, we used a stronger prior
359 with $\lambda_j = 1$.

360 The transition matrix and initial state distribution were each given Dirichlet priors with symmetric
361 concentration parameters. For the transition matrix A , we used a structured Dirichlet prior over
362 each row A_k that encourages self-transitions (i.e., persistence within states) where the concentration
363 parameters $\alpha_k \in \mathbb{R}^K$ were set as $\alpha_k = \alpha \cdot \mathbf{1}_K + \kappa \cdot \mathbf{e}_k$. Here, $\alpha = 1.1$ is a weakly informative
364 base concentration applied to all transitions, $\kappa = 100$ is a stickiness parameter that adds mass to
365 the diagonal (self-transition) entry, $\mathbf{1}_K$ is a vector of ones, and \mathbf{e}_k is a one-hot vector indicating the
366 k -th state. This form biases the prior toward self-transitions while still allowing transitions to other
367 states. For the initial state distribution $\alpha_\pi = \alpha \cdot \mathbf{1}_K$, with $\alpha = 1.1$. This weakly informative prior
368 encourages a broadly uniform initial state distribution while still allowing the model to learn the
369 estimate of π from the data.

370 **Fitting using Expectation-Maximization (EM) algorithm.** We used the EM algorithm to maxi-
371 mize the log-posterior given in Eq. 3 with respect to the GLM-HMM parameters. As the sum involves
372 an exponential number of terms— $\mathcal{O}(K^T)$ to be specific—we do not maximize this expression directly.
373 Instead, the EM algorithm provides an efficient way to compute this term using a single forward and
374 backward pass over the data. During the E-step of the EM algorithm, we compute the ‘expected
375 complete data log-likelihood’ (ECLL), which is a lower bound on the right-hand side of Eq. 3. Then,
376 during the ‘maximization’ or M-step of the algorithm, we maximize the ECLL with respect to the
377 model parameters θ . It can be shown that this procedure has the effect of always improving the
378 log-posterior in each step of the algorithm and converges to a local optimum of the log-likelihood
379 [12, 14].

380 The ‘complete data log-likelihood’ (CLL) for a session is written as $\log P(Y, Z|S; \theta)$:

$$\begin{aligned}
CLL(\theta) &= \log P(Y, Z|S; \theta) \\
&= \log P(y_{1:T}, z_{1:T} | s_{1:T}; \theta) \\
&= \log \left[P(z_1 \mid \pi) \prod_{t=2}^T P(z_t | z_{t-1}, A) \prod_{t=1}^T P(y_t | z_t, s_t, w_k, b_k, \Sigma_k) \right] \\
&= \log \pi_{z_1} + \sum_{t=2}^T \log A_{z_{t-1}, z_t} + \sum_{t=1}^T \log B_{z_t}(y_t, s_t)
\end{aligned}$$

381 where $B_{z_t}(y_t, s_t)$ is the Gaussian and Bernoulli distribution given by the emission model equation
382 (Eq. 2).

383 The Expected-CLL or the ECLL for a session, where the expectation is with respect to the distribution
384 over the latents $\sum_Z p(Z \mid Y, S; \theta_{old})$ computed during the E-step, can now be written as:

$$\begin{aligned}
ECLL(\theta) &= \sum_Z P(Z | Y, S; \theta^{old}) CLL \\
&= \sum_Z P(Z | Y, S; \theta^{old}) \log P(Y, Z | S; \theta) \\
&= \sum_Z \log \pi_{z_1} P(Z | Y, S; \theta^{old}) + \sum_Z \sum_{t=2}^T \log A_{z_{t-1}, z_t} P(Z | Y, S; \theta^{old}) + \sum_Z \sum_{t=1}^T \log B_{z_t}(y_t, s_t) P(Z | Y, S; \theta^{old}) \\
&\vdots \\
&= \sum_{k=1}^K \log \pi_k \gamma_k(1) + \sum_{j=1}^K \sum_{k=1}^K \sum_{t=2}^T \log A_{jk} \xi_{j,k}(t) + \sum_{k=1}^K \sum_{t=1}^T \log B_k(y_t, s_t) \gamma_k(t) \tag{5}
\end{aligned}$$

Here, we denote $\gamma_k(t) = P(z_t = k | Y_{1:T}, s_{1:T}, \theta^{old})$ for the posterior state probability of being in state k at time point t , and $\xi_{j,k}(t) = P(z_{t-1} = j, z_t = k | Y, S; \theta^{old})$ is the joint posterior state distribution for two consecutive latents z_t and z_{t-1} . We compute these two posterior distributions γ and ξ in the E-step as below:

E-step

The E-step of the EM algorithm involves computing the posterior distribution $P(Z | Y, \theta_{old})$ over the hidden variables given the data and the current setting of the GLM-HMM parameters θ_{old} using the forward-backward algorithm. The forward-backward algorithm makes use of recursion and memoization to allow these posterior probabilities to be calculated efficiently, with the forward and backward passes of the algorithm each requiring just a single pass through the whole session.

The goal of the forward pass is to obtain, for each time point t within a session and each state k , the quantity $a_i(t) = P(Y_1 = y_1, Y_2 = y_2, \dots, Y_t = y_t, z_t = i | s_{1:t})$ of observing $Y = y_1, y_2, \dots, y_t$ which represents the posterior probability of the female behavior data up until time t and the latent state at time t being state k . Assuming there are K total states, it can be recursively computed as:

$$a_j(t+1) = \sum_{k=1}^K a_k(t) A_{jk} B_j(y_t, s_t)$$

where $a_j(1) = \pi_j P(y_1 | z_t = j, s_1)$ and $B_j(y_t, s_t) = P(y_t | z_t = j, s_t)$ is the usual Gaussian or Bernoulli GLM distribution.

During the backward pass, the goal is to calculate the posterior probability of the future behavior data given the latent state $b_j(t) = P(Y_{t+1} = y_{t+1}, \dots, Y_T = y_T | z_t = j, s_{t+1:T})$, as follows:

$$b_j(t) = \sum_{k=1}^K b_k(t+1) A_{jk} B_k(y_{t+1}, s_{t+1})$$

where $b_j(T) = 1$.

From the $a_j(t)$ and $b_j(t)$ quantities obtained from the forward-backward algorithm, we can compute the posterior state distribution γ over the latent state at every time step (this uses data from the whole session):

$$\begin{aligned}
\gamma_k(t) &= P(z_t = k | Y_{1:T}, s_{1:T}, \theta_{old}) \\
&= \frac{P(z_t = k, Y_{1:T} | s_{1:T}, \theta_{old})}{P(Y_{1:T} | s_{1:T}, \theta_{old})} \\
&= \frac{P(Y_{1:t}, z_t = k | s_{1:t}, \theta_{old}) \cdot P(Y_{t+1:T} | z_t = k, s_{t+1:T}, \theta_{old})}{P(Y_{1:T} | s_{1:T}, \theta_{old})} \\
&= \frac{a_k(t) \cdot b_k(t)}{\sum_{i=1}^K a_i(t) \cdot b_i(t)} \tag{6}
\end{aligned}$$

407 Similarly, we can obtain the joint posterior state distribution ξ for the consecutive latents:

$$\begin{aligned}\xi_{j,k}(t) &= P(z_{t-1} = j, z_t = k \mid Y, S; \theta_{old}) \\ &= \frac{a_j(t) A_{jk} b_j(t+1) B_k(y_{t+1}, s_{t+1})}{\sum_{i=1}^K a_i(t) \cdot b_i(t)}\end{aligned}\quad (7)$$

408 Having now computed γ and ξ , ECLL is now a simply a function of model parameters θ with every
409 other term known (Eq. 5).

410 **M-step**

411 After running the forward-backward algorithm, we can compute the total ECLL by summing over
412 the per-session ECLLs (Eq. 5) and adding the log-prior (Eq. 4). During the M-step, we maximize the
413 ECLL with respect to the GLM-HMM parameters θ . This uses the smoothed state probabilities $\gamma_t(k)$
414 and $\xi_{j,k}(t)$ computed during the E-step (Eqs. 6 and 7). For the initial state distribution π , transition
415 matrix A and GLM weights for continuous emissions, this results in closed-form updates. The initial
416 state probability π_k is updated as:

$$\pi_k^{new} = \frac{\sum_{e=1}^E \gamma_1(k)}{E},$$

417 The updated transition probabilities A_{ij} are given by the mode of the posterior Dirichlet distribution:

$$A_{ij}^{new} = \frac{\alpha - 1 + \sum_{e=1}^E \sum_{t=2}^T \xi_{i,j}(t)}{\sum_{j'=1}^K (\alpha - 1 + \sum_{e=1}^E \sum_{t=2}^T \xi_{i,j'}(t))},$$

418 Because these GLMs contribute independently to the terms A , π and emission terms B , we can
419 optimize the filters for each output dimension separately.

420 In case of continuous emissions (forward, lateral and angular velocity), each state-specific emission
421 model assumes a Gaussian distribution over the output $y_t^{(j)} \in \mathbb{R}$ ($j = 1 \dots 3$) with mean linearly
422 dependent on the input vector $s_t \in \mathbb{R}^M$:

$$y_t^{(j)} \mid z_t = k, s_t, \theta \sim \mathcal{N}(w_k^{(j)} \cdot s_t, \sigma_k^{2(j)}) \quad (\text{for forward, lateral and angular velocity emissions})$$

423 To estimate the GLM weights $\{w_k^{(j)}, b_k^{(j)}, \sigma_k^{2(j)}\}$, we solve the weighted multivariate linear regression
424 problem for each state k . For notational simplicity, we assume $w_k^{(j)}$ includes the bias term $b_k^{(j)}$,
425 with s_t augmented by a constant 1. To estimate the weights for each state k in the GLM-HMM, we
426 pooled sufficient statistics across all sessions; thus, the variables $\gamma_k(t)$, y_t , and s_t below represent
427 data concatenated across sessions.

428 The state-specific linear weights $w_k^{(j)} \in \mathbb{R}^{M+1}$ and the emission covariances $\sigma_k^{2(j)} \in \mathbb{R}$ have
429 closed-form solution for the updates given by:

$$\begin{aligned}w_k^{new(j)} &= \left(\sum_{t=1}^T \gamma_t(k) y_t s_t^\top \right) \left(\sum_{t=1}^T \gamma_t(k) s_t s_t^\top + \lambda_j I \right)^{-1} \\ \sigma_k^{2new(j)} &= \frac{1}{\sum_{t=1}^T \gamma_t(k)} \sum_{t=1}^T \gamma_t(k) \left(y_t^{(j)} - w_k^{(j)} \cdot s_t \right) \left(y_t^{(j)} - w_k^{(j)} \cdot s_t \right)^\top\end{aligned}$$

430 For numerical stability, we added a small constant to the estimated covariance $\sigma_k^{2new(j)} \leftarrow \sigma_k^{2new(j)} +$
431 10^{-8} .

432 For binary outputs such as wing flicking ($y_t^{(j)} \in \{0, 1\}$), Bernoulli GLM weights have no such
433 closed-form update.

$$P(y_t^{(j)} \mid z_t = k, s_t, \theta) = \sigma(w_k^{(j)} \cdot s_t) \quad (\text{for wing flicking})$$

We use the Dynamax Python package [16] to minimize the negative ECLL for all emission models. For continuous emissions, closed-form updates are implemented in the package. For Bernoulli emissions, Dynamax performs gradient-based optimization using the Adam optimizer, which is implemented via the Optax library.

Initializing GLM-HMM weights. We first fit a single-state linear regression model (i.e., a GLM without latent states) to each behavioral output. The estimated weights were then used to initialize the emission parameters of the GLM-HMM, with small random noise added independently to each state’s parameters to break symmetry and encourage state specialization.

D Assessing model performance

Normalized Test LogLikelihood. We assessed model performance by calculating the log-likelihood of data held-out from training. We held-out entire sessions of courtship data for assessing test set performance. That is, when fitting the model, the ECLL in Eq.5 are modified to include only 80% of sessions (because we use five-fold cross-validation throughout this work); and the log-likelihood of the held-out 20% of sessions E' is calculated using the fit parameters θ and a single run of the forward pass on the held-out sessions. In particular, we assessed how well the model predicted the next output given knowledge of all the data up to the present moment. In practice, it can be computed as:

$$\text{LL}_{\text{forward}}(\text{model}) = \sum_{\text{Test set}} \log \sum_{k=1}^K a_k(T) \quad (8)$$

that is, the sum of the last column of the a matrix obtained after doing a single forward pass on a test session.

To report the log-likelihood in more interpretable units, we normalized by subtracting the log-likelihood under the Chance model (described above; Figure 3), as follows:

$$\text{LL}_{\text{norm}}(\text{model}) = \text{LL}_{\text{forward}}(\text{model}) - \text{LL}_{\text{forward}}(\text{chance})$$

The chance model was drawn from the full distribution of behavior across all courtship recordings (Figure 3). To express this in interpretable units, we report $\text{LL}_{\text{norm}}(\text{model})$ as bits per second, by dividing it by the total duration of courtship in seconds. The normalized log-likelihood of the forward model thus reports the improvement in predicting female behavior over the Chance model, based on knowledge of her history to better estimate the current state.

State Inference. Latent states were inferred using forward filtering in the GLM-HMM framework. For each time point t , we computed the predictive state distribution $\hat{\gamma}_t$ over states using all observations up to the previous time point $t - 1$:

$$\hat{\gamma}_k(t) = P(z_t = k \mid s_{1:t}, y_{1:t-1})$$

This procedure is applied after training the GLM-HMM, using the learned transition and emission parameters to decode state sequences on the training and held-out data. It can be computed and stored using intermediate values during the calculations of the matrix a and γ during a forward pass in the E-step. For visualization purposes (e.g., state sequences over time), we assigned each time point to the most probable state ($\arg \max_k \hat{\gamma}_t(k)$). However, for model predictions, we used the soft state probabilities to compute a weighted sum of outputs across all states.

Behavior prediction. Rather than relying on hard state assignments, the model prediction at each time point was taken as a weighted sum over the predictions from all latent states, with weights given by the predictive state probabilities (from the forward filtering algorithm). Formally, for a behavioral output y_t , our model prediction \hat{y}_t is given by:

$$\begin{aligned}\hat{y}_t &= \sum_{k=1}^K p(z_t = k \mid s_{1:t}, y_{1:t-1}) \cdot \hat{y}_{t,k} & \text{where } \hat{y}_{t,k} &= w_k \cdot s_t + b_k \\ &= \sum_{k=1}^K \hat{\gamma}_k(t) \cdot \hat{y}_{t,k}\end{aligned}\tag{9}$$

where $\hat{y}_{t,k}$ is the GLM prediction from state k for one of the continuous velocity emission predictions. However, for wing flick predictions, we used hard state assignments (obtained from the forward filtering step and using the state with the maximum probability at each time step; importantly, we did not use the Viterbi algorithm for state inference) and included an additional sigmoid nonlinearity (assuming animal in state k at time t):

$$P(\hat{y}_t) = \sigma(w_k \cdot s_t + b_k)$$

Pearson correlation coefficient To evaluate model performance for the continuous velocity emissions, we computed the Pearson correlation score between the observed behavioral output y and the model's soft predictions \hat{y} as defined above. We compute:

$$r = \frac{\sum_t (y_t - \bar{y})(\hat{y}_t - \bar{\hat{y}})}{\sqrt{\sum_t (y_t - \bar{y})^2} \sqrt{\sum_t (\hat{y}_t - \bar{\hat{y}})^2}}.$$

where \bar{y} and $\bar{\hat{y}}$ denote the mean of the observed y_t and predicted \hat{y}_t outputs, respectively. This correlation captures the linear relationship between predicted and observed signals while incorporating uncertainty in latent state identity. The correlation score was computed separately for each fly and for each behavioral output variable (forward, lateral, and angular velocity). In Figure 3c, we report the mean Pearson correlation score averaged across these three output dimensions.

Pearson correlation coefficient per state. To evaluate how well each latent state predicts continuous behavioral outputs, we computed a state-specific version of the Pearson correlation coefficient $\hat{r}^{(k)}$ using soft assignments (Figure 3g). Specifically, we computed a soft-assignment weighted Pearson correlation coefficient between the true behavioral output y_t and the state-specific prediction $\hat{y}_{t,k}$, for each state k . This method incorporates the posterior state probabilities $\hat{\gamma}_t(k)$ (predictive state distribution using data up to time point $t - 1$) as weights. For each state k , the steps were as follows:

- Compute the weighted means:

$$\mu_y^{(k)} = \frac{\sum_t \hat{\gamma}_t(k) y_t}{\sum_t \hat{\gamma}_t(k)}, \quad \mu_{\hat{y}}^{(k)} = \frac{\sum_t \hat{\gamma}_t(k) \hat{y}_{t,k}}{\sum_t \hat{\gamma}_t(k)}$$

- Compute the weighted covariance:

$$\text{Cov}^{(k)} = \frac{\sum_t \hat{\gamma}_t(k) (y_t - \mu_y^{(k)}) (\hat{y}_t - \mu_{\hat{y}}^{(k)})}{\sum_t \hat{\gamma}_t(k)}$$

- Compute the weighted variances:

$$\text{Var}_y^{(k)} = \frac{\sum_t \hat{\gamma}_t(k) (y_t - \mu_y^{(k)})^2}{\sum_t \hat{\gamma}_t(k)}, \quad \text{Var}_{\hat{y}}^{(k)} = \frac{\sum_t \hat{\gamma}_t(k) (\hat{y}_t - \mu_{\hat{y}}^{(k)})^2}{\sum_t \hat{\gamma}_t(k)}$$

- Finally, compute the weighted Pearson correlation:

$$\hat{r}^{(k)} = \frac{\text{Cov}^{(k)}}{\sqrt{\text{Var}_y^{(k)} \cdot \text{Var}_{\hat{y}}^{(k)} + \varepsilon}}$$

where ε is a small constant added for numerical stability.

497 **F1 score** To evaluate model performance on wing flicking (Figure 3d), we computed the F1 score
498 between the observed binary behavioral output y and the model’s predicted output \hat{y} . The predicted
499 output was thresholded at 0.5 to yield a binary classification, and F1 score was computed across all
500 time points.

501 **F1 score per state.** To evaluate model performance on wing flicking in each state (Figure 3h),
502 we used hard state assignments obtained by taking the most probable state at each time point
503 ($\arg \max_k \hat{\gamma}_t(k)$). For each state, we then computed the standard F1 score between the observed and
504 predicted binary outputs, using only the time points assigned to that state.

505 **Cross-validation.** To select the appropriate number of latent states, we performed cross-validation
506 by splitting the dataset into training and test sets. For each candidate model (with a different number
507 of states), we fit the model parameters on the training data and evaluated performance on held-out
508 test data Figure 3b. We used multiple random train/test splits to ensure robustness of the model
509 comparison Figure 6.

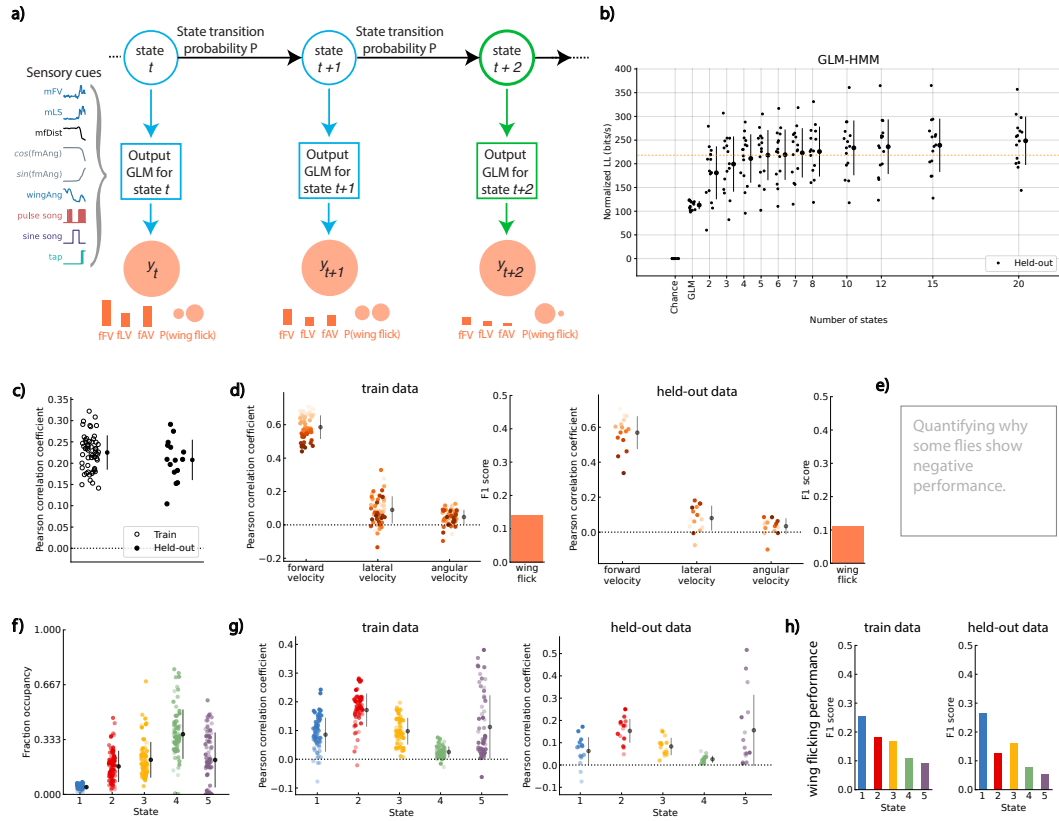


Figure 3: Model architecture and performance. (a) Schematic of the GLM-HMM framework. At each time point t , the female is assumed to be in a latent internal state that determines how multimodal male sensory cues are linearly weighted to predict her motor outputs: forward (fFV), lateral (fLV), and angular velocity (fAV), along with wing flicks. (b) Normalized log-likelihood (LL) on test data (in bits/s; see Appendix). Each circle represents one courtship pair. (c) Pearson correlation between predicted and observed female behavioral outputs, for both training (60 fly pairs) and held-out test data (15 fly pairs), using the five-state GLM-HMM. Each open or filled circle represents one courtship pair. (d) Model performance (Pearson correlation) broken down by behavioral variable—forward velocity (fFV), lateral velocity (fLV), and angular velocity (fAV)—for training (left) and held-out (right) fly pairs. Each dot represents one courtship pair. A subset of fly pairs also exhibit near-zero or negative performance, suggesting inter-fly variability. Within each panel, flies are color-coded consistently: the same shade of orange denotes the same fly pair across behaviors. (e) Quantifying why some flies show negative performance. (f) Fractional occupancy of each latent state across fly pairs in the 5-state GLM-HMM. Each dot denotes the proportion of time a given fly pair spent in a particular state; black markers indicate the mean \pm sem. Together with the variability in predictive performance across states shown in g, differences in fractional occupancy suggest that some states are more behaviorally informative or more frequently drive motor output than others. (g) Female behavioral predictability by latent state, using soft state assignments from the GLM-HMM, on training data (left) and held-out (right). Pearson correlation between predicted and observed outputs is computed within each state by weighting time points proportionally to their inferred state probabilities (See Appendix). States vary in how strongly female behavior is predicted by male cues, suggesting that certain internal states reflect periods of stronger sensory coupling. Black markers indicate the mean \pm 1 s.d. (h) Wing flicking F1 scores by state.

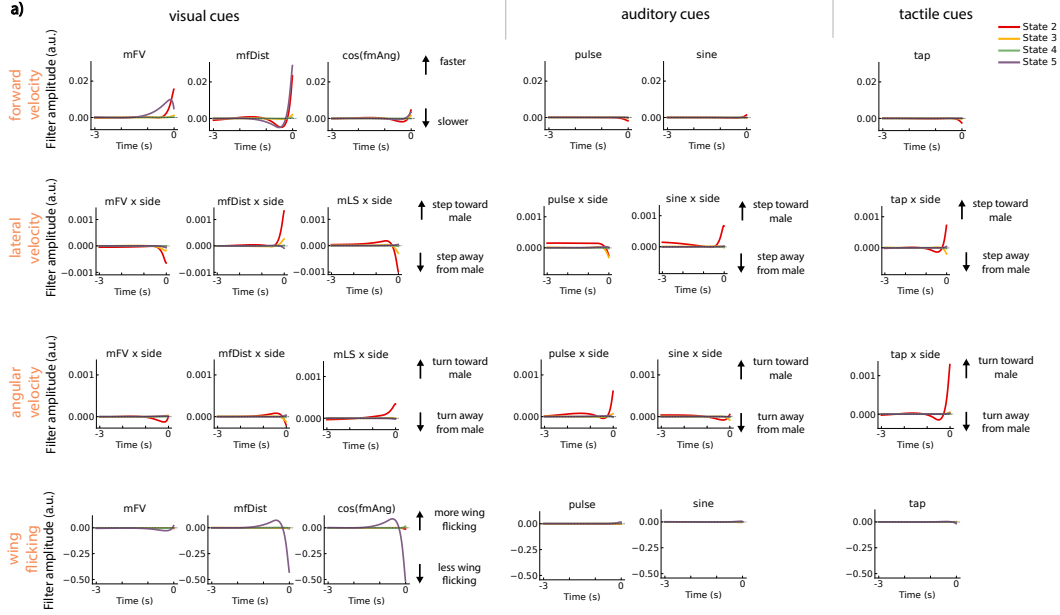


Figure 4: State-dependent sensory filters reveal modality-specific drive of female locomotion. (a) Output filters that predict female velocities for some of the feedback cues. For each latent state (State 2-5) inferred by the five-state GLM-HMM, we plot the GLM filter that converts a 3-second history of male cues into predicted female velocities. Columns are grouped by visual (left), auditory (center), and tactile (right) cues; rows show the effect on female forward (top), lateral (middle), and angular (bottom) velocity. Positive filter amplitude predicts faster forward motion, steps or turns toward the male (\uparrow); negative amplitude predicts slowing or motion away (\downarrow). States 2 and 5 show strong filters, whereas the states 3 and 4 filters are nearly flat, indicating minimal sensorimotor coupling. Cue abbreviations are in Figure 1. Interaction terms ($\times side$) in the lateral and angular filters capture cue laterality: for mFV, mfDist, mLS and tap cues, *side* indicates which side of the female the male's thorax occupies, whereas for pulse and sine song, it specifies whether the singing wing is on her left or right.

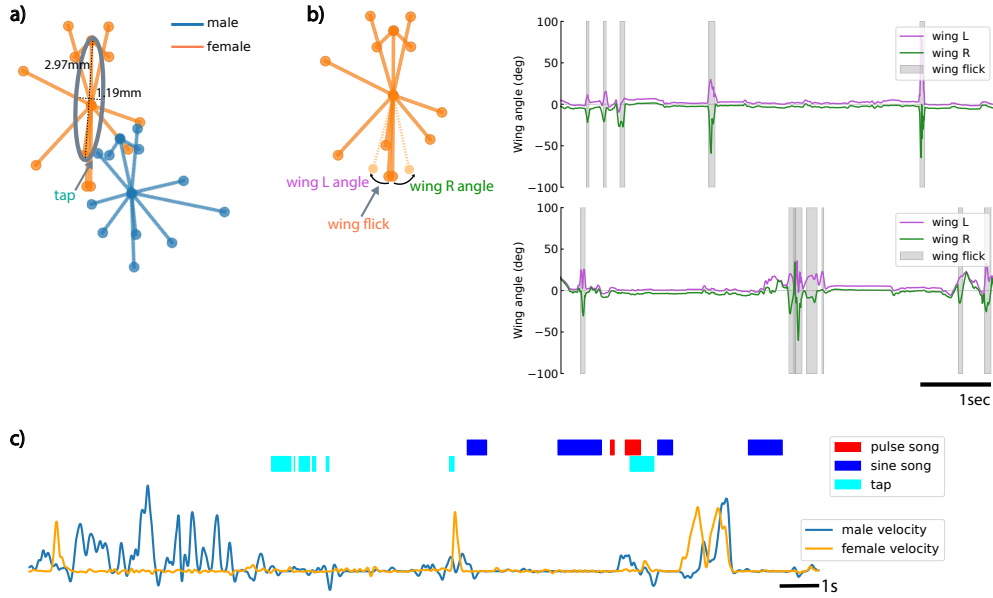


Figure 5: Sensory features. (a) Taps are identified using a heuristic based on proximity: a tap is registered when the male foreleg tip enters the ellipsoidal region containing the female from head to abdomen. Shown is an example frame where a male tap (cyan) is detected. (b) Left: Female wing flicking is detected when the left and right wing angles deviate by more than 20° from each other. Wing angles are measured relative to the body axis using tracked the wing tip positions. Right: Top and bottom panels each show a 5-second segment of left wing (magenta) and right wing (green) angle traces from two different flies. Gray shaded bars mark detected wing flick events. (c) A representative bout showing male and female forward velocity over time, overlaid with male sensory events (pulse and sine song, and taps). Female responses do not show consistent moment-to-moment coupling with male cues, motivating the need for a latent state model to explain longer-timescale structure.

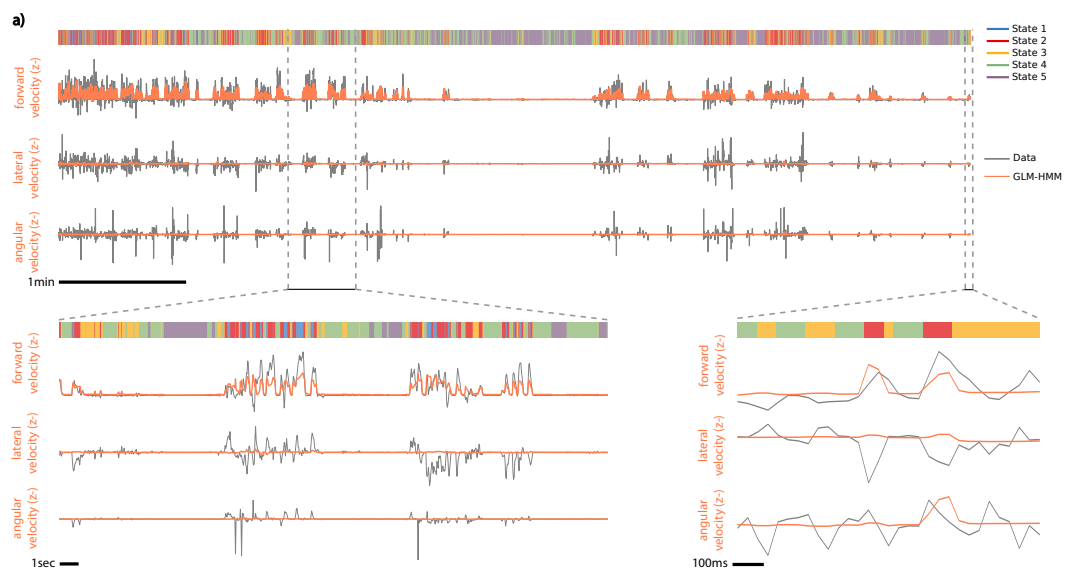


Figure 6: GLM-HMM predictions of female velocity dynamics across multiple timescales. (a) Top row: continuous latent-state assignments (five colors) over a full courtship session of 6-minutes. Below, raw z-scored forward, lateral, and angular velocity traces (gray) are overlaid with GLM-HMM predictions (orange). Two dashed-box insets show progressively shorter epochs: the left inset (expanded below) spans 30-seconds, revealing rapid transitions and high-frequency fluctuations; the right inset (expanded at bottom right) spans 1-second. For this analysis, we assigned each time point to its most likely state during the filtering E-step, shown on top. The velocity predictions however use "soft" state assignments combining probabilistic contributions from all states.

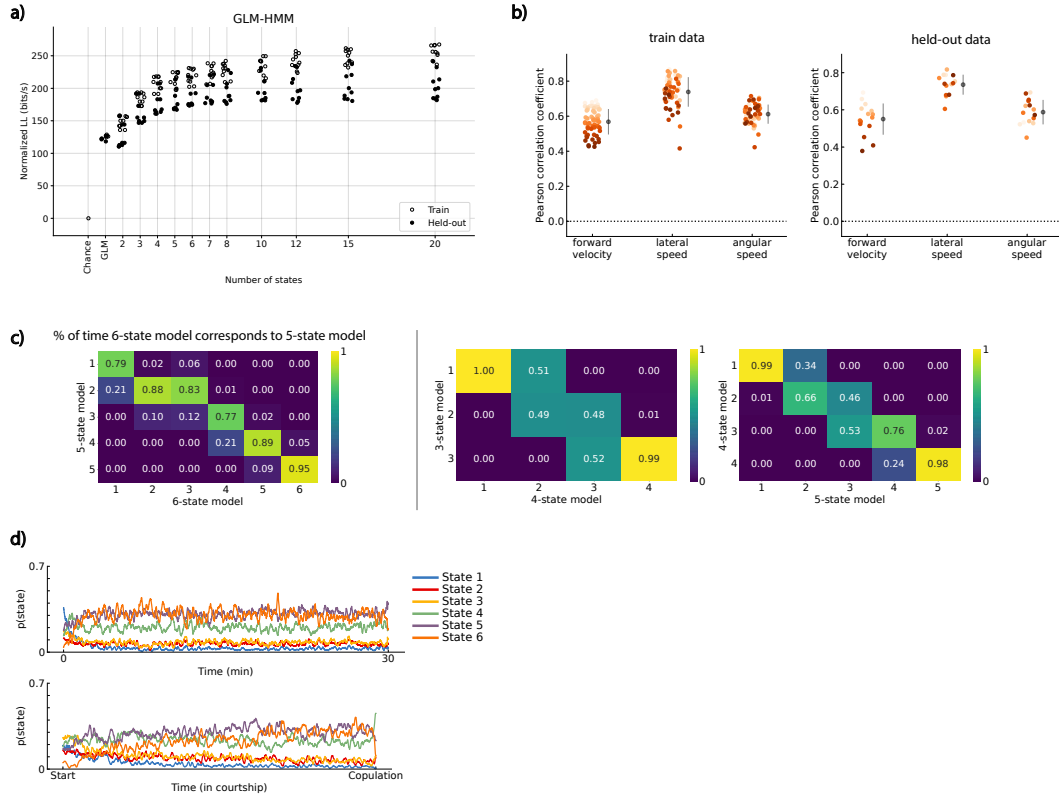


Figure 7: Evaluating the states of the GLM-HMM. **(a)** CV for different train-test splits. **(b)** Performance (Pearson’s correlation) of the model trained to predict forward velocity (fFV), lateral speed (fLS), and angular speed (fAS), shown separately for training pairs (left) and held-out pairs (right). Lateral and angular speeds—rather than signed velocities—are being predicted here that do not have information about the direction of female turning. Prediction performance is substantially higher for lateral and angular speeds than for signed velocities (see Figure 3), implying that directional movements are more variable and harder to predict, whereas changes in movement magnitude are more consistently driven by male cues and internal state. Each dot represents one courtship pair. Within each panel, flies are color-coded consistently: the same shade of orange denotes the same fly pair across behaviors. Black markers indicate the mean \pm 1 s.d. **(c)** Left: The correspondence between the 5-state GLM-HMM and the 6-state GLM-HMM. Shown is the conditional probability of the 5-state model being in the one of its states given the state of the 6-state model. State 2 and State 3 in the 6-state model both correspond to State 2 of the 5-state model most of the time. Right: Same analysis showing correspondences across multiple model comparisons: 3 \rightarrow 4 and 4 \rightarrow 5. **(d)** Same as Figure 2 for 6-state model. State 2 (red) and State 3 (yellow) overlap with each other for most of the session, indicating the two aren’t distinguishable and are identical to State 2 from the 5-state model.

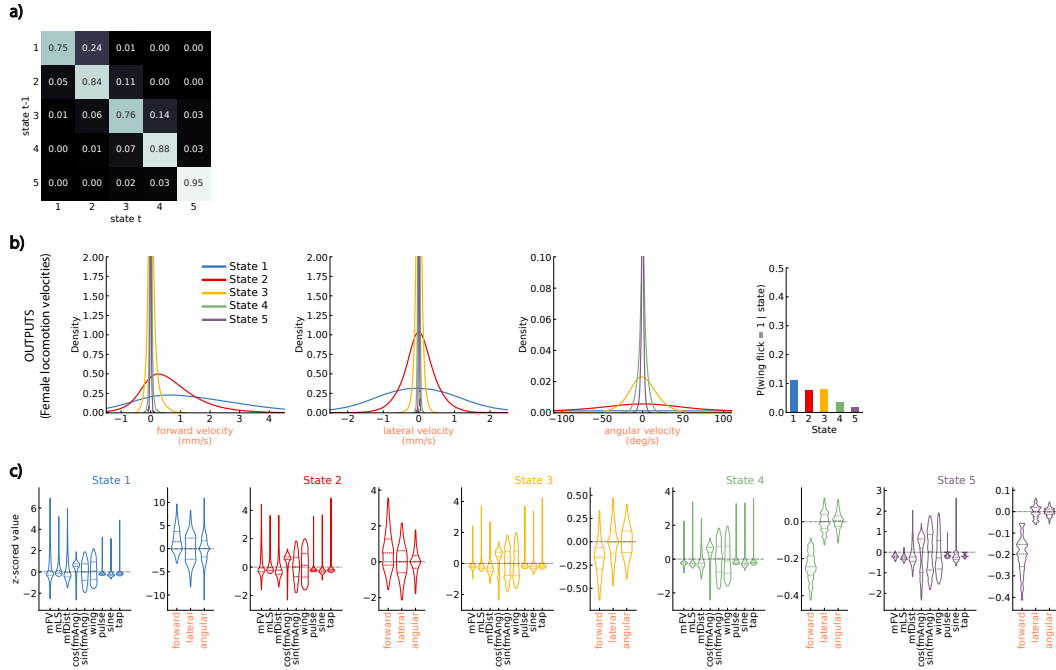


Figure 8: State-wise distributions of behavioral outputs and sensory inputs. (a) Inferred full transition matrix of the 5-state GLM–HMM, showing the probability of transitioning from each state (rows) to every other state (columns). Large entries along the diagonal indicate a high probability of remaining in the same state ("sticky"). State 1, interpreted as a "noisy" or "chamber introduction" state, is the least stable and primarily shows one-way transition into State 2 (also see Figure 2). (b) Distributions of each female behavioral output variable within each latent state. (c) Distributions of male sensory cues within each latent state, complementing Figure 2 where only the mean z-scored value per cue was shown.