

# Supporting Information

Giret et al. 10.1073/pnas.1317087111

## SI Methods

All experiments were done using methods approved by the Veterinary Office of the Canton of Zurich, Switzerland.

**Loop Delay Estimation.** Motor and auditory delays were estimated as the time lags of stimulation at which song frequency modulation (FM) or evoked neural responses deviated from baseline levels by more than 3 SDs. To estimate lateral magnocellular nucleus of the anterior nidopallium (LMAN) motor delays, we chronically implanted bipolar 50- $\mu\text{m}$  stainless steel wires separated by about 0.5 mm into LMAN and stimulated using monophasic 0.2-ms current pulses in the range of 100–500  $\mu\text{A}$ . We either used single pulses (Fig. S2A) or double pulses separated by 1 ms (Fig. 2B).

**LMAN Extracellular Recordings in Freely Moving Zebra Finches.** We identified LMAN in microdrive implantation surgeries by its antidromic response to electrical stimulation in the robust nucleus of the arcopallium (RA). In each recording session, we let birds first produce undirected songs with lights on, after which we exposed birds to playback of their own songs in the dark.

**Spectrograms, Raster Plots, and Spike-Triggered Averages.** Sounds were band-pass filtered in the range of 0.3–13 kHz and digitized at 32-kHz sampling rate. They were then Fourier transformed in 16-ms Hamming windows with 12-ms overlap and were displayed as log-power sound spectrograms in Figs. 2, and 3 and Figs. S2 and S4. We produced raster plots in Fig. 3 and Fig. S4 using the real-time display format corresponding to right alignment, i.e., we plotted a spike at time  $t$  in between spectrogram columns associated with Fourier windows  $[t_0 - 16 \text{ ms}, t_0]$  and  $[t_0 - 12 \text{ ms}, t_0 + 4 \text{ ms}]$ , where  $t_0$  is the largest integer multiple of 4 ms smaller than  $t$ . We plotted the spike triggered average (STA) log sound amplitude in Fig. S4 using centered alignment: the STA log amplitude at zero time lag corresponded to the log RMS sound waveform in the window  $[t_0 - 8 \text{ ms}, t_0 + 8 \text{ ms}]$  relative to a spike at time  $t_0$ .

**Cross-Covariance Analysis.** To compute motif cross-covariance (CC) functions, song-related (mean-subtracted) spike trains  $\rho_S(t)$  were taken from windows  $[0, T]$  starting 32 ms before onset of the stereotyped motif, to include hypothetical motor commands that initiate song motifs, and ending coincidentally with the last syllable of the motif  $C_{SP}(\tau) = \frac{1}{T} \int_0^T \rho_S(t) \rho_P(t + \tau) dt$ . Playback spike trains  $\rho_P(t)$  were taken from corresponding windows  $[\tau, T + \tau]$  of the same motif rendition.

Bout CC functions were computed analogously, with the interval  $[0, T]$  starting 32 ms before onset of the song bout and ending coincidentally with the last syllable of the bout.

To explore robustness of our covariance analysis, we varied the preonset interval of 32 ms from 4 to 100 ms: our conclusions were not affected by such changes. In one alternative analysis of motif covariance, we zero padded playback spike trains  $\rho_P(t)$  outside the boundary  $[0, T]$  of the stereotyped song motif, in which case  $C_{SP}(\tau)$  was normalized by  $T - |\tau|$  instead of  $T$  to yield an unbiased estimate of covariance. Results were not qualitatively different when using this definition of covariance function.

Peaks in both motif and bout CC functions were purely stimulus driven and therefore could in principle be observed at arbitrary large positive and negative time lags. Conceptually, negative dips in CC functions (instead of local maxima as in Fig. 1) can also be supportive of inverse models, i.e., models in which the feedback is

negative and not positive (neurons are inhibited by the sounds they produce and not excited).

## LMAN Extracellular Recordings in Anesthetized Zebra Finches

Experiments were performed in two head-fixed anesthetized adult male zebra finches (>90 dph). A few days (typically 1 wk) before the surgery, birds were caught from our colony and placed alone in an isolation chamber. Their song was recorded with a custom Labview (National Instruments) software. For playback, a song bout containing at least three motifs was selected for each bird. Birds were deeply anesthetized with isoflurane (1.5–3%) and received up to four urethane injections in intervals of at least 30 min. Injections consisted of 20–40  $\mu\text{L}$  of 20% urethane in saline injected i.m. or s.c. Birds were placed in a stereotaxic apparatus, and the head angle formed by the flat part of the skull above the beak was set at 65°. Windows were made into the skull near HVC, RA, LMAN, and Area X. A metal plate was fixed on the skull with dental acrylic. A stimulation electrode made of two coated stainless steel wires (diameter: 51  $\mu\text{m}$ ) connected to copper wires was implanted into RA [average coordinates from the bifurcation of the midsagittal sinus ( $\lambda$ ): 2.4 mm medio-lateral (ML); –2.8 mm anterior-posterior (AP); –1.7 mm dorso-ventral (DV) with a manipulator angle of –38°]. LMAN and HVC were located based both on stereotaxic coordinates and antidromic activation to RA electrical stimulation (average LMAN coordinates: 1.7 mm ML; 5.0 mm AP; 2.0–2.7 mm DV with a manipulator angle of 55°; average HVC coordinates: 2.0 mm ML; –0.5 mm AP; 0.3–0.6 mm DV with a manipulator angle of –45°). Area X was identified based on steady baseline neuronal activity (average coordinates: 1.9 mm ML; 6.0 mm AP; 2.5–3.5 mm DV with a manipulator angle of 45°). To perform extracellular electrophysiological recordings of single neurons in LMAN, we used glass pipettes filled in with a 3 M saline solution. Glass pipettes were pulled with a P-97 micropipette puller (Sutter Instruments) and had an impedance of 10–14 M $\Omega$ . The signal was amplified with an Axoclamp 2B, filtered (300 Hz to 10 kHz) with a LHBF-48X NPI filter, and recorded with a custom Labview software. We used tungsten metal electrodes with impedances of 1.5–4.0 M $\Omega$  (Microprobes) to perform extracellular recordings in Area X. The signal was amplified with a custom-made amplifier, filtered (300 Hz to 10 kHz) with a LHBF-48X NPI filter, and recorded with a custom Labview software. A glass pipette filled in with 250 M GABA in saline was implanted into HVC.

While simultaneously recording single units in LMAN and Area X, we performed the following steps. First, RA was stimulated to assess the latency of antidromic spike responses in LMAN; typical stimulation currents evoking responses in LMAN were ~60  $\mu\text{A}$ , and typical antidromic response latencies were ~2.0 ms. Second, we played back the bird's own song (BOS) for 5 min with an inter-BOS interval of 2 s. Third, while still recording the same LMAN and Area X units and playing back BOS, we inactivated HVC via pressure injection of GABA (using a Picospritzer). We injected ~0.1  $\mu\text{L}$  of GABA between each BOS playback during roughly 10 min. Fourth, we stopped the GABA injections but kept recording the same units while playing back BOS for another 10 min.

At the end of the experiment, we made burns at the position of the glass pipette by delivering a 10- $\mu\text{A}$  current for 10 s via a metal electrode (inserted into HVC after removal of the pipette). Similar electrolytic lesions were performed at the last recording sites in LMAN and Area X and at the RA stimulation site. The bird was then euthanized with a 0.05 mL i.m. injection of 20%

pentobarbital (50 mg/mL). Thereafter, the brain was removed and immersed into a 4% solution of paraformaldehyde. The brain was rinsed with a 0.01 M Phosphate Buffer Solution, the hemispheres were separated from each other, and the hemisphere of interest was glued on a metal plate and embedded in 2% agar. Sagittal slices of 80  $\mu\text{m}$  were cut with a Thermo Microm HM650V microtome and mounted on slides. Regions of interest were imaged and inspected for correct placement of electrolytic burns.

### Mirroring with Multiple Latencies

We can quantitatively work out a prediction for the mirroring CC function in a simple but conceptually instructive toy model, in the case in which a motor neuron's firing may elicit multiple auditory consequences at a range of latencies. Consider a single motor neuron, as in Fig. S1A and B, with time varying firing rate activity  $m(t)$ . Also, let  $a_i(t)$ , with  $i = 1, \dots, N$  being the time-varying firing rate activity of a population of  $N$  auditory neurons presynaptic to the motor neuron (Fig. S1B). Now suppose that each time the motor neuron is active at time  $t$ , it may generate multiple song features, which in turn generate subsequent auditory responses of different latencies in each auditory neuron. If we denote the impulse response (to the firing of the motor neuron) of each auditory neuron  $i$  by  $K_i(u)$  (which is nonzero only for  $u > 0$ ), then during any bout of motor exploration  $m(t)$ , the response of each auditory neuron, in a simple linear response model, will be

$$a_i(t) = \int_0^{\infty} du K_i(u) m(t-u). \quad [\text{S1}]$$

Of course, if there are multiple motor neurons firing, there will be additional contributions to the auditory feedback response. However, when explorations between different motor neurons are uncorrelated, the auditory feedback due to this one motor neuron is the only part of the auditory response that can correlate with the motor neuron's past activity and therefore is the only part of the response that determines the final weights from the presynaptic auditory population to the given motor neuron under a learning rule that has previously been reported to generate inverse models (1, 2).

This learning rule was derived as gradient descent on the following energy function on the vector of synaptic weights  $\mathbf{w}$

$$E(\mathbf{w}) = \left\langle \int_0^{\infty} ds e(s) [m(t-s) - \mathbf{w} \cdot \mathbf{a}(t)]^2 \right\rangle_m, \quad [\text{S2}]$$

where  $e(s)$  is an eligibility trace, and  $\langle \cdot \rangle_m$  denotes an average over the statistics of the motor explorations  $m(t)$ . This rule can be thought of as attempting to tune the inverse model synaptic weight vector  $\mathbf{w}$  so that the input to the motor neuron using auditory activity pattern  $\mathbf{a}(t)$  postdicts any past motor exploration  $m(t-s)$ , with the importance of postdiction at time lag  $s$  weighted by the eligibility trace  $e(s)$ . We assume the eligibility trace is normalized so that  $\int_0^{\infty} ds e(s) = 1$ . Gradient descent on this energy function yields a biologically plausible learning rule in which synapses from auditory to motor neurons undergo Hebbian association and hetero-synaptic competition.

Because the energy function is quadratic in the weights, with a unique minimum, it is straightforward to determine the outcome of learning:  $\mathbf{w} = \mathbf{C}^{-1}\mathbf{p}$  where

$$\mathbf{p}_i = \int_0^{\infty} ds e(s) \langle m(t-s) a_i(t) \rangle_m, \quad [\text{S3}]$$

is the average integrated, eligibility weighted correlation between auditory neuron  $i$  and the motor neuron and

$$\mathbf{C}_{ij} = \langle a_i(t) a_j(t) \rangle_m, \quad [\text{S4}]$$

is the equal time correlation between auditory neurons.

Now assume that motor exploration while inverse model synapses  $\mathbf{w}$  are being learned is temporally uncorrelated:  $\langle m(t_1) m(t_2) \rangle_m = \delta(t_1 - t_2)$ . Then using Eq. S1, we find

$$\langle m(t_1) a_i(t_2) \rangle_m = K_i(t_2 - t_1), \quad [\text{S5}]$$

and substituting this into Eq. S3 yields

$$\mathbf{p}_i = \int_0^{\infty} ds e(s) K_i(s). \quad [\text{S6}]$$

Also,  $\mathbf{C}_{ij}$  in Eq. S4 reduces to

$$\mathbf{C}_{ij} = \int_0^{\infty} du K_i(u) K_j(u). \quad [\text{S7}]$$

Now that we have the learned synaptic weight vector  $\mathbf{w}$ , we can compute the mirroring CC function. Consider for example a simple case where the motor neuron fires at an isolated time during motor production, as in Fig. S1A, Upper. Call this time  $t = 0$ . Thus, during motor production, motor activity is  $m_M(t) = \delta(t)$ . Now consider the motor response during playback of the song. Each auditory neuron  $i$  will have a firing rate response  $K_i(t)$ . This auditory activity will propagate to the motor neuron through the learned inverse model synapses  $\mathbf{w}$ , yielding the motor response to playback,  $m_P(t) = \sum_i \mathbf{w}_i K_i(t)$ . The mirroring CC function is then

$$CC(\tau) \equiv \int dt m_M(t) m_P(t+\tau) = \sum_i \mathbf{w}_i K_i(\tau), \quad [\text{S8}]$$

where again  $\mathbf{w} = \mathbf{C}^{-1}\mathbf{p}$ , and  $\mathbf{p}$  and  $\mathbf{C}$  are given in Eqs. S6 and S7, respectively. This equation yields a general prediction for the mirroring CC function directly in terms of the arbitrary, diverse latency, auditory responses  $K_i(t)$  and the eligibility trace  $e(t)$ .

To understand this prediction intuitively, imagine a situation where each auditory neuron  $i$  responds only during a narrow window of time, centered at  $\tau_i$ , as in Fig. S1A, Lower. If these windows of time do not overlap significantly, then the auditory correlation matrix  $\mathbf{C}$  in Eq. S4 is approximately diagonal. Furthermore, if the auditory responses have roughly equal strength, then  $\mathbf{C}$  is simply proportional to the identity matrix. Thus, the synaptic weights  $\mathbf{w}_i$  are simply proportional to the eligibility weighted auditory motor correlations  $\mathbf{p}_i$ . Now suppose the eligibility trace is a monotonically decaying exponential with time constant  $\tau_e$ :  $e(s) = \frac{1}{\tau_e} e^{-s/\tau_e}$ . Furthermore, suppose the width of the auditory response  $K_i(s)$  is narrow relative to  $\tau_e$ . Then  $\mathbf{p}_i$  in Eq. S6 is approximately proportional to  $e^{-\tau_i/\tau_e}$  for all  $i$ . Thus, the weight from an auditory neuron to the motor neuron is exponentially suppressed by the latency of the auditory neuron response relative to the motor neuron's firing, shown schematically in Fig. S1B. Therefore, in this scenario, the CC mirroring function is

$$CC(\tau) \approx \sum_i e^{-\tau_i/\tau_e} K_i(\tau). \quad [\text{S9}]$$

Therefore, the mirroring CC function is a weighted sum of the auditory responses, where long latency responses are exponentially suppressed by their latency. Thus, if a motor neuron has a complex playback response represented by a diversity of latencies, the mirroring CC function will tend to peak near the earliest









