

Multimodal system for recording individual-level behaviors in songbird groups

L. Rüttimann^{1*}, J. Rychen¹, T. Tomka^{1,2}, H. Hörster¹, M. D. Rocha¹, and R.H.R. Hahnloser^{1,2*}

¹Institute of Neuroinformatics, University of Zurich and ETH Zurich, 8057 Zurich, Switzerland.

²Neuroscience Center Zurich (ZNZ), University of Zurich and ETH Zurich, 8057 Zurich Switzerland.

*) corresponding authors: rlinus@ini.ethz.ch, rich@ini.ethz.ch

Abstract

In longitudinal observations of animal groups, the goal is to identify individuals and to reliably detect their interactive behaviors including their vocalizations. However, to reliably extract individual vocalizations from their mixtures and other environmental sounds remains a serious challenge. Promising approaches are multi-modal systems that make use of animal-borne wireless sensors and that exploit the inherent signal redundancy. In this vein, we designed a modular recording system (BirdPark) that yields synchronized data streams and contains a custom software-defined radio receiver. We record pairs of songbirds with multiple cameras and microphones and record their body vibrations with custom low-power frequency-modulated (FM) radio transmitters. Our custom multi-antenna radio demodulation technique increases the signal-to-noise ratio of the received radio signals by 6 dB and reduces the signal loss rate by a factor of 87 to only 0.03% of the recording time compared to standard single-antenna demodulation techniques. Nevertheless, neither a single vibration channel nor a single sound channel is sufficient by itself to signal the complete vocal output of an individual, with each sensor modality missing on average about 3.7% of vocalizations. Our work emphasizes the need for high-quality recording systems and for multi-modal analysis of social behavior.

Introduction

Acoustic communication is vital for many social behaviors. However, to study interactions among animals that are kept in groups entails many measurement challenges beyond the already considerable challenges of analyzing longitudinal data of isolated animals¹⁻³. One of the key difficulties of group-level behavior research is to perform automatic recognition of individuals and their actions.

Action recognition is the task of detecting behaviors from video or audio, or from signals collected with animal-borne sensors. Video-based action recognition has traditionally been based on posture tracking⁴⁻⁶ to avoid data-hungry training of classifiers on high-dimensional video data. Recently, posture tracking has greatly improved thanks to deep-learning approaches⁷⁻¹¹. Action recognition from video requires good visibility of focal animals because visual obstructions tend to hamper recognition accuracy. Given that freely moving animals may occlude one another, e.g., in birds during nesting, there seems to be a limitation to the usefulness of pure vision-based approaches.

Sound recordings have also been instrumental in action recognition. Sounds can be informative about both vocal and non-vocal behaviors. For example, wing flapping during flying produces a characteristic sound signature. But also preening, walking, and shaking can be recognized from sounds¹². The task of classifying sounds is known as acoustic scene classification¹³. Similar to vision, microphones record not just the focal animal but also background sounds. Therefore, action recognition and actor identification from sounds alone are challenging tasks, especially when many animals interact with one another. Possible workarounds to these issues are multi-modal approaches that combine multiple cameras and microphones, for example to assign vocalizations from dairy cattle in a barn to individual animals¹⁴. Other examples are systems that include also motion-tracking devices, for example to quantify gesture–speech synchrony in humans¹⁵.

In general, limitations due to sight occlusions and sound superpositions can be overcome with animal-borne sensors such as accelerometers^{16,17}, gyroscopes, microphones¹⁸, and global positioning systems (GPS)¹⁷. In combination with wireless transmitters¹⁸ and loggers¹⁶, these sensors enable the detection of behaviors such as walking, grooming, eating, drinking, and flying, for example, in birds¹⁹, cats²⁰, and dogs²¹, though often with low reliability because of noisy and ambiguous sensor signals¹². In general, animal-borne transmitter devices are designed to achieve high reliability, low weight, small size, and long battery life, giving rise to a complex tradeoff. Among the best transmitters, in terms of battery life, size, and weight, are analog frequency-modulated (FM) radio transmitters. Their low power requirement minimizes animal handling frequency and associated handling stress, making them an excellent choice for longitudinal observations of small vertebrates^{18,22,23}.

Among the challenges associated with FM radio reception are radio signal fadings due to relative movements of animal-borne transmitters and stationary receivers. Fading arises when electromagnetic waves arrive over multiple paths and interfere destructively (channel fading)²⁴ e.g., by reflection of metallic walls. Fading also occurs because every receiver has a direction of zero gain, which may affect signal reception from a moving transmitter. Signal fading can be addressed with antenna diversity, i.e., the use of several antennas. In multi-antenna approaches of diversity combining, either the strongest signal is selected, all signals are summed up, or signals are first weighted by their strength and then summed²⁵. However, these approaches do not guarantee protection from fading when signals annihilate. Alternatively, diversity combining is possible with phase compensation, which is the technique of shifting signal phases such that shifted signals align and sum constructively^{26,27}. Phase shifting reduces fading and increases the signal-to-noise of the received signal, and it provides cues for localizing a transmitter^{28,29}. We set out to bring the benefits of antenna diversity and of phase-shifting techniques to ethology research.

We focus on systems that make use of more than one sensor modality and recognize actions with higher accuracy than would be possible from one sensor modality alone. One particular challenge with multi-modal systems is the synchronization of the multimodal data streams. Usually, each data modality is recorded with a dedicated recording device that uses its own internal sampling clock. Furthermore, clocks tend to drift apart, and often the recordings cannot be started at exactly at the same time on all devices. Therefore, the individual data streams must be aligned post-recording using either markers in the sensor signals or auxiliary synchronization channels, which are labor intense and error-prone processes^{15,30,31}.

To perform individual-level longitudinal observations of social behaviors and to record high-quality multimodal data sets suitable for action recognition, we present a custom recording system for groups of vocalizing animals (BirdPark). We built a naturalistic environment inside a soundproof enclosure that features a set of microphones to record sounds and several video cameras to capture the entire scene. Moreover, all animals wear a miniature low-power transmitter device that transmits body movements from a firmly attached accelerometer via analog, frequency-modulated (FM) radio that we receive with several antennas.

Our system is optimized for robust longitudinal recordings of vocal interactions in songbirds. The on-animal accelerometers enable week-long monitoring of vocalizations without a change of battery. The combination of multiple antennas minimizes signal losses. All sensor signals are perfectly synchronized, which we achieve by routing dedicated sample triggers to all recording devices (radio receiver, stationary microphone digitizer, and cameras) derived from a central quartz clock using clock dividers. We release our custom recording software and we demonstrate the high data quality and redundant signaling of vocal gestures.

Results

The Recording Arena

We built an arena optimized for audiovisual recordings, minimizing acoustic resonances and visual occlusions (Figure 1A). It provides space for up to 8 songbirds and contains nest boxes, perches, sand baths, food, and water sources. To record the sounds inside the chamber, we installed five microphones. Three video cameras capture the overall scene from three orthogonal viewpoints. In addition, we installed a camera and a microphone in each of two nest boxes. To simplify video analysis, we combined the images from all five cameras into a single video image (Figure 1B). The camera resolutions are sufficient to resolve key points on birds even in midflight (Figure 1D).

To record vocalizations, we mounted transmitter devices to birds' backs. On each device there is an accelerometer that picks up body vibrations from vocalizations and body movements such as hopping and wing flapping¹⁶. The devices transmit the acceleration signals as frequency-modulated (FM) radio waves to four antennas inside the recording chamber.

We demodulated the FM radio signals with a custom eight-channel radio receiver. To ensure that the data streams from the video cameras, microphones, and transmitter devices are well synchronized, the clock of the radio receiver triggers the video frames and audio samples (Figure 1C): We divided the frequency of this 200 MHz clock by 2^{13} to generate the audio clock rate of 24.414 kHz, and with a further division by 2^9 , we generated the 47.684 Hz video frame rate.

On the host computer that controlled the acquisition system, we ran two custom applications: BirdVideo, that writes the video data to a file, and BirdRadio, that acquires the microphone and sensor data to a file (see Methods). The generated file pairs are synchronized such that with each video frame there are 512 audio and sensor samples. While the recording is gapless, it is split into files of typically 7 minutes duration.

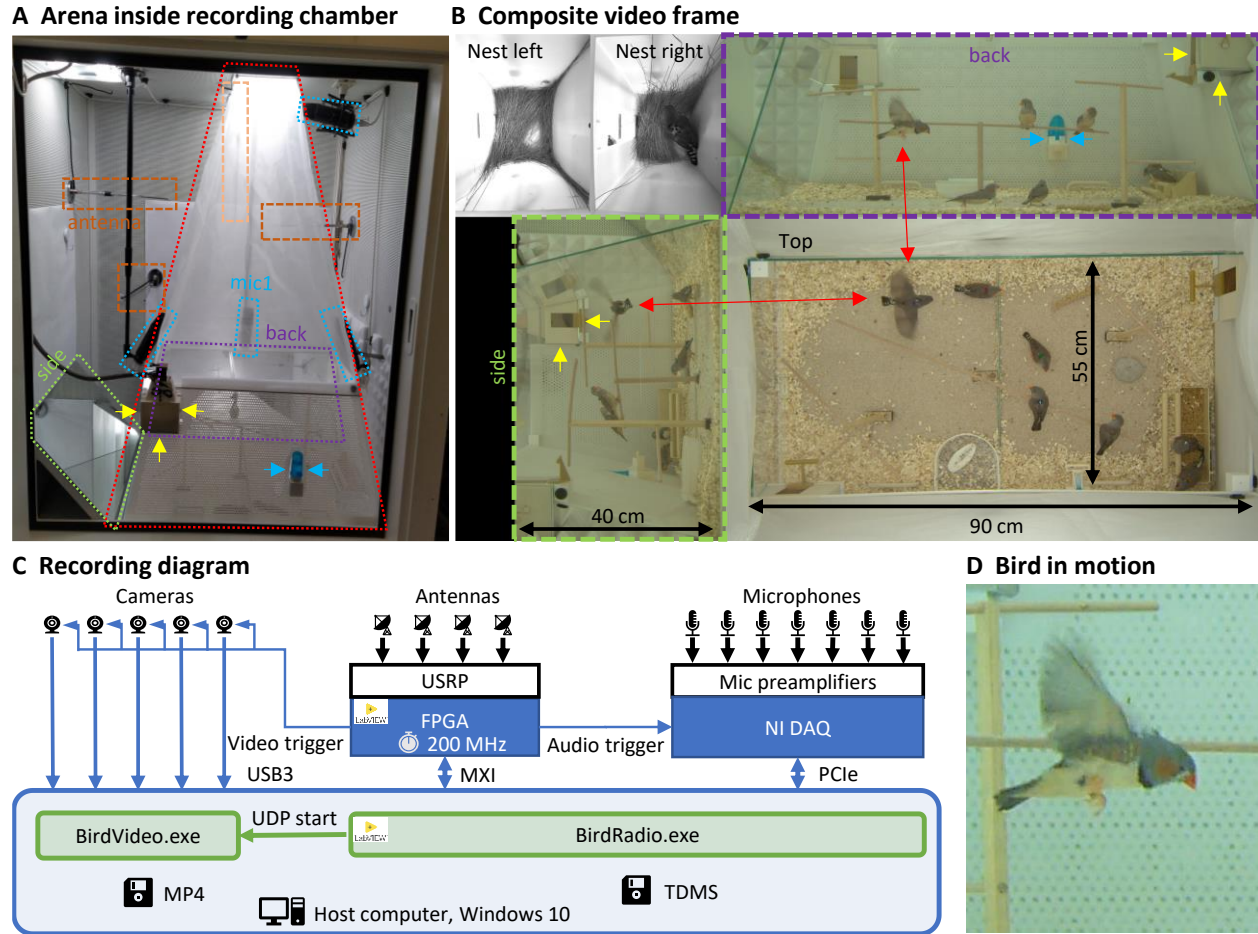


Figure 1: Recording arena and schematic of recording system. **A:** Inside a soundproof chamber, we built a recording arena (red dotted line) for up to 8 birds. We record the animals' behaviors with three cameras mounted through the ceiling. These provide a direct top view and indirect side and back views via two mirrors (delimited by green and magenta dotted lines). To record the sounds in the chamber, we installed five microphones (blue dotted lines) among all four sides of the cage (one attached to the front door is not visible) and the ceiling, and two small microphones in the nest boxes. The radio signals from the transmitter device are received with four radio antennas (orange dotted lines) mounted on three side walls and the ceiling. One nest box is indicated with yellow arrows and a water bottle with blue arrows. **B:** A composite still image of all camera views shows two monochrome nest box views (top left) and three views of the arena (top, side, back) with 8 birds among which one is flying (red arrows). Yellow and blue arrows as in A. **C:** Schematic of the recording system for gapless and synchronized recording of audio (microphones), radio (accelerometer sensors), and video channels (cameras). The radio receiver is implemented on a universal software radio peripheral (USRP) with a large field programmable gate array (FPGA) that runs at the main clock frequency of 200 MHz. Clock dividers on the FPGA provide the sample trigger for audio recordings and the frame trigger for the cameras. The data streams are collected on a host computer that runs two custom programs, one (BirdRadio) for streaming audio and sensor signals to disk and one (BirdVideo) for encoding video data. **D:** Zoom-in on an airborne bird, illustrating the spatial and temporal resolution of the camera.

Transmitter device

Our transmitter devices are based on the FM radio transmitter circuit described in Ter Maat et al.¹⁸ (Figure 2A). To distinctly record birds' vocalizations irrespective of external sounds, we replaced the microphone with an accelerometer¹⁶ that acts as a contact microphone. The radio circuit uses a single transistor to frequency modulate the sensor signal $a(t)$ onto a radio carrier frequency ω_c (that is set by the resonator circuit properties) with an inductor coil as an emitting antenna. As a result, the acceleration signal $a(t)$ is encoded as momentary radio transmitter frequency $\omega_T(t) \approx \omega_c(t) + ca(t)$

with c some constant. We found that the carrier frequency ω_c is not constant but modulated by the proximity of body parts (proximity effect, Figure 5). Additionally, it is subject to temperature and end-of-battery-life drift (see Methods: Transmitter device). While the instability of the carrier frequency ω_c is a disadvantage of the single-transistor circuit design, its advantages are its low weight (1.5 g) and small power consumption (the battery lifetime is 12 days).

Before mounting a device on a bird, we adjusted the coil to the desired carrier frequency in the range 250–350 MHz by slightly bending the wires. Thereafter, as a means of protection, we fixed the coil and the electronics in dyed epoxy. The purpose of the dyed epoxy is to help identify the birds in the video images and to stabilize the coil's inductance (Figure 2B,C). We mounted the devices on birds using a rubber-string harness adapted from Alarcón-Nieto et al.³² (Figure 2D).

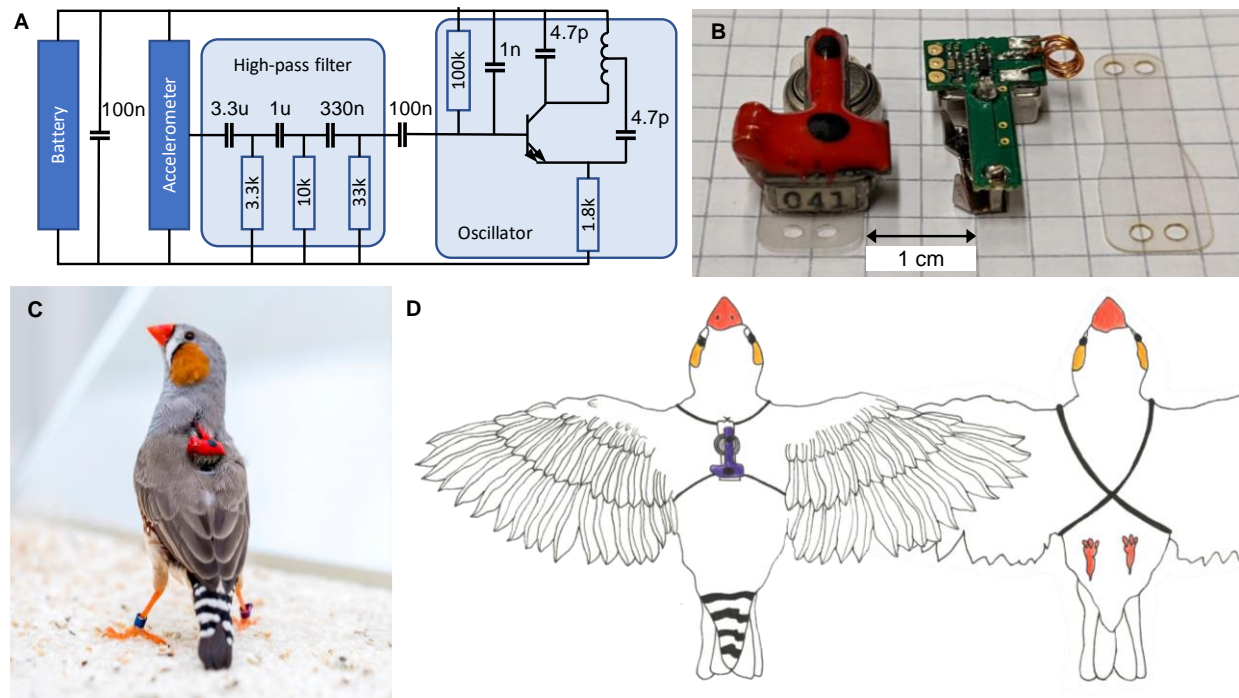


Figure 2: Transmitter device. **A:** Schematic of the electronic circuit (adapted from Ter Maat et al.¹⁸). The analog FM radio transmits the vibration transducer (accelerometer) signal via a high-pass filter (cutoff frequency: 15Hz) followed by a radiating oscillator. **B:** A fully assembled transmitter (left), and another one without epoxy and battery (middle), and a piece of mounting foil (right). **C:** Picture of the device mounted on a bird. The transmitters are color coded to help identify the birds in video data. **D:** Schematic of a bird wearing the device. The harness (adapted from Alarcón-Nieto et al.³²) is made of rubber string (black) that goes around the wings and the chest.

Radio receiver

To reconstruct the acceleration signal $a(t)$ of a given vibration transmitter device from the received multi-antenna FM signals, we demodulated the latter using a Phase Locked Loop (PLL), which measures the momentary transmitter frequency $\omega_T(t)$. A PLL generates an internal oscillatory signal of variable instantaneous frequency $\omega(t)$. It adjusts that frequency to maintain a zero relative phase with respect to the received (input) signal. As long as the zero relative phase condition is fulfilled, the PLL's instantaneous frequency $\omega(t)$, after high-pass filtering, forms our estimate of the bird's acceleration $a(t)$ (knowing that proximity effect degrade this relationship).

In our diversity combining approach, we construct the PLL's input signal as a combination of all four antenna signals. Namely, we compensate the individual phase offsets of the four antenna signals in such a way that all phases are aligned (we achieve this phase shifting with phase-compensation circuits — one for each antenna). We then form the desired mixture signal by summing the phase-shifted signals; the phase of the summed signal serves as the PLL's error signal that we use to adjust the instantaneous frequency $\omega(t)$. The variable phase offsets on the four antennas arise from the variable locations and orientations of transmitters relative to receivers.

In summary, our approach to minimizing fading in wireless ethology research is to use a demodulator comprising a PLL and a diversity combining approach based on phase-compensation. Our FM radio demodulation is described in more detail in the following.

We implemented our custom demodulation technique by installing four antennas labeled A, B, C, and D perpendicular to the walls and the ceiling of the chamber (see Methods: Radio reception). We fed the antenna signals into a universal software radio peripheral (USRP) containing a large field programmable gate array (FPGA) (Figure 3A). The four input stages of the USRP filter out an 80 MHz wide band around the local oscillator frequency ω_{LO} , which we typically set to 300 MHz. These four signals then become digitally available on the FPGA as complex valued signal of 100 MHz sampling rate. We call them the intermediate signals $z^a(t)$, $a \in \{A, B, C, D\}$ (intermediate band in the frequency domain, Methods: Intermediate band, Figure 3B).

On the FPGA, we instantiated eight demodulators indexed by $i \in \{1, \dots, 8\}$. Each demodulator contained four digital downconverters (DDC) that cut out from the four intermediate signals $z^a(t)$ the four baseband signals $u_i^a(t)$, $a \in \{A, B, C, D\}$ around the eight instantaneous frequencies $\omega_i(t)$, sampled at 781.25 kHz (Methods: Baseband, Figure 3C and D). All baseband signals and the main signal are complex valued and so we interchangeably call them vectors and their phases we refer to as angles.

The PLL is driven by the main vector $u_i^M(t) = u_i^A(t) + u_i^B(t) + u_i^C(t) + u_i^D(t)$ that we form as the sum of the four baseband vectors. The error signal for PLL's feedback controller is given by the phase of the main vector $\theta_i(t) = \arg(u_i^M(t))$ (Methods: PLL, Figure 3D). When that phase is approximately zero (PLL is locked), the instantaneous frequency $\omega_i(t)$ tracks the frequency of transmitter i . We updated the PLL's instantaneous frequency at a rate of 781.25 kHz.

To avoid destructive interferences in the summation of baseband vectors, we compensated their phases $\alpha_i^a(t) = \angle(u_i^M(t), u_i^a(t))$ relative to the main vector. We introduced individual phase offset $\Delta\varphi_i^a(t)$ that were set in feedback loops to drive the phases $\alpha_i^a(t)$ towards zero (Methods: Phase compensation, Figure 3E). Phase offsets were updated at a rate of 1.5 kHz and logged at every video frame.

The PLL and the phase compensation form independent control loops. When the PLL is unlocked (e.g., off), the instantaneous frequency does not match the momentary transmitter frequency and the baseband vectors rotate at the difference frequency (Figure 3F left). When the PLL is switched on and locked, the baseband vectors do not rotate, and the main signal displays a phase $\theta_i(t) \simeq 0$ (Figure 3F middle). When the phase alignment is switched on, the baseband signals align and their sum maximizes the magnitude of the main vector (Figure 3F right).

Because birds' locomotion is slower than their rapid vocalization-induced vibratory signals, the phases $\alpha_i^a(t)$ change more slowly than the instantaneous frequency $\omega_i(t)$ and therefore we updated phase offsets less often than the PLL's internal frequency.

Operation

The intermediate band of $z^a(\omega)$ is wide enough to accommodate up to eight FM transmitters. Provided the transmitters' FM carrier frequencies are roughly evenly spread, the momentary transmitter frequencies do not cross, even during very large frequency excursions such as when a bird pecks on its sensor. Nevertheless, we limited the instantaneous frequencies ω_i to the range $[\Omega_i - \Delta\Omega, \Omega_i + \Delta\Omega]$, where Ω_i is the center frequency and $\Delta\Omega = 1$ MHz is the common limiting range of all channels, Figure 3B. The center frequency Ω_i of a channel we manually set at the beginning of an experiment to the associated FM carrier frequency. The limiting range of 1 MHz we found narrow enough such that the PLL rapidly re-locked after brief signal losses.

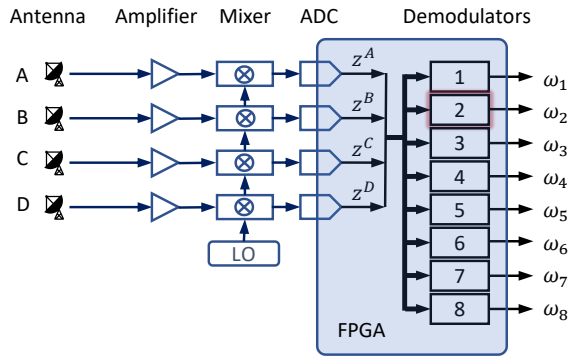
During operation, we often observed large excursions of the instantaneous frequency. These excursions occurred while birds pecked on the device or while they preened their feathers near the sensor, or when one bird sat on top of another, such as during copulations. The magnitude of these jumps could reach 1 MHz, which is much larger than the maximally 1-kHz shifts induced by vocalizations, Figure 5. The large excursions likely were caused by capacitive-inductive effects of movements on the resonator circuit. These observations demonstrate the large bandwidth and robustness of our PLLs.

Performance of diversity combining technique

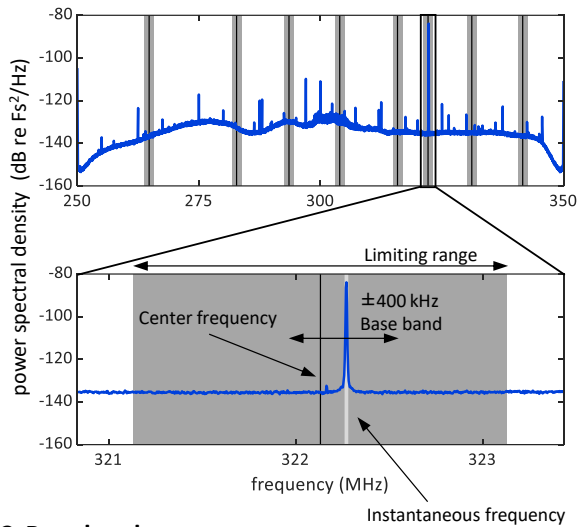
We validated the robustness of our multi-antenna demodulation method by analyzing the signal-to-noise ratio of the received transmitter signal in a zebra finch pair over a full day (2 x 14 h measurement period). We quantified the diversity combining performance in terms of the fraction of times at which signal fading occurred, both with and without diversity combining. We defined the radio signal-to-noise ratio (RSNR) of $u_i^M(t)$ as the signal power (mean square) divided by the median of the signal variance (the median is taken over the duration of 7-min long data files). In theory, the noise (variance) of $u_i^M(t)$ is twice as large as that of either baseband signals $u_i^a(t)$ ($a=A..D$), assuming independence of radio amplifier noises. The signal power of $u_i^M(t)$ is upper bounded by 16 times that of its largest constituent, i.e., the best-placed antenna signal (the base-band signal $u_i^a(t)$ of largest amplitude).

Indeed, we found that the signal-to-noise ratio of $u_i^M(t)$ was about 6 dB higher on average than that of the best single-antenna signal, Figure 4b. The total reception time during which the RSNR was critically low (<13 dB, our operational definition of signal fading) was 87 times shorter when using the multi-antenna demodulation than when using the best-antenna demodulation (Fig. 4c). Thus, our diversity combining technique is useful and results in a coherent constructive addition of the four antenna signals and to very robust frequency tracking even when the signal in one antenna vanishes.

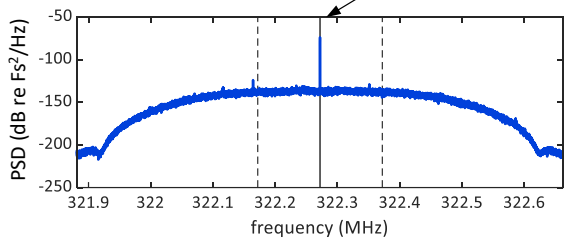
A Radio receiver



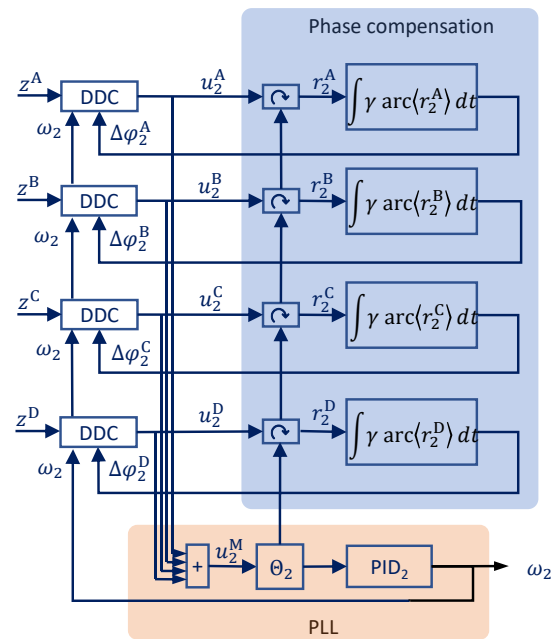
B Intermediate band



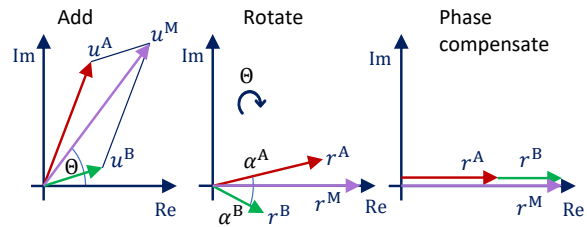
C Base band



D Demodulator: PLL & phase compensation



E Vector diagram



F Scatter plot

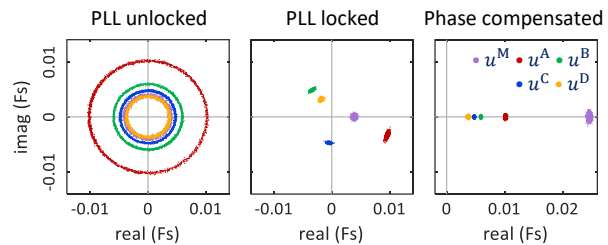
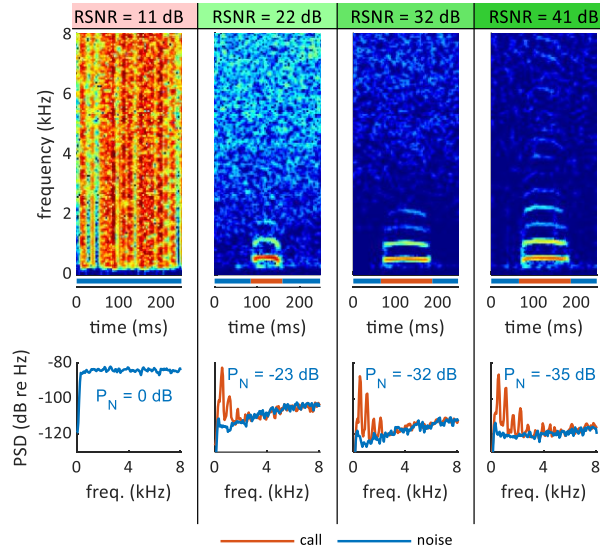


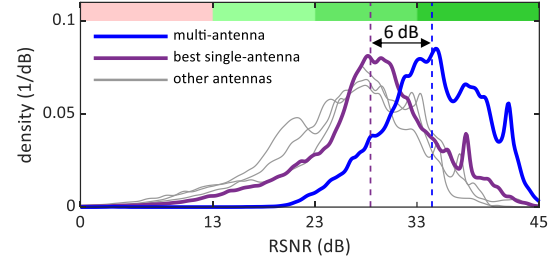
Figure 3: Radio receiver with PLL demodulators and phase compensation. **A:** Schematic of the software-defined radio receiver. Analog antenna signals are first down converted by mixing the amplified signals with a local oscillator of frequency ω_{LO} , followed by low-pass filtering. The resulting radio frontend (RF) outputs have digital representations $z^A(t)$, $z^B(t)$, $z^C(t)$, and $z^D(t)$. From these, eight FM demodulators extract the instantaneous frequencies $\omega_i(t)$. **B:** The power spectrum of the 80 MHz wide intermediate band $z^A(\omega)$ centered on the oscillator frequency $\omega_{LO} = 300$ MHz. The limiting ranges (gray bars) of the 8 channels (transmitters) are typically set to ± 1 MHz of their manually set center frequencies (black vertical lines). A zoom-in (lower graph) to the limiting range of one active channel reveals the large peak associated with the momentary radio transmitter frequency. **C:** The baseband power spectral density $u_i^A(\omega)$ is a down-converted, ± 100 kHz (dashed vertical lines) frequency flat band around the instantaneous frequency (solid vertical line) that tracks the momentary radio transmitter frequency (the large peak). **D:** For each channel (transmitter), a demodulator using a PLL (orange shading) and four phase compensation circuits (blue shading) computes the instantaneous frequency ω_i (shown for $i = 2$). The baseband signals $u_i^A(t)$ are derived from intermediate signals with digital downconverters (DDCs) that operate on a common instantaneous frequency ω_2 . The phase compensation circuits drive the phases $\alpha_i^A(t) = \angle(u_i^M(t), u_i^A(t)) = \arg(r_i^A(t))$, with the rotated baseband

signals $r_i^a = u_i^a e^{-i\theta_i}$, to zero. **E:** Vector diagram in the complex plane illustrating the effects of alignment by the PLL (middle) and of phase compensation (right) on main and baseband vectors. **F:** The measured phases of the four baseband vectors u_i^a are shown without phase compensation (left), after aligning the main vector with the PLL (middle), and after additional phase compensation (right). The result is that all vectors are aligned, and the master vector is of maximal amplitude.

A Radio signal-to-noise ratio (RSNR) examples



B RSNR histogram



C RSNR statistics

	RSNR ≤ 13 dB	13 dB < RSNR ≤ 23 dB	23 dB < RSNR ≤ 33 dB	33 dB < RSNR
Multi-antenna	0.03 %	1.31 %	36.19 %	62.46 %
Best single-antenna	2.62 %	14.73 %	59.87 %	22.79 %

Figure 4: **A:** Spectrogram examples of demodulated acceleration signals generated by single zebra finch calls with radio signal-to-noise ratios of 11, 22, 32, and 41 dB (top, from left to right). The bottom plots show the spectra of the calls (red lines) and of noise (blue lines), computed as time averages from the spectrograms above (time windows indicated by red and blue horizontal bars). The relative noise power P_N (integral of blue line relative to the noise power of the first example) is decreasing with increasing RSNR. When the RSNR is below 13 dB, the noise power spectral density is above the signal power of most vocalizations (these become invisible). **B:** Histogram of RSNR over time for the multi-antenna signal (blue line) and the single-antenna signals for the antenna with the greatest mean RSNR (purple line) and all other antennas (grey lines). The mean of the multi-antenna RSNR (dashed blue line) is 6 dB larger than the mean RSNR of the best single antenna (dashed purple line). **C:** Multi-antenna demodulation is significantly better than best-antenna demodulation, as demonstrated by the significantly longer reception periods above a given signal-to-noise ratio. The time below the critical RSNR of 13 dB is reduced by a factor of 87.

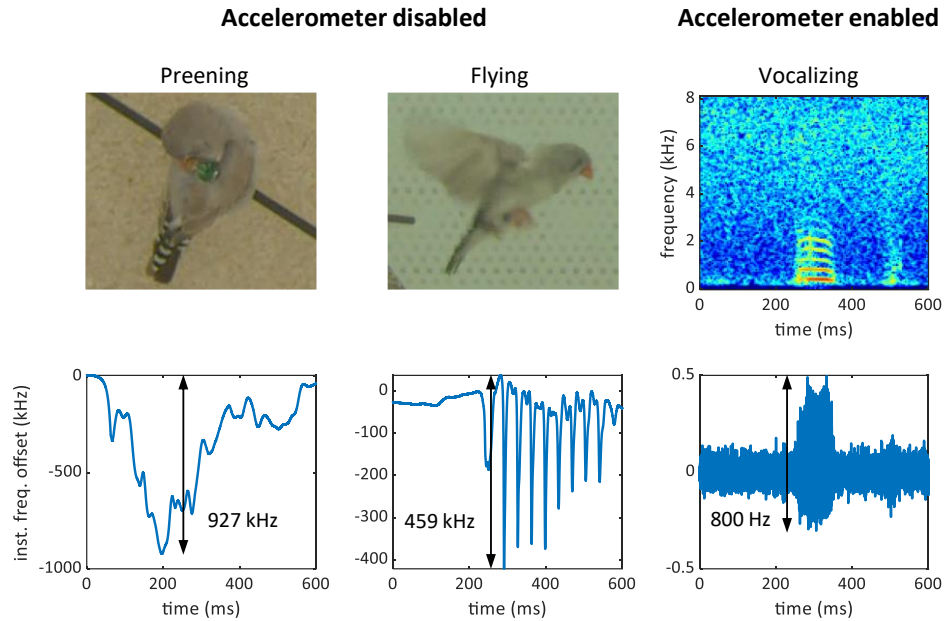
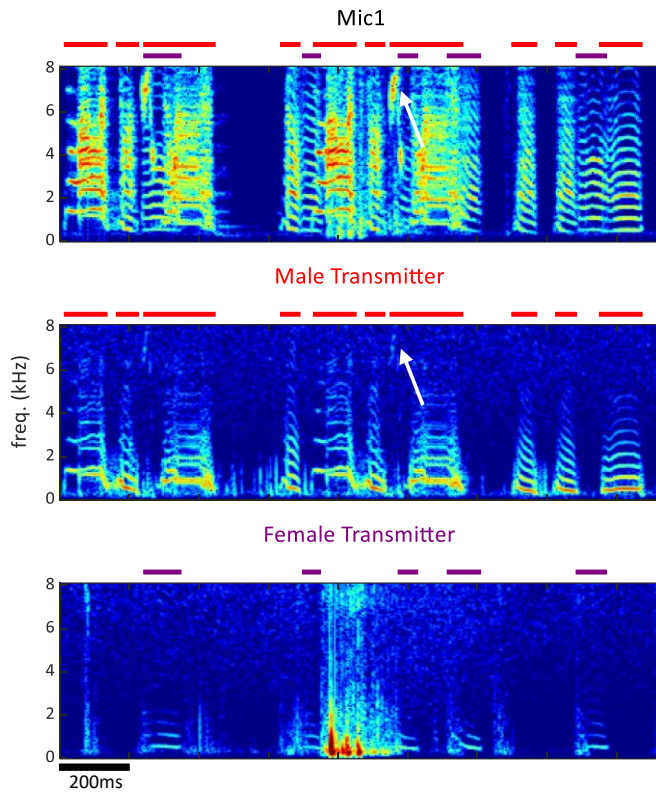


Figure 5: Proximity affects the instantaneous frequency much more than vocalizations. On a transmitter device with a disabled (short circuited) accelerometer, the instantaneous frequency is strongly modulated by the proximity of the head during preening (left) and by wing movements during flight (middle). In contrast, the modulation of instantaneous frequency is much weaker for vocalizations (right) on a transmitter with enabled accelerometer. The modulation amplitudes (black arrows) due to proximity are about 1000 times larger than the modulation amplitudes due to vocalizations.

A Example spectrograms



B Misses

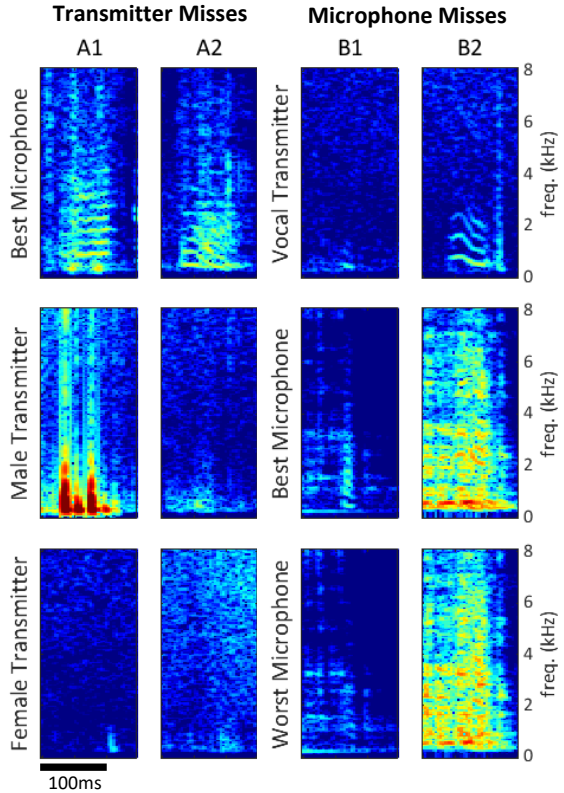


Figure 6: Vocal signals and their segmentation. **A:** Example log-power spectrograms of vocalizations produced by a mixed-sex zebra finch pair. The songs and calls overlap on the microphone channel (Mic1, top) but appear well separated on spectrograms of the male's momentary frequency $\omega_1(t)$ (male transmitter, middle) and the female momentary frequency $\omega_2(t)$ (female transmitter, bottom). Distinct transmitter-based vocal segments are indicated by red (male) and purple (female) horizontal bars on top of the spectrograms. High-frequency vibrations appear attenuated, but even a high-pitched 7 kHz song note (white arrow) by the male is still visible. **B:** Example vocalizations that are not captured by the transmitter device (A1: syllable masked by wing flaps, A2: low signal level) or not visible on either some channels (B1: low signal level) or on all microphone channels (B2: syllable masked by loud noise).

	Experiment 1		Experiment 2		Total
	male	female	male	female	
# vocalizations	447	91	487	83	1108
# voc. missed on tr. channels	26 (5.82 %)	5 (5.49 %)	6 (1.23 %)	4 (4.82 %)	41 (3.70 %)
# voc. missed on all mic. channels	1 (0.22 %)	0 (0.00 %)	3 (0.62 %)	0 (0.00 %)	4 (0.36 %)
# voc. missed on Mic1 channel	11 (2.46 %)	9 (10.00 %)	23 (4.72 %)	3 (3.61 %)	46 (4.16 %)
vocal overlaps	2.73 %		3.51 %		3.15 %

Table 1: Statistics of missed vocalizations on a transmitter channel (second row) and on all microphone channel (third row) and on an individual microphone channel (Mic1, fourth row). The 'syllable overlaps' row quantifies the percentage of times in which both birds vocalize simultaneously (as in the example in Figure 6A) expressed as a fraction of total vocalization duration.

Analysis of vocalizations

We observed that signals picked up from vocalizations reached up to 7 kHz in optimal cases, but only up to 1 kHz in the worst case. This variability stems from differences in vocalization amplitude, RSNR levels, and presumably, also from how well the sensor is in contact with the skin of the bird and at which position. The vibration sensor also measures accelerations from body movements like wing flaps or hopping. To quantify our ability to detect vocalizations in the multimodal data stream, we manually segmented vocalizations by inspection of log-power sound spectrograms from 2 vibration and 6 microphone channels (see Methods for details). We were particularly interested in missed vocalizations that are invisible on a given channel but whose presence is revealed on another channel. Compared to the miss rate associated with a single microphone channel of 4.2% (range 2.5%-10.0%, 7 mins recordings of n=2 bird pairs), all microphone channels combined produced a reduction in miss rate by a factor of about 10, which reveals a large benefit of the multi-microphone approach. The percentage of missed vocalizations on a given vibration channel was similar to that on a single microphone channel, 3.7% (range 1.2%-5.8%, n=4 birds), although these numbers are not directly comparable because the miss rate on a microphone channel stems from two birds and the miss rate on a vibration channel from a single bird. A very small number of vocalizations was missed on all six microphone channels, 0.4% (range 0.0%-0.6%, n=2 bird pairs), these were all very faint calls that were masked by loud noises.

We release our open-source radio receiver software and our graphical user interface (GUI) for monitoring life data streams on a host computer. Also, we share a sample of our dataset together with the annotations of all vocal segments, these may be useful for benchmarking machine learning systems on the task of vocal segmentation.

Discussion

We designed and validated a behavioral recording system for up to 8 songbirds, yielding perfectly synchronized multi-modal data. Our custom multi-antenna FM demodulation technique, compared with single antenna demodulation increases radio SNRs by 6 dB and reduces signal-loss events by a factor of 87. The wireless devices transmit well-separated vocalizations unless these are masked by large body movements such as wing flaps or because of weak mechanical coupling between the bird and the device, which we were unable to completely get rid of. We segmented vocalizations in two 7-mins recordings of pair-housed birds and found that each signal channel by itself is of limited value for reliably extracting vocal segments. We believe this ability to evaluate the inherent quality of a sensor modality can be very valuable in ethology research.

Ideally, the quality of a data set of vocalizations should be evaluated relative to ground truth, i.e., the true segmentation of vocalizations and background noise. However, such a ground truth data set does not (yet) exist in practice. For example, to perfectly measure the vocal output of a bird would require simultaneous measurements of syringeal labia and muscle activity, of sub-syringeal air pressure, and of tracheal air flow³³, which are measurements that have not been performed simultaneously in freely moving group-housed birds. Without such a data set or at least an approximation thereof, quantifications of vocal output will be biased.

We estimated vocal output by visual inspection of sound spectrograms from five microphones and two animal-borne accelerometers. With respect to this approximate ground truth, we find that no sensor

channel by itself achieves a rate of misses (false negative vocalizations) of less than 1%, with an average of 3.7%, which is a large number given that we took our measurements in ideal settings using the minimal group size of two instrumentalized birds housed in a relatively small, acoustically well isolated environment. Since we might have missed some vocalizations even considering all the channels used, our estimated single-channel miss rates constitute lower bounds, implying that the true miss rate of vocalizations must be higher. And, assuming that published datasets are of worse quality than ours (less instrumentalization), it is likely that findings on animal vocal communication in the literature exhibit a bias due to a miss rate higher than 3.7%. Moreover, since we expect the miss rate to increase with the number of interacting birds studied, the benefits of our recording system likely increase with the size of the social group.

Our system could promote research on the meaning of vocalizations for social behavior. Of key interest are courtship and reproductive behaviors, which have attracted much attention in the past. For songbirds, selecting a partner for copulation and subsequent offspring rearing involves complex courtship displays that include varied vocalizations and coordinated body movements on very rapid time scales. The reproductive behavior of zebra finches has been studied thoroughly^{34–37}, with efforts to define comprehensive ethograms^{38,39}. The copulation behavior of zebra finches has attracted much attention, partly because of the phenomenon of extra-pair copulation: the tendency of zebra finches to copulate with mates that are not their partner^{40–42}. However, not much is known about the roles of vocalizations in signaling copulation readiness and in reflecting the subsequent pair bond.

Similarly, multimodal group-level studies using our system could help us better understand the learning strategies young birds use while they modify their immature vocalizations to match a sensory target provided by a tutor. Much of our knowledge on song learning stems from research of isolated animal in controlled environments^{43–45}. However, the isolated-animal paradigm is impoverished compared to the natural setting in which animals live in groups, because vocal learning is subject to social influences⁴⁶. For example, the song learning success in juveniles is positively influenced by social interactions^{47,48} including by non-singing adult females⁴⁹. Furthermore, the directed songs of males produced towards females are different and subserved by different brain mechanisms than the undirected songs produced while alone⁵⁰. To study these processes, it would be valuable to acquire longitudinal high-quality audio and video data sets of freely-behaving and vocalizing animals in complex social settings. As vocalizations are communicative signals, the information they contain can only be fully captured by considering both the environment and the social group of the studied individual.

Methods

The chamber

All our experiments were performed in a sound-isolation chamber (Industrial Acoustic Company, UK) of dimensions 124 cm (width) x 90 cm (depth) x 130 cm (height). The isolation chamber has a special silent ventilation system for air exchange at a rate of roughly 500 l/min. Two types of light sources made of light emitting diodes (LEDs, Waveform Lighting) are mounted on aluminum plates on the ceiling of the chamber: 1) LEDs with natural frequency spectrum (consuming a total electric power of 80 W); and 2) ultraviolet (UV) LEDs of wavelength 365 nm and consuming a total electric power of 13 W.

The ceiling of the chamber contains three circular holes of 7 cm diameter through which we inserted aluminum tubes of 5 mm wall thickness that serve as holders for three cameras (Basler acA2040-120uc).

The tubes also conduct the heat of the LEDs (93 W) and the cameras (13 W) to the outside of the chamber where we placed silent fans to keep the cameras below 55 deg C.

With the cameras, we filmed the arena directly from the top and indirectly from two sides via two tilted mirrors in front of the glass side panels. This setup yields an object distance of about 1 – 1.5 m, which allows for a depth of field that covers the entire arena. Furthermore, the large object distance results in a small perspective distortion.

We installed four microphones on the four side walls of the isolation chamber and one microphone on the ceiling. We mounted two further microphones in the nest boxes and one microphone outside the isolation chamber. Wherever possible, we covered the sidewalls of the chamber with sound-absorbing foam panels. A door sensor measures whether the door of the chamber is open or closed. The temperature inside the chamber we kept at roughly 26 °C and the humidity was roughly 24 %.

The arena

Inside the chamber, we placed the bird arena of dimensions 90 cm x 55 cm (floor). To minimize acoustic resonances and optical reflections, the 40-cm high side panels of the arena are tilted by 11° and 13° towards the inside. Two of the side panels are made of glass for videography and the two opposite panels are made of perforated plastic plates. The floor of the arena is covered with a sound-absorbing carpet. A pyramidal tent made of a fine plastic net covers the upper part of the arena reaching up to 125 cm above ground. At a height of 35 cm, we attached two nest boxes to the side panels, each equipped with a microphone and a camera. Furthermore, the arena is equipped with perches, a sand-bath, and a food and water tray.

Video acquisition system

To visualize the arena, we used industrial cameras (3 Megapixel, MP) with zoom lenses (opening angles: top view: 45° x 26°, back view: 55° x 26°, side view: 26° x 35°) and exposure times of 3 ms. To visualize nests in the dimly lit nest boxes, we used monochrome infrared cameras (2 MP, Basler daA1600-60um) and fisheye lenses (143° x 112°).

The uncompressed camera outputs (ca. 400 MB/s in total) are relayed to a host computer over a USB3 Vision interface. Each camera receives frame trigger signals generated by the USRP (Figure 1C). A custom program (BirdVideo), written in C++ and using the OpenCV library⁵¹ and FFMPEG⁵², undistorts the nest camera images with a fisheye lens model and transforms all five camera images to a single 2976 x 2016 pixel-sized image (Figure 1B). The composite images are then encoded with the h264 codec on a NVIDIA GPU and stored in an MP4 (ISO/IEC 14496-14) file. We used constant-quality compression with a variable compression ratio in the range 150-370, resulting in a data rate of 0.72 MB/s – 1.8 MB/s, depending on the number of birds and their activity in the arena. Compression ratios in this range do not significantly decrease key point tracking performance⁵³. The frame rate of the video is about 48 frames per second (frame period 21 ms). The spatial resolution of the main cameras is about 2.2 pixels/mm.

Transmitter device

The vibration transducer (Knowles BU21771) senses acceleration with a sensitivity of 1.4 mV/m/s² in the 30 Hz – 6 kHz frequency range⁵⁴. Inspired by Ter Maat et al.¹⁸, we performed the frequency modulation of the vibration signal onto a radio carrier using a simple Hartley oscillator transmitter stage with only

one transistor and a coil that functions both as the inductor in the LC resonator and as the antenna. The LC resonator determines the carrier frequency ω_c of the radio signal. We set ω_c of our transmitters in the vicinity of 300 MHz, which corresponds to an electromagnetic wavelength of about 1 m and is roughly equal to the dimension of our soundproof chambers, implying that our radio system operates in the near field.

We found that ω_c depends on temperature at an average of -73 kHz/ $^{\circ}$ C (range -67 kHz/ $^{\circ}$ C to -87 kHz/ $^{\circ}$ C, $n=3$ transmitters). Towards the end of the battery life (over the course of the last 3 days), we observed an increase of ω_c by an average of $500 - 800$ kHz. These slow drifts can easily be accounted for by tracking and high-pass filtering the momentary frequency. The measured end-to-end sensitivity of the frequency modulation is 5 kHz/g, with g being the gravitational acceleration constant. The transmitters are powered by a zinc-air battery (type LR41), which lasts at least 12 days. The total weight of the transmitter including the harness and battery is 1.5 g, which is ca. 10% of the body weight of a zebra finch.

We moderately tightened the harness during attachment in a tradeoff between picking up vibratory signals from the singer and preserving the wearer's natural behaviors. It is known that the act of mounting devices can transiently affect zebra finch behavior: right after attachment, the singing rate and amount of locomotion both tend to decrease, they return to baseline within 3 days for sensors weighing 0.75 g⁵⁵ and within 2 weeks for sensors weighing 3 g⁵⁶.

Radio reception

We mounted four whip antennas of 30 cm length (Albrecht 6157) perpendicular to the metallic sidewalls and ceiling of the chamber (using magnetic feet). We fed the antenna signals to a universal software radio interface (USRP-2940, National Instruments, USA), which comprises four independent antenna amplifiers the gains of which we set to 68 dB, adjustable in the range -20 dB to $+90$ dB.

Intermediate band

The USRP generates from the amplified radio signal a down-converted signal with an analog bandwidth of 80 MHz around the local oscillator frequency (ω_{LO}). After digitization, this intermediate signal $z^a(t)$ for antenna $a \in \{A, B, C, D\}$, is a 100 MS/s complex-valued signal (IQ) of 2×18 bits precision and is relayed to the FPGA. For details about the analog down conversion and the sampling of complex valued signals, see the NI manual of the USRP and the documentation of our custom radio software.

Digital downconverters (DDCs) and baseband

Every PLL, $i \in \{1, \dots, 8\}$, makes use of four digital downconverters (DDCs, Figure 3D), one for each antenna a , that cut from the intermediate band a narrow baseband around the instantaneous frequency ω_i . The 3-stage decimation filter of this down-conversion has a flat frequency response within a ± 100 kHz band. Given the decimation factor of 128, the resulting complex base-band signals $u_i^a(t)$ (the outputs of the DDCs) have a sample rate of 781.25 kHz and a precision of 2×25 bits.

Phase locked loop

Each PLL generates its instantaneous frequency ω_i by direct digital synthesis (DDS) with a 48-bit phase accumulation register and a lookup table. This frequency is dynamically adjusted to keep the phase of the main vector (i.e., the summed baseband signals) close to zero. We calculated the angle $\theta_i = \arg(u_i^M)$ of the main vector with the CORDIC algorithm⁵⁷ and unwrapped the phase up to ± 128 turns

before using it as the error signal for a proportional-integral-derivative (PID) controller. The PID controller adjusts the instantaneous frequency of the PLL.

The PID controller was implemented on the FPGA in single precision floating point (32-bit) arithmetic. The controller included a limiting range and an anti-windup mechanism⁵⁸. The unwrapping of the error phase was crucial for the PLL to quickly lock-in and to keep the lock even during large and fast frequency deviations of the transmitter. To tune the PID parameters, we measured the closed-loop transfer function of the PLL by adding a white-noise signal to the control signal (instantaneous frequency), and then adjusted the PID parameters until we observed a low-pass characteristic without overshoot. We achieved a closed loop bandwidth of about 30 kHz.

Phase compensation

For a given PLL i , we compensated the relative phases under which a radio signal arrives at the four antennas a in order to align the four baseband vectors u_i^a . The alignment was achieved by providing a phase offset $\Delta\varphi_i^a$ to each downconverter, where it acts as offset to the phase accumulation register of the direct digital synthesis. To compute the angle of baseband vector relative to the main vector, we rotated the baseband vectors u_i^a by the phase θ_i of the main vector to result in the rotated vector $r_i^a = u_i^a e^{-i\theta_i}$. After averaging the rotated vectors across 512 samples, we computed its angle $\alpha_i^a = \arg(\langle r_i^a \rangle_{512})$. We then compensated that angle by iteratively adding a fraction $\gamma \in [0,1]$ of it to the phase offset: $\Delta\varphi_i^a := \Delta\varphi_i^a + \gamma\alpha_i^a$. The parameter γ is the phase compensation gain (Figure 3D), typically set to $\gamma = 0.2$. This iterative update was performed at a rate of about 1.5 kHz (781.25 kHz / 512), which is faster than birds' locomotion (changes in physical position or orientation of the transmitter).

Central control software

On the host we run our Central control software (BirdRadio) programmed with LabView, that acquires the microphone and transmitter signals and writes them to a TDMS file (Figure 1C). Furthermore, BirdRadio sends UDP control signals to BirdVideo, it automatically starts and stops the recording in the morning and evening respectively, it controls the light in the sound-isolation chamber with dimming that simulates sunrise and sunset, it triggers an email alarm when the radio signal from a transmitter device is lost and it automatically adjusts the center frequency of each radio channel every morning to adjust for carrier frequency drift.

Data management

The BirdPark is designed for continuous recordings over multiple months, producing data at a rate of 60 GB/day for 2 birds and 130 GB/day for 8 birds. We implemented the FAIR (findable, accessible, interoperable, and reusable) principles of scientific data management⁵⁹ as follows:

Our recording software splits the data into gapless files of 20'000 video frames (ca. 7 mins duration). At the end of a recording day, all files are processed by a script that converts the TDMS files into HDF5 files and augments them with rich metadata. The HDF5 files are self-descriptive in that every data field (main data and metadata) contains its description as a property. We use the lossless compression feature of HDF5 to obtain a compression ratio of typically 2.5 for the audio and accelerometer data. The script also adds two AAC-compressed microphone channels into the video files. Although this step introduces redundancy, the availability of sound in the video files is very useful during manual annotation of the videos. Furthermore, the script also exports the metadata as a JSON file and copies the processed data

onto a NAS server. At the end of an experiment, the metadata is uploaded onto an openBIS⁶⁰ server and is linked with the datafiles on the NAS.

Manual Segmentation of vocalizations

We manually segmented vocalizations in recordings of mixed-sex zebra finch pairs (7 mins recordings of n=2 bird pairs). We high-pass filtered (FIR filter with a stopband frequency of 100 Hz and passband frequency of 200 Hz) the raw accelerometer and microphone signals and produced spectrograms (window size=384 samples, hop size=96 samples) that we manually segmented using Raven Pro 1.6⁶¹. We performed three types of segmentations in the following order:

1. Transmitter-based vocal segments: We separately segmented all vocalizations on either transmitter channel, precisely annotating their onsets and offsets. In our released data sets, these segments are referred to as transmitter-based vocal segments. When we were uncertain whether a sound segment was a vocalization or not, we also looked at spectrograms of microphone channels and listened to sound playbacks. If we remained uncertain, the segment was tagged with the label 'Unsure', such segments were treated as (non-vocal) noises and were excluded from further analysis.
2. Microphone-based vocal segments: We simultaneously visualized all microphone spectrograms (from Mic1 to Mic6) using the multi-channel view in Raven Pro (we ignored Mic7, which was located in the second nest that was not accessible to the birds). On those, we annotated each vocal segment on the first microphone channel on which it was visible (e.g., a syllable that is visible on all microphone channels is only annotated on Mic1). Overlapping vocalizations were annotated as a single vocal segment. When we were uncertain whether a sound segment was a vocalization or not, we also looked at spectrograms of accelerometer channels and listened to sound playbacks. If we remained uncertain, the segment was tagged with the label 'Unsure', such segments were treated as (non-vocal) noises and were excluded from further analysis. These segments are referred to as microphone-based vocal segments.
3. Consolidated vocal segments: All consistent (perfectly overlapping) accelerometer- and microphone-based vocal segments, we labelled as consolidated vocal segments. We then inspected all inconsistent (not perfectly overlapping) segments by visualizing all channel spectrograms. We fixed inconsistencies that were caused by human annotation errors (e.g. lack of attention) by fixing the erroneous or missing transmitter- and microphone-based segments. From the inconsistent (partially overlapping) segments that were not caused by human error, we generated one or several consolidated segments by trusting the modality that more clearly revealed the presence of a vocalization (hence our reporting of 'misses' in Table 1).

In our released comma-separated value (CSV) files, we give each consolidated vocal segment a Bird Tag (e.g., either 'b15p5_m' or 'b14p4_f') that identifies the bird that produced the vocalization, a Transmitter Tag that identifies the transmitter channel on which the vocalization was identified (either 'b15p5_m' or 'b14p4_f' or 'None'), and a FirstMic Tag that identifies the first microphone channel on which the segment was visible ('Mic1' to 'Mic6', or 'None'). We resolved inconsistencies and chose these tags as follows:

- If a microphone (-based vocal) segment was paired (partially overlapping) with exactly one transmitter segment, a consolidated segment was generated with the onset time set to the minimum onset time and the offset time set to the maximum offset time of the segment pair. The Bird and Transmitter Tags were set to the transmitter channel name, and the FirstMic Tag was set to the microphone channel name.

- If a transmitter segment was unpaired, a consolidated segment was created with the same onset and offset times. The Bird and TrCh Tags were set to the transmitter channel name, and the FirstMic Tag was set to 'None'.
- If a microphone segment was unpaired, a transmitter channel name was inferred based on the vocal repertoire and noise levels on both transmitter channels. We visually verified that the vocal segment was not the result of multiple overlapping vocalizations (which was never the case). Then we created a consolidated segment with the same on- and offset, set the FirstMic Tag to the microphone Id, the Bird Tag to the guessed transmitter channel name, and the Transmitter Tag to 'None'.
- If a microphone segment was paired with more than one transmitter segment, a consolidated segment was created for each of the accelerometer segments. The onsets and offsets were manually set based on visual inspection of all spectrograms. Bird and Transmitter Tags were set to the transmitter channel name and the FirstMic Tag was set to the microphone channel name.
- We never encountered the case where a transmitter segment was paired with multiple microphone segments.

The statistics in Table 1 were calculated as follows:

vocalizations = (number of consolidated segments)

voc. missed on tr. channels = (number of consolidated segments with Transmitter=None)

voc. missed on all mic. channels = (number of consolidated segments with FirstMic=None)

voc. missed on Mic1 = (number of consolidated segments with FirstMic≠Mic1 and FirstMic≠None)

The assignment to the female/male column was determined by the 'Bird' tag.

The vocal overlap statistic was calculated based on the number of spectrogram bins in the consolidated segmentation as follows:

vocal overlap = $2 * (\text{number of spectrogram bins with vocal overlap}) / (\text{number of spectrogram bins with vocalization})$

Animals and Experiments

Mixed-sex zebra finch (*Taeniopygia guttata*) pairs bred and raised in our colony (University of Zurich and ETH Zurich) were kept in the BirdPark on a 14/10 h light/dark daily cycle for multiple days, with food and water provided ad libitum. All experimental procedures were approved by the Cantonal Veterinary Office of the Canton of Zurich, Switzerland (license numbers ZH054/2020). All methods were carried out in accordance with relevant guidelines and regulations (Swiss Animal Welfare Act and Ordinance, TSchG, TSchV, TVV).

Data and code availability

The recording software (BirdRadio and BirdVideo) and the segment annotations including the raw data files (HDF5 files with transmitter and microphone signals and MP4 videos) will be made available for download.

Author contributions

Conceptualization, J.R., L.R., M.D.R. and R.H.R.H.; System and software development, L.R. and J.R.; animal experiments, L.R. and H.H.; data annotation and data analysis, T.T. and L.R.; writing—original draft, J.R., L.R. and R.H.R.H.; writing—reviewing and editing, J.R., L.R., T.T., M.D.R. and R.H.R.H.

Acknowledgments

We thank Aymeric Nager for his help with the design and construction of the BirdPark. Financial support was provided by Swiss national Science Foundation 31003A_182638 (The roles of vocal communication in pair formation and cultural learning in songbirds) and Swiss national Science Foundation NCCR Evolving Language (Agreement #51NF40_180888).

Competing interests

The author(s) declare no competing interests.

References

- [1] D. Lipkind, A. T. Zai, A. Hanuschkin, G. F. Marcus, O. Tchernichovski, and R. H. R. Hahnloser, “Songbirds work around computational complexity by learning song vocabulary independently of sequence,” *Nat Commun*, vol. 8, no. 1, pp. 1–11, Nov. 2017, doi: 10.1038/s41467-017-01436-0.
- [2] O. Tchernichovski, P. P. Mitra, T. Lints, and F. Nottebohm, “Dynamics of the vocal imitation process: How a zebra finch learns its song,” *Science (1979)*, vol. 291, no. 5513, pp. 2564–2569, Mar. 2001, doi: 10.1126/science.1058522.
- [3] S. Kollmorgen, R. H. R. Hahnloser, and V. Mante, “Nearest neighbours reveal fast and slow components of motor learning,” *Nature*, vol. 577, no. 7791, pp. 526–530, Jan. 2020, doi: 10.1038/s41586-019-1892-x.
- [4] C. Segalin *et al.*, “The mouse action recognition system (MARS) software pipeline for automated analysis of social behaviors in mice,” *Elife*, vol. 10, Jul. 2021, doi: 10.7554/eLife.63720.
- [5] S. Fujimori, T. Ishikawa, and H. Watanabe, “Animal behavior classification using DeepLabCut,” in *2020 IEEE 9th Global Conference on Consumer Electronics, GCCE 2020*, Oct. 2020, pp. 254–257. doi: 10.1109/GCCE50665.2020.9291715.
- [6] A. Perkes, B. Pfrommer, K. Daniilidis, D. J. White, and M. Schmidt, “Variation in female songbird state determines signal strength needed to evoke copulation,” *bioRxiv*, p. 2021.05.19.444794, May 2021, doi: 10.1101/2021.05.19.444794.
- [7] T. Nath, A. Mathis, A. C. Chen, A. Patel, M. Bethge, and M. W. Mathis, “Using DeepLabCut for 3D markerless pose estimation across species and behaviors,” *Nat Protoc*, vol. 14, no. 7, pp. 2152–2176, Jul. 2019, doi: 10.1038/s41596-019-0176-0.
- [8] T. D. Pereira *et al.*, “SLEAP: Multi-animal pose tracking,” *BioRxiv*, Sep. 2020, doi: 10.1101/2020.08.31.276246.

- [9] A. Mathis, S. Schneider, J. Lauer, and M. W. Mathis, “A Primer on Motion Capture with Deep Learning: Principles, Pitfalls, and Perspectives,” *Neuron*, vol. 108, no. 1. pp. 44–65, Oct. 2020. doi: 10.1016/j.neuron.2020.09.017.
- [10] T. Walter and I. D. Couzin, “Trex, a fast multi-animal tracking system with markerless identification, and 2d estimation of posture and visual elds,” *Elife*, vol. 10, pp. 1–73, Feb. 2021, doi: 10.7554/eLife.64000.
- [11] M. Badger *et al.*, “3D Bird Reconstruction: A Dataset, Model, and Shape Recovery from a Single View,” in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2020, vol. 12363 LNCS, pp. 1–17. doi: 10.1007/978-3-030-58523-5_1.
- [12] D. Stowell, E. Benetos, and L. F. Gill, “On-Bird Sound Recordings: Automatic Acoustic Recognition of Activities and Contexts,” *IEEE/ACM Trans Audio Speech Lang Process*, vol. 25, no. 6, pp. 1193–1206, Jun. 2017, doi: 10.1109/TASLP.2017.2690565.
- [13] D. Barchiesi, D. D. Giannoulis, D. Stowell, and M. D. Plumbley, “Acoustic Scene Classification: Classifying environments from the sounds they produce,” *IEEE Signal Process Mag*, vol. 32, no. 3, pp. 16–34, May 2015, doi: 10.1109/MSP.2014.2326181.
- [14] G. H. Meen, M. A. Schellekens, M. H. M. Slegers, N. L. G. Leenders, E. van Erp-van der Kooij, and L. P. J. J. Noldus, “Sound analysis in dairy cattle vocalisation as a potential welfare monitor,” *Comput Electron Agric*, vol. 118, pp. 111–115, Oct. 2015, doi: 10.1016/j.compag.2015.08.028.
- [15] W. Pouw, J. P. Trujillo, and J. A. Dixon, “The quantification of gesture–speech synchrony: A tutorial and validation of multimodal data acquisition using device-based and video-based motion tracking,” *Behav Res Methods*, vol. 52, no. 2, pp. 723–740, Apr. 2020, doi: 10.3758/s13428-019-01271-9.
- [16] V. N. Anisimov, J. A. Herbst, A. N. Abramchuk, A. v Latanov, R. H. R. Hahnloser, and A. L. Vyssotski, “Reconstruction of vocal interactions in a group of small songbirds,” *Nat Methods*, vol. 11, no. 11, pp. 1135–1137, Nov. 2014, doi: 10.1038/nmeth.3114.
- [17] E. Eisenring *et al.*, “Quantifying song behavior in a free-living, light-weight, mobile bird using accelerometers,” *Ecol Evol*, vol. 12, no. 1, p. e8446, Jan. 2022, doi: 10.1002/ece3.8446.
- [18] A. ter Maat, L. Trost, H. Sagunsky, S. Seltmann, and M. Gahr, “Zebra finch mates use their forebrain song system in unlearned call communication,” *PLoS One*, vol. 9, no. 10, p. e109334, Oct. 2014, doi: 10.1371/journal.pone.0109334.
- [19] A. G. Laich, R. P. Wilson, F. Quintana, and E. L. C. Shepard, “Identification of imperial cormorant *Phalacrocorax atriceps* behaviour using accelerometers,” *Endanger Species Res*, vol. 10, no. 1, pp. 29–37, May 2010, doi: 10.3354/esr00091.
- [20] S. Watanabe, M. Izawa, A. Kato, Y. Ropert-Coudert, and Y. Naito, “A new technique for monitoring the detailed behaviour of terrestrial animals: A case study with the domestic cat,” *Appl Anim Behav Sci*, vol. 94, no. 1–2, pp. 117–131, Oct. 2005, doi: 10.1016/j.applanim.2005.01.010.

- [21] L. Gerencsér, G. Vásárhelyi, M. Nagy, T. Vicsek, and A. Miklósi, "Identification of Behaviour in Freely Moving Dogs (*Canis familiaris*) Using Inertial Sensors," *PLoS One*, vol. 8, no. 10, p. e77814, Oct. 2013, doi: 10.1371/journal.pone.0077814.
- [22] L. F. Gill, P. B. D'Amelio, N. M. Adreani, H. Sagunsky, M. C. Gahr, and A. ter Maat, "A minimum-impact, flexible tool to study vocal communication of small animals with precise individual-level resolution," *Methods Ecol Evol*, vol. 7, no. 11, pp. 1349–1358, Nov. 2016, doi: 10.1111/2041-210X.12610.
- [23] L. F. Gill, W. Goymann, A. ter Maat, and M. Gahr, "Patterns of call communication between group-housed zebra finches change during the breeding cycle," *Elife*, vol. 4, no. OCTOBER2015, Oct. 2015, doi: 10.7554/eLife.07770.
- [24] D. Tse and V. Pramod, *Fundamentals of Wireless Communication*. Cambridge University Press, 2005. [Online]. Available: <https://web.archive.org/web/20070810052329/http://www.eecs.berkeley.edu/~dtse/book.html>
- [25] R. S. Shatara, "Combined Switched and Phase Aligned Multi-Antenna Diversity System for Signal-Error-Reduction in Mobile Receiving Systems," Universität der Bundeswehr München, 2003.
- [26] O. Voitsun, S. Senega, and S. Lindenmeier, "Multi-Antenna Diversity Set for Transmission and Reception in Car-to-Car and Car-to-X Communication," in *GeMIC 2020 - Proceedings of the 2020 German Microwave Conference, 2020*, pp. 80–83.
- [27] S. Senega, A. Nassar, and S. Lindenmeier, "Automotive antenna diversity system for satellite radio with high phase accuracy in low SNR-scenarios," *Clay Miner*, vol. 10, no. 5–6, pp. 578–586, Jun. 2018, doi: 10.1017/S1759078718000296.
- [28] A. Haniz *et al.*, "A Novel Phase-Difference Fingerprinting Technique for Localization of Unknown Emitters," *IEEE Trans Veh Technol*, vol. 66, no. 9, pp. 8445–8457, 2017, doi: 10.1109/TVT.2017.2696049.
- [29] C. A. Berdanier and Z. Wu, "A novel RF emitter localization method through phase information," *IEEE National Radar Conference - Proceedings*, pp. 1–5, 2013, doi: 10.1109/RADAR.2013.6586137.
- [30] P. H. Zimmerman, J. E. Bolhuis, A. Willemsen, E. S. Meyer, and L. P. J. J. Noldus, "The observer XT: A tool for the integration and synchronization of multimodal signals," *Behav Res Methods*, vol. 41, no. 3, pp. 731–735, Aug. 2009, doi: 10.3758/BRM.41.3.731.
- [31] T. C. Dolmans, M. Poel, J. Klooster, and B. P. Veldkamp, *Data Synchronisation and Processing in Multimodal Research*, vol. 1, no. October. 2021.
- [32] G. Alarcón-Nieto, J. M. Graving, J. A. Klarevas-Irby, A. A. Maldonado-Chaparro, I. Mueller, and D. R. Farine, "An automated barcode tracking system for behavioural studies in birds," *Methods Ecol Evol*, vol. 9, no. 6, pp. 1536–1547, Jun. 2018, doi: 10.1111/2041-210X.13005.
- [33] D. Morris, "The Reproductive Behaviour of the Zebra Finch (*Poephila Guttata*), With Special Reference To Pseudofemale Behaviour and Displacement Activities," *Behaviour*, vol. 6, no. 1, pp. 271–322, Jan. 2008, doi: 10.1163/156853954x00130.

- [34] D. MORRIS, "THE COMPARATIVE ETHOLOGY OF GRASSFINCHES (ERYTHRURAE) AND MANNIKINS (AMADINAE)," *Proceedings of the Zoological Society of London*, vol. 131, no. 3, pp. 389–439, Nov. 1958, doi: 10.1111/j.1096-3642.1958.tb00695.x.
- [35] P. G. Caryl, "Sexual behaviour in the zebra finch *Taeniopygia guttata*: response to familiar and novel partners," *Anim Behav*, vol. 24, no. 1, pp. 93–107, Feb. 1976, doi: 10.1016/S0003-3472(76)80103-0.
- [36] R. Ullrich, P. Norton, and C. Scharff, "Waltzing *Taeniopygia*: Integration of courtship song and dance in the domesticated Australian zebra finch," *Anim Behav*, vol. 112, pp. 285–300, Feb. 2016, doi: 10.1016/j.anbehav.2015.11.012.
- [37] A. J. Figueredo, D. M. Ross, and L. Petrinovich, "The Quantitative Ethology of the Zebra Finch: A Study in Comparative Psychometrics," *Multivariate Behav Res*, vol. 27, no. 3, pp. 435–458, Jul. 1992, doi: 10.1207/s15327906mbr2703_7.
- [38] J. Hau, H. Jacobs, N. Smith, P. Smith, L. Smyth, and P. Yew, "Zebra finch behaviour in standard cages and effect of simple enrichment," *Scandinavian Journal of Laboratory Animal Science*, vol. 23, no. SUPPL. 1, pp. 129–134, 1996.
- [39] T. R. Birkhead, J. Pellatt, and F. M. Hunter, "Extra-pair copulation and sperm competition in the zebra finch," *Nature*, vol. 334, no. 6177, pp. 60–62, 1988, doi: 10.1038/334060a0.
- [40] T. R. Birkhead, F. M. Hunter, and J. E. Pellatt, "Sperm competition in the zebra finch, *Taeniopygia guttata*," *Anim Behav*, vol. 38, no. 6, pp. 935–950, Dec. 1989, doi: 10.1016/S0003-3472(89)80135-6.
- [41] W. Forstmeier, K. Martin, E. Bolund, H. Schielzeth, and B. Kempenaers, "Female extrapair mating behavior can evolve via indirect selection on males," *Proc Natl Acad Sci U S A*, vol. 108, no. 26, pp. 10608–10613, Jun. 2011, doi: 10.1073/pnas.1103195108.
- [42] S. Kollmorgen, R. H. R. Hahnloser, and V. Mante, "Nearest neighbours reveal fast and slow components of motor learning," *Nature 2020 577:7791*, vol. 577, no. 7791, pp. 526–530, Jan. 2020, doi: 10.1038/s41586-019-1892-x.
- [43] D. Lipkind, A. T. Zai, A. Hanuschkin, G. F. Marcus, O. Tchernichovski, and R. H. R. Hahnloser, "Songbirds work around computational complexity by learning song vocabulary independently of sequence," *Nature Communications 2017 8:1*, vol. 8, no. 1, pp. 1–11, Nov. 2017, doi: 10.1038/s41467-017-01436-0.
- [44] O. Tchernichovski, P. P. Mitra, T. Lints, and F. Nottebohm, "Dynamics of the vocal imitation process: How a zebra finch learns its song," *Science (1979)*, vol. 291, no. 5513, pp. 2564–2569, Mar. 2001, doi: 10.1126/SCIENCE.1058522/ASSET/7C0AA79A-E81E-4A48-9284-A22A9A1ACDE1/ASSETS/GRAPHIC/SE1119293005.JPEG.
- [45] I. Ljubičić, J. Hyland Bruno, and O. Tchernichovski, "Social influences on song learning," *Curr Opin Behav Sci*, vol. 7, pp. 101–107, Feb. 2016, doi: 10.1016/J.COBEHA.2015.12.006.

- [46] Y. Chen, L. E. Matheson, and J. T. Sakata, "Mechanisms underlying the social enhancement of vocal learning in songbirds," *Proc Natl Acad Sci U S A*, vol. 113, no. 24, pp. 6641–6646, Jun. 2016, doi: 10.1073/pnas.1522306113.
- [47] S. Carouso-Peck and M. H. Goldstein, "Female Social Feedback Reveals Non-imitative Mechanisms of Vocal Learning in Zebra Finches," *Current Biology*, vol. 29, no. 4, pp. 631-636.e3, Feb. 2019, doi: 10.1016/j.cub.2018.12.026.
- [48] D. Y. Takahashi, D. A. Liao, and A. A. Ghazanfar, "Vocal Learning via Social Reinforcement by Infant Marmoset Monkeys," *Current Biology*, vol. 27, no. 12, pp. 1844-1852.e6, Jun. 2017, doi: 10.1016/j.cub.2017.05.004.
- [49] S. C. Woolley and A. J. Doupe, "Social context-induced song variation affects female behavior and gene expression," *PLoS Biol*, vol. 6, no. 3, pp. 0525–0537, Mar. 2008, doi: 10.1371/journal.pbio.0060062.
- [50] Gary Bradski, "The OpenCV Library," *Dr. Dobb's Journal of Software Tools*, 2000.
- [51] S. Tomar, "Converting Video Formats with FFmpeg | Linux Journal," *Linux Journal*, vol. 2006, no. 146, p. 10, 2006, [Online]. Available: <https://www.linuxjournal.com/article/8517>
- [52] A. Mathis and R. Warren, "On the inference speed and video-compression robustness of DeepLabCut," *bioRxiv*, p. 457242, Oct. 2018, doi: 10.1101/457242.
- [53] Knowles Electronics, "TB-26: The use of BU series accelerometers." 2017.
- [54] L. F. Gill, P. B. D'Amelio, N. M. Adreani, H. Sagunsky, M. C. Gahr, and A. ter Maat, "A minimum-impact, flexible tool to study vocal communication of small animals with precise individual-level resolution," *Methods Ecol Evol*, vol. 7, no. 11, pp. 1349–1358, Nov. 2016, doi: 10.1111/2041-210X.12610.
- [55] V. N. Anisimov, J. A. Herbst, A. N. Abramchuk, A. v Latanov, R. H. R. Hahnloser, and A. L. Vyssotski, "Reconstruction of vocal interactions in a group of small songbirds," *Nat Methods*, vol. 11, no. 11, pp. 1135–1137, Nov. 2014, doi: 10.1038/nmeth.3114.
- [56] J. E. Volder, "The CORDIC Trigonometric Computing Technique," *IRE Transactions on Electronic Computers*, vol. EC-8, no. 3, pp. 330–334, 1959, doi: 10.1109/TEC.1959.5222693.
- [57] K. J. Astrom and L. Rundqwist, "Integrator windup and how to avoid it," 1989, pp. 1693–1698. doi: 10.23919/acc.1989.4790464.
- [58] M. D. Wilkinson *et al.*, "Comment: The FAIR Guiding Principles for scientific data management and stewardship," *Sci Data*, vol. 3, no. 1, pp. 1–9, Mar. 2016, doi: 10.1038/sdata.2016.18.
- [59] A. Bauch *et al.*, "OpenBIS: A flexible framework for managing and analyzing complex data in biology research," *BMC Bioinformatics*, vol. 12, no. 1, pp. 1–19, Dec. 2011, doi: 10.1186/1471-2105-12-468.
- [60] Bioacoustics Research Program and B. R. Program, "Raven Pro: Interactive Sound Analysis Software (Version 1.5)," *Ithaca, NY: The Cornell Lab of Ornithology*. The Cornell Lab of Ornithology, Ithaca, NY, p. <http://www.birds.cornell.edu/raven>, 2014. [Online]. Available:

papers3://publication/uuid/DE0F17E7-8210-4F82-907B-
B923E8F9AE7A%5Cnhttp://www.birds.cornell.edu/raven