

1 Amplitude Increases of Vocalizations are Associated with Body Accelerations in Siamang
2 (*Symphalangus syndactylus*)
3
4

5 Wim Pouw¹, Mounia Kehy², Marco Gamba³, Andrea Ravignani^{4, 5, 6}

6 1. Donders Institute for Brain, Cognition, and Behaviour, Radboud University Nijmegen,
7 Netherlands

8 2. Equipe de Neuro-Ethologie Sensorielle, Université Jean Monnet, France

9 3. Dipartimento di Scienze della Vita e Biologia dei Sistemi, Università di Torino, Torino, Italy

10 4. Comparative Bioacoustics Research Group, Max Planck Institute for Psycholinguistics,
11 Nijmegen, The Netherlands

12 5. Center for Music in the Brain, Department of Clinical Medicine, Aarhus University & The
13 Royal Academy of Music, Aarhus, Denmark

14 6. Department of Human Neurosciences, Sapienza University of Rome, Rome, Italy
15

16 **Author note:** Correspondences can be addressed to Wim Pouw
17 (wim.pouw@donders.ru.nl). We would like to thank the ‘Jaderpark Tier- und Freizeitpark an der
18 Nordsee’ for allowing us to record audiovisual data of the Siamang family at their facility. This
19 work has been supported by the Max Planck Institute (MPI) for Psycholinguistics Nijmegen, and
20 the Donders Institute for Cognition, Brain, and Behavior. We would like to thank Jeroen Geerts
21 of the MPI, for support of organizing the audiovisual equipment, and Maarten Snellen with his
22 help setting up a dedicated server to run our dashboard application. We would like to thank
23 Diandra Düngen for her support in making the recordings possible. WP is funded by a VENI
24 grant (VI.Veni 0.201G.047: PI Wim Pouw) and was further supported by a Donders Postdoctoral
25 Development fund. The Comparative Bioacoustics Group is supported by Max Planck
26 Independent Research Group Leader funding to A.R. Center for Music in the Brain is funded by
27 the Danish National Research Foundation (DNRF117). AR is funded by the European Union
28 (ERC, TOHR, 101041885).

29 **Open data:** All data and code supporting this manuscript can be found on

30 https://github.com/WimPouw/siamang_physical_constraints_code_repo

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

Abstract

Siamangs (*Symphalangus syndactylus*), one of the few singing apes, vocalize loudly, often while they move. We hypothesize that movement and vocalization coordinate, possibly due to vigorous thorax-loading movements such as brachiation affecting vocal-respiratory dynamics. To assess this vocal-motor coordination we recorded more than a hundred stereotypical vocalizations combined with movement from two captive Siamang (isolated from 7 hours of singing). We observed that stereotypical calls coincided with a movement display and were performed by juvenile individuals during solo singing (which allowed for isolation of the calls). Investigating these vocal-motor events, we found that body acceleration estimated using computer vision was statistically associated with the nearest peak in the amplitude envelope of the call, and that body acceleration timeseries contained mutual information about the amplitude envelope timeseries during these events. By confirming via quantitative methods that singing and movement are coordinated, the current report invites further mechanistic investigation on vocal-locomotor coupling in siamang.

Keywords: Siamang, Locomotor-vocal coupling, Locomotion, Respiration, Vocalization, Multimodal Communication

1 Introduction

2 Human and non-human animals often coordinate whole-body movement with
3 vocalizations. Multiple bat species (e.g., *Phyllostomus hastatus*) flexibly synchronize their echo-
4 locating pulses with their wingbeats while flying, often in 1:1 or polyrhythmic fashion [1–4].
5 Twelve species of North- and South-American birds show an allometrically coupled wingbeat
6 duration with vocal unit durations [5]. Gerbils (*Meriones unguiculatus*) locomote in a saltating
7 way, when they hop and hit the ground with their forelimbs they synchronously emit a
8 vocalization [6]. The male brown-headed cowbird (*Molothrus ater*) uses vigorous wingbeats in
9 courtship displays, and such activity affects respiratory-vocal activity [7]. Finally, humans also
10 show synchronization of impulses of their limb movements with their vocalization (Pearson and
11 Pouw, 2022; Pouw et al., 2020; Pouw et al., 2023; Serré et al., 2022; for an overview see Pouw
12 and Fuchs, 2022).

13 Surprisingly, it is virtually unknown whether non-human primates synchronize their
14 whole-body movements with vocalization. One taxon seems the perfect model to study vocal-
15 motor interactions, building on anatomical and neural circuitry shared across primates, including
16 us: the Gibbons and Siamang (*Hylobatidae*). They are highly vocal species and they load their
17 entire body weight on their pectoral system during their primary mode of locomotion –
18 *brachiation*. Here we focus on the heavier-weight Siamang (*Syndactylus symphalangus*).

19 Siamang and Gibbons diverge in several ways from great apes. They are a highly vocal
20 species, performing extremely loud (sometimes > 100 Decibels [13,14], rhythmically
21 coordinated songs [15], supporting family-bonding and territory-marking, and on occasion
22 alarming. While usually understood as containing stereotyped vocalizations, Siamang song
23 contain considerable variability for certain call phrases [16]. Siamang and Gibbons apes also
24 move at extremely fast 45km/h speed using hand-over-hand grasps, also known as brachiation.
25 These Small Asian Apes also move *while* they sing (e.g., [see here](#)). Interestingly, Haimoff
26 (Haimoff, 1981), *p.* 135) observes a temporal coordination of locomotion and vocalization in the
27 wild Siamang they studied. These and other widespread *qualitative* observations [18,19] suggest
28 that singing in Siamang and Gibbons may at times be a combined display¹, such that two
29 behaviors (vocalization & movement) that in principle can operate alone, structurally operate

¹ Note that some Gibbon species (but not Siamang) are known to also dance outside of the context of singing [20], but that refers to a different sort of behavior.

1 together, much like other animals that move and vocalize in coordinated ways [1,5,7,21].
2 However, no quantitative evidence exists to support the idea of a combined display. Perhaps in
3 part because it is difficult to isolate calls often produced in group singing, and in part because
4 movements of these apes are difficult to track – at present it seems virtually impossible to track
5 movements of the Siamang in the wild as these apes move through the canopies with high speeds
6 [22].

7 One interesting possible reason for vocal-motor coupling in Siamang is biomechanics. A
8 range of animals that include their pectoral limbs for locomotion (e.g., bats, dogs, horses,
9 rhinoceros) synchronize their locomotor cycles at increasing gait speeds with respiratory cycles
10 [23]. This synchronization is held to occur because of a piston-like effect where the visceral
11 organs push forward on the lungs during accelerations and decelerations of the body during
12 locomotion strides (Daley et al., 2013; Lafortuna et al., 1996; for an overview see Pouw and
13 Fuchs, 2022), or because muscle activation for locomotion driving movement accelerations
14 simultaneously compress the rib-cage [1]. Interestingly, it has been generally acknowledged that
15 brachiation in primates likely affects respiratory control and thus vocalization [26–29]. Thus
16 from biomechanics of brachiation alone we would hypothesize some respiratory interactions,
17 suggesting some potential for vocal-motor coupling. Other non-mutually exclusive hypotheses
18 for coordination of voice and whole-body movements hold that combining visual and auditory
19 signals increases the likelihood of communicative success [30,31].

20 In this study we opportunistically observed captive Siamang engaged in solo singing to
21 assess vocal-motor coupling, whereby we audio-visually recorded singing for a period of 21
22 days. We noticed characteristic vocalizations combined with pulse-like movements (e.g., for
23 examples see [here](#)). These movement-accompanied calls contained ululating screams as main
24 units [32] and were produced by juveniles who engaged in solo singing after the duetting
25 singing performed by the entire group was completed or was winding down. To our surprise
26 these specific stereotypical calls were *always* produced with a pulse-like movement. These
27 pulse-like movements were seemingly synchronized with the calls and we will assess whether
28 these movements associate with the loudness of these solo-songs in two analyses. We know from
29 biomechanics in humans and other animals that the physical impact of a movement on the
30 musculo-skeletal system is during acceleration or deceleration (as forces are a function by mass

1 [a constant] and acceleration) [6,9,24,25]. Therefore, we test here whether thorax accelerations
2 during vocalizations statistically relate to the amplitude of concurrent vocalization in Siamang.

4 **Materials and methods**

5 **Data recording**

6 Audiovisual recordings of a family of Siamang (6 members; female adult, male adult,
7 two male juveniles, one infant, and a newborn) were collected in the June and August of 2022
8 over two visits at the Jaderpark Zoo in Lower Saxony, Germany. This yielded over 7 hours of
9 recorded singing, collected by the first and second authors. The Siamang sang primarily in the
10 morning around 9-10am, or after their fruit and vegetable lunchtime, around 1pm, and
11 occasionally around 5pm. Only two juvenile/young adult individuals performed the stereotypical
12 movement-accompanied vocalizations we consider here, and these events were ideal for acoustic
13 analysis as there was no overlap compared to the typical collective singing of Siamang. Due to
14 the adults and juveniles singing together, there is almost constant overlap in calls, which made us
15 focus on the solo singing of the juveniles in this first investigation of siamang singing and
16 movement. The two individuals were Baju (7 years and 8m) and Fajar (4 years and 11m). Baju
17 and Fajar were both born in the Jaderpark zoo. During data collection, Baju was separated from
18 the family due to risk of injury after a fight where all family members attacked Baju. This also
19 means we could not collect *more* data from Baju during our second testing period as he was
20 transferred to a smaller facility and stopped singing during our visit when the transfer was made.

21 **Audio-visual recording**

22 Four GOPRO Hero9 were installed, set at 1080 quality, sampling at 59.74fps, with linear
23 lens settings. We then cropped frames and recompressed the video to 50fps. The camera
24 positions were positioned as orthogonally from each other as the site allowed. Figure S1, panel a,
25 shows a sketched map with geometrically estimated distances based on a laser-pointer
26 measurement device. We further use in this study four audio sources from Sennheiser
27 microphones with windjammers (MKE400), plugged into the GO-PRO ensuring audio-visual
28 synchronization, sampling at 48kHz.

29 Recording was set at similar gains across microphones, and we checked for clipping
30 during pilot recordings. The four acoustic waveforms were combined using an ‘Adobe Premiere
31 Pro 2019 CC’ waveform alignment, which uses a cross-correlation approach to find the optimal

1 lag to synchronize peaks in the audio. After synchronizing the waveforms of the four sources, we
2 recompressed the audio as a single-channel audio source which thereby contains the combined
3 time-aligned sources (48kHz). Therefore, we always estimate amplitude using four combined
4 audio sources that collected from multiple locations to increase our measurement accuracy to
5 track the sound's amplitude. Specifically, we placed the audiovisual recorders at four different
6 angles; this strategy minimized problems with differences in sound radiation and differences in
7 sender-recording distances across different events, which may otherwise affect measurement of
8 amplitude peaks. However, the amplitude measurements will be flawed in open environments
9 such as these. For this reason, we also report a second time series analysis unaffected by
10 differences in amplitude measurements across events (see lagged mutual information).

11 **Identification of movement-accompanied vocalizations and related features**

12 The first and the second author identified opportunistically as many movement-
13 accompanied vocalizations by going through all the recorded songs and annotating these events
14 in ELAN [33]. These vocalizations were easy to identify because they all consisted of a ululating
15 scream as a main unit, and any variability in the call structure was stereotypical within each of
16 the two individuals. Interestingly, these stereotypical vocalizations always co-occurred with a
17 pulse-like movement, though with variable intensity and different types of locomotion. The
18 annotators (second and first author) drew a boundary around the movement-vocalization event,
19 such that the movement and the vocalization sequence was contained. We will refer to these
20 annotated events as movement-accompanied vocalizations throughout.

21 Importantly, we did not observe any stereotypical vocalizations that occurred with no
22 observable movement at all (though some contained small amplitude movements, and these are
23 part of the variability in our dataset; e.g., see supplemental table S1, Example 1). Please also see
24 a longer segment of a juvenile's singing in which it becomes clear that the stereotypical
25 movement-vocal calls that comprise our data are embedded in other types of calls such as
26 sequences of barks: see supplemental table S1, Example 2).

27 Also note that we focused on Juveniles because they would sing on their own. The adult
28 female and male produced movement and vocalizations too. However, their movement-
29 accompanied vocalizations almost always occurred when all individuals were singing. This lead
30 to frequent, multiple overlaps. Analyses of overlaps that would require another set of acoustic
31 post-processing steps that would greatly complicate the analysis relative to the current approach.

1 The degree of physical interactions expected are likely dependent on the type of action
2 performed during vocalizations; therefore, we attempted to apply a standardized description to
3 locomotor actions. We used Hunt's typology of locomotion types [34] to characterize the action
4 that occurred during the vocalization (Figure 1). In some instances, we slightly deviated from
5 Hunt's [35] categories to accommodate for a particular locomotion action (e.g., 'drop fore limb
6 swing': the individual sits on top of horizontal structure, and then scoots backwards or forward to
7 drop and then swing forward with two extended arms). A common mode of locomotion for
8 Siamang, ricochet brachiation (e.g., see supplemental table S1, Example 3), is absent in our
9 dataset, possibly due to the facility having more ropes than rigid and connected supports, and
10 thus being more tailored towards swinging movements rather than ricochet brachiation. Some
11 movement-accompanied vocalizations remained undefined as they did not fall into a clear
12 category (i.e., mixed locomotion modes). We will make a crude binary distinction between
13 locomotion types which load the entire weight on the thorax via the shoulder girdle(s) (forelimb
14 only²), or those that involve distribution of weight or support via the lower limbs (other). In the
15 case of forelimb only loads, one would particularly expect accelerations to constrain respiratory-
16 vocal interactions.

17

18 **Video preprocessing and post-processing**

19 Each event can potentially be recorded via four camera angles. We first checked all these
20 camera angles to see whether the individual was visible. If the individual was not visible, it was
21 excluded as a potential camera angle submitted for analysis. Video processing was performed in
22 Python, the specific steps discussed below. Further processing to prepare the dataset for
23 statistical analyses was performed in R and R-studio.

24 **Cutting scenes and performing initial motion detection with OpenCv2.** Firstly, a
25 custom Python script automatically cut the videos based on the ELAN annotations begin and end
26 times using moviepy, ffmpeg and pydub (see supplemental table S1, Script 1).

27 Secondly, we determined regions of interest for each scene. The installed cameras had a
28 field of view to opportunistically capture behavior at different locations. This makes tracking
29 with supervised computer vision more computationally costly as there may be many individuals

² Strictly speaking, since siamang are primarily bipedal, we could have also referred to the pectoral limbs as upper limbs (rather than forelimbs). But we decided to follow Hunt's categories here.

1 moving in a complex structured site. As a pre-processing step we therefore created a Python
2 pipeline (see supplemental table S1, Script 2) using OpenCv2 that determined pixel deviations
3 from the median to ascertain a key area of movement per frame. After smoothing and obtaining
4 maxima, this key area was used to determine a static bounding box that would further serve to
5 crop the video to primarily contain the movement of the siamang and exclude the rest of the
6 complex scene (for an example see Figure S1 panel b).

7 **Kinematics**

8 **Tracking DeepLabCut.** A convolutional neural network (Resnet-50) was trained using
9 DeepLabCut (version 2)[36] with 250 hand-labeled frames (see Supplemental table S1, Resource
10 1). Two key points were used for the training set. The first key point effectively tracked the
11 upper thorax region targeted at the most posterior region (i.e., which was used for acceleration).
12 The second key point was targeted at the sacrum of the individual (which was used for the body
13 normalization of the acceleration magnitudes).

14 We trained the model with half a million iterations, reaching an average error rate of 1
15 pixel (for keypoints with 60% confidence rates) in the training set, and 25 pixels in the test set. If
16 we normalize these errors by the original frame-sizes (1500x1080), then we get an error of
17 0,0015%.

18 When using the trained model to extract position traces for the video recordings, we
19 applied the model to all events with DLC's native filtering option to remove noise-related jitter,
20 yielding x,y position traces for the two key points (see Supplemental table S1, Example 4) and
21 likelihoods. Since derivatives increase power of noise-related jitter relative to slower frequencies
22 [37], we also applied extra smoothing of the resultant position traces with a 9th order
23 Kolmogorov Zurbenko filter with a span of 110ms (R-package `kza`).

24 Likelihoods were further used for data quality curation. Specifically, if a camera angle
25 had tracking that dropped for more than 5% below a threshold of .80, than we did not submit that
26 particular camera for further analyses; the remaining 5% of the data, that had low confidences,
27 were linearly interpolated (`na.approx` function using R-package `zoo`) using the surrounding
28 high-confidence tracking samples. Note that the DLC team has recommended a likelihood of at
29 least .60 for good tracking, and our pipeline is slightly more conservative than that.

30 **Normalization.** The thorax accelerations were calculated by differentiating speed
31 measured in pixels over time. However, different camera positions, and different locations of the

1 Siamang, make pixel-units problematic: a pixel as a unit of space will differ per distance of the
 2 camera to the individual. Therefore, we normalized the kinematics to a dimensionless quantity
 3 by scaling the position traces by the mean body size of the Siamang detected (for all frames that
 4 had a DeepLabCut confidence estimate of 100%). This means that all kinematics in this report
 5 are always first normalized by body size units.

6 **Approximating kinematics from multiple camera angles.** Depending on the location
 7 of the individual, we had one to four camera angles that recorded the movement-accompanied
 8 vocalizations. The many objects on the site, the distances of the cameras, the lack of access to the
 9 site, combined did not allow us to perform stereoscopic reconstruction of the camera angles to
 10 estimate 3D postures using a Charuco board [38,39]; this means the cameras could not be
 11 spatially calibrated for angle intrinsics and extrinsics. Note further that it would not be
 12 inappropriate to simply combine the 2D accelerations determined for each camera (by taking the
 13 Euclidean norm of all the 2D accelerations recorded) as this would yield an overestimation given
 14 that camera angles have correlated information (for example, because they all capture vertical
 15 acceleration of the individual). Therefore, to combine the information we need to find the non-
 16 correlated information in each of the camera views.

17 To still make use of multiple cameras (in the case when more than one camera had
 18 sufficient-quality data) that were impossible in the zoo to calibrate, we devised another solution
 19 to utilize multi-angle camera information using Principal Component Analyses (see e.g., [40],
 20 which formed the basis for our approximation of the 3D acceleration of the Siamang. We
 21 combined position traces of n cameras (c) available $\mathbf{M} = [x_{c=1}, y_{c=1}, x_{c=2}, y_{c=2} \dots x_{c=n}, y_{c=n}]$ to a
 22 principal component analysis, to calculate the eigenvectors (v), where we extract the three
 23 highest loading principal components, $\mathbf{PC}_{1-3} = [\mathbf{M}v(:,1), \mathbf{M}v(:,2), \mathbf{M}v(:,3)]$, which will
 24 approximate positional information about orthogonal planes. Subsequently, we take the
 25 Euclidean norm of the derivatives of the first three principal components to get an approximation
 26 of 3D speed vector s , such that $\vec{s} = \sqrt{(\Delta PC_1^2 + \Delta PC_2^2 + \Delta PC_3^2)}$. This norm is then differentiated
 27 once more, and absolutized, to yield an approximated 3D acceleration magnitude vector,
 28 ($\vec{a} = |\Delta \vec{s}|$). Note after each differentiation we smooth the data with a 5th order Kolmogorov
 29 Zurbenko filter with a span of 50ms.

30 **Acoustics**

1 **Smoothed amplitude envelope.** We extracted a smoothed amplitude envelope of the
2 waveform from the combined audio sources using a common approach, by first applying a
3 Hilbert transformation [41], and then taking the complex modulus. This resulted in a one-
4 dimensional time series, which was downsampled to 100Hz and further smoothed with a 12 Hz
5 Hanning window.

6 **Aggregation of multiple data streams**

7 We synchronized the motion tracking data and the acoustic data by first aligning the data
8 samples in time, and upsampling the motion tracking to 100Hz to preserve the high-sampling
9 rate of the acoustics. We upsampled the motion tracking data by linearly interpolating the data
10 along a time vector using R-package `zoo` (function `na.approx`; Zeileis and Grothendieck,
11 2005) so as to have a regular sampling rate at 100Hz that exactly matches the sampling times of
12 the amplitude envelope. This yields a combined time series with acoustics and motion tracking
13 that could be further processed for analysis. R-code performing the above-mentioned processing
14 steps is available online (see supplemental table S1, Script 3).

15 **Final processed dataset**

16 After having to exclude 29 events that had low-confidence tracking (trackings that
17 dropped more than 5% below a threshold of 80% likelihood determined by DeepLabCut) the
18 final dataset consisted of 83 events (of which; 26 Bajus and 57 Fajar).

19 **Analysis 1: Peak analyses**

20 For each movement-accompanied vocal event we determined the global maximum in
21 absolutized acceleration (max acc), i.e., the largest magnitude of either acceleration or
22 deceleration. Then we obtained the local maxima in the smoothed amplitude envelope, with a
23 `findpeaks` function using R-package `pracma`. This allowed us to select the magnitude of the
24 local peak in the envelope nearest to max acceleration. Since the acoustic and acceleration peak
25 data were long-tailed distributed, we log transformed the variables (which also improved model
26 fits), and z-normalized. As additional information, we also report change in amplitude in terms
27 of Decibels by scaling the raw envelope signal by $\text{dB} = 10 \cdot \log_{10}(x)$. This way, we can estimate
28 how body acceleration tends to increase the amplitude by a certain amount of dB.

29 **Analysis 2: Lagged Mutual Information Analyses**

1 We used the function `cmi` from R-package `mpmi`[43], to compute for each vocal-motor
 2 event the continuous mutual information (cmi) at different lags, where cmi at some lag is defined
 3 as $I(X;Y)_\tau$ in eq. 1:

$$4 \quad I(X;Y)_\tau = \int \int p(x_{t-\tau}, y_t) \cdot \log \left(\frac{p(x_{t-\tau}, y_t)}{p(x_{t-\tau}) \cdot p(y_t)} \right) dx_{t-\tau} dy_t \quad (1)$$

6
 7 where X_t is the z-normalized timeseries of body acceleration, Y_t is the z-normalized
 8 timeseries of amplitude envelope, and τ is a predefined lag capturing the time delay between the
 9 two variables. Note, that the $p(x, y)$ refers to the joint probability density function of X and Y,
 10 and $p(x)$ and $p(y)$ the marginal probabilities.

11 Mutual information thus provides an estimate of the information gained in Y given X (at
 12 a particular time lag). This analysis helps uncover linear and non-linear covarying information
 13 and temporal dependencies between body acceleration and amplitude envelope. By z-
 14 normalizing the time series within each event, this analysis is not dependent on the relative
 15 differences in absolute amplitude of the sound and movement between events, but rather is a test
 16 of whether movement and sound within an event are mutually coupled.

17 **False random pair.** A strong comparison condition was constructed to determine
 18 whether mutual information was higher in the observed time series relative to some baseline.
 19 Lagged mutual information was produced by comparing the amplitude envelope with a randomly
 20 simulated time series, together forming a *false random pair*. For each event comparison, the
 21 simulated body acceleration time series had the same mean, standard deviation, and
 22 autocorrelation as the original event's body acceleration time series (using an autoregressive
 23 integrated moving average, or ARIMA, procedure). We would expect that mutual information is
 24 lower in false random pairs relative to real paired movement-vocalization events.

25 **Main measures.** To summarize the data, we obtained the maximum observed mutual
 26 information between movement and sound for each event over the different lags (max mutual
 27 information). We also collected the lag at which the maximum mutual information was observed
 28 to understand whether movement is more predictive of sound or vice versa.

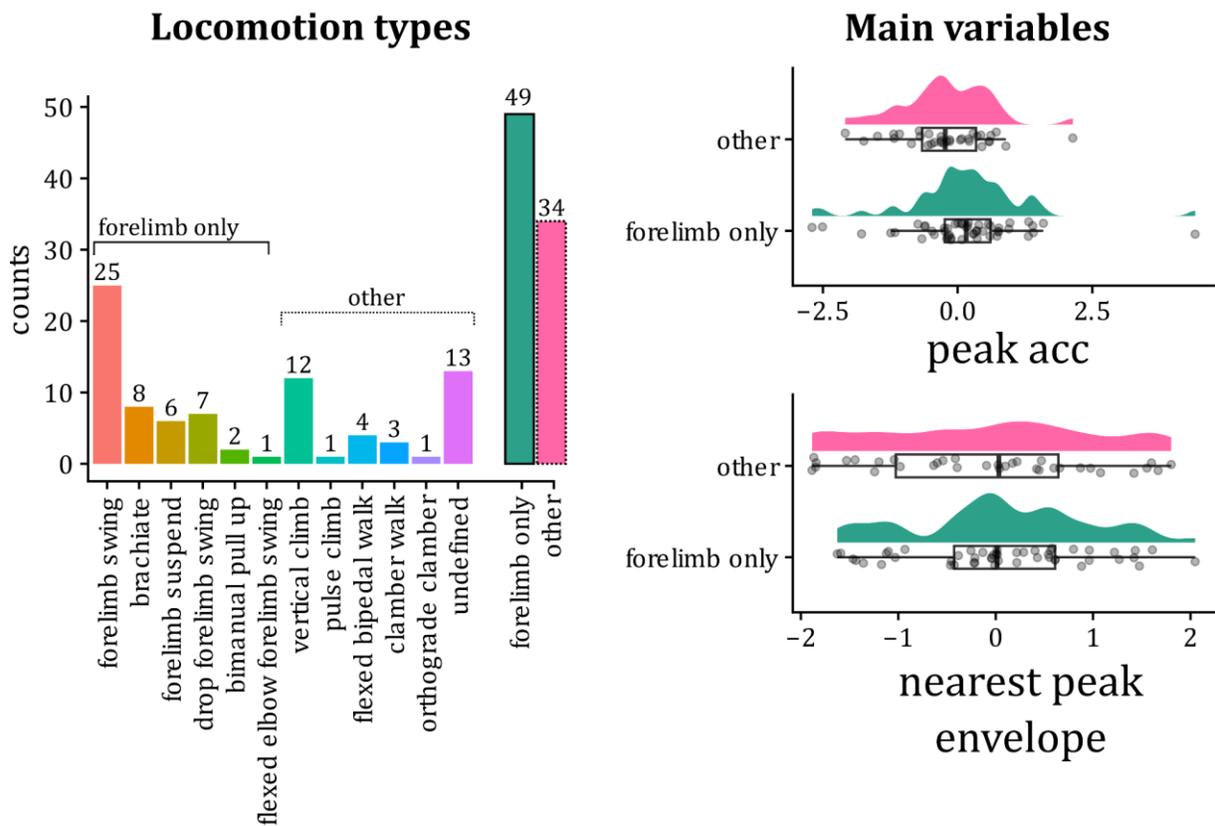
1 Descriptive statistics

2 An R markdown-script supporting the statistical analyses can be found online (see
3 Supplemental table S1, Script 4). Figure 1 and Table S2 provide information about the main
4 variables, as well as the time windows of the annotations (approximate length of the multimodal
5 events). The temporal inter-peak distance between the global maximum in acceleration and the
6 nearest local maximum in amplitude was on average 120 milliseconds ($SD = 160$). This low
7 temporal distance provides confidence that the two point-estimates for kinematics and acoustics
8 occurred sufficiently close in time to be possibly coupled. Based on 95% confidence intervals,
9 the older juvenile Bajju seemed to generate higher peaks in the amplitude envelope (nearest to
10 peak acceleration) as compared to his younger brother Fajar, while being comparable in their
11 body accelerations peaks (see Table S2). Further, forelimb only locomotion tends to have higher
12 accelerations than other locomotion types. The main variables did not dramatically differ either
13 by individual or locomotion type.

14

15

1 **Figure 1. Frequency distributions for the locomotion types, peak acceleration, and limbs.**
 2 The left panel shows the number of different locomotion types that we categorized for each
 3 vocal-motor event. Since there are too many categories, with few instances, we created super
 4 categories that indicate locomotion actions ($N = 49$ forelimb only, $N = 34$ for other types of
 5 locomotion) that only included the fore/upper limbs (forelimb only), versus those that included
 6 another limb (other). We use this super category to check whether our continuous kinematic-
 7 acoustic analyses may give different results when we also consider what type of locomotor type
 8 was performed. On the right panel, the distributions are shown for the magnitude (z-scaled units
 9 per individual) in the global peak of acceleration, per binary locomotion category, showing for
 10 example higher mean acceleration for forelimb only locomotor actions as opposed to other
 11 actions that include the hindlimbs.



12

13

Results and Discussion

14 To assess whether body accelerations constrain vocalizations amplitude, we assess
 15 whether both parameters reliably scale in magnitude. Figure 2 provides a summary of the
 16 procedure and the main results.

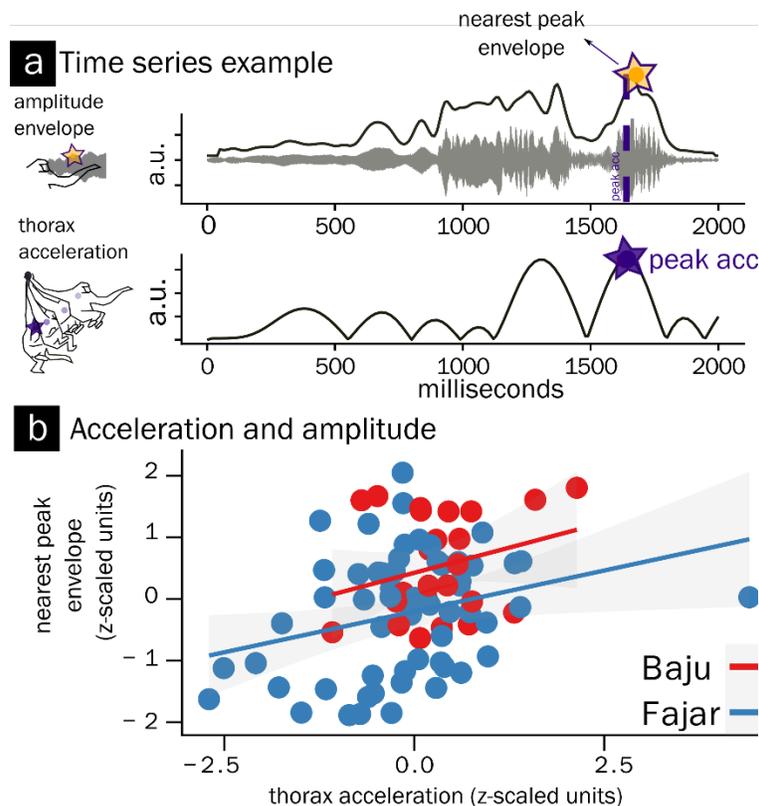
17 A mixed linear regression model was performed associating peak acceleration with
 18 nearest peak of the call amplitude envelope (using maximum likelihood from R-package `nlme`;
 19 Pinheiro et al., 2019). Individual (Fajar, Baju) was set as random intercept, but a model with
 20 random slopes for individual did not converge. By setting Individual as random intercept we

1 assess a relationship within individual, and we therefore do not simply lump together the data in
2 establishing an association between body acceleration and amplitude envelope. The model
3 statistically predicting peak envelope from peak acceleration reliably outperformed a base model
4 predicting the overall mean in amplitude; change in $\chi^2(1) = 7.68, p = .006$. The model
5 coefficients indicated that a higher magnitude peak in acceleration reliably associated with
6 higher magnitude envelope peaks, $b = 0.29, b\ 95\%CI = [0.08, 0.49], t(80) = 2.82, p = .006$,
7 intercept $b = 0.09, t(80) = 0.41, p = 0.682$. Since the models are performed on a log-log scale,
8 we can interpret the coefficients as indicating that for a unit of increase in acceleration there is a
9 30% increase in the magnitude of the nearest peak envelope (i.e., $\frac{1}{3}$ power relationship). If we
10 rescale the amplitude envelope to dB, and perform the same model we yield that per body scaled
11 change in acceleration, there is an increase in the peak envelope of about 1.4dB, $b = 1.43, b$
12 $95\%CI = [0.043, 1.435], t(80) = 2.82, p = .006$ (see [supplemental online figure](#)). A simple
13 regression analysis (which lumps the data together) yields a similar conclusion for the main
14 effect of body acceleration, $r = .33, t(81) = 3.18, p = .002$. Also, note that if we remove the
15 possible outlier shown in Figure 2 panel b (see [supplemental online figure](#) without outlier) the
16 conclusions remain unchanged, $r = .38, t(80) = 3.1829, p < .0001$.

17 We explored possible confounds (e.g., number of camera angles available) or moderators
18 (e.g., locomotion type) of this kinematic-acoustic coupling via interactions, but such interaction
19 models were not reliably outperforming our main model, $\chi^2(1) < 8.44, p's > .21$ (see
20 [supplemental online figure](#)). These checks for confounds or conditional factors provide
21 confidence that effects of acceleration and vocalization are stable among different measurement
22 conditions, individuals, and locomotion types.

23 Thus, we can conclude that the magnitude of vocalization amplitude peaks that occur
24 around moments when the thorax undergoes its maximum acceleration or deceleration is
25 positively associated with the magnitude of that acceleration. For an inspection of the movement-
26 accompanied vocalizations underlying figure 2 see our [dynamic data dashboard](#) (for code, see
27 supplemental Table S1, Script 5).

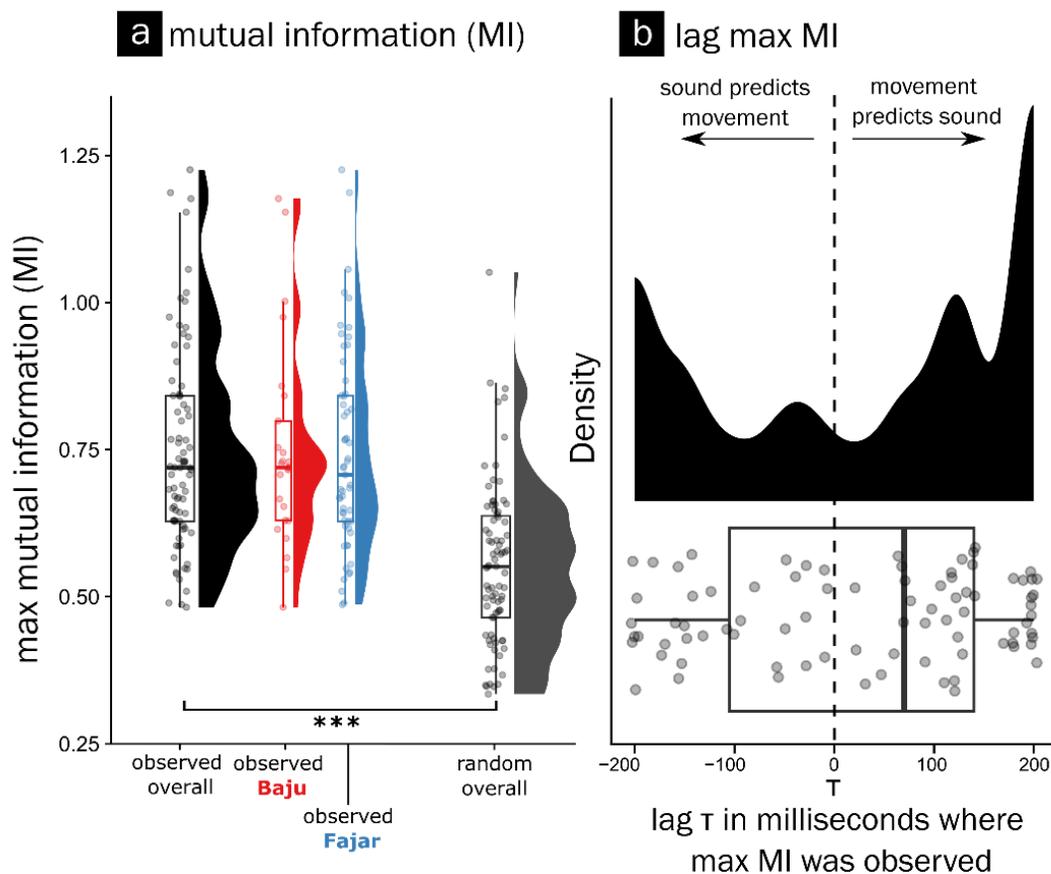
1 **Figure 2. Nearest peak analyses and results for acceleration and amplitude envelope** The
 2 time series example (a) shows data for a single event (of 83 in total, 26 Baju and 57 Fajar). The
 3 upper panel of (a) shows the smoothed amplitude envelope of the call; the lower panel of (a)
 4 shows the acceleration of the thorax, both given in arbitrary units (a.u.). The purple star shows
 5 the global maximum for acceleration. This maximum is then used to determine the nearest local
 6 maximum in the amplitude envelope (yellow star, with purple outline), which in this case also
 7 happens to be the global maximum in amplitude. The point-estimates are then submitted for
 8 linear mixed regression with random intercepts for each individual, Baju and Fajar (b). Panel (b)
 9 shows the nearest peak in body-scaled z-normalized thorax acceleration on the x-axis related to
 10 the nearest z-normalized peak in the amplitude envelope on the y-axis (this relation was found to
 11 be statistically reliable in a mixed regression analysis, $p = .006$). See the [online dynamic](#)
 12 [dashboard to explore the audiovisual data](#) underlying these data points which allows for
 13 exploring each underlying audiovisual event.
 14



15 While the peak magnitude analysis provides a promising indication that body
 16 accelerations are associated with the loudness of the call, there is always a risk of unreliable
 17 results given noisy amplitude measurement conditions, which may vary over the different
 18 occasions the calls were produced. Therefore, we devised another test to statistically confirm an
 19 association of movement and sound based on lagged mutual information calculations that assess
 20 the continuous co-variance between two continuous variables. We assess mutual information at
 21 different lags (-200 to +200 ms) in movement (body-scaled acceleration, z-normalized for each
 22

1 event) and sound (amplitude envelope, z-normalized for each event). We set the lags at a
 2 maximum of 200 milliseconds as any longer delays are unlikely to reflect a close coordination of
 3 sound and movement (but rather some sequential organization). The lagged analyses allow for
 4 possible coupling delays, indicating whether movement predicts sound in time rather than vice
 5 versa.

6 **Figure 3. Mutual information between normalized movement and sound** a) The maximum
 7 observed mutual information (in bits) in body acceleration and amplitude envelope over lags from
 8 -200 milliseconds to 200 milliseconds is shown for observed movement-accompanied vocal events
 9 (overall [N=83]; as well as split over individuals, Bajju [N=26] and Fajar [N=57]) and false random
 10 pair time series (random overall [N=83]). It is clear in a) that observed time series have much
 11 higher mutual information than random pairs which was confirmed in a linear mixed regression (p
 12 $< .0001$), and it can be seen that Bajju and Fajar show this high mutual information. b) shows the
 13 lag at which the maximum mutual information is found for all actually observed events as
 14 expressed in milliseconds. We obtain that relatively more often a maximum in mutual information
 15 (MI) is found for when movement predicts sound in time, but our mixed regression analyses did
 16 not find that a lag at max MI is reliably higher than 0ms ($p = .08$).



1 A mixed linear regression model assessed whether the observed maximum mutual
2 information is higher than a false random pair. Individual variability (Fajar, Baju) is accounted for
3 as a random intercept (random slopes did not converge). A model contrasting a false random pair
4 with the observed time series outperformed a model predicting the overall mean of mutual
5 information, $\chi^2(1) = 54.49$, $p < 0001$. The mutual information for false random pairs was indeed
6 reliably lower than the observed mutual information, $b = 0.20$, b 95%CI = [-0.25, -0.15], $t(163) =$
7 7.982 , $p < .0001$. We further assessed whether the lags at which we obtained the max mutual
8 information were reliably different from 0, thereby indicating that either sound or movement was
9 predictive of the other in time. Though movement seemed more likely to be maximally predictive
10 of sound at a lag between 0 to 200 milliseconds (see Figure 3) rather than the other way around,
11 the intercept of a mixed regression model (with individual as random intercept) predicting the
12 overall mean of the lag at max mutual information was not reliably different from 0 with a two-
13 sided test, $b = +27\text{ms}$, $t(81) = -1.80$, $p = .08$. To conclude, we find clear additional evidence that
14 there is mutual information in movement and sound, even if we ignore absolute measurements of
15 sound and movement (as we z-normalize each time series per event). This report has therefore to
16 remain inconclusive about whether movement predicts sound in time (though the data pattern
17 toward movement predicting sound rather than vice versa).

18

19 The evolution of vocal production can be framed as a continuous stabilization of multiple
20 interacting subsystems originally evolved for different purposes [45–49]. Many bodily systems
21 must fall in line to produce such skilled and energetic vocal duet singing in Gibbons and Siamang.
22 This duet singing is certainly shaped, as hypothesized till now, by several adaptations, such as
23 complex vocal tract shaping [13] and the presence of laryngeal air sacs [50–52]. However, here
24 we show that that co-vocal movement and vocalization are mutually coupled in time and
25 amplitude. We firstly obtain striking stereotypical pulse-like movements that co-occur with
26 stereotypical vocalizations in two juvenile Siamang engaged in solo singing. In fact, these
27 stereotypical vocalizations always occurred with (some) movements (few of these vocalizations
28 occurred with very little movement, and these are included in our data, see the left region of the
29 [dynamic plot to inspect the events with low body accelerations](#)). Using unsupervised and
30 supervised computer vision and data science methods we approximated the magnitude of peaks in

1 3D acceleration/deceleration of the thorax, which we then related to the magnitude of the peak in
2 the smoothed amplitude envelope nearest to the moment of peak acceleration. We observed that
3 the more the Siamang's thorax undergoes acceleration, the higher the amplitude of the
4 vocalization, i.e., the louder the call. We further find that body accelerations are associated in time
5 with changes in loudness, as established through finding higher mutual information of z-
6 normalized movement and vocalization timeseries as compared against false pairs.

7 We speculate that the non-overlapping and loud solo songs may serve an 'honest signaling'
8 function towards prospective mates (see e.g., Raemaekers et al., 1984). We would suspect that if
9 these high-amplitude solo calls are designed to inform others about one's adaptive fitness, all
10 available bodily resources will likely be recruited. This may mean that body movements stabilize
11 and amplify said calls through some biomechanical cooperation (i.e., a synergy), or it may simply
12 mean that more intense movements are coordinated as a visual display – a visual honest signal –
13 that is coordinated with vocalization (i.e., a systematic combination of two signals). The current
14 study merely shows an association between singing and movement, and more research is needed
15 to understand the mechanisms and possible adaptive benefits of mutual coordination of song and
16 movement.

17 Our results have several implications. Firstly, vigorous movements during singing in the
18 Siamang have generally been described as combined 'movement *displays*'. However, our results
19 suggest that these behaviors should also be understood as a coordination, where whole-body
20 movement and sound combined in a structural way [9,30,31]. This further has essential
21 implications for the large field of primate gesture studies, specifically as it tailors to the need for a
22 better understanding of how different modalities contribute to communicative signaling [54,55] -
23 after all, it seems that Siamang signal by using their body and their vocalization in a synchronized
24 way. This synchrony reminds us of how professional human singers couple their upper limb
25 movements during vocalizations [8]. It may be mechanistically analogous to how other non-human
26 animals couple their pectoral limb activity to vocalization [1,5–7,23]. Our findings further tie in
27 with recent research suggesting that species with more arboreal locomotion repertoires also have
28 increased vocal singing abilities [56].

29 Importantly however, coupling of the respiratory-vocal system with peripheral bodily
30 movements is not necessarily something that might drive vocal flexibility over evolutionary time

1 [5,57]. Several avian vocal learners tend to have a looser allometric scaling of their wingbeat
2 duration with their vocalization duration, as compared to birds who are not vocal learners [5].
3 Mammalian vocal learners also tend to escape allometric scaling laws that relate vocalization and
4 vocal tract size [58]. In humans, the increased flexibility in respiratory control has been in part
5 attributed to a weakening of biomechanical constraint of locomotion with respiration (as the
6 thorax was no longer impacted by locomotion; Bramble and Carrier, 1983). Furthermore, when
7 mammals, including primates, need to do forceful activities that load onto the thorax, the glottis
8 might need to be temporarily closed off to trap air to stiffen the thorax [6,9,27,59].

9 Thus, the possible functional interactions between locomotion and vocalization in
10 primates is an open question. In avian species, for example, it has been shown that vocal
11 development can be dependent on locomotor development [21,60]. Developmental research with
12 marmoset monkeys has however shown that vocal development can occur even when locomotor
13 development is delayed - Only an association of locomotor-postural development and vocal
14 development was found (Gustison et al., 2019). Thus, more research is needed that not only
15 assesses whether there is a link between locomotion and vocalization, but also how this then may
16 play out over ontogenetic and evolutionary development.

17 There are several limitations to the current study. It is based on data from a single zoo
18 obtained from captive juveniles, rather than wild adult Siamangs. It only considers certain
19 vocalizations. These issues can only be resolved by collecting and analyzing more data and
20 behaviors. Further, our measurements do not allow, at present, for more fine-grained 3D tracking
21 of the animals as the cameras were impossible to calibrate at this zoo; the study setting also did
22 not allow for perfect vocal amplitude measurements since that would require much more
23 controlled parameters (e.g., the constant distance between source, the constant direction of
24 radiation). We have reduced these issues through measuring from multiple directions and using
25 signal processing techniques to obtain unique movement information from the cameras.

26 We should further highlight some of the hypothesized mechanisms that we put forward
27 for vocal-motor interaction is on the level of biomechanics (kinematics) but our analyses and
28 findings limited to kinematics, and accordingly more work is needed on the level of
29 biomechanics [1,22,62]. Such biomechanical work would need to evaluate how locomotion
30 affects the surrounding respiratory system so that sub-glottal pressure increases due to some
31 compression of the lungs. For example, activity of muscles that insert into the human rib cage

1 (e.g., pectoralis major), measured via electromyography, predicts fluctuations in loudness in
2 steady-state vocalizations during upper limb movements [11]. In Siamang similar but less
3 invasive investigations could consist of measuring branch reaction forces during locomotion [22]
4 and directly relating this to modulations in singing.

5 **Conclusions**

6 The potential coordination of singing with movement has fascinated gibbon researchers
7 for a long time [17]. At the same time, current leading primatologists see the coordination of the
8 voice and communicative movement as something that is not common in non-human primates
9 [63]. The current report provides quantitative support that primate vocalization and whole-body
10 movement is coordinating. Siamangs fluctuate their song amplitude together with body
11 accelerations. They thus coordinate movement and singing in a way that is currently unknown.
12 We hope therefore that this research further invites inquiry into the impressive movement-
13 entangled vocalizations of Siamang and Gibbons.

15 **References**

- 16 1. Lancaster WC, Henson OW, Keating AW. Respiratory muscle activity in relation to
17 vocalization in flying bats. *J Exp Biol.* 1995;198: 175–191.
- 18 2. Suthers RA, Thomas SP, Suthers BJ. Respiration, wing-beat and ultrasonic pulse emission
19 in an echo-locating bat. *J Exp Biol.* 1972;56: 37–48.
- 20 3. Stidsholt L, Johnson M, Goerlitz HR, Madsen PT. Wild bats briefly decouple sound
21 production from wingbeats to increase sensory flow during prey captures. *iScience.*
22 2021;24: 102896. doi:10.1016/j.isci.2021.102896
- 23 4. Koblitz JC, Stilz P, Schnitzler H-U. Source levels of echolocation signals vary in
24 correlation with wingbeat cycle in landing big brown bats (*Eptesicus fuscus*). *J Exp Biol.*
25 2010;213: 3263–3268. doi:10.1242/jeb.045450
- 26 5. Berg KS, Delgado S, Mata-Betancourt A. Phylogenetic and kinematic constraints on avian
27 flight signals. *Proc R Soc B Biol Sci.* 2019;286: 20191083. doi:10.1098/rspb.2019.1083
- 28 6. Blumberg M. Rodent ultrasonic short calls: locomotion, biomechanics, and communication.
29 *J Comp Psychol.* 1992;106: 360–365. doi:10.1037/0735-7036.106.4.360
- 30 7. Cooper BG, Goller F. Multimodal signals: enhancement and constraint of song motor
31 patterns by visual display. *Science.* 2004;303: 544–546. doi:10.1126/science.1091099

- 1 8. Pearson L, Pouw W. Gesture–vocal coupling in Karnatak music performance: A neuro–
2 bodily distributed aesthetic entanglement. *Ann N Y Acad Sci.* 2022;n/a.
3 doi:10.1111/nyas.14806
- 4 9. Pouw W, Fuchs S. Origins of vocal-entangled gesture. *Neurosci Biobehav Rev.* 2022;141:
5 104836. doi:10.1016/j.neubiorev.2022.104836
- 6 10. Pouw W, Paxton A, Harrison SJ, Dixon JA. Acoustic information about upper limb
7 movement in voicing. *Proc Natl Acad Sci.* 2020 [cited 12 May 2020].
8 doi:10.1073/pnas.2004163117
- 9 11. Pouw W, Werner R, Burchardt L, Selen L. The human voice aligns with whole-body
10 kinetics. *bioRxiv*; 2023. p. 2023.11.28.568991.
11 doi:https://doi.org/10.1101/2023.11.28.568991
- 12 12. Serré H, Dohen M, Fuchs S, Gerber S, Rochet-Capellan A. Leg movements affect speech
13 intensity. *J Neurophysiol.* 2022 [cited 23 Sep 2022]. doi:10.1152/jn.00282.2022
- 14 13. Koda H, Nishimura T, Tokuda IT, Oyakawa C, Nihonmatsu T, Masataka N. Soprano
15 singing in gibbons. *Am J Phys Anthropol.* 2012;149: 347–355. doi:10.1002/ajpa.22124
- 16 14. McAngus Todd NP, Merker B. Siamang gibbons exceed the saccular threshold: Intensity of
17 the song of *Hylobates syndactylus*. *J Acoust Soc Am.* 2004;115: 3077–3080.
18 doi:10.1121/1.1736273
- 19 15. Raimondi T, Di Panfilo G, Pasquali M, Zarantonello M, Favaro L, Savini T, et al.
20 Isochrony and rhythmic interaction in ape duetting. *Proc R Soc B Biol Sci.* 2023;290:
21 20222244. doi:10.1098/rspb.2022.2244
- 22 16. D’Agostino J, Spehar S, Abdullah A, Clink DJ. Evidence for Vocal Flexibility in Wild
23 Siamang (*Symphalangus syndactylus*) Ululating Scream Phrases. *Int J Primatol.* 2023 [cited
24 20 Sep 2023]. doi:10.1007/s10764-023-00384-5
- 25 17. Haimoff EH. Video Analysis of Siamang (*Hylobates syndactylus*) Songs. *Behaviour.*
26 1981;76: 128–151.
- 27 18. Mitani JC. Gibbon Song Duets and Intergroup Spacing. *Behaviour.* 1985;92: 59–96.
- 28 19. Redmond J, Lamperez A. Leading limb preference during brachiation in the gibbon family
29 member, *Hylobates syndactylus* (siamangs): A study of the effects of singing on
30 lateralisation. *Laterality.* 2004;9: 381–396. doi:10.1080/13576500342000211
- 31 20. Fan P-F, Ma C-Y, Garber PA, Zhang W, Fei H-L, Xiao W. Rhythmic displays of female
32 gibbons offer insight into the origin of dance. *Sci Rep.* 2016;6: 34606.
33 doi:10.1038/srep34606
- 34 21. Berg KS, Beissinger S, Bradbury J. Factors shaping the ontogeny of vocal signals in a wild
35 parrot. *J Exp Biol.* 2013;216: 338–345. doi:10.1242/jeb.073502

- 1 22. Cheyne SM. Gibbon Locomotion Research in the Field: Problems, Possibilities, and
2 Benefits for Conservation. In: D’Août K, Vereecke EE, editors. *Primate Locomotion:
3 Linking Field and Laboratory Research*. New York, NY: Springer; 2011. pp. 201–213.
4 doi:10.1007/978-1-4419-1420-0_11
- 5 23. Boggs DF. Interactions between locomotion and ventilation in tetrapods. *Comp Biochem
6 Physiol A Mol Integr Physiol*. 2002;133: 269–288. doi:10.1016/S1095-6433(02)00160-5
- 7 24. Daley MA, Bramble DM, Carrier DR. Impact loading and locomotor-respiratory
8 coordination significantly influence breathing dynamics in running humans. *PLOS ONE*.
9 2013;8: e70752. doi:10.1371/journal.pone.0070752
- 10 25. Lafortuna CL, Reinach E, Saibene F. The effects of locomotor-respiratory coupling on the
11 pattern of breathing in horses. *J Physiol*. 1996;492 (Pt 2): 587–596.
12 doi:10.1113/jphysiol.1996.sp021331
- 13 26. Harrison DFN. *The anatomy and physiology of the mamalian larynx*. Cambridge:
14 Cambridge University Press; 1995.
- 15 27. Hayama S. The origin of the completely closed glottis. Why does not the monkey fall from
16 a tree? *Primate Res*. 1996;12: 179–206. doi:10.2354/psj.12.179
- 17 28. Mott F. A Study by Serial Sections of the Structure of the Larynx of *Hylobates syndactylus*
18 (Siamang Gibbon). *Proc Zool Soc Lond*. 1924;94: 1161–1170. doi:10.1111/j.1096-
19 3642.1924.tb03336.x
- 20 29. Negus VE. *The comparative anatomy and physiology of the larynx*. London: William
21 Heinemann Medical Books; 1949.
- 22 30. Partan SR, Marler P. Communication Goes Multimodal. *Science*. 1999;283: 1272–1273.
23 doi:10.1126/science.283.5406.1272
- 24 31. Halfwerk W, Varkevisser J, Simon R, Mendoza E, Scharff C, Riebel K. Toward Testing for
25 Multimodal Perception of Mating Signals. *Front Ecol Evol*. 2019;7.
26 doi:10.3389/fevo.2019.00124
- 27 32. Geissmann T. Duet Songs of the Siamang, *Hylobates Syndactylus*: II. Testing the Pair-
28 Bonding Hypothesis during a Partner Exchange. *Behaviour*. 1999;136: 1005–1039.
- 29 33. Wittenburg P, Brugman H, Russel A, Klassmann A, Sloetjes H. ELAN: a Professional
30 Framework for Multimodality Research. 2006; 4.
- 31 34. Hunt KD. Why are there apes? Evidence for the co-evolution of ape and monkey
32 ecomorphology. *J Anat*. 2016;228: 630–685. doi:10.1111/joa.12454
- 33 35. Hunt KD, Cant JGH, Gebo DL, Rose MD, Walker SE, Youlatos D. Standardized
34 descriptions of primate locomotor and postural modes. *Primates*. 1996;37: 363–387.
35 doi:10.1007/BF02381373

- 1 36. Mathis A, Mamidanna P, Cury KM, Abe T, Murthy VN, Mathis MW, et al. DeepLabCut:
2 markerless pose estimation of user-defined body parts with deep learning. *Nat Neurosci.*
3 2018;21: 1281–1289. doi:10.1038/s41593-018-0209-y
- 4 37. Winter DA. *Biomechanics and Motor Control of Human Movement.* John Wiley & Sons;
5 2009.
- 6 38. Karashchuk P, Rupp KL, Dickinson ES, Walling-Bell S, Sanders E, Azim E, et al. Anipose:
7 A toolkit for robust markerless 3D pose estimation. *Cell Rep.* 2021;36: 109730.
8 doi:10.1016/j.celrep.2021.109730
- 9 39. Theriault DH, Fuller NW, Jackson BE, Bluhm E, Evangelista D, Wu Z, et al. A protocol
10 and calibration method for accurate multi-camera field videography. *J Exp Biol.* 2014;
11 jeb.100529. doi:10.1242/jeb.100529
- 12 40. Mena-Chalco JP, Macêdo I, Velho L, Cesar RM. 3D face computational photography using
13 PCA spaces. *Vis Comput.* 2009;25: 899–909. doi:10.1007/s00371-009-0373-x
- 14 41. He L, Dellwo V. Amplitude envelope kinematics of speech: Parameter extraction and
15 applications. *J Acoust Soc Am.* 2017;141: 3582–3582. doi:10.1121/1.4987638
- 16 42. Zeileis A, Grothendieck G. zoo: S3 Infrastructure for Regular and Irregular Time Series. *J*
17 *Stat Softw.* 2005;14: 1–27. doi:10.18637/jss.v014.i06
- 18 43. Pardy C. mpmi: Mixed-pair mutual information estimators version 0.4 from R-Forge. 2012.
19 Available: <https://rdr.io/rforge/mpmi/>
- 20 44. Pinheiro J, Bates D, DebRoy S, Sarkar D, R Team RC. *nlme: Linear and nonlinear mixed*
21 *effects models.* 2019.
- 22 45. Pouw W, Proksch S, Drijvers L, Gamba M, Holler J, Kello C, et al. Multilevel rhythms in
23 multimodal communication. *Philos Trans R Soc B Biol Sci.* 2021.
24 doi:10.1098/rstb.2020.0334
- 25 46. Deacon TW. *The symbolic species: The co-evolution of language and the brain.* W.W.
26 Norton; 1998.
- 27 47. Deacon TW. *Incomplete Nature: How Mind Emerged from Matter.* W.W. Norton; 2013.
- 28 48. MacLarnon AM, Hewitt GP. The evolution of human speech: the role of enhanced
29 breathing control. *Am J Phys Anthropol.* 1999;109: 341–363. doi:10.1002/(SICI)1096-
30 8644(199907)109:3<341::AID-AJPA5>3.0.CO;2-2
- 31 49. MacNeilage PF. *The origin of speech.* Oxford University Press; 2010.
- 32 50. Burchardt LS, Sande Y van de, Kehy M, Gamba M, Ravignani A, Pouw W. A computer
33 vision toolkit for the dynamic study of air sacs in Siamang with a general application to the

- 1 study of elastic kinematics in other animals. 2023 [cited 15 Oct 2023]. Available:
2 <https://ecoevorxiv.org/repository/view/6080/>
- 3 51. de Boer B. Acoustic analysis of primate air sacs and their effect on vocalization. *J Acoust*
4 *Soc Am.* 2009;126: 3329–3343. doi:10.1121/1.3257544
- 5 52. Dunn JC. Sexual selection and the loss of laryngeal air sacs during the evolution of speech.
6 *Anthropol Sci.* 2018;126: 29–34. doi:10.1537/ase.180309
- 7 53. Raemaekers JJ, Raemaekers PM, Haimoff EH. Loud Calls of the Gibbon (*Hylobates lar*):
8 Repertoire, Organisation and Context. *Behaviour.* 1984;91: 146–189.
- 9 54. Liebal K, Slocombe KE, Waller BM. The language void 10 years on: multimodal primate
10 communication research is still uncommon. *Ethol Ecol Evol.* 2022;0: 1–14.
11 doi:10.1080/03949370.2021.2015453
- 12 55. Slocombe KE, Waller BM, Liebal K. The language void: the need for multimodality in
13 primate communication research. *Anim Behav.* 2011;81: 919–924.
14 doi:10.1016/j.anbehav.2011.02.002
- 15 56. Schruth DM, Templeton CN, Holman DJ, Smith EA. Evolution of primate protomusicality
16 via locomotion. *bioRxiv*; 2021. p. 2020.12.29.424766. doi:10.1101/2020.12.29.424766
- 17 57. Bramble D, Carrier DR. Running and breathing in mammals. *Science.* 1983;219: 251–256.
18 doi:10.1126/science.6849136
- 19 58. Garcia M, Ravignani A. Acoustic allometry and vocal learning in mammals. *Biol Lett.*
20 2020;16: 20200081. doi:10.1098/rsbl.2020.0081
- 21 59. Orlikoff RF. Voice Production during a Weightlifting and Support Task. *Folia Phoniatr*
22 *Logop.* 2008;60: 188–194. doi:10.1159/000128277
- 23 60. Liu W-C, Landstrom M, Cealie M, MacKillop I. A juvenile locomotor program promotes
24 vocal learning in zebra finches. *Commun Biol.* 2022;5: 573. doi:10.1038/s42003-022-
25 03533-3
- 26 61. Gustison ML, Borjon JI, Takahashi DY, Ghazanfar AA. Vocal and locomotor coordination
27 develops in association with the autonomic nervous system. Tchernichovski O, Calabrese
28 RL, Goller F, editors. *eLife.* 2019;8: e41853. doi:10.7554/eLife.41853
- 29 62. Bertram JEA. New perspectives on brachiation mechanics. *Am J Phys Anthropol.* 2004;39:
30 100–117. doi:10.1002/ajpa.20156
- 31 63. Fröhlich M, Sievers C, Townsend SW, Gruber T, Schaik CP. Multimodal communication
32 and language origins: integrating gestures and vocalizations. *Biol Rev.* 2019;94: 1809–
33 1829. doi:10.1111/brv.12535

Supplemental materials

1

2 Table S1. Links to online materials

Reference	Information	Hyperlink
Example 1	Example of data of small amplitude movements	link
Example 2	Example of how stereotypical calls analyzed in this research are embedded in longer singing sequences	link
Example 3	Example ricochetal brachiation	link
Example 4	Example of video tracking DeepLaBcut	link
Script 1	Python code for snipping videos from annotations	link
Script 2	Python OpenCv2 automatic bounding box pre-processing script	link
Script 3	R Processing script	link
Script 4	R Statistical analysis	link
Script 5	Python reproducible code for data dashboard	link
Resource 1	Deeplabcut trained model + info	link

3

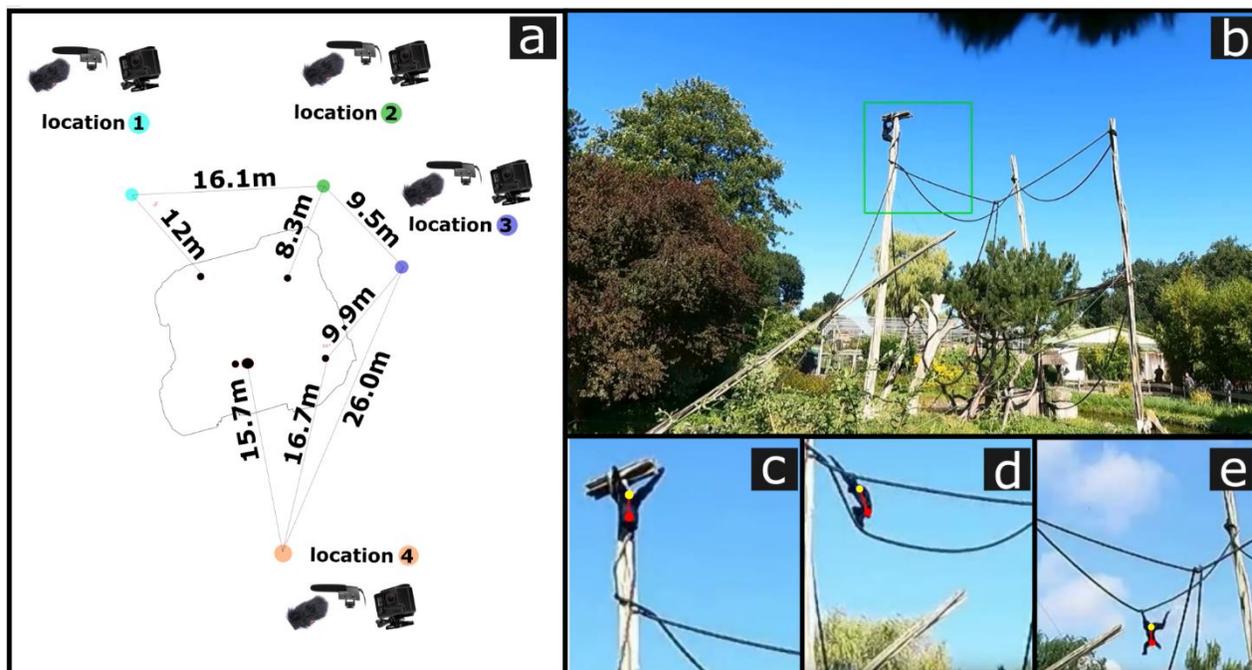
4

5 Table S2. Descriptives stats main variables

	Overall	Baju	Fajar	Forelimb only	Other
	<i>M</i> (<i>SD</i>) 95%CI[]	<i>M</i> (<i>SD</i>) 95%CI[]	<i>M</i> (<i>SD</i>) 95%CI[]	<i>M</i> (<i>SD</i>) 95%CI[]	<i>M</i> (<i>SD</i>) 95%CI[]
					
Time annotation	2851 ms (491) [2744, 2959]	2495 (268) [2386, 2603]	3014 (485) [2886, 3143]	2815 (485) [2675, 2954]	2905 (504) [2729,3080]
Nearest peak envelope (<i>z</i>)		0.52 (0.82) [0.19, 0.85]	-0.24 (.99) [-0.50, 0.03]	0.07 (0.91) [-0.19, 0.33]	-0.10 (1.12) [-0.49, 0.29]
Peak acceleration (<i>z</i>)		0.29 (0.69) [0.01, 0.57]	-0.13 (1.09) [-0.42, 0.16]	0.15 (1.09) [-0.21, 0.46]	-0.21 (0.83) [-0.50, 0.08]
Inter-peak distance	121 ms (160) [86, 156]	157ms (165) [90, 224]	104 (156) [62, 146]	107 (154) [62, 151]	141 (169) [83, 200]

6 *Note.* Max nearest envelope is the *z*-scaled magnitude of log peak smoothed amplitude envelope. Note that there are
7 no overall descriptives for variables that have been *z*-scaled (amounting to $M = 0$, $sd = 1$).

1 Figure S1. Keypoints tracked for different frames



2
3 *Note.* Panel a) shows a sketch of the location with four different recording sites with camera and microphones and
4 estimated distances obtained via laser-based distance estimation. b) Shows a frame from location 4 and an automatic
5 selection of the frame denoted by the green rectangle. c-d) These subframes containing movement in the frame were
6 then used for training DLC model to detect the two keypoints, shown in yellow (upper thorax, at T1) and red
7 (sacrum). The Euclidean distance of the upper thorax to the sacrum was used to normalize kinematics expressed in
8 pixel units to units relative to body size.
9

10