

# Emergent intensity invariance in a physiologically inspired model of the grasshopper auditory system

Jona Hartling, Jan Benda

## 1 Exploring a grasshopper’s sensory world

Our scientific understanding of sensory processing systems results from the distributed accumulation of anatomical, physiological and ethological evidence. This process is undoubtedly without alternative; however, it leaves us with the challenge of integrating the available fragments into a coherent whole in order to address issues such as the interaction between individual system components, the functional limitations of the system overall, or taxonomic comparisons between systems that process the same sensory modality. Any unified framework that captures the essential functional aspects of a given sensory system thus has the potential to deepen our current understanding and facilitate systematic investigations. However, building such a framework is a challenging task. It requires a wealth of existing knowledge of the system and the signals it operates on, a clearly defined scope, and careful reduction, abstraction, and formalization of the underlying structures and mechanisms.

One sensory system about which extensive information has been gathered over the years is the auditory system of grasshoppers (*Acrididae*). Grasshoppers rely on their sense of hearing primarily for intraspecific communication, which includes mate attraction (D. v. Helversen 1972) and evaluation (Stange and Ronacher 2012), sender localization (D. v. Helversen and Rheinlaender 1988), courtship display (Elsner 1968), rival deterrence (Greenfield and Minckley 1993), and loss-of-signal predator alarm (SOURCE). In accordance with this rich behavioral repertoire, grasshoppers have evolved a variety of sound production mechanisms to generate acoustic communication signals for different contexts and ranges using their wings, hindlegs, or mandibles (Otte 1970). Among the most conspicuous acoustic signals of grasshoppers are their species-specific calling songs, which broadcast the presence of the singing individual — mostly the males of the species — to potential mates within range. These songs are usu-

ally more characteristic of a species than morphological traits (Tishechkin and Vedenina 2016, Tarasova et al. 2021), which can vary greatly within species (Rowell 1972, Köhler et al. 2017). The reliance on songs to mediate reproduction represents a strong evolutionary driving force, that resulted in a massive species diversification (Vedenina and Muge 2011, Sevastianov et al. 2023), with over 6800 recognized grasshopper species in the *Acrididae* family (Cigliano et al. 2024). It is this diversity of species, and the crucial role of acoustic communication in its emergence, that makes the grasshopper auditory system an intriguing candidate for attempting to construct a functional model framework. As a necessary reduction, the model we propose here focuses on the pathway responsible for the recognition of species-specific calling songs, disregarding other auditory functions such as directional hearing (D. v. Helversen 1984, Ronacher, D. v. Helversen, and Helversen 1986, D. v. Helversen and Rheinlaender 1988).

To understand the functional challenges faced by the grasshopper auditory system, one has to understand the properties of the songs it is designed to recognize. Grasshopper songs are amplitude-modulated broad-band acoustic signals. Most songs are produced by stridulation, during which the animal pulls the serrated stridulatory file on its hindlegs across a resonating vein on the forewings (O. v. Helversen and Elsner 1977, Stumpner and Helversen 1994, D. v. Helversen and O. v. Helversen 1997). Every tooth that strikes the vein generates a brief pulse of sound. Multiple pulses make up a syllable; and the alternation of syllables and relatively quiet pauses forms a characteristic, through noisy, waveform pattern. Song recognition depends on certain temporal and structural parameters of this pattern, such as the duration of syllables and pauses (D. v. Helversen 1972), the slope of pulse onsets (D. v. Helversen 1993), and the accentuation of syllable onsets relative to the preceeding pause (Balakrishnan et al. 2001, D. v. Helversen, Balakrishnan, and Helversen 2004). The amplitude modulation, or envelope, of the song is sufficient for recognition (D. v. Helversen and O. v. Helversen 1997). However, the essential recognition cues can vary considerably with external physical factors, which requires the auditory system to be invariant to such variations in order to reliably recognize songs under different conditions. For instance, the temporal structure of grasshopper songs warps with temperature (Skovmand and Boel Pedersen 1983). The auditory system can compensate for this variability by reading out relative temporal relationships rather than absolute time intervals (Creutzig, Wohlgemuth, et al. 2009, Creutzig, Benda, et al. 2010), as those remain relatively constant across different temperatures (D. v. Helversen 1972). Another, perhaps even more fundamental external source of song variability lays in the attenuation of sound intensity with increasing distance to the sender. Sound attenuation depends on both the frequency content of the signal and the

vegetation of the habitat (Michelsen 1978). For the receiving auditory system, this has two major implications. First, the amplitude dynamics of the song pattern are steadily degraded over distance, which limits the effective communication range of grasshoppers to 1-2 m in their typical grassland habitats (Lang 2000). Second, the overall intensity level of songs at the receiver's position varies depending on the location of the sender, which should ideally not affect the recognition of the song pattern. This necessitates that the auditory system achieves a certain degree of intensity invariance — a time scale-selective sensitivity to faster amplitude dynamics and simultaneous insensitivity to slower, more sustained amplitude dynamics. Intensity invariance in different auditory systems is often associated with neuronal adaptation (Benda and Hennig 2008, Barbour 2011, Ozeri-Engelhard et al. 2018), which represents an important principle of dynamic sensory systems in general (Benda 2021). In the grasshopper auditory system, a number of neuron types along the processing chain exhibit spike-frequency adaptation in response to sustained stimulus intensities (Römer 1976, Gollisch et al. 2002, Hildebrandt et al. 2009, Clemens, Weschke, et al. 2010) and thus likely contribute to the emergence of intensity-invariant song representations. This means that intensity invariance is not the result of a single processing step but rather a gradual process, in which different neuronal populations contribute to varying degrees (Clemens, Weschke, et al. 2010) and by different mechanisms (Hildebrandt et al. 2009). Approximating this process within a functional model framework thus requires a considerable amount of simplification. In this work, we demonstrate that even a small number of basic physiologically inspired signal transformations — specifically, pairs of nonlinear and linear operations — is sufficient to achieve a meaningful degree of intensity invariance. Due to the critical role of intensity-invariant representations for reliable song recognition, these transformations are at the core of the proposed model framework.

Multi-species, multi-individual communally inhabited environments

- Temporal overlap: Simultaneous singing across individuals/species common
  - Frequency overlap: Little speciation into frequency bands (likely unused)
  - "Biotic noise": Hetero-/conspecifics ("Another one's songs are my noise")
  - "Abiotic noise": Wind, water, vegetation, anthropogenic
  - Effects of habitat structure on sound propagation (landscape - soundscape)
- Sensory constraints imposed by the (acoustic) environment

Cluster of auditory challenges (interlocking constraints → tight coupling):

From continuous acoustic input, generate neuronal representations that...

- 1)...allow for the separation of relevant (song) events from ambient noise floor
- 2)...compensate for behaviorally non-informative song variability (invariances)

3)...carry sufficient information to characterize different song patterns, recognize the ones produced by conspecifics, and make appropriate behavioral decisions based on context (sender identity, song type, mate/rival quality)

How can the auditory system of grasshoppers meet these challenges?

- What are the minimum functional processing steps required?
- Which known neuronal mechanisms can implement these steps?
- Which and how many stages along the auditory pathway contribute?
- What are the limitations of the system as a whole?

How can a human observer conceive a grasshopper's auditory percepts?

- How to investigate the workings of the auditory pathway as a whole?
- How to systematically test effects and interactions of processing parameters?
- How to integrate the available knowledge on anatomy, physiology, ethology?
- Abstract, simplify, formalize → Functional model framework

### **Precursor work for model construction (special thanks to authors):**

Linear-nonlinear modelling of behavioral responses to artificial songs

- Feature expansion as implemented in our model: Major contribution!
- Bank of linear filters, nonlinearity, temporal integration, feature weighting
- Clemens and Hennig 2013 (crickets)
- Clemens and Ronacher 2013 (grasshoppers)
- Ronacher, Hennig, and Clemens 2015

### **Own advancements/key differences:**

- 1) Used boxcar functions as artificial "songs" (focus on few key parameters)  
→ Now actual, variable songs (as naturalistic as possible)
- 2) Fitted filters to behavioral data  
→ More general, simpler, unfitted formalized Gabor filter bank

## **2 Developing a functional model of the grasshopper auditory pathway**

The grasshopper auditory system has been studied extensively over the past decades; and a corresponding number of involved neuron types has been described (Rehbein et al. 1974; Kalmring 1975; Rehbein 1976; Eichendorf and Kalmring 1980). The functional model we propose here focuses on the pathway responsible for song recognition and assumes a strict feed-forward organization of three consecutive neuronal populations: Peripheral auditory

receptor neurons (1st order), local interneurons of the metathoracic ganglion (2nd order), and ascending neurons (3rd order) projecting towards the supraesophageal ganglion.

Previous authors have reported a marked increase in response heterogeneity within the population of ascending neurons compared to receptors and local interneurons, which exhibit almost identical filter characteristics, respectively (Clemens, Kutzki, et al. 2011). Based on these findings, the model pathway can be divided into two distinct portions (Fig. 1c+d). In the preprocessing portion, generated

The preprocessing portion comprises the tympanal membrane, receptors, and local interneurons. The different signal representations

Due to the similar response properties within the involved

1) "Pre-split portion" of the auditory pathway:

Tympanal membrane → Receptor neurons → Local interneurons

Similar response/filter properties within receptor/interneuron populations (Clemens, Kutzki, et al. 2011)

→ One population-wide response trace per stage (no "single-cell resolution")

2) "Post-split portion" of the auditory pathway:

Ascending neurons (AN) → Central brain neurons

Diverse response/filter properties within AN population (Clemens, Kutzki, et al. 2011)

- Pathway splitting into several parallel branches

- Expansion into a decorrelated higher-dimensional sound representation

→ Individual neuron-specific response traces from this stage onwards

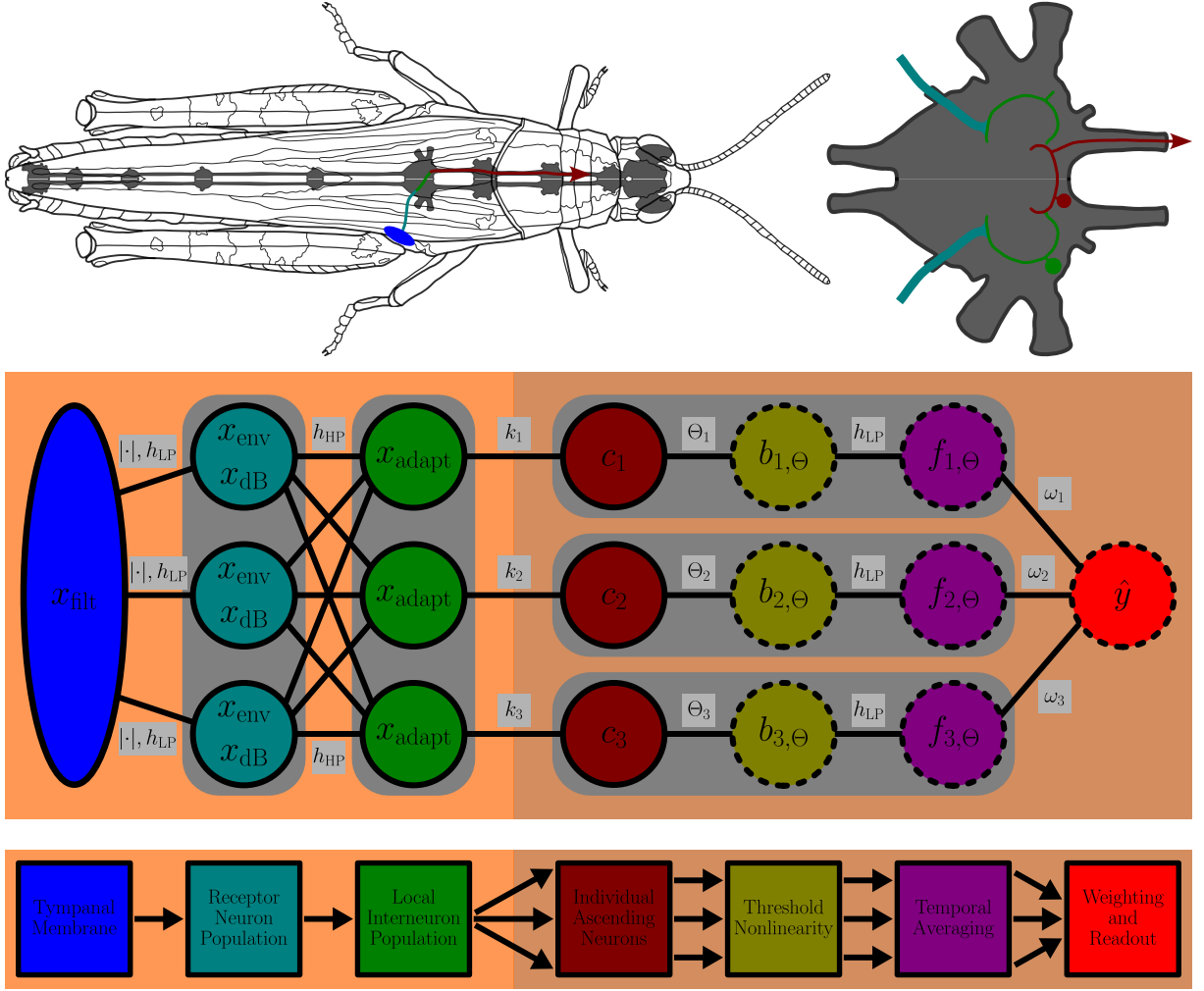


Figure 1: The auditory system of grasshoppers.

## 2.1 Population-driven signal pre-processing

Grasshoppers receive airborne sound waves by a tympanal organ at each side of the thorax (Fig. 1a). The tympanal membrane acts as a mechanical resonance filter: Vibrations that fall within specific frequency bands are focused on different membrane areas, while others are attenuated (Michelsen 1971; Windmill et al. 2008; Malkin et al. 2014). This processing step can be approximated by an initial bandpass filter

$$x_{\text{filt}}(t) = x(t) * h_{\text{BP}}(t), \quad f_{\text{cut}} = 5 \text{ kHz}, 30 \text{ kHz} \quad (1)$$

applied to the acoustic input signal  $x(t)$ . The auditory receptor neurons connect directly to the tympanal membrane (Fig. 1a). Besides performing the mechano-electrical transduction, the receptor population is substrate to several known processing steps. First, the receptors extract the signal envelope (Machens, Prinz, et al. 2001), which likely in-

volves a rectifying nonlinearity (Machens, Stemmler, et al. 2001). This can be modelled as full-wave rectification followed by lowpass filtering

$$x_{\text{env}}(t) = |x_{\text{filt}}(t)| * h_{\text{LP}}(t), \quad f_{\text{cut}} = 500 \text{ Hz} \quad (2)$$

of the tympanal signal  $x_{\text{filt}}(t)$ . Furthermore, the receptors exhibit a sigmoidal response curve over logarithmically compressed intensity levels (Suga 1960; Gollisch et al. 2002). In the model, logarithmic compression is achieved by conversion to decibel scale

$$x_{\text{dB}}(t) = 10 \cdot \log_{10} \frac{x_{\text{env}}(t)}{x_{\text{ref}}}, \quad x_{\text{ref}} = \max[x_{\text{env}}(t)] \quad (3)$$

relative to the maximum intensity  $x_{\text{ref}}$  of the signal envelope  $x_{\text{env}}(t)$ . Next, the axons of the receptor neurons project into the metathoracic ganglion, where they synapse onto local interneurons (Fig. 1b). Both the local interneurons (Hildebrandt et al. 2009; Clemens, Weschke, et al. 2010) and, to a lesser extent, the receptors themselves (Fisch et al. 2012) display spike-frequency adaptation in response to sustained stimulus intensity levels. This mechanism allows for the robust encoding of faster amplitude modulations against a slowly changing overall baseline intensity. Functionally, this processing step resembles a highpass filter

$$x_{\text{adapt}}(t) = x_{\text{dB}}(t) * h_{\text{HP}}(t), \quad f_{\text{cut}} = 10 \text{ Hz} \quad (4)$$

over the logarithmically scaled envelope  $x_{\text{dB}}(t)$ . The projections of the local interneurons remain within the metathoracic ganglion and synapse onto a small number of ascending neurons (Fig. 1b), which marks the transition between the preprocessing stream and the parallel processing stream of the model pathway.

## 2.2 Feature extraction by individual neurons

The small population of ascending neurons

### Stage-specific processing steps and functional approximations:

Template matching by individual ANs

- Filter base (STA approximations): Set of Gabor kernels
- Gabor parameters:  $\sigma_i, \phi_i, f_i \rightarrow$  Determines kernel sign and lobe number

$$k_i(t, \sigma_i, f_i, \phi_i) = e^{-\frac{t^2}{2\sigma_i^2}} \cdot \sin(2\pi f_i \cdot t + \phi_i) \quad (5)$$

→ Separate convolution with each member of the kernel set

$$c_i(t) = x_{\text{adapt}}(t) * k_i(t) = \int_{-\infty}^{+\infty} x_{\text{adapt}}(\tau) \cdot k_i(t - \tau) d\tau \quad (6)$$

Thresholding nonlinearity in ascending neurons (or further downstream)

- Binarization of AN response traces into "relevant" vs. "irrelevant"

→ Shifted Heaviside step-function  $H(c_i - \Theta_i)$  (or steep sigmoid threshold?)

$$b_i(t, \Theta_i) = \begin{cases} 1, & c_i(t) > \Theta_i \\ 0, & c_i(t) \leq \Theta_i \end{cases} \quad (7)$$

Temporal averaging by neurons of the central brain

- Finalized set of slowly changing kernel-specific features (one per AN)

- Different species-specific song patterns are characterized by a distinct combination of feature values → Clusters in high-dimensional feature space

→ Lowpass filter 1 Hz

$$f_i(t) = b_i(t) * h_{\text{LP}}(t), \quad f_{\text{cut}} = 1 \text{ Hz} \quad (8)$$

### 3 Two mechanisms driving the emergence of intensity-invariant song representation

**Definition of invariance (general, systemic):**

Invariance = Property of a system to maintain a stable output with respect to a set of relevant input parameters (variation to be represented) but irrespective of one or more other parameters (variation to be discarded) → Selective input-output decorrelation

**Definition of intensity invariance (context of neurons and songs):**

Intensity invariance = Time scale-selective sensitivity to certain faster amplitude dynamics (song waveform, small-scale AM) and simultaneous insensitivity to slower, more sustained amplitude dynamics (transient baseline, large-scale AM, current overall intensity level)

→ Without time scale selectivity, any fully intensity-invariant output will be a flat line

#### 3.1 Logarithmic scaling & spike-frequency adaptation

Envelope  $x_{\text{env}}(t) \xrightarrow{\text{dB}}$  Logarithmic  $x_{\text{dB}}(t) \xrightarrow{h_{\text{HP}}(t)}$  Adapted  $x_{\text{adapt}}(t)$

- Rewrite signal envelope  $x_{\text{env}}(t)$  (Eq. 2) as a synthetic mixture:



- 1) Song signal  $s(t)$  ( $\sigma_s^2 = 1$ ) with variable multiplicative scale  $\alpha \geq 0$
- 2) Fixed-scale additive noise  $\eta(t)$  ( $\sigma_\eta^2 = 1$ )

$$x_{\text{env}}(t) = \alpha \cdot s(t) + \eta(t), \quad x_{\text{env}}(t) > 0 \quad \forall t \in \mathbb{R} \quad (9)$$

- Signal-to-noise ratio (SNR): Ratio of variances of synthetic mixture  $x_{\text{env}}(t)$  with ( $\alpha > 0$ ) and without ( $\alpha = 0$ ) song signal  $s(t)$ , assuming  $s(t) \perp \eta(t)$

$$\text{SNR} = \frac{\sigma_{s+\eta}^2}{\sigma_\eta^2} = \frac{\alpha^2 \cdot \sigma_s^2 + \sigma_\eta^2}{\sigma_\eta^2} = \alpha^2 + 1 \quad (10)$$

### Logarithmic component:

- Simplify decibel transformation (Eq. 3) and apply to synthetic  $x_{\text{env}}(t)$
- Isolate scale  $\alpha$  and reference  $x_{\text{ref}}$  using logarithm product/quotient laws

$$\begin{aligned} x_{\text{dB}}(t) &= \log \frac{\alpha \cdot s(t) + \eta(t)}{x_{\text{ref}}} \\ &= \log \frac{\alpha}{x_{\text{ref}}} + \log b_i g[s(t) + \frac{\eta(t)}{\alpha} b_i g] \end{aligned} \quad (11)$$

- In log-space, a multiplicative scaling factor becomes additive
- Allows for the separation of song signal  $s(t)$  and its scale  $\alpha$
- Introduces scaling of noise term  $\eta(t)$  by the inverse of  $\alpha$
- Normalization by  $x_{\text{ref}}$  applies equally to all terms (no individual effects)

### Adaptation component:

- Highpass filter over  $x_{\text{dB}}(t)$  (Eq. 4) can be approximated as subtraction of the local signal offset within a suitable time interval  $T_{\text{HP}}$  ( $0 \ll T_{\text{HP}} < \frac{1}{f_{\text{cut}}}$ )

$$x_{\text{adapt}}(t) \approx x_{\text{dB}}(t) - \log \frac{\alpha}{x_{\text{ref}}} = \log b_i g[s(t) + \frac{\eta(t)}{\alpha} b_i g] \quad (12)$$

### Implication for intensity invariance:

- Logarithmic scaling is essential for equalizing different song intensities
  - Intensity information can be manipulated more easily when in form of a signal offset in log-space than a multiplicative scale in linear space
- Scale  $\alpha$  can only be redistributed, not entirely eliminated from  $x_{\text{adapt}}(t)$ 
  - Turn initial scaling of song  $s(t)$  by  $\alpha$  into scaling of noise  $\eta(t)$  by  $\frac{1}{\alpha}$
- Capability to compensate for intensity variations, i.e. selective amplification of output  $x_{\text{adapt}}(t)$  relative to input  $x_{\text{env}}(t)$ , is limited by input SNR (Eq. 10):

- $\alpha \gg 1$ : Attenuation of  $\eta(t)$  term  $\rightarrow s(t)$  dominates  $x_{\text{adapt}}(t)$
- $\alpha \approx 1$ : Negligible effect on  $\eta(t)$  term  $\rightarrow x_{\text{adapt}}(t) = \log[s(t) + \eta(t)]$
- $\alpha \ll 1$ : Amplification of  $\eta(t)$  term  $\rightarrow \eta(t)$  dominates  $x_{\text{adapt}}(t)$
- $\rightarrow$  Ability to equalize between different sufficiently large scales of  $s(t)$
- $\rightarrow$  Inability to recover  $s(t)$  when initially masked by noise floor  $\eta(t)$
- Logarithmic scaling emphasizes small amplitudes (song onsets, noise floor)
- $\rightarrow$  Recurring trade-off: Equalizing signal intensity vs preserving initial SNR

### 3.2 Threshold nonlinearity & temporal averaging

Convolved  $c_i(t) \xrightarrow{H(c_i - \Theta_i)}$  Binary  $b_i(t) \xrightarrow{h_{\text{LP}}(t)}$  Feature  $f_i(t)$

#### Thresholding component:

- Within an observed time interval  $T$ ,  $c_i(t)$  follows probability density  $p(c_i, T)$
- Within  $T$ ,  $c_i(t)$  exceeds threshold value  $\Theta_i$  for time  $T_1$  ( $T_1 + T_0 = T$ )
- Threshold  $H(c_i - \Theta_i)$  splits  $p(c_i, T)$  around  $\Theta_i$  in two complementary parts

$$\int_{\Theta_i}^{+\infty} p(c_i, T) dc_i = 1 - \int_{-\infty}^{\Theta_i} p(c_i, T) dc_i = \frac{T_1}{T} \quad (13)$$

$\rightarrow$  Semi-definite integral over right-sided portion of split  $p(c_i, T)$  gives ratio of time  $T_1$  where  $c_i(t) > \Theta_i$  to total time  $T$  due to normalization of  $p(c_i, T)$

$$\int_{-\infty}^{+\infty} p(c_i, T) dc_i = 1 \quad (14)$$

#### Averaging component:

- Lowpass filter over binary response  $b_i(t)$  (Eq. 8) can be approximated as temporal averaging over a suitable time interval  $T_{\text{LP}}$  ( $T_{\text{LP}} > \frac{1}{f_{\text{cut}}}$ )
- Within  $T_{\text{LP}}$ ,  $b_i(t)$  takes a value of 1 ( $c_i(t) > \Theta_i$ ) for time  $T_1$  ( $T_1 + T_0 = T_{\text{LP}}$ )

$$f_i(t) \approx \frac{1}{T_{\text{LP}}} \int_t^{t+T_{\text{LP}}} b_i(\tau) d\tau = \frac{T_1}{T_{\text{LP}}} \quad (15)$$

$\rightarrow$  Temporal averaging over  $b_i(t) \in [0, 1]$  (Eq. 7) gives ratio of time  $T_1$  where  $c_i(t) > \Theta_i$  to total averaging interval  $T_{\text{LP}}$

$\rightarrow$  Feature  $f_i(t)$  approximately represents supra-threshold fraction of  $T_{\text{LP}}$

#### Combined result:

- Feature  $f_i(t)$  can be linked to the distribution of  $c_i(t)$  using Eqs. 13 & 15

$$f_i(t) \approx \int_{\Theta_i}^{+\infty} p(c_i, T_{LP}) dc_i = P(c_i > \Theta_i, T_{LP}) \quad (16)$$

→ Because the integral over a probability density is a cumulative probability, the value of feature  $f_i(t)$  (temporal compression of  $b_i(t)$ ) at every time point  $t$  signifies the probability that convolution output  $c_i(t)$  exceeds the threshold value  $\Theta_i$  during the corresponding averaging interval  $T_{LP}$

### **Implication for intensity invariance:**

- Convolution output  $c_i(t)$  quantifies temporal similarity between amplitudes of template waveform  $k_i(t)$  and signal  $x_{adapt}(t)$  centered at time point  $t$

→ Based on amplitudes on a graded scale

- Feature  $f_i(t)$  quantifies the probability that amplitudes of  $c_i(t)$  exceed threshold value  $\Theta_i$  within interval  $T_{LP}$  around time point  $t$

→ Based on binned amplitudes corresponding to one of two categorical states → Deliberate loss of precise amplitude information

→ Emphasis on temporal structure (ratio of  $T_1$  over  $T_{LP}$ )

- Thresholding of  $c_i(t)$  and subsequent temporal averaging of  $b_i(t)$  to obtain  $f_i(t)$  constitutes a remapping of an amplitude-encoding quantity into a duty cycle-encoding quantity, mediated by threshold function  $H(c_i - \Theta_i)$

- Different scales of  $c_i(t)$  can result in similar  $T_1$  segments depending on the magnitude of the derivative of  $c_i(t)$  in temporal proximity to time points at which  $c_i(t)$  crosses threshold value  $\Theta_i$

→ The steeper the slope of  $c_i(t)$ , the less  $T_1$  changes with scale variations

→ If  $T_1$  is invariant to scale variation in  $c_i(t)$ , then so is  $f_i(t)$

- Suggests a relatively simple rule for optimal choice of threshold value  $\Theta_i$ :

→ Find amplitude  $c_i$  that maximizes absolute derivative of  $c_i(t)$  over time

→ Optimal with respect to intensity invariance of  $f_i(t)$ , not necessarily for other criteria such as song-noise separation or diversity between features

- Nonlinear operations can be used to detach representations from graded physical stimulus (to fasciliate categorical behavioral decision-making?):

1) Capture sufficiently precise amplitude information:  $x_{env}(t)$ ,  $x_{adapt}(t)$

→ Closely following the AM of the acoustic stimulus

2) Quantify relevant stimulus properties on a graded scale:  $c_i(t)$

- More decorrelated representation, compared to prior stages
- 3) Nonlinearity: Distinguish between "relevant vs irrelevant" values:  $b_i(t)$
- Trading a graded scale for two or more categorical states
- 4) Represent stimulus properties under relevance constraint:  $f_i(t)$
- Graded again but highly decorrelated from the acoustic stimulus
- 5) Categorical behavioral decision-making requires further nonlinearities
- Parameters of a behavioral response may be graded (e.g. approach speed), initiation of one behavior over another is categorical (e.g. approach/stay)

## **4 Discriminating species-specific song patterns in feature space**

## **5 Conclusions & outlook**