# "Wow!" Bayesian Surprise for Salient Acoustic Event Detection

Boris Schauerte & Rainer Stiefelhagen

Karlsruhe Institute of Technology

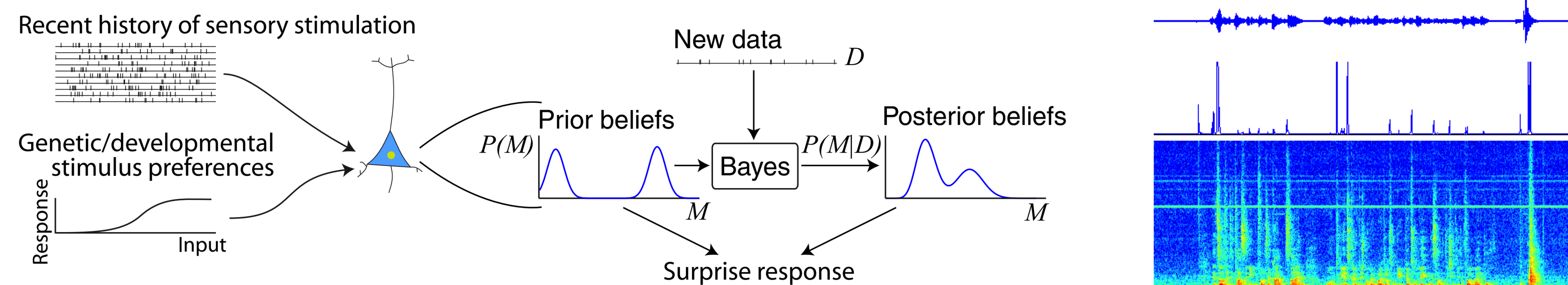{boris.schauerte, rainer.stiefelhagen}@kit.edu

## ABSTRACT

We propose the use of Bayesian surprise to detect arbitrary, salient acoustic events. We use Gaussian or Gamma distributions to model the spectrogram distribution and use the Kullback-Leibler divergence of the posterior and prior distribution to calculate how "unexpected" and thus surprising newly observed audio samples are. This way, we efficiently detect arbitrary surprising/salient acoustic events.

## MOTIVATION

- identify subsets within sensory inputs that are likely to contain important information

- focus complex processing operations on the selected, potentially relevant information

- in general, drastically reduce the computational requirements to process data

- real-time processing and reflex-like reactions despite computational restrictions

## PRINCIPLE

- an observed spectrogram element $G(t,\omega)$ is "surprising" if the updated (using Bayes' rule) distribution $P_{\mathrm{post}}^\omega$ differs significantly from the prior distribution $P_{\mathrm{prior}}^\omega$

$$S_{\mathrm{A}}(t,\omega) = D_{\mathrm{KL}}(P_{\mathrm{post}}^\omega \| P_{\mathrm{prior}}^\omega) \quad (1)$$

$$= \int P_{\mathrm{post}}^\omega \log \frac{P_{\mathrm{post}}^\omega}{P_{\mathrm{prior}}^\omega} dg \quad (2)$$

with Kullback-Leiber divergence $D_{\mathrm{KL}}$

- surprise at time $t$ over all frequencies

$$S_{\mathrm{A}}(t) = \frac{1}{|\Omega|} \sum_{\omega \in \Omega} S_{\mathrm{A}}(t,\omega) \quad (3)$$

- the unit of surprise is called "wow" [4]

## MODELS

**Gaussian:**

$$S_{\mathrm{A}}(t,\omega) = \frac{1}{2}[\log \frac{|\Sigma_{\mathrm{prior}}^\omega|}{|\Sigma_{\mathrm{post}}^\omega|} + \mathrm{Tr}\left[\Sigma_{\mathrm{prior}}^{\omega^{-1}} \Sigma_{\mathrm{post}}^\omega\right] - I_{\mathrm{D}} + (4)$$

$$(\mu_{\mathrm{post}}^\omega - \mu_{\mathrm{prior}}^\omega)^T \Sigma_{\mathrm{prior}}^{\omega^{-1}} (\mu_{\mathrm{post}}^\omega - \mu_{\mathrm{prior}}^\omega)]$$

- mean $\mu$ and variance $\Sigma$ of the data in the considered time window, i.e. history

- advantage: exact closed form solution exists (highly efficient to calculate)

**Gamma:**

$$S_{\mathrm{A}}(t,\omega) = \alpha' \log \frac{\beta}{\beta'} + \log \frac{\Gamma(\alpha')}{\Gamma(\alpha)} \quad (5)$$

$$+\beta' \frac{\alpha}{\beta} + (\alpha - \alpha')\psi(\alpha) \quad (6)$$

- $\alpha, \beta > 0$, and Gamma function $\Gamma$ and Digamma function $\psi$

- advantage: better control over the history using the decay/forgetting factor $0 < \zeta < 1$ and update rule

$$\alpha' = \zeta\alpha + G(t,\omega) \quad (7)$$
$$\beta' = \zeta\beta + 1 \quad (8)$$

## EVALUATION

**Idea:**

- we can not simply observe humans to provide a measure of acoustic saliency

- pragmatic, application-oriented approach: use existing acoustic event detection and classification datasets

- salient acoustic event detection has to suppress "uninteresting" audio data while highlighting potentially relevant and thus salient data segments

CLEAR2007 acoustic event detection **dataset:**

- recordings of meetings in a smart room

- a human user marked and classified (14 classes) acoustic events

- not all events could be classified by the human user (i.e., "unknown" class)

$F_\beta$ score as **evaluation measure:**

$$F_\beta = (1 + \beta^2) \cdot \frac{\mathrm{precision} \cdot \mathrm{recall}}{(\beta^2 \cdot \mathrm{precision}) + \mathrm{recall}} \quad (9)$$

- we want to detect all prominent events, i.e. a high recall is most important

- we can tolerate false positives as long as we achieve a net run-time benefit when focusing subsequent algorithms, i.e. a high precision is of secondary interest

- "$\beta$ times as much importance to recall as precision"

**Results:**

| | $F_1$ | $F_2$ | $F_4$ |
|---|---|---|---|
| STFT + Gamma | 0.7668 | 0.8924 | 0.9665 |
| STCT + Gamma | 0.7658 | 0.8916 | 0.9655 |
| MDCT + Gamma | 0.7644 | 0.8894 | 0.9647 |
| STFT + Gaussian | 0.7604 | 0.8832 | 0.9531 |
| STCT + Gaussian | 0.7612 | 0.8813 | 0.9529 |
| MDCT + Gaussian | 0.7613 | 0.8805 | 0.9538 |

## APPLICATIONS

**Robotics [1]:**

- the robot can efficiently detect, investigate, and react on arbitrary, unexpected - i.e., interesting - events

- focus and make better use of the robot's limited computational ressources

**Intensive Care [2]:**

- use Gaussian surprise to detect (sudden) patient agitation based on facial features

## BIOLOGICAL MOTIVATION

- spectrogram $\sim$ basilar membrane [3]

- surprise $\sim$ early sensory neurons [4]

## CODE?

- Gamma and Gaussian surprise implementation public (BSD license) at http://bit.ly/ZjzXqr

- comes with a ready to go audio example

**References**

[1] B. Schauerte, B. Kühn, K. Kroschel, R. Stiefelhagen, Multimodal Saliency-based Attention for Object-based Scene Analysis, in IROS, 2011

[2] M. Martinez, R. Stiefelhagen, Automated Multi-Camera System for Long Term Behavioral Monitoring in Intensive Care Units, in MVA, 2013

[3] Schnupp J., Nelken I., King A, Auditory Neuroscience, MIT Press, Cambridge, MA, 2011.

[4] L. Itti and P. F. Baldi, Bayesian surprise attracts human attention, in NIPS, 2006.

ICASSP 2013
Vancouver Convention & Exhibition Centre
May 26 - 31, 2013 • Vancouver, Canada