

Learning Sparse Representations for Audiovisual Signals



REDWOOD CENTER
for Theoretical Neuroscience

Gianluca Monaci and Friedrich T. Sommer,

Redwood Center for Theoretical Neuroscience, University of California, Berkeley

Email: {gmonaci,fsommer}@berkeley.edu



Motivation

BACKGROUND

- Audio-visual interactions in the brain at many levels
 - *Perceptual evidence*: McGurk effect, bounce-stream illusion, sound-induced flashing... [1]
- Evidence for **early fusion mechanisms**
 - ERP, MEG, BOLD dynamics studies in humans [2]
 - Anatomical studies in monkeys [3]

SPARSE CODING PARADIGM

- Successful in describing auditory and visual coding
 - A1: Smith&Lewicki 2006 - V1: Olshausen&Field 1996

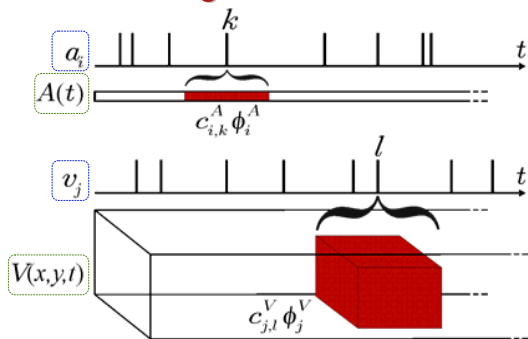
OBSERVATION

- No computational model of early crossmodal interactions

GOAL

- Design a sparse, biologically plausible, audio-visual signal model accounting for early fusion mechanisms in the brain

Audio-visual signal model



audio and visual signals audio and visual spikes audio and visual kernels

$$S = \begin{pmatrix} A(t) \\ V(x, y, t) \end{pmatrix} \approx \begin{pmatrix} \sum_i \sum_k a_{i,k} c_{i,k}^A \phi_i^A(t - k \cdot \tau^A) \\ \sum_j \sum_l v_{j,l} c_{j,l}^V \phi_j^V(t - l \cdot \tau^V) \end{pmatrix}$$

Adaptive crossmodal projections code

Learning sparse codes

- Maximizing log-likelihood

$$O(Z, \Phi, s) = \log(p(Z|\Phi)) = \log(\int p(Z|\Phi, s)p(s)ds)$$
- Approximation: $\int p(Z|\Phi, s)p(s)ds = \int p(Z, s|\Phi)ds \approx p(Z, s^*|\Phi)$
- Iteratively solving two nested steps
 - ▶ **Code**: Find s^* optimizing $O(Z, \Phi, s)$ w.r.t. s
 - ▶ **Learn**: Find Φ^* optimizing $O(Z, \Phi, s^*)$ w.r.t. Φ

Learning sparse audiovisual codes [4]

- Enforce synchronicity: $a = v = s$

$$O = \log(\int p(A, V|\Phi^A, \Phi^V, s)p(s)ds)$$

Learning adaptive crossmodal projections

- Codes a and v are dependent

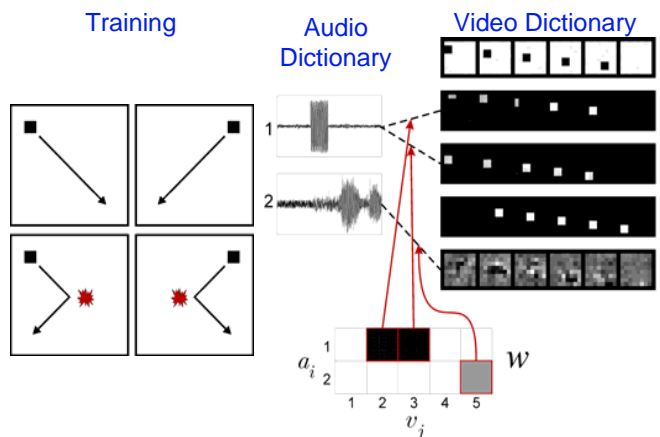
$$p(a, v) \propto \exp(-\lambda_a \|a\|_0 - \lambda_v \|v\|_0 + \sum_i \sum_j a_i w_{ij} v_j)$$

spikes in a # spikes in v crossmodal projections
- $O = \log(\int \int p(A, V|\Phi^A, \Phi^V, a, v)p(a, v)dadv)$

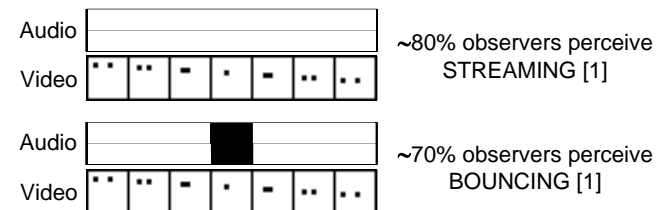
Learning done iteratively solving three nested steps

- ▶ **Code**: greedy approximation, AV-Matching Pursuit [4]
- ▶ **Learn Dictionary**: Gradient Descent on Φ
- ▶ **Learn Projections**: Hebbian learning on w

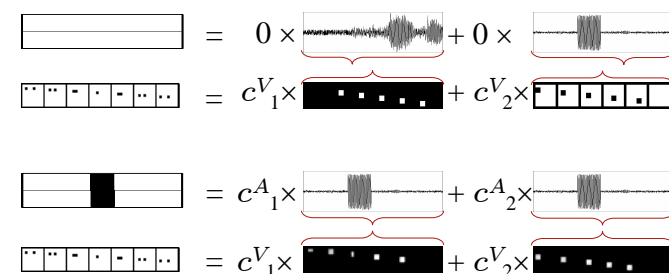
Training



The bounce-stream illusion



Encoding using the learned audiovisual dictionary with adaptive crossmodal projections



Conclusions

- New model for early audio-visual fusion
- Model based on joint sparse coding
- New method to learn basis functions and cross-modal associations
- Model "suffers" from bounce-stream illusion

- [1] Sensory modalities are not separate modalities: plasticity and interactions. Shimoho, S. and Shams, L., Current Opinion in Neurobiology, 2001, 11:505-509
- [2] Sound alters activity in human V1 in association with illusory visual perception. Watkins, S. et al., NeuroImage, 2006, 31:1247-1256
- [3] Multisensory convergence in calcarine visual areas in macaque monkey. Rockland, K.S. and Ojima, H., Int. Journal of Psychophysiology 50:19-26
- [4] Learning sparse generative models of audiovisual signals. Monaci, G. and Sommer, F.T., submitted to European Signal Processing Conference (EUSIPCO), 2008