# AUDIO-VISUAL FUSION AND CONDITIONING WITH NEURAL NETWORKS FOR EVENT RECOGNITION

Mathilde Brousmiche [1] [2]    Jean Rouat [2]    Stéphane Dupont [1]
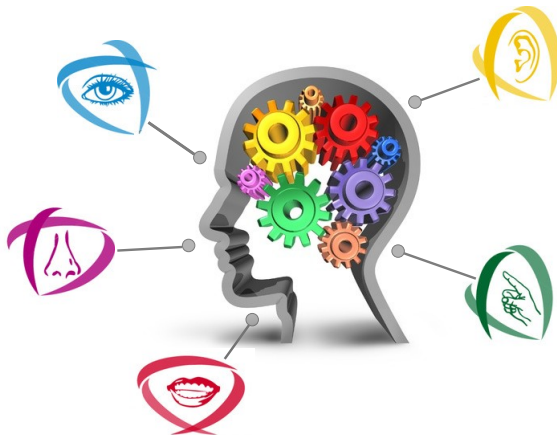
[1]Numediart Institute
University of Mons
Belgium

[2]Necotis Lab
University of Sherbrooke
Canada

MLSP, October 2019
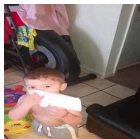
UMONS
University of Mons

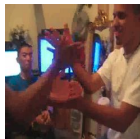UNIVERSITÉ DE
SHERBROOKE

## Multimodality

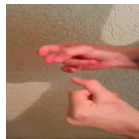# Problem setting : Audio-visual Event Classification
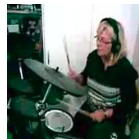
Subset of kinetics[1]:



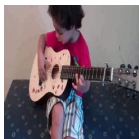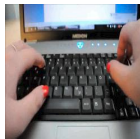blowing_nose    clapping    crying    finger_snapping    playing_drums
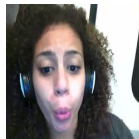
playing_guitar    sneezing    using_computer    whistling    yawning

---

[1] W. Kay et al. "The kinetics human action video dataset". In: *arXiv preprint arXiv:1705.06950* (2017).

# Problem setting : Audio-visual Event Classification

# Problem setting : Audio-visual Event Classification

# Problem setting : Audio-visual Event Classification

**Introduction**
○●

Audio-visual Fusion
○○○

Audio-visual Conditioning
○○○○○○

Conclusion
○○○

# Problem setting : Audio-visual Event Classification



Visual

Audio

Playing guitar

- Fusion ?
- Conditioning ?

Introduction
oo

Audio-visual Fusion
●oo

Audio-visual Conditioning
oooooo

Conclusion
ooo

# Techniques of fusion

Concatenation

Element-wise addition

Multimodal Compact Bilinear pooling (MCB)[1]



---

[1]Y. Gao et al. "Compact bilinear pooling". In: *IEEE Proc. CVPR.* 2016, pp. 317–326.

Introduction
oo

Audio-visual Fusion
o●o

Audio-visual Conditioning
oooooo

Conclusion
ooo

# Fusion levels : Early fusion

Introduction
oo

Audio-visual Fusion
○●○

Audio-visual Conditioning
○○○○○○

Conclusion
○○○

# Fusion levels : Middle fusion

Introduction
oo

Audio-visual Fusion
○●○

Audio-visual Conditioning
○○○○○○

Conclusion
○○○

# Fusion levels : Late fusion

Introduction
oo

**Audio-visual Fusion**
ooo●

Audio-visual Conditioning
oooooo

Conclusion
ooo

# Audio-visual Fusion efficiency

Introduction
oo

Audio-visual Fusion
ooo

Audio-visual Conditioning
●ooooo

Conclusion
ooo

## What is conditioning ?

Introduction
oo

Audio-visual Fusion
ooo

**Audio-visual Conditioning**
●ooooo

Conclusion
ooo

# What is conditioning ?

Introduction
oo

Audio-visual Fusion
ooo

**Audio-visual Conditioning**
o●oooo

Conclusion
ooo

# Modalities conditioning with a attention model[2]

[2]Y. Tian et al. "Audio-visual event localization in unconstrained videos". In: Proc. of ECCV. 2018, pp. 247–263.

# Proposal : Visual feature map modulation with audio information

Introduction
oo

Audio-visual Fusion
ooo

**Audio-visual Conditioning**
ooo●ooo

Conclusion
ooo

# Proposal : Visual feature map modulation with audio information

# Feature-wise Linear Modulation (FiLM)[3]

$\gamma_{i,c}$ and $\beta_{i,c}$ modulate the activations $\mathbf{F}_{i,c}$ :

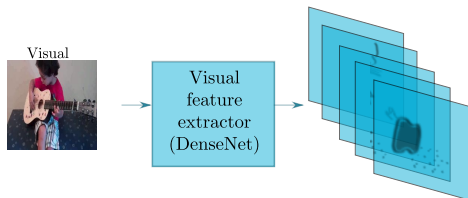$$FiLM(\mathbf{F}_{i,c}|\gamma_{i,c}, \beta_{i,c}) = \gamma_{i,c}\mathbf{F}_{i,c} + \beta_{i,c}$$

where

$$\gamma_{i,c} = f_c(x_i) \qquad \beta_{i,c} = h_c(x_i)$$

$f$ and $h$ can be arbitrary functions

---

[3]E. Perez et al. "Film: Visual reasoning with a general conditioning layer".
In: *Thirty-Second AAAI Conference on Artificial Intelligence*. 2018.

# Contribution of FiLM in audio-visual event classification

| Accuracy [%]            | Image             | Sound             |
|-------------------------|-------------------|-------------------|
| Without FiLM modulation | $61.00 \pm 5.11$  | $66.67 \pm 4.60$  |
| With FiLM modulation    | $75.75 \pm 5.35$  | $75.75 \pm 3.14$  |

Introduction
○○

Audio-visual Fusion
○○○

**Audio-visual Conditioning**
○○○○○○●

Conclusion
○○○

# Better embedding clustering with FiLM



(a) Without FiLM    (b) With FiLM

Introduction
○○

Audio-visual Fusion
○○○

**Audio-visual Conditioning**
○○○○○○●

Conclusion
○○○

# Better embedding clustering with FiLM



(a) Without FiLM  (b) With FiLM

Introduction
oo

Audio-visual Fusion
ooo

**Audio-visual Conditioning**
ooooo●

Conclusion
ooo

# Better embedding clustering with FiLM



(a) Without FiLM    (b) With FiLM

Introduction
○○

Audio-visual Fusion
○○○

**Audio-visual Conditioning**
○○○○○●

Conclusion
○○○

# Better embedding clustering with FiLM



(a) Without FiLM    (b) With FiLM

## Conclusion

- Relevant information for event recognition exists both in visual and audio modalities.
- Exploiting both audio and visual modalities through fusion or conditioning improves event recognition performance
- The use of FiLM layers allows exploiting both audio and visual modalities without an explicit implementation of the fusion

# Future Work

- Test another conditioning method based on multimodal Long Short-Term Memory (LSTM) neural networks
- Analyze the robustness of all methods in the presence of noise as well as in the absence of one modality.

Thank you !