# Learnable pooling with Context Gating for Video Classification

Antoine Miech, Ivan Laptev, Josef Sivic
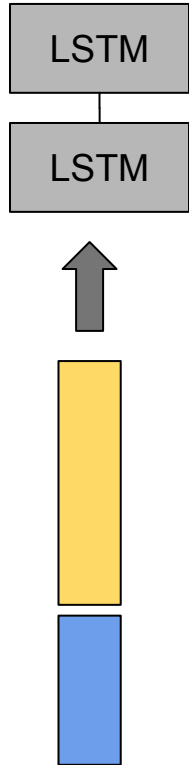
ENS
ÉCOLE NORMALE SUPÉRIEURE

1794
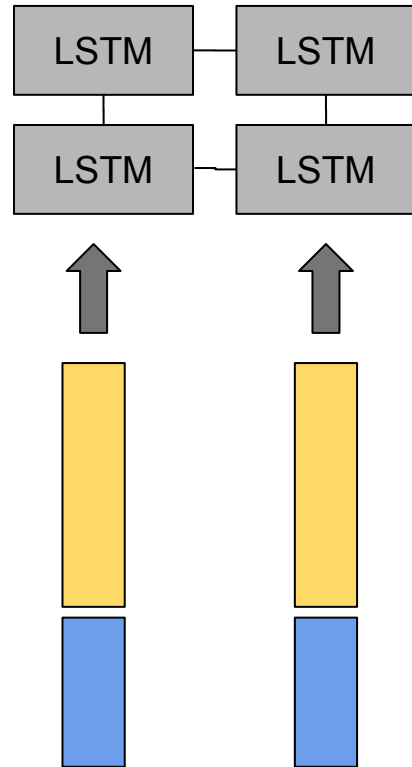
*Inria*
INVENTEURS DU MONDE NUMÉRIQUE

CZECH TECHNICAL UNIVERSITY IN PRAGUE

1

# Goal: Multi-modal features pooling



**POOLING MODULE**

**CLASSIFICATION**
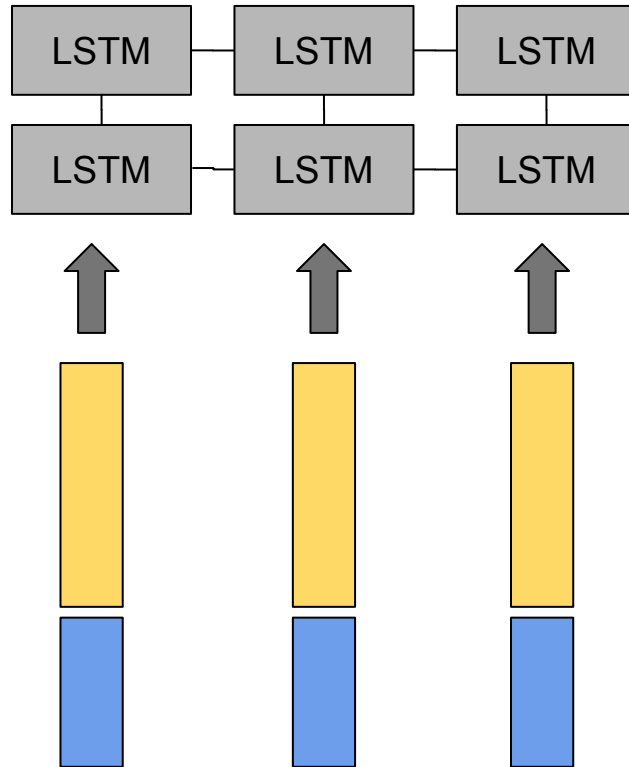
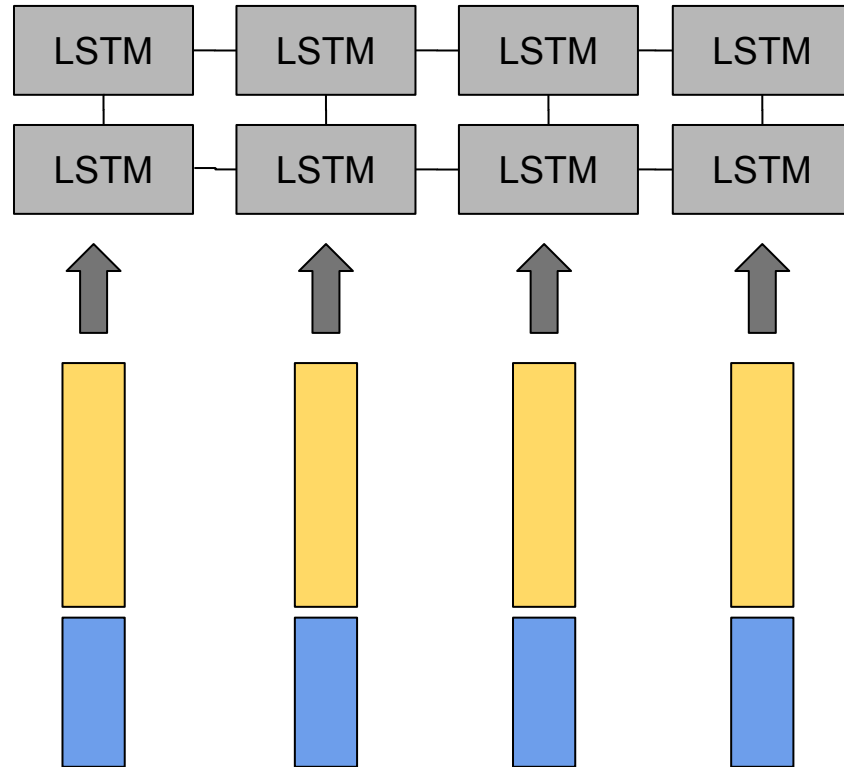# Recurrent model (e.g LSTM)

LSTM
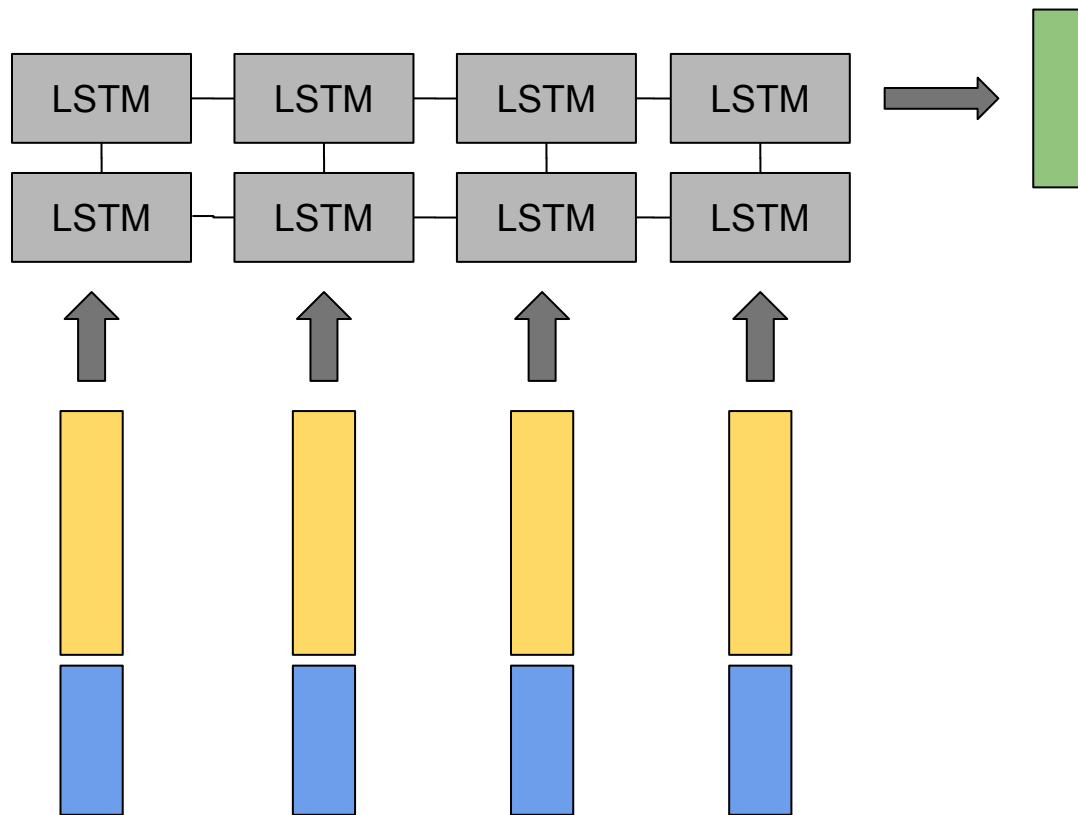
LSTM

# Recurrent model (e.g LSTM)

# Recurrent model (e.g LSTM)
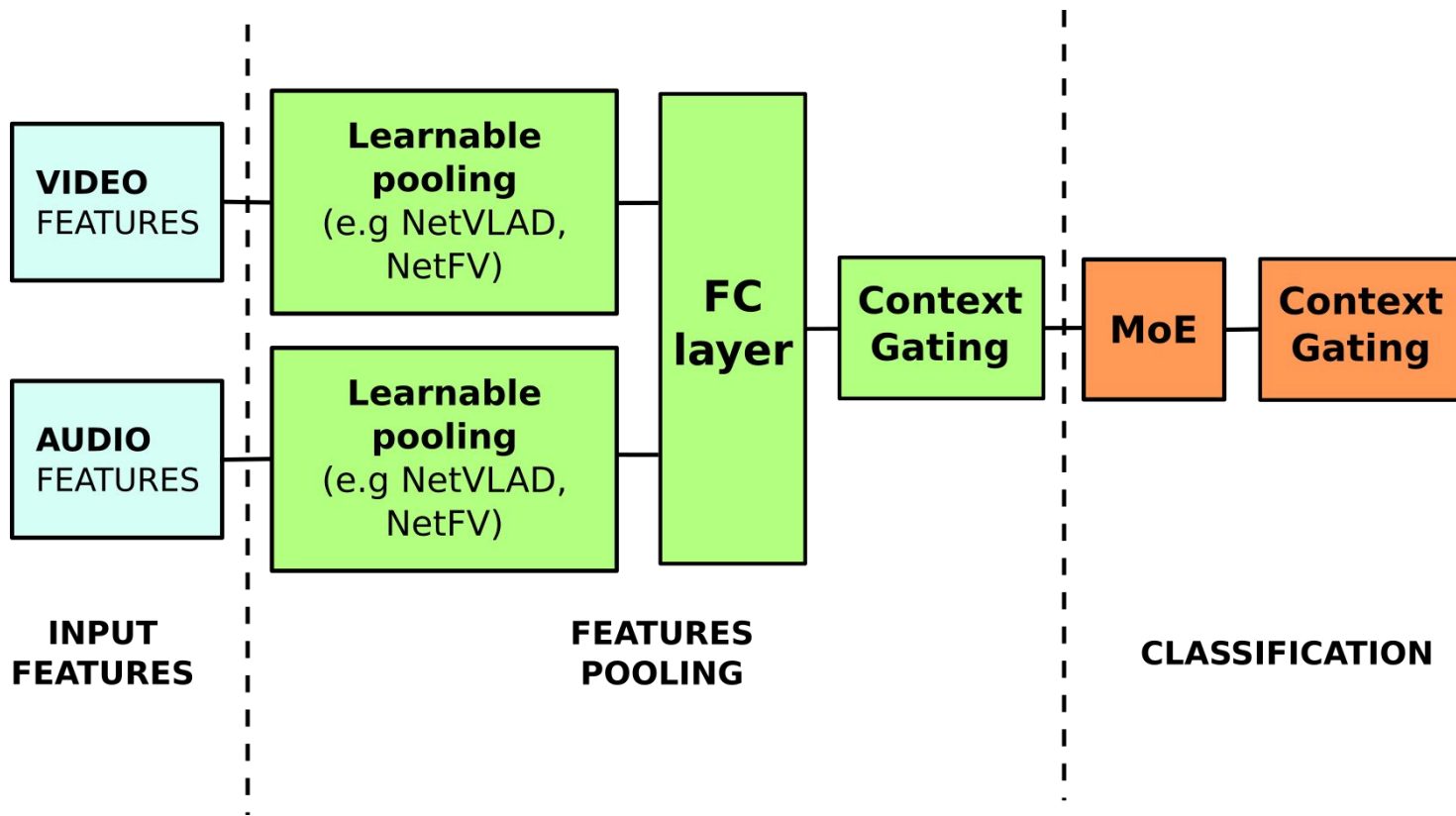
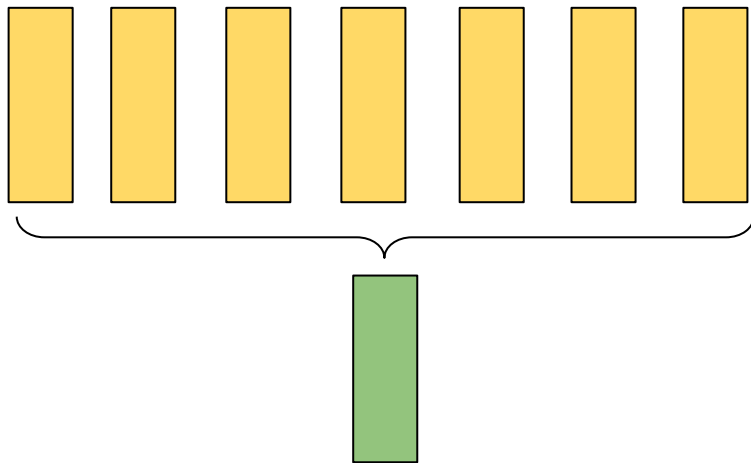# Recurrent model (e.g LSTM)

# Recurrent model (e.g LSTM)

# Problem ?

▷ Slow for inference/training
▷ This is NOT a sequential problem
▷ Needs lots of data for training
▷ How about very long videos ?

But surprisingly good results !

# Our approach

# Traditional pooling

▷ **Bag-of-visual-words** [Sivic and Zisserman, 2003][Csurka et al., 2004]

▷ **VLAD** [Jégou et al., 2010]

▷ **Fisher Vector** [Perronnin et al., 2007]

# Learnable pooling

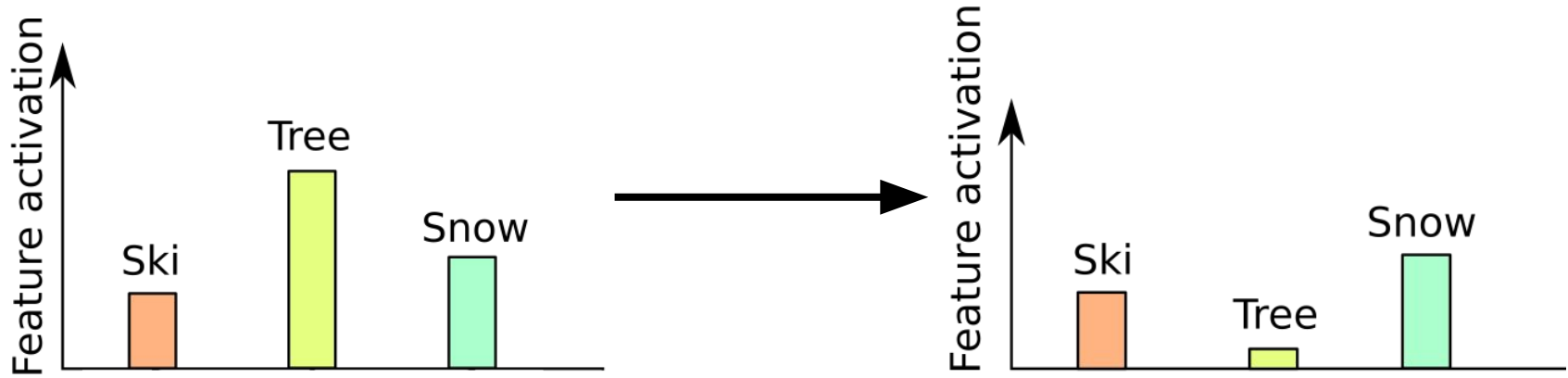| Unsupervised | End-to-End |
|---|---|
| Bag-of-visual-Words [Sivic and Zisserman, 2003] | Soft-DBoW |
| VLAD [Jégou et al., 2010] | NetVLAD [Arandjelović et al. CVPR 2016] |
| Fisher Vector [Perronnin et al., 2007] | NetFV |

# Model overview



**INPUT FEATURES**

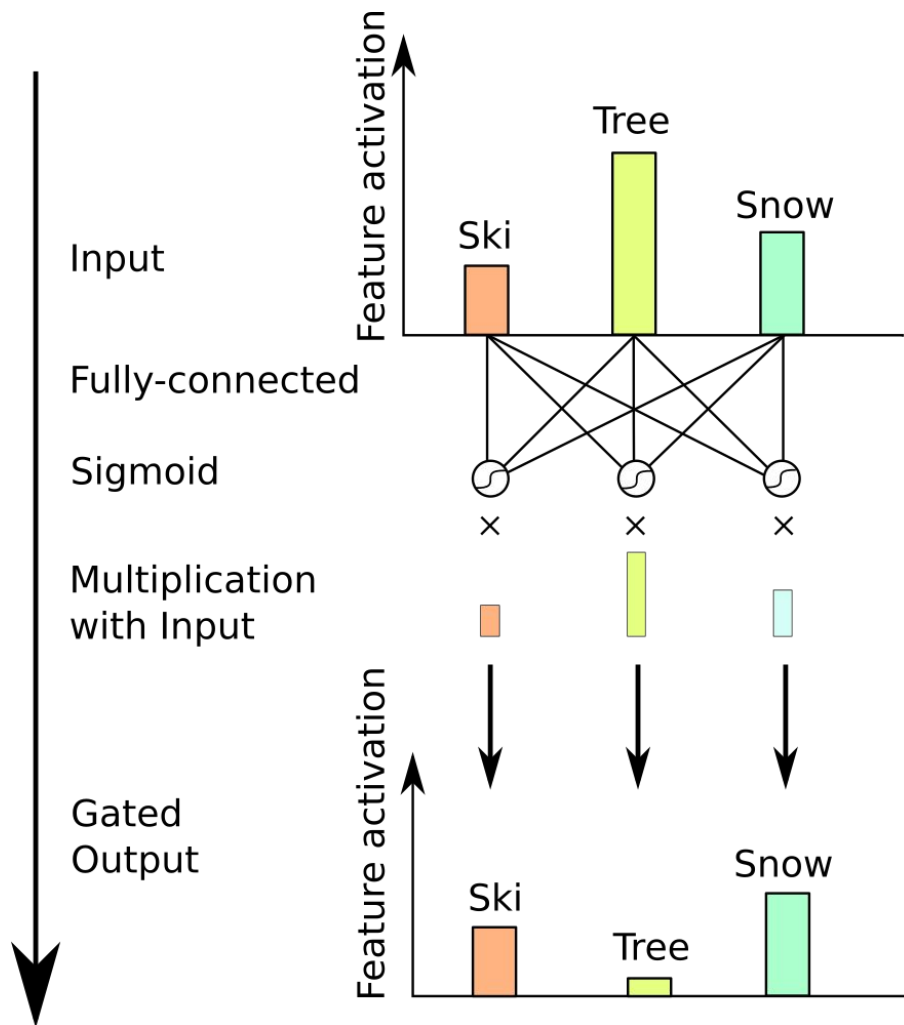**FEATURES POOLING**

**CLASSIFICATION**

# Context Gating

# Context Gating

# Context Gating

Equation:

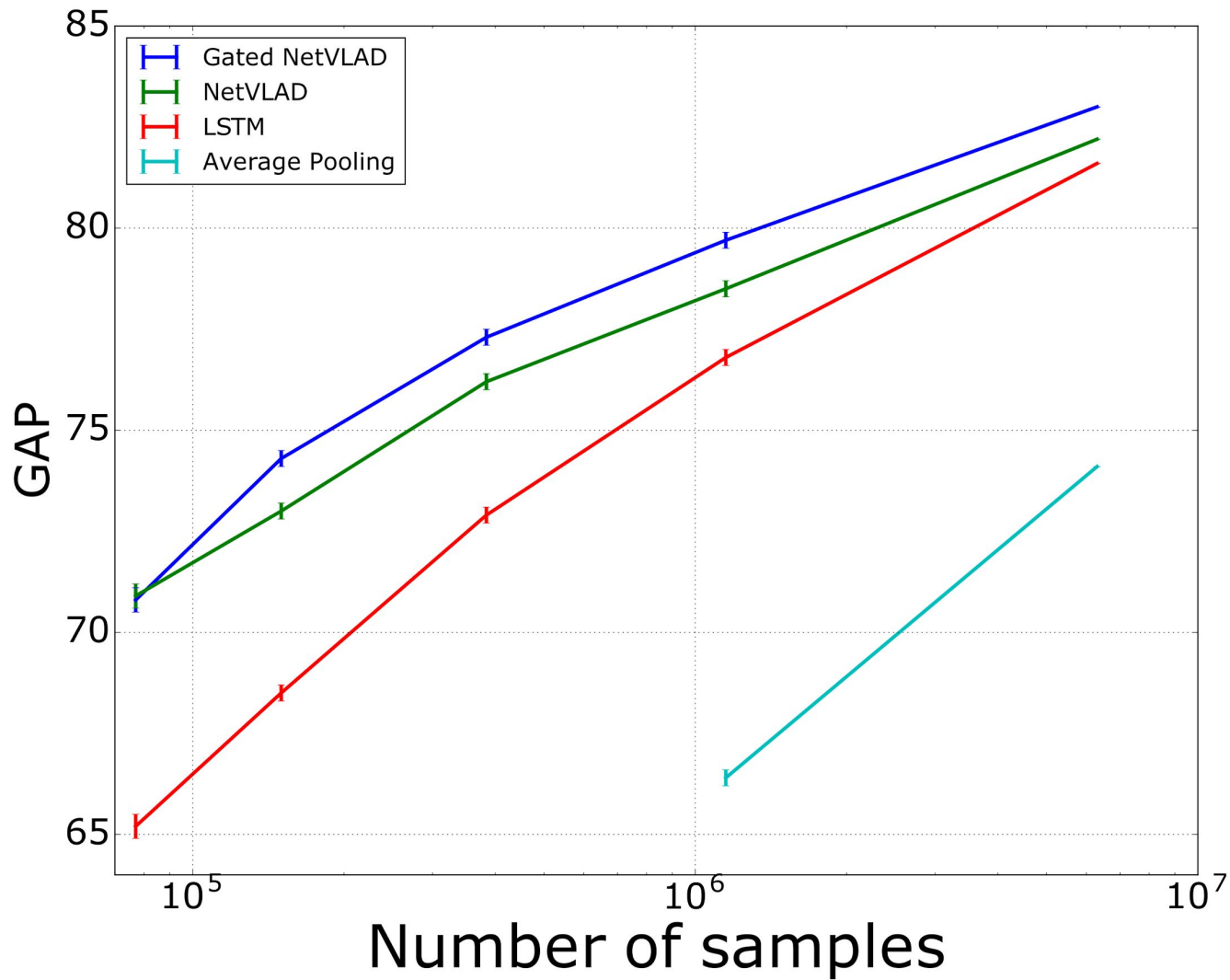$$Y = \sigma(WX + b) \circ X$$

Input

Fully-connected

Sigmoid

Multiplication
with Input

Gated
Output

Feature activation

Ski

Tree

Snow

Feature activation

Ski

Tree

Snow

# Model overview



**VIDEO** FEATURES

**AUDIO** FEATURES

**Learnable pooling** (e.g NetVLAD, NetFV)

**Learnable pooling** (e.g NetVLAD, NetFV)

**FC layer**

**Context Gating**

**MoE**

**Context Gating**

**INPUT FEATURES**

**FEATURES POOLING**

**CLASSIFICATION**

# Results

# Generalization

# Bonus

🏆 Winning the kaggle competition

Effects of ensembling

# LOUPE Tensorflow toolbox

**GITHUB LOUPE repo:**
github.com/antoine77340/LOUPE

**GITHUB Kaggle code repo:**
github.com/antoine77340/Youtube-8M-WILLOW

# *Questions ?*

,,