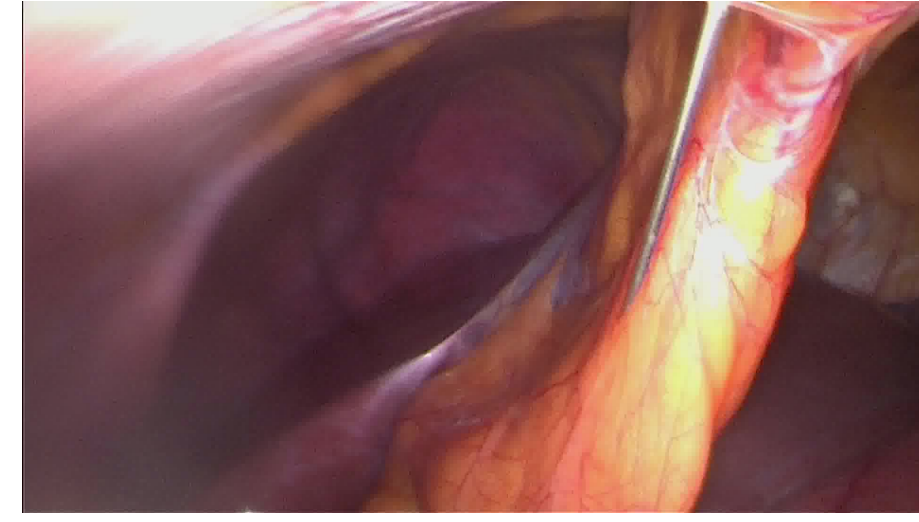


CONTEXT

Goal: Surgical video frames classification

- ▷ Videos of size 1920x1080 Shot at 25 frames per second at IRCAD research center in Strasbourg, France
- ▷ 27 training videos
- ▷ 15 testing videos
- ▷ 8 classes



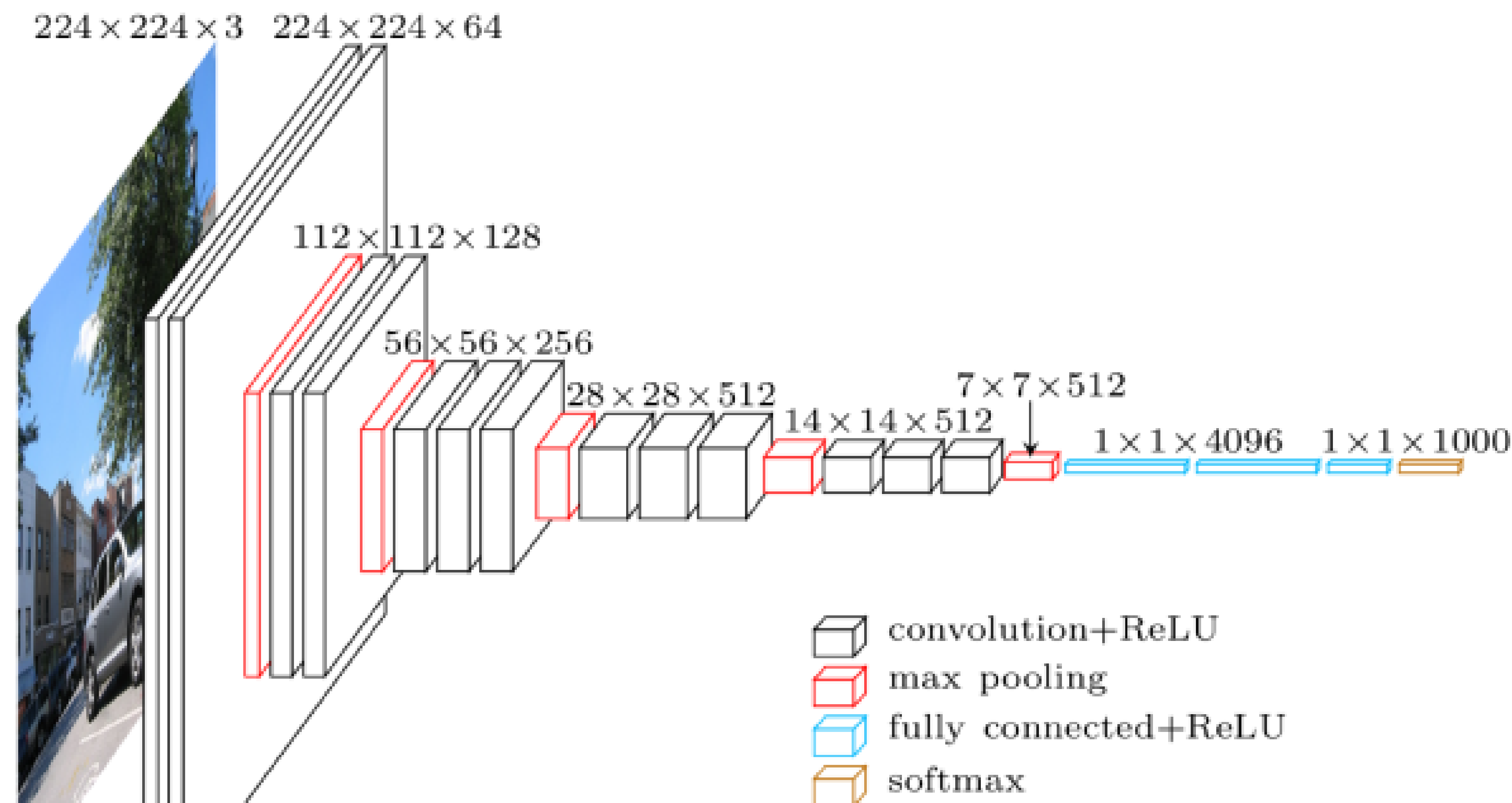
OK for centered object KO for "natural" image

- ▷ Online prediction: $P(y|x_i, x_{i-1}, x_{i-2}, \dots)$
- ▷ Usefull to
 - ▷ Monitor surgeons
 - ▷ Trigger automatic actions

Our approach (usable in production)

- ▷ Firstly, we train a classifier using Deep Learning
- ▷ Secondly, we smooth its predictions
 - ▷ Averaging predictions over last 15 frames
 - ▷ Then, we use a Hidden Markov Model (HMM) as a "denoizer"

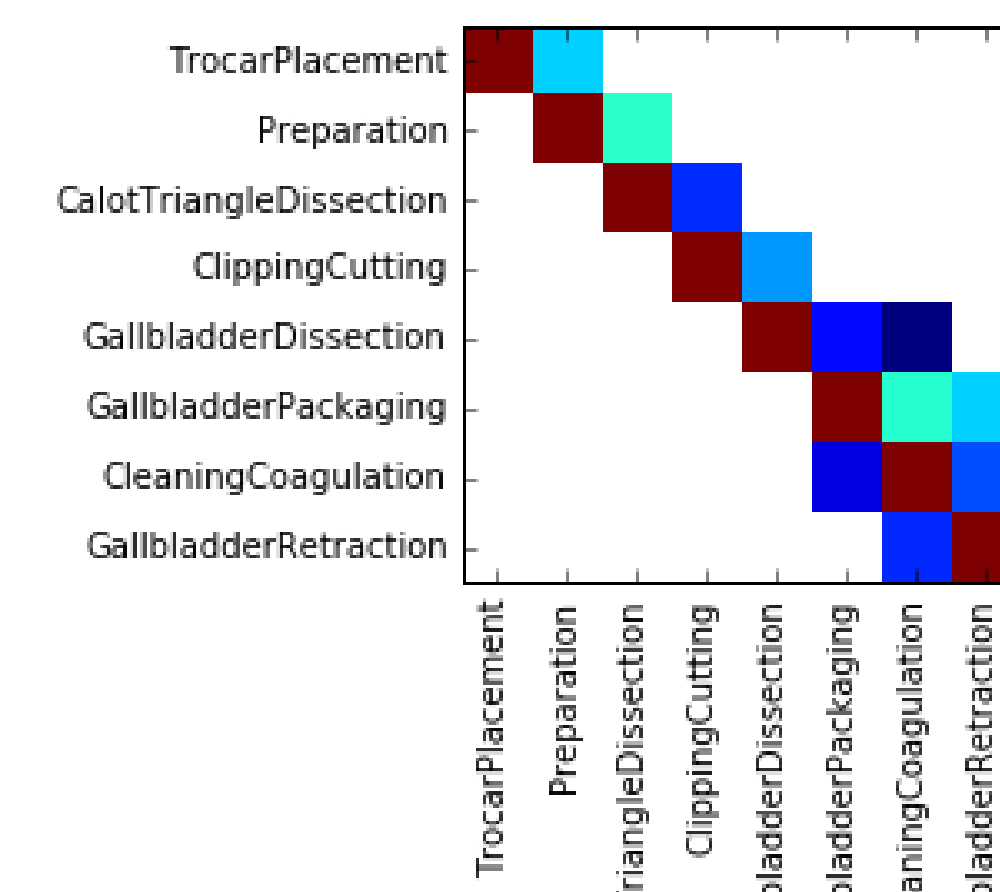
DEEP LEARNING ARCHITECTURE AND METHODS



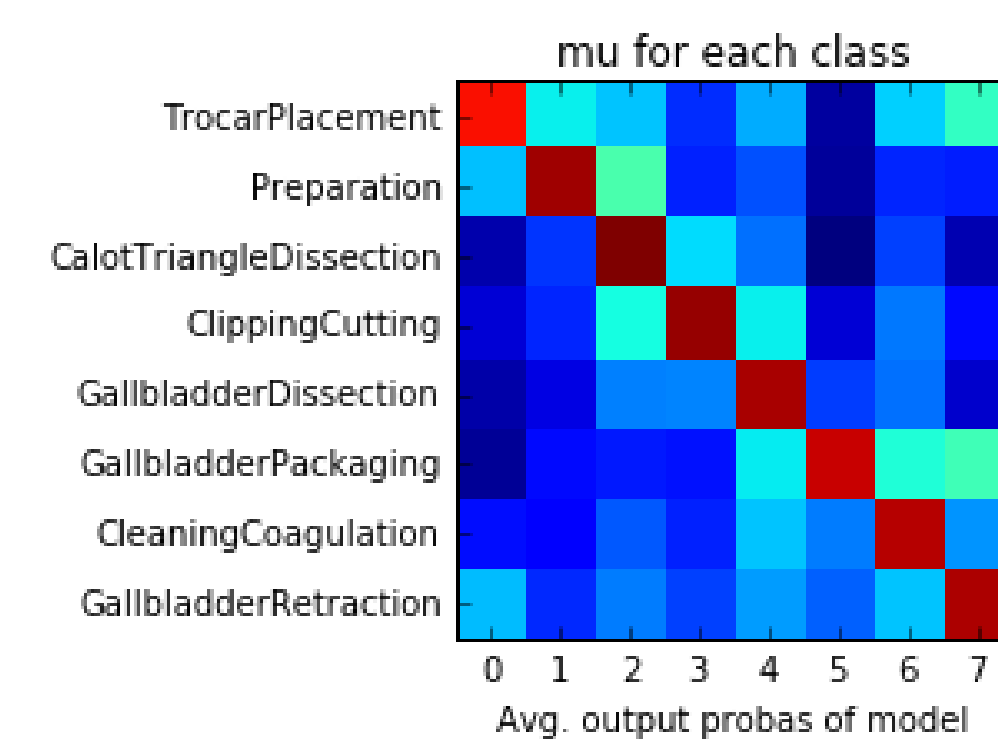
- ▷ Fully connected layer → convolution layer
- ▷ Sliding window approach / shared features
- ▷ Spatial aggregation
- ▷ Object localization prediction

HIDDEN MARKOV MODEL

- ▷ Initial state probabilities
- ▷ Matrix of probabilities of transition between states
- ▷ Gaussian parameters for emissions of observations (mean and co-variance matrix)

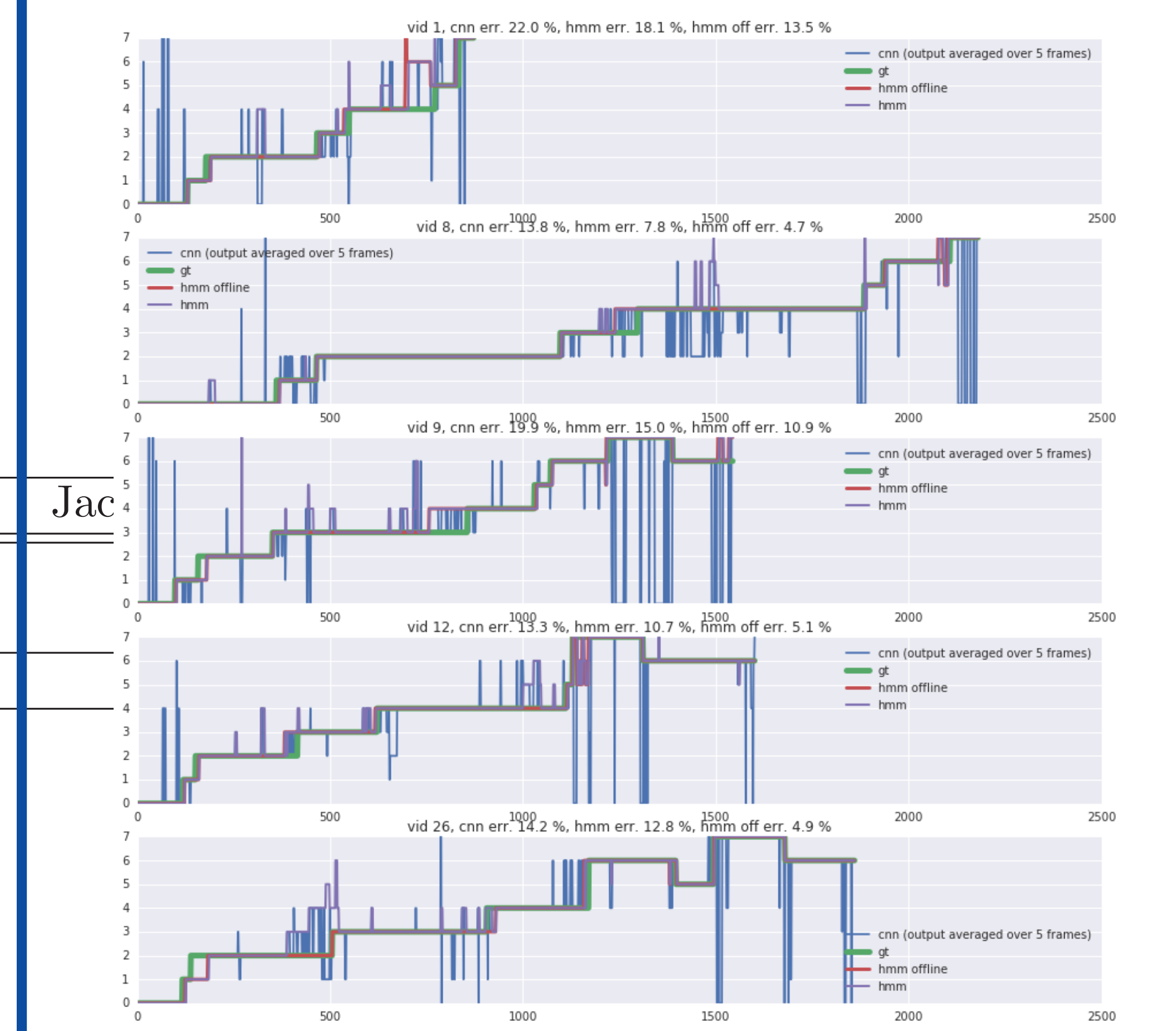


(b)



(c mean)

VISUAL RESULTS



FRAMES CLASSIFIER

Creating validation set by random split:

- ▷ Training set: 22 videos ▷ Validation set: 5 videos {2, 9, 10, 13, 27}

Extracting one frame every 25 frames (1 f/s):

- ▷ Training set: 59,493 images
- ▷ Validation set: 8,062 images
- ▷ Testing set: 28,732 images
- 3 approaches using Deep Learning:**
 - ▷ Training CNN From Scratch
 - ▷ Good: learning representations for this task
 - ▷ Bad: maybe not enough data + lot of hyperparameters
 - ▷ Pre-trained CNN as Features Extractor + SVM
 - ▷ Good: using (good) representations learned on ImageNet
 - ▷ Bad: large semantic gab between ImageNet and this task
 - ▷ Fine Tuning a pre-trained CNN
 - ▷ Good: adapting the representations learned on ImageNet for this task

EXPERIMENTS

| Model | Type | Accuracy Test (%) | Accuracy Val (%) | Jaccard Val |
|------------------|---|-------------------------|------------------|--------------|
| InceptionV3 | Extraction (repres. of ImageNet) | Avg Smoothing 60.53 | 85.97 | 74.67 |
| InceptionV3 | From Scratch (repres. of M2CAI) | Avg + HMM Online 69.13 | 88.90 | 81.60 |
| InceptionV3 | Fine-tuning (both representations) | Avg + HMM Offline 79.06 | 93.47 | 87.59 |
| ResNet200 | Fine-tuning (both representations) | 79.24 | | |

- [1] Ouab et al. Is object localization for free? *CVPR*, 2015.
- [2] Durand et al. MANTRA. *ICCV*, 2015.
- [3] Parizi et al. Automatic discovery of parts. *ICLR*, 2015.
- [4] Gong et al. Multi-scale orderless pooling. *ECCV*, 2014.

CONCLUSION

- ▷ Fine Tuning most accurate approach ▷ HMM is usefull to smooth the predictions ▷ Production ready !
- ▷ Future Work:
 - ▷ Fine Tuning CNN on full train-set
 - ▷ Ensembling several fine tuned CNNs