

Diana Tulpene - Debunking Myths of Mental Illness and Behavioral Issues from Online Communities (University of Ottawa)

- NLP-based user modeling → Big 5 model: (vs. message level)
  - (1) extraversion, (2) emotional stability, (3) agreeableness, (4) conscientiousness, (5) openness to experiences
- Social media provides unique info source! → 'natural setting'
- Mental disorders: insomnia, disorders, posttraumatic disorder, PTSD → Debunked!
  - a lot of old school NLP → but more recently deep learning end-to-end
- Datasets:
  - 'Bell Let's Talk' (Twitter)
  - UPsych 2017 shared dataset (Twitter)
  - Georgeform Dataset (Reddit posts)
- Problem of Deep Learning: Little data on user-level
  - ↳ sleek features → Multi-Task Deep Learning (different mental disorders)
  - ↳ Train models independently and extract layer by layer!
- FWT → Embedding → Lm → predict/act // population-based LSTM
 

```

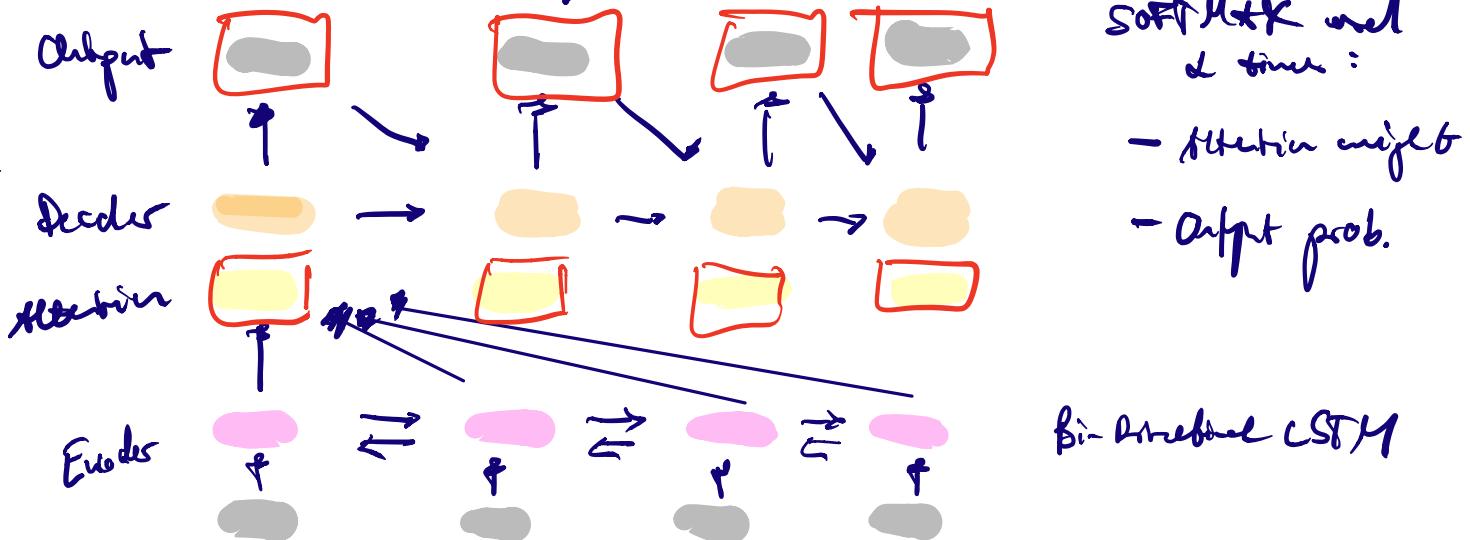
graph TD
    FWT[FWT] --> Embedding[Embedding]
    Embedding --> Lm[Lm]
    Lm --> predictAct[predict/act // population-based LSTM]
    
```

  - ↳ word level
  - ↳ character level
  - ↳ Multi-Lingual
  - often combining tools can improve performance for all tools!
  - ↳ form of delexicalization as in FLTR case! More labels per point!
- Compared different pre-trained embedding techniques! → How (F) end-to-end optimizes with message level models are hard due to noise / little data!
  - ↳ but: Ensures user privacy setting! // Population-based: Facilitates to analyze / extract insights for policy devs!

# Vlad Niculae - Sparse Sequence-to-Squence Models

- S-to-S w. attention  $\rightarrow$  Bahdanau et al 2015

- Architecture classically:



- Usually beam search used to efficiently reduce spending at form of post-processing.
  - HERE  $\rightarrow$  Roberts and Way 2019  $\Rightarrow$  Sparse Attention (lights/flights)
- $\hookrightarrow$  Softmax:  $p_i = \exp z_i / \sum \exp z_i$   $\Rightarrow$  Complex: Discrete probability over choices
- $\rightarrow$  minimizes expected dev + entropy
- $\rightarrow$  Generalize different entropy measure!

$\Delta$   $\rightarrow$  Spherical caps to set

$$D_{\alpha}(z) = \min_{p \in \Delta} p^T z + H(p)$$

$$\hookrightarrow H_0(p) = 0 \rightarrow \text{softmax}$$

$$\hookrightarrow H^S(p) = - \sum_j p_j \log p_j \rightarrow \text{softmax}$$

$$\hookrightarrow H^G(p) = \frac{1}{2} \sum_j p_j (1 - p_j) \rightarrow \text{soft entropy} \rightarrow P & W \text{ diverg}$$

$$\hookrightarrow \text{Here: } H_{\alpha}^T(\phi) = \frac{1}{\alpha(\alpha-1)} \sum_j (\phi_j - \phi_j^{\alpha}) \rightarrow \alpha=1: \text{softmax} \quad \alpha=\infty: \text{discrete}$$

$$\hookrightarrow \text{if } \alpha \neq 1 \text{ else } H^S(p)$$

- Grid search:  $\alpha$  for expected  $\rightarrow$  learn the other

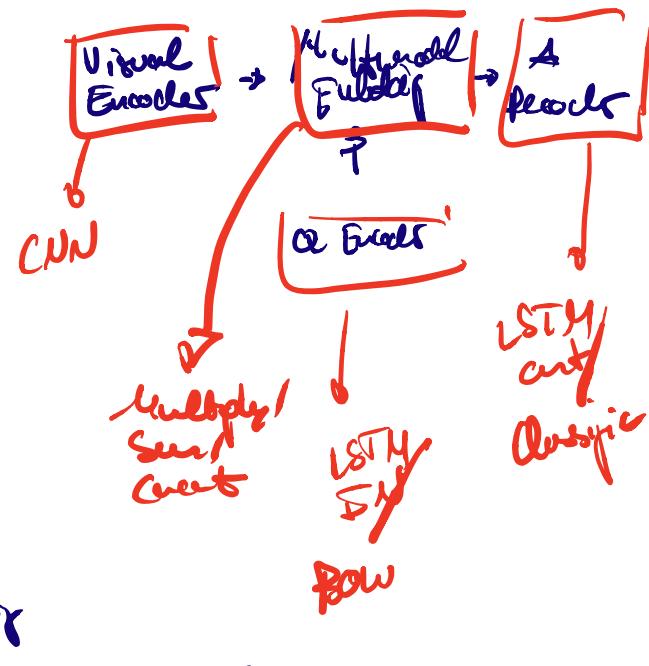
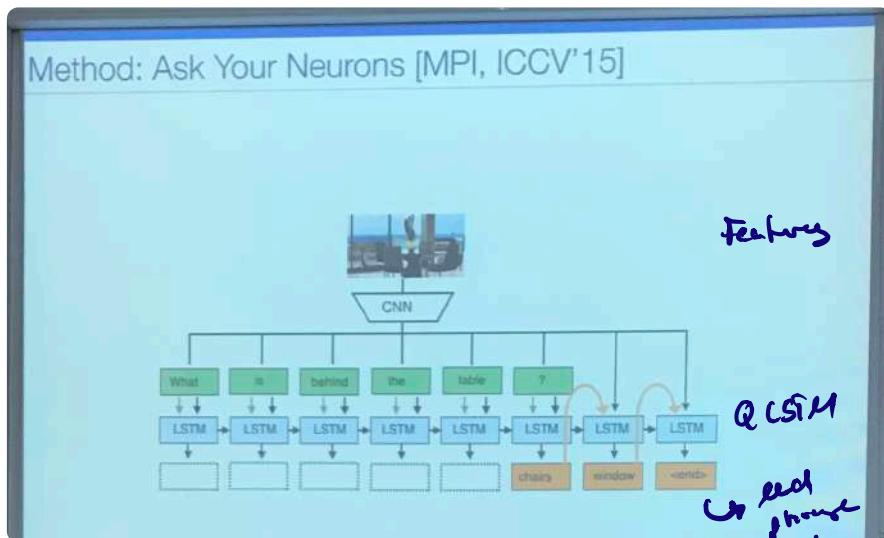
$\hookrightarrow$  try to learn end-to-end!  $\rightarrow$  might be better to take care of beam size!

# LTHL Conference - Belovest - Pg 2

(DeepML)

Matusz & Malinowski - Grounded Language: Answer Qs about Images + Questions

- Visual Turing Test - 1 image  $\leftrightarrow$  many Qs  $\rightarrow$  multi-modal Task
- Visual Models!  $\rightarrow$  Downstream behavior vs. learned representations
  - $\hookrightarrow$  Have to test algo on both in order to achieve better understanding
- Grounding language into percepts  $\rightarrow$  association! Motivation
- S. Gershman  $\rightarrow$  (Fitts' Similarity in Humans and Machines)
- looking different Qs  $\rightarrow$  What is in image?  $\rightarrow$  Object Recognition
  - $\hookrightarrow$  How many or how?  $\Rightarrow$   $\oplus$  Category
- Architecture design:



- Multimodal Functionality: Mapping of all mod.  $\rightarrow$  vector  $\rightarrow$  unimodal signal

$\hookrightarrow$  Concatenation: Form of Integration!

$$\text{Hadamard: } f = \phi^T (U^{(c)})^T \odot V^{(c)} y$$

$\oplus$  Sum:  $FIM$

1st Modality      2nd Modality  $\rightarrow$  Qs

$\mapsto$  Relation  $\rightarrow$  bilinear pooling

$$x^T U^{(c)} V^{(c)} y = 1^T (U^{(c)T} \odot V^{(c)}) y$$

$\hookrightarrow$  Quadratic  $\leftrightarrow$  Hadamard

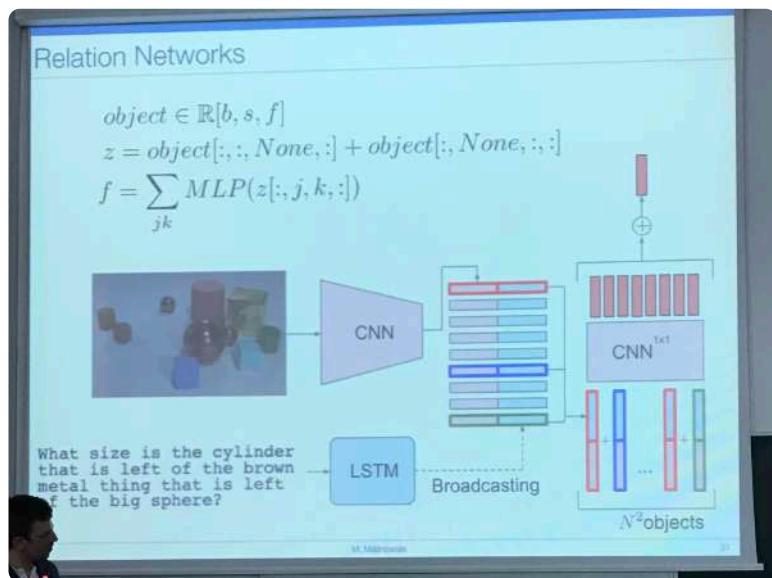
via bilinear  $\hookrightarrow$

$\hookrightarrow$  Can lead to class. red.

- Could theoretically add auditory percepts to  $\times$  LSTM streams
- Usually only performs the vision component but not the embedding  
↳ otherwise vision component feeds to our  $f$ !

• Data: CLEVR  $\rightarrow$  shapes  $\begin{pmatrix} \text{WT} \\ \text{O} \end{pmatrix} \rightarrow$  Facebook  $\Rightarrow$  Searched Soln!

### • Relation Networks:



→ Problem of combinatoriality  $\Rightarrow N^2$  objects

↳ Reduce to  $k^2$  subset

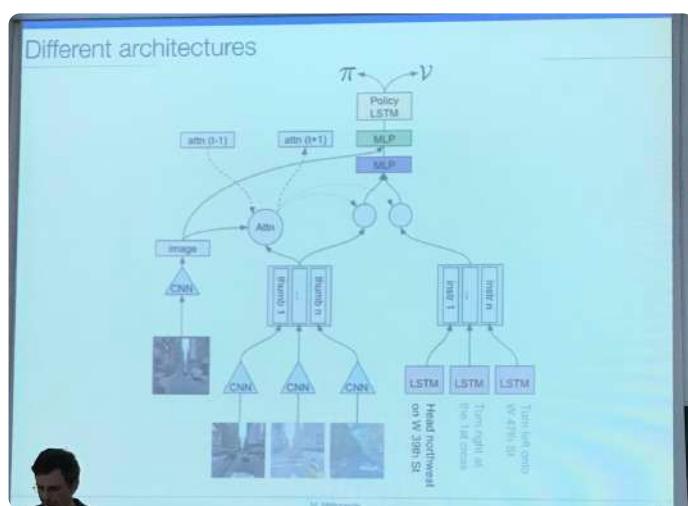
||

(Hard Attention)

↳ Sutskever et al (2018) - Relational Reasoning

↳ Part need to attend to full object

- Want to work with ambiguous naturalistic images and not abstract 3D block renderings!  
↳ Map navigation based on instructions  $\rightarrow$  Street View + Google Maps  
↳ Easy robot on different city environments!



- Different several steps

↳ step-by-step, Decentral, Only local

→ agent has to close relevant network!

Skills  
need fit.

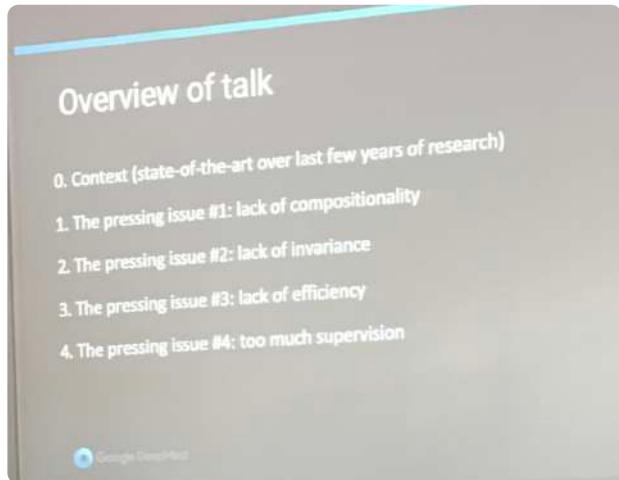
→ Local steps  
w. intermodular steps

- Supervised vs RL goals

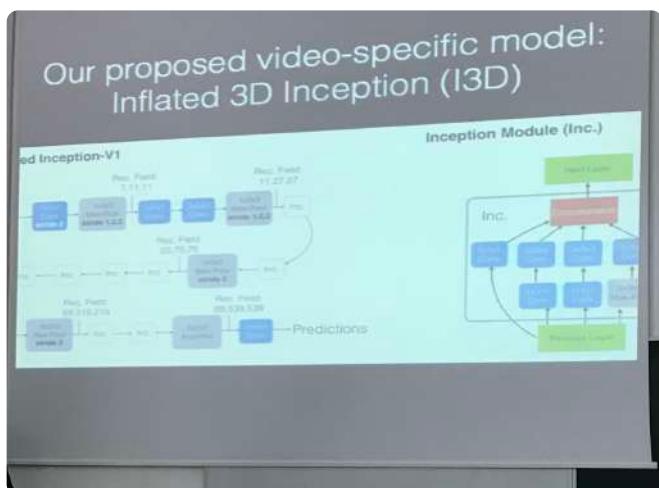
- Variables vs Options?  $\rightarrow$  Hierarchy in layers!

# Joao Carreira - 'Fixing problems in video understanding'

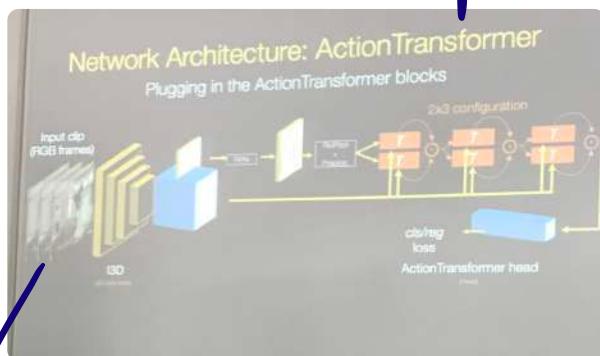
- Initial video understanding → mostly visual! Not visual!
- Social - IQ datasets AFR long → solved video understanding
  - Pascal VOC / Coco → segmentation test! datasets



- Kinetics dataset (800k videos) → aim to be 'fingerset for video'  
↳ 700 classes → 680k videos 2019!



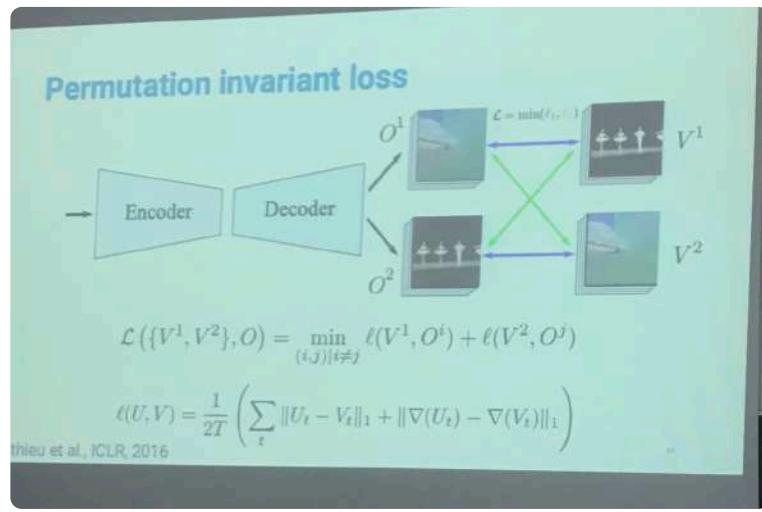
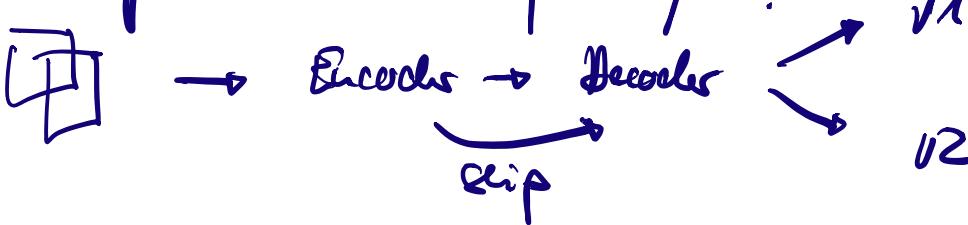
- ① Compositionality: Problem → Multiple objects + multiple actions  
→ Action Transformer Blocks
  - ↳ requires all objects to interact!



- ↳ what do transformers do?  
↳ Map layer-specific heads to 3D  
⇒ look at RGB images  
↳ Detect people // Track people

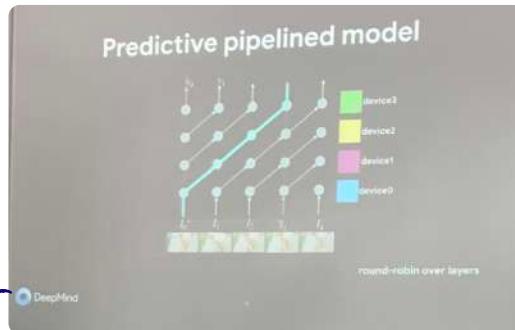
detecting full video!

- ② Invariance → Get rid of nuisance features (shadows, etc.)  
 → Layer decomposition of invariance! → Reflection/Decomposition = Lin. Comb. Ftrs  
 → Usually optimization & what is good way to learn? → Self-Supervised  
 ↳ merge two videos in pixel space!



- Works on synthetic data  
 ↳ How to transfer to the wild?
- Get more invariance to features of naturalistic images!  
 ↳ Go to ocean!
- Test different state-of-the-art methods  
 ↳ Not only linear last patches or GAN generated

- ③ Efficiency ⇒ Fast + Online (no buffer!) processing  
 → Train = Massively parallel ⇒ Not separable  
 ↳ Reptile backtracking necessary ⇒ Improve throughput by processing on different devices!  
 ↳ Go parallelize to improve latency issues

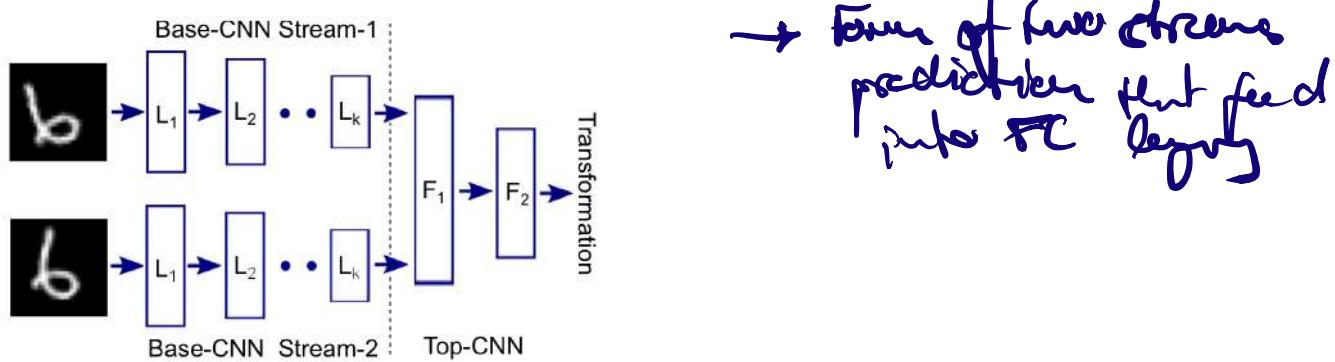


↳ Parallelize layers!

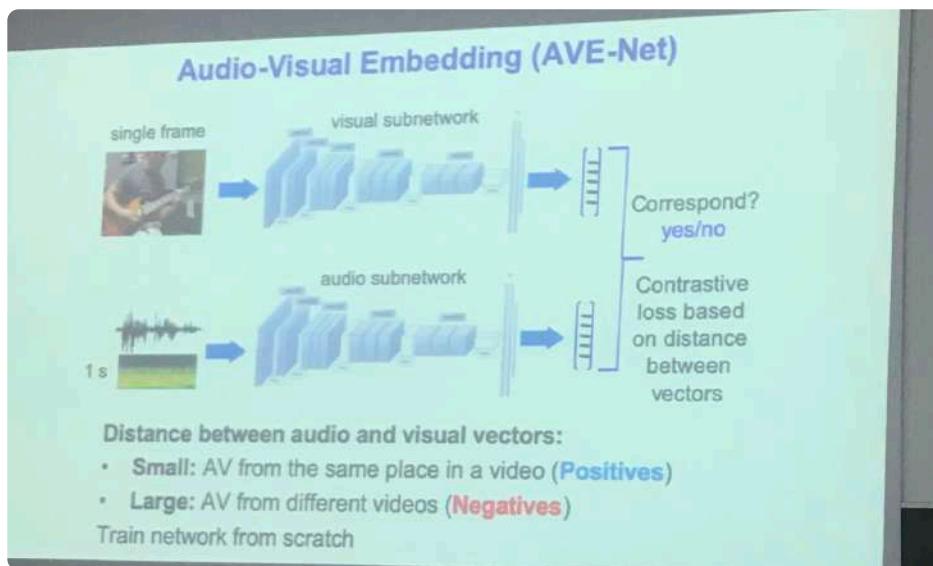
- ⇒ Trade-off between latency & accuracy!  
 ↳ Fully Separable ⇒ Offline & off-line work!

## ④ Human Supervision $\Rightarrow$ Self-Supervised Learning

- Video passive  $\Rightarrow$  device infers perception  $\rightarrow$  Not attorney and!
- J.J. Gibson - 'We move in order to see and we see in order to move.'
- $\hookrightarrow$  Visual features should predict motion!



→ Form of two streams  
predicted that feed  
into FC layer



→ With audio  
 $\Rightarrow$  Objects that sand  
?

Zissmann, Welt  
↓  
also by traps take to  
negatives! ?  
Multi-learning  $\Leftrightarrow$  few  
models

→ For NLP self-supervised is already better than supervised + BERT

- General Idea: Train in self-supervised fashion on massive data  
↳ perform few shot learning or supervised dataset  
↳ Is this the future? Bayesian Model  $\rightarrow$  feature space?!
- Putting all 4 developments together?

# ICML Conference - Bucharest - Day 3

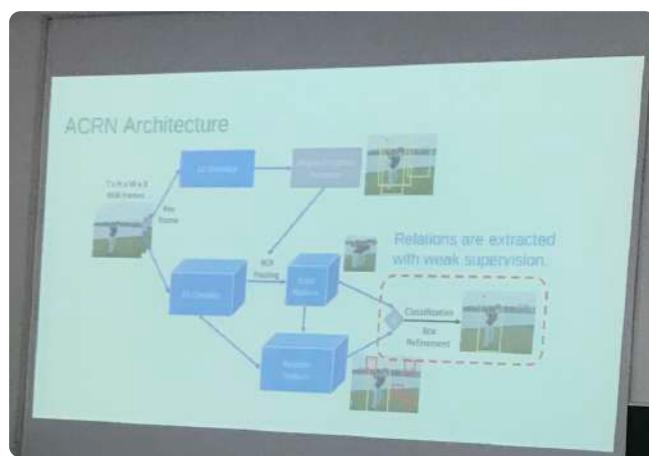
## Rahul Srikantiahar - Recent Research in Machine Perception

- NN architectures ① Labeled data ② Infrastructure
- Comp. photography → High dynamic range imaging
- Privacy hiding and blurring on YouTube
- Connective with @ Google #



### Large-Scale Video Annotation for YouTube #

- Importance of Meta-Data
- Transfer from ET descriptors to YouTube
- Domain adaptation



### Actor-Centric Relation

Defined for action detection,

- Faster R-CNN for object detection segmentation
  - ↳ Need more spatio-temporal context
  - ↳ Iter: Focus on actors + other objects over time + space
  - ↳ Graph ConvNet → we are working straight edges associated with specific actors

### Pixel Sharpness ↴

### Open Scene datasets

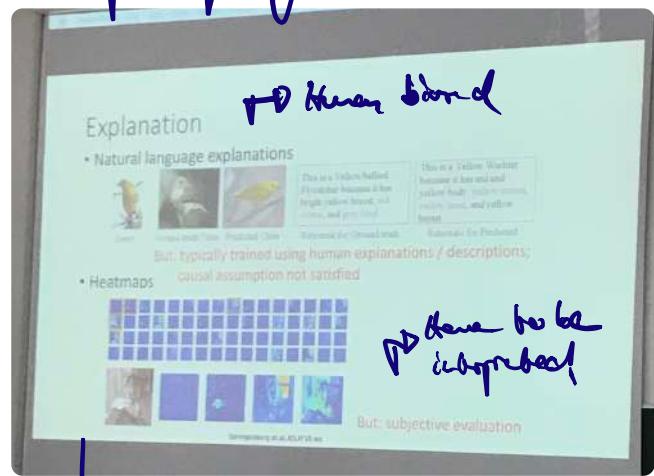
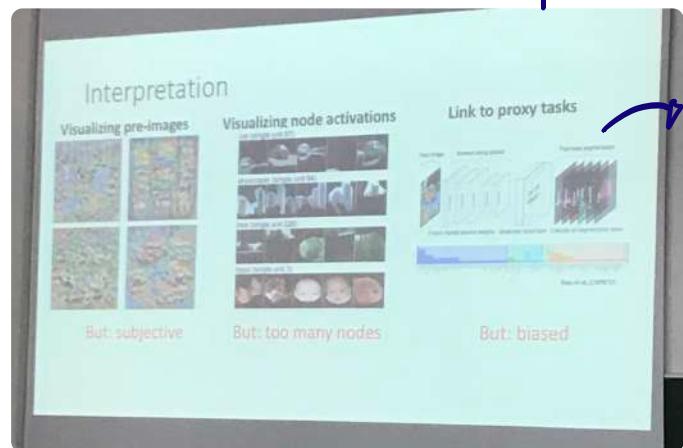
- TF Object Detection API
- Open Images vs
- Audio #

### Weakly/Self-Supervised learning on YouTube DB,

- Form of view synthesis / scene depth + pose of camera
- LDI - layered depth image
- NN = differentiable measure ↳
- Monocular Challenge → Saliency + depth cues

# Tinne Tuytelaars - Explanation by Interpretability

- Interpretability is the degree to which a human can understand the cause of a decision' - Miller, 2018
  - subjective → Need to adapt to the recipient of explanation
  - explain vs. understand → Teacher might be able to understand a concept but not be good at explaining it → Lai et al., 2018
- EU guidelines → predictive systems must provide explanation of internal logic.
  - ↳ Ease of interpretability → White boxes / simple / linear logic
  - INTERPRETATION: explain model → interpretable
  - EXPLANATION: explain decision → input-specific



Problems: ① Too many arrows

↳ Which neurons? ..

↳ Pictures # visualizations

via guided background + guided output

↳ Select top h-things with highest response relevant features

↳ proxy task: Zhou et al  
CVPR 16

↳ user study: Feuer & Lampert  
ECCV 14

↳ visualize only face?

② Visualizing heatmaps ↳ focused BP not for class labeled but for relevant features

Evaluation:

↳ roughly balanced to what users expect! → Go to synthetic dataset!

↳ Synthetic & Have access to ground truth labels! → Can compare heatmaps