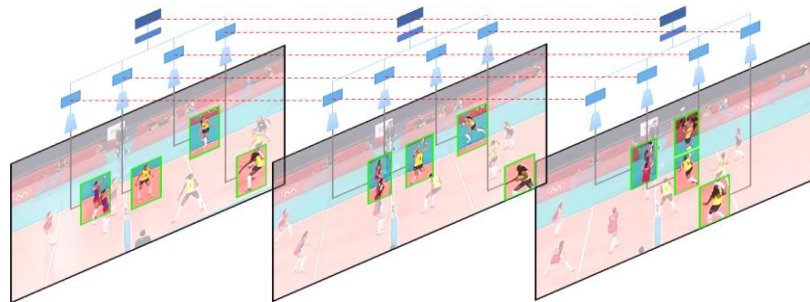# Deep Structured Models for Group Activities and Label Hierarchies

Greg Mori

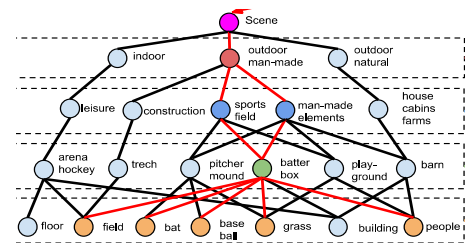Simon Fraser University

# Outline



- Temporal structured models for group activities
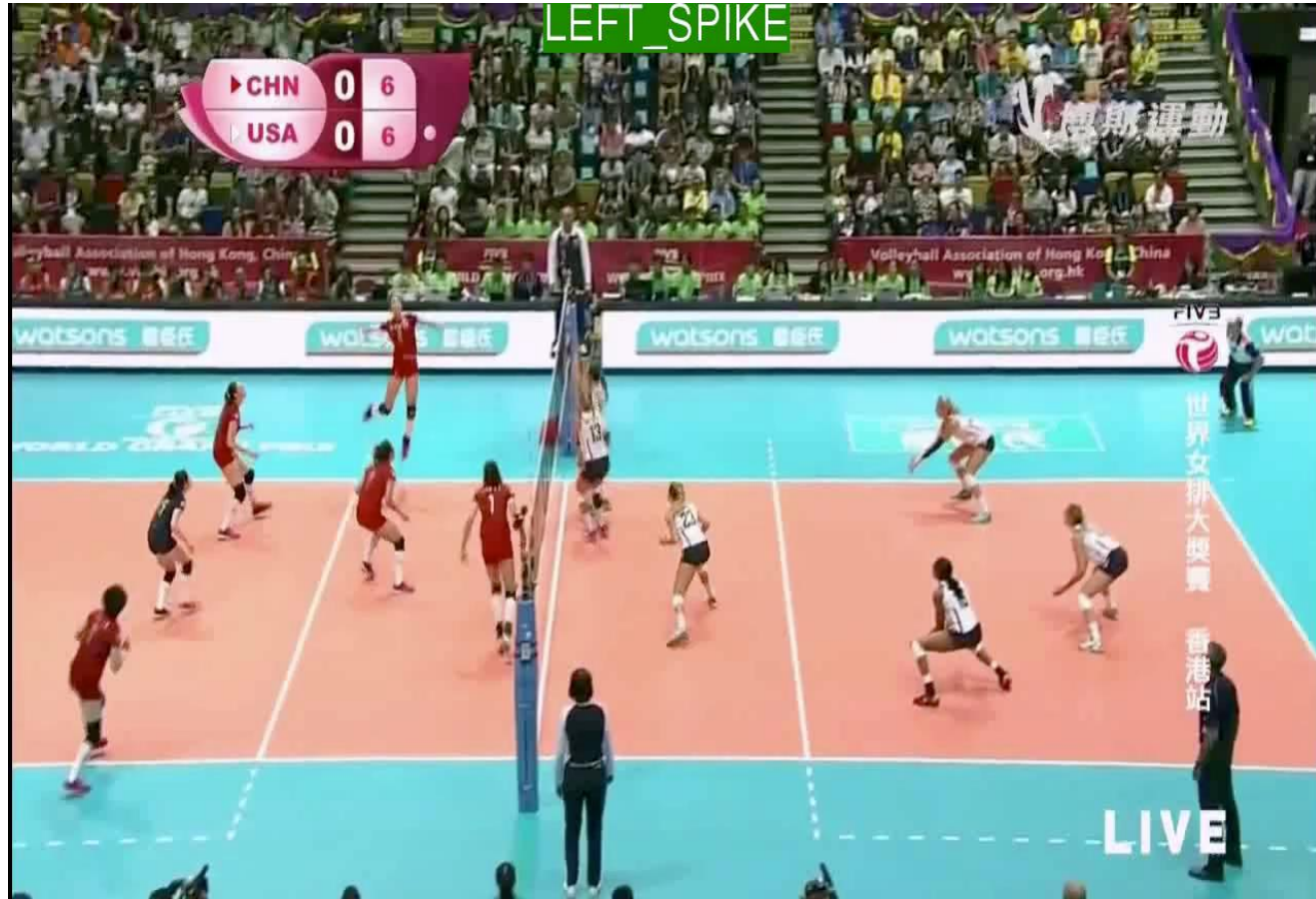  - Ibrahim et al. CVPR 2016

- Image annotation with label hierarchies
  - Hu et al. CVPR 2016

# Example: Rally in a Volleyball Game

Image Classifier → Group activity label

Challenge:
- high level description
- aggregate information over whole scene
- focus on relevant people



Group Activity = Majority's Activity

Group Activity = Key Player's Activity

Group Activity – Right spike

Intuitive fix: use person-centric representation

# Person Tracks

- Extract trajectories by tracking each person forward/backward in time

# Stage 1 : Learning Individual Action Features

# Stage1 : Learning Individual Action Features

# Stage 2: Learning Frame Representations

# Summary

# Volleyball Dataset – Frame Labels

- 4830 frames annotated from 55 volleyball videos
- 2/3 videos for training, 1/3 testing
- 9 player action labels
- 4 scene labels



**Spiking**   **Setting**   **Passing**   **Win point**

Left/right team variants

# Volleyball Dataset – People Labels

**Waiting**



**Digging**



**Setting**



**Spiking**



**Falling**



**Jumping**



**Moving**



**Standing**



**Blocking**

# Experimental results on Volleyball Dataset



| Method | Accuracy |
|---|---|
| Image Classification | 66.7 |
| Person Classification | 64.5 |
| Person - Fine tuned | 66.8 |
| Temp Model - Person | 67.5 |
| Temp Model - Image | 63.1 |
| Our Model w/o LSTM1 | 73.3 |
| Our Model w/o LSTM2 | 80.9 |
| Our Model | 81.6 |



Dense trajectories: 73.4-78.7

# Visualization of results

# Summary

- A two stage hierarchical model for group activity recognition

- LSTMs as a highly effective temporal model and temporal feature source

- People-relation modeling with simple pooling

# Outline

- Temporal structured models for group activities
  - Ibrahim et al. CVPR 2016



- Image annotation with label hierarchies
  - Hu et al. CVPR 2016

# Image Classification

- A natural image can be categorized with labels at different concept layers



Indoor ◯ ◯ outdoor man-made ◯ outdoor natural

leisure ◯ ◯ sports field ◯ man-made elements ◯ cabins houses

trench ◯ ◯ pitcher mound ◯ batter's box ◯ play-ground ◯ barn

field ◯ ◯ bat ◯ base ball ◯ grass ◯ person ◯ building

# Label Correlation Helps

- Such categorization at different concept layers can be modeled with label graphs

- It is natural and straightforward to leverage label correlation

# Goal: A generic label relation model

- Infer the entire label space from visual input
- Infer missing labels given a few fixed provided labels



| Metadata or Partial Label | | | |
| --- | --- | --- | --- |
| Visual Architecture | Prior Activation | Inference Machine on Knowledge Graph | Refined Probability |

An End-to-end Trainable System

Back-propagate Gradient from Loss Function

# Top-down Inference Neural Network

- Refine activations for each label
- Pass messages top-down and within each layer of label graph

Produce Visual Prior activation from CNN

$$x_t^i = W_t \cdot CNN(I^i) + b_t$$

Top-down inference



Visual Architecture

Activation at current concept layer

Vertical weight propagates information across concept layers

Horizontal weight propagates information within concept layers

$$a_t^i = V_{t-1,t} \cdot a_{t-1}^i + H_t \cdot x_t^i + b_t$$

Activation at last concept layer

# Bidirectional Inference Neural Network (BINN)

- Bidirectional inference to make information propagate across entire label structure

- Inference in each direction independently and blend results

**Top down Inference**

$$\overrightarrow{a}^i_t = \overrightarrow{V}_{t-1,t} \cdot \overrightarrow{a}^i_{t-1} + \overrightarrow{H}_t \cdot x^i_t + \overrightarrow{b}_{t,}$$

$$\overleftarrow{a}^i_t = \overleftarrow{V}_{t+1,t} \cdot \overleftarrow{a}^i_{t+1} + \overleftarrow{H}_t \cdot x^i_t + \overleftarrow{b}_{t,}$$

$$a^i_t = \overrightarrow{U}_t \cdot \overrightarrow{a}^i_t + \overleftarrow{U}_t \cdot \overleftarrow{a}^i_t + b_t$$

**Bottom-up Inference**

**Combine bidirectional inference result**

Visual Architecture

Bidirectional inference

# Structured Inference Neural Network (SINN)

- BINN is **hard** to train well

- **Regularize** connections

  with prior knowledge

  about label correlations

- Decompose connections

  into **Positive correlation** +

  **Negative correlation**



$$V^{+}_{<t-1,t>} = \begin{bmatrix} w_{0,0} & w_{0,1} & w_{0,2} & 0 \\ 0 & 0 & w_{1,2} & w_{1,3} \\ 0 & w_{2,1} & 0 & w_{2,3} \\ 0 & 0 & 0 & w_{3,3} \end{bmatrix}$$

Positive Correlation

$$V^{-}_{<t-1,t>} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ w_{1,0} & 0 & 0 & 0 \\ w_{2,0} & 0 & w_{2,1} & 0 \\ 0 & w_{3,1} & w_{3,2} & 0 \end{bmatrix}$$

Negative Correlation

Class: Cat, Zebra, Leopard, Hound
Attributes: Domestic, Spotted, Striped, Fast

# Structured Inference Neural Network (SINN)

- Evolve BINN formulation with regularization in connections

# Prediction from Purely Visual Input

- Visual architecture (e.g. Convolutional Neural Networks ) produces **visual activation**
- SINN implements **Information propagation** bidirectionally and produces refined **output activation**

# Prediction with Partially Observed Labels

- **Reverse Sigmoid** neuron produces activation from Partial labels

- **SINN** adapts both **visual activation** and **activation from partial labels** to infer the remaining labels



SINN Prediction with Partial Human Labels

# Datasets

- Evaluate method with two types of experiments on three datasets



### Animals with Attributes
[Lampert et al. 2009]

<u>Labels</u>
28 taxonomy terms
50 animal classes
85 attributes

**Task**: predict entire label set

- Taxonomy terms are constructed from Word Net as [Hwang et al. 2012]
- Knowledge graph constructed by combining class-attributes graph with taxonomy graph

### NUS-WIDE
[Chua et al. 2009]



<u>Labels</u>
698 image groups
81 concepts
1000 tags

**Task**: predict 81 concepts with observing tags/image groups

- Knowledge graph produced by Word Net using <u>semantic similarity</u>
- 698 image groups constructed from image meta data

### SUN 397
[Xiao et al. 2012]



<u>Labels</u>
3 coarse
16 general
397 fine-grained

**Task 1**: predict entire label set
**Task 2**: predict fine-grained scene given coarse scene category

- Knowledge graph provided by dataset

# Ex2: Inference from partial labels (NUS-WIDE)

- Produce predictions given partial 1k tags and 698 image groups



**Ground Truth**: railroad
**CNN + Logistic**: statue buildings person
**Our Predictions**: railroad person sky

**Ground Truth**: animal grass water dog
**CNN + Logistic**: grass person animal
**Our Predictions**: water animal dog

**Ground Truth**: rainbow clouds sky
**CNN + Logistic**: clouds water sky
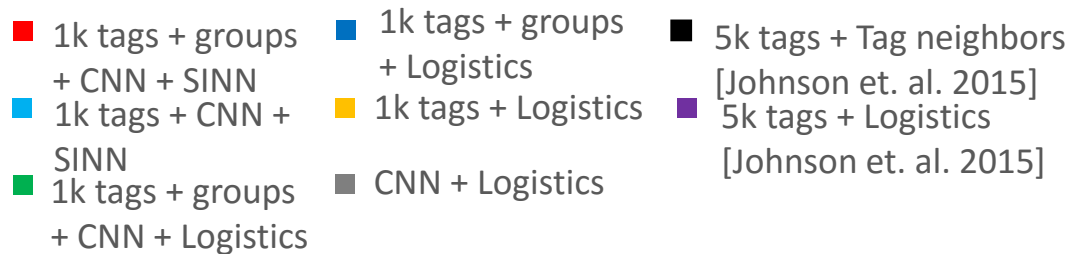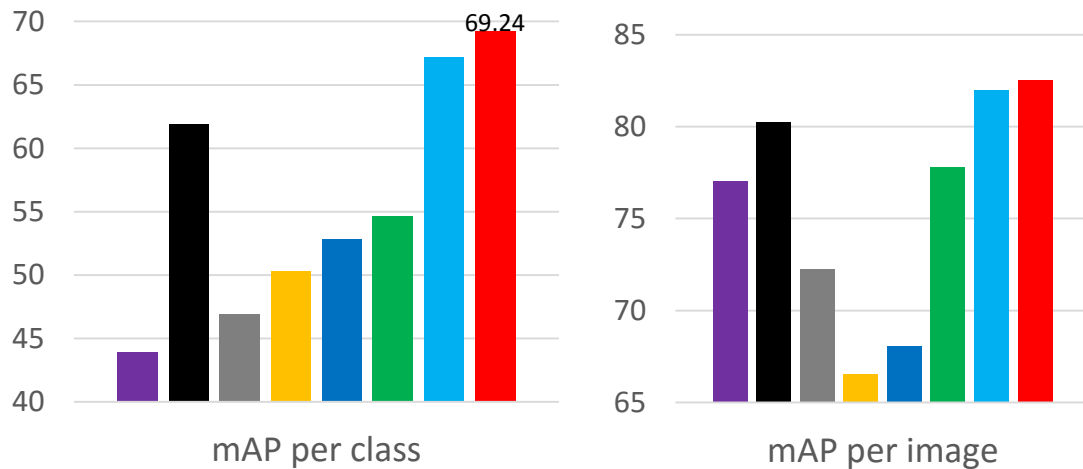**Our Predictions**: rainbow clouds sky

**Ground Truth**: food water
**CNN + Logistic**: food plants flower
**Our Predictions**: food plants water

Correct predictions are marked in **blue** while incorrect are marked in **red**

# Ex2: Inference from partial labels (NUS-WIDE)

- Evaluate on standard 81 ground truth classes of NUSWIDE
- **Outperform all baselines by large margin**



mAP per class — 69.24

mAP per image

- ■ 1k tags + groups + CNN + SINN
- ■ 1k tags + CNN + SINN
- ■ 1k tags + groups + CNN + Logistics
- ■ 1k tags + groups + Logistics
- ■ 1k tags + Logistics
- ■ CNN + Logistics
- ■ 5k tags + Tag neighbors [Johnson et. al. 2015]
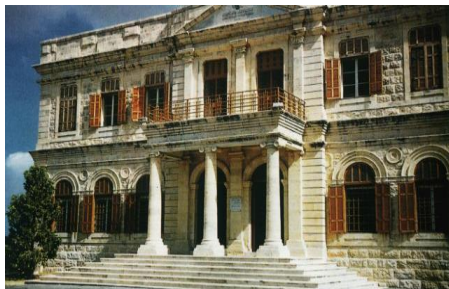- ■ 5k tags + Logistics [Johnson et. al. 2015]

# Ex2: Inference with partial labels (SUN397)

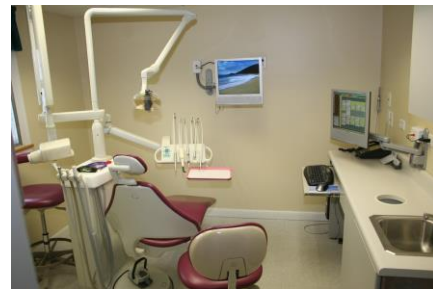- Produce predictions given coarse-level labels (3 coarse categories)



**CNN + Logistic**: campus
**Observed Label**:
outdoor/man-made
**Our Predictions**: abbey
**Ground Truth**: abbey

**CNN + Logistic**: building facade
**Observed Label**:
outdoor/man-made
**Our Predictions**: library/outdoor
**Ground Truth**: library/outdoor

**CNN + Logistic**: patio
**Observed Label**:
outdoor/natural;
outdoor/man-made
**Our Predictions**: picnic area
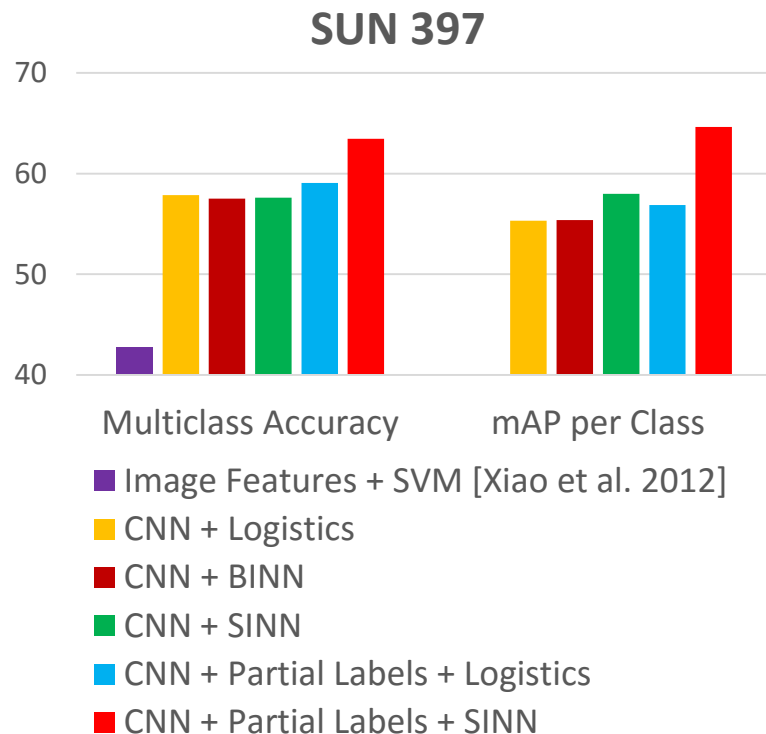**Ground Truth**: picnic area

**CNN + Logistic**: operating room
**Observed Label**: indoor
**Our Predictions**: dentists office
**Ground Truth**: dentists office

Correct predictions are marked in **blue** while incorrect are marked in **red**

# Ex2: Inference with partial labels (SUN397)

- Evaluate on 397 fine-grained scene categories
- **Significantly improved performance**



**SUN 397**

Legend:
- Image Features + SVM [Xiao et al. 2012]
- CNN + Logistics
- CNN + BINN
- CNN + SINN
- CNN + Partial Labels + Logistics
- CNN + Partial Labels + SINN

# Summary

- Inference in structured label space

- Relations within and across levels of a label space

- Model positive and negative correlations between labels in end-to-end trainable model

# Conclusion

- Methods for handling *structures* in deep networks
  - Spatial structure: learning gating functions to connect people for group activity recognition [Deng, Vahdat, Hu, Mori CVPR 2016]

  - Temporal structure: hierarchies of long short-term memory models for group activities [Ibrahim, Muralidharan, Deng, Vahdat, Mori CVPR 2016]

  - Label structure: message passing algorithms for multi-level image labeling; purely from image data or with partial labels [Hu, Zhou, Deng, Liao, Mori CVPR 2016]

# Acknowledgement