



<http://dmirlab.com>

VALSE2019 Sharing

Speaker : 章浩

Time : 2019/05/17



<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲



<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲



Outline

VALSE 2019 合肥

VALSE年度研讨会的主要目的是为计算机视觉、图像处理、模式识别与机器学习研究领域内的中国青年学者（以80后为主）提供一个深层次学术交流的舞台。在这个舞台上，我们恪守并倡导理性批判、勇于探索、实证、创新等科学精神；在这个舞台上，我们倡导自由平等原则下、理性而纯学术的百家争鸣和思想交锋；这个舞台上，我们期望欣赏到国内青年学者越来越优美的学术华尔兹（VALSE）。通过这个舞台，我们期望促进国内青年学者的思想交流和学术合作，从而在相关领域做出重量级学术贡献，提升中国学者在国际学术舞台上的学术影响力。

截至目前，VALSE已成功举办8届，分别为 VALSE 2011（杭州），VALSE 2012（西安），VALSE 2013（南京），VALSE 2014（青岛），VALSE 2015（成都），VALSE 2016（武汉），VALSE 2017（厦门），VALSE 2018（大连）。VALSE 2019 将于2019年4月在合肥举行，由中国人工智能学会与安徽大学汤进教授团队联合合肥工业大学、中国科技大学等团队共同承办。



<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲



<http://dmirlab.com>

Weakly | Active Learning

- What Weakly ?
- Mingming Chen, NKU
 - Multi-scale Backbone Architecture
 - 通用视觉基元属性感知
 - 互联网大数据自主学习
- Qixiang Ye, CAS
 - 粗粒度的弱监督标记
 - Towards Self-Learning



<http://dmirlab.com>

Weakly | Active Learning

- Shengjun Huang, NHU
 - Active Learning
- Summary (个人观点, 真伪自辩)



<http://dmirlab.com>

Weakly | Active Learning

- What Weakly ?
- Mingming Chen, NKU
 - Multi-scale Backbone Architecture
 - 通用视觉基元属性感知
 - 互联网大数据自主学习
- Qixiang Ye, CAS
 - 粗粒度的弱监督标记
 - Towards Self-Learning



Weakly | Active Learning

- 弱监督分类:
 - 不完全监督: 只有一部分训练数据具备标签
 - 不确切监督: 训练数据只具备粗粒度标签
 - 不准确监督: 给出的标签并不总是真值（标签有噪声）
- 弱监督的含义:
 - 弱监督给出的标签会在某种程度上弱于我们面临的任务所要求的输出
- 研究背景:
 - 数据集很重要, 但是标注成本很大
 - 研究方法难选: 深度神经网络共性技术, 视觉基元属性感知。



<http://dmirlab.com>

Weakly | Active Learning

- What Weakly ?
- Mingming Chen, NKU
 - Multi-scale Backbone Architecture
 - 通用视觉基元属性感知
 - 互联网大数据自主学习
- Qixiang Ye, CAS
 - 粗粒度的弱监督标记
 - Towards Self-Learning



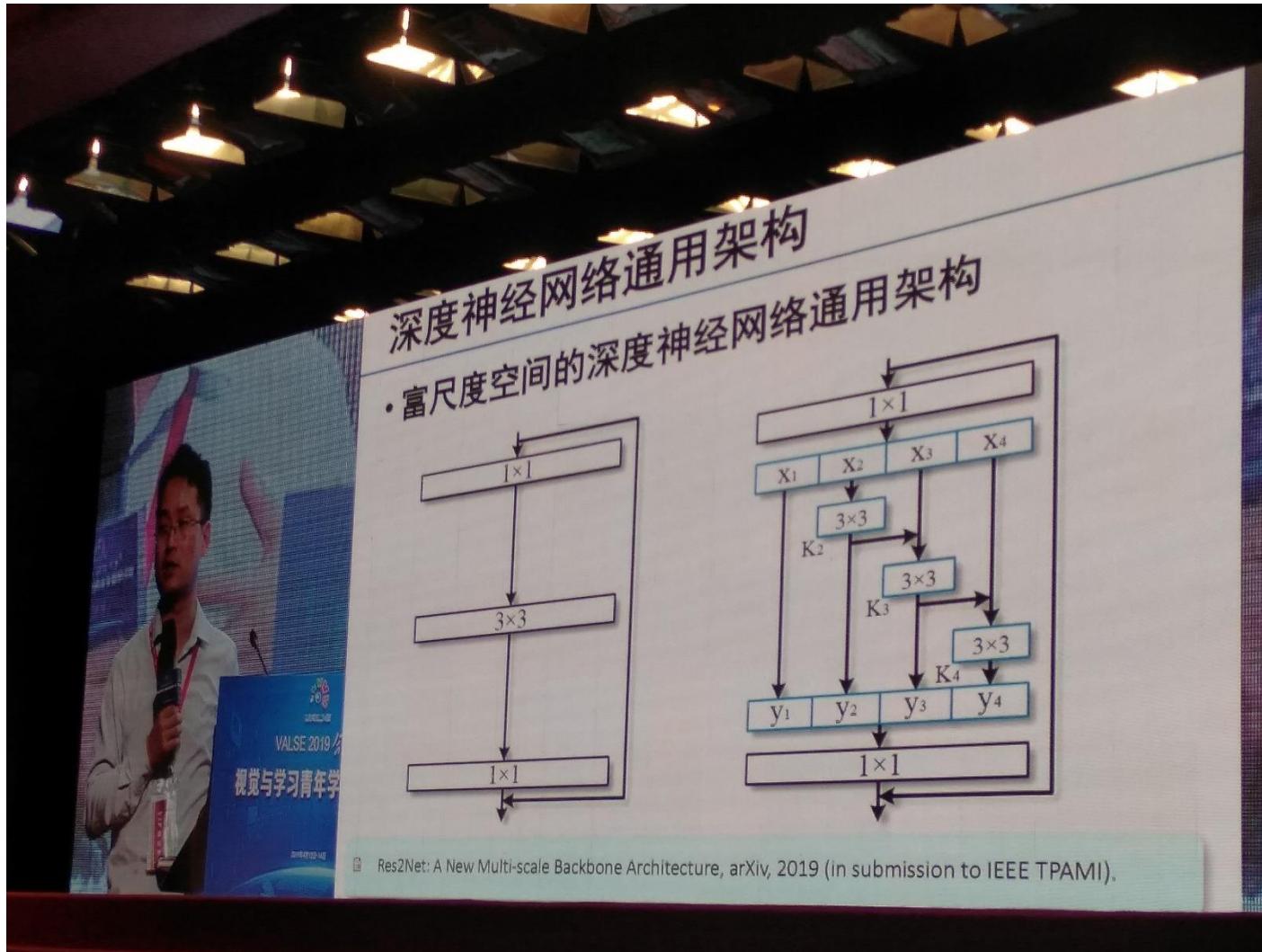
Weakly | Active Learning

<http://dmirlab.com>

- Motivation:
 - 当前各种深度网络的进步得益于网络多尺度信息综合能力的提升
- 报告核心:
 - 富尺度空间神经网络架构：多任务协同求解，鲁棒性提高
 - 显著性物体检测：预设基元属性感知能力，减少数据依赖
 - 互联网大数据自主学习：减少人工标注，自动学习

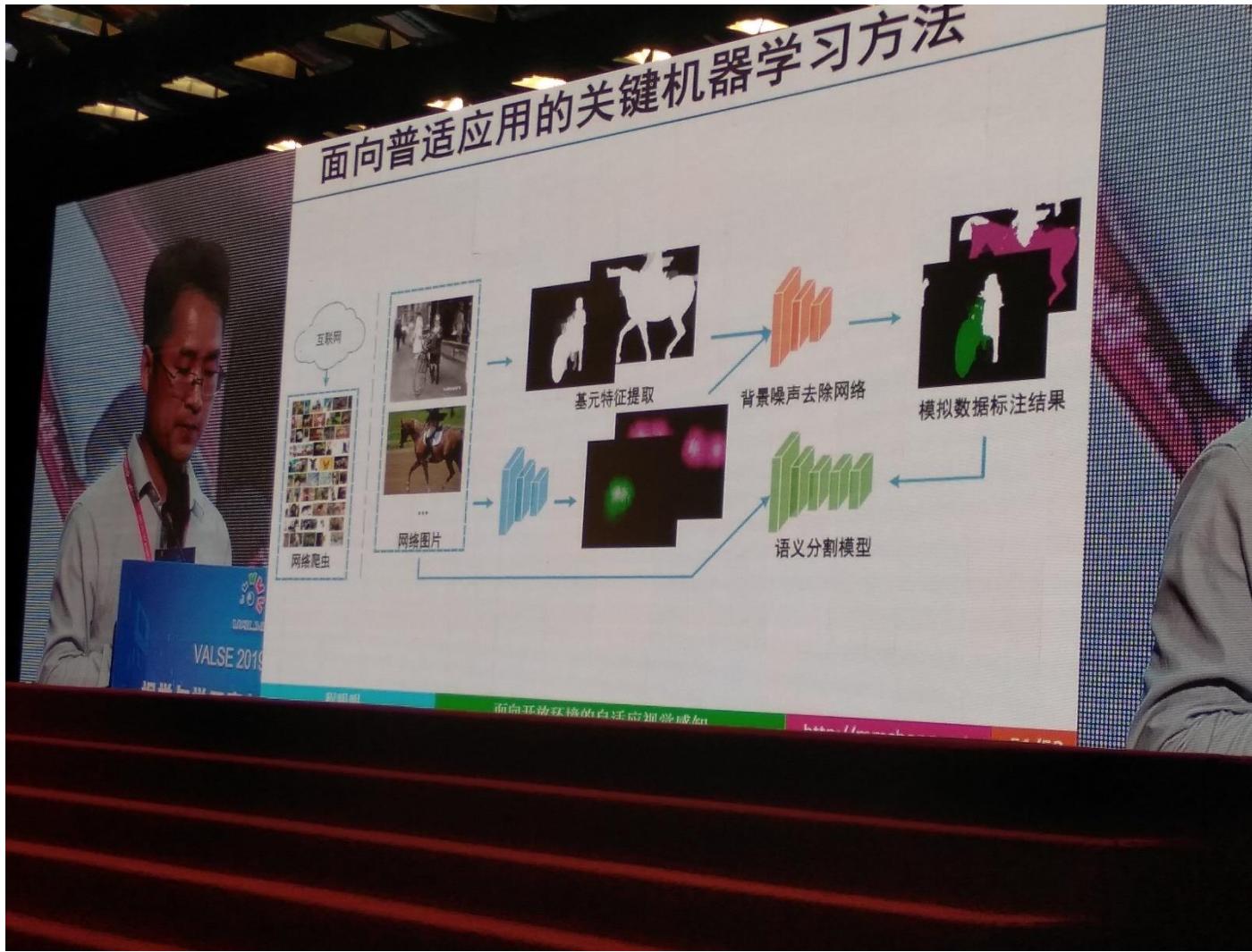
Weakly | Active Learning

<http://dmirlab.com>



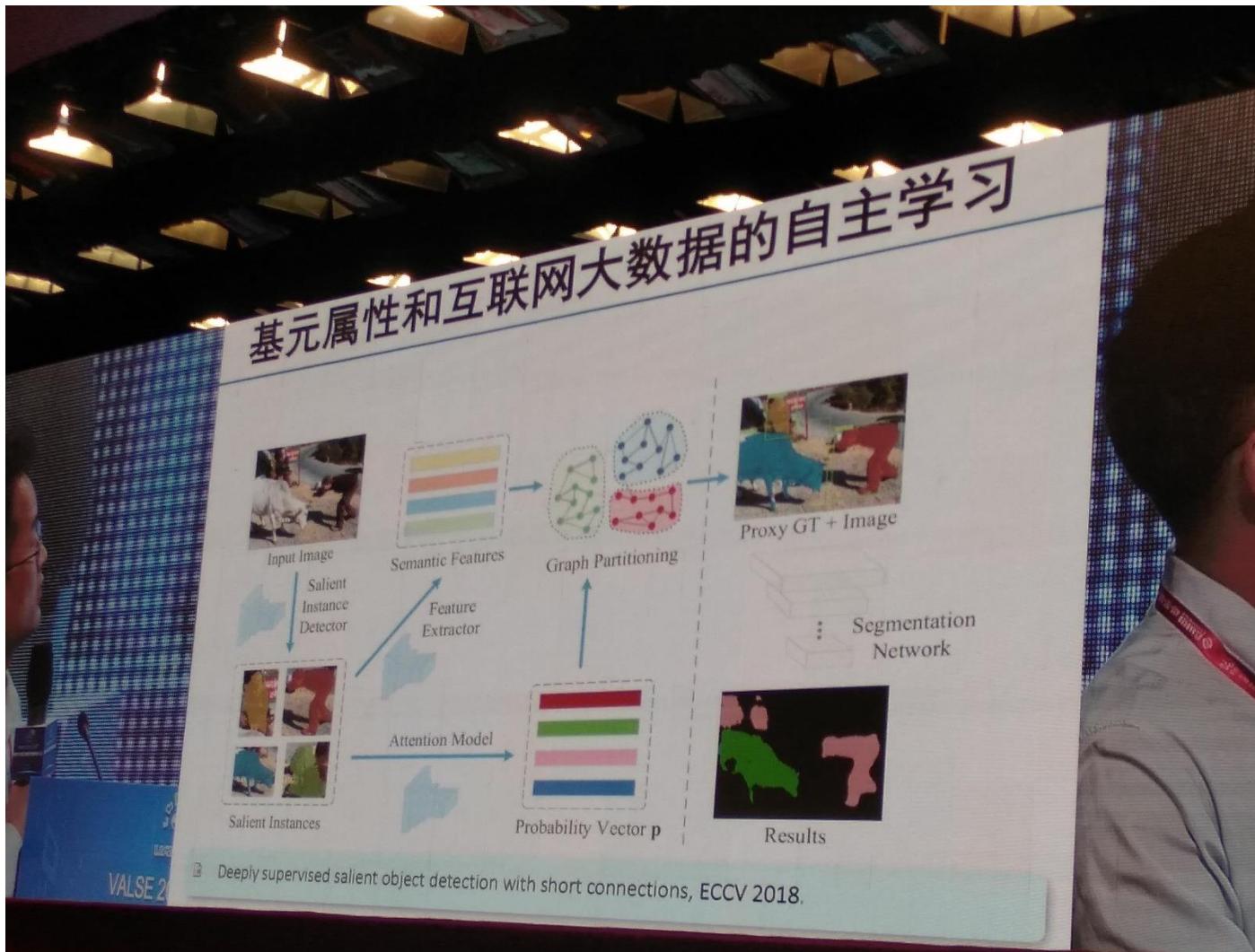
Weakly | Active Learning

<http://dmirlab.com>



Weakly | Active Learning

<http://dmirlab.com>



Weakly | Active Learning





<http://dmirlab.com>

Weakly | Active Learning

- What Weakly ?
- Mingming Chen, NKU
 - Multi-scale Backbone Architecture
 - 通用视觉基元属性感知
 - 互联网大数据自主学习
- Qixiang Ye, CAS
 - 粗粒度的弱监督标记
 - Towards Self-Learning



<http://dmirlab.com>

Weakly | Active Learning

- What Weakly ?
- Mingming Chen, NKU
 - Multi-scale Backbone Architecture
 - 通用视觉基元属性感知
 - 互联网大数据自主学习
- Qixiang Ye, CAS
 - 粗粒度的弱监督标记
 - Towards Self-Learning



Weakly | Active Learning

- Motivation:
 - 目前detection / segmentation需要为每一个任务制定详细而具体的标注，而且无法共享，故而标注成本极大
- 核心观点:
 - 粗粒度的弱监督标记：比如，
 - 只给目标物体上画一条线，
 - 只在目标物体上打一个点，
 - 仅仅告诉模型一系列图片中包含什么而不给位置，让模型自己学习找到这些目标

Weakly | Active Learning

<http://dmirlab.com>



Weakly | Active Learning

<http://dmirlab.com>





Weakly | Active Learning

<http://dmirlab.com>

- Min-entropy Latent Model for Weakly Supervised object Detection, (CVPR2018)
- CMIL: Continuation Multiple Instance Learning for Weakly Supervised Detection (CVPR2019 **Oral**)
- SPN: Soft Proposal Network for Weakly Supervised Object Localization (ICCV2017)
- Learning Instance Activation Maps for Weakly Supervised Instance Segmentation (CVPR2019)
- PAMI2019: Recurrent Learning(MELM+RecurrentLearning)
- PeakResponseMapping(PRM) (CVPR2018)



<http://dmirlab.com>

Weakly | Active Learning

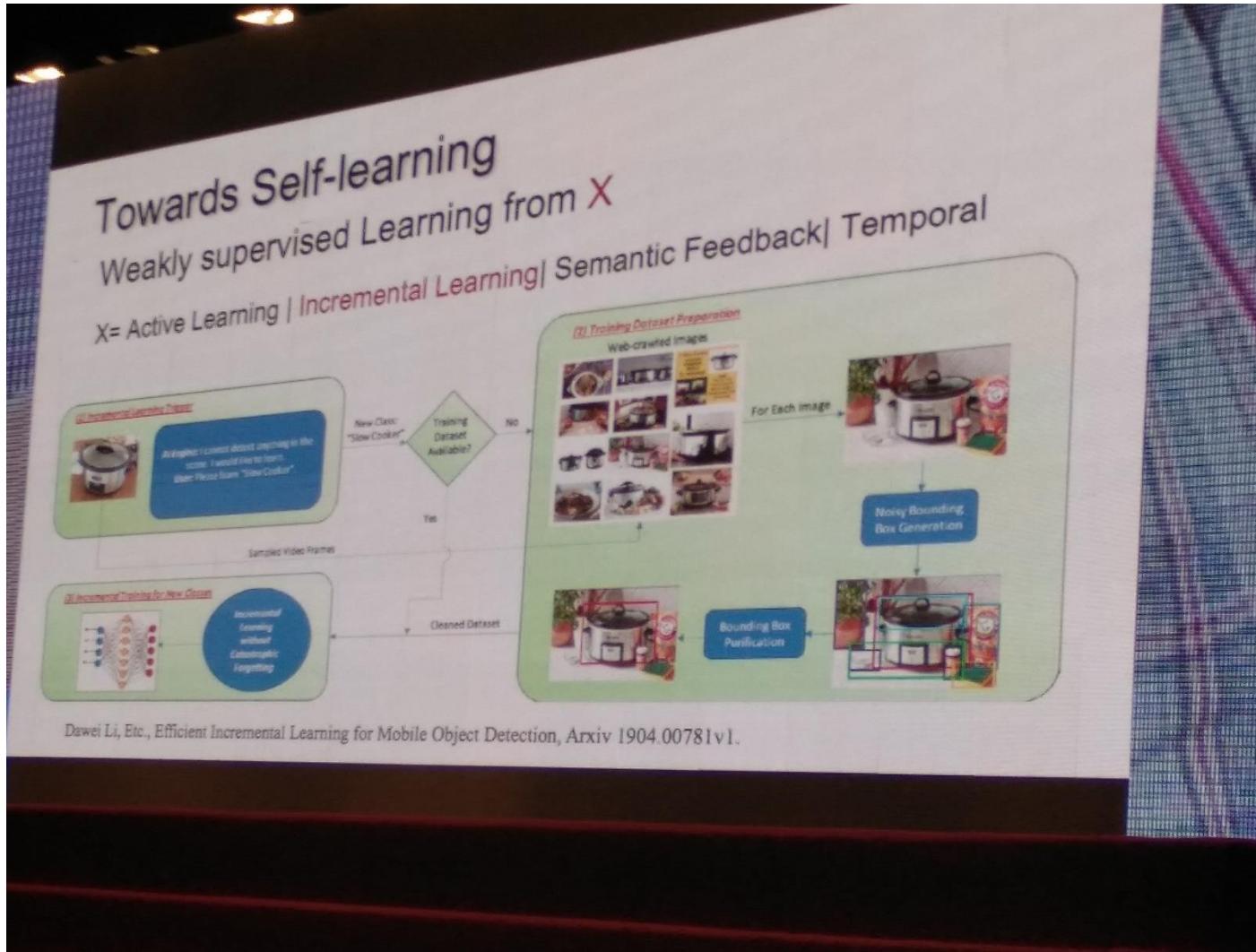
- What Weakly ?
- Mingming Chen, NKU
 - Multi-scale Backbone Architecture
 - 通用视觉基元属性感知
 - 互联网大数据自主学习
- Qixiang Ye, CAS
 - 粗粒度的弱监督标记
 - Towards Self-Learning

Weakly | Active Learning



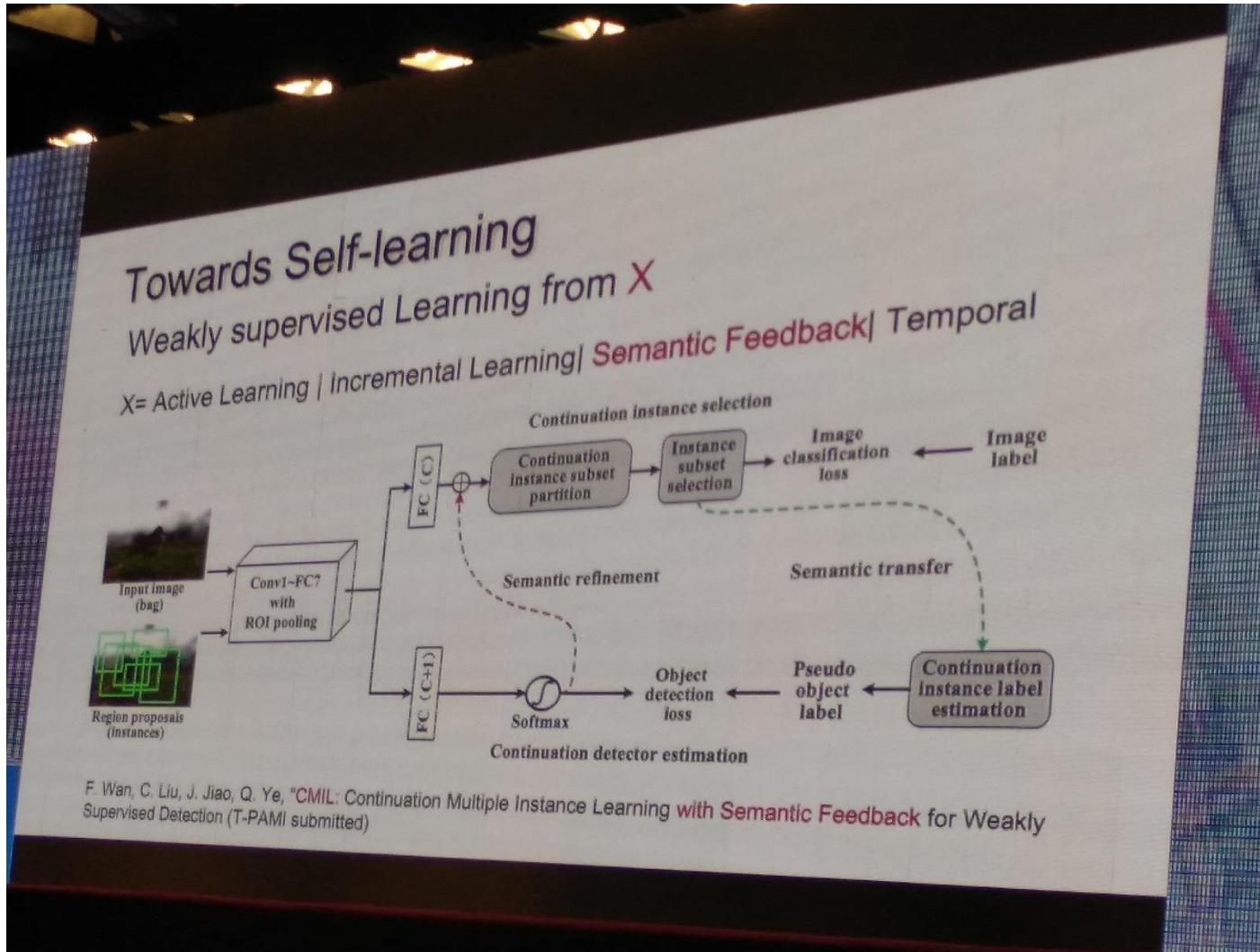
Weakly | Active Learning

<http://dmirlab.com>



Weakly | Active Learning

<http://dmirlab.com>



<http://dmirlab.com>



<http://dmirlab.com>

Weakly | Active Learning

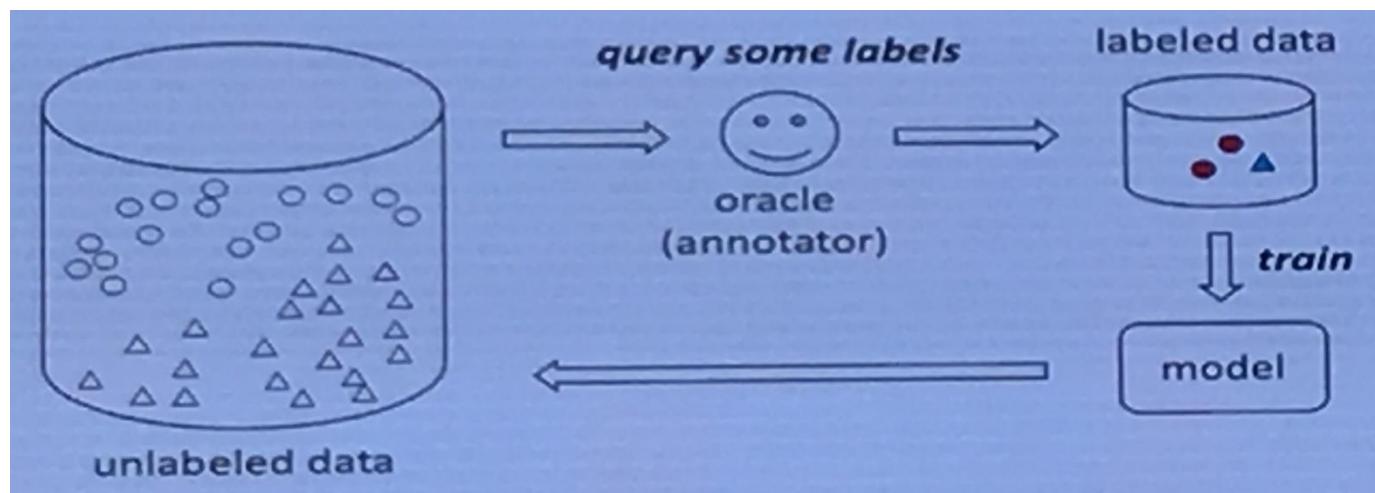
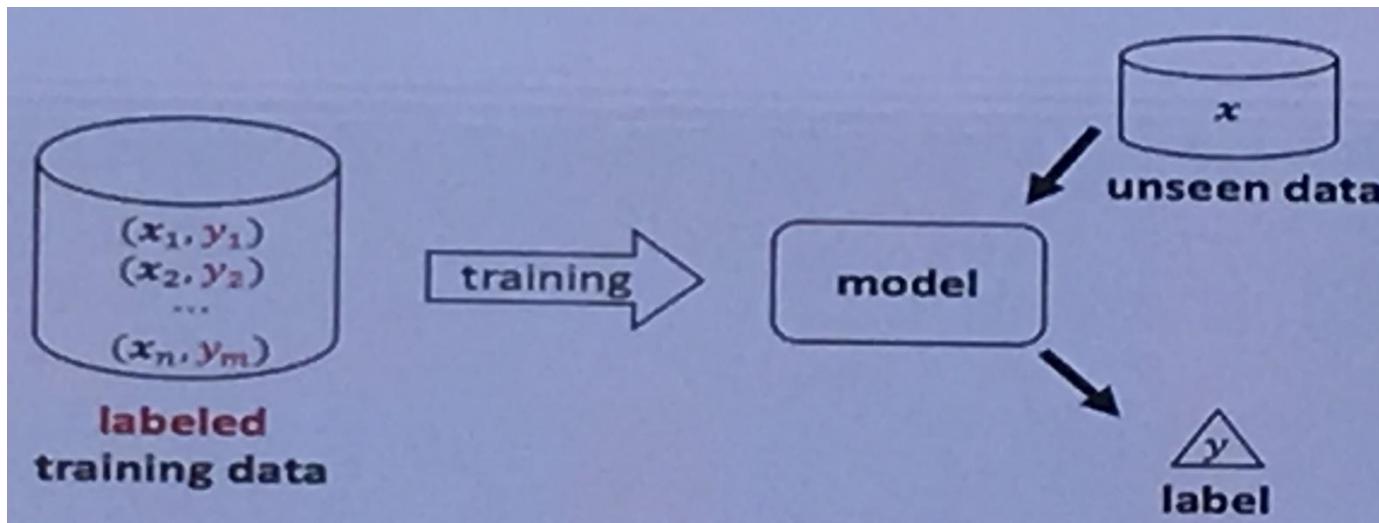
- Shengjun Huang, NHU
 - Active Learning
- Summary (个人观点, 真伪自辩)



Weakly | Active Learning

- Motivation:
 - 大量的数据标记成本巨大，有一些甚至是不可得的
 - 比如医院的患者信息，异常检测的异常样本，几年才发生一次异常
 - 标记成本的判定非常复杂：“**Label 生而不平等**”
- Solution:
 - Active Learning: 用**最小的标记代价**获得最大的performance
- Non-Active vs. Active Learning
 - 如下图

Weakly | Active Learning





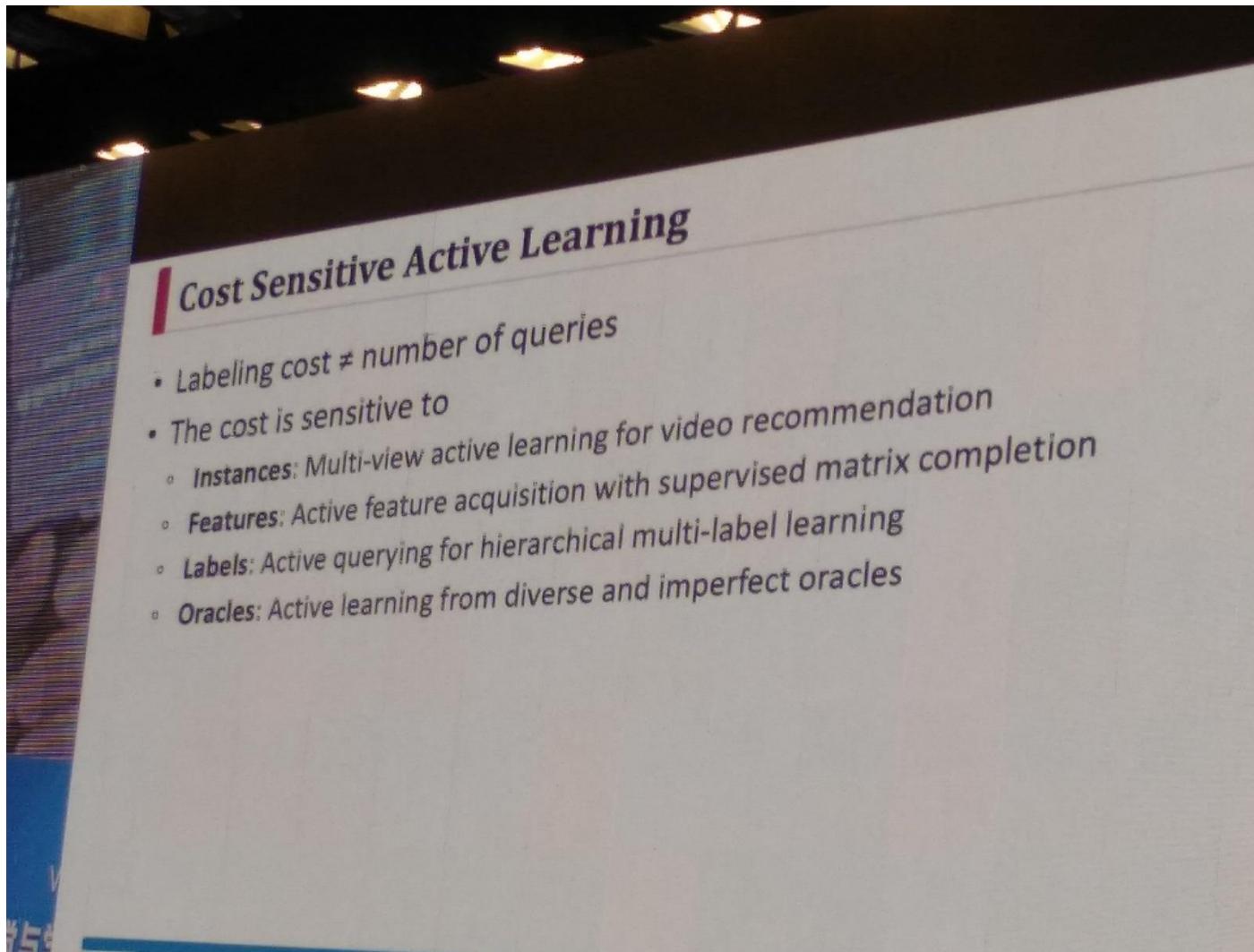
<http://dmirlab.com>

Weakly | Active Learning

- **Cost Sensitive** Active Learning
 - goal: train an effective model with **least** labeling cost
 - **细致定义 least** 就是 Active Learning 的 核心
- The cost is **not** number of queries (too naive too young),
 - It is about...
 - E.g. Multi-view Active learning for Video Recommendation

Weakly | Active Learning

<http://dmirlab.com>



A photograph of a presentation slide titled "Cost Sensitive Active Learning". The slide content includes:

- Labeling cost ≠ number of queries
- The cost is sensitive to
 - Instances: Multi-view active learning for video recommendation
 - Features: Active feature acquisition with supervised matrix completion
 - Labels: Active querying for hierarchical multi-label learning
 - Oracles: Active learning from diverse and imperfect oracles

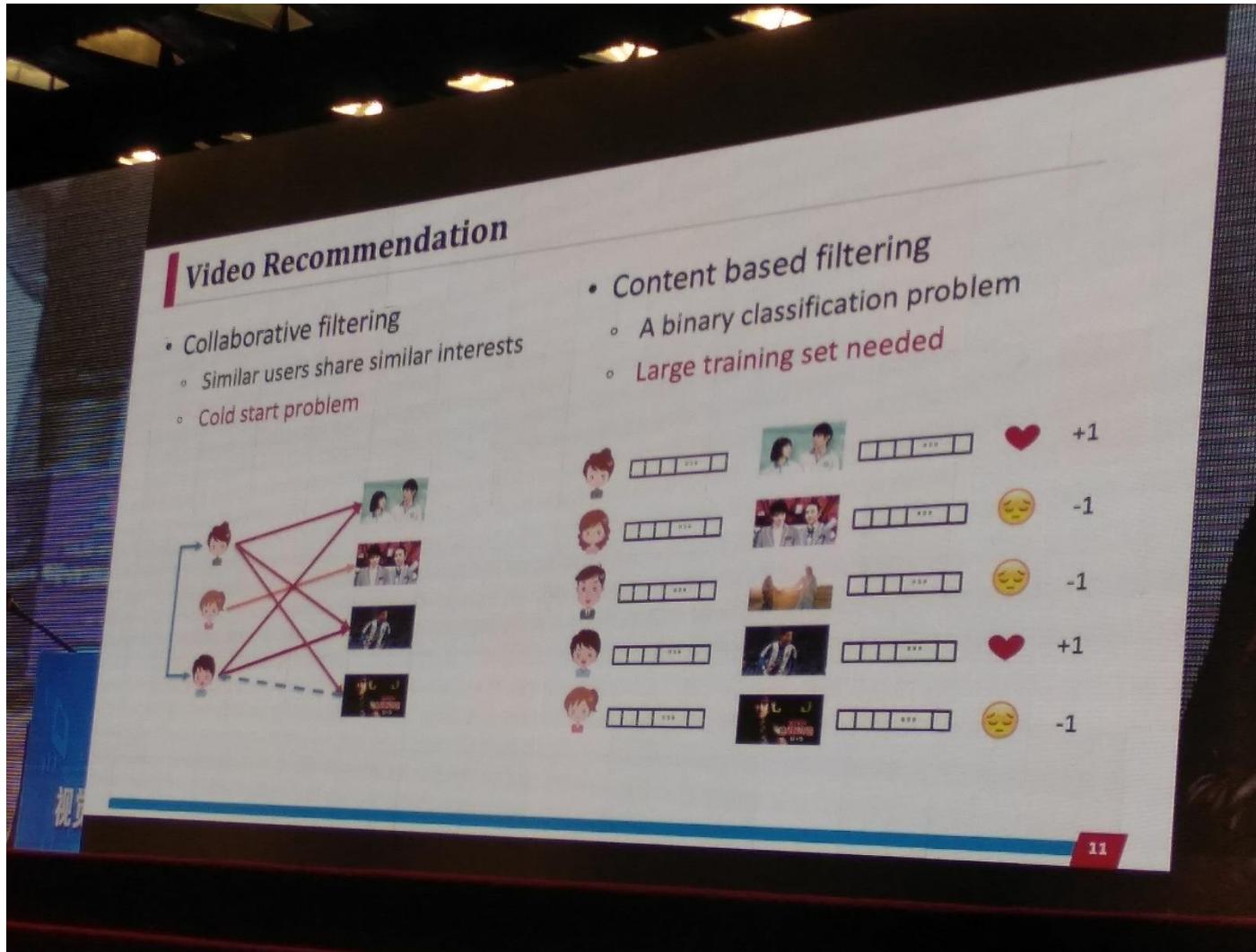


Weakly | Active Learning

- Instances
 - 核心：怎么选**最好的**instance？
- 视频推荐：
 - 协同过滤（冷门启动问题）、基于内容的过滤（需要大量数据训练）
- 多视角视频表示：
 - 视觉特征、文本特征、用户特征、标签
- Motivation：
 - 在视频推荐任务中，文本特征（即评论）获取需要很大代价，视觉特征不需要人力代价。
- Idea: Visual to text Mapping

Weakly | Active Learning

<http://dmirlab.com>



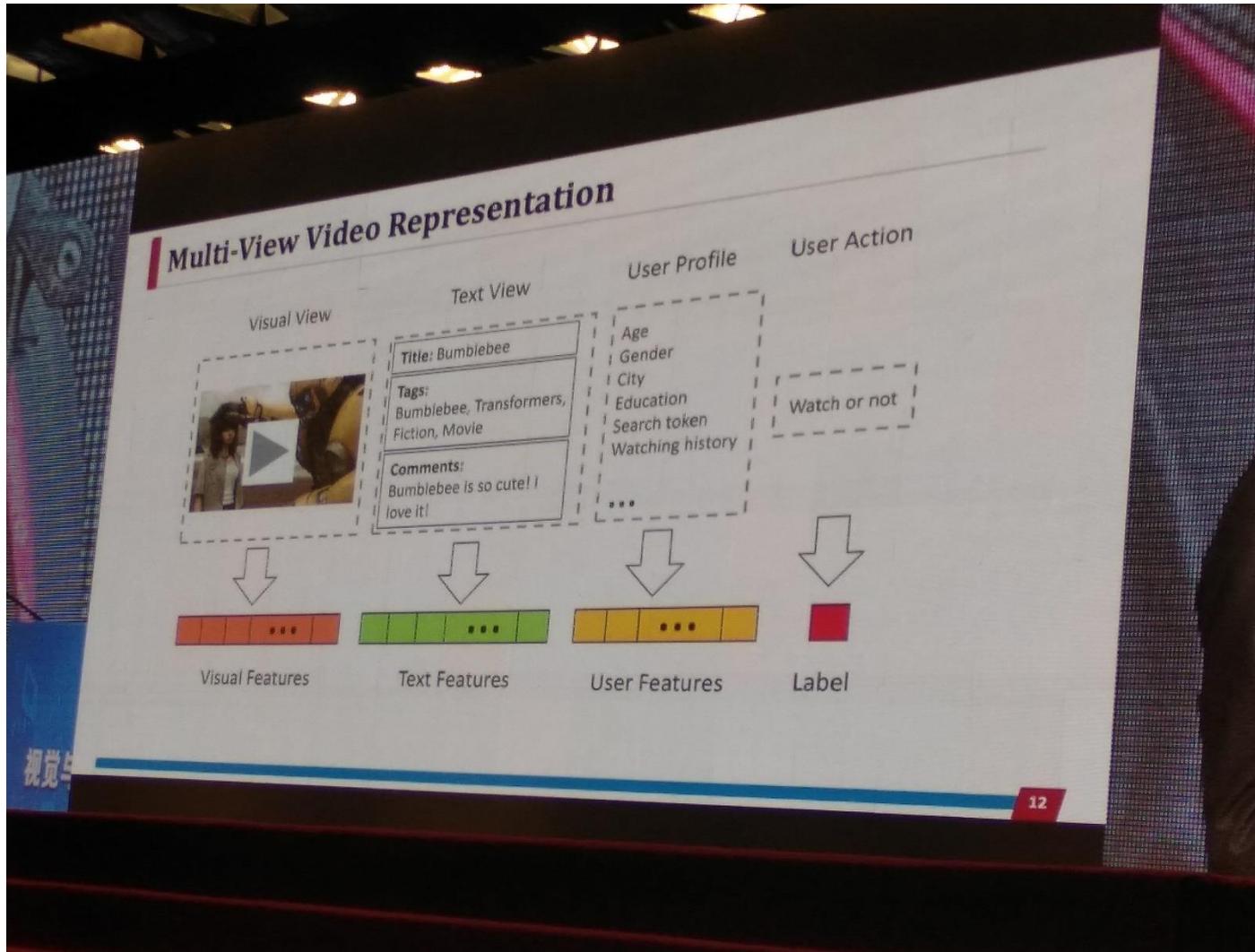
Video Recommendation

- Collaborative filtering
 - Similar users share similar interests
 - Cold start problem
- Content based filtering
 - A binary classification problem
 - Large training set needed

The slide illustrates two recommendation paradigms:

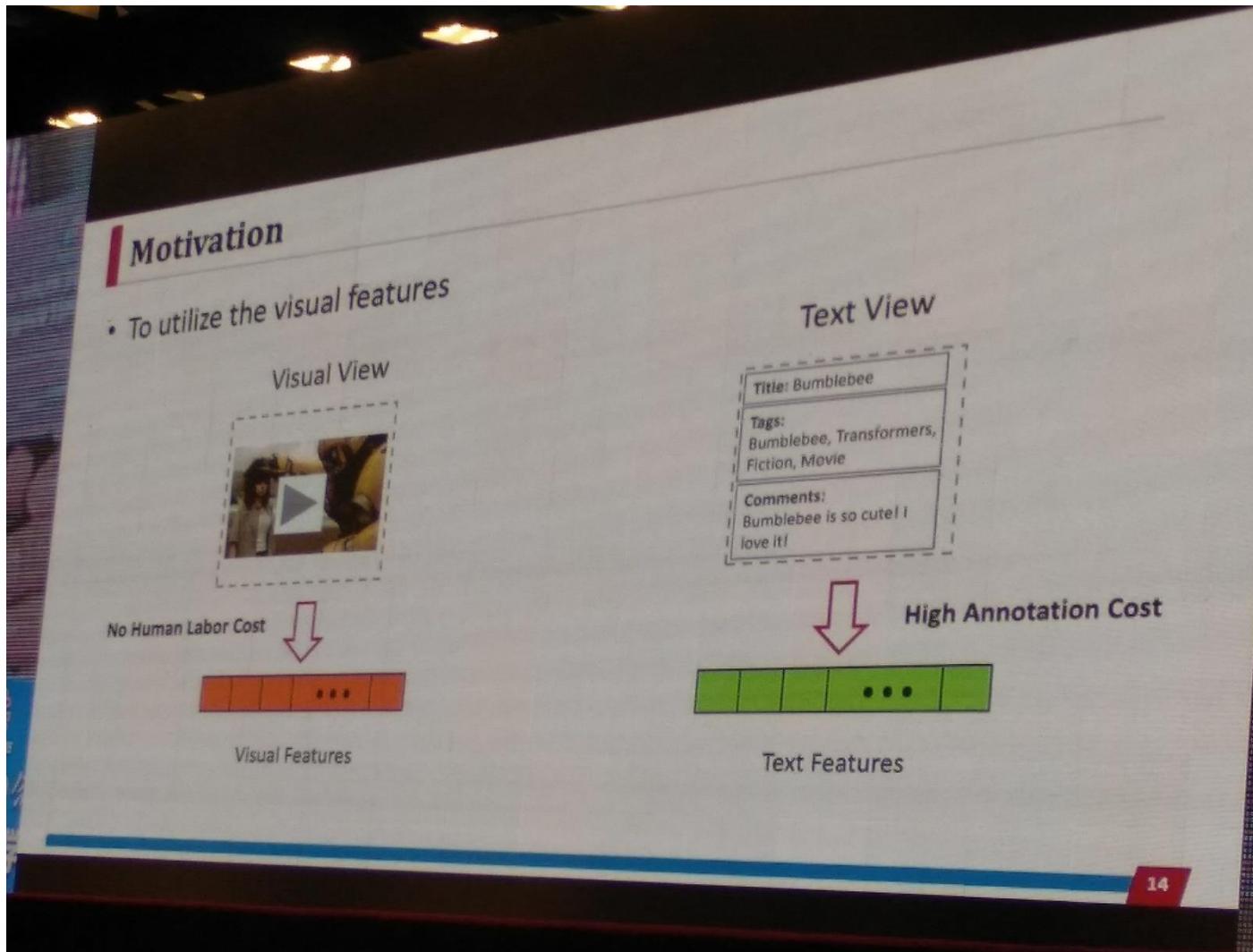
- Collaborative filtering:** Shows a diagram where four users are connected by arrows to five video items. One user has a dashed arrow pointing to a video item they have not interacted with yet.
- Content based filtering:** Shows a grid of 8 users and 8 video items. Each user has a profile represented by a row of 5 binary values (0 or 1). To the right of each video item is a user icon and a sentiment rating:
 - Row 1: User 1 (heart) +1
 - Row 2: User 2 (neutral face) -1
 - Row 3: User 3 (neutral face) -1
 - Row 4: User 4 (heart) +1
 - Row 5: User 5 (neutral face) -1

Weakly | Active Learning

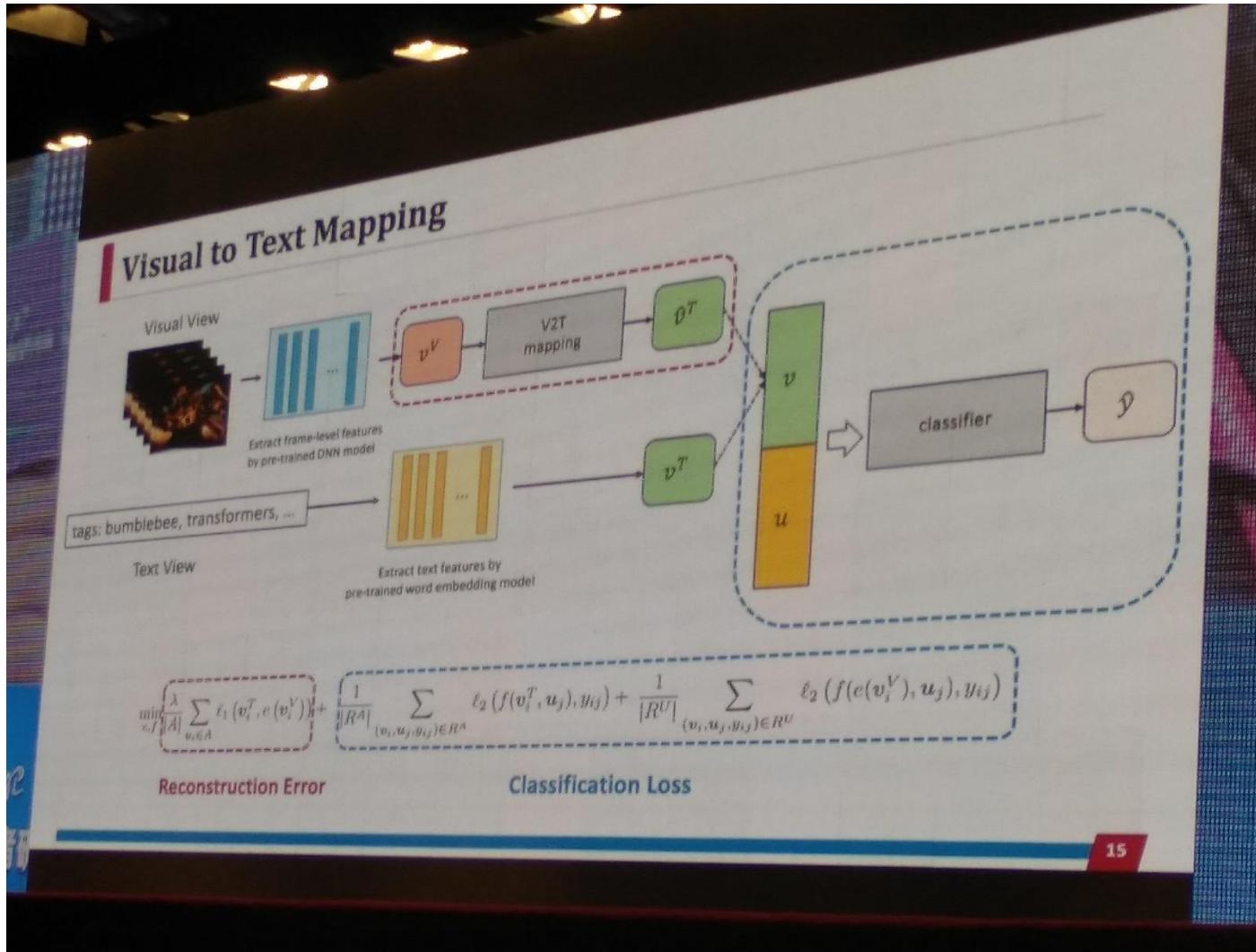


Weakly | Active Learning

<http://dmirlab.com>



Weakly | Active Learning

<http://dmirlab.com>




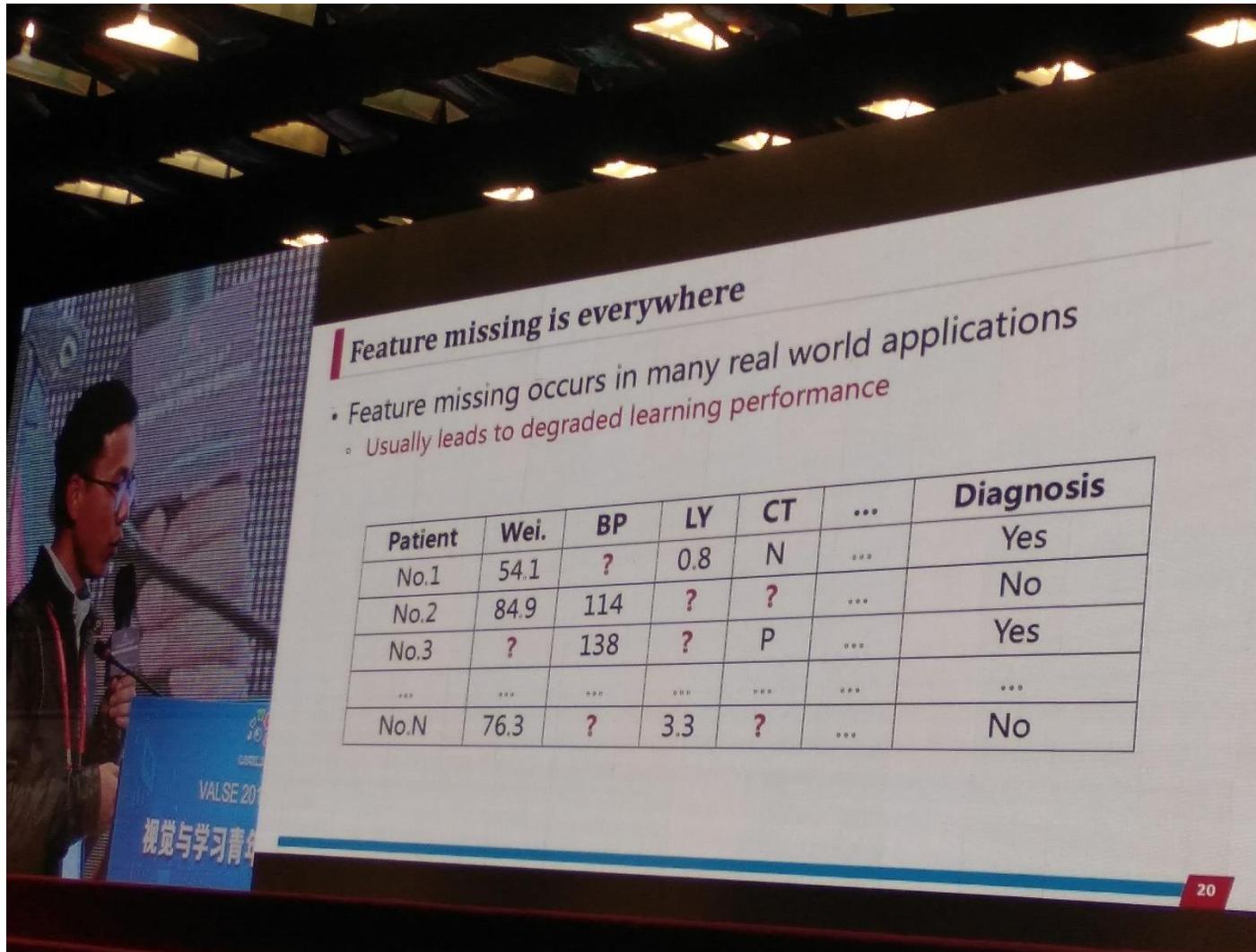
<http://dmirlab.com>

Weakly | Active Learning

- Features

- 问题：现实应用中往往会出现特征丢失现象，通常导致学习性能下降
 - 核心：如何选**更好的**特征？

Weakly | Active Learning

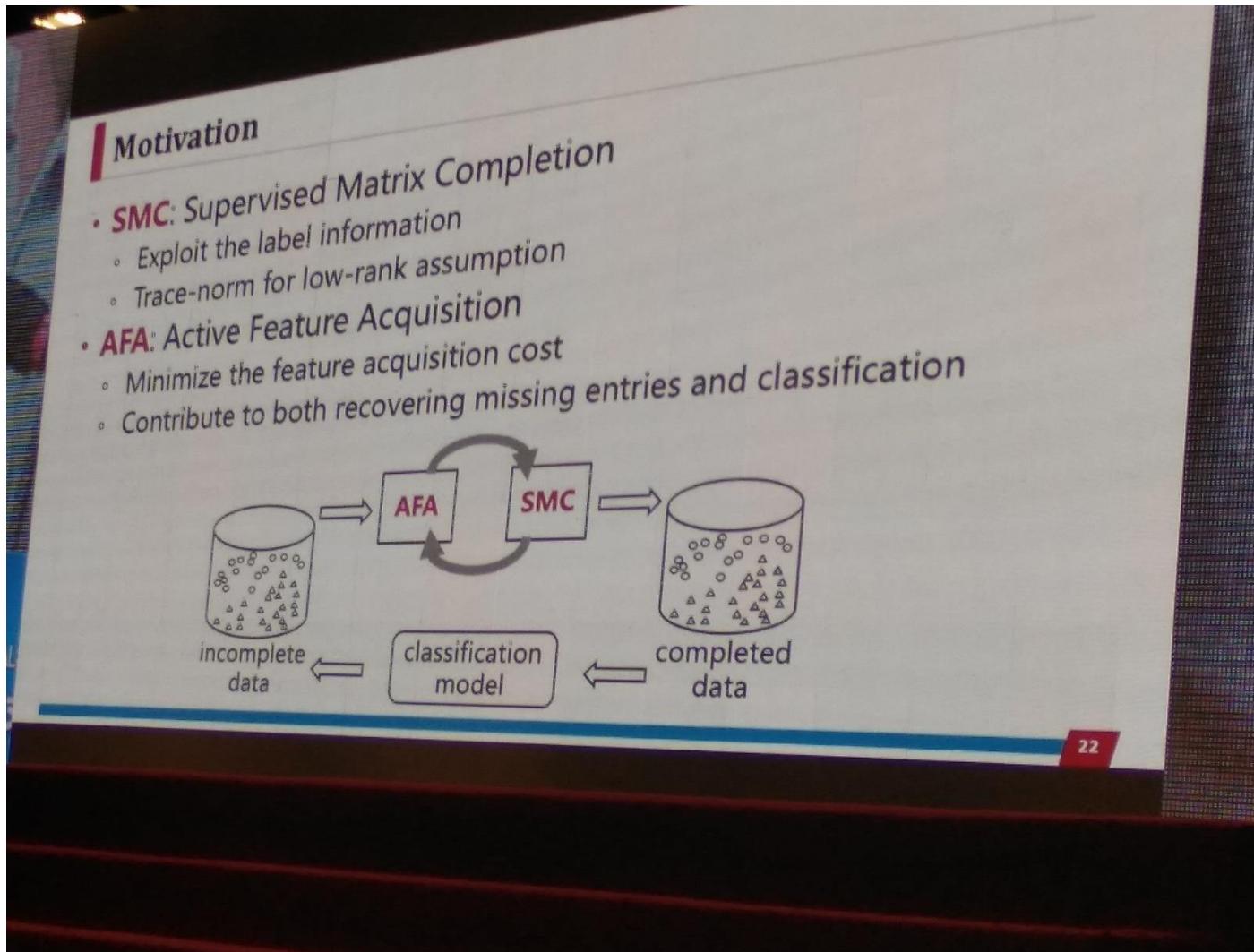


The image shows a presentation slide with a background video of a speaker. The slide has a title and a table. The title is "Feature missing is everywhere". Below the title is a bulleted list: "• Feature missing occurs in many real world applications" and "◦ Usually leads to degraded learning performance". The table has columns labeled "Patient", "Wei.", "BP", "LY", "CT", "...", and "Diagnosis". The rows contain data for patients No.1, No.2, No.3, and No.N. The "Wei." column has values 54.1, 84.9, ?, and 76.3. The "BP" column has values ?, 114, ?, and ?. The "LY" column has values 0.8, ?, ?, and 3.3. The "CT" column has values N, ?, P, and ?. The "..." column has values ..., ..., ..., and The "Diagnosis" column has values Yes, No, Yes, and No. There are also some small numbers (e.g., 0.00, 0.000) in the empty cells of the table.

Patient	Wei.	BP	LY	CT	...	Diagnosis
No.1	54.1	?	0.8	N	...	Yes
No.2	84.9	114	?	?	...	No
No.3	?	138	?	P	...	Yes
...
No.N	76.3	?	3.3	?	...	No

Weakly | Active Learning

<http://dmirlab.com>



Weakly | Active Learning

Supervised Matrix Completion

- Supervised matrix completion
 - Jointly optimize the missing entries and the classification model
 - Class labels guide the matrix completion
 - Recovered values contribute the classification

$$\min_{\widehat{\mathbf{X}}, f} \frac{1}{2} \|\mathcal{R}_\Omega(\widehat{\mathbf{X}} - \mathbf{X})\|_F^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}} + \lambda_2 \ell(\widehat{\mathbf{X}}, f)$$

Recovery error on
observed entries
Trace-norm
for low-rank
Classification
loss

$$[\mathcal{R}_\Omega(X)]_{i,j} = \begin{cases} X_{i,j} & \text{if } (i,j) \in \Omega \\ 0 & \text{otherwise} \end{cases}$$

Weakly | Active Learning

<http://dmirlab.com>

Active Feature Acquisition

- Which entries are most important?
 - Useful for recovering missing entries
 - Useful for the classifier training
- A variance-based criterion for active selection
 - Informativeness

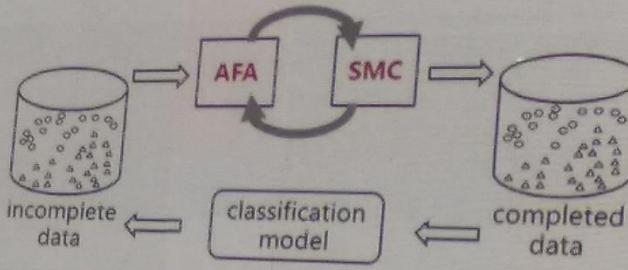
$$I_{i,j} = \sum_{t=1}^T (X_{i,j}^t - \bar{X}_{i,j})^2$$

recovered value
at t -th iteration

$$\bar{X}_{i,j} = \frac{1}{T} \sum_{t=1}^T X_{i,j}^t$$

$$\min_{\widehat{\mathbf{X}}, f} \frac{1}{2} \|\mathcal{R}_\Omega(\widehat{\mathbf{X}} - \mathbf{X})\|_F^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}} + \lambda_2 \ell(\widehat{\mathbf{X}}, f)$$

Patient	Wei.	BP	LY	CT	...	Diagnosis
No.1	54.1	?	0.8	N	...	Yes
No.2	84.9	114	?	?	...	No
No.3	?	138	?	P	...	Yes
...
No.N	76.3	?	3.3	?	...	No



28



<http://dmirlab.com>

Weakly | Active Learning

- Labels
 - 标签有层次结构，要平衡成本和信息

Weakly | Active Learning

<http://dmirlab.com>

Active Feature Acquisition

- Which entries are most important?
 - Useful for recovering missing entries
 - Useful for the classifier training
- A variance-based criterion for active selection
 - Informativeness

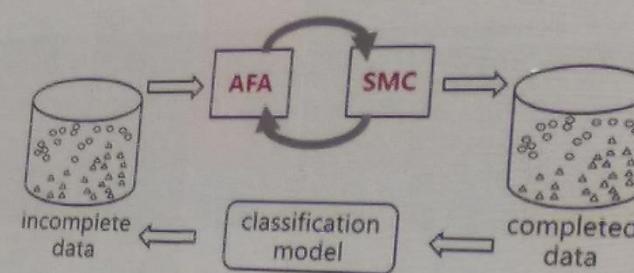
$$I_{i,j} = \sum_{t=1}^T (X_{i,j}^t - \bar{X}_{i,j})^2$$

recovered value
at t -th iteration

$$\bar{X}_{i,j} = \frac{1}{T} \sum_{t=1}^T X_{i,j}^t$$

$$\min_{\widehat{\mathbf{X}}, f} \frac{1}{2} \|\mathcal{R}_\Omega(\widehat{\mathbf{X}} - \mathbf{X})\|_F^2 + \lambda_1 \|\widehat{\mathbf{X}}\|_{\text{tr}} + \lambda_2 \ell(\widehat{\mathbf{X}}, f)$$

Patient	Wei.	BP	LY	CT	...	Diagnosis
No.1	54.1	?	0.8	N	...	Yes
No.2	84.9	114	?	?	...	No
No.3	?	138	?	P	...	Yes
...
No.N	76.3	?	3.3	?	...	No



incomplete data classification model completed data

28



Weakly | Active Learning

<http://dmirlab.com>

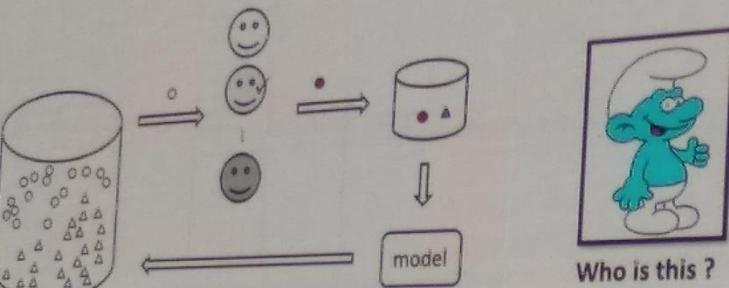
- Oracles
 - 不同的oracles有不同的价格

Weakly | Active Learning

<http://dmirlab.com>

Cost-Sensitive Active Learning

- Oracles are cost-sensitive
 - Different oracles have diverse prices
 - Selecting both instance and oracle
 - Accurate yet cheap annotations



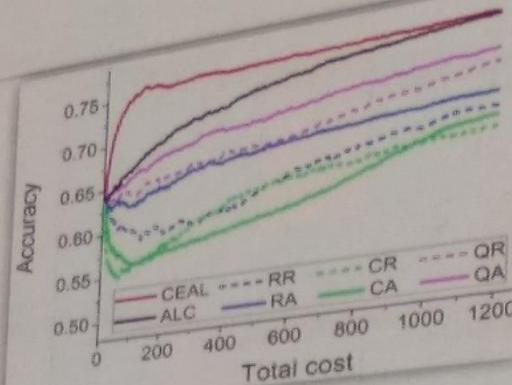
Who is this ?

• Low overall quality
 • Low price
 • Expert for this query

• High overall quality
 • High price
 • Less familiar with it

[Huang et al., IJCAI 2017]

37





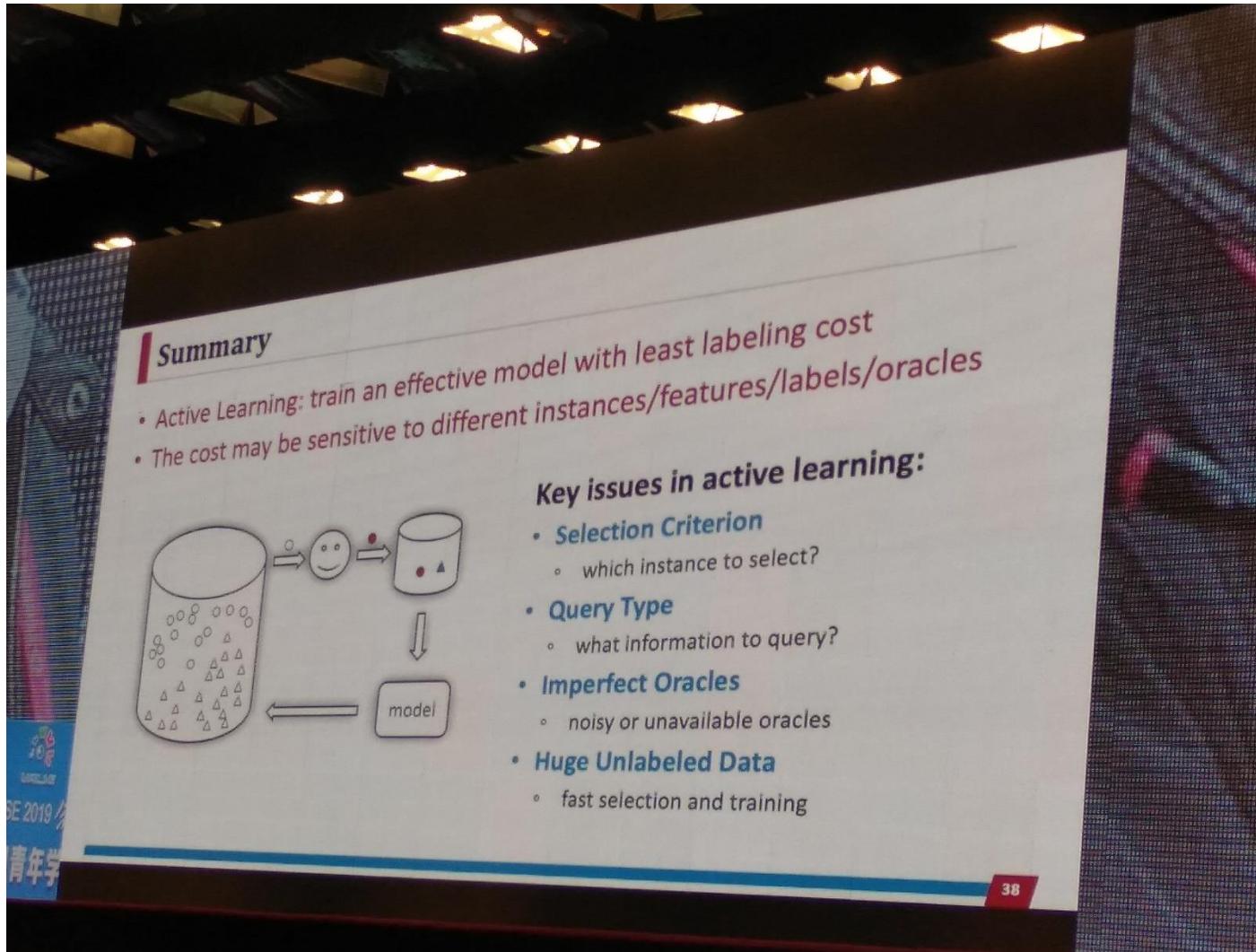
<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲

Weakly | Active Learning

<http://dmirlab.com>



Summary

- Active Learning: train an effective model with least labeling cost
- The cost may be sensitive to different instances/features/labels/oracles

Key issues in active learning:

- **Selection Criterion**
 - which instance to select?
- **Query Type**
 - what information to query?
- **Imperfect Oracles**
 - noisy or unavailable oracles
- **Huge Unlabeled Data**
 - fast selection and training

SE 2019 青年学

38



<http://dmirlab.com>

Weakly | Active Learning

- Shengjun Huang, NHU
 - Active Learning
- Summary (个人观点, 真伪自辩)



Weakly | Active Learning

- 陈明明
 - 基元属性感知：一种特殊的特征共享，挖掘各任务间的共享性，减少重复数据标记
- 叶齐祥：
 - 粗粒度标记：只给目标物体上画一条线、打一个点、图像级别的标记，让模型自动学习完整的Label



Weakly | Active Learning

<http://dmirlab.com>

- Idea:
 - 卷积激励的显著分布估计: 对使原图产生激活大的区域进行分析
 - GAN生成边界周围的样本: 得到更细粒度的分类边界
 - 多模态数据对齐: 为了衡量图像和文本的相似度, 将图像编码到文本的特征空间中, 或者将文本编码到图像的特征空间中, 以此衡量, 而不是将二者编码到第三个特征空间中衡量
 - 设计好的**loss**是发表论文的好方法(尤其是人脸分析)



<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲

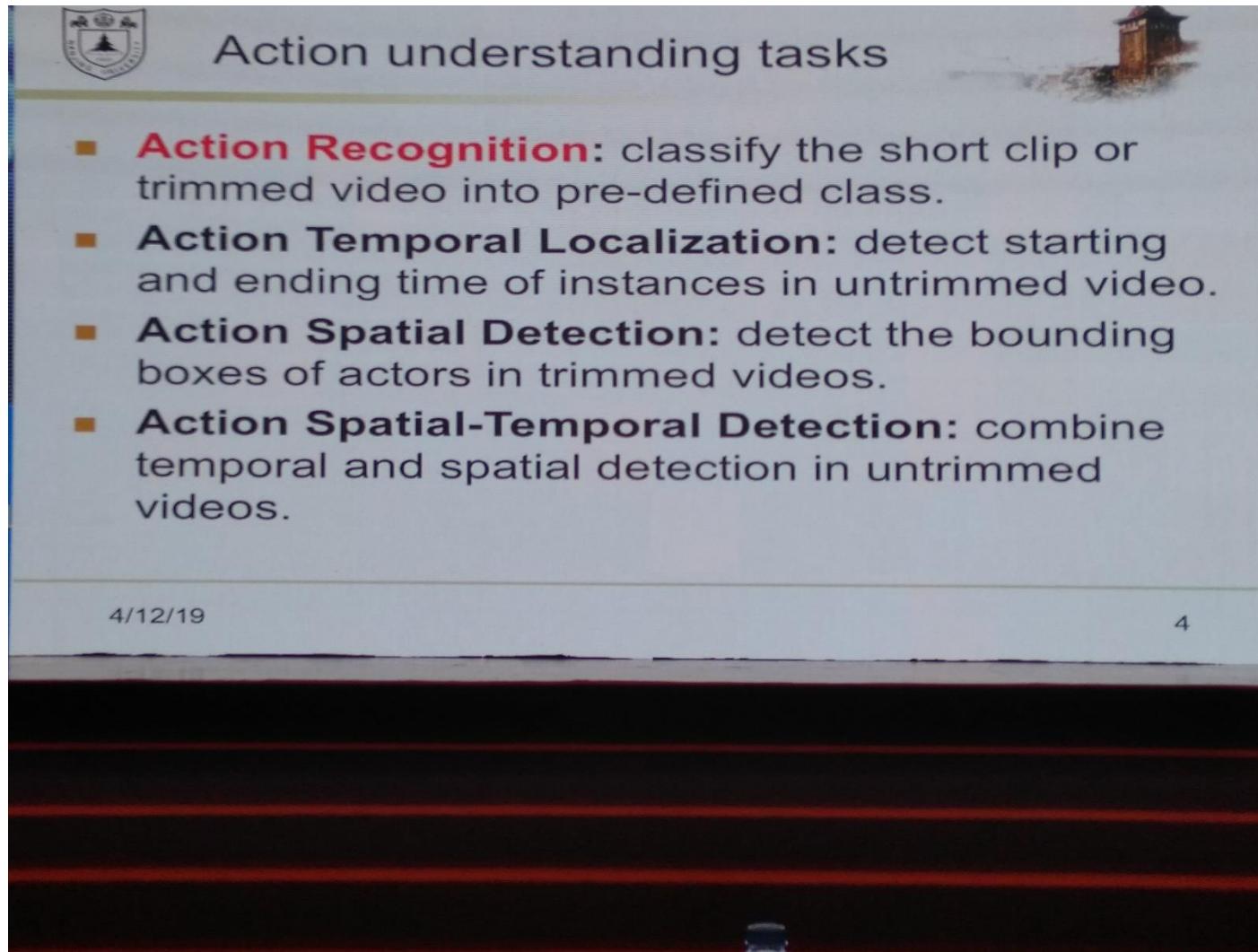


<http://dmirlab.com>

Person

- 通用视频的时序建模与动作识别 (liming wang, NJU)
 - 通用视频的动作识别的深度网络全家福
 - Liming wang的三个工作

Person



Action understanding tasks

- **Action Recognition:** classify the short clip or trimmed video into pre-defined class.
- **Action Temporal Localization:** detect starting and ending time of instances in untrimmed video.
- **Action Spatial Detection:** detect the bounding boxes of actors in trimmed videos.
- **Action Spatial-Temporal Detection:** combine temporal and spatial detection in untrimmed videos.

4/12/19

4

Person

<http://dmirlab.com>



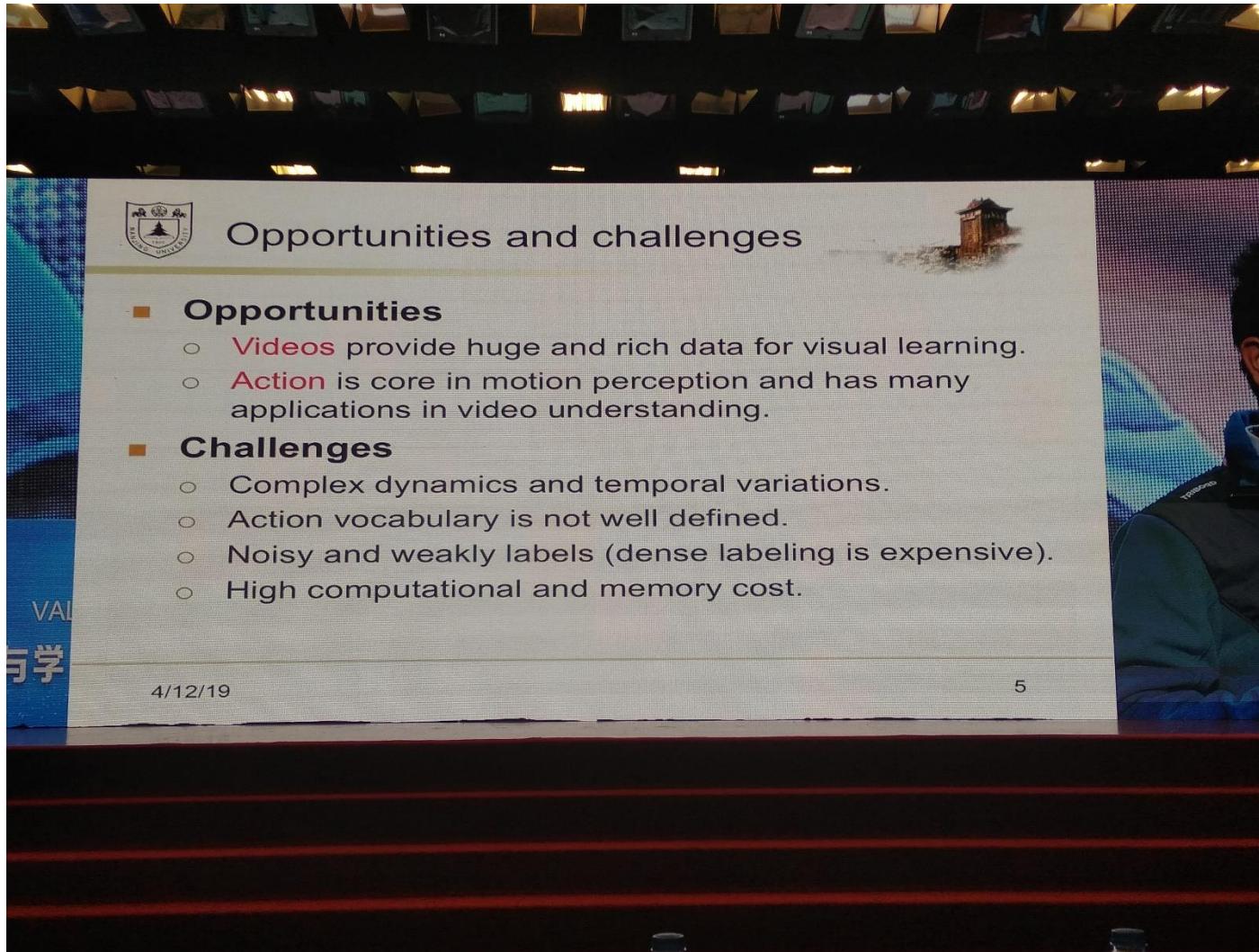
A photograph of a presentation slide titled "Action recognition in videos". The slide features a collage of numerous small video frames showing various actions like running, jumping, and playing sports. Below the collage is a bulleted list of three categories of action recognition datasets:

- 1. Action recognition “in the lab”: KTH, Weizmann etc.
- 2. Action recognition “in TV, Movies”: UCF Sports, Hollywood etc.
- 3. Action recognition “in Web Videos”: HMDB, UCF101, THUMOS, ActivityNet etc.

At the bottom left, there is a small vertical logo for "Al". At the bottom center, the date "4/12/19" is displayed. On the right side, the number "2" is visible, likely indicating the slide number. A small caption at the bottom right reads: "Haroon Idrees et al. The THUMOS Challenge on Action Recognition for Videos "in the Wild", in Computer Vision and Image Understanding (CVIU), 2017."

Person

<http://dmirlab.com>



The image shows a presentation slide titled "Opportunities and challenges". The slide is framed by a decorative border featuring a repeating pattern of colored shapes (triangles and rectangles). At the top left is a university crest. On the right side, there is a small illustration of a traditional Chinese building. The slide content is organized into two main sections: "Opportunities" and "Challenges", each with a bulleted list of points. The footer of the slide includes the date "4/12/19" and the number "5".

Opportunities and challenges

- Opportunities
 - Videos provide huge and rich data for visual learning.
 - Action is core in motion perception and has many applications in video understanding.
- Challenges
 - Complex dynamics and temporal variations.
 - Action vocabulary is not well defined.
 - Noisy and weakly labels (dense labeling is expensive).
 - High computational and memory cost.

VAL
与学

4/12/19

5



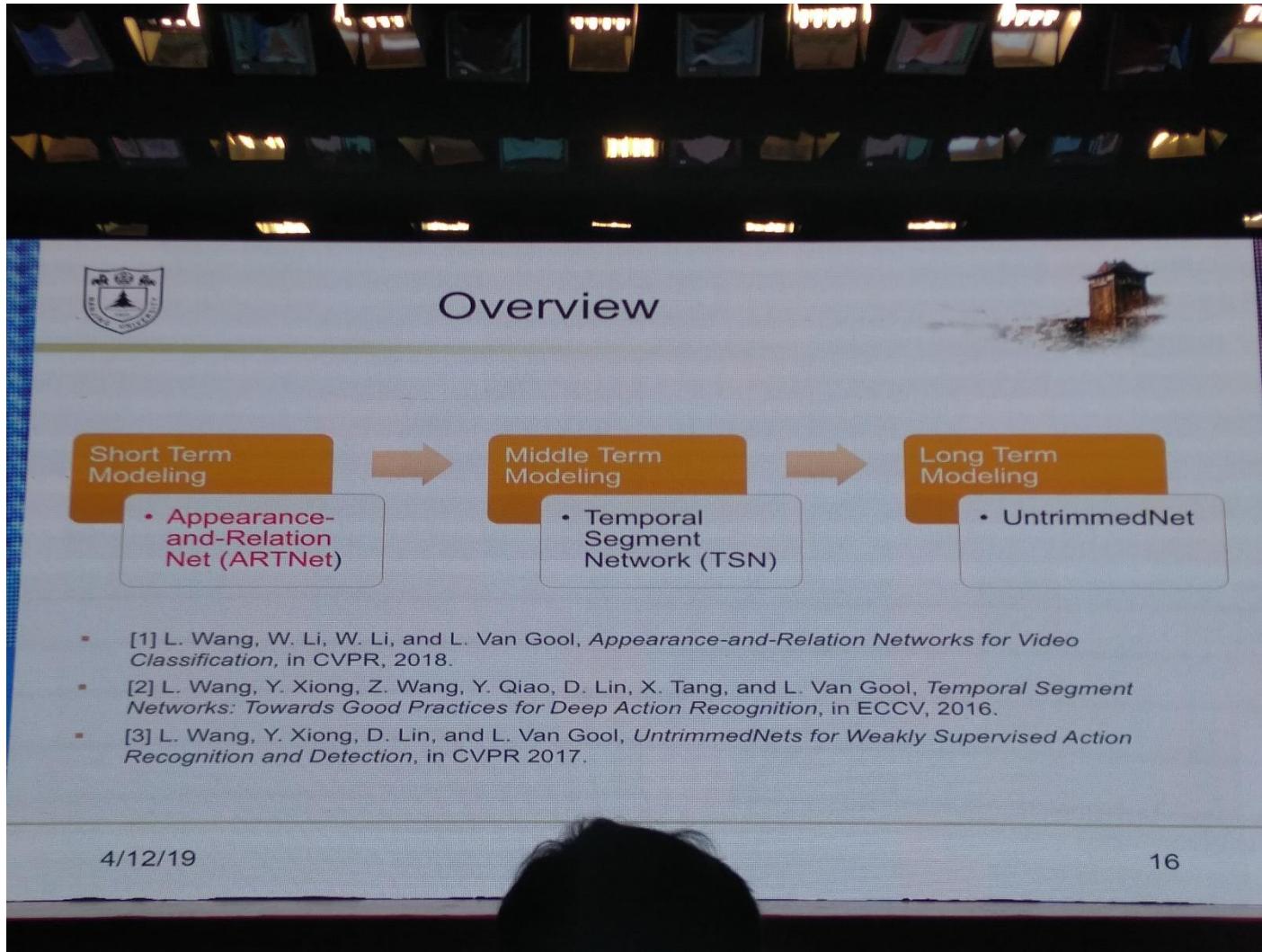
<http://dmirlab.com>

Person

- 代表性论文
 - Large-scale video classification with CNN (FeiFei Li, CVPR2014)
 - Two-Stream CNN for action recognition in videos (NIPS2014)
 - learning spatiotemporal features with 3D CNN (ICCV2015)
 - TDD (liming wang, CVPR2015)
 - Real-time action recognition with enhanced motion vector CNNs (CVPR2016)
 - Two Stream I3D (CVPR2017)
 - R(2+1)D (CVPR2018)
 - SlowFast Networks (kaiming he, CVPR2019)

Person

<http://dmirlab.com>



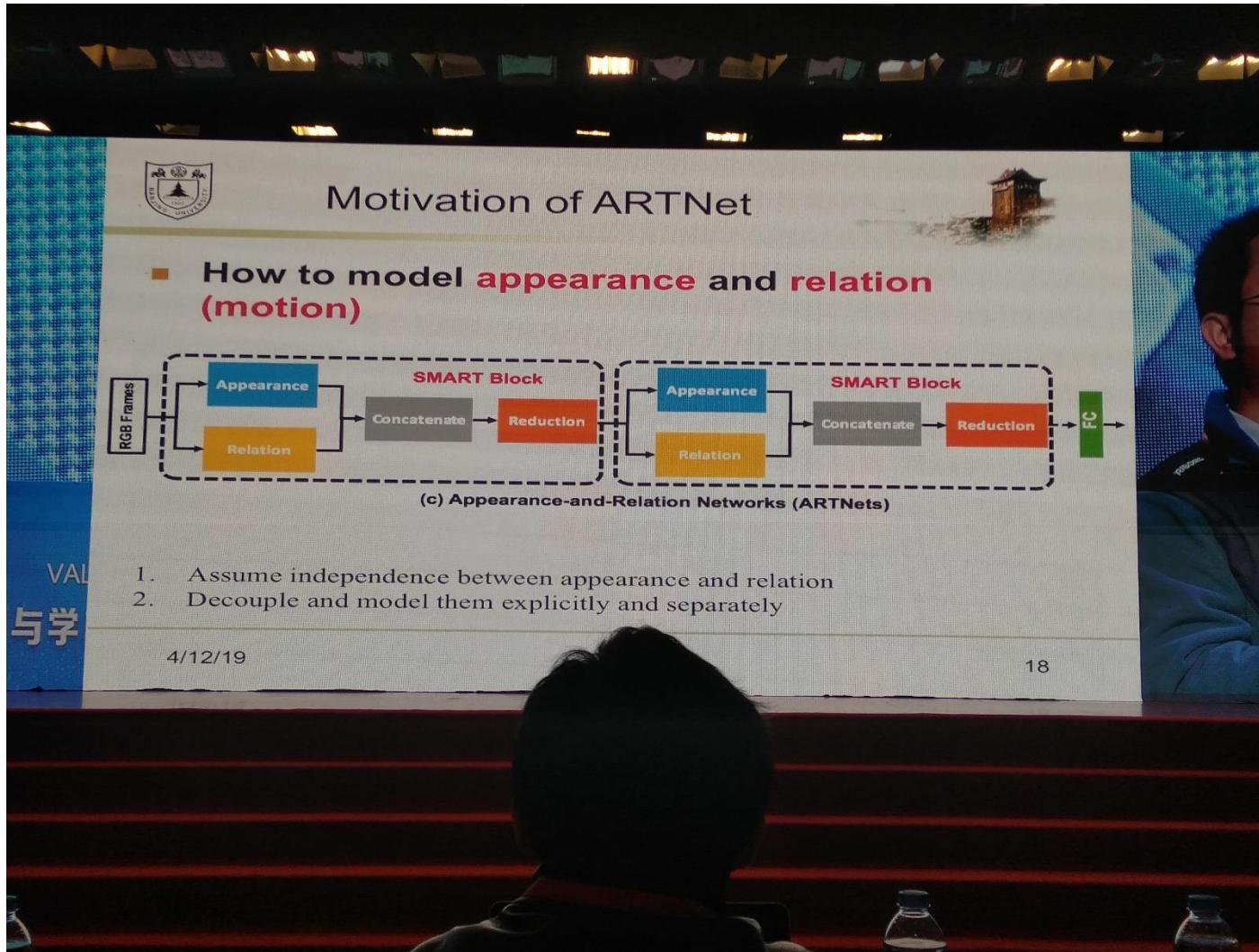
The slide has a decorative header with a university crest on the left and a traditional Chinese painting of a pavilion on the right. The title "Overview" is centered above three orange rectangular boxes connected by arrows. The first box contains "Short Term Modeling" and "Appearance-and-Relation Net (ARTNet)". The second box contains "Middle Term Modeling" and "Temporal Segment Network (TSN)". The third box contains "Long Term Modeling" and "UntrimmedNet". Below the boxes is a list of three academic references:

- [1] L. Wang, W. Li, W. Li, and L. Van Gool, *Appearance-and-Relation Networks for Video Classification*, in CVPR, 2018.
- [2] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool, *Temporal Segment Networks: Towards Good Practices for Deep Action Recognition*, in ECCV, 2016.
- [3] L. Wang, Y. Xiong, D. Lin, and L. Van Gool, *UntrimmedNets for Weakly Supervised Action Recognition and Detection*, in CVPR 2017.

At the bottom left is the date "4/12/19" and at the bottom right is the page number "16".

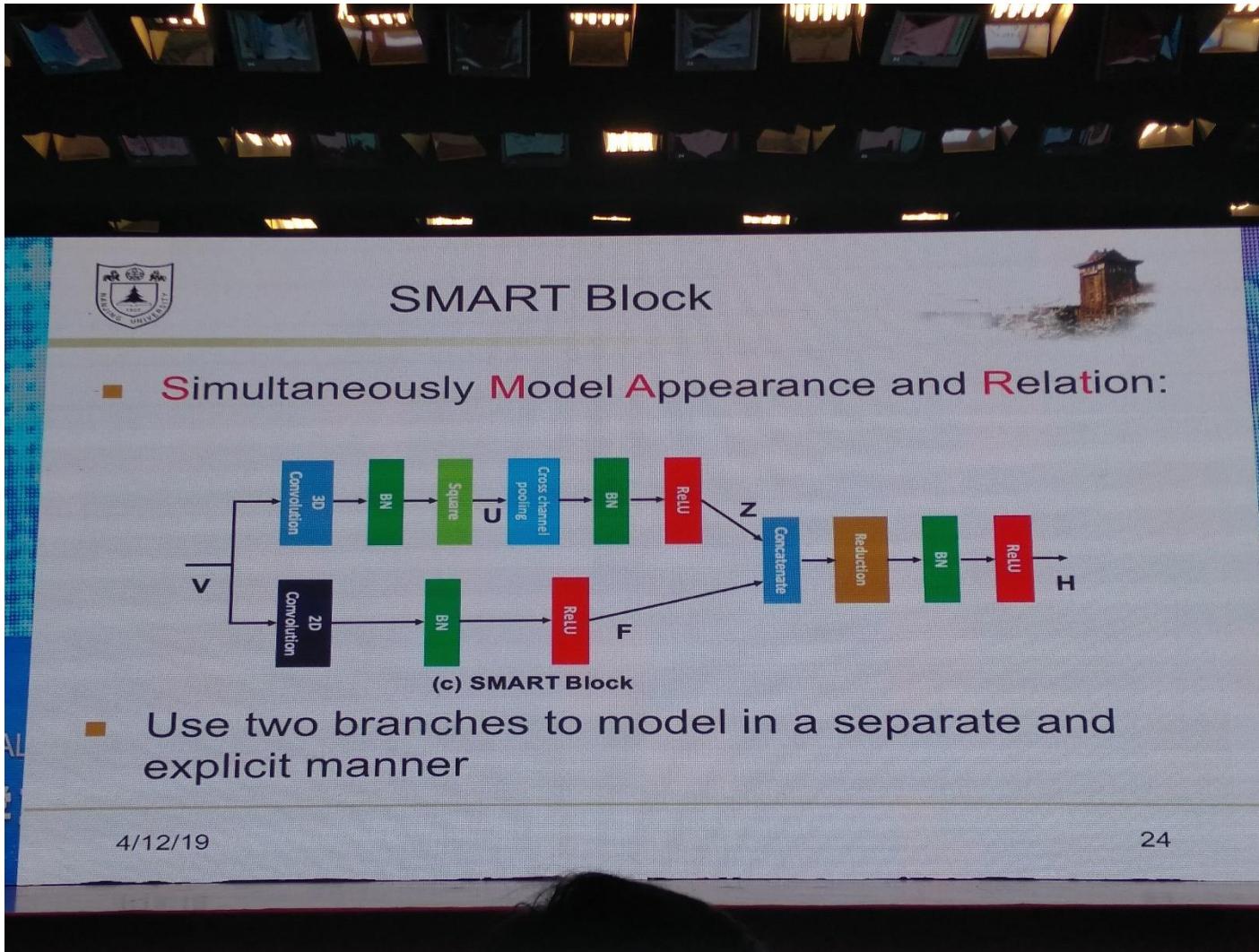
Person

<http://dmirlab.com>



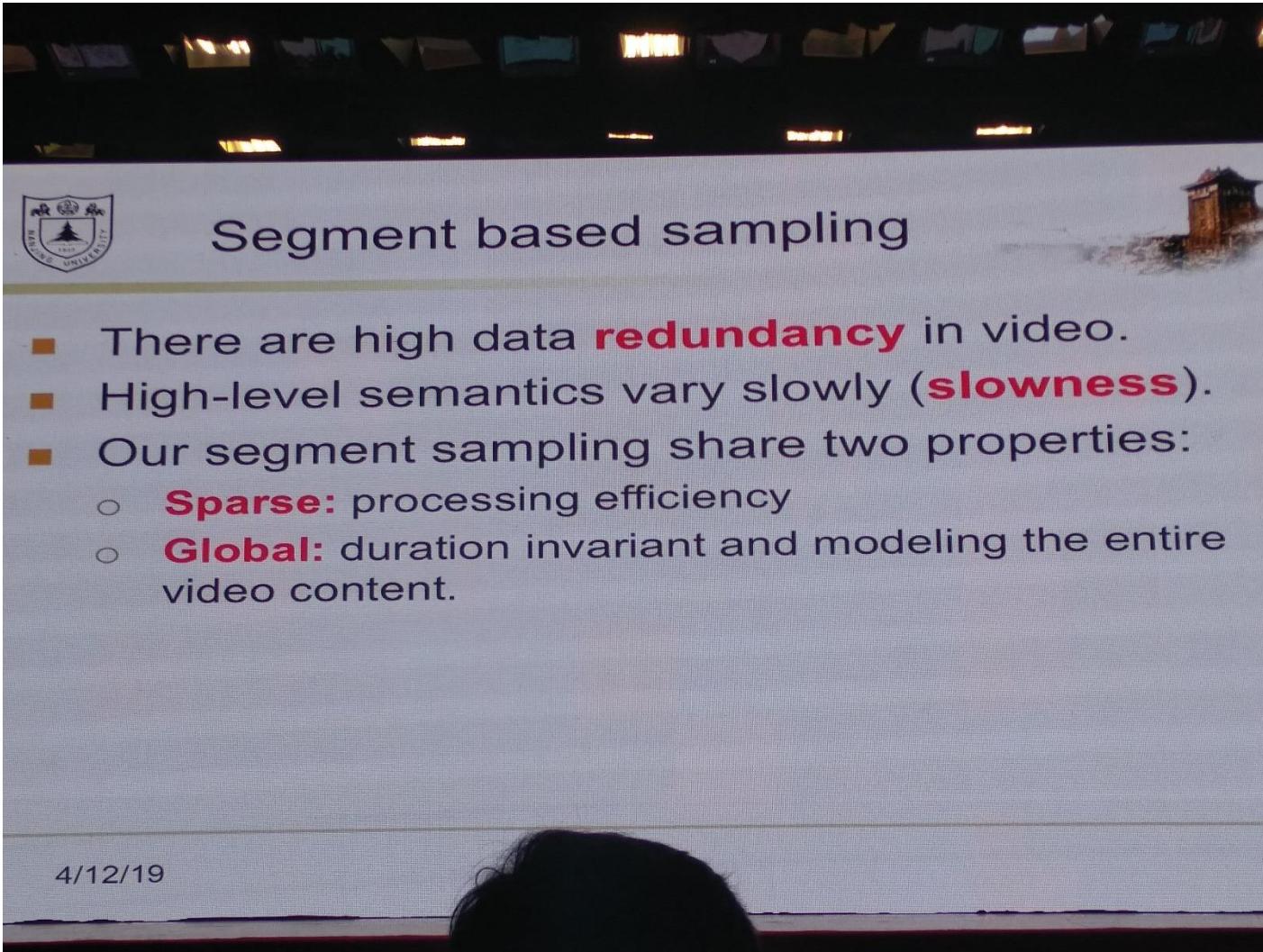
Person

<http://dmirlab.com>



Person

<http://dmirlab.com>



The image shows a presentation slide titled "Segment based sampling". The slide has a dark background with a decorative border at the top featuring a repeating pattern of small, glowing yellow and orange shapes. In the upper left corner, there is a logo for Tsinghua University, which includes a shield with a tree and the university's name in Chinese and English. The main title "Segment based sampling" is centered in a large, dark font. To the right of the title, there is a small, faint illustration of a traditional Chinese building with a tiled roof. The slide contains a bulleted list of five points, each preceded by a square bullet. The text is in a dark font, except for the words "redundancy", "slowness", "Sparse", and "Global", which are highlighted in red. The bottom left corner of the slide displays the date "4/12/19".

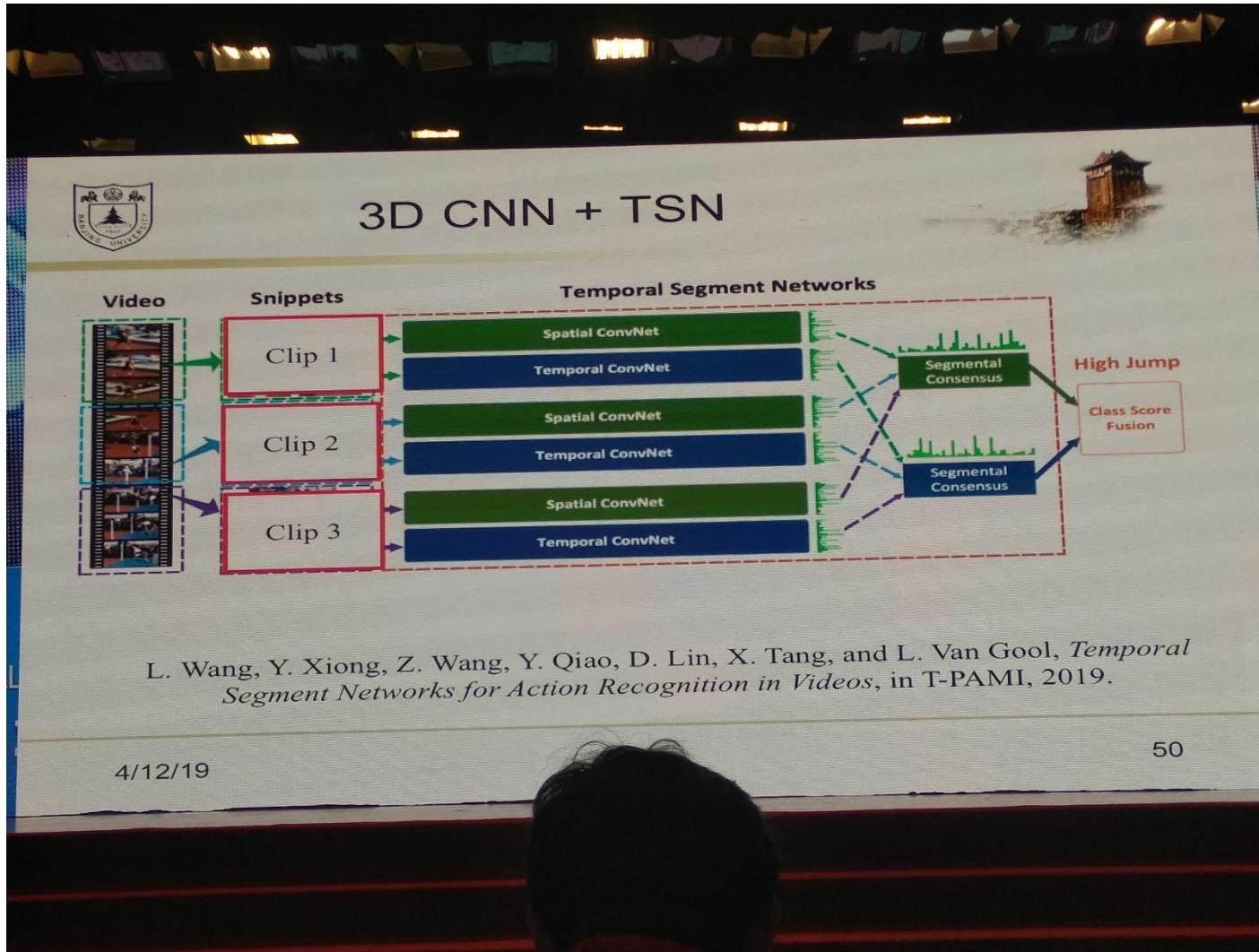
- There are high data **redundancy** in video.
- High-level semantics vary slowly (**slowness**).
- Our segment sampling share two properties:
 - **Sparse**: processing efficiency
 - **Global**: duration invariant and modeling the entire video content.

4/12/19

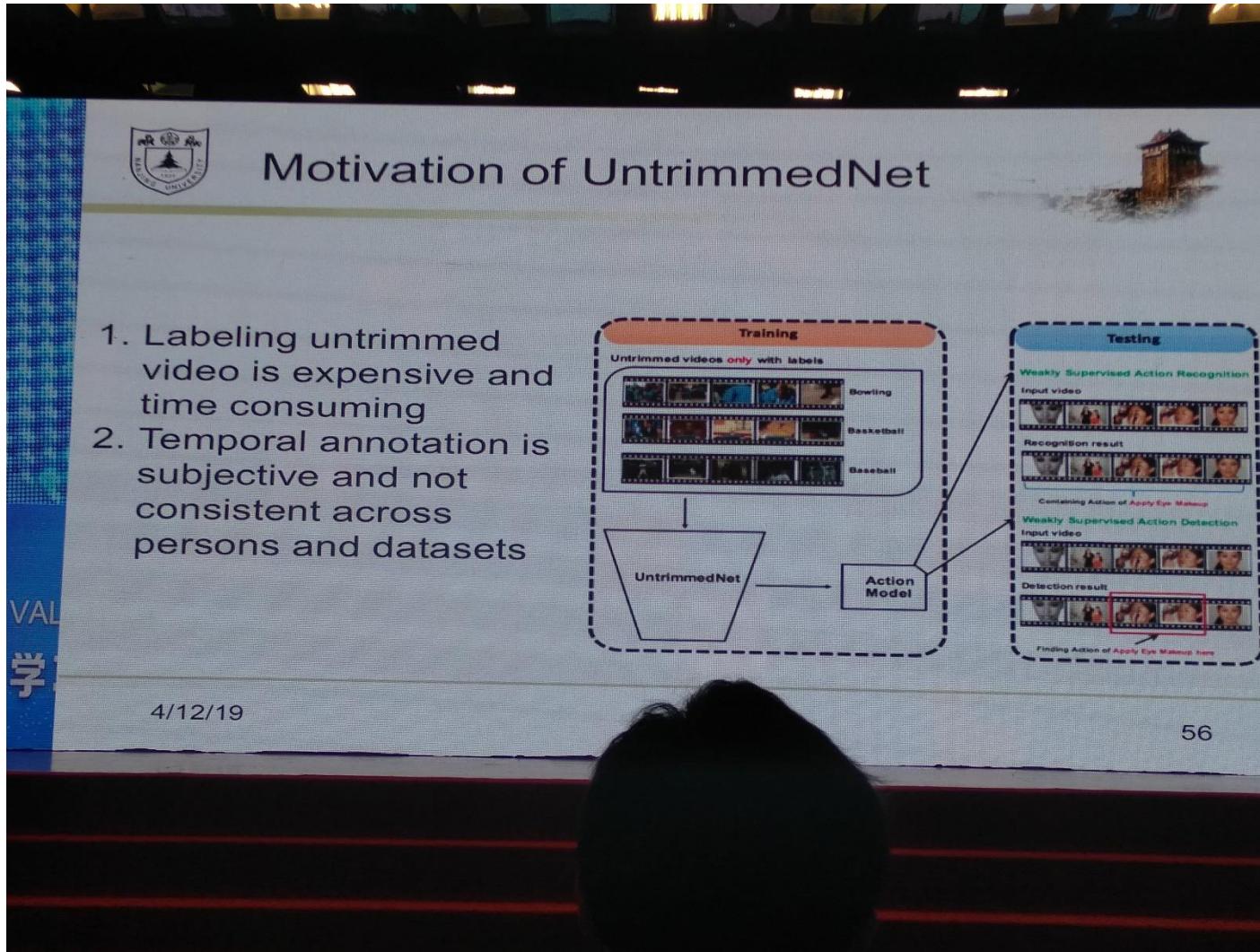
3

Person

<http://dmirlab.com>



Person





<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲



Knowledge | Resoning

<http://dmirlab.com>

- 基于知识驱动的行为理解(Ceiwu Lu, SJU)
- Towards X visual reasoning (hanwang zhang, NTU)



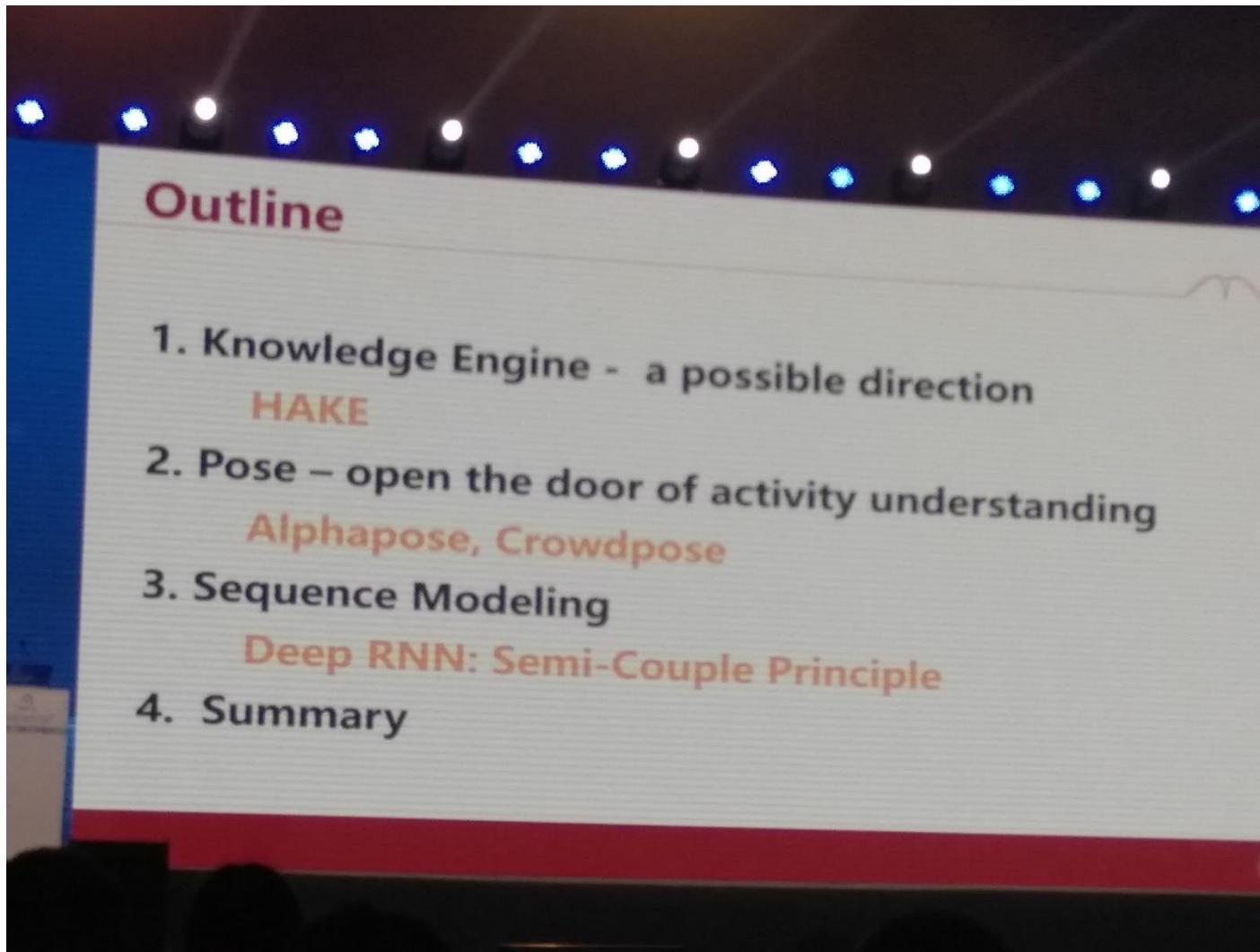
Knowledge | Resoning

<http://dmirlab.com>

- 基于知识驱动的行为理解(Ceiwu Lu, SJU)
- Towards X visual reasoning (hanwang zhang, NTU)

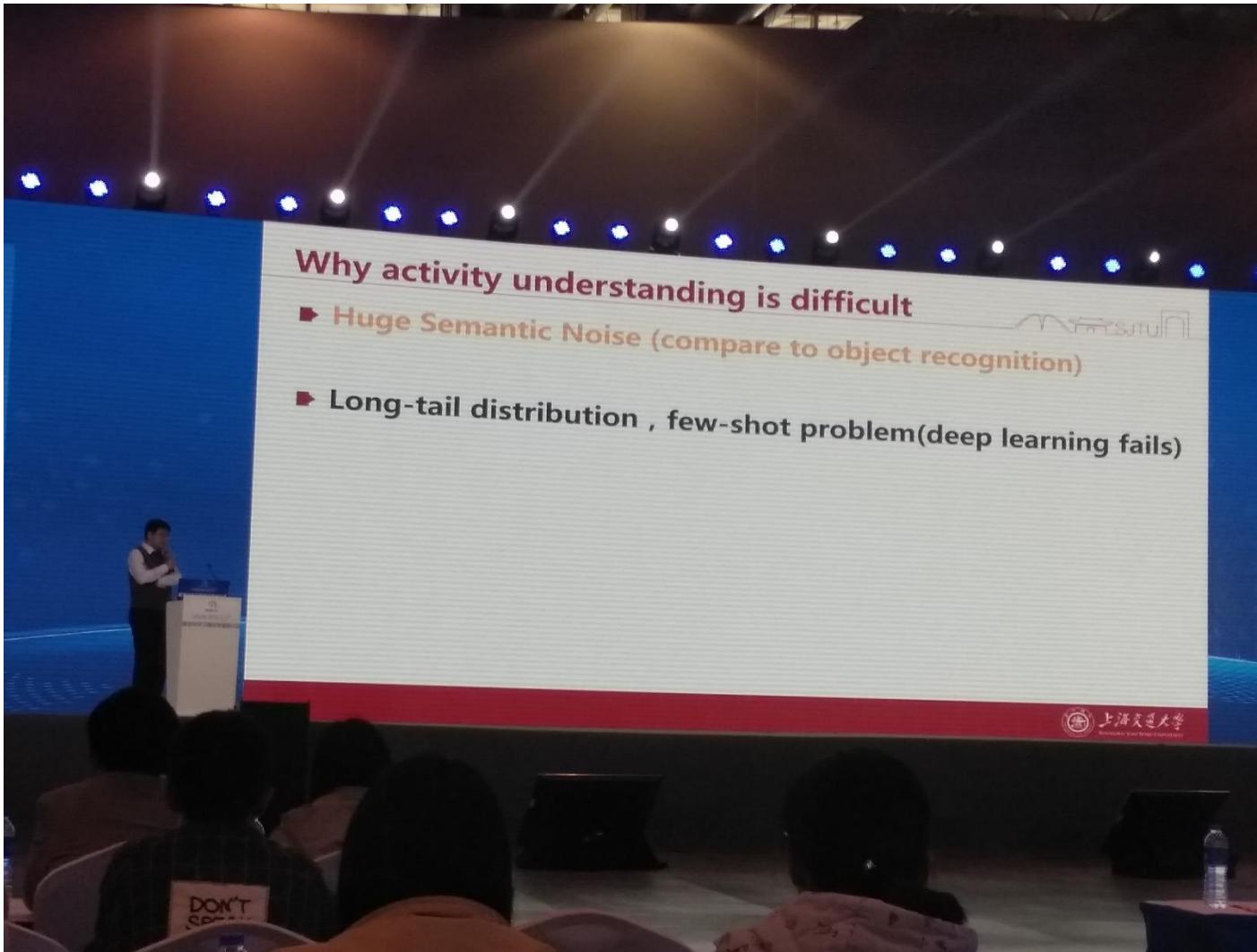
Knowledge | Resoning

<http://dmirlab.com>



Knowledge | Resoning

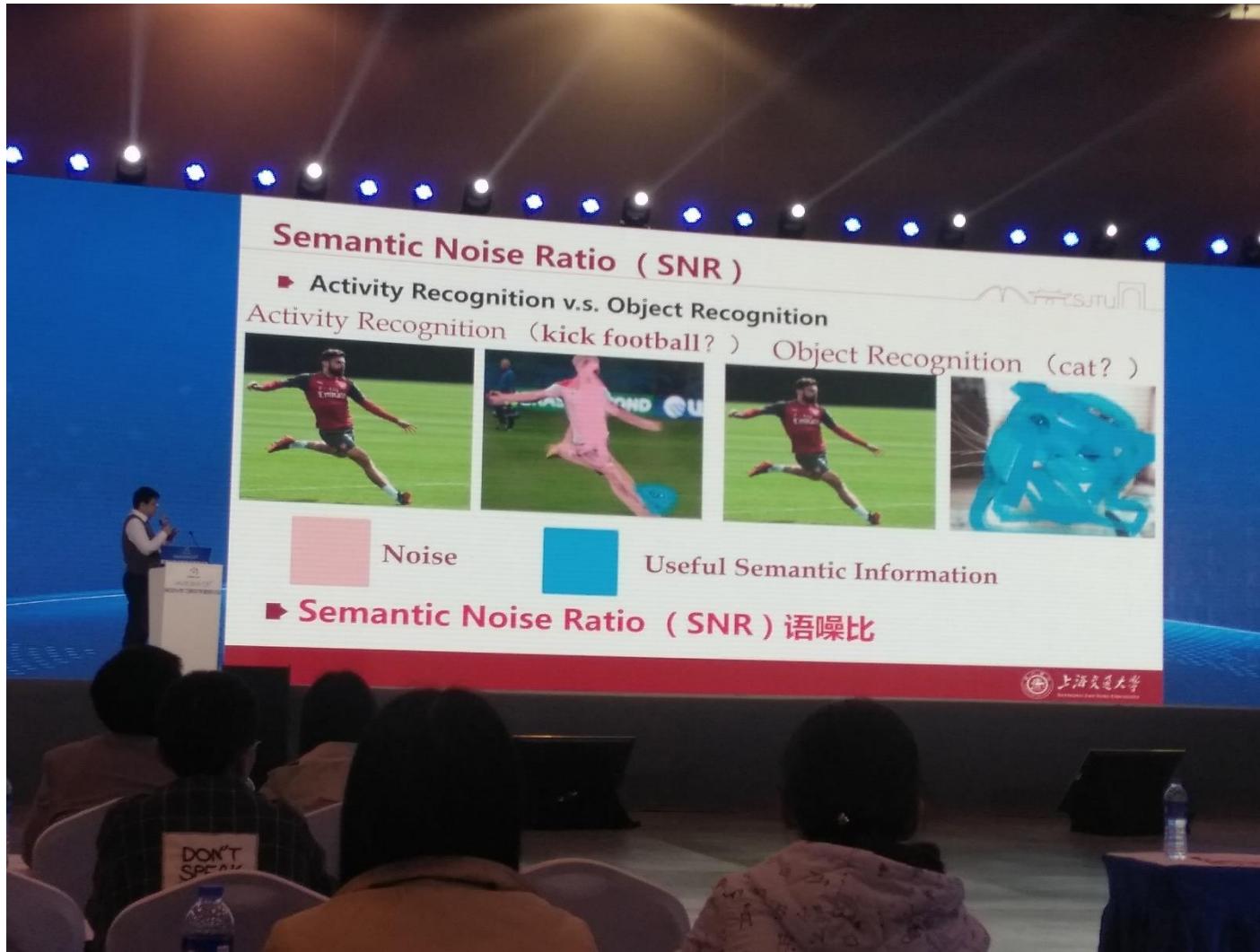
<http://dmirlab.com>



<http://dmirlab.com>

Knowledge | Resoning

<http://dmirlab.com>



<http://dmirlab.com>

Analogy: A Toy Model

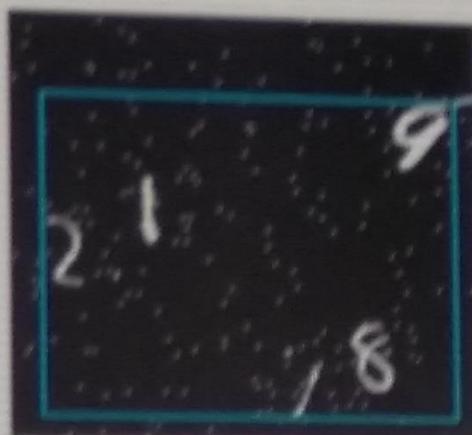
Instance-level Understanding



$f(I) \approx$

Labels:
person-kick-football

Analogy



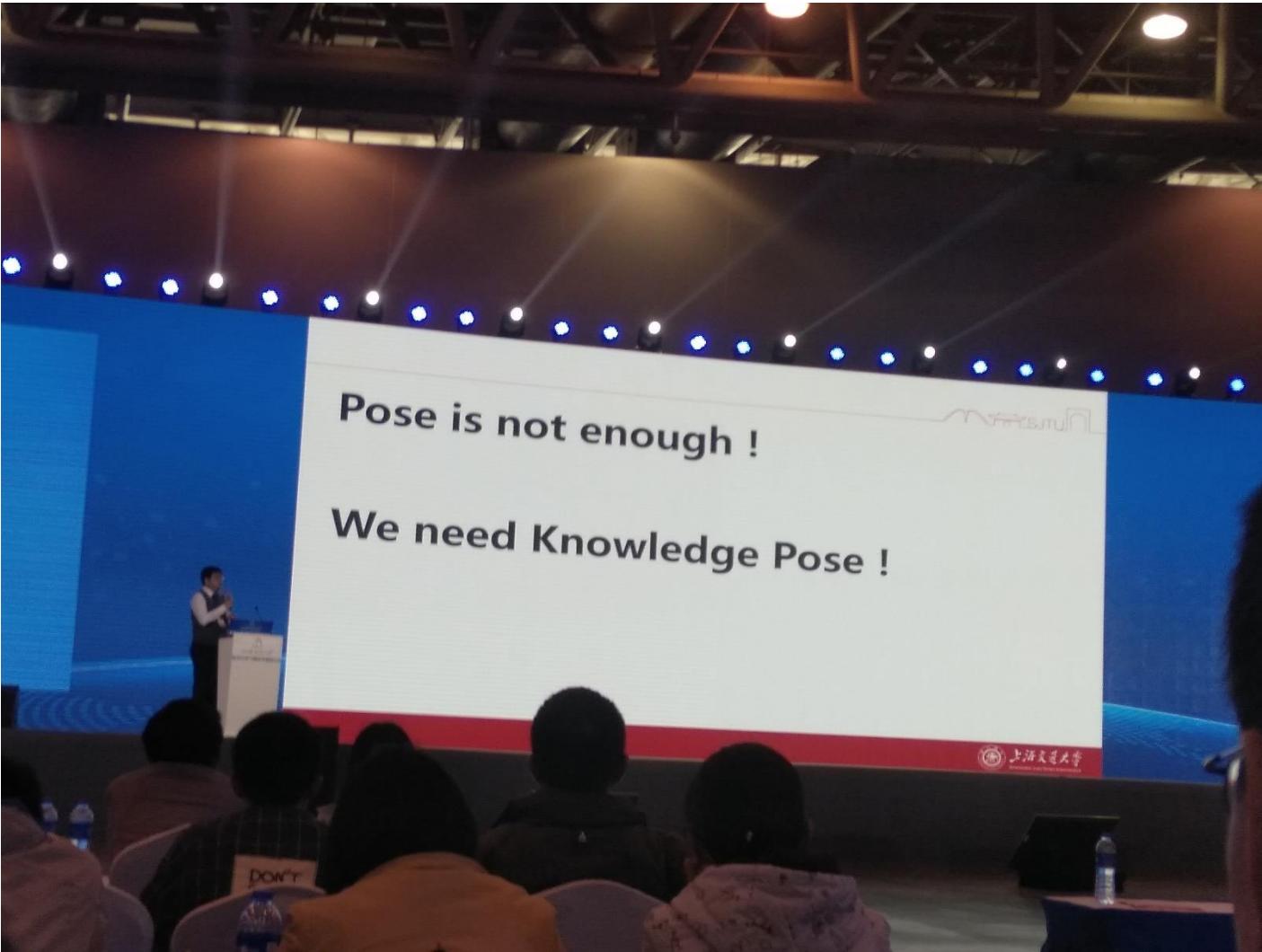
$g(x) \approx$

Labels:
17

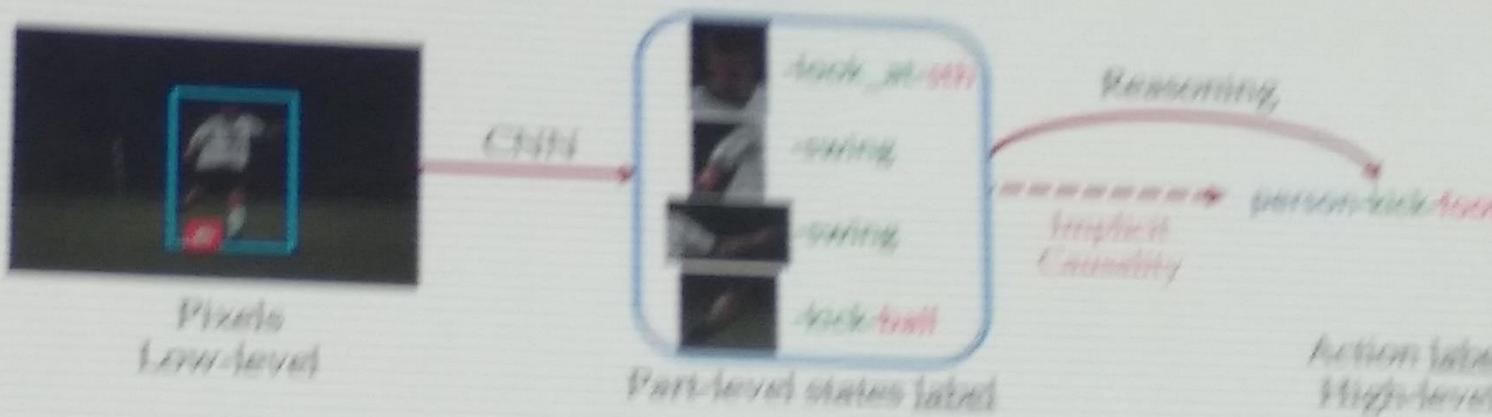
Assumption: $g(\star)$ is
"The sum of top-2 values, e.g. $9+8=17$ "
model \rightarrow fitting $\rightarrow g(\star)$

Knowledge | Resoning

<http://dmirlab.com>



Two-stage Paradigm: Overview



Analogy: A Toy Model

A analogy to human activity recognition:

Part-level Understanding



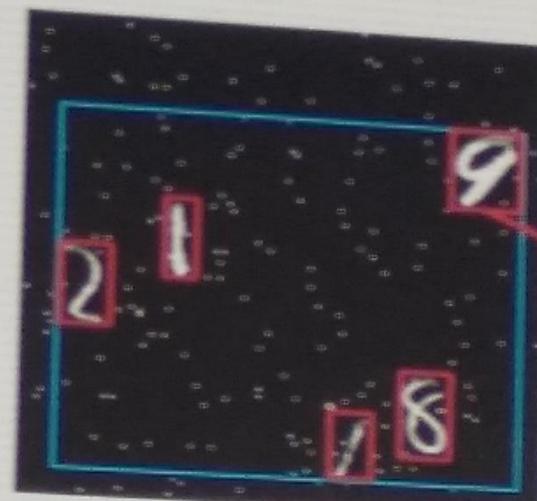
"Part Bounding Box"
from pose keypoints,
Annotators → Part States

Analogy

$f(I)$

Part State Labels:

head-look_at_sth right_arm-swing
left_arm-swing right_foot-h



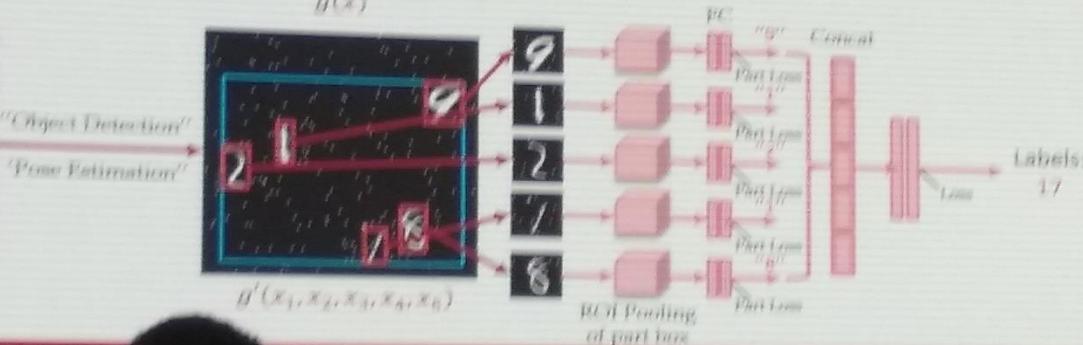
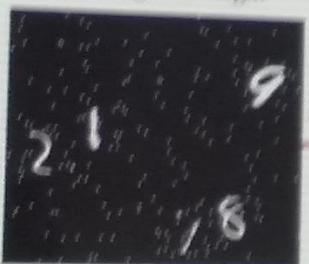
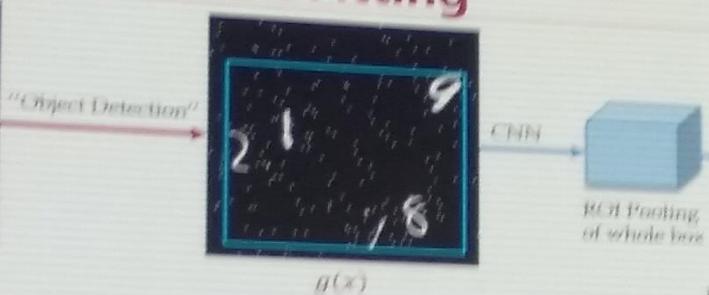
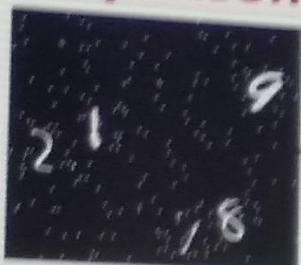
"Part Bound
Annotators -"

$g(x) \rightarrow g'(x_1, x_2, x_3, x_4, x_5)$

Number Labels:
9 8 1 1 2

Labels:
17

Comparison---Function Fitting



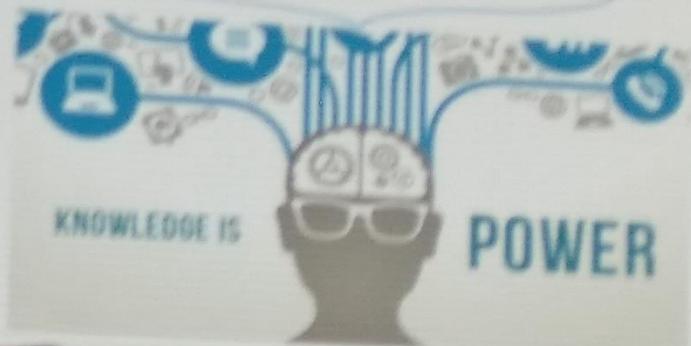
Human Activity Knowledge Engine



HAKE

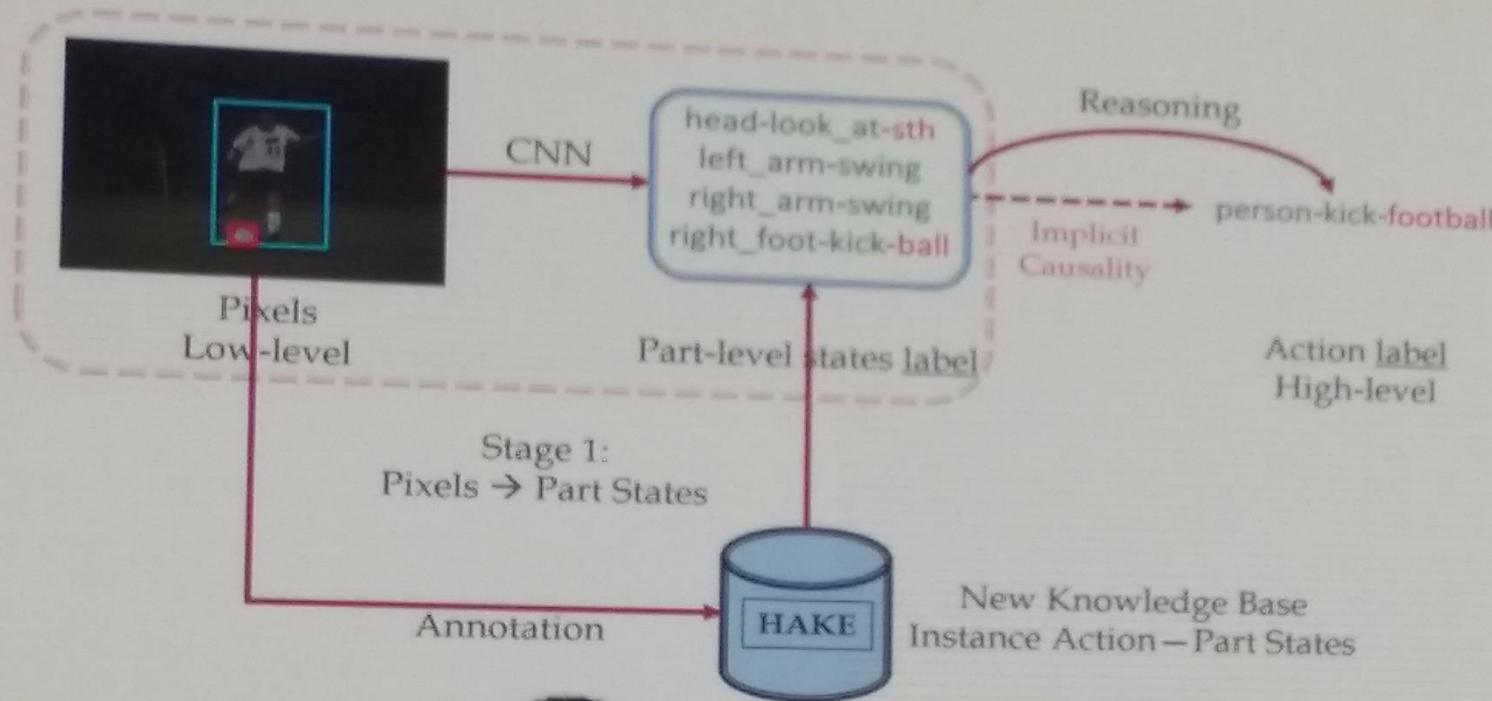
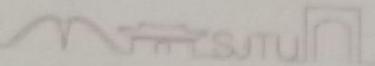


- To see the activity
- To parse the activity
- To understand the activity



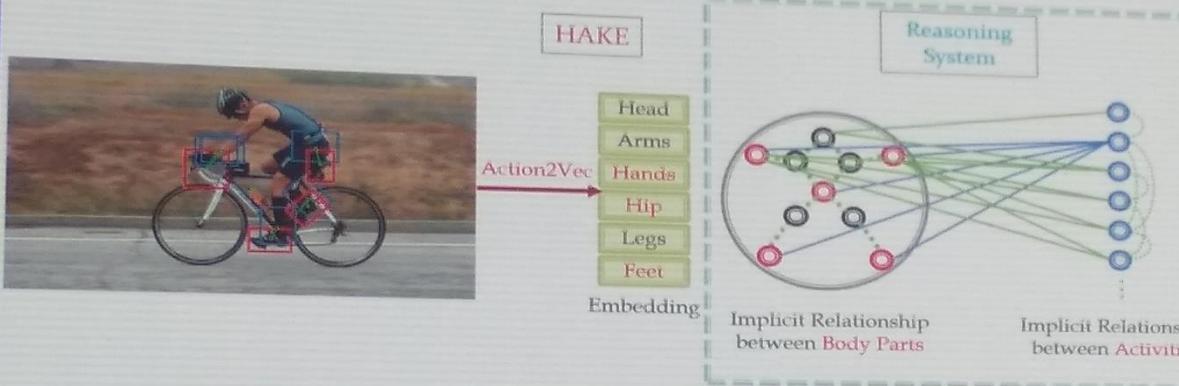
© JBL GROUP

Knowledge Engine Construction



Reasoning via Part States

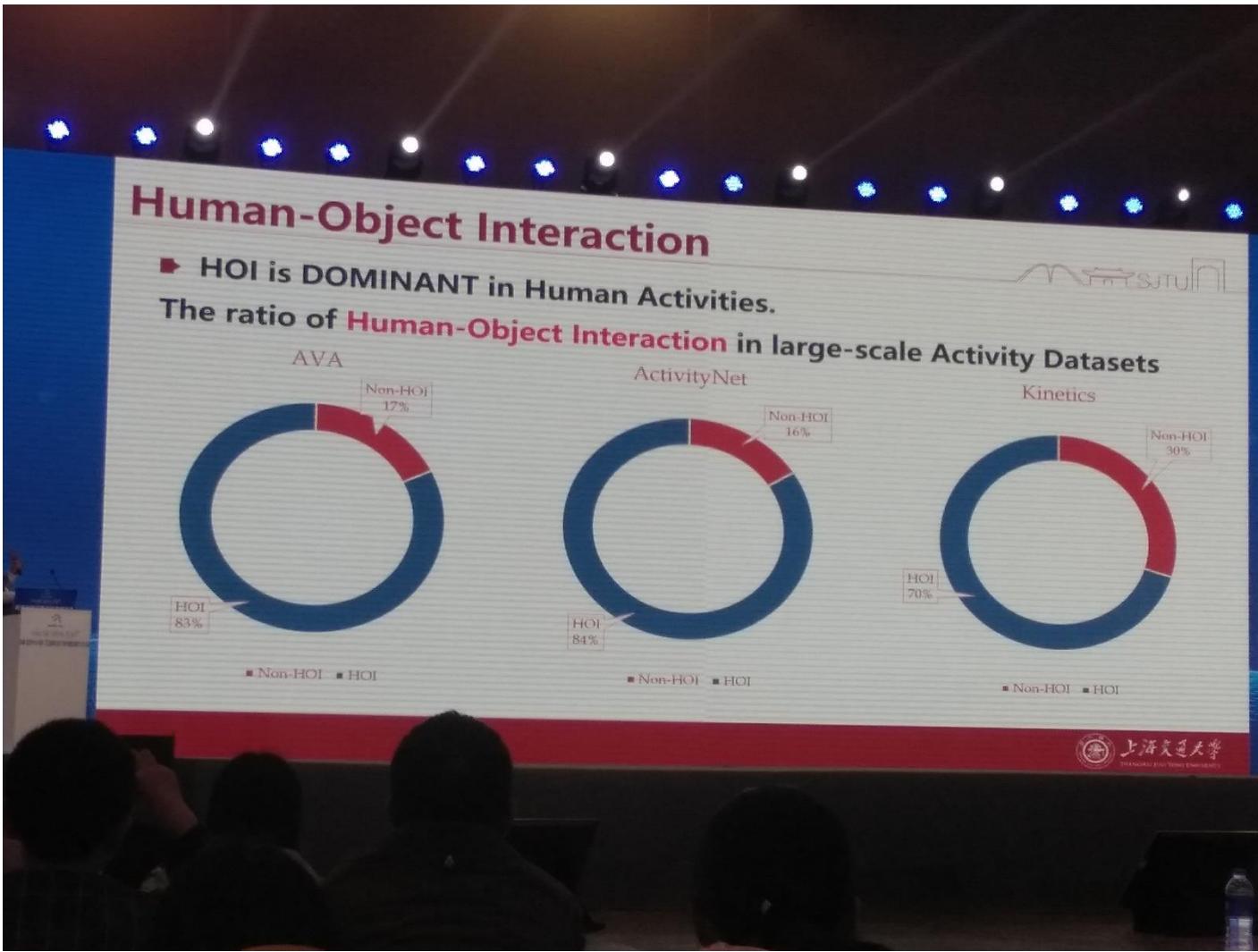
- Activity Knowledge based Reasoning
 - Implicit relationship between activity and part states



- Reasoning from Part States to Instance Activities $R(\{p_i\}_{i=1,2,\dots,10}) = \{A_j\}_{j=N}$
- Taking GNNs, RNNs, FCs as $R(\bullet)$ all show prominent performance improvement

Knowledge | Resoning

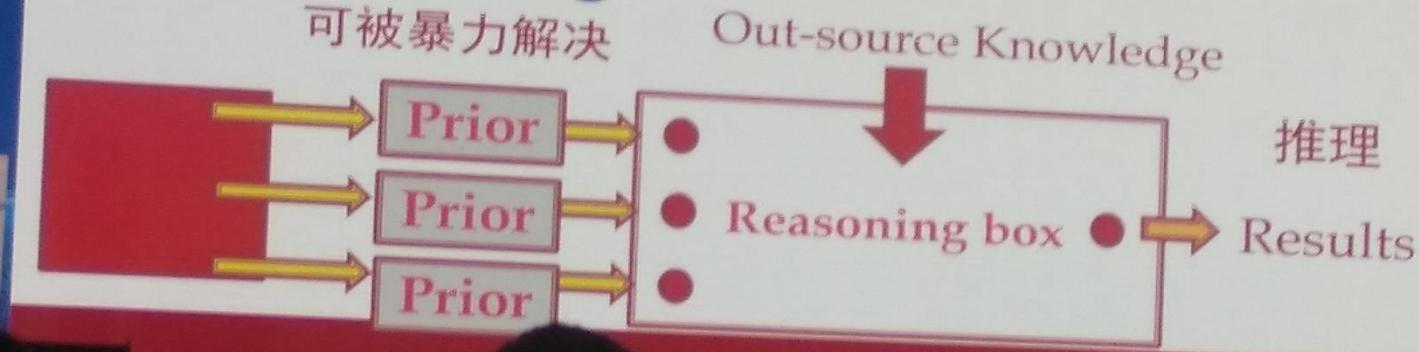
<http://dmirlab.com>



Our Insight : Knowledge Equation

Prior-Reasoning

可被暴力解决



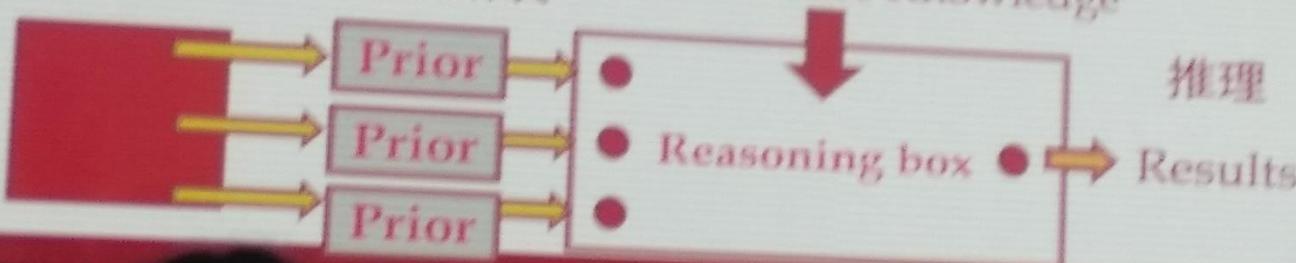
DON'T
SPEAK

Our Insight : Knowledge Equation



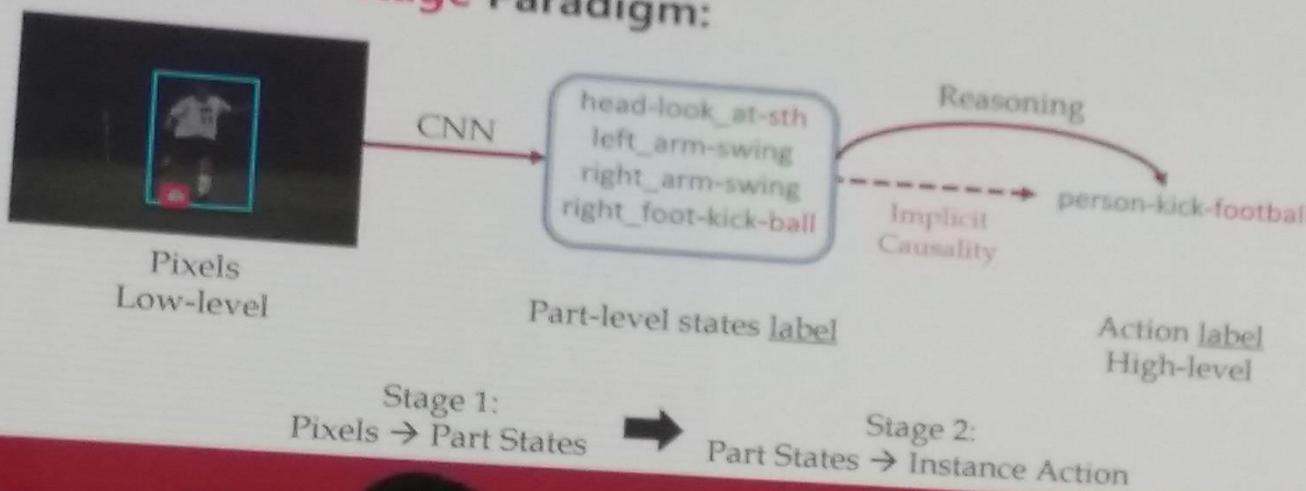
Prior-Reasoning

可被暴力解决



Conclusion

- Activity data is Semantically NOISY!
- Knowledge at Body Part can help to denoise
- HAKE: Human Activity Knowledge Engine
- HAKE based Two-stage Paradigm:





Knowledge | Resoning

<http://dmirlab.com>

- 基于知识驱动的行为理解(Ceiwu Lu, SJU)
- Towards X visual reasoning (hanwang zhang, NTU)

X Visual Reasoning (eXplainable & eXplicit)

Hanwang Zhang 张含望

MReaL Lab (mreallab.github.io)

hanwangzhang@ntu.edu.sg



School of Computer Science and Engineering

NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE



Explainable

- The inductive bias is human-understandable language
- The network structure is interpretable
- e.g., Neural Architecture Search (NAS) is not explainable

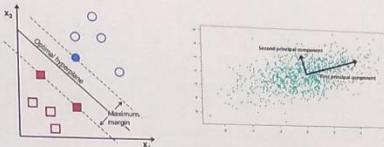


Inductive Bias



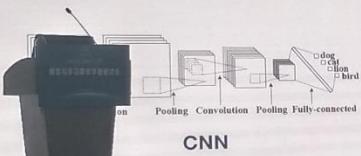
K-NN

Class A
Class B
Unknown class

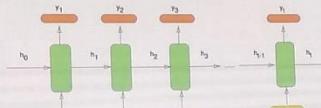


Large Margin

PCA



CNN



RNN



Inductive Bias for Visual Reasoning

- How to define the bias?
 - Where is it from?
 - What is it like?
- How to use the bias?
 - Where to impose it?
 - How to make it computationally feasible?



NANYANG
TECHNOLOGICAL
UNIVERSITY

XVR: Structure Building from 3 Modalities

- Visual Data
 - Re-Engineer the visual feature, visual attention is a special case
- Language
 - Re-Comprehend the language
- Dialog History
 - Re-Compile the history



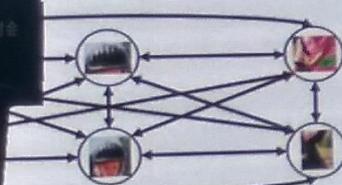
Tree



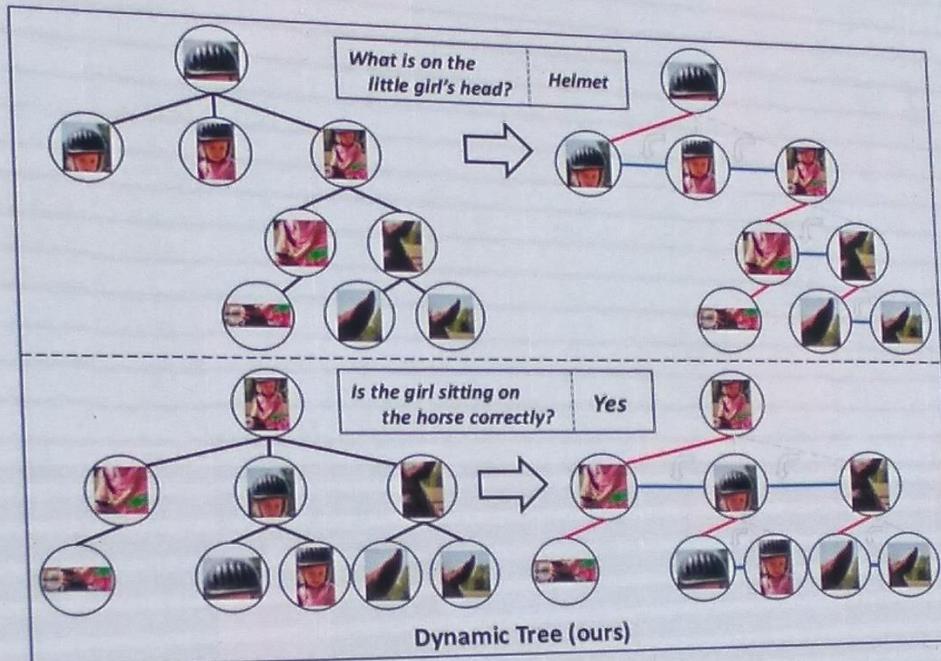
Independent Nodes
[Anderson et al 2018]



Chain
[Zellers et al 2018]

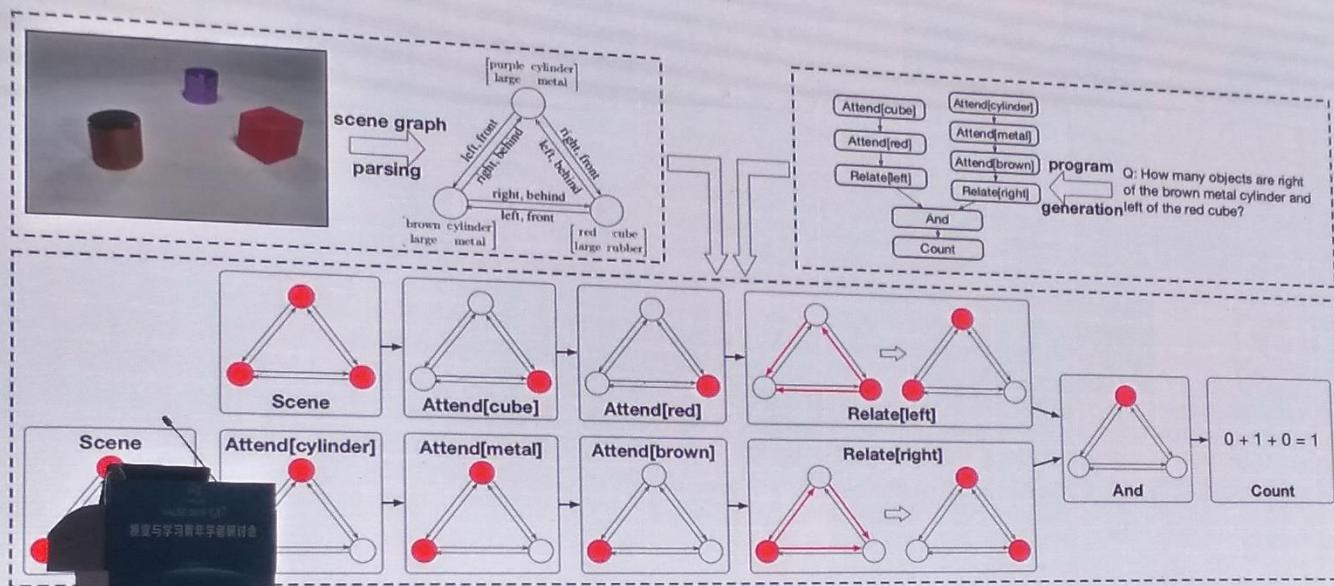


Graph
[Xu et al 2017]



Tang et al. **CVPR'19 oral.** Learning to Compose Dynamic Tree Structures for Visual Contexts

Scene Graph

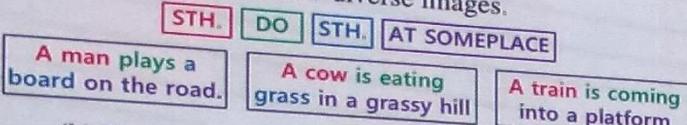


Shi et al. CVPR'19 Explainable and Explicit Visual Reasoning over Scene Graphs

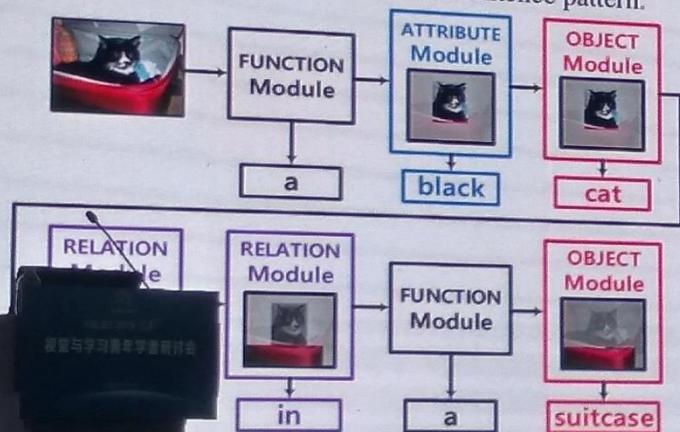
Modular Attention



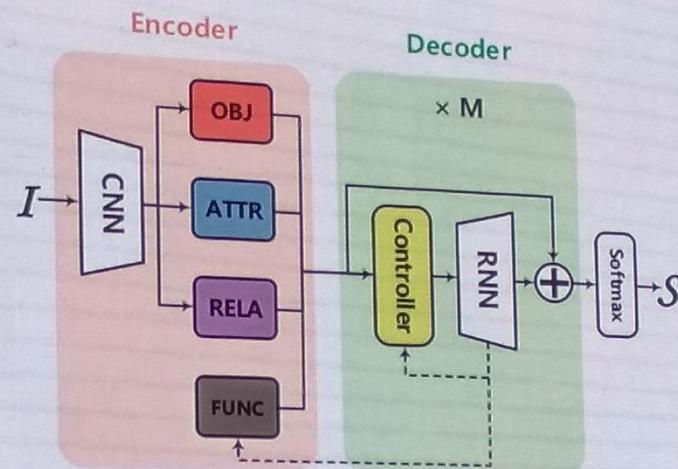
(a) Three diverse images.



(b) Three captions with the same sentence pattern.



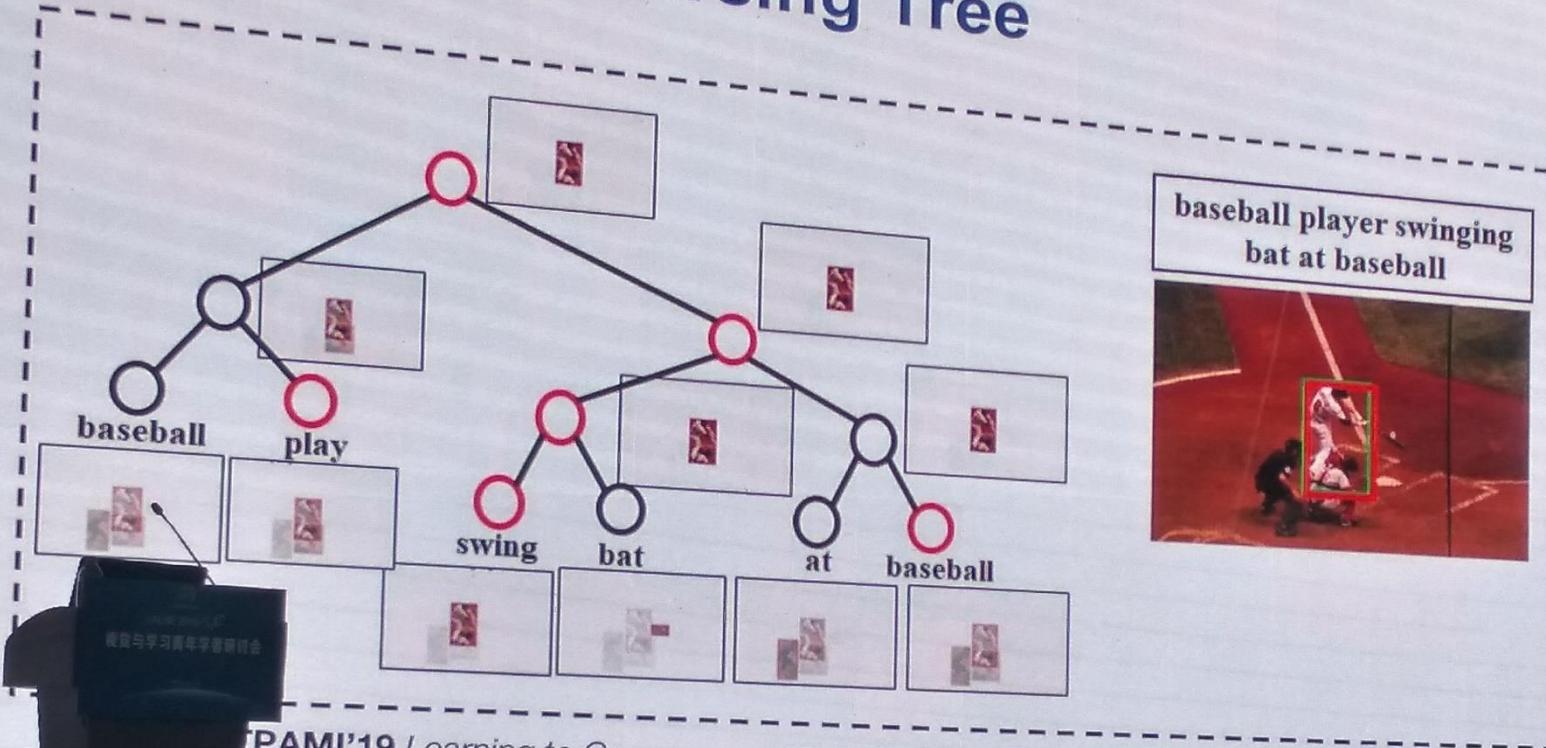
caption generation process of CNM.



C: 127.9 Kaparthy, 126.0 server

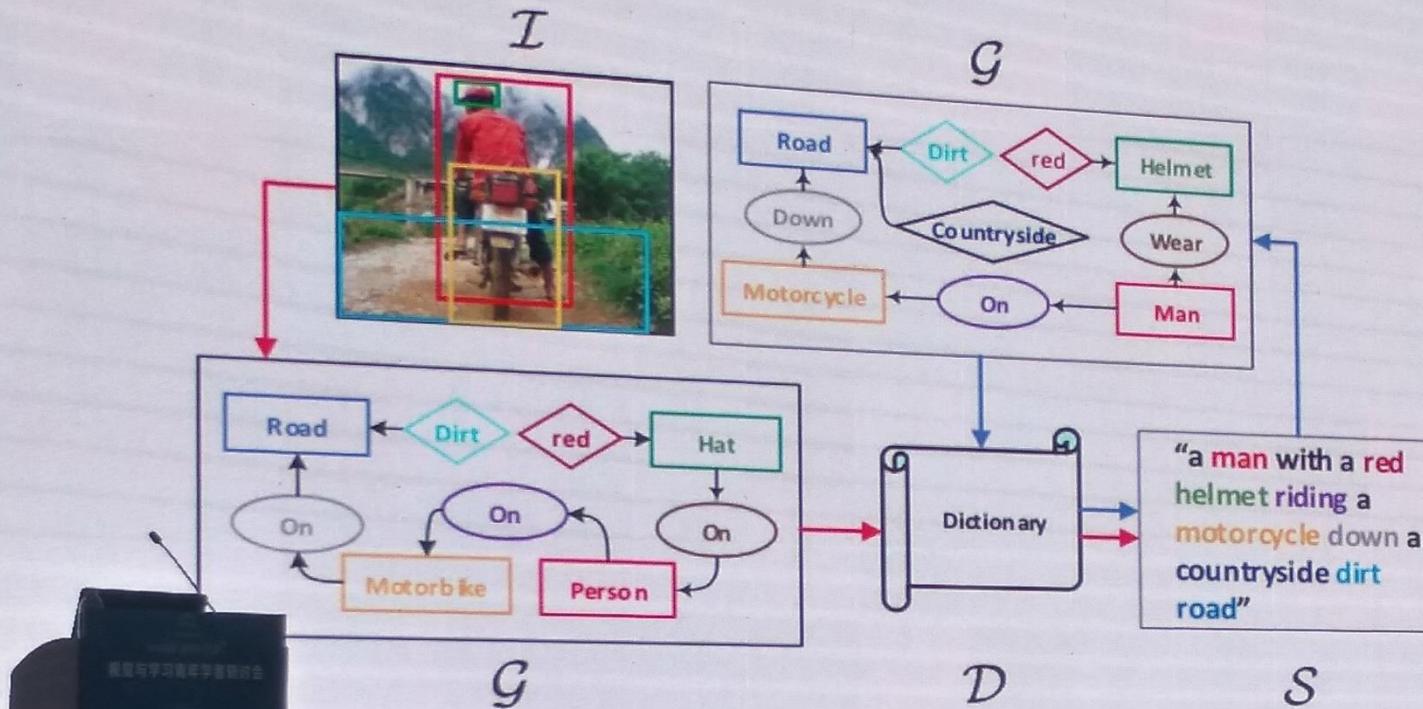
Yang et al. arxiv'19 Learning to Collocate Neural Modules for Image Captioning

Constituency Parsing Tree



PAMI'19 Learning to Compose and Reason with Language Tree Structures for Visual Grounding

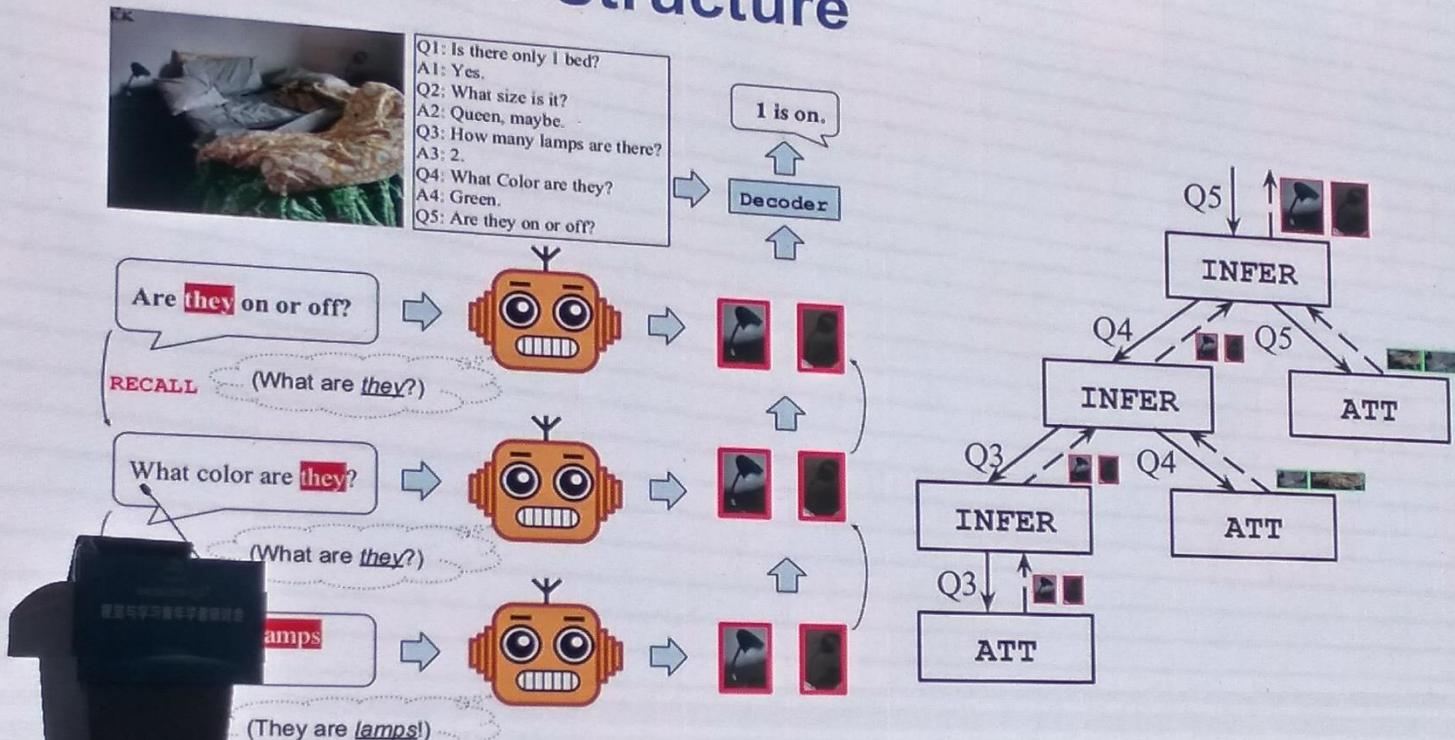
The Graph



C: 129.1 Kaparthy, 126.5 server (code available)

Yang et al. **CVPR'19 oral** Auto-Encoding Scene Graphs for Image Captioning

History-aware Structure



Niu et al. CVPR'19 oral Recursive Visual Attention in Visual Dialog

Summary

Motivation

- Learning → Reasoning
- Reasoning ← Inductive Bias

This talk

- XVR ← Structure from the bias
- Vision, Language, Dialog History

Future

- Comprehension (dataset-independent)
- Novel induction (knowledge discovery)



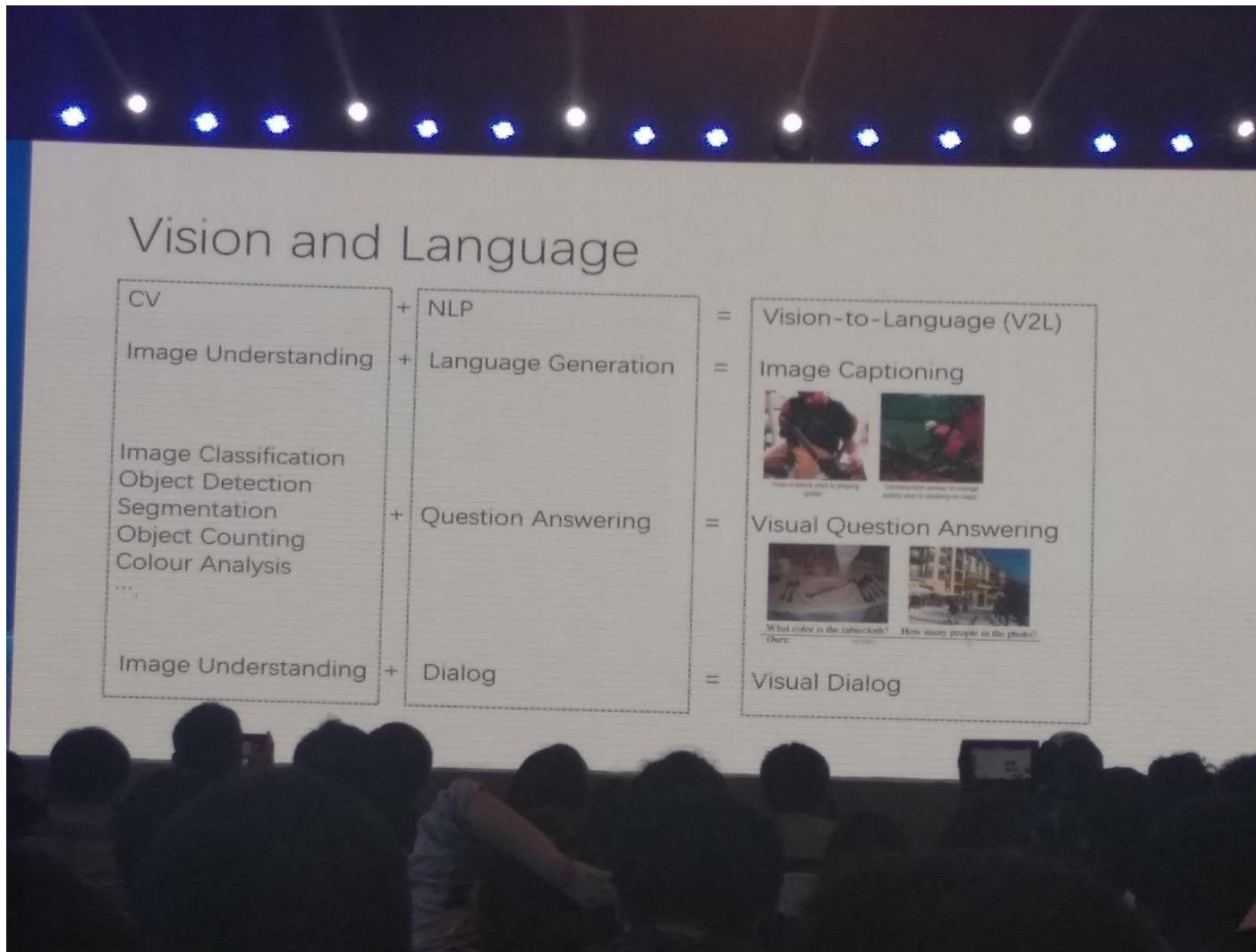


<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲





More Vision-and-Language Tasks

- Referring Expression
 - ReferIt Game, *EMNLP 2014*
 - RefCOCO, *ECCV 2016*
 - GuessWhat?!*, CVPR 2017*
- Visual Dialog
 - VisDial, *CVPR 2017*
 - Image Grounded Conversation, *ACL 2017*
 - Dialog-based Image Retrieval, *NIPS 2018*
- Text 2 image/video



Query from
GuessWhat!:

is it a person? Yes
is the person a man? No
is it the left woman? No
does she wear glasses? Yes

Query from ReferCOCO:
woman in black

Query from ReferCOCO+:
woman wearing glasses with smile
Query from ReferCOCOg:
a women in black playing a game
with her friends

Visual Dialog

A screenshot of a "Visual Dialog" application. At the top, there's a small image of a coffee mug with the caption "A red drinking water cup of a coffee mug". Below it, a user message says "What color is the mug?". The system response is "The mug only has one color". Another user message asks "Are there other items on the table?". The system response is "Yes, magazines, books, toaster and basket, and a plane". On the left, there's a sidebar with a "VisDial" button and a "GuessWhat!" button.



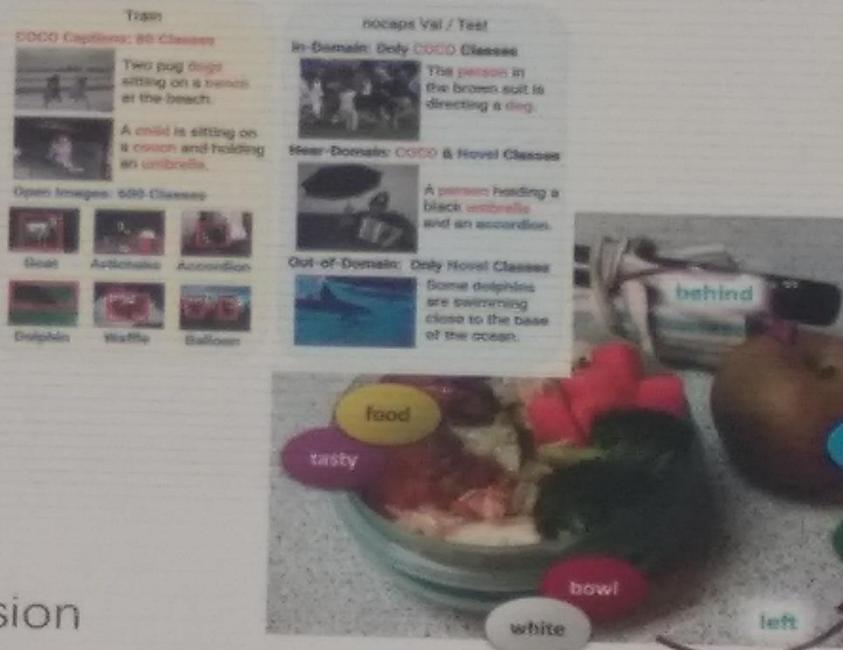
User1: My son is ahead and surprised!

User2: Did he end up winning the race?

User1: Yes he won, he can't believe it!

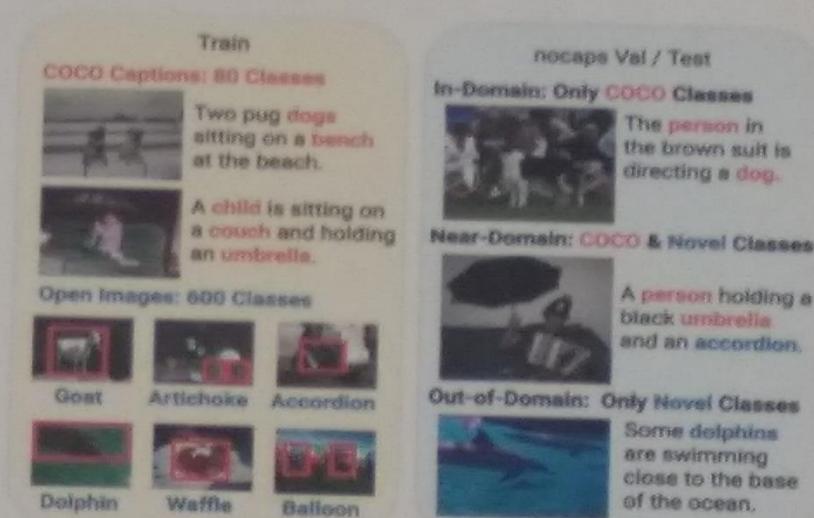
New Challenges – V&L 2.0

- Diversity and Controllable
 - Novel object captioning
 - Captioning with styles
- Reasoning
 - CLEVR dataset
 - GQA – CLEVR in the real world
 - Visual Commonsense Reasoning
- Embodied
 - Embodied VQA, Interactive QA
 - Language-guided visual navigation
 - Remote Embodied Referring Expression



Novel Object Captioning

- In-domain
 - objects have been described in the training data
- Near-domain
 - the most salient objects in the image are novel objects
- Out-of-domain
 - do not contain any COCO classes



The **nocaps** benchmark for novel object captioning (at scale).

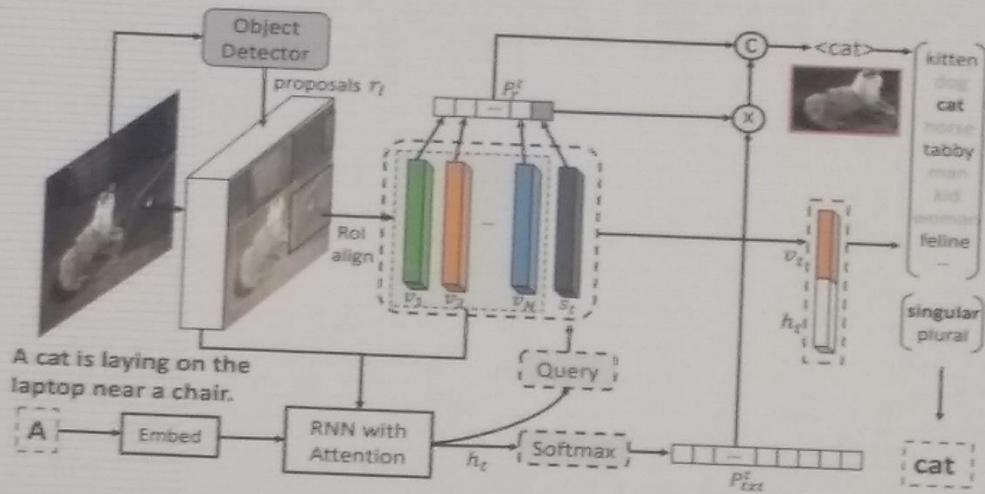
nocaps: novel object captioning at scale
Harsh Agrawal*, Karan Desai*, Xinlei Chen,
Parikh, Stefan Lee, Trevor Anderson, arXiv

uv Batra, Devi

Neural Baby Talk

Template' generation

- at each time step, the model decides whether to generate a word from the textual vocabulary or generate a "visual" word.
- A <region-17> is sitting at a <region-123> with a <region-3>, filling in the slots
- Classifies each of the indicated regions into object categories.
- A puppy is sitting at a table with a cake

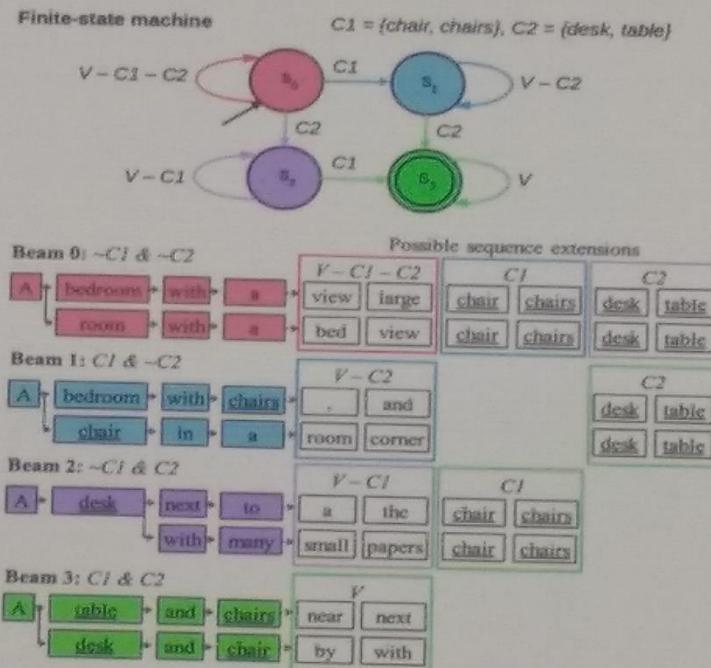
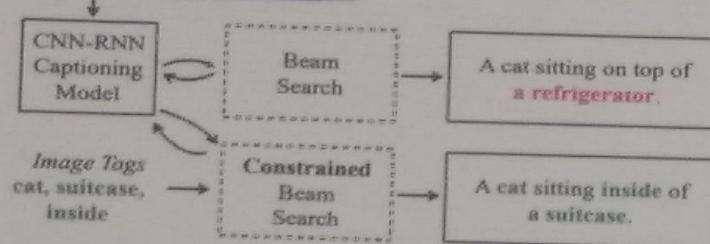


Lu, Jiasen, Jianwei Yang, Dhruv Batra, and Devi Parikh. "Neural baby talk." In CVPR. 2018.

Constrained beam search



Input image containing previously unseen object ('suitcase')



Anderson, Peter, Stephen Gould, and Mark Johnson. "Partially-Supervised Image Captioning." In NIPS, 2018.

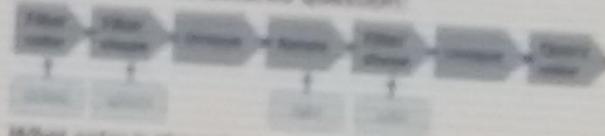
Visual Reasoning - CLEVR

Questions in CLEVR test various aspects of visual reasoning including **attribute identification**, counting, comparison, spatial relationships, and logical operations

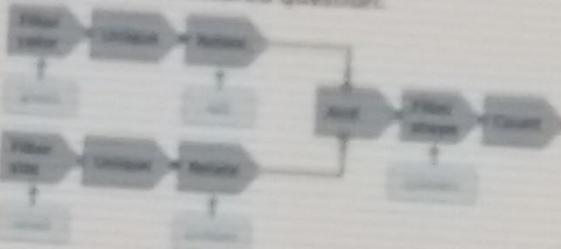


equal number of large things and metal spheres?
the cylinder that is left of the brown metal thing that is left of the
here with the same size as the metal cube, is it made of the
as the small red sphere?
objects are either small cylinders or red things?

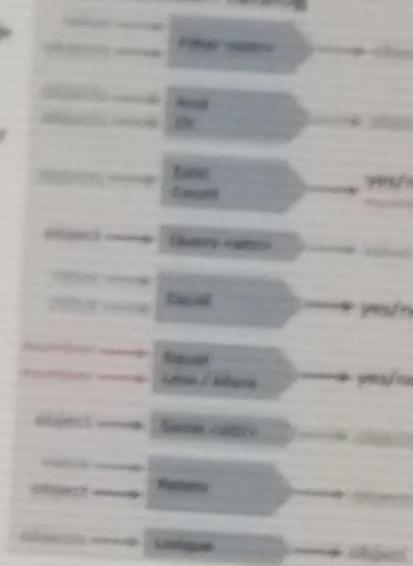
Sample chain-structured question:



Sample tree-structured question:



CLEVR function catalog



Johnson, Justin, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. "Clevr: A diagnostic dataset for compositional language and elementary visual reasoning." In CVPR. 2017.

GQA – Graph QA



Pattern: What/which types do you think the objects, colors or decoys?

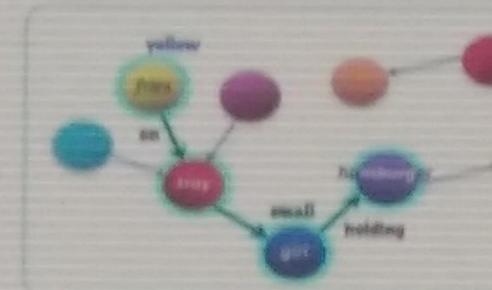
Program: Select :objects → Choose :types; :values; :decoys

Reference: The food on the red object left of the small girl that is holding a hamburger

Decoy: brown

What color is the food on the red object left of the small girl that is holding a hamburger, yellow or brown?

Select: hamburger → Relation: girl, holding → Filter size: small → Relation: object, left → Filter color: red → Relation: food, on → Choose: color: yellow | brown



Graph Normalization

- Ontology construction
- Edge Pruning
- Object Augmentation
- Global Properties

Question Generation

- Patterns Collection
- Compositional References
- Decoys Selection
- Probabilistic Generation

Sampling and Balancing

- Distribution Balancing
- Type Based Sampling
- Deduplication

Entailments Relations

- Functional Programs
- Entailment Relations
- Recursive Reachability

New Metrics

- Consistency
- Validity & Plausibility
- Distribution
- Grounding

Hudson, Drew A., and Christopher D. Manning. "GQA: a new dataset for compositional question answering over real-world images." CVPR 2019

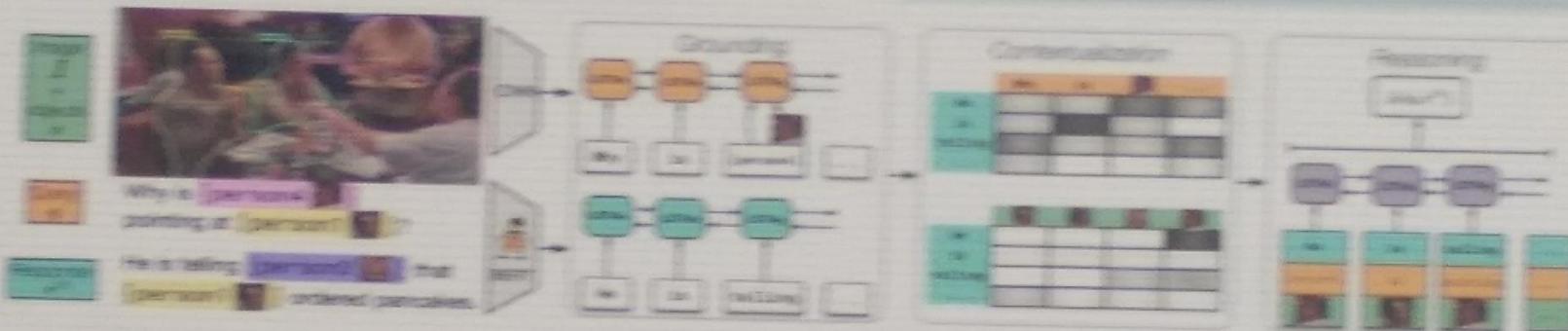
VCR – Visual Commonsense Reasoning



Why is [person4] pointing at [person1]?

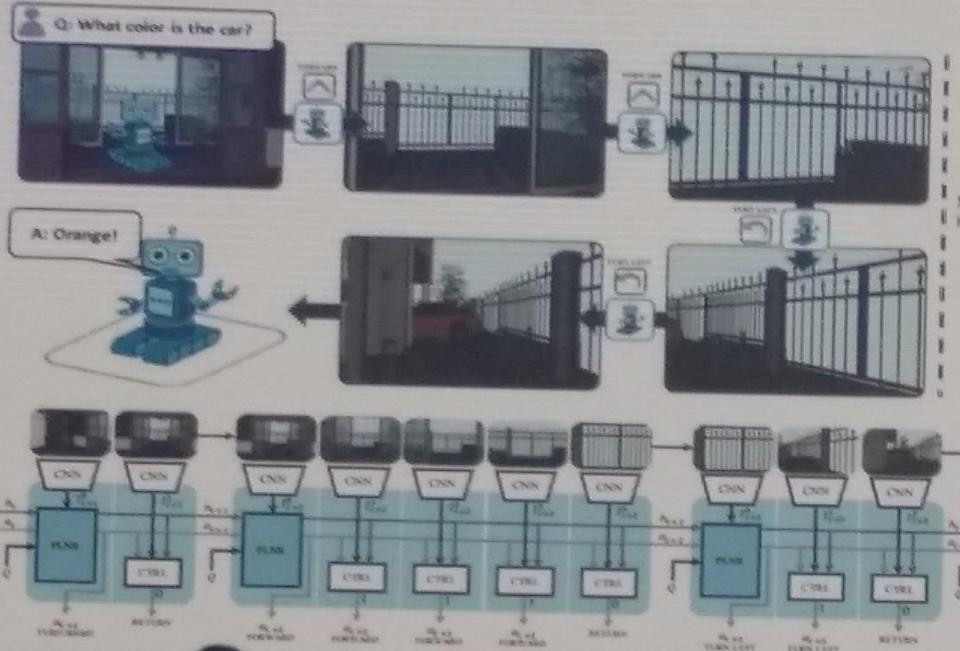
- (A) He is telling [person1] what [person1] should do next.
- (B) He is too nervous.
- (C) He is feeling competitive towards [person1].
- (D) He is giving [person1] directions.

- (E) [person1], [person2], [person3], [person4] had the pancakes in front of them.
- (F) [person1], [person2], [person3], [person4] are taking turns to order and asked for confirmation.
- (G) [person1], [person2], [person3], [person4] are looking at the pancakes and both sides and are smiling slightly.
- (H) [person1], [person2], [person3], [person4] are delivering food to the table and take right one which is closer.



Yers, Rowan, Yannick Biel, Ali Farhadi, and Yejin Choi. "From Recognition to Cognition: Visual Commonsense Reasoning." CVPR 2019.

Embodied VQA & MT-Embodied VQA



EQQA-v1: What color is the car?

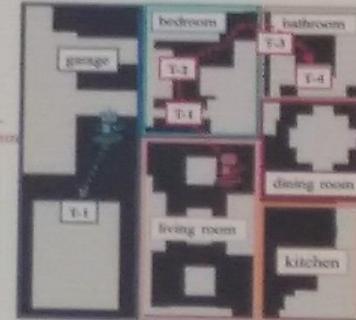


Answer: Orange

MT-EQA: Does the dressing table in the bedroom have same color as the sink in the bathroom?



Answer: No



Das, Abhishek, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra.
Embodied question answering. CVPR 2018.

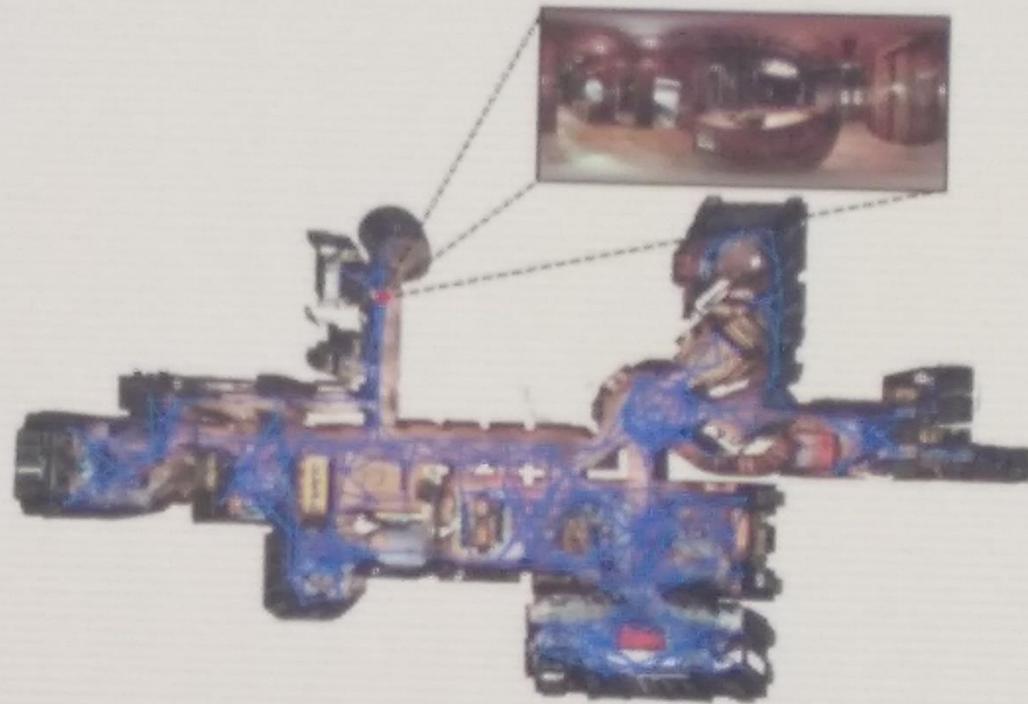
Licheng Yu, Xinlei Chen, Georgia Gkioxari, Mohit Bansal, Tamara L Berg, Dhruv Batra. Multi-Target Embodied Question Answering. CVPR 2019

Vision-and-Language Navigation

Goal: 3.7m



...om. Turn left and exit the room on the left. Wait there.

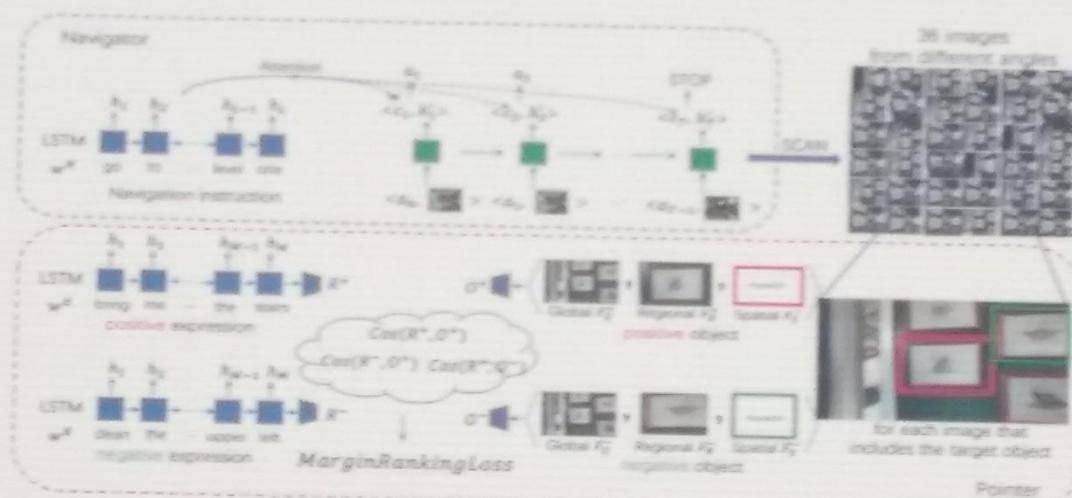


Anderson, Peter, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. "Vision-and-language navigation: interpreting visually-grounded navigation instructions in real environments." In *CVPR 2018*.

Remote Embodied Referring Expression



Instruction: Go to the stairs on level one and bring me the
picture that is next to the top of the stairs.



Conclusions

- Limited improvements on classical image captioning and VQA
- Face to the real challenges!
 - Diversity
 - Controllable
 - Reasoning
 - Deploy the model in the real environment
 - Ask, Answer and Act
- Work on our **embodied visual-navigation & referring expression** tasks
- Two Fully-Funded International PhD positions!



<http://dmirlab.com>

Outline

- About VALSE
- 弱监督视觉理解与主动学习
- 以人为中心的视觉理解
- 视觉中的知识推理
- CV+NLP 2.0
- Poster 选讲



Poster

<http://dmirlab.com>

- IndRNN: Independently Recurrent Neural Networks (CVPR2018)
- FishNet: A Versatile Backbone for Image, Region, and Pixel Level Prediction (NIPS2018)
- HRNet: Deep High-Resolution Representation Learning for Human Pose Estimation (CVPR 2019, MSRA)



<http://dmirlab.com>

Poster

- GLCN: semi-supervised learning with graph learning-convolutional networks (CVPR2019)
- Data Representation and Learning with graph diffusion-embedding networks (CVPR 2019)
- multi-label image recognition with GCN (CVPR2019)
- two-stream adaptive GCN for skeleton-based action recognition (CVPR2019)

Proposed IndRNN

gradient backpropagation through time (BPTT)

- For n-th neuron: $h_{n,t} = \sigma(\mathbf{w}_n x_t + u_n h_{n,t-1} + b_n)$
- $\frac{\partial J_n}{\partial h_{n,t}} = \frac{\partial J_n}{\partial h_{n,T}} \frac{\partial h_{n,T}}{\partial h_{n,t}} = \frac{\partial J_n}{\partial h_{n,T}} \prod_{k=t}^{T-1} \sigma'_{n,k+1} u_n = \frac{\partial J_n}{\partial h_{n,T}} u_n^{T-t} \prod_{k=t}^{T-1} \sigma'_{n,k+1}$

To keep long and short-term memory: $|u_n| \in [0, \sqrt[T-t]{\frac{\gamma}{\prod_{k=t}^{T-1} \sigma'_{n,k+1}}})$

Easily solving the gradient exploding and vanishing problem

Robust with ReLU: $|u_n| \in [0, \sqrt[T-t]{\gamma})$

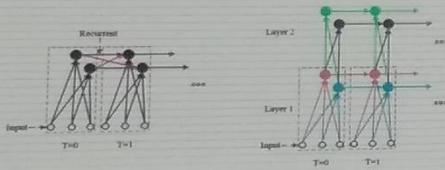
Easier to interpret: Neurons independent, weight for range of memory

New perspective of RNN

- Independently aggregating spatial patterns (through W) over time (through u)

Proposed Independently Recurrent Neural Networks (IndRNN)

- Form
 - $\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{u} \odot \mathbf{h}_{t-1} + \mathbf{b})$
 - Using element-wise product to process recurrent input
- Independently recurrent in one layer
 - For n -th neuron: $h_{n,t} = \sigma(\mathbf{w}_n x_t + u_n h_{n,t-1} + b_n)$
- Cross-channel information explored in the following layers



FishNet

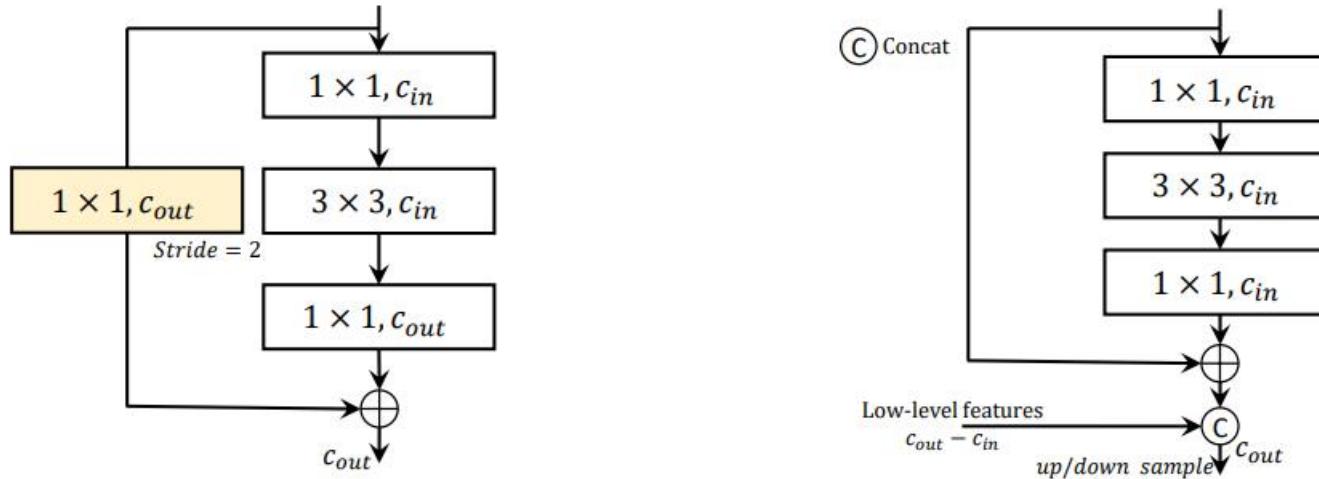


Figure 1: The up/down-sampling block for ResNet (left), and FishNet (right). The 1×1 convolution layer in yellow indicates the *Isolated convolution (I-conv, see Section 2)*, which makes the direct BP incapable and degrades the gradient from the output to shallow layers.

FishNet

<http://dmirlab.com>

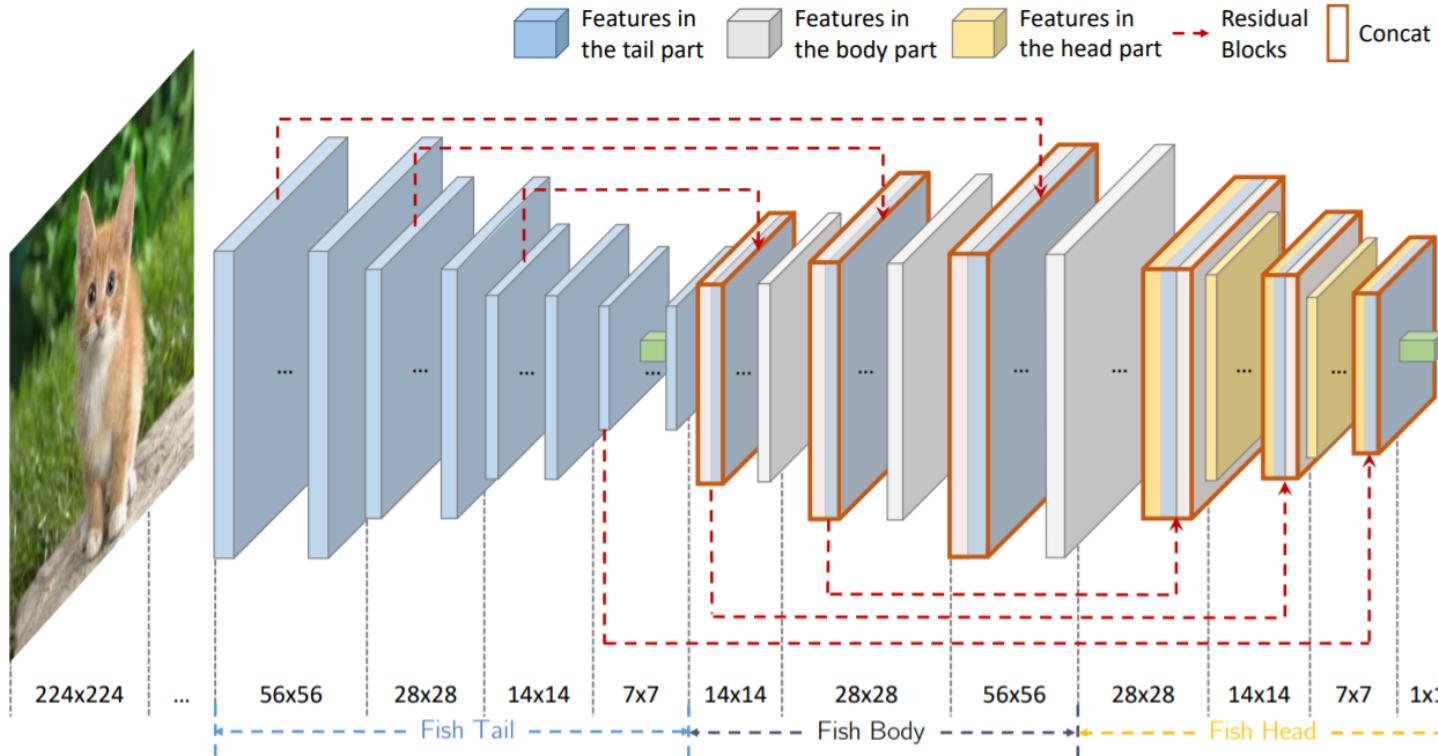
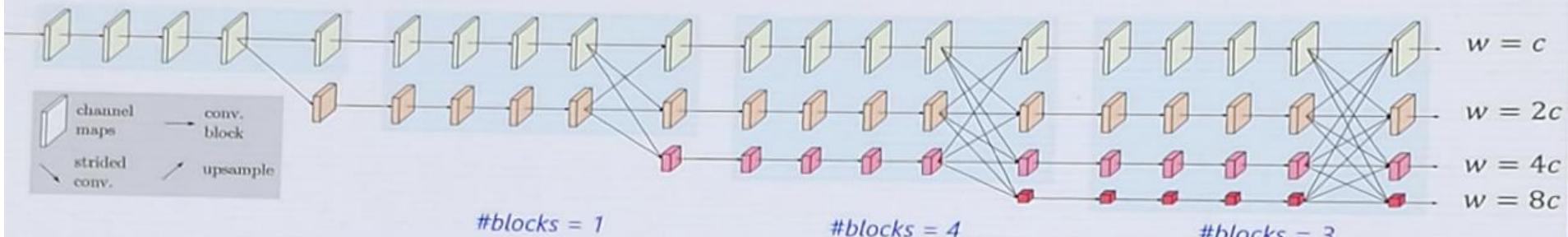


Figure 2: Overview of the FishNet. It has three parts. *Tail* uses existing works to obtain deep low-resolution features from the input image. *Body* obtains high-resolution features of high-level semantic information. *Head* preserves and refines the features from the three parts.

FishNet: A Versatile Backbone for Image, Region, and Pixel Level Prediction (NIPS2018)

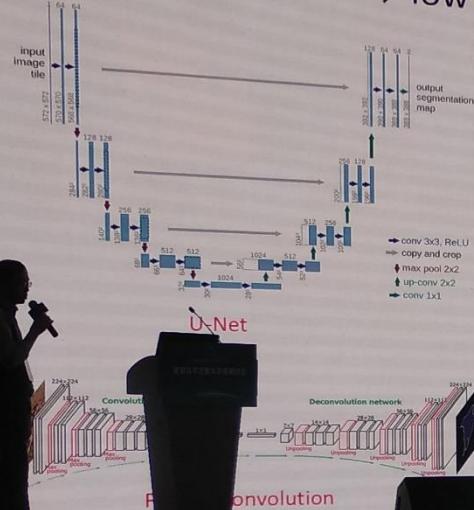
<http://dmirlab.com>

HRNet instantiation

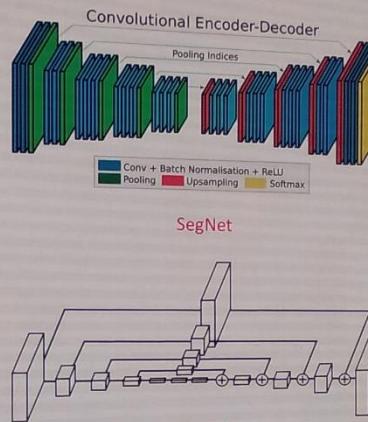


HRNet: Deep High-Resolution Representation Learning for Human Pose Estimation (CVPR 2019, MSRA)

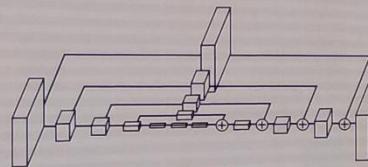
High resolution → low resolution → high resolution



Residual convolution



SegNet



Hourglass

6



Motivation

Traditional Graph CCNS

- Two-stage framework;
- Use a *fixed* graph.

The disadvantage of fixing graph

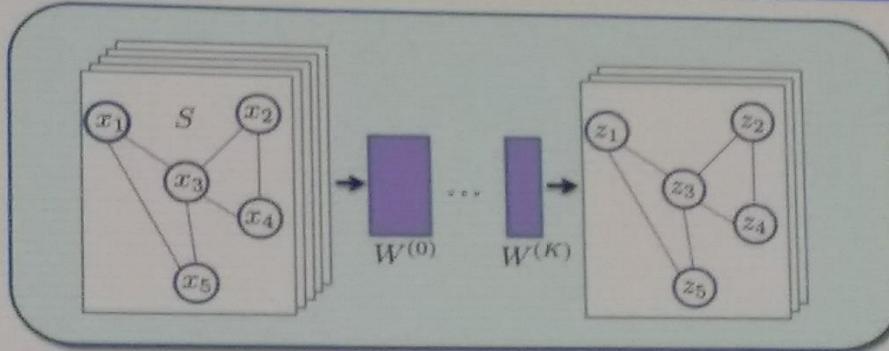
- The fixed graph may contain *noise*;
- The fixed graph is obtained from a separate network which is also not *guaranteed* to best serve the graph CNNs. That ignores the correlation between *graph construction* and *graph convolutional*, which may lead to weak suboptimal solution.



The Overall Architecture

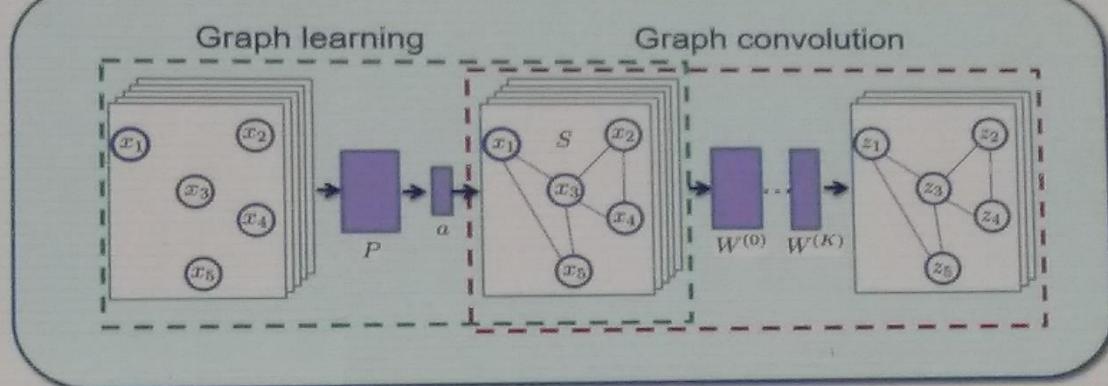
GCN

Graph convolution



GLCN

Graph learning +
Graph convolution

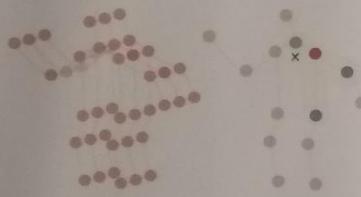


Introduction

- In existing GCN-based methods for skeleton-based action recognition task, the topology of the graph is set manually, and it is fixed over all layers and input samples. This may not be optimal for the hierarchical GCN and diverse samples in action recognition tasks. In this work, an adaptive graph convolutional network is proposed to adaptively learn the topology of the graph for different GCN layers and skeleton samples in an end-to-end manner, which can better suit the action recognition task and the hierarchical structure of the GCNs.
- The second-order information (the lengths and directions of bones) of the skeleton data, which is naturally more informative and discriminative for action recognition, is rarely investigated in existing methods. In this work, the second-order information of the skeleton data is explicitly formulated and combined with the first-order information using a two-stream framework, which brings notable improvement for the recognition performance.
- On two large-scale datasets for skeleton-based action recognition, the proposed 2s-AGCN exceeds the state-of-the-art by a significant margin. Code: <https://github.com/lshiwjx/2s-AGCN>.

Proposed Methods

Graph Convolutional Network (GCN)



$$f_{out}(v_i) = \sum_{v_j \in \mathcal{B}_i} \frac{1}{Z_{ij}} f_{in}(v_j) \cdot w(l_i(v_j))$$

$$f_{out} = \sum_k^K \mathbf{W}_k (\mathbf{f}_{in} \mathbf{A}_k) \odot \mathbf{M}_k$$

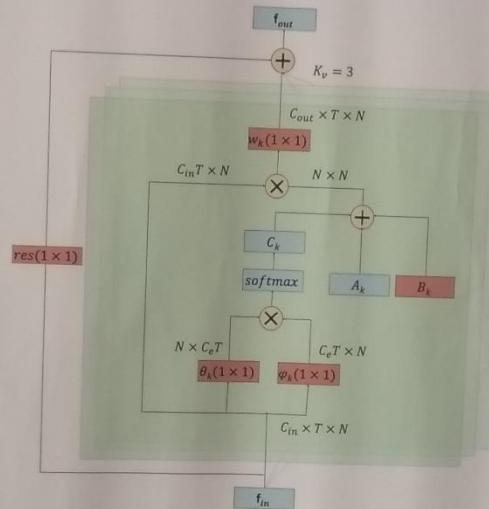
Two-Stream Adaptive Graph Convolutional Networks for Skeleton-Based Action Recognition

Lei Shi, Yifan Zhang, Jian Cheng, Hanqing Lu
Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China



Proposed Methods

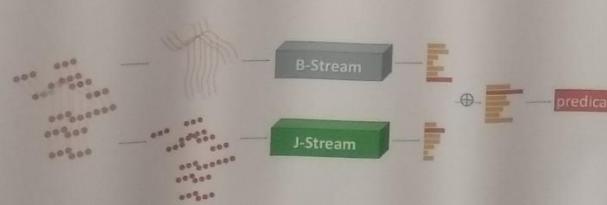
Adaptive Graph Convolutional Network (AGCN)



$$\mathbf{f}_{out} = \sum_k^K \mathbf{W}_k \mathbf{f}_{in} (\mathbf{A}_k + \mathbf{B}_k + \mathbf{C}_k)$$

$$\mathbf{C}_k = \text{softmax}(\mathbf{f}_{in}^T \mathbf{W}_{\theta k}^T \mathbf{W}_{\phi k} \mathbf{f}_{in})$$

Two-Stream Network



Experiments

Preprocessing methods

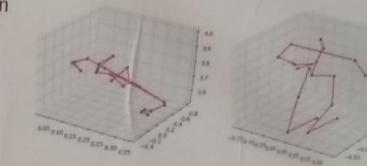


Figure 2. Example of the data preprocessing on the NTU-dataset. The left is the original skeleton, and the right is the processed skeleton.

Results on NTU-RGB-D

Methods	Accuracy (%)	Methods	Accuracy (%)
ST-GCN	92.7	Js-AGCN	93.7
ST-GCN wo/M	91.1	Bs-AGCN	93.2
AGCN wo/A	93.4	2s-AGCN	95.1
AGCN wo/B	93.3		
AGCN wo/C	93.4		
AGCN	93.7		

Visualization

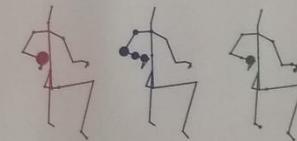


Figure 8. Visualization of the graphs for different layers.

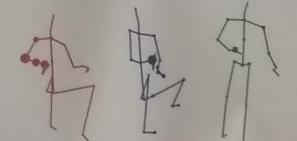


Figure 9. Visualization of the graphs for different samples.



<http://dmirlab.com>

Thank you !
Q&A