

Open Elective Course [OE]

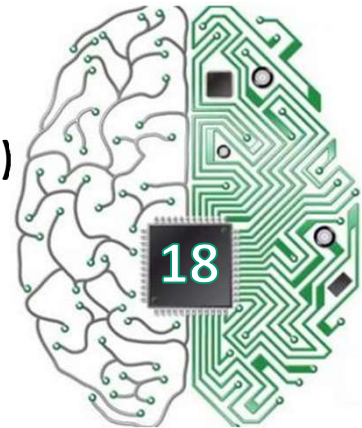
Course Code: CSO507

Winter 2023-24

Lecture#

Deep Learning

Unit-4: Convolutional Neural Networks (Part-VI)

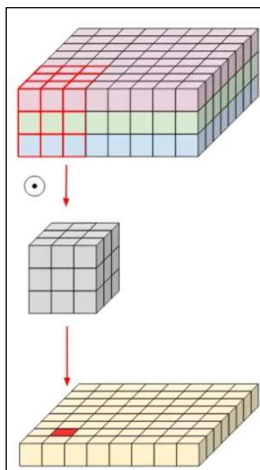
**Course Instructor:****Dr. Monidipa Das**

Assistant Professor

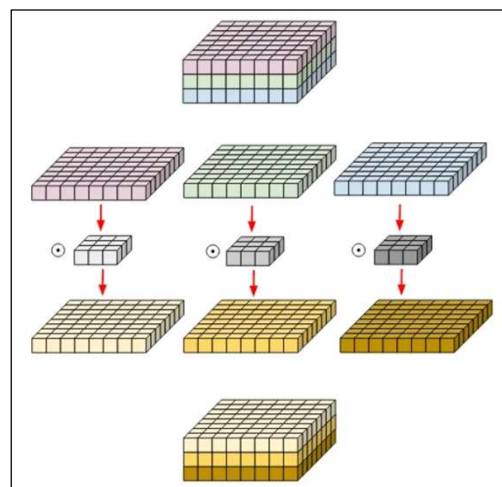
Department of Computer Science and Engineering

Indian Institute of Technology (Indian School of Mines) Dhanbad, Jharkhand 826004, India

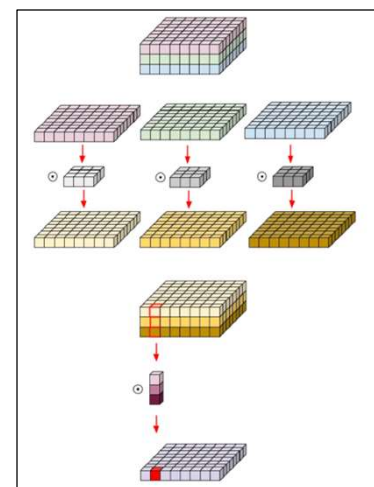
Illustrations for Standard, Depth-wise, and Depth-wise Separable Convolution



Standard convolution



Depth-wise convolution



Depth-wise separable convolution

Source: Internet

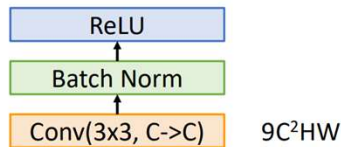
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

MobileNets: Tiny Networks (For Mobile Devices)



Standard Convolution Block

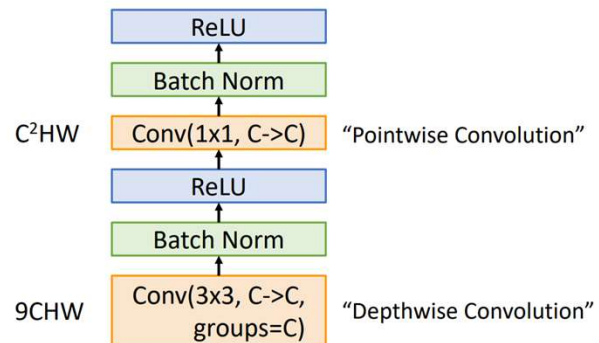
Total cost: $9C^2HW$



$$\begin{aligned}\text{Speedup} &= 9C^2 / (9C + C^2) \\ &= 9C / (9 + C) \\ &\Rightarrow 9 \text{ (as } C \rightarrow \infty)\end{aligned}$$

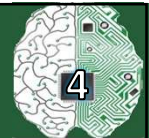
Depthwise Separable Convolution

Total cost: $(9C + C^2)HW$



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

ConvNets Application

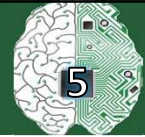


- ConvNets are today ubiquitous in **computer vision**!

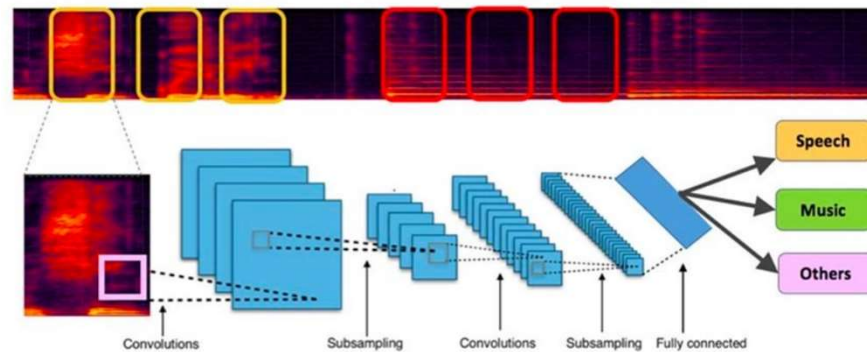


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

ConvNets beyond vision



- CNNs are not only useful for image tasks!
- They are becoming the standard in audio tasks and very competitive in text processing tasks (e.g., sentiment classification).

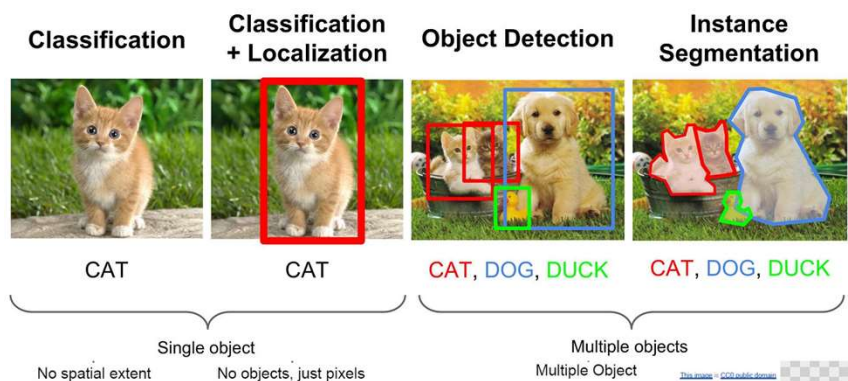


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

ConvNets Application for Computer Vision

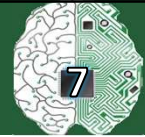


- Four major computer vision problems:
 - Image Classification
 - Image Classification with Localization
 - Object Detection
 - Object Segmentation

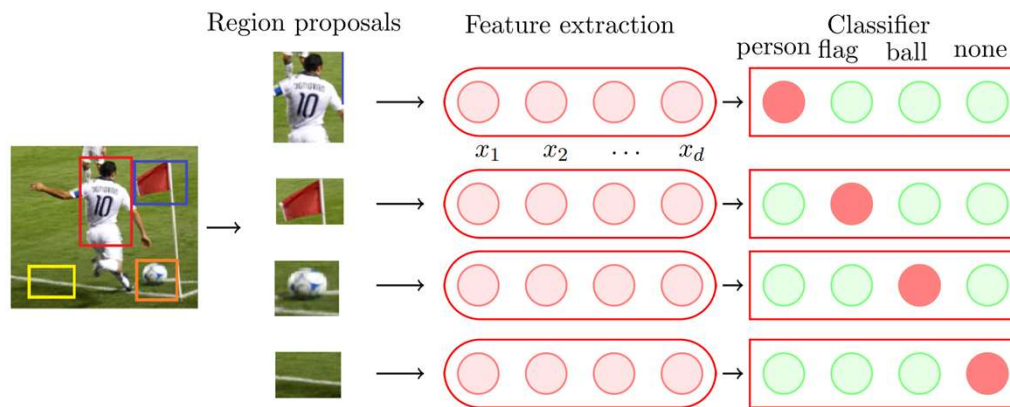


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pipeline for Object Detection



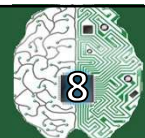
- Starts with a region proposal stage where we identify potential regions which may contain objects



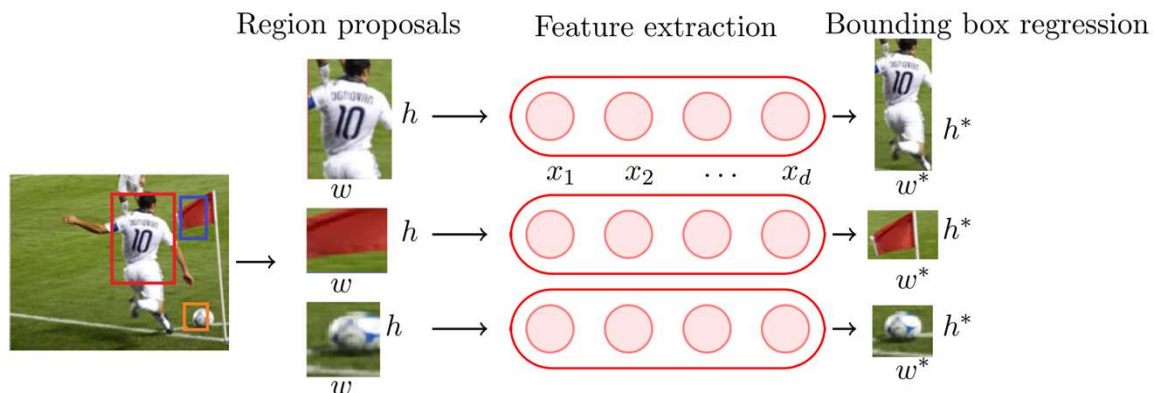
Acknowledgement: Prof. Mitesh M. Khapra

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pipeline for Object Detection

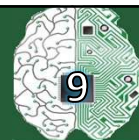


- In addition we would also like to correct the proposed bounding boxes This is posed as a regression problem (for example, we would like to predict w^* , h^* from the proposed w and h)



Acknowledgement: Prof. Mitesh M. Khapra

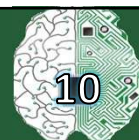
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



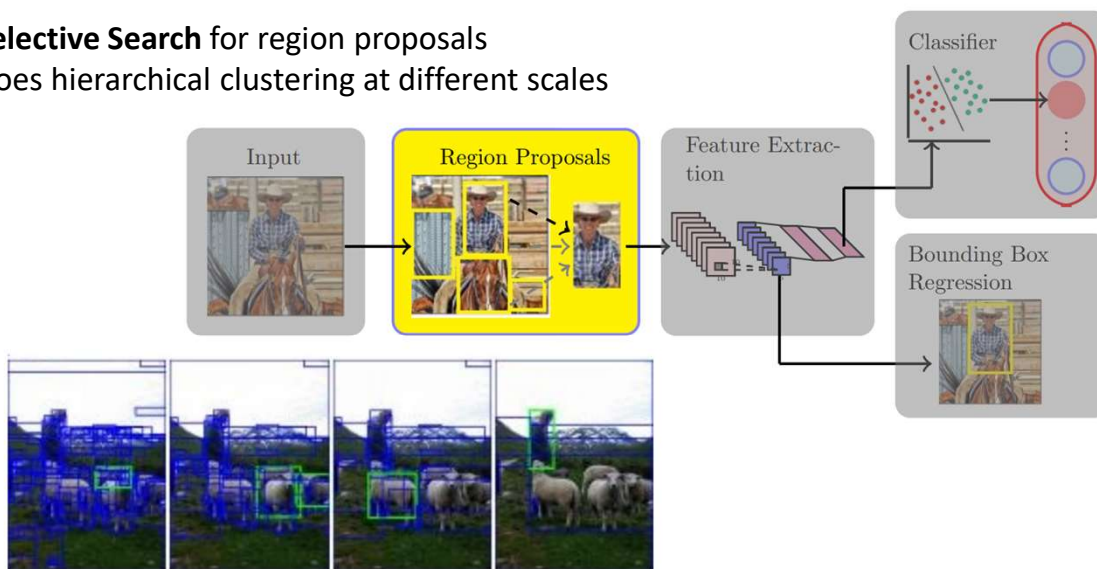
Region-based CNNs (RCNNs)

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection

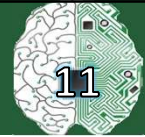


- **Selective Search** for region proposals
- Does hierarchical clustering at different scales

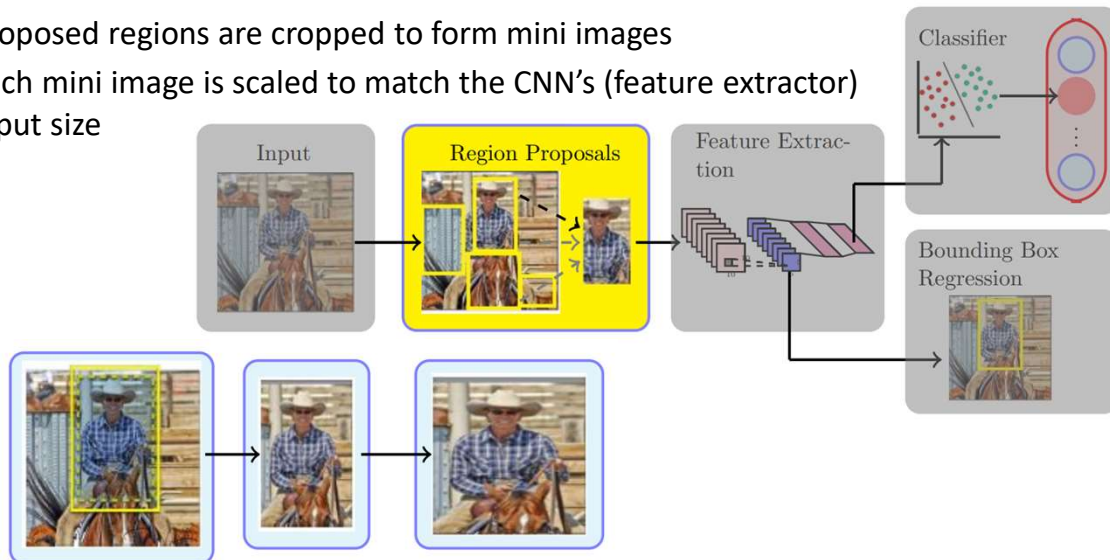


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection



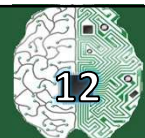
- Proposed regions are cropped to form mini images
- Each mini image is scaled to match the CNN's (feature extractor) input size



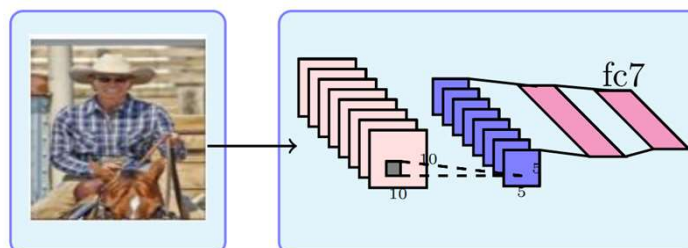
Acknowledgement: Prof. Srinivas Aravamudan

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection



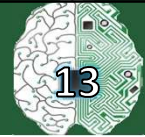
- For **feature extraction** any CNN trained for Image Classification can be used (AlexNet/ VGGNet etc.)
- Outputs from fc7 layer are taken as features
- CNN is fine tuned using ground truth (cropped) object images



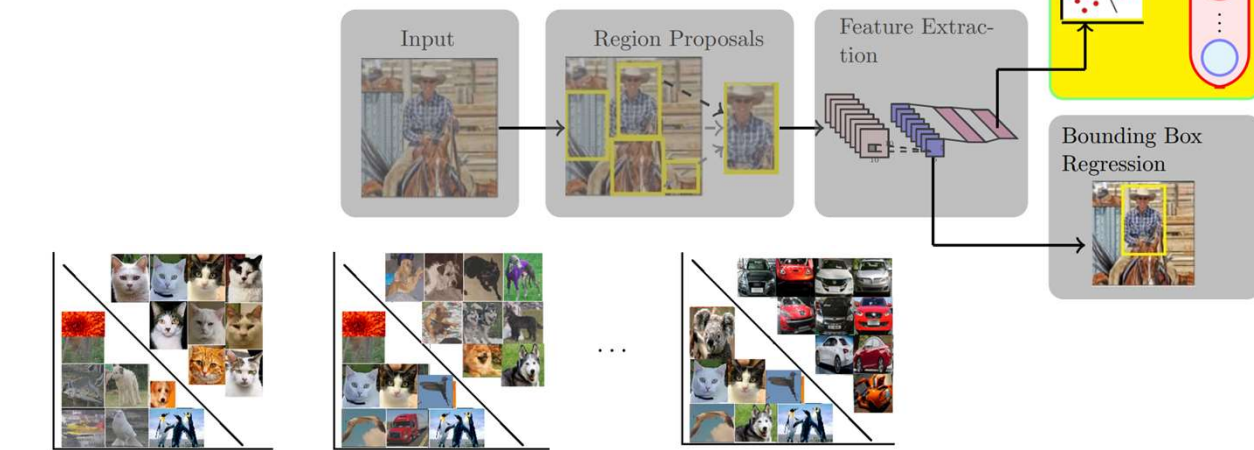
Acknowledgement: Prof. Srinivas Aravamudan

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection



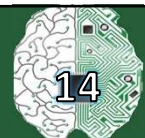
- Linear models (SVMs) are used for classification (1 model per class)



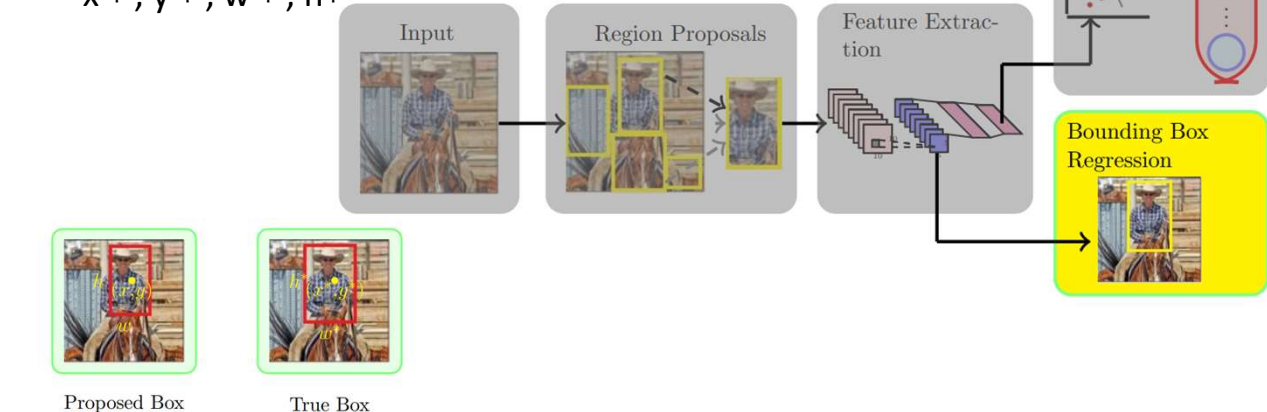
Adapted from [1] and [2]

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection



- The proposed regions may not be perfect
- Learn four regression models which will learn to predict x^* , y^* , w^* , h^*



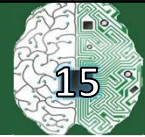
Proposed Box

True Box

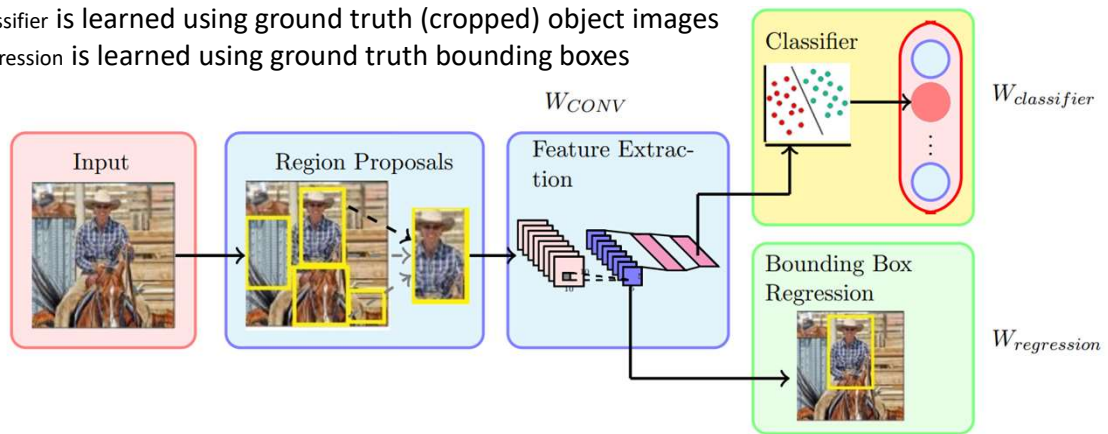
z : features from pool5 layer of the network

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection



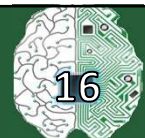
- What are the parameters of this model?
- W_{CONV} is taken as it is from a CNN trained for Image classification (say on ImageNet)
- W_{CONV} is then fine tuned using ground truth (cropped) object images
- $W_{classifier}$ is learned using ground truth (cropped) object images
- $W_{regression}$ is learned using ground truth bounding boxes



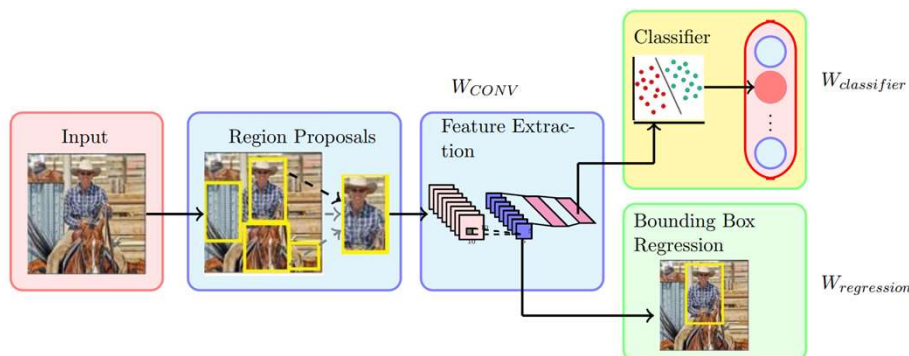
Acknowledgement: Prof. Mitesh M. Khapra

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection

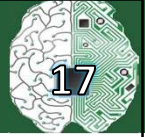


- **What is the computational cost for processing one image at test time?**
 - **Inference Time** = Proposal Time + # Proposals \times Convolution Time + # Proposals \times classification + # Proposals \times regression



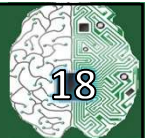
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

RCNN model for object detection



- No joint learning
- Use ad hoc training objectives
 - Fine tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hinge loss)
 - Train post-hoc bounding-box regressors (squared loss)
- Training (≈ 3 days) and testing (47s per image) is slow
- Takes a lot of disk space

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



Questions?

Acknowledgement: Prof. Alex C. Berg, University of Washington, Microsoft Research

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad