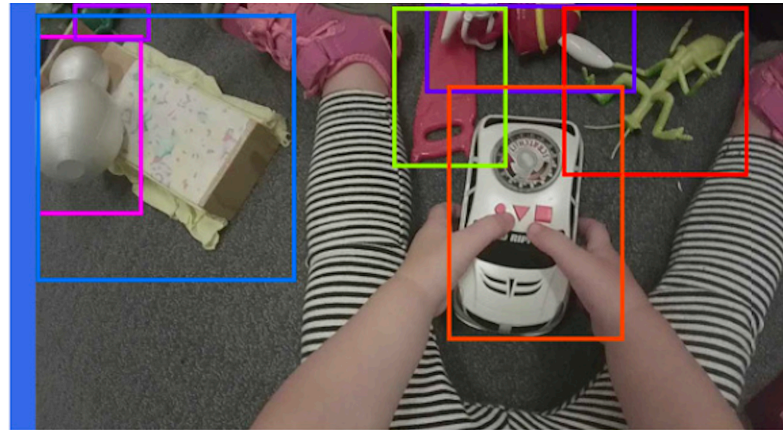# An Egocentric Perspective on Active Vision and Visual Object Learning in Toddlers

S. Bambach, D. Crandall, L. Smith, C. Yu.
ICDL 2017
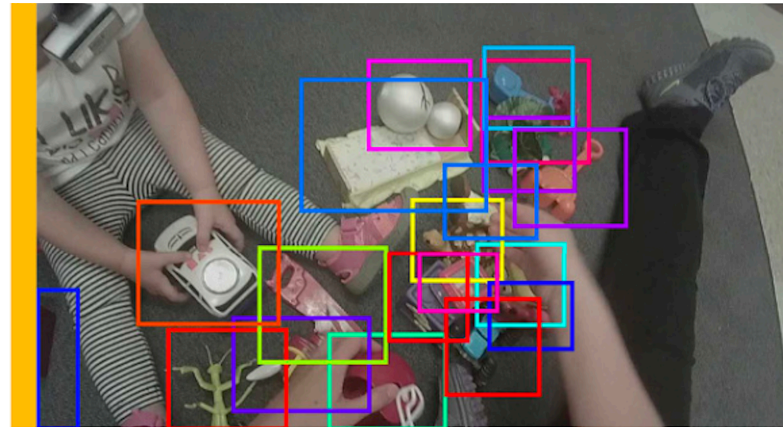
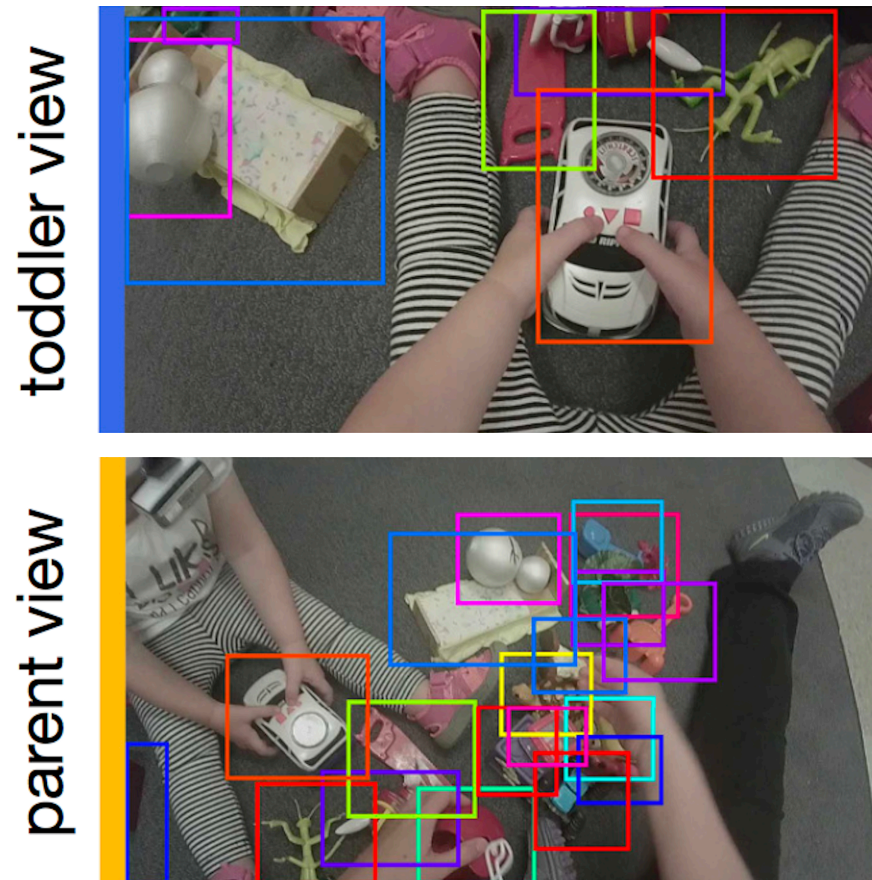Experiment presenters: Arjun, Ginevra

# Their Experiments

# Their Experiments



toddler view

parent view

Authors could not control training set

# Our Experiments

- We generate images where
  - Labeled object occupies fixed percentage of view
  - Background objects do not move



| Scale 0 10% of view | Scale 1 20% of view | Scale 2 30% of view | Scale 3 60% of view | Clean Image |

Image source: collages we made from Caltech 256 database

# Our Experiments

- Simulate toddler bringing object to face
  - We control scale to measure its effect on testing accuracy



| Scale 0 10% of view | Scale 1 20% of view | Scale 2 30% of view | Scale 3 60% of view | Clean Image |

Image source: collages we made from Caltech 256 database

# Our Dataset

- 5 classes, 3633 images

- Collages
  - Construct 'scenes of toys' using Caltech-256
  - 1 positive image amongst many negatives
  - Simulate toddler perspective

Image source: Caltech 256 database

# Scene Generation

- ## Scene dim: **224 x 224**
  - Scale largest image dim to 70
  - Rotate randomly from -15° to 15°

- ## **10** negatives
  - Select uniformly from Caltech-256 negatives
  - Placed randomly in within scene boundary

- ## **1** positive
  - Scale 0 (1x), 1 (1.5x), 2 (2x), 3 (3x)
  - Place randomly within scene boundary (at scale 1)
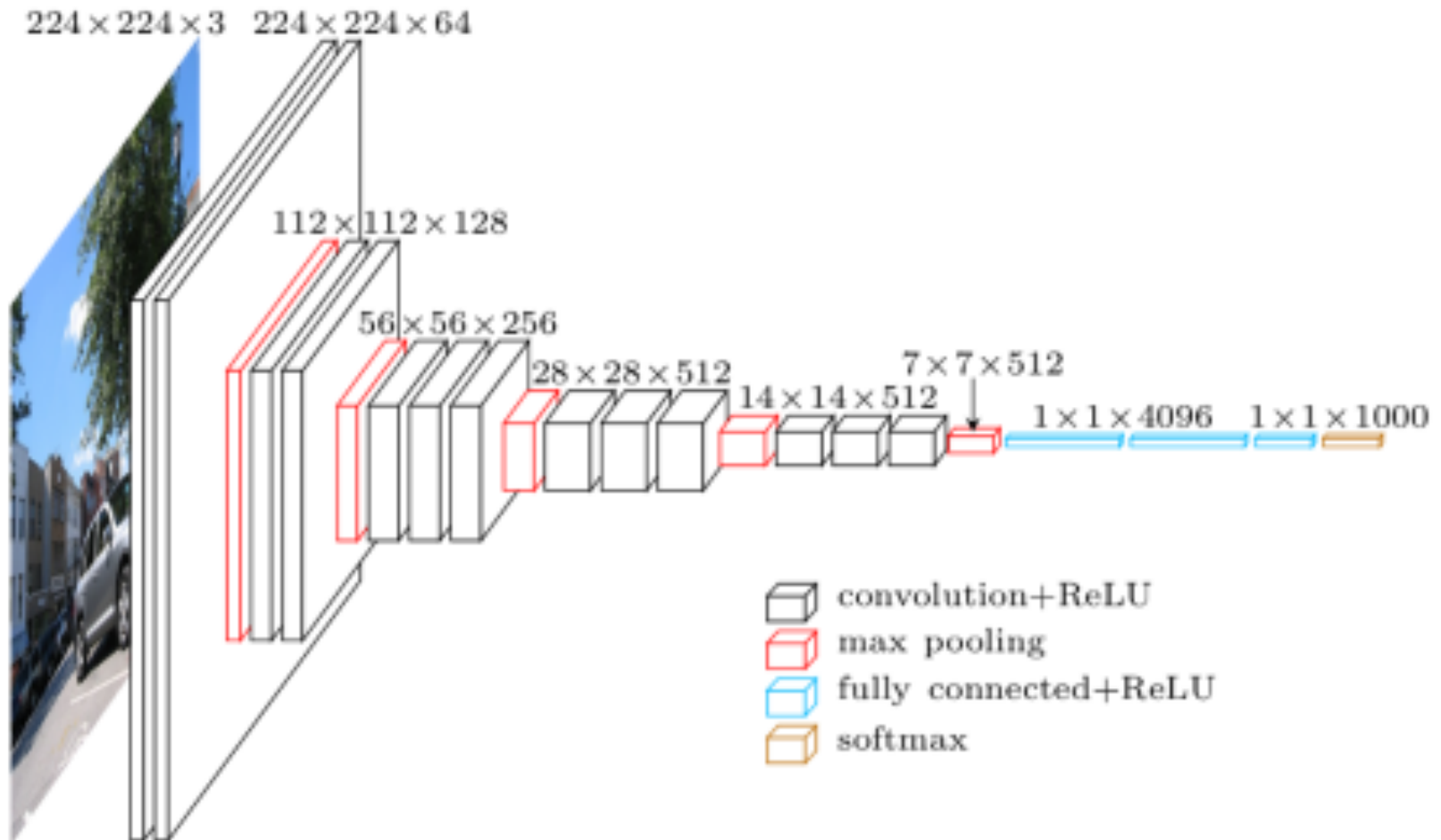
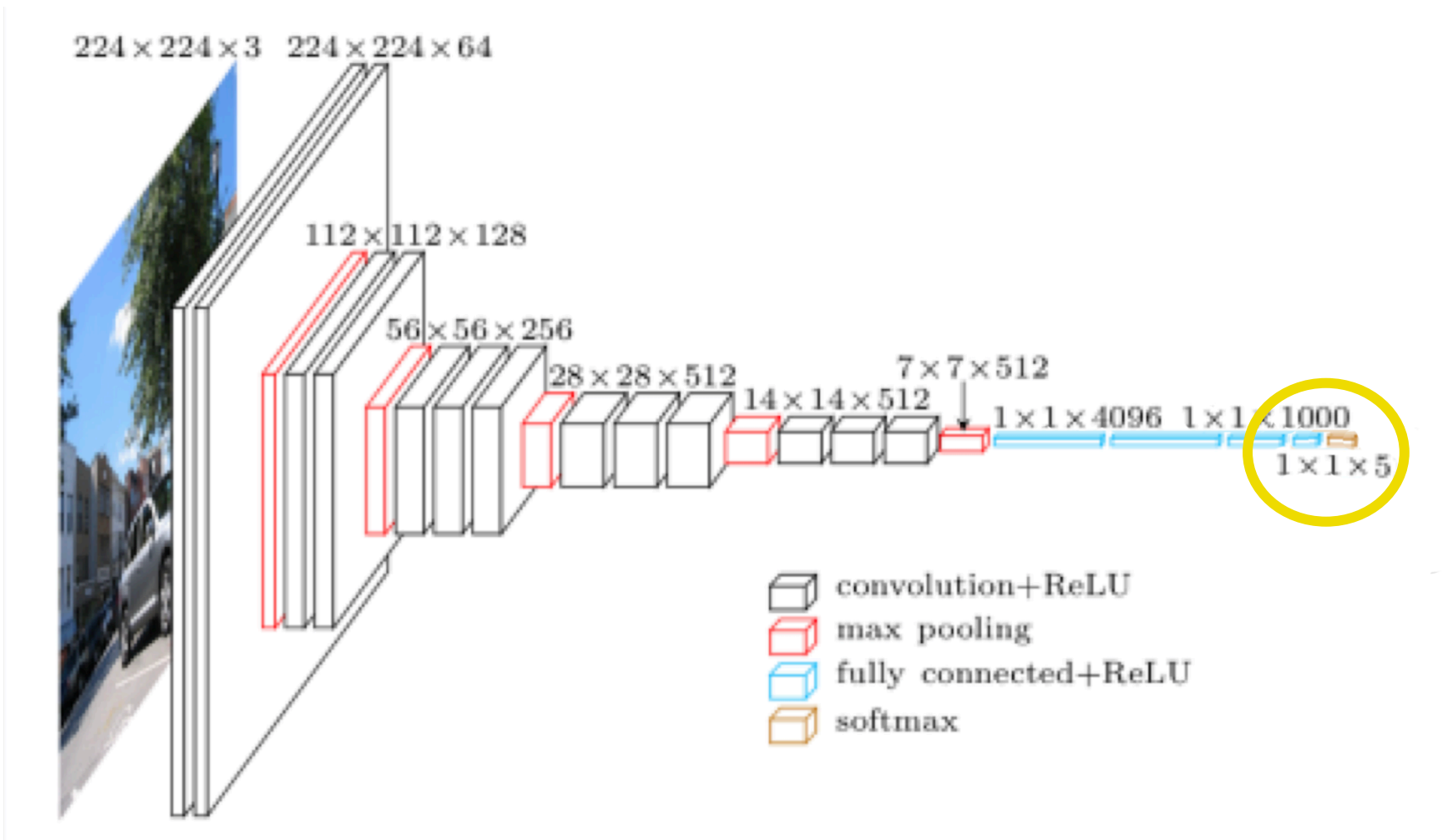- ## **2** scenes per training instance

# VGG 16



224 × 224 × 3  224 × 224 × 64

112 × 112 × 128

56 × 56 × 256

28 × 28 × 512

14 × 14 × 512

7 × 7 × 512

1 × 1 × 4096  1 × 1 × 1000

convolution+ReLU
max pooling
fully connected+ReLU
softmax

# VGG 16 for 5 classes



$224 \times 224 \times 3$  $224 \times 224 \times 64$

$112 \times 112 \times 128$

$56 \times 56 \times 256$

$28 \times 28 \times 512$

$14 \times 14 \times 512$

$7 \times 7 \times 512$

$1 \times 1 \times 4096$  $1 \times 1 \times 1000$

$1 \times 1 \times 5$

convolution+ReLU
max pooling
fully connected+ReLU
softmax

Image source: https://www.cs.toronto.edu/~frossard/post/vgg16/, modified by us

# Experiment Setup

- Experiment 1
  - Train on different scales, test on clean image

- Experiment 2
  - Train on different scales and clean, test on different scales



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view

Clean
Image

Image source: collages we made from Caltech 256 database

# Experiment Setup

- ## Experiment 1
  - Train on different scales, test on clean image
- ## Experiment 2
  - Train on different scales and clean, test on different scales



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view

Clean
Image

Image source: collages we made from Caltech 256 database

# Experiment 1 - objective

- Test effect of 'bringing object to face' for isolated classification

- Questions to consider
  - Effect of viewing at multiple scales?
  - Single ideal scale or result of multiple scales?

# Experiment 1 - data



Scale 0
10% of view
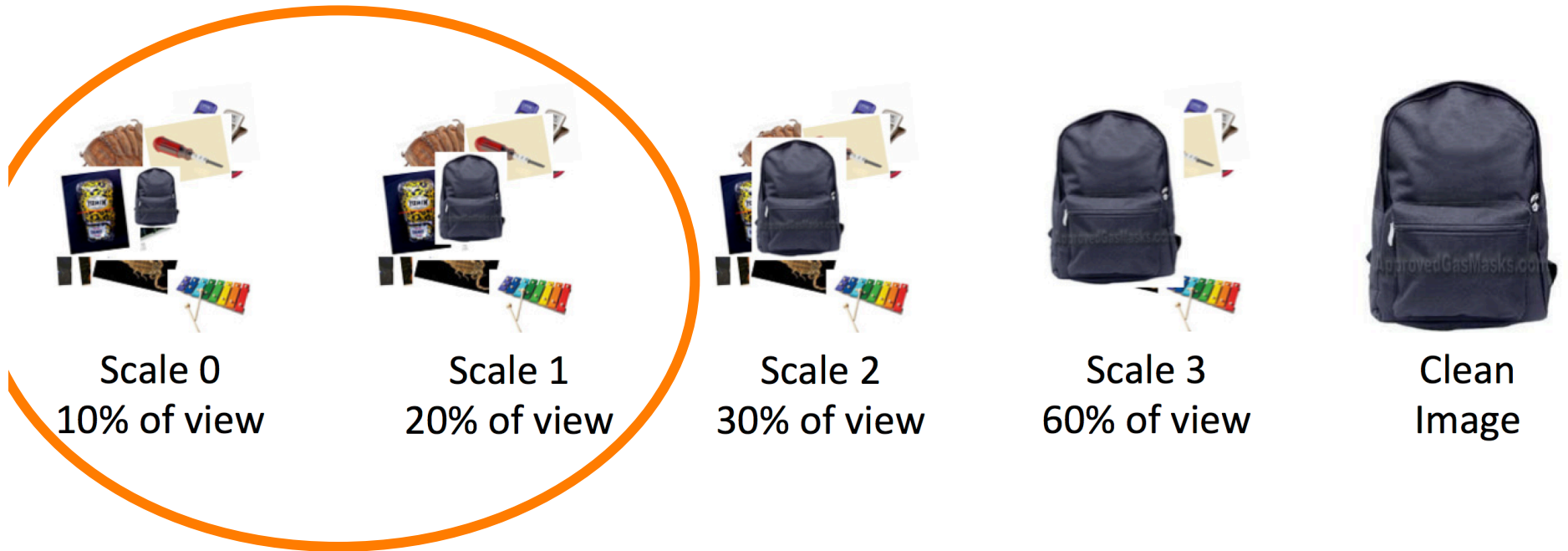
Scale 1
20% of view

Scale 2
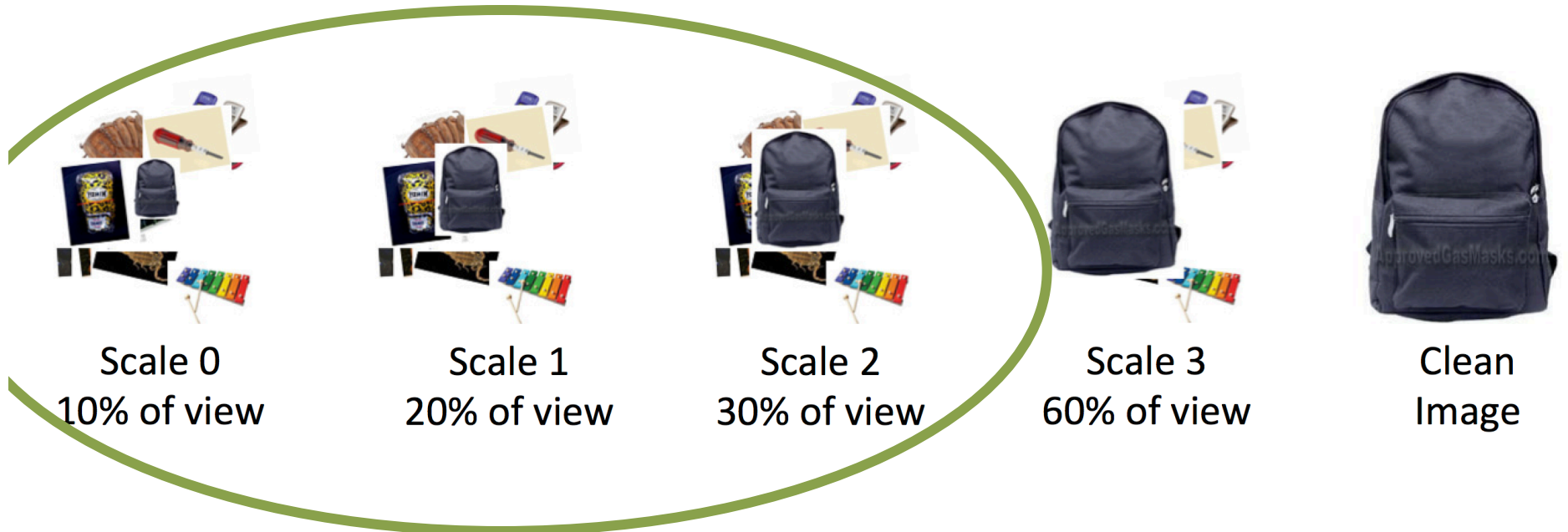30% of view

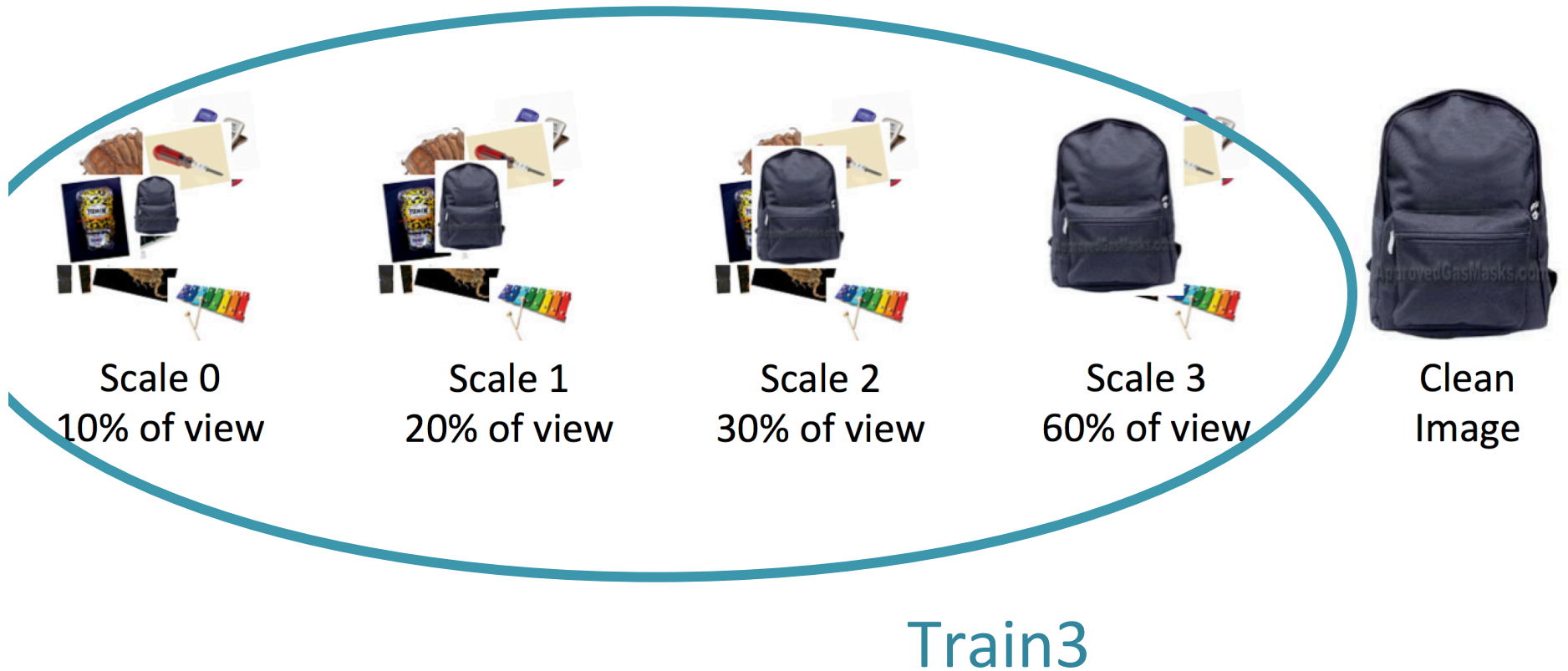Scale 3
60% of view

Clean
Image

Train0

# Experiment 1 - data



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view

Clean
Image

Train1

Image source: collages we made from Caltech 256 database

# Experiment 1 - data



| Scale 0<br>10% of view | Scale 1<br>20% of view | Scale 2<br>30% of view | Scale 3<br>60% of view | Clean<br>Image |

**Train2**

# Experiment 1 - data



| Scale 0 10% of view | Scale 1 20% of view | Scale 2 30% of view | Scale 3 60% of view | Clean Image |

Train3

# Experiment 1 - data



| Scale 0 | Scale 1 | Scale 2 | Scale 3 | Clean |
| 10% of view | 20% of view | 30% of view | 60% of view | Image |

Train3only

# Experiment 1 - data



| Scale 0 | Scale 1 | Scale 2 | Scale 3 | Clean |
|---------|---------|---------|---------|-------|
| 10% of view | 20% of view | 30% of view | 60% of view | Image |

## Correct number of epochs to compensate for more training examples

# Experiment 1 - data



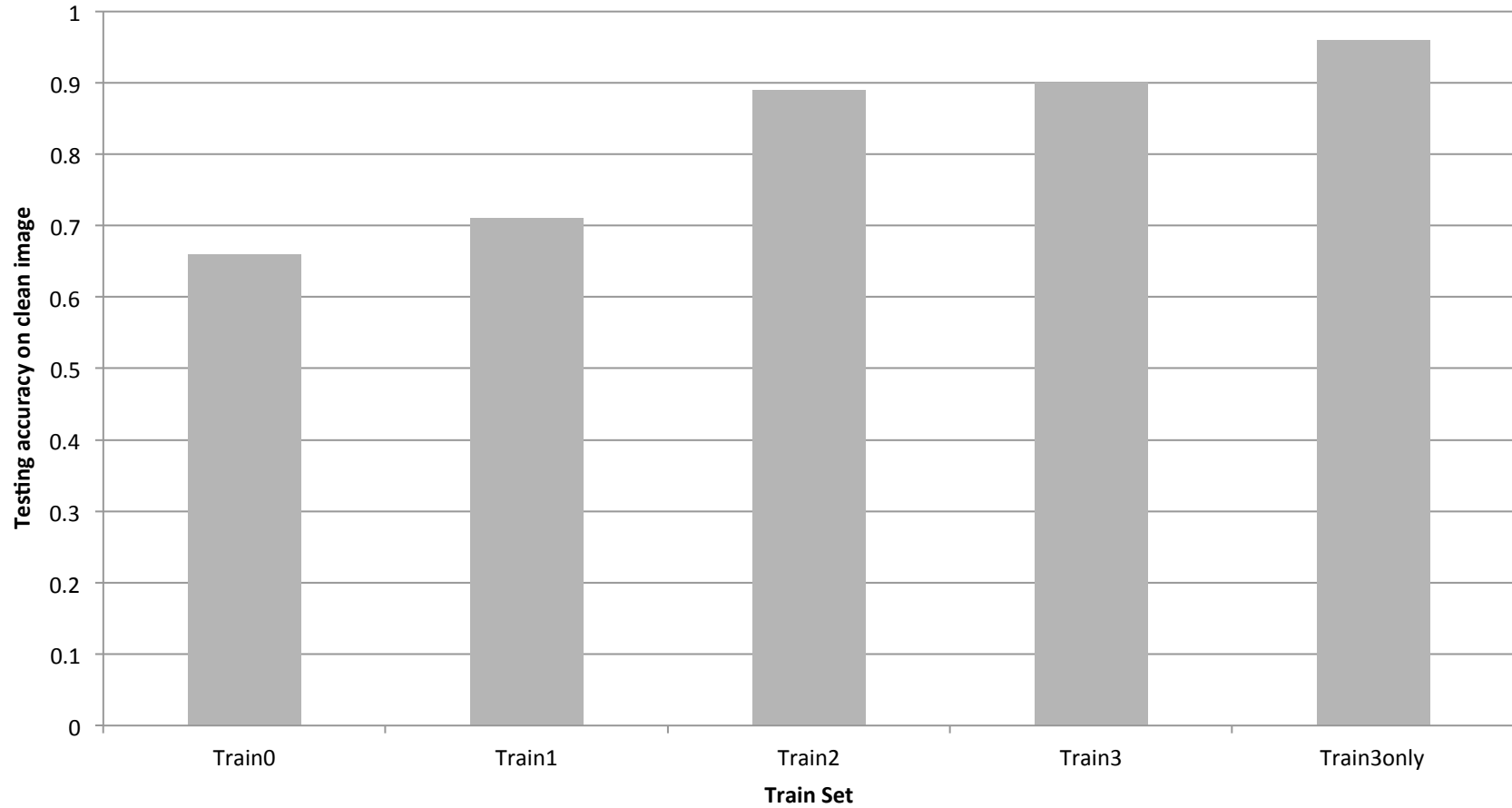Scale 0
10% of view

Scale 1
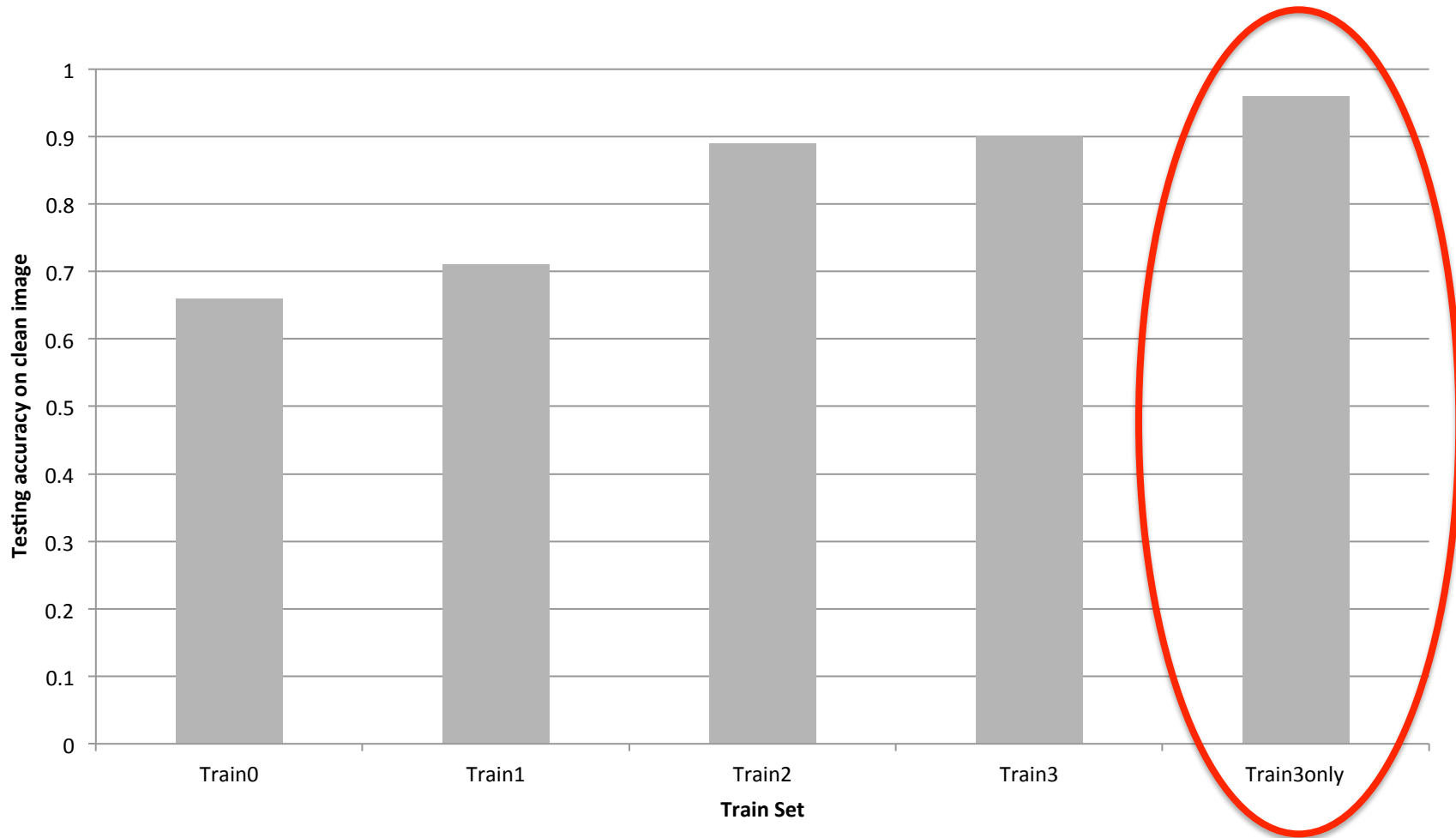20% of view

Scale 2
30% of view
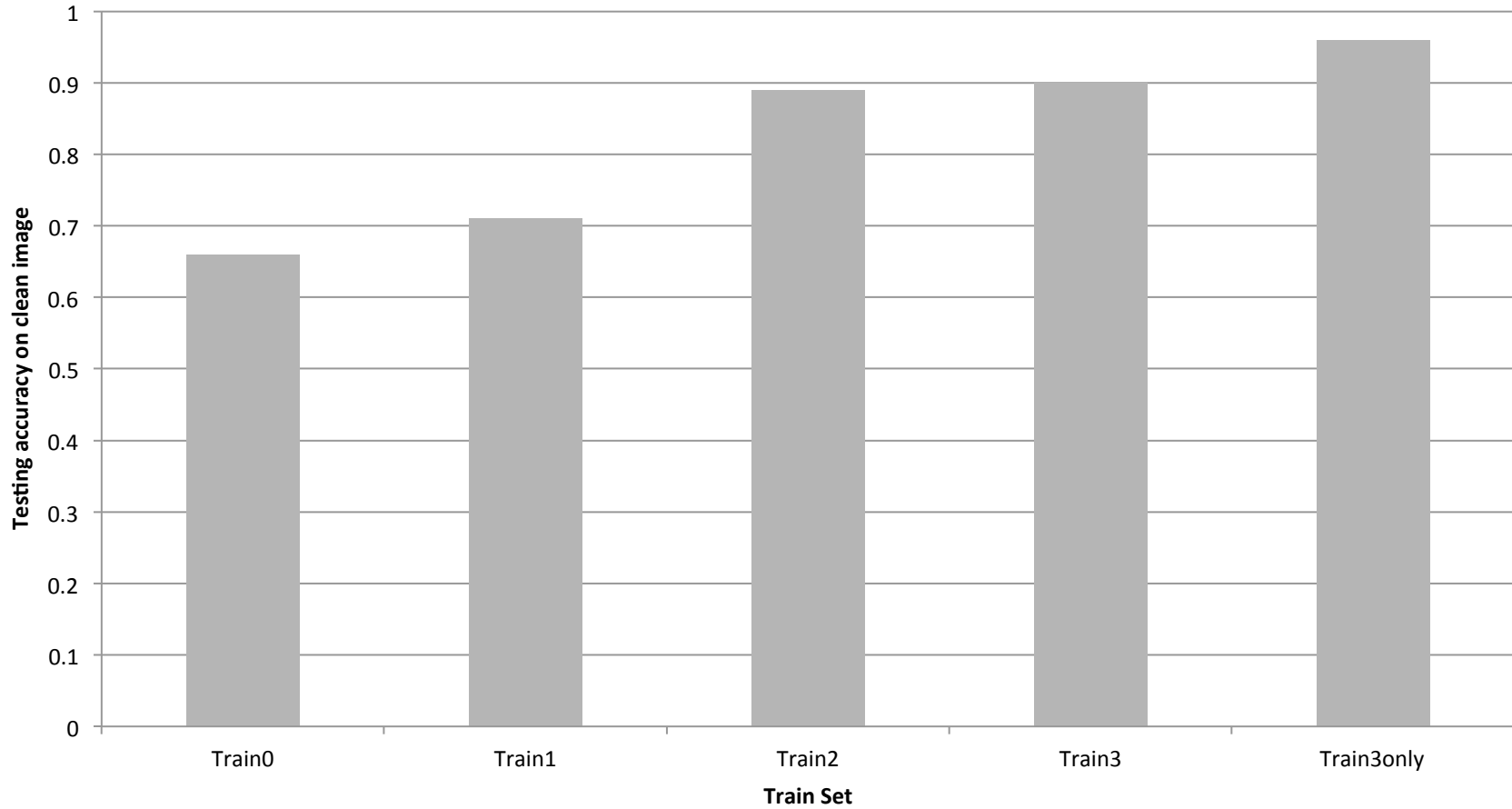
Scale 3
60% of view

Clean
Image

Test

Image source: collages we made from Caltech 256 database
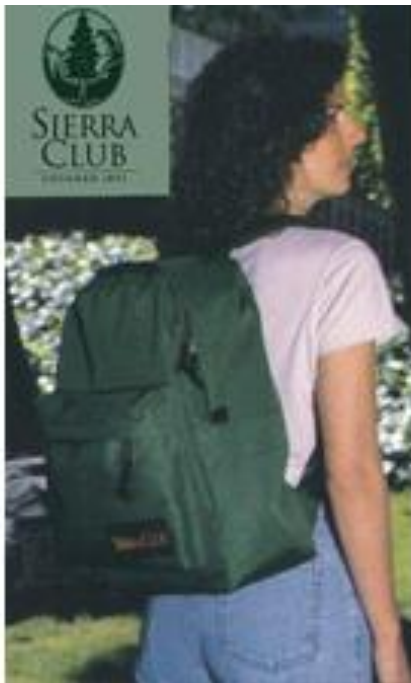
# Experiment 1 - results

# Experiment 1 - results



Training on larger scale images only yields to best test accuracy.

# Experiment 1 - results

- Images misclassified when network trained in low scales benefit from training in higher scales



Misclassified after train0, train1, train2

Correctly classified after train3 and train3only

(Category: bag)

Image source: Caltech 256 database

# Experiment 1 - results

- Images misclassified when network trained in low scales benefit from training in higher scales



Misclassified after train0, train1, train2, train3

Correctly classified only after train3only

(Category: plane)

# Experiment 1 - results

- Images misclassified after train3only were misclassified after all other trainings

Bag

Plane

Plane

# Experiment 1 - conclusions

- Toddler's data gives better training because object is closer, not because it is 'brought to face'

- Significant jump in accuracy if object occupies >30% of view in training

- Training images where object occupies <30% of view do more harm than good



| Scale 0 10% of view | Scale 1 20% of view | Scale 2 30% of view | Scale 3 60% of view | Clean Image |

Image source: collages we made from Caltech 256 database

# Experiment Setup

- Experiment 1
  - Train on different scales, test on clean image

- **Experiment 2**
  - Train on different scales and clean, test on different scales



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view

Clean
Image

Image source: collages we made from Caltech 256 database

# Experiment 2 - objective

- Effect of 'bringing to face' for object-in-scene detection

- Questions to consider
  - Does 'cleaning' the scene decrease detection in cluttered environment?

# Experiment 2 - data



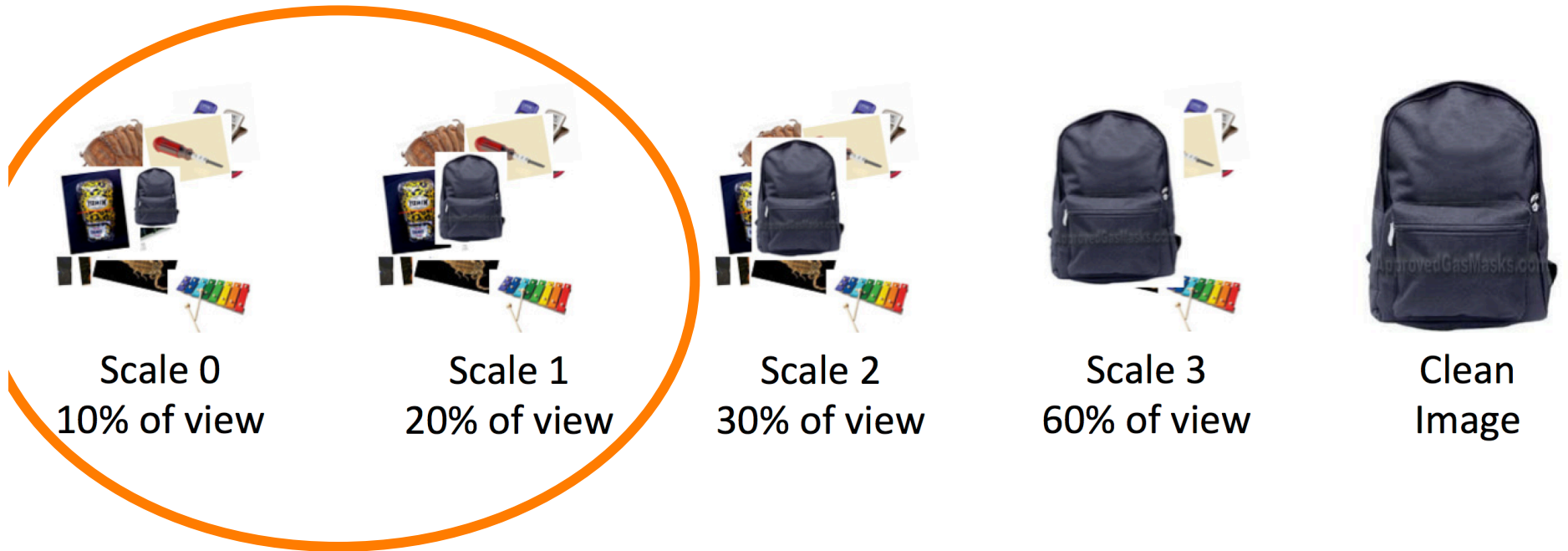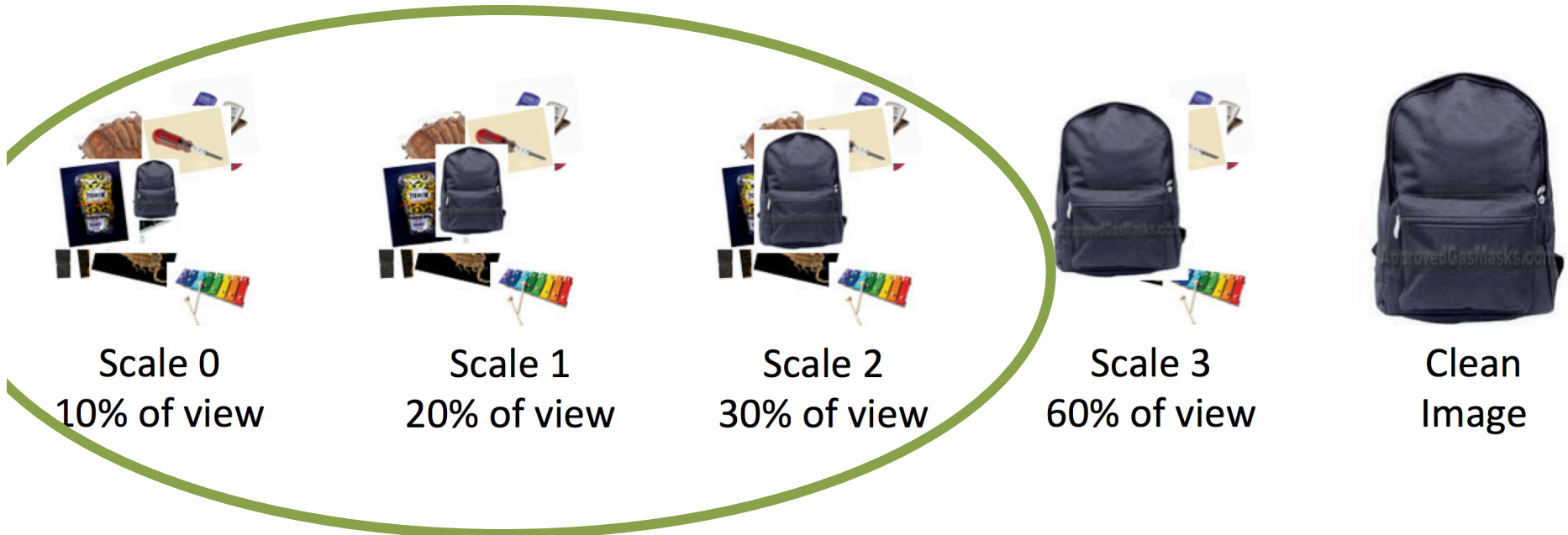| Scale 0 10% of view | Scale 1 20% of view | Scale 2 30% of view | Scale 3 60% of view | Clean Image |

Train0

Image source: collages we made from Caltech 256 database

# Experiment 2 - data



| Scale 0 10% of view | Scale 1 20% of view | Scale 2 30% of view | Scale 3 60% of view | Clean Image |

Train1

# Experiment 2 - data



| Scale 0 | Scale 1 | Scale 2 | Scale 3 | Clean |
| 10% of view | 20% of view | 30% of view | 60% of view | Image |

Train2

# Experiment 2 - data



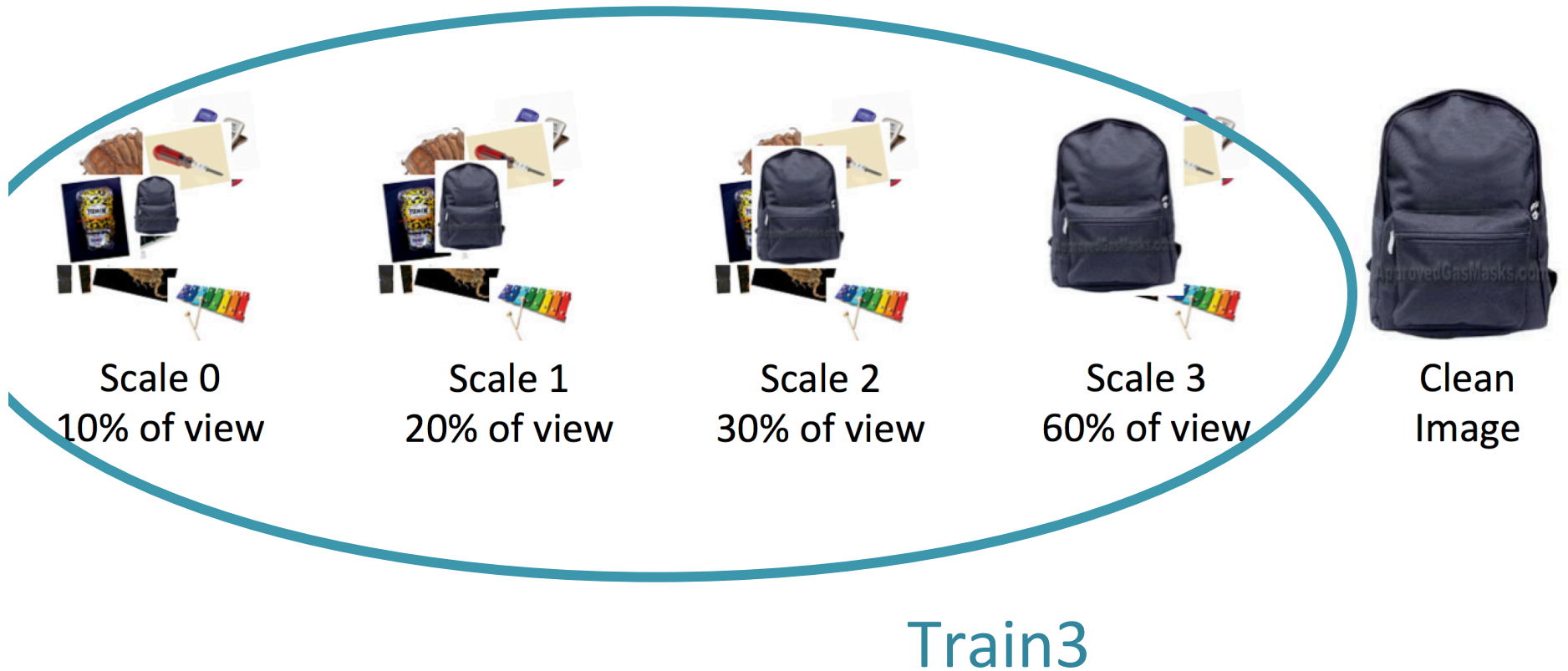| Scale 0<br>10% of view | Scale 1<br>20% of view | Scale 2<br>30% of view | Scale 3<br>60% of view | Clean<br>Image |

Train3

# Experiment 2 - data



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view

Clean
Image

TrainClean

Image source: collages we made from  Caltech 256 database

# Experiment 2 - data



| Scale 0<br>10% of view | Scale 1<br>20% of view | Scale 2<br>30% of view | Scale 3<br>60% of view | Clean<br>Image |

Correct number of epochs to compensate for more training examples

Image source: collages we made from Caltech 256 database

# Experiment 2 - data



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view

Clean
Image

Test0

On different images compared to train sets

Image source: collages we made from  Caltech 256 database

# Experiment 2 - data



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view
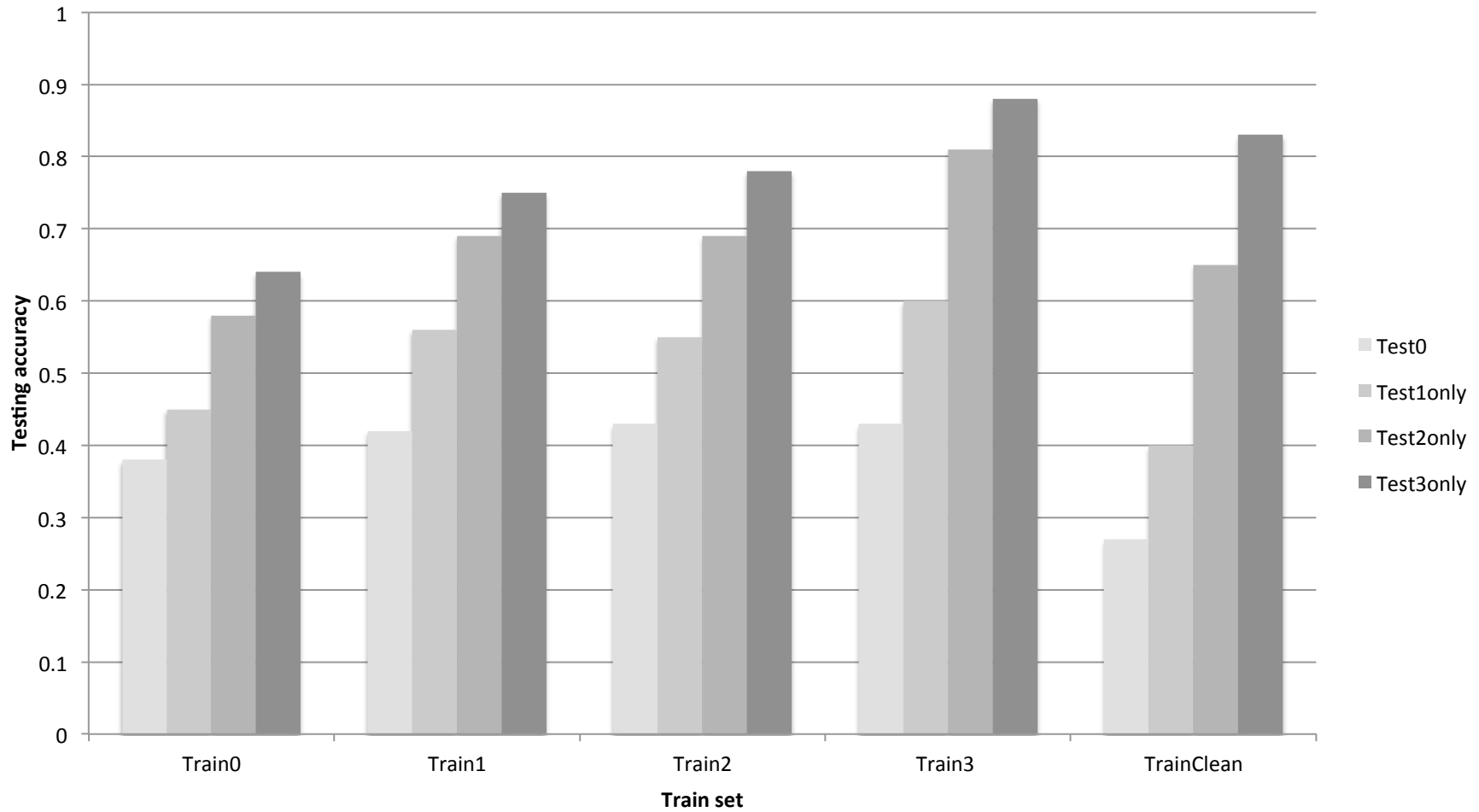
Clean
Image

Test1only

On different images compared to train sets

Image source: collages we made from Caltech 256 database

# Experiment 2 - data



| Scale 0<br>10% of view | Scale 1<br>20% of view | Scale 2<br>30% of view | Scale 3<br>60% of view | Clean<br>Image |

Test2only

On different images compared to train sets

Image source: collages we made from Caltech 256 database

# Experiment 2 - data



Scale 0
10% of view

Scale 1
20% of view
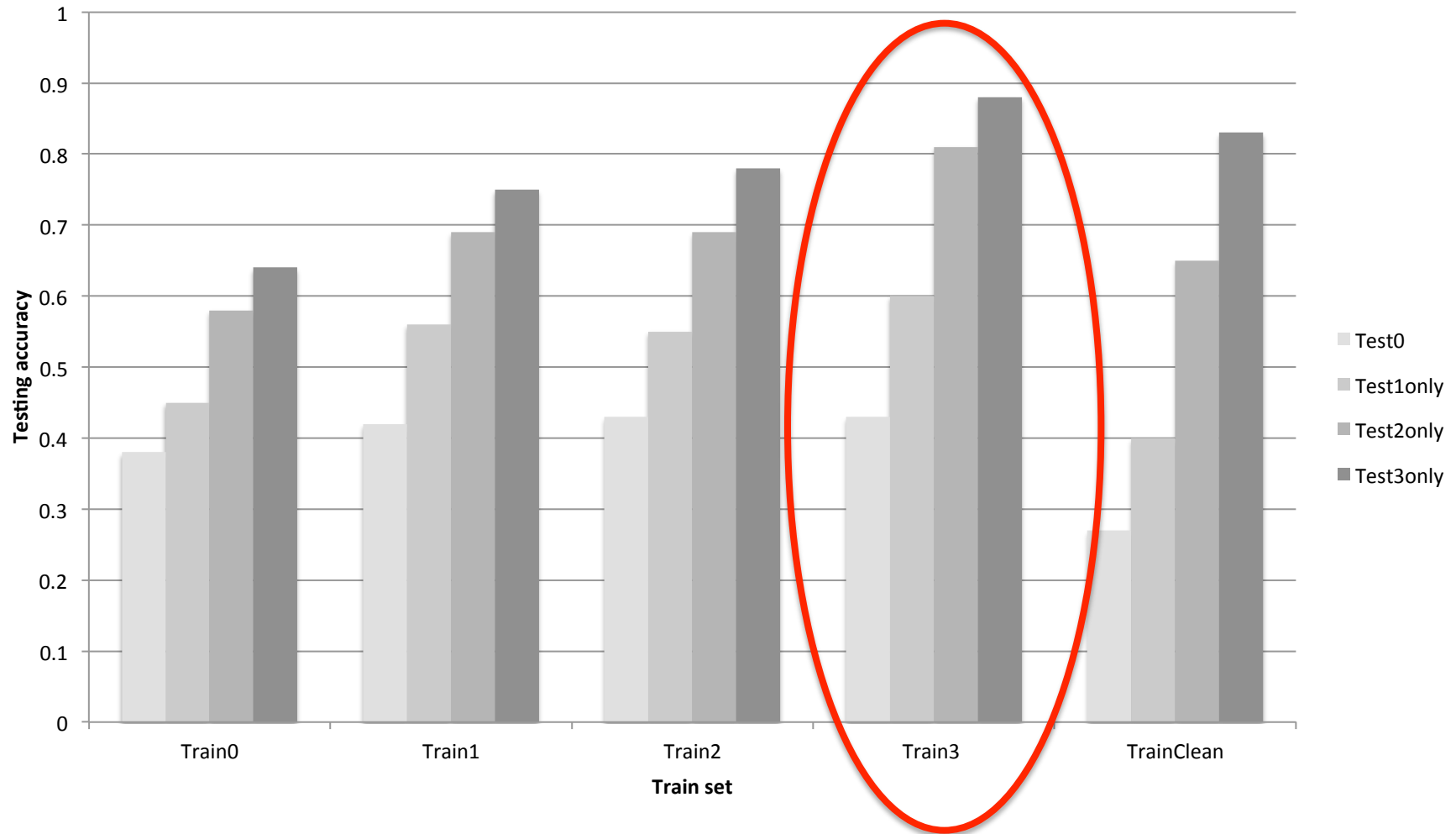
Scale 2
30% of view

Scale 3
60% of view

Clean
Image

Test3only

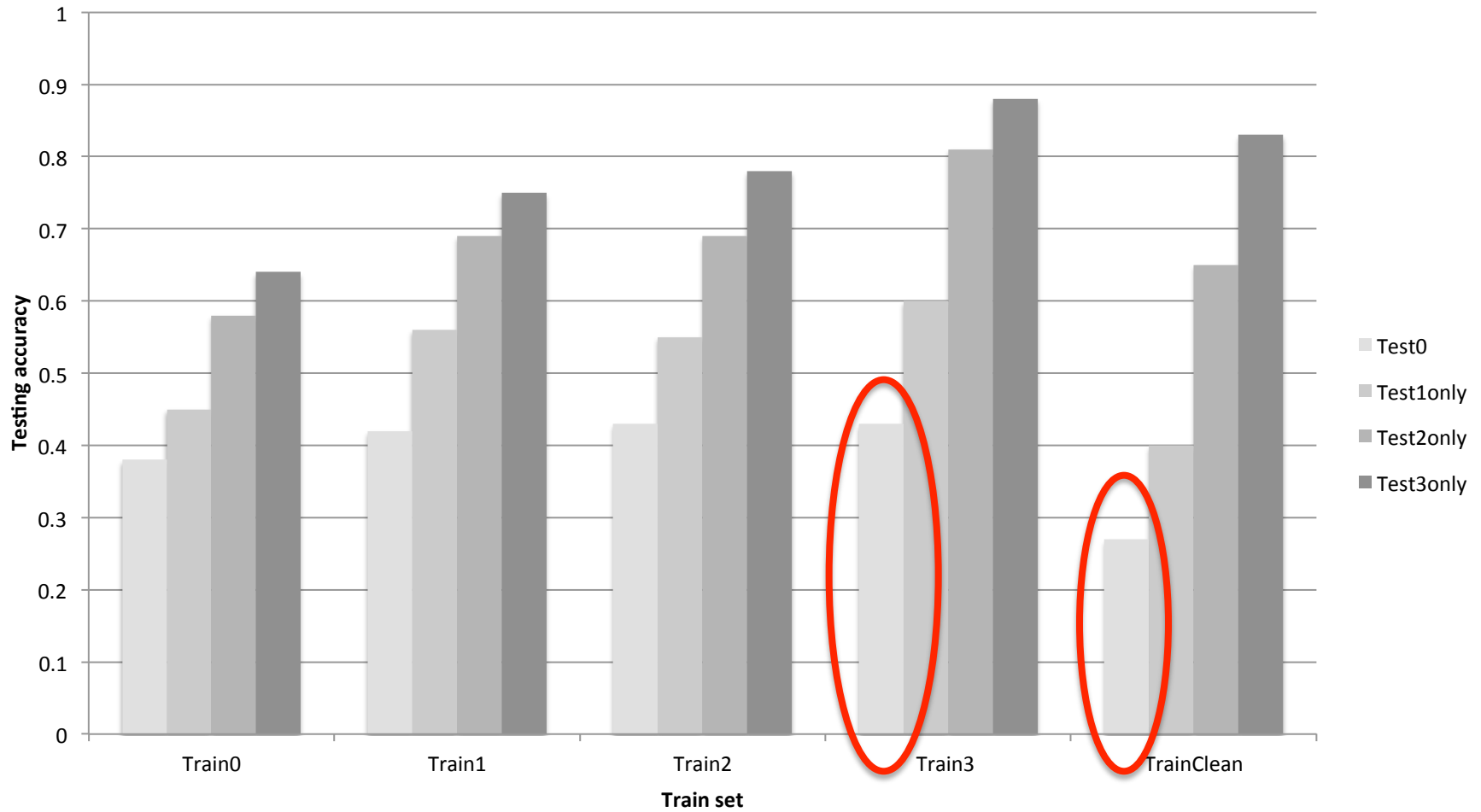On different images compared to train sets

Image source: collages we made from Caltech 256 database
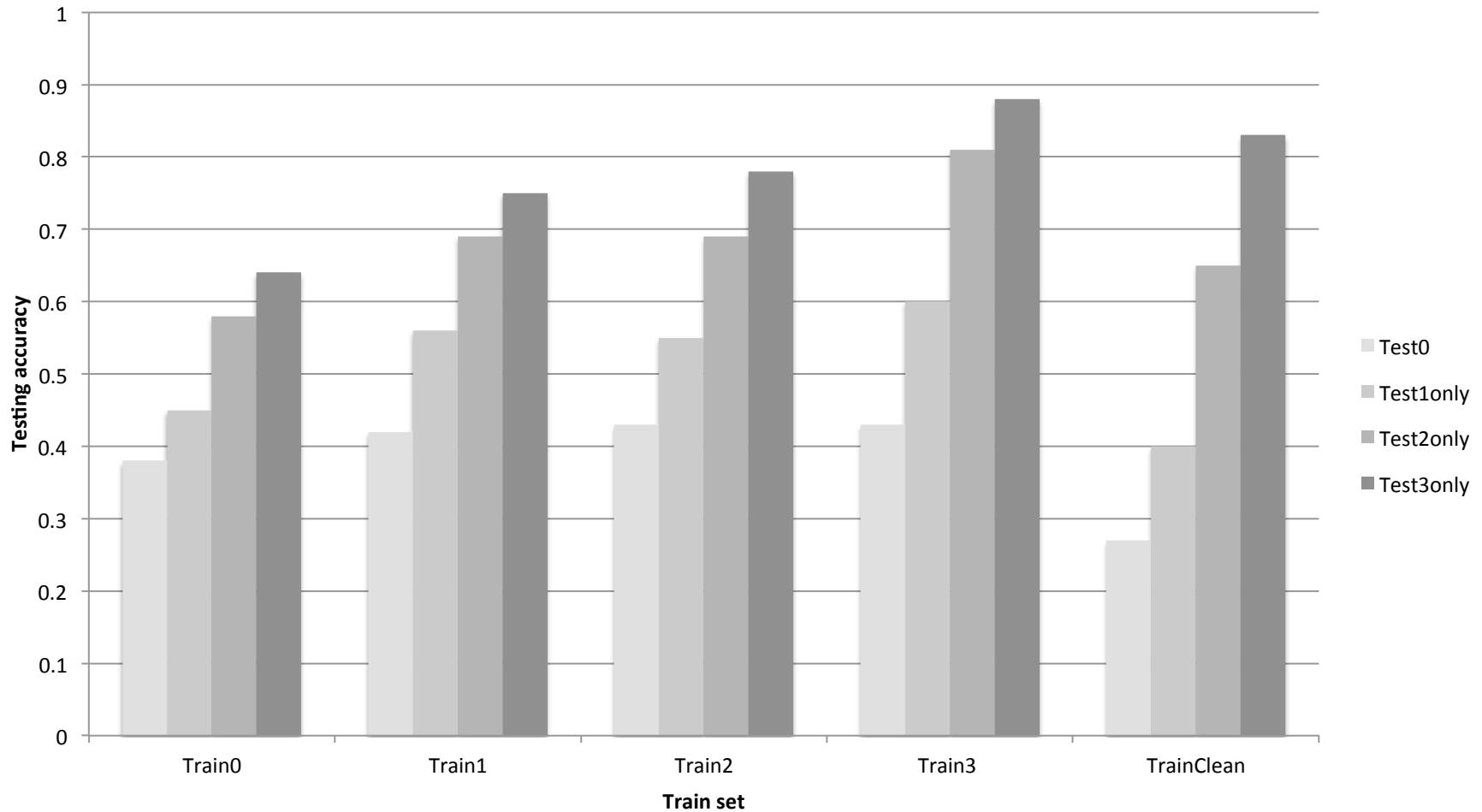
Experiment 2 - results

Experiment 2 - results

# Experiment 2 - results

# Experiment 2 - results



Training by 'bringing to face' yields to best accuracy

# Experiment 2 - conclusions

- Can learn more from different scales than from clean, as long as scale 3 is included

- Learning from different scales gives better accuracies when tested on lower scales

- Test on clean much better than test on scales



Scale 0
10% of view

Scale 1
20% of view

Scale 2
30% of view

Scale 3
60% of view

Clean
Image

Image source: collages we made from Caltech 256 database

# Conclusions

- With our controlled datasets, we could verify that network learns better from larger scale

- Testing needs to be done on clean images, no matter which scales were used in training

- Training on scales >30% gives more robustness when testing on all scales

- Training on scales <30% hurts accuracy