



# Learning with Side Information through Modality Hallucination

*Judy Hoffman Saurabh Gupta Trevor Darrell*  
*CVPR 2016*

Nan Yang

Supervisor: Benjamin Busam

Computer Aided Medical Procedures  
Technische Universität München

15.12.2016



# Outline

---

- 1 Introduction
- 2 Related work
- 3 Modality Hallucination Model
- 4 Experiments
- 5 Conclusion
- 6 Sources

## Introduction

Use available paired RGB-d training data to learn to hallucinate mid-level convolutional features from a RGB image for **object detection** task.

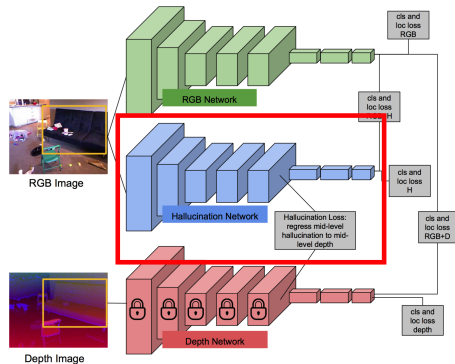


Figure: Network Structure<sup>1</sup>



# Introduction

---

- Task: Object Detection
- Framework: Deep Learning (Convolutional Neuron Networks)
- What's new: Additional ConvNets to extract depth features from RGB image (**modality hallucination**)

## Results

RGB object detection:

**34.0 mAP (this work)** vs 29.9 mAP (Fast R-CNN)

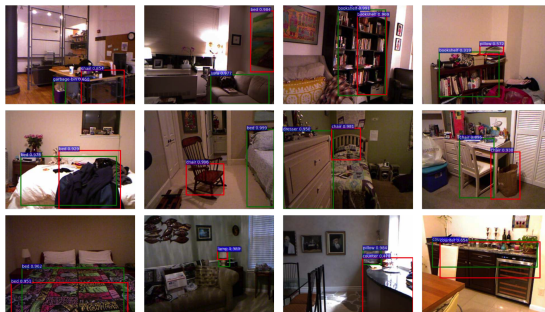


Figure 4: Example Detections on the NYUD2 *test set* where our RGB hallucination network's (green box) top scoring detection for the image is correct while the baseline RGB detector's (red box) top scoring detection is incorrect.

## Related Work

---

- Object detection using deep learning
  - RGB<sup>[6]</sup>
  - RGB-d<sup>[7]</sup>
- Transfer Learning<sup>[8][9]</sup>
- Learning using side information<sup>[10]</sup>
- Network distillation<sup>[11][12]</sup>

## Related Work

---

- **Object detection using deep learning**
  - RGB<sup>[6]</sup>
  - RGB-d<sup>[7]</sup>
- Transfer Learning<sup>[8][9]</sup>
- Learning using side information<sup>[10]</sup>
- Network distillation<sup>[11][12]</sup>



## Object Detection

---

- Extract region of interest (region proposal)
- Classification of this region
- Result: Region (bounding box coordinates) + Category



# Object Detection

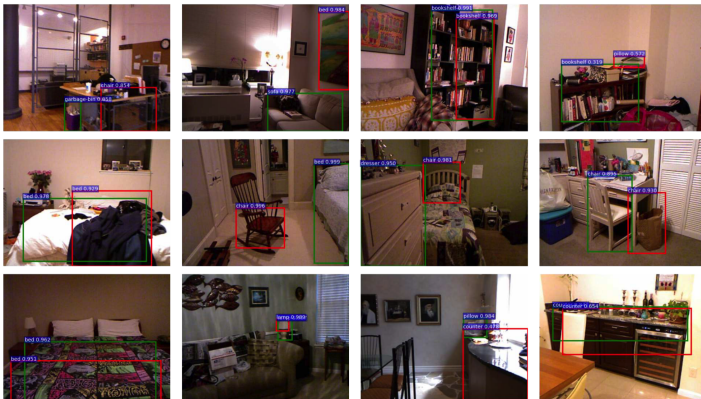
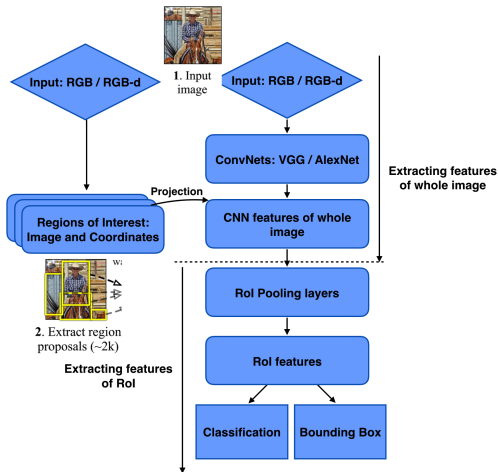


Figure 4: Example Detections on the NYUD2 *test set* where our RGB hallucination network's (green box) top scoring detection for the image is correct while the baseline RGB detector's (red box) top scoring detection is incorrect.

# Object Detection Pipeline (Fast R-CNN<sup>[6]</sup>)



# Object Detection Pipeline (Fast R-CNN<sup>[6]</sup>)

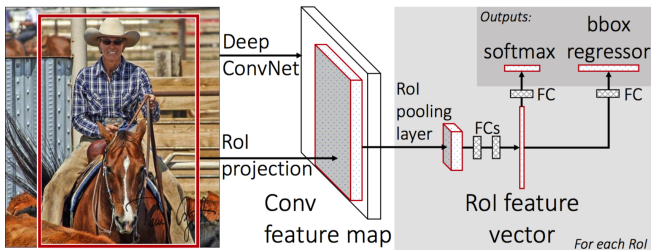


Figure: F-RCNN Structure<sup>6</sup>



## State of the art

---

- Fast R-CNN (Mean Average Precision: 29.9%)<sup>[6]</sup>
- RGB-d Fast R-CNN (Mean Average Precision: **44.4%**)<sup>[7]</sup>

## Observation

---

Object detection can be improved significantly with **RBD-d images** (additional modality, complementary information).

## Problem

---

Not all RGB images have corresponding depth information.

# Idea

---

How about extracting depth features from RGB images?

## Related Work

---

- Object detection using deep learning
  - RGB<sup>[6]</sup>
  - RGB-d<sup>[7]</sup>
- **Transfer Learning**<sup>[8][9]</sup>
- **Learning using side information**<sup>[10]</sup>
- **Network distillation**<sup>[11][12]</sup>



## Transfer Learning<sup>[8][9]</sup>

---

- Transfer learning and domain adaptation which learns to share information from one task (**Obj detection using depth features**) to another (**Obj detection with RGB features**).
- Previous work used depth information at training time to inform RGB test time detection by learning transformations into a **single representation** for the joint modality space.
- This work focuses on learning an **additional RGB representation**, which is the **hallucination network**.



## Learning using side information<sup>[10][13]</sup>

---

- Learning algorithm has additional knowledge at training time, whether meta data or an **additional modality**.
- Use this extra information to inform training of a **stronger model** than could be produced otherwise, e.g. surface normals at training time could produce detection improvement.

## Network distillation<sup>[11][12]</sup>

---

- Network Distillation: The output from one network (**depth network**) is used as the target probability distribution for a new network (**hallucination network**). Input can be different.
- This approach can also be seen as using distillation to learn representations on RGB images by **transferring supervision** from paired depth images, by employing joint training.

## Hallucination Network (Training time)

Sharing depth modality information with RGB through hallucination network.

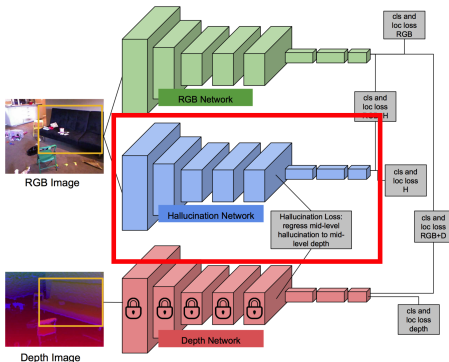


Figure: Network Structure<sup>1</sup>

## Hallucination Network (Test time)

Given **only a RGB image** and regions of interest, pass through both the RGB network and the hallucination network to produce **two scores** per category, per region, which we **average** and take the softmax to produce our final predictions.

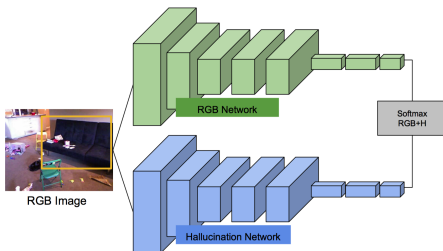


Figure: Test Network Structure<sup>2</sup>

## Explanation

1. A convolutional neuron network (**Hallucination network**) whose input is a **RGB image**.

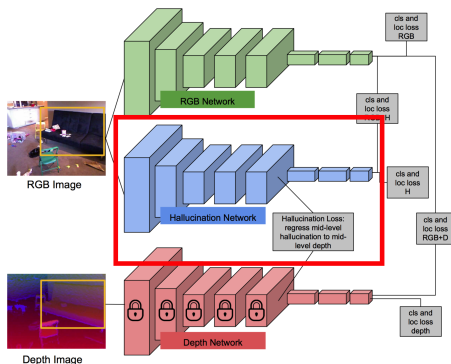


Figure: Network Structure<sup>1</sup>

## Explanation

2. The **features maps** generated by this network is as identical as possible to the ones generated by depth ConvNets.

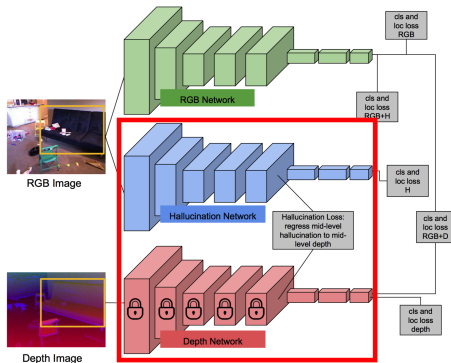


Figure: Network Structure<sup>1</sup>

## Explanation

3. A network whose input is RGB image generates depth features.

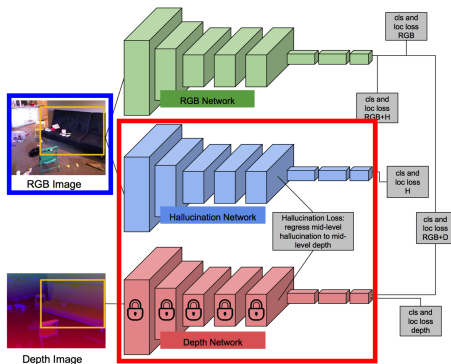


Figure: Network Structure<sup>1</sup>



## Explanation

4. Similar feature maps but different convolutional kernels ( weights ) because of different input modalities.

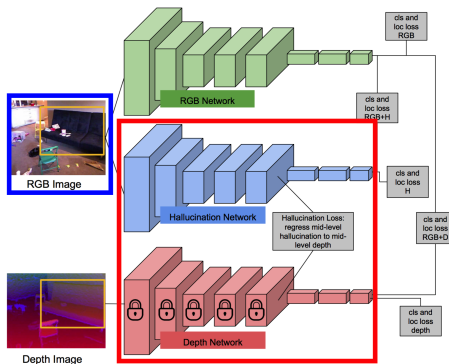


Figure: Network Structure<sup>1</sup>

## Architecture Definition

Multi channels ConvNets:

- RGB Network
- Hallucination Network
- Depth Network

The base network structures can be AlexNet / VGG (to extract features).

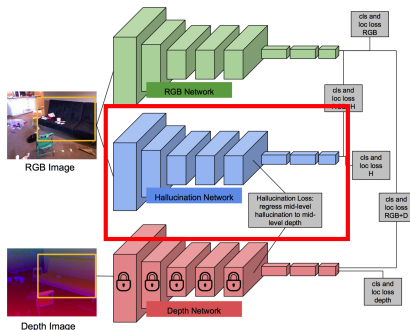


Figure: Network Structure<sup>1</sup>

# Modality Hallucination Model

## Training time

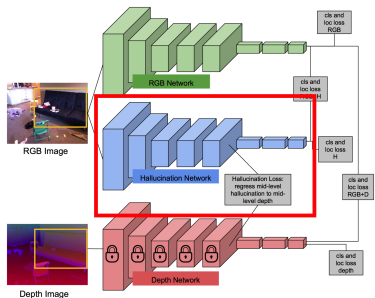


Figure: Network Structure<sup>1</sup>

## Testing time

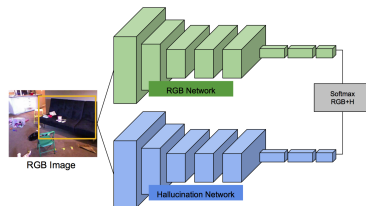


Figure: Test Network Structure<sup>2</sup>

# Training Pipeline

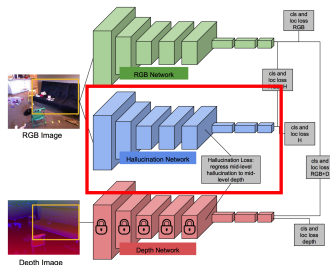
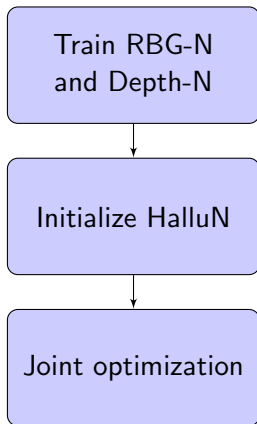
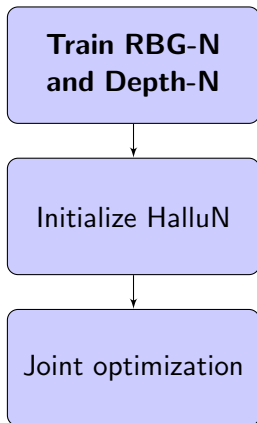


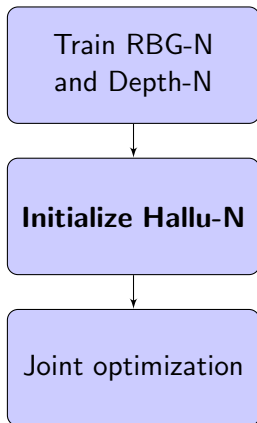
Figure: Network Structure<sup>1</sup>

## Training Pipeline



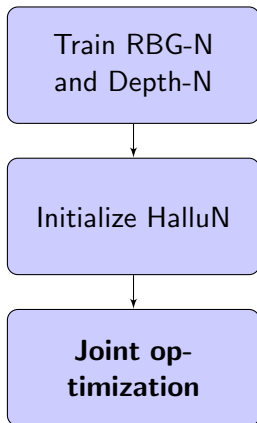
RGB and depth network are independently trained using the Fast R-CNN algorithm with the corresponding image input.

## Training Pipeline



The hallucination network parameters are initialized with the learned **depth network weights**.

## Training Pipeline



Multi-task Optimization with **hallucination loss**.

## Hallucination Loss

Activations after some layer,  $l$ , should be similar between the hallucination and depth networks.

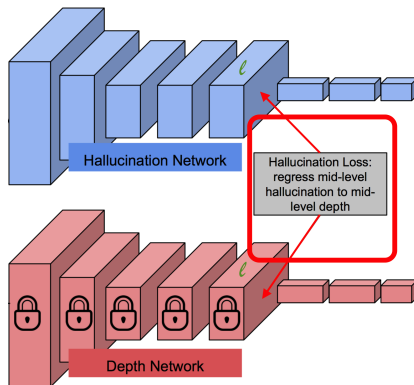


Figure: Network Structure<sup>1</sup>



## Hallucination Loss

Euclidean loss:  $L_{hallucinate}(l) = \|\sigma(A_l^{dNet}) - \sigma(A_l^{hNet})\|_2^2$   
 where  $\sigma(x) = 1/(1 + \exp(-x))$  (sigmoid function) and  $l = layer$

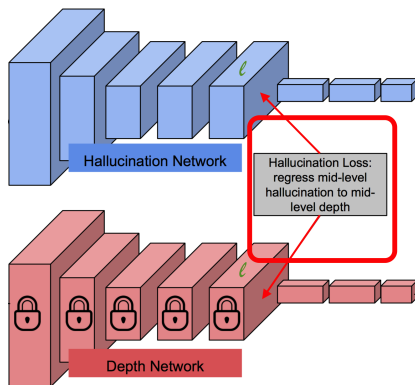


Figure: Network Structure<sup>1</sup>

## Hallucination Loss

Asymmetric transfer of information: **0** learning rate lower than  $l$

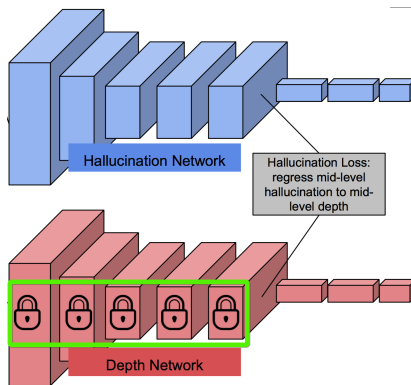
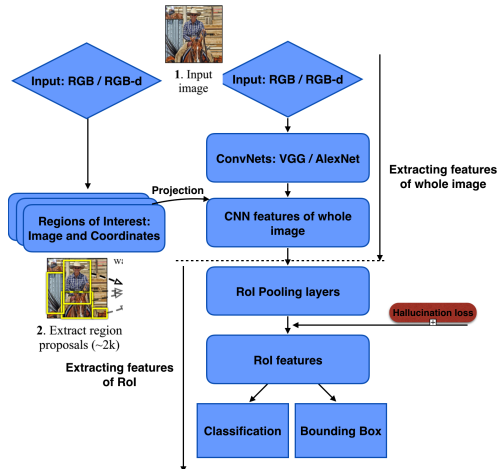


Figure: Network Structure<sup>1</sup>

# Hallucination Loss





## Multi-task Optimization

---

11 total losses:

- 5 softmax cross-entropy losses for **classification**
- 5 smooth L1 losses for **bounding box** coordinates regression
- 1 hallucination loss

# Multi-task Optimization

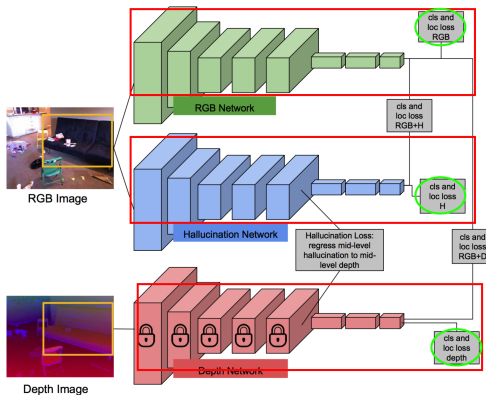


Figure: Network Structure<sup>1</sup>

# Multi-task Optimization

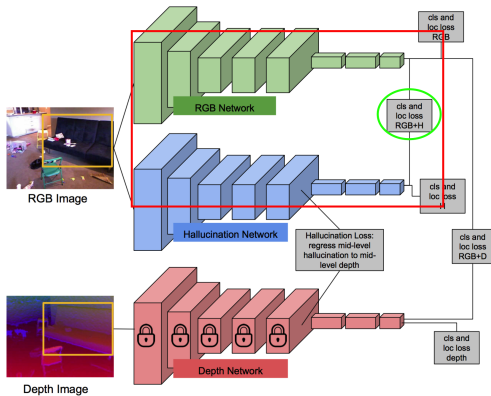


Figure: Network Structure<sup>1</sup>

# Multi-task Optimization

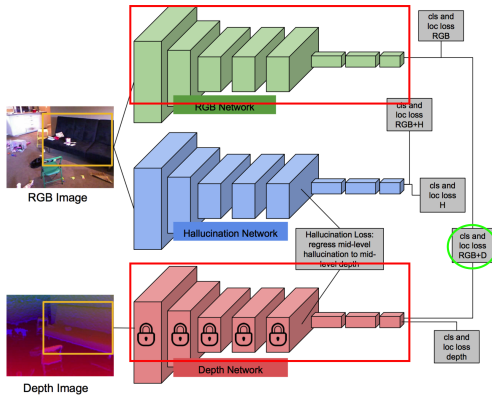


Figure: Network Structure<sup>1</sup>

# Multi-task Optimization

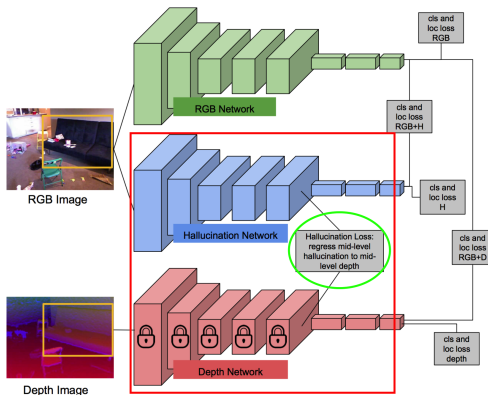


Figure: Network Structure<sup>1</sup>



## Balancing objectives

---

$$\begin{aligned}\mathcal{L} = & \gamma \mathcal{L}_{\text{hallucinate}} \\ & + \alpha \left[ \mathcal{L}_{\text{loc}}^{\text{dNet}} + \mathcal{L}_{\text{loc}}^{\text{rNet}} + \mathcal{L}_{\text{loc}}^{\text{hNet}} + \mathcal{L}_{\text{loc}}^{\text{rdNet}} + \mathcal{L}_{\text{loc}}^{\text{rhNet}} \right] \\ & + \beta \left[ \mathcal{L}_{\text{cls}}^{\text{dNet}} + \mathcal{L}_{\text{cls}}^{\text{rNet}} + \mathcal{L}_{\text{cls}}^{\text{hNet}} + \mathcal{L}_{\text{cls}}^{\text{rdNet}} + \mathcal{L}_{\text{cls}}^{\text{rhNet}} \right]\end{aligned}$$

$$\alpha = 0.5, \beta = 1, \gamma = ?$$

$\gamma$  is heuristically set to **10x** the contribution of any of the other losses

## Experiments

---

- Base Network: AlexNet, VGG
- Region Proposals: Multiscale Combinatorial Grouping (MCG)
- SGD Hyper-parameters: LR = 0.001, Momentum = 0.9, Weight decay = 0.0005, T = 10 (Threshold of gradient clipping)

## NYUD2 Detection Evaluation

A - AlexNet, V - VGG

method	btub	bed	bshelf	box	chair	counter	desk	door	dresser	gbin	lamp	monitor	nstand	pillow	sink	sofa	table	tv	toilet	mAP
RGB only [10] (A)	7.5	50.6	<b>36.8</b>	1.4	30.2	34.9	10.8	21.5	27.8	16.9	26.0	32.6	20.6	<b>25.1</b>	31.6	<b>36.7</b>	14.8	25.1	54.6	26.6
RGB ensemble (A-A)	10.5	53.7	33.6	1.6	32.0	34.8	12.2	20.8	34.5	19.6	28.6	45.7	28.5	24.4	31.4	34.7	14.5	<b>34.0</b>	<b>56.1</b>	29.0
Our Net (A-RGB, A-H)	<b>13.9</b>	<b>56.1</b>	34.4	<b>1.9</b>	<b>32.9</b>	<b>40.5</b>	<b>12.9</b>	<b>22.6</b>	<b>37.4</b>	<b>22.0</b>	<b>28.9</b>	<b>46.2</b>	<b>31.9</b>	22.9	<b>34.2</b>	34.2	<b>19.4</b>	33.2	53.6	<b>30.5</b>
RGB only [10] (V)	15.6	59.4	38.2	1.9	33.8	36.3	12.1	24.5	31.6	18.6	25.5	46.5	30.1	20.6	30.3	40.5	19.5	<b>37.8</b>	45.7	29.9
RGB ensemble (A-V)	14.8	60.4	<b>43.1</b>	<b>2.1</b>	36.4	40.7	13.3	<b>27.1</b>	35.5	20.8	29.9	<b>52.9</b>	<b>33.5</b>	26.2	33.0	44.4	19.9	36.7	50.2	32.7
Our Net (A-RGB, V-H)	<b>16.8</b>	<b>62.3</b>	41.8	<b>2.1</b>	<b>37.3</b>	<b>43.4</b>	<b>15.4</b>	24.4	<b>39.1</b>	<b>22.4</b>	<b>30.3</b>	46.6	30.9	<b>27.0</b>	<b>42.9</b>	<b>46.2</b>	<b>22.2</b>	34.1	<b>60.4</b>	<b>34.0</b>

Figure: NYUD2 evaluation<sup>3</sup>

## NYUD2 Detection Evaluation

A - AlexNet, V - VGG

method	mAP
RGB only [10] (A)	26.6
RGB ensemble (A-A)	29.0
<b>Our Net (A-RGB, A-H)</b>	<b>30.5</b>
RGB only [10] (V)	29.9
RGB ensemble (A-V)	32.7
<b>Our Net (A-RGB, V-H)</b>	<b>34.0</b>

Figure: NYUD2 evaluation<sup>3</sup>

## How to initialize the hallucination net?

Initial Weights	bathub	bed	bshelf	box	chair	counter	desk	door	dresser	gbin	lamp	monitor	nstand	pillow	sink	sofa	table	tv	toilet	mAP
RGB	7.5	50.4	9.9	<b>0.9</b>	<b>26.2</b>	<b>24.9</b>	5.8	<b>15.8</b>	13.0	<b>29.8</b>	12.0	43.1	20.9	14.7	17.9	25.3	<b>15.1</b>	32.5	59.1	22.4
depth	9.9	<b>52.4</b>	<b>14.9</b>	<b>0.9</b>	24.9	24.4	4.3	15.3	18.1	24.1	<b>14.8</b>	<b>45.8</b>	<b>27.2</b>	<b>18.5</b>	<b>21.3</b>	<b>29.0</b>	13.7	<b>33.6</b>	<b>66.4</b>	<b>24.2</b>
random	<b>10.5</b>	47.6	12.3	0.6	23.5	20.2	<b>6.0</b>	13.0	<b>19.3</b>	12.0	13.3	42.8	12.8	12.1	13.6	23.0	13.9	28.6	61.5	20.3

Figure: Initialization evaluation<sup>4</sup>

## How to initialize the hallucination net?

Initial Weights	mAP
RGB	22.4
depth	<b>24.2</b>
random	20.3

Figure: Initialization evaluation<sup>4</sup>

## What did the hallucination net learn?

roi-pool5 activations from corresponding regions in the test image  
which have highest final detection scores

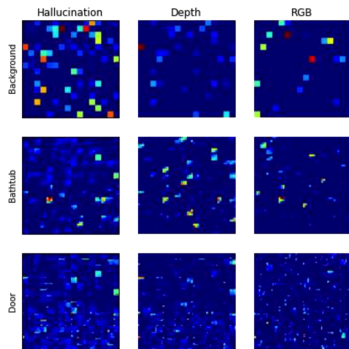


Figure: Activation of Rol Pool5<sup>5</sup>

## What did the hallucination net learn?

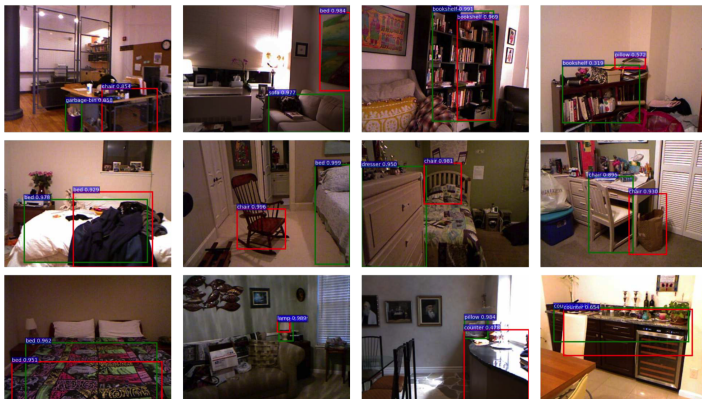


Figure 4: Example Detections on the NYUD2 *test set* where our RGB hallucination network's (green box) top scoring detection for the image is correct while the baseline RGB detector's (red box) top scoring detection is incorrect.



## What did the hallucination net learn?



## What did the hallucination net learn?



## But it's not perfect.

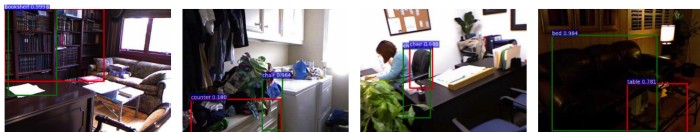


Figure 5: Example Detections on the NYUD2 *test set* where our RGB hallucination network's (green box) top scoring detection for the image is a false positive while the baseline RGB detector's (red box) top scoring detection is a true positive.

But it's not perfect.

---





## Conclusion

---

- A novel technique for extracting **unseen modality** features.
- Outperforms the corresponding Fast R-CNN RGB detection models on the NYUD2 dataset.

## Sources

---

- [1][2][3][4][5]: J. Hoffman, S. Gupta, T. Darrell, Learning with Side Information through Modality Hallucination
- [6]: R. Girshick, Fast R-CNN
- [7]: S. Gupta, R. Girshick, P. Arbelaez, and J. Malik. Learning rich features from rgb-d images for object detection and segmentation.
- [8]: L. Chen, W. Li, and D. Xu. Recognizing rgb images by learning from rgb-d data.
- [9]: L. Spinello and K.O. Arras, Leveraging rgb-d data: Adaptive fusion and domain adaptation for object detection.
- [10]: V. Vapnik and A. Vashist, A new learning paradigm: Learning using privileged information.
- [11]: G. Hinton, O. Vinyals, and J. Dean, Distilling the knowledge in a neural network.
- [12]: S. Gupta, J. Hoffman, and J. Malik. Cross modal distillation for supervision transfer.
- [13]: A. Shrivastava and A. Gupta. Building part-based object detectors via 3d geometry. In ICCV, 2013.



Thanks.