

# Zero-shot Task Transfer



భారతదేశమీ సిద్ధా విభాగ  
Indian Institute of Technology Hyderabad

**Vineeth N Balasubramanian**

(Joint work with Arghya Pal, PhD student)

Indian Institute of Technology, Hyderabad, INDIA

Presented as Oral at



# Our Group's Research

## Algorithmic

- Non-convex optimization for DL\*
- Explainable ML§
- Learning with low supervision⌘

## Applied

- Recognition of Expressions, Poses, Gestures, Actions
- Vision on UAVs/Drones
- Computer Vision for Agriculture
- Autonomous Navigation

§ Neural Network Attributions: A Causal Perspective, **ICML 2019**

⌘Zero-shot Task Transfer, **CVPR 2019**

\* Submodular Batch Selection for Training Deep Neural Networks, **IJCAI 2019**

\* On Noise and Optimality in Neural Networks, **ICML 2018 Workshops**

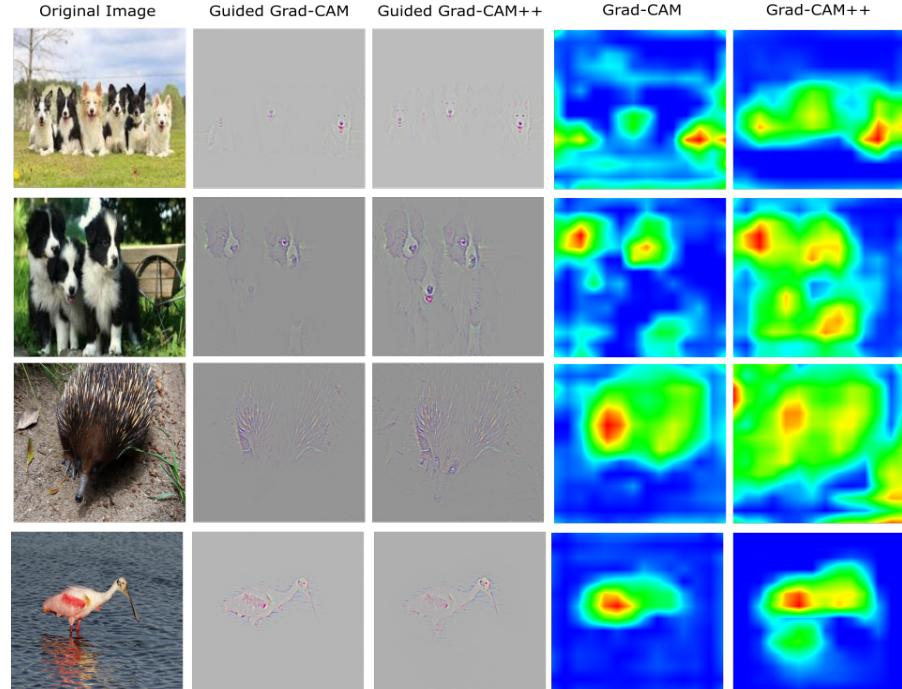
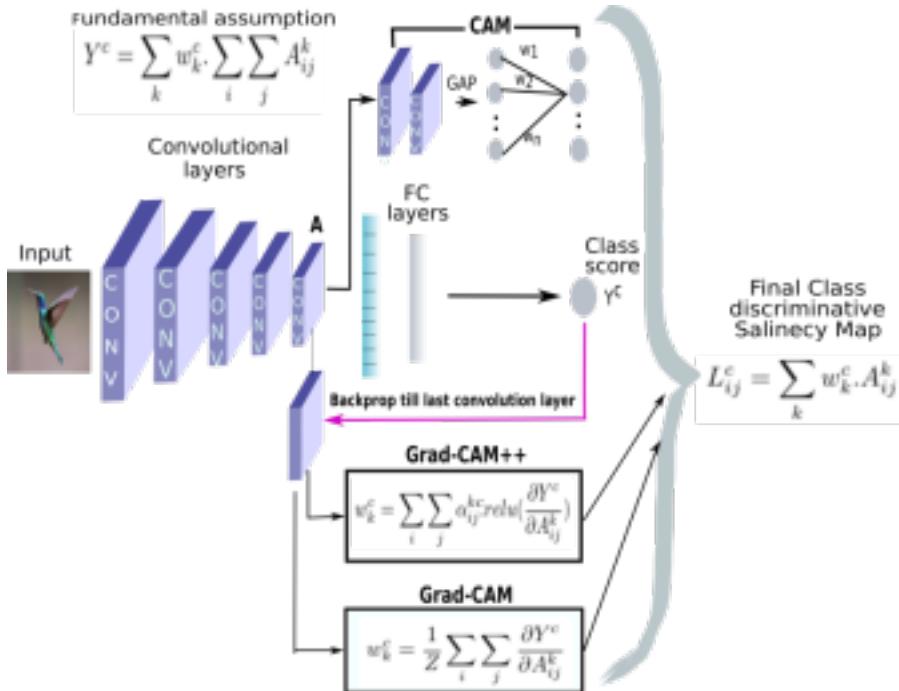
⌘Adversarial Data Programming, **CVPR 2018**

§ Grad-CAM++: Generalized Gradient-based Visual Explanations for Convolutional Networks, **WACV 2018**

⌘Attentive Semantic Video Generation using Captions, **ICCV 2017, ACM MM 2017**



# Grad-CAM++



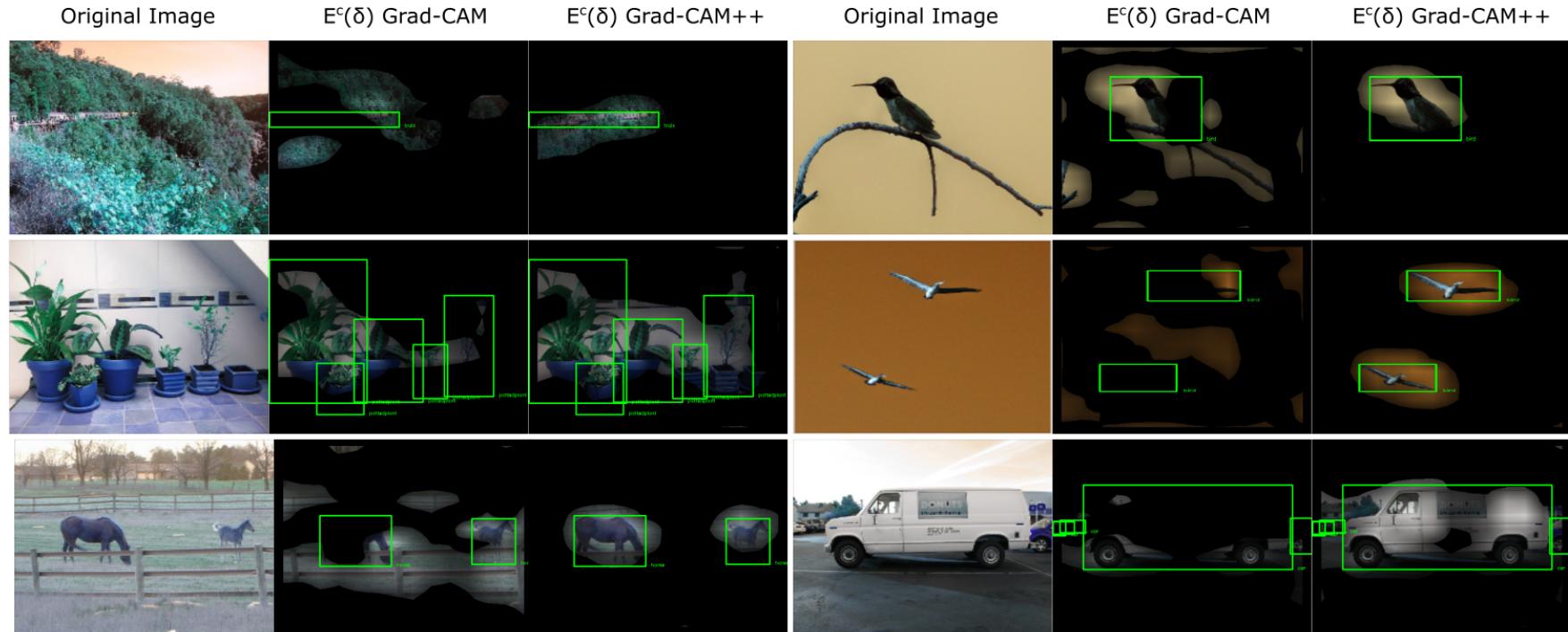
Chattopadhyay, Sarkar, Howlader, Balasubramanian, WACV 2018



Vineeth N B  
Indian Institute of Technology, Hyderabad, INDIA

Zero-shot Task  
Transfer

# Grad-CAM++



Chattopadhyay, Sarkar, Howlader, Balasubramanian, WACV 2018



Vineeth N B  
Indian Institute of Technology, Hyderabad, INDIA

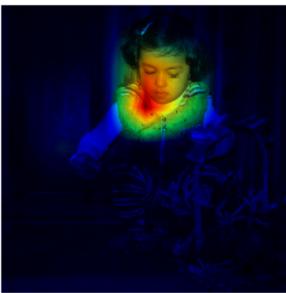
Zero-shot Task  
Transfer

# Grad-CAM++

Original Image



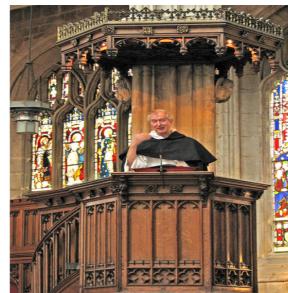
Grad-CAM



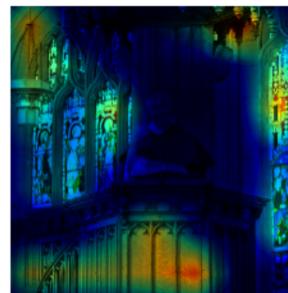
Grad-CAM++



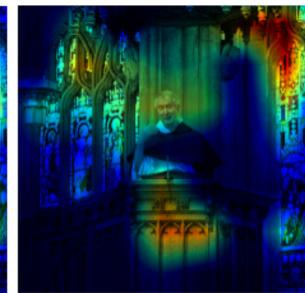
Original Image



Grad-CAM

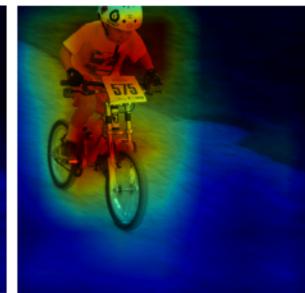
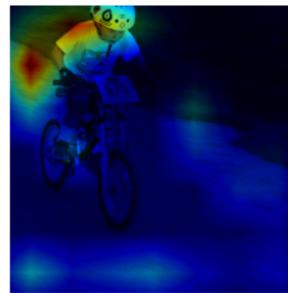
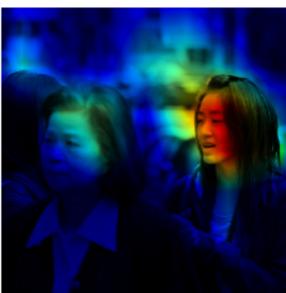


Grad-CAM++



A young girl accompanied by a small plant.

Photo of men act under colored pillars in a museum.



Two girls focussed on their faces on a sunny day

A motocross bike race four little kids are riding a bike race.

Chattopadhyay, Sarkar, Howlader, Balasubramanian, WACV 2018

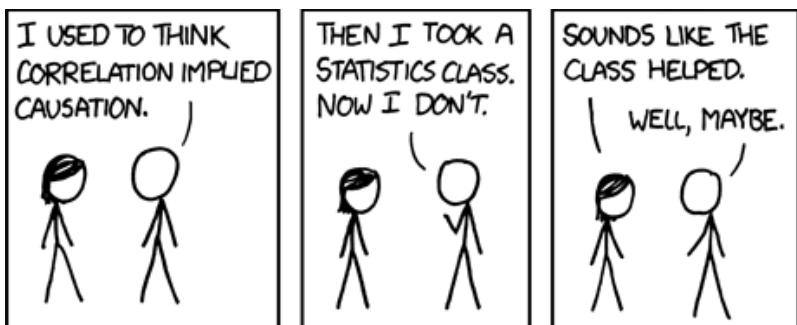


Vineeth N B  
Indian Institute of Technology, Hyderabad, INDIA

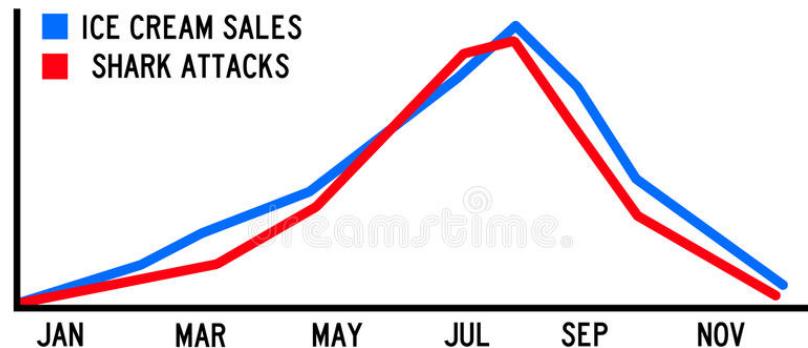
Zero-shot Task  
Transfer

# Causal Attributions in Neural Networks

- Is feature correlation to output a true indicator of explainability?
- Or do we need to find causal relationships in the analyzed data-output pairs?



**CORRELATION IS NOT CAUSATION!**

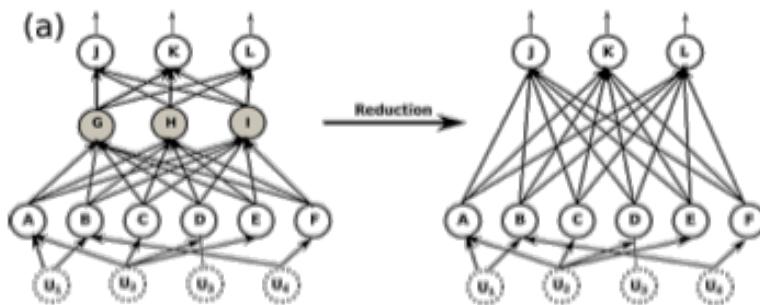


Both ice cream sales and shark attacks increase when the weather is hot and sunny, but they are not caused by each other (they are caused by good weather, with lots of people at the beach, both eating ice cream and having a swim in the sea)

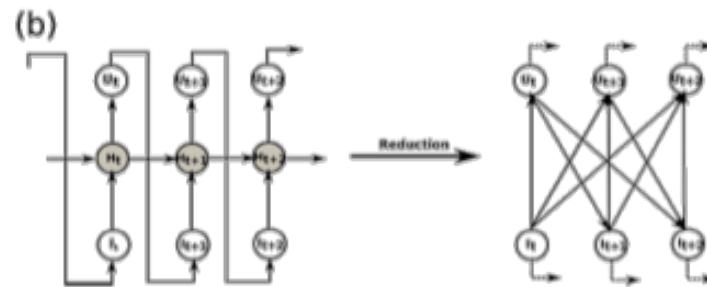
Chattopadhyay, Manupriya, Sarkar, Balasubramanian, ICML 2019

# Causal Attributions in Neural Networks

Neural network as a SCM



Feedforward neural network



Recurrent neural network

Chattopadhyay, Manupriya, Sarkar, Balasubramanian, ICML 2019



Vineeth N B

Indian Institute of Technology, Hyderabad, INDIA

Zero-shot Task  
Transfer

# Causal Attributions in Neural Networks

$$\mathbb{E}[y|do(x = 1)] - \mathbb{E}[y|do(x = 0)]$$

For continuous variables:

$$ACE_{do(x_i=\alpha)}^y = \mathbb{E}[y|do(x_i = \alpha)] - baseline_{x_i}$$

where baseline is defined as:

$$\mathbb{E}_{x_i}[\mathbb{E}_y[y|do(x_i = \alpha)]] \quad \text{the average ACE across all } x_i$$

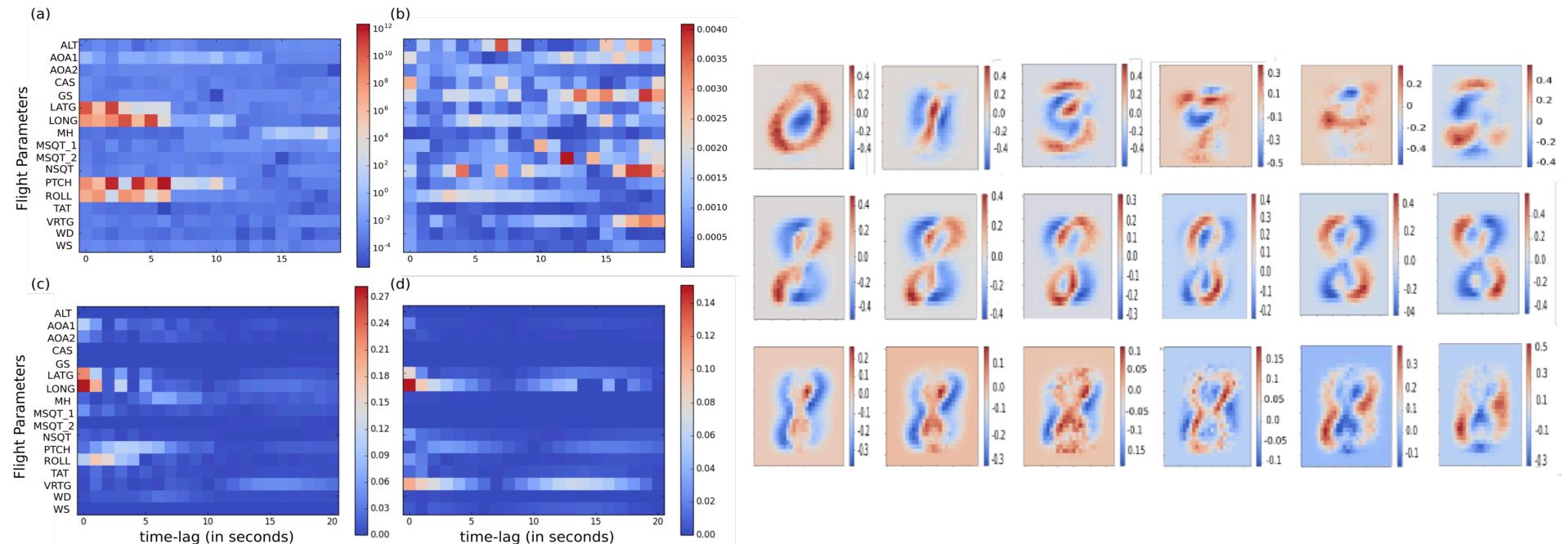
Chattopadhyay, Manupriya, Sarkar, Balasubramanian, ICML 2019



Vineeth N B  
Indian Institute of Technology, Hyderabad, INDIA

Zero-shot Task  
Transfer

# Causal Attributions in Neural Networks



Chattopadhyay, Manupriya, Sarkar, Balasubramanian, ICML 2019



Vineeth N B  
Indian Institute of Technology, Hyderabad, INDIA

Zero-shot Task  
Transfer

# Our Group's Research

## Algorithmic

- Non-convex optimization for DL\*
- Explainable ML§
- Learning with low supervision⌘

## Applied

- Recognition of Expressions, Poses, Gestures, Actions
- Vision on UAVs/Drones
- Computer Vision for Agriculture
- Autonomous Navigation

§ Neural Network Attributions: A Causal Perspective, **ICML 2019**

⌘Zero-shot Task Transfer, **CVPR 2019**

\* Submodular Batch Selection for Training Deep Neural Networks, **IJCAI 2019**

\* On Noise and Optimality in Neural Networks, **ICML 2018 Workshops**

⌘Adversarial Data Programming, **CVPR 2018**

§ Grad-CAM++: Generalized Gradient-based Visual Explanations for Convolutional Networks, **WACV 2018**

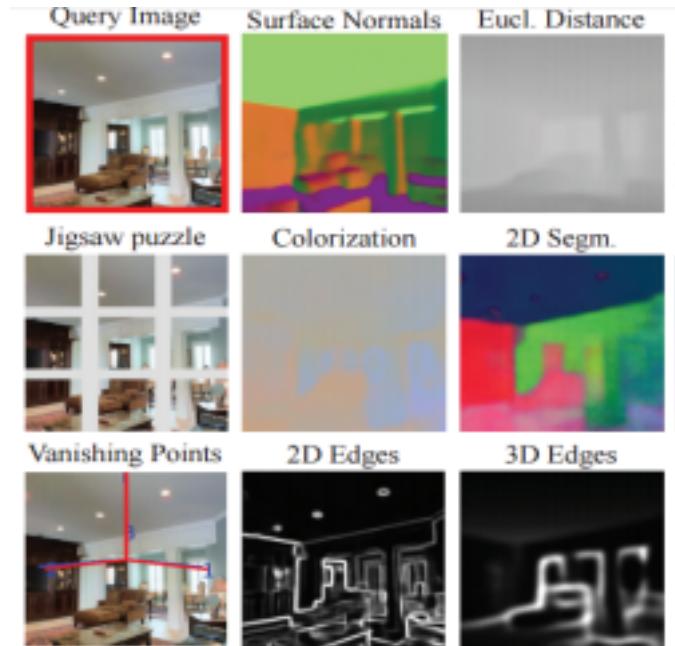
⌘Attentive Semantic Video Generation using Captions, **ICCV 2017, ACM MM 2017**



# Introduction

- Consider  $K$  vision tasks to accomplish:  $\tau_1, \dots, \tau_K$ 
  - Object recognition
  - Depth estimation
  - Edge detection
  - Pose estimation
  - ... etc etc...

How to automatically predict model parameters  
for tasks without ground truth?



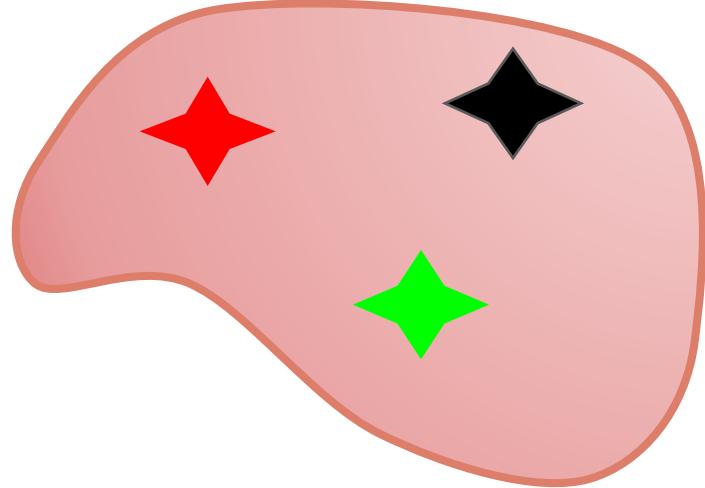
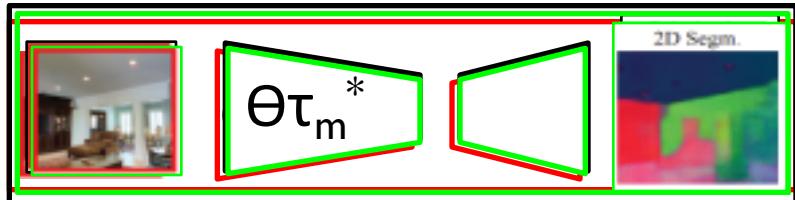
# Zero-shot Tasks vs Known Tasks

- Of  $K$  tasks:
- Ground-truth available for  $m$  tasks,  $\{\tau_1, \dots, \tau_m\}$  → **Known Tasks**
- No ground truth available for tasks  $\{\tau_{(m+1)}, \dots, \tau_K\}$  **Zero-shot Tasks**



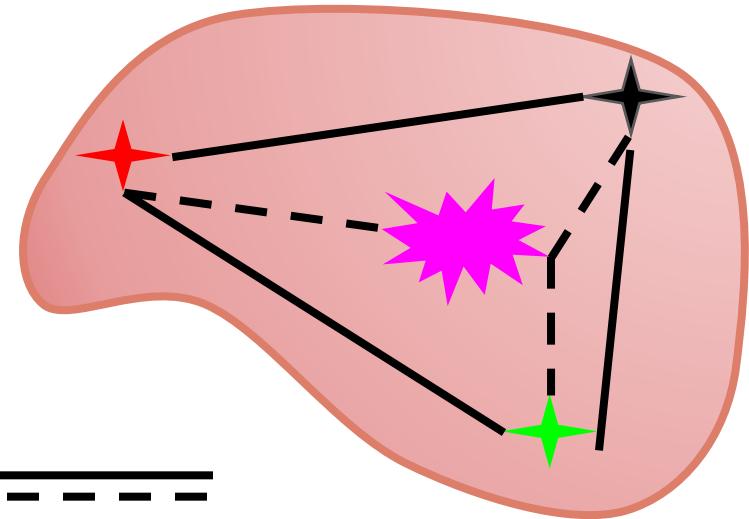
# Task Meta-Manifold

- Known Tasks:  $\{\tau_1, \dots, \tau_m\}$
- We can learn their corresponding model parameters on task meta-manifold
- Not possible for **zero-shot task**



# Zero-shot Task Transfer: Our Objective

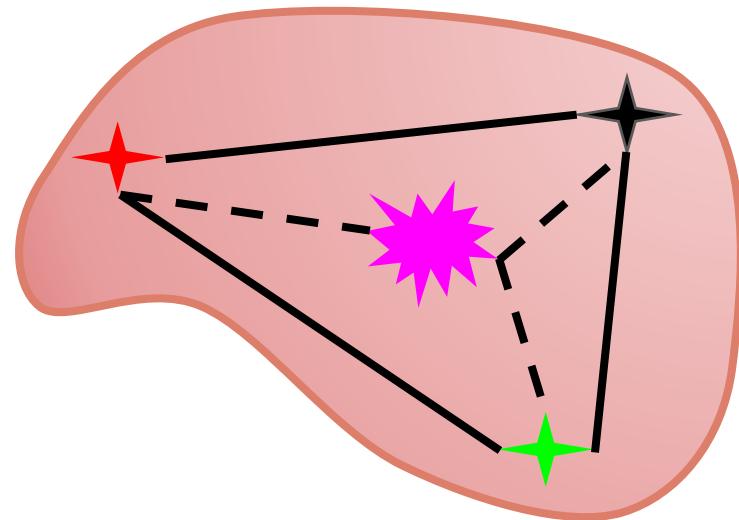
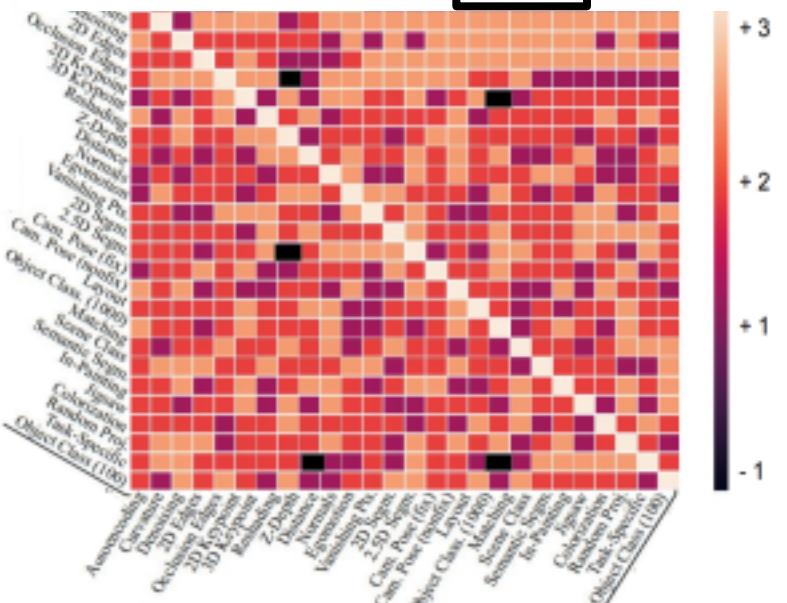
- Learn a meta-learning function  $F$  to **regress the unknown zero-shot model parameters**; given
  - Model parameters of known tasks 
  - Pair-wise task correlations 



$$\mathcal{F}(\theta_{\tau_1}, \dots, \theta_{\tau_m}, \Gamma) = \theta_{\tau_j}$$

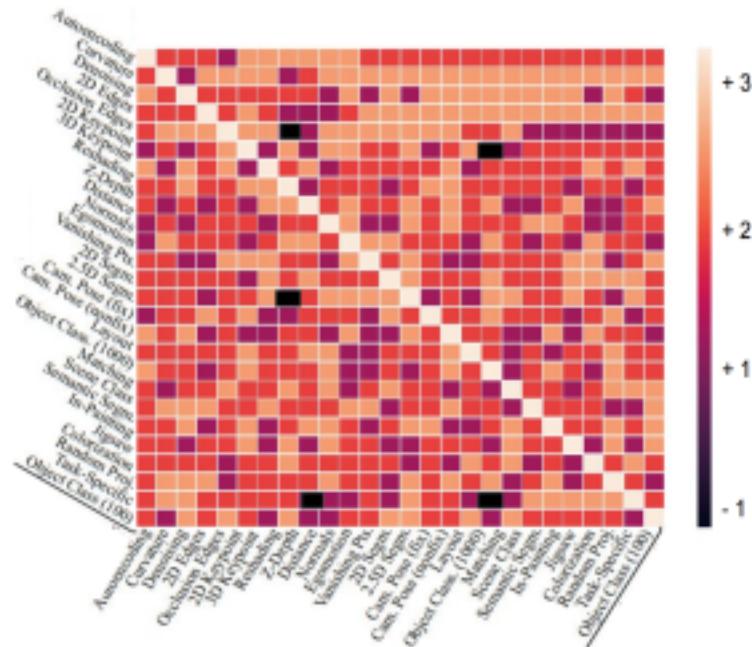
# Pairwise Task Correlation

$$\mathcal{F}(\theta_{\tau_1}, \dots, \theta_{\tau_m}, \Gamma) = \theta_{\tau_j} \quad \text{A matrix of pairwise correlations}$$



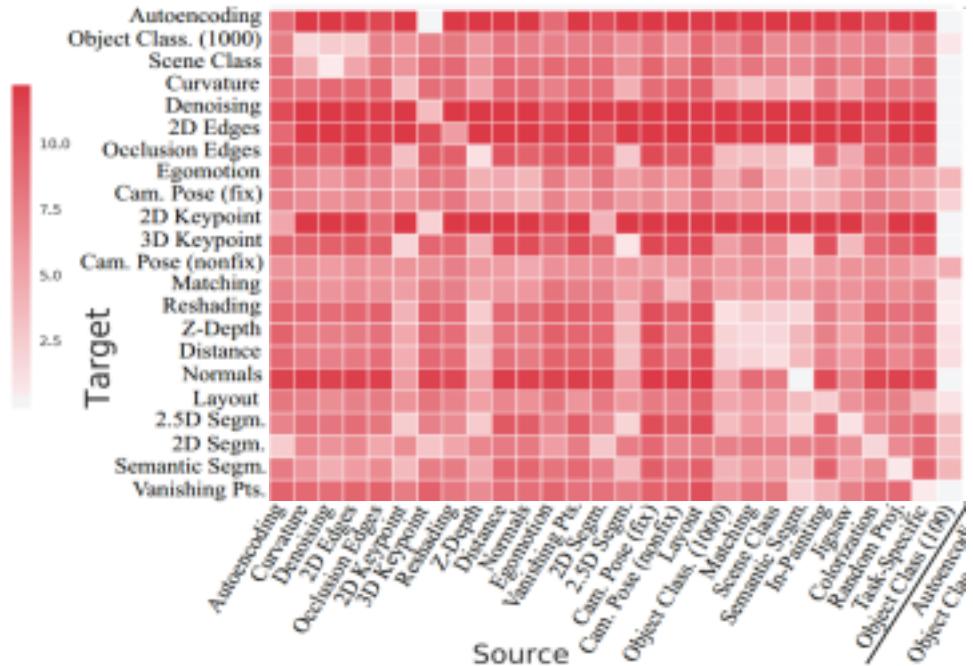
# Pairwise Task Correlation Matrix

- Wisdom-of-crowd
- A set of experts (30 in our case) asked to rate relation between task pairs  $\tau_i$  and  $\tau_j$ , across known and zero-shot tasks
- Note that our framework can admit any other source of task correlation beyond crowdsourcing



# Pairwise Task Correlation Matrix

Taskonomy CVPR 2018 (Best Paper): 26 vision tasks



Sampled set of tasks and not an exhaustive list

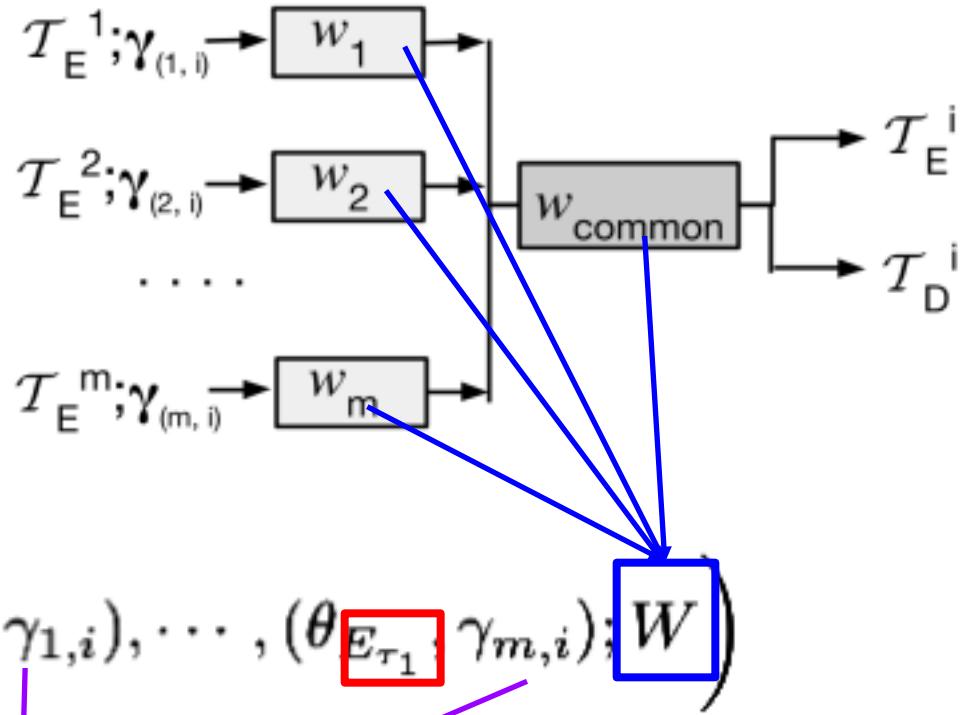
Zero-shot Task Transfer

# Methodology

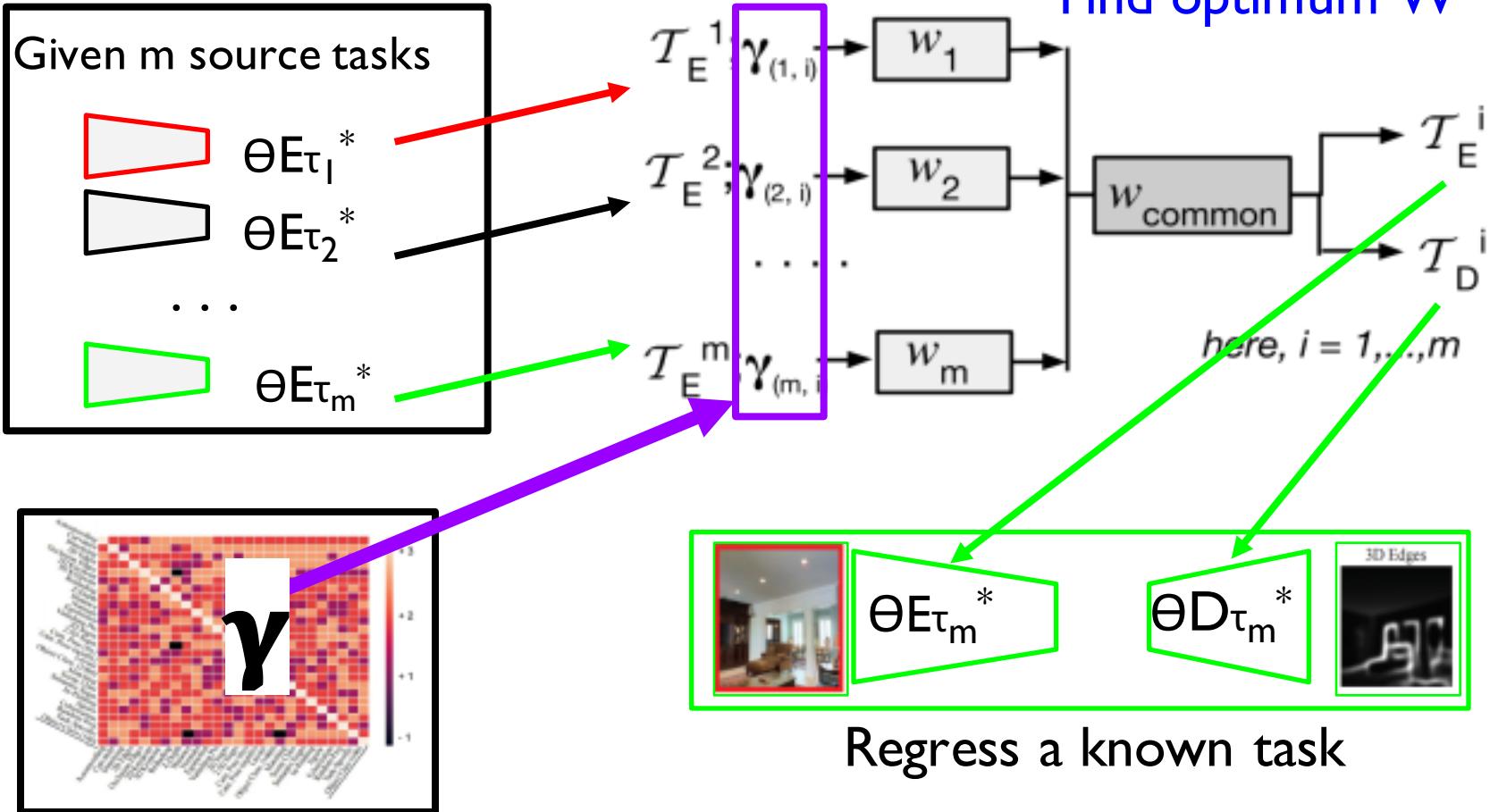
- We introduce TTNet, an architecture, to learn our meta-learner  $F(\cdot)$  parameterized by the weight  $\mathbf{W}$

$$\min_{\mathbf{W}} \sum_{i=1}^m \left| \left| \mathcal{F} \left( (\theta_{E_{\tau_1}}, \gamma_{1,i}), \dots, (\theta_{E_{\tau_1}}, \gamma_{m,i}); \mathbf{W} \right) - (\theta_{E_{\tau_i}}^*, \theta_{D_{\tau_i}}^*) \right| \right|^2$$

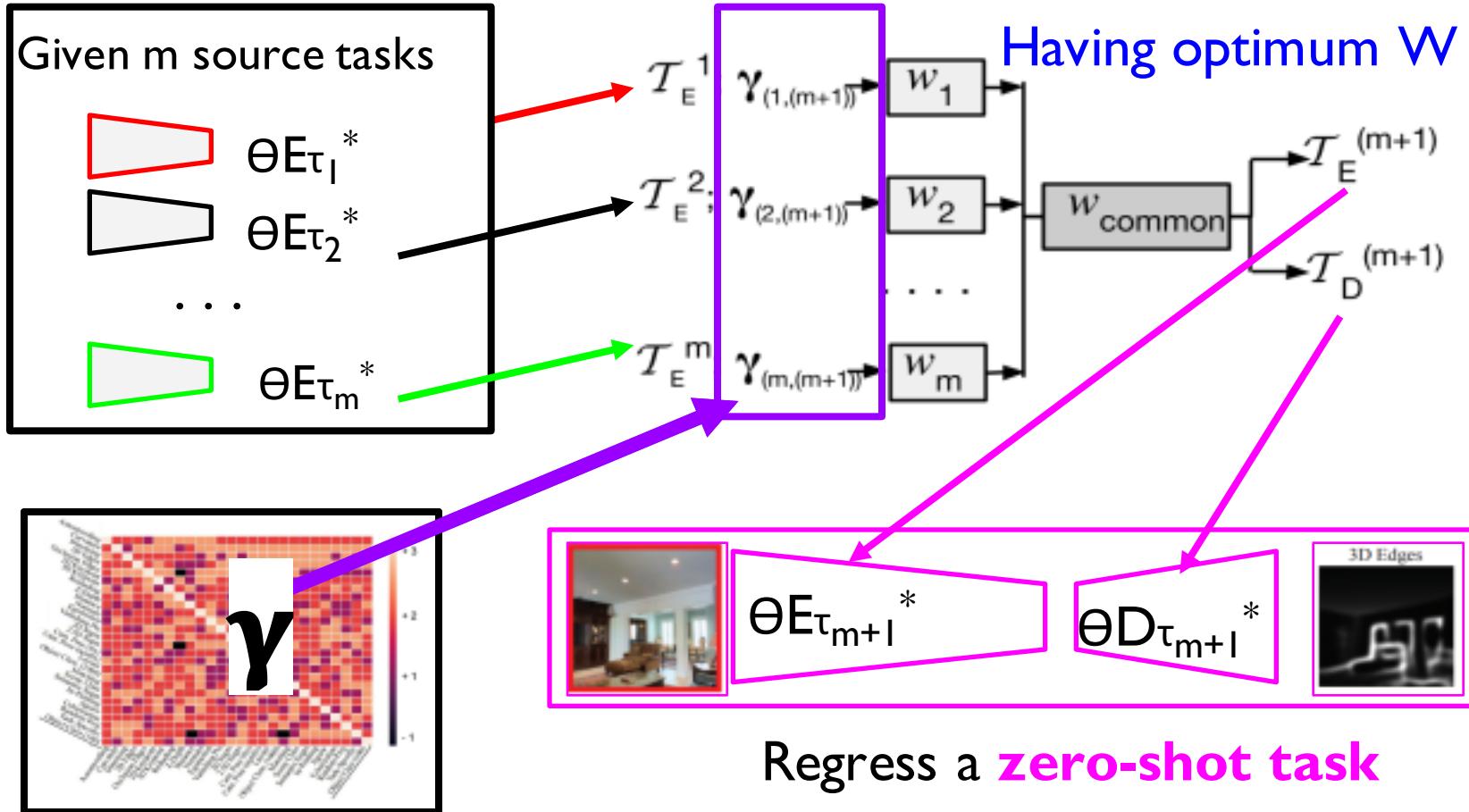
Pairwise task correlations



# TTNet: Training

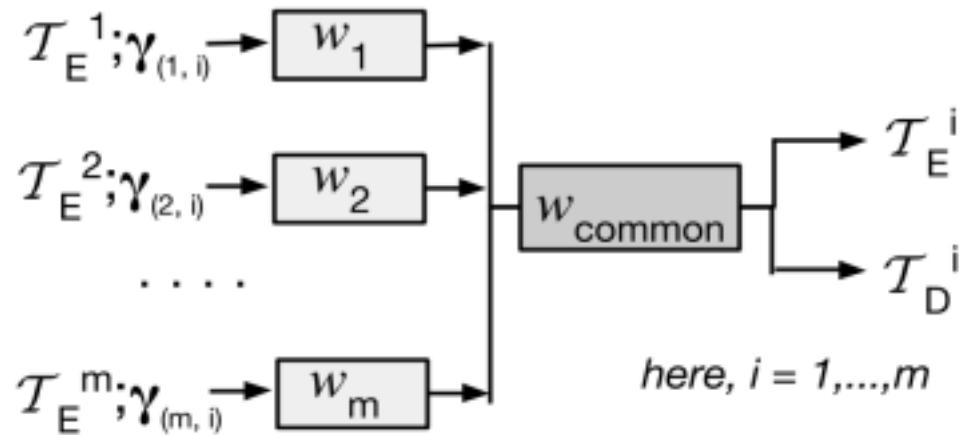


# TTNet: Zero-shot Parameter Regression



# Loss Function

- We add a **data model consistency loss**



$$\min_W \sum_{i=1}^m \left| \left| \mathcal{F} \left( (\theta_{E_{\tau_1}}, \gamma_{1,i}), \dots, (\theta_{E_{\tau_1}}, \gamma_{m,i}); W \right) \right. \right.$$

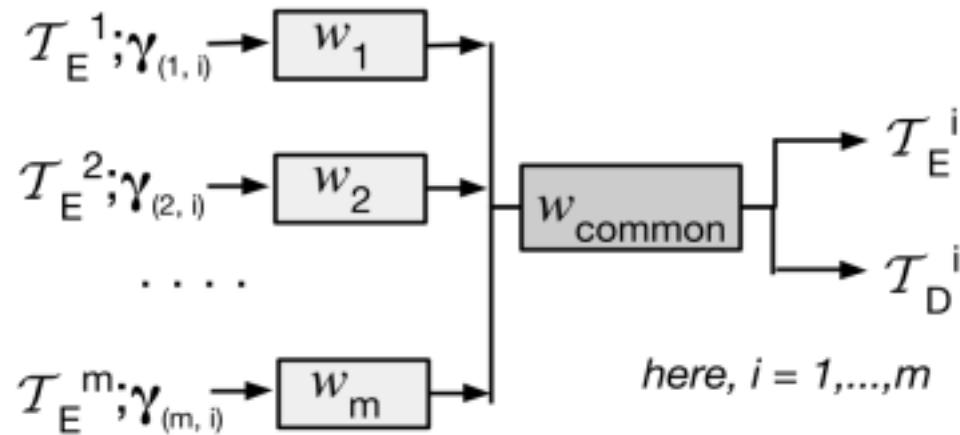
$$- (\theta_{E_{\tau_i}}^*, \theta_{D_{\tau_i}}^*) \|^2$$

$$+ \lambda \sum_{\substack{x \in X_{\tau_i} \\ y \in \mathbf{y}_{\tau_i}}} \mathcal{L} \left( \mathcal{D}_{\tilde{\theta}_{D_{\tau_i}}}(\mathcal{E}_{\tilde{\theta}_{E_{\tau_i}}}(x)), y \right)$$

Where  $\mathcal{L}$  is a loss appropriately chosen for a particular task

# Loss Function

- We add a **data model consistency loss**



$$\min_W \sum_{i=1}^m \left| \left| \mathcal{F} \left( (\theta_{E_{\tau_1}}, \gamma_{1,i}), \dots, (\theta_{E_{\tau_1}}, \gamma_{m,i}); W \right) \right. \right.$$

$$- (\theta_{E_{\tau_i}}^*, \theta_{D_{\tau_i}}^*) \|^2$$

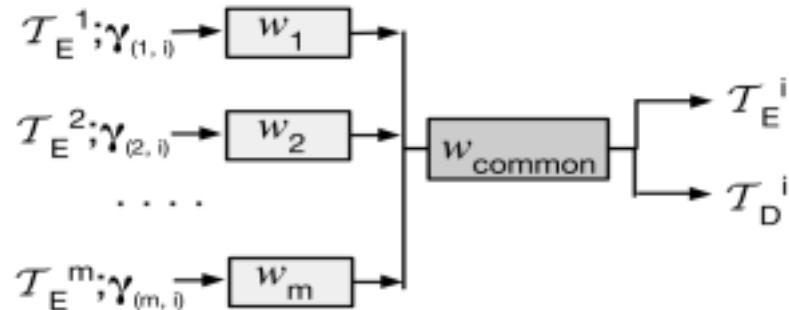
$$+ \lambda \sum_{\substack{x \in X_{\tau_i} \\ y \in \mathbf{y}_{\tau_i}}} \mathcal{L} \left( \mathcal{D}_{\tilde{\theta}_{D_{\tau_i}}}(\mathcal{E}_{\tilde{\theta}_{E_{\tau_i}}}(x)), y \right)$$

Where  $\mathcal{L}$  is a loss appropriately chosen for a particular task

# Methodology: Other Details

## Input:

- $P$  subsets of labeled data obtained from each task
- Model parameters learned for each of  $p$  subsets
- For  $m$  tasks, we get a  $p \times m$  data samples as input



## Training:

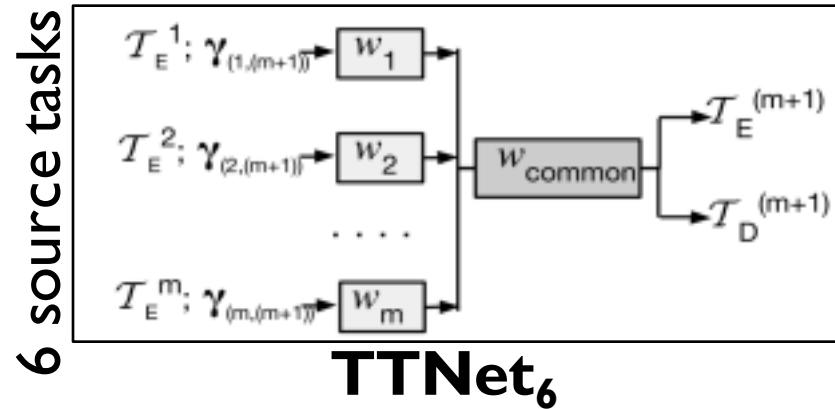
- *Self mode*: Train each  $W_i$  individually (and  $W_{common}$ )
- *Transfer mode*: Train other  $W_i$ s in the network

## Inference:

- Runs in transfer mode for target zero-shot task

# Experiments

- **Taskonomy dataset\*\***
  - \* 150k RGB data of indoor scenes
- We studied three different TTNets with 6, 10 and 20 source tasks
- Also, considered TTNet<sub>LS</sub> (fine-tuned on small amount of target task data)

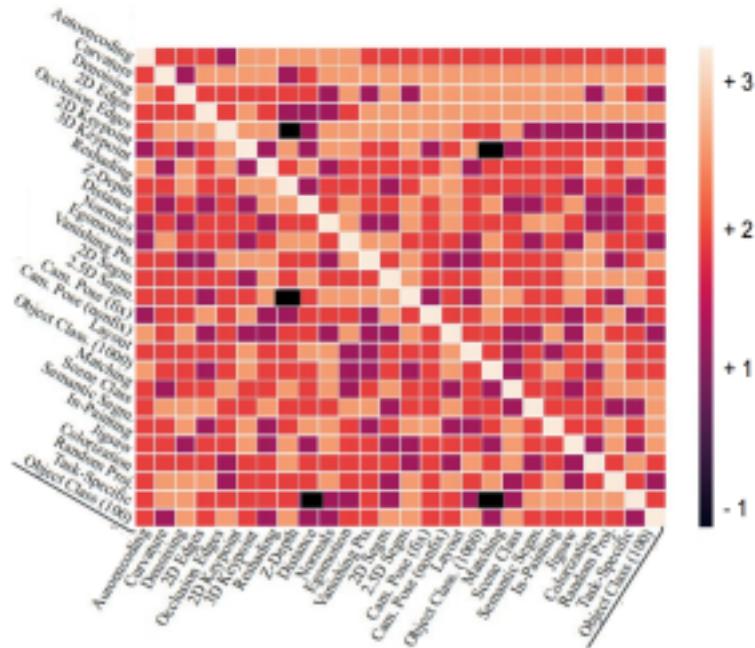


\*\* Taskonomy, Amir R Zamir et al., CVPR 2018



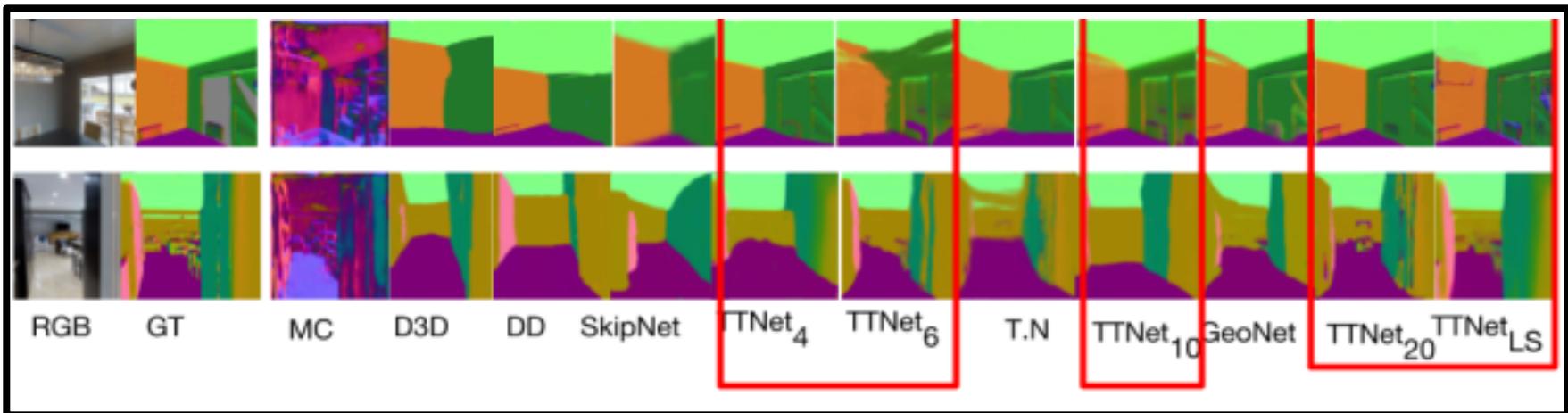
# Task Correlation

- Scale of -1 to +2 (negative correlation to strong positive correlation)
- Crowd votes aggregated from 30 users using the widely used Dawid-Skene algorithm (EM-based)



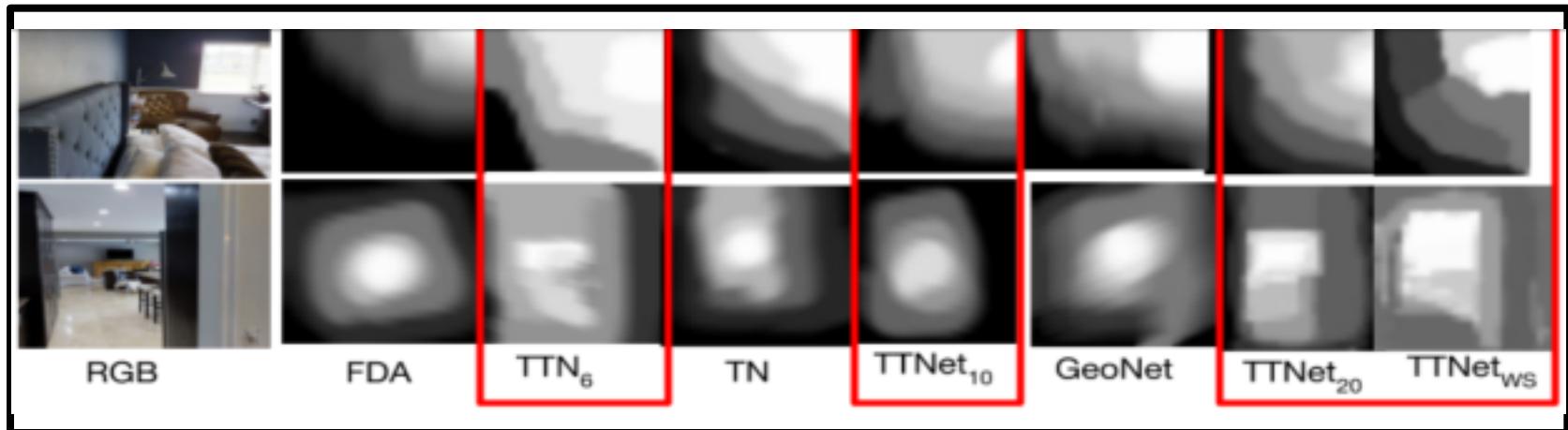
# Results

<b>Source (Known) Tasks</b>	Autoencoding, Scene Class, 3D key point, Reshading, Vanishing Point, Colorization
<b>Zero-shot Task</b>	Surface Normal Estimation

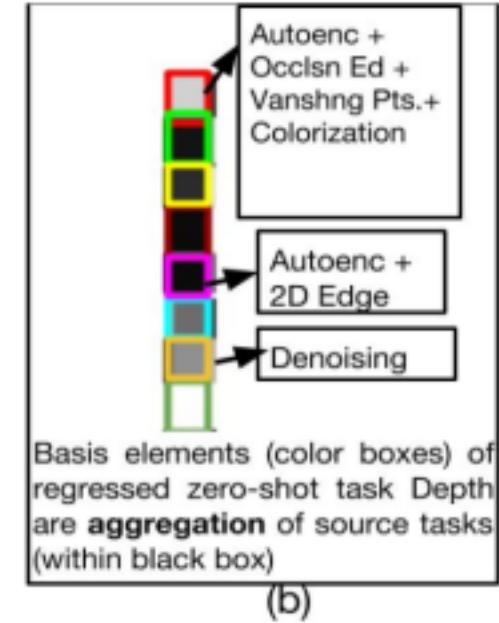
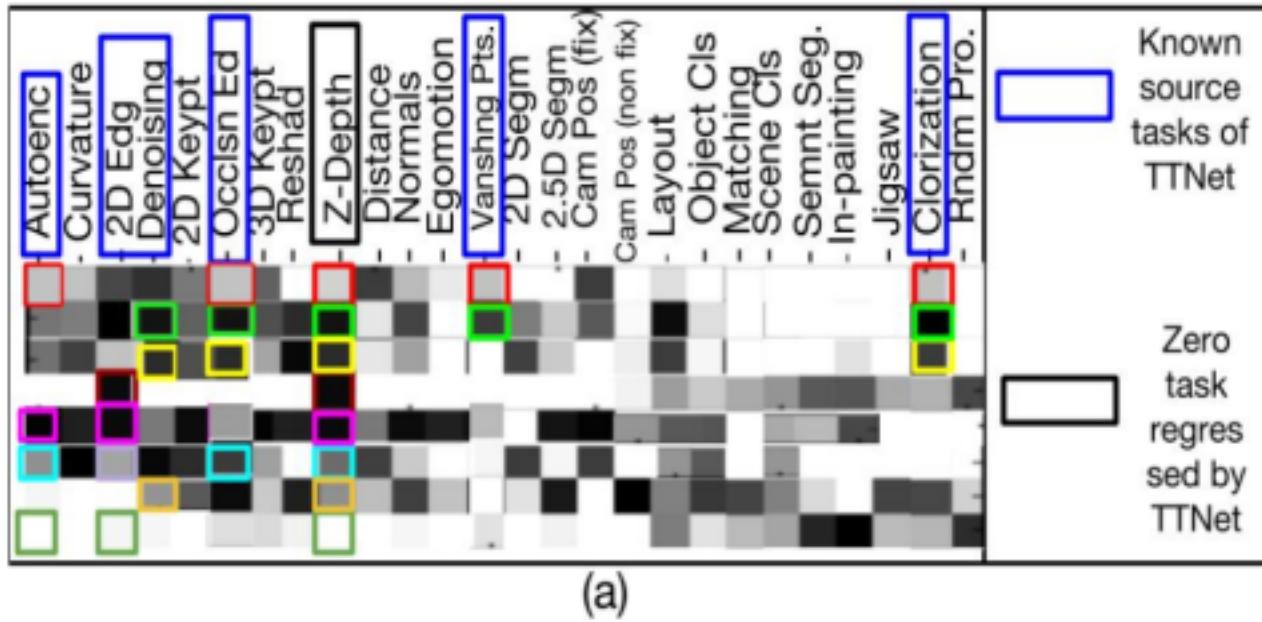


# Results

<b>Source (Known) Tasks</b>	Autoencoding, Scene Class, 3D key point, Reshading, Vanishing Point, Colorization
<b>Zero-shot Task</b>	Z-Depth Estimation

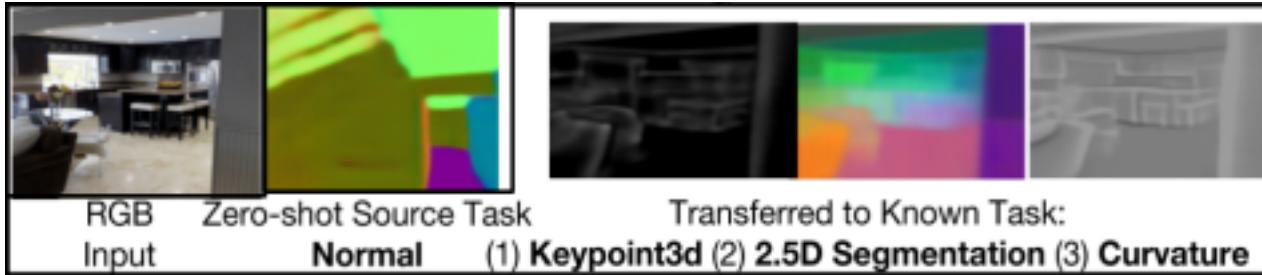


# Why better than Supervised Learning?

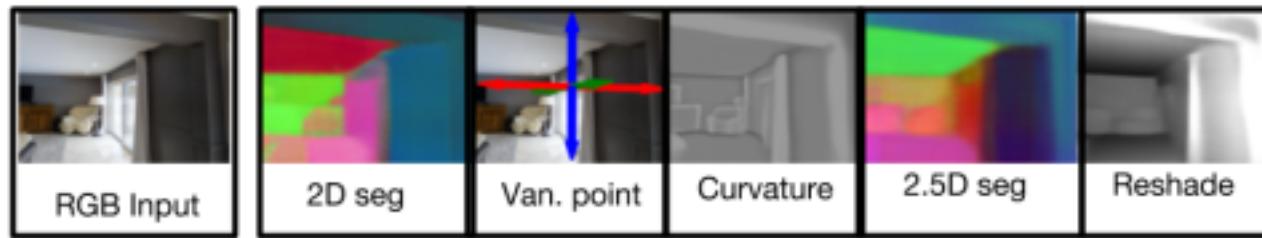


# More Results

- Zero-shot to Known Task Transfer



- Different Choice of Zero-Shot Tasks



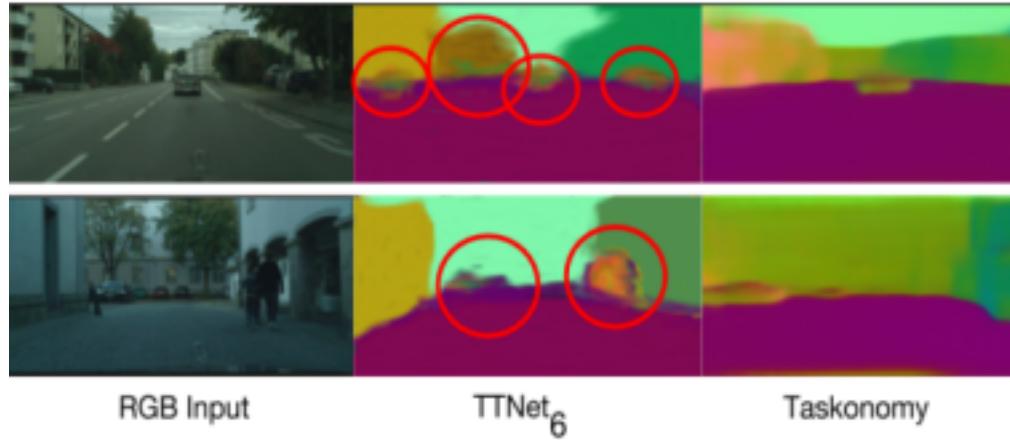
# How many source tasks do we need?

	Win Rate (Camera Pt,fixed) (%)	Win Rate (Depth) (%)	Win Rate (Room Layout) (%)	Win Rate (Normal) (%)	Colorization	Random Projection	Jig-Saw Puzzle	Vanishing Point	Semantic Segmentation	2D Segmentation	2.5D Segmentation	Room Layout	Normals	Distance	Z-Depth	Reshading	Matching	Cam Pose (non-fixed)	3D Key Point	2D Key Point	Cam Pose (fixed)	Ego motion	Occlusion Edges	2D Edges	Denoising	Curvature	Scene Class	Object Class	Autoencoding				
4	✓	x	x	x	✓	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	79%	62%	71%	71%			
	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	71%	58%	61%	59%			
	✓	x	x	x	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	75%	79%	79%	52%			
6	✓	x	x	x	✓	✓	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	88%	85%	87%	89%			
	✓	x	✓	x	x	x	x	x	x	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	87%	86%	86%	89%			
	✓	✓	✓	✓	✓	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	85%	88%	86%	82%			
10	✓	✓	✓	✓	✓	✓	x	x	x	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	85%	84%	87%	85%			
	✓	✓	✓	✓	✓	✓	x	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	87%	88%	91%	92%			
	✓	✓	✓	✓	✓	✓	x	x	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	88%	83%	81%	89%			
15	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	88%	85%	91%	93%				
	✓	x	x	✓	✓	✓	x	✓	x	x	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			
18	✓	✓	✓	✓	✓	✓	x	✓	x	✓	✓	✓	✓	✓	✓	✓	x	✓	x	✓	✓	✓	✓	✓	✓	✓	✓	✓	93%	91%	97%	91%	
	✓	x	✓	✓	✓	✓	x	✓	x	✓	✓	✓	✓	✓	✓	✓	x	✓	x	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	95%	91%	93%
20	✓	x	✓	✓	✓	✓	✓	✓	✓	✓	x	✓	✓	✓	✓	✓	✓	x	✓	x	✓	✓	✓	✓	✓	✓	✓	✓	✓	94%	91%	93%	89%



# More Results

- Performance on Other Datasets (Cityscapes)



# Questions?



[cs15resch11001@iith.ac.in](mailto:cs15resch11001@iith.ac.in)  
[vineethnb@iith.ac.in](mailto:vineethnb@iith.ac.in)

Department of Computer Science  
and Engineering, IIT-Hyderabad  
<http://www.iith.ac.in/~vineethnb>

