# Taskonomy: Disentangling Task Transfer Learning

Written by Amir. R Zamir, Alexander Sax, William Shen, Leonidas Guibas, Jitendra Malik and Silvio Savarese
Presented by Peyman Bateni, Lucas Porto, Tanzila Rahman, Vibudh Agrawal

*unless otherwise stated, all figures come from the paper

# Timeline for Today

— — —

- **Introduction**
- **Related Work**
- **Method**
  - **Stage I: Task Specific Modelling**
  - **Stage II: Transfer Modelling**
  - **Stage III: Ordinal Normalization using AHP**
  - **Stage IV: Computing the Global Taxonomy**
- **Experiments**
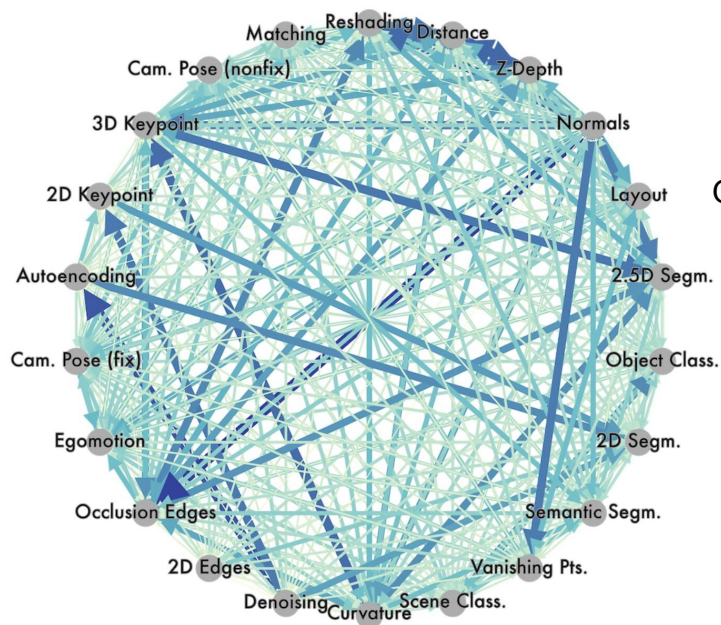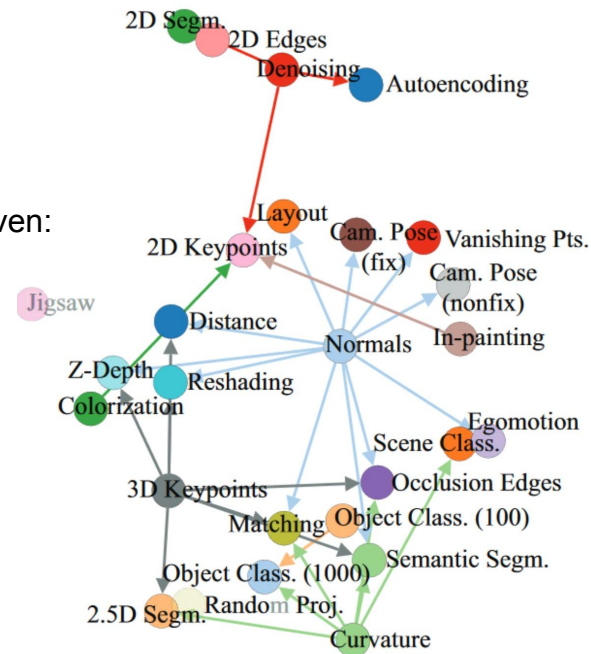- **Discussion**
- **Question**

# Introduction

— — —

- **Motivation:** transferability among visual tasks
  - Are visual tasks related? Should we learn a task from scratch every time?
  - Saving time and data: reducing supervision and computation
  - Efficient learning of comprehensive perception models

- **Taskonomy (task taxonomy)**: a structure that describes efficient task transfers
  - Directed hypergraph
  - Nodes: tasks
  - Edges: transfers

# Taskonomy: Task Taxonomy



Optimize for performance given:
- Supervision budget
- Task importance

# Related Work

— — —

- **Self-supervised learning:** learning from labels derived directly from the input data

- **Multi-task learning:** learning representations for multiple tasks

- **Domain adaptation:** robustness to different input domains

- **Meta-learning:** similar motivations; higher-level understanding of the learning process

# Method - Overview

— — —

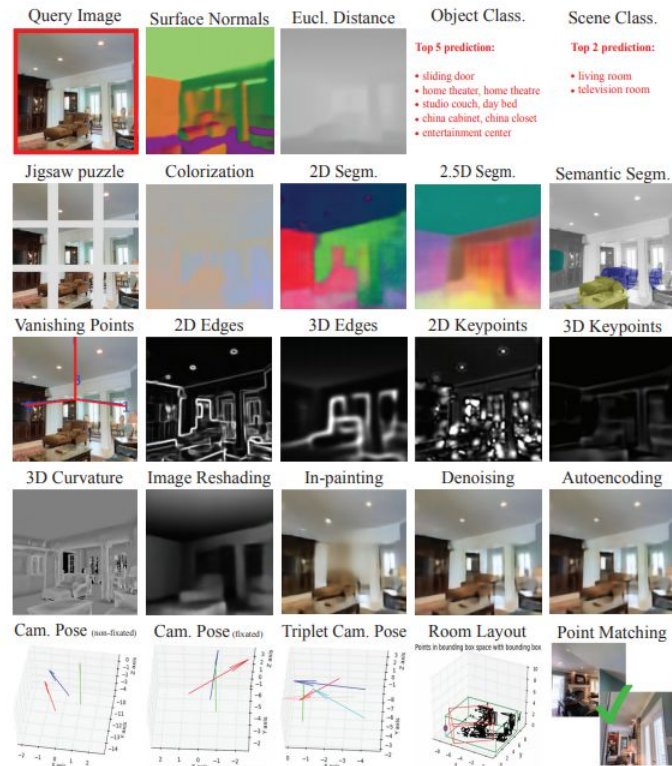**Task dictionary/bank (V = $\mathcal{T} \cup$ S)**

- 26  2D, 2.5D, 3D, and semantic tasks

**Dataset**

- 4 million images of indoor scenes from about 600 buildings.

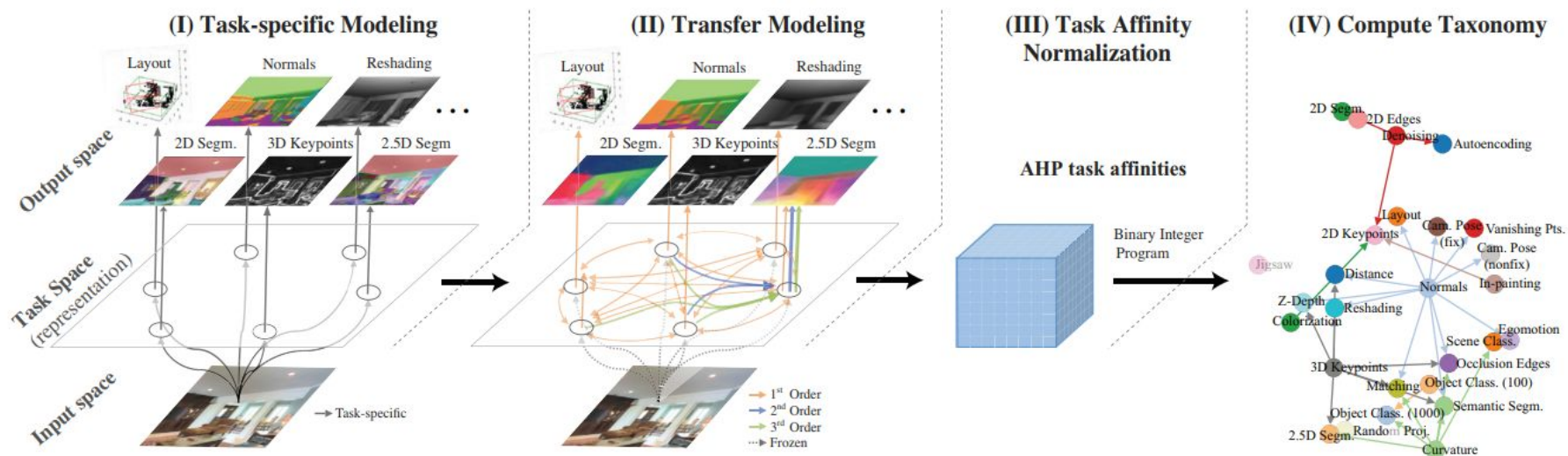- Each image has GT label for all the tasks.

**Task specific network**

- 26x Network



**Task Dictionary**

# Method - Overview



(I) Task-specific Modeling

(II) Transfer Modeling

(III) Task Affinity Normalization
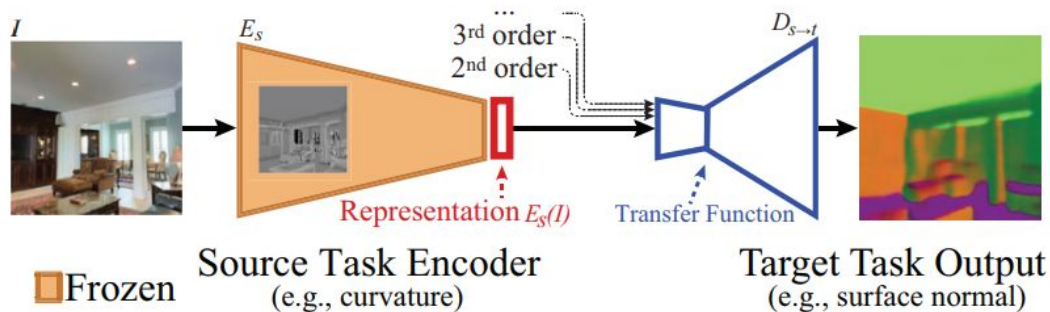
(IV) Compute Taxonomy

# Method - Stage I: Task-Specific Modelling

— — —

- Use an encoder-decoder architecture homogeneous across all tasks.

- Encoder: A fully convolutional ResNet-50 without pooling and extract powerful representations.

- Decoder: Different architecture depending on tasks, but smaller than encoder.

# Method - Stage II: Task-Specific Modelling

— — —

- Transfer network learns a small transfer/readout function ($D_{s \to t}$).

- $D_{s \to t}$ is parameterized by $\theta_{s \to t}$ minimizing the loss $L_t$:

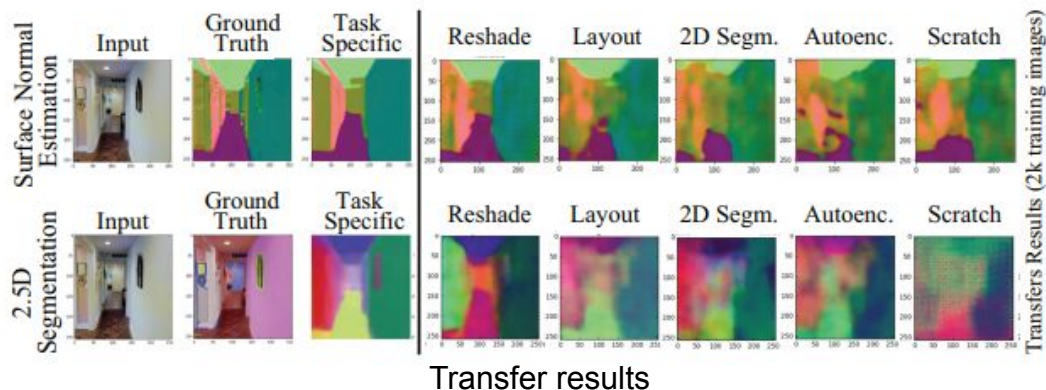$$D_{s \to t} := \arg\min_{\theta} \mathbb{E}_{I \in \mathcal{D}} \Big[ L_t \Big( D_\theta \big( E_s(I) \big), f_t(I) \Big) \Big]$$
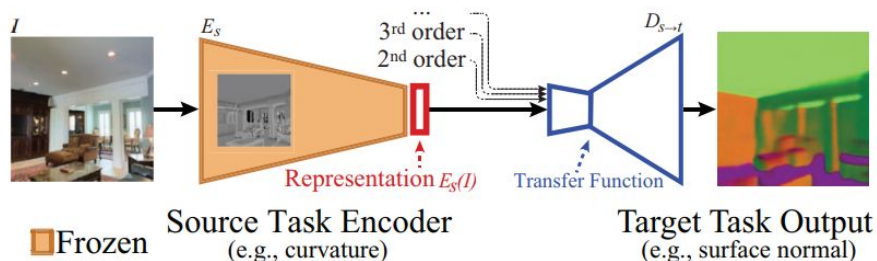
# Method - Stage II: Task-Specific Modelling

— — —

- Learn a readout function ($D_{s \rightarrow t}$).

- $D_{s \rightarrow t}$ is parameterized by $\theta_{s \rightarrow t}$ minimizing the loss $L_t$:

$$D_{s \rightarrow t} := \arg\min_{\theta} \mathbb{E}_{I \in \mathcal{D}} \left[ L_t \Big( D_\theta \big( E_s(I) \big), f_t(I) \Big) \right]$$



Transfer results

# Method - Stage II: Task-Specific Modelling

— — —

- Learn a readout function ($D_{s \to t}$).

- $D_{s \to t}$ is parameterized by $\theta_{s \to t}$ minimizing the loss $L_t$:

$$D_{s \to t} := \arg\min_{\theta} \mathbb{E}_{I \in \mathcal{D}} \Big[ L_t \Big( D_\theta \big( E_s(I) \big), f_t(I) \Big) \Big]$$

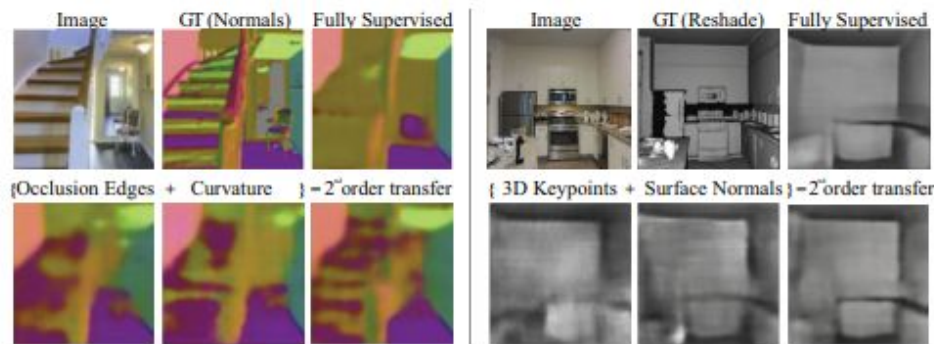- Include higher-order transfers to receive multiple representations in the input.

# Method - Stage II: Task-Specific Modelling

- Learn a readout function ($D_{s\rightarrow t}$).

- $D_{s\rightarrow t}$ is parameterized by $\theta_{s\rightarrow t}$ minimizing the loss $L_t$:

$$D_{s\rightarrow t} := \arg\min_{\theta} \mathbb{E}_{I\in\mathcal{D}}\left[L_t\left(D_{\theta}(E_s(I)), f_t(I)\right)\right]$$
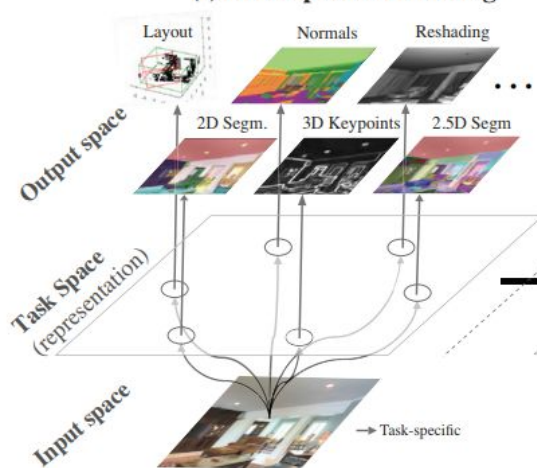
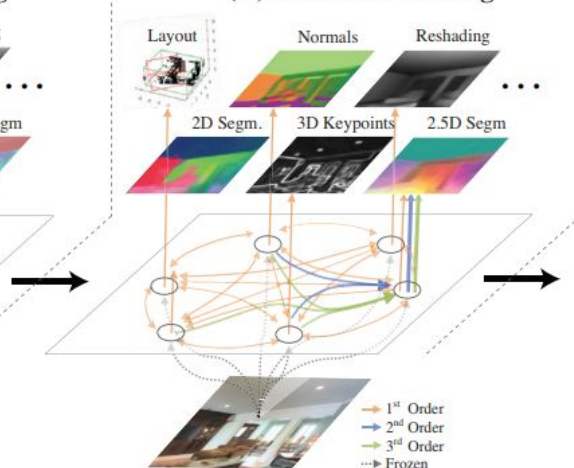- Include higher-order transfers to receive multiple representations in the input.



Higher order transfers

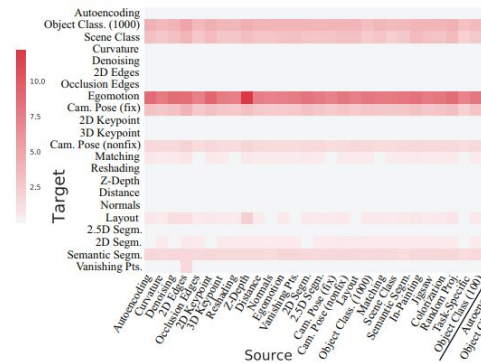# Method - Stage II: Task-Specific Modelling

# Method - Stage III: Ordinal Normalization using AHP

— — —

- A naive normalization can apply to linearly rescale each row of the matrix to the range [0, 1]. But this approach fails when the actual output quality increases at different speeds w.r.t. the loss.
- Apply ordinal normalization derived from Analytic Hierarchy Process (AHP).

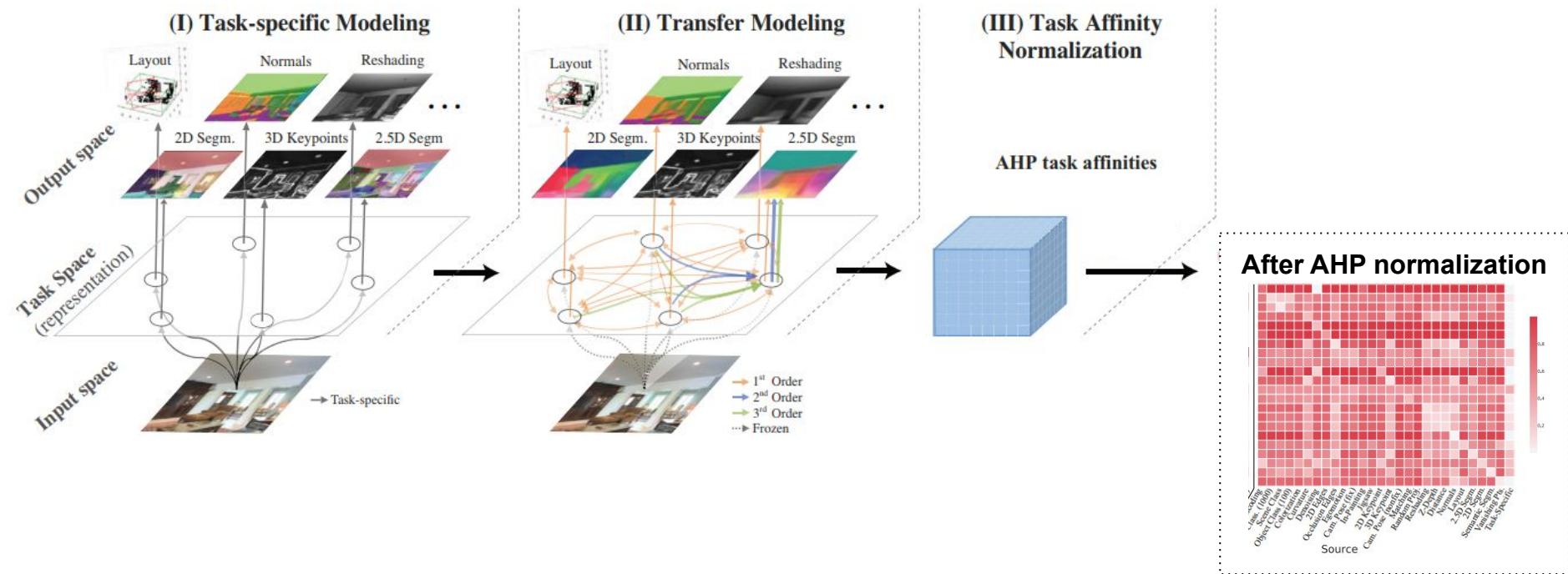# Method - Stage III: Ordinal Normalization using AHP

— — —

Ordinal Normalization using AHP :

- Construct a pairwise tournament matrix $W_t$ .
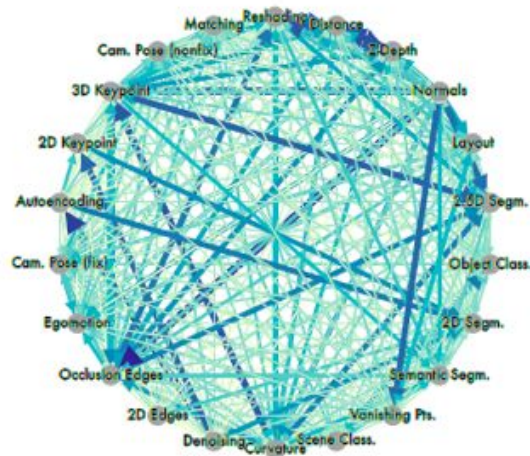
- Take win-lose ratio between transfer $s_i$ and $s_j$.

$$w'_{i,j} = \frac{\mathbb{E}_{I \in \mathcal{D}_{test}}[D_{s_i \to t}(I) > D_{s_j \to t}(I)]}{\mathbb{E}_{I \in \mathcal{D}_{test}}[D_{s_i \to t}(I) < D_{s_j \to t}(I)]}$$

- Take the 1st principle component (normalized to sum to 1) of the matrix.

- Create the final matrix by stacking the 1st principle components for all t.

# Method - Stage III: Ordinal Normalization using AHP

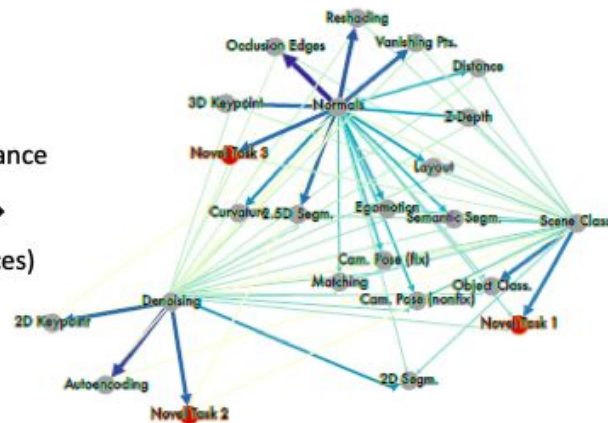# Step IV: Computing the Global Taxonomy (Boolean Integer Programing (BIP))



Maximize collective performance

Minimize supervision (sources)

Post Analytic Hierarchy Process

Post Boolean Integer Programing

# Problem formulation

- E (edges in a graph, transfers):

$$(\{s_1^i, \ldots, s_{m_i}^i\}, t^i)$$
$$\{s_1^i, \ldots, s_{m_i}^i\} \subset S$$
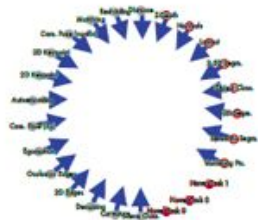$$t^i \subset T$$

- V (vertices in a graph, tasks):

$$S \cup T$$

- Operations:

$$(\{s_1^i, \ldots, s_{m_i}^i\}, t^i) \xrightarrow{\text{sources}} \{s_1^i, \ldots, s_{m_i}^i\}$$

$$(\{s_1^i, \ldots, s_{m_i}^i\}, t^i) \xrightarrow{\text{target}} t^i$$

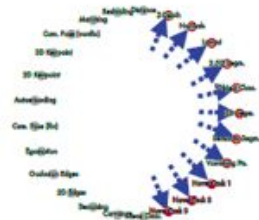$(|E| + |V|, 1)$                        $(|E| + |V|, 1)$

Maximize: $c^t z$

Subject to: $Ax \leq b \; and \; x \in \{0,1\}^{|E|+|V|}$

$(|E| + 2|V| + 1, |E| + |V|)$         $(|E| + 2|V| + 1, 1)$
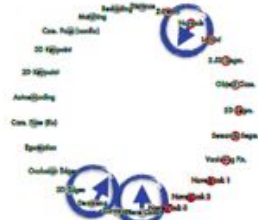
- Constrains:



**Constraint I:**
only transfer from sources.

**Constraint II:**
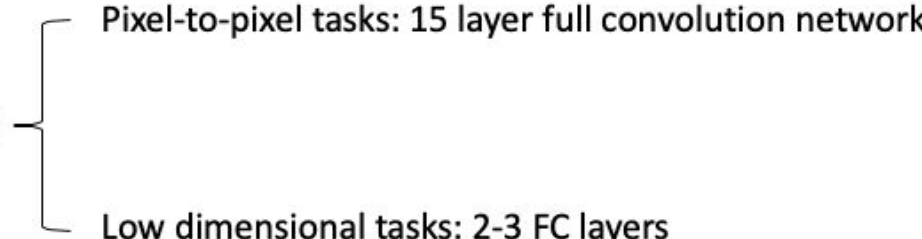all targets are transferred to.

**Constraint III:**
not exceed budget.

# Experiments

---

- Total tasks in dictionary: 26
- Source only tasks: 4 (colorization, jigsaw puzzle, in-painting, random projection)
- Total transfer functions: ~3000
- GPU hours: 47,886

# Network Architecture

− − −

- Encoder: Fully connected ResNet-50 without pooling

Pixel-to-pixel tasks: 15 layer full convolution network

- Decoders:

Low dimensional tasks: 2-3 FC layers

- Transfer: 2 convolution layers

- Same hyperparameters as well as architecture across different tasks.

# Data splits

———

- Dataset: 4 million images
- Training: 120k
- Validation: 16k
- Test: 17k
- Task specific network (Training dataset)
- Transfer network (validation dataset, ranging from 1k – 16K)

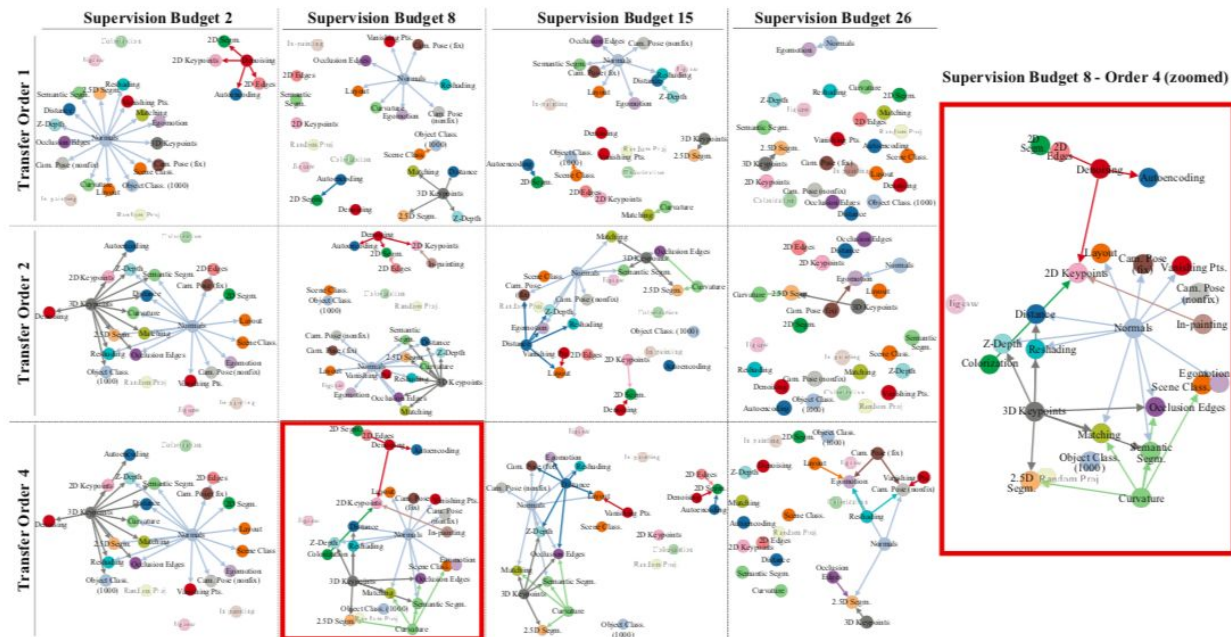# Experiments I - Evaluation of Computed Taxonomies



Figure 8: **Computed taxonomies** for solving 22 tasks given various supervision budgets (x-axes), and maximum allowed transfer orders (y-axes). One is magnified for better visibility. Nodes with incoming edges are target tasks, and the number of their incoming edges is the order of their chosen transfer function. Still transferring to some targets when tge budget is 26 (full budget) means certain transfers started performing better than their fully supervised task-specific counterpart. See the interactive solver website for color coding of the nodes by *Gain* and *Quality* metrics. Dimmed nodes are the source-only tasks, and thus, only participate in the taxonomy if found worthwhile by the BIP optimization to be one of the sources.

# Experiments I - Evaluation of Computed Taxonomies

- - - -

**Gain:** win rate (%) against a network trained from scratch using the same training data as transfer networks'. That is, the best that could be done if transfer learning was not utilized. This quantifies the *gained* value by transferring.

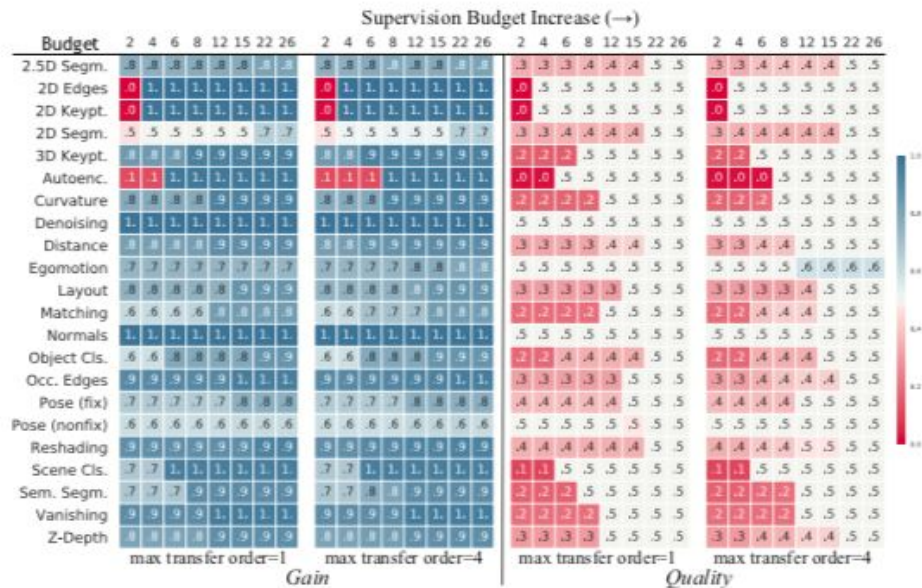**Quality:** win rate (%) against a fully supervised network trained with 120k images (gold standard).



Figure 9: **Evaluation of taxonomy computed for solving the full task dictionary.** Gain (left) and Quality (right) values for each task using the policy suggested by the computed taxonomy, as the supervision budget increases(→). Shown for transfer orders 1 and 4.

# Experiments II - Generalization to Novel Tasks

- - -

- **Carefully transferring policies depending on the target is superior to fixed transfers**
- Although taxonomy transfer policies lose to fully supervised networks, **in most cases results get close to the gold standard** with win rate at 40%



Figure 10: **Generalization to Novel Tasks.** Each row shows a novel test task. Left: Gain and Quality values using the devised "all-for-one" transfer policies for novel tasks for orders 1-4. Right: Win rates (%) of the transfer policy over various self-supervised methods, ImageNet features, and scratch are shown in the colored rows. Note the large margin of win by taxonomy. The uncolored rows show corresponding loss values.

# Experiments III - Significance Test of the Structure

‒ ‒ ‒

- **Outperforms all other randomly assigned connective networks, indicating the significant and existence of an underlying connective structure between the tasks**
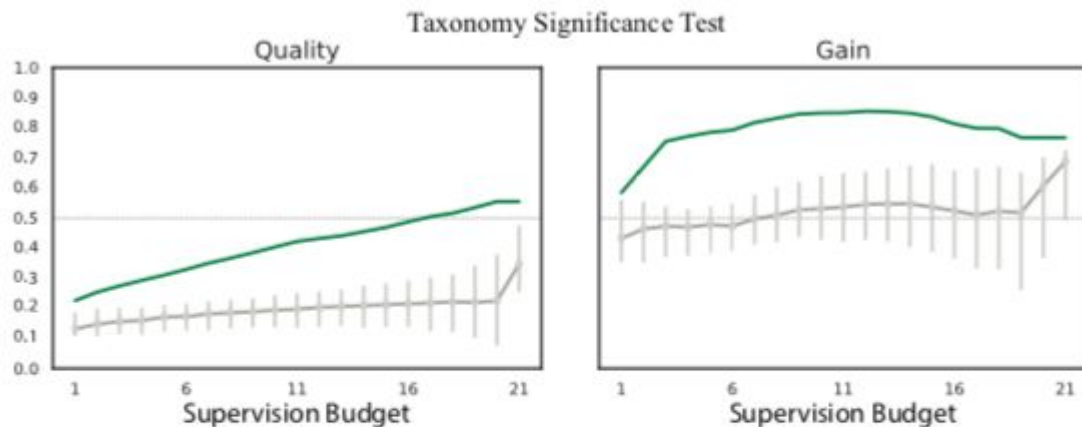


Figure 11: **Structure Significance.** Our taxonomy compared with random transfer policies (random feasible taxonomies that use the maximum allowable supervision budget). Y-axis shows *Quality* or *Gain*, and X-axis is the supervision budget. Green and gray represent our taxonomy and random connectivities, respectively. Error bars denote $5^{th}$–$95^{th}$ percentiles.

# Experiments IV - Evaluation on MIT Places & ImageNet

─ ─ ─

- **Spearman's rho between taxonomy ranking and the Top-1 ranking is 0.857 on Places and 0.823 on ImageNet showing a notable correlation**
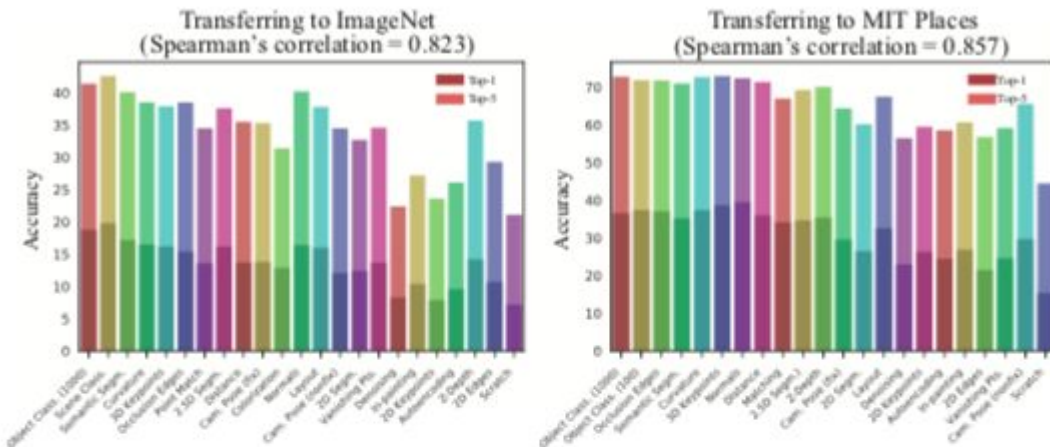


Figure 12: **Evaluating the discovered structure on other datasets: ImageNet [78] (left) for object classification and MIT Places [104] (right) for scene classification.** Y-axis shows accuracy on the external benchmark while bars on x-axis are ordered by taxonomy's predicted performance based on our dataset. A monotonically decreasing plot corresponds to preserving identical orders and perfect generalization.

# Experiments V - Universality of the Structure

— — —

**System choices**

      I.    **Architecture of Task-Specific Networks**
     II.    **Architecture of Transfer Function Networks**
    III.    **Amount of Data Available for Training Networks**
    IV.    **Datasets**
     V.    **Data Splits**
    VI.    **Choice of Dictionary**

**Remarkably stable leading to almost no change in the output taxonomy**

# Experiments VI - Task Similarity Tree



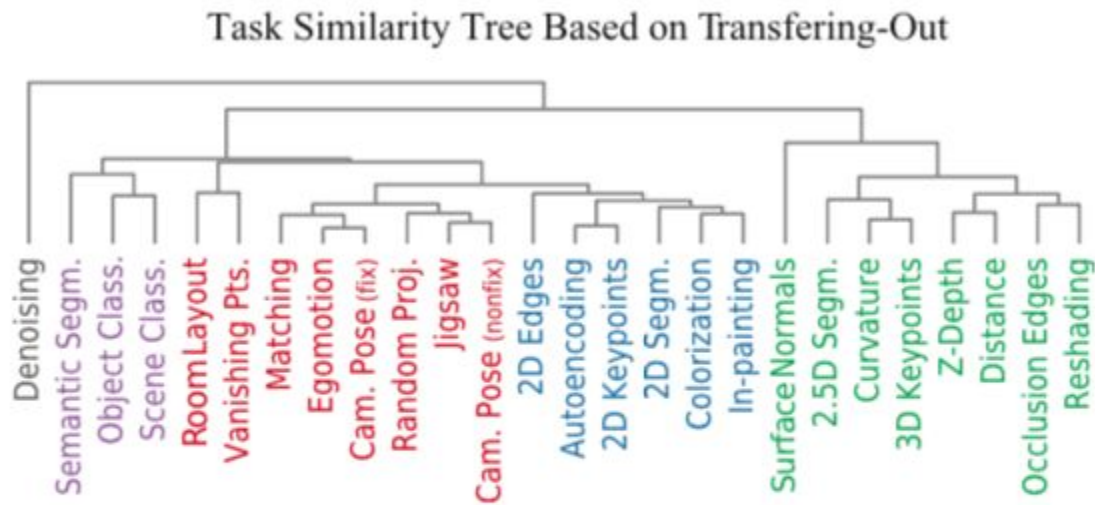Task Similarity Tree Based on Transfering-Out

Figure 13: **Task Similarity Tree.** Agglomerative clustering of tasks based on their transferring-out patterns (i.e. using columns of normalized affinity matrix as task features). 3D, 2D, low dimensional geometric, and semantic tasks clustered together using a fully computational approach.

# Limitations and Discussion

— — —

- **Model and data specific**
  - Taxonomy is computed w.r.t. a particular model and dataset

- **Compositionality**
  - Tasks as composition of subtasks

- **Non-visual tasks**
  - Does this method apply to non-visual tasks?

- **Lifelong/Continuous learning**
  - Expanding the taxonomy after computing it

# Thank You

— — —

- **We're happy to answer any questions :)**