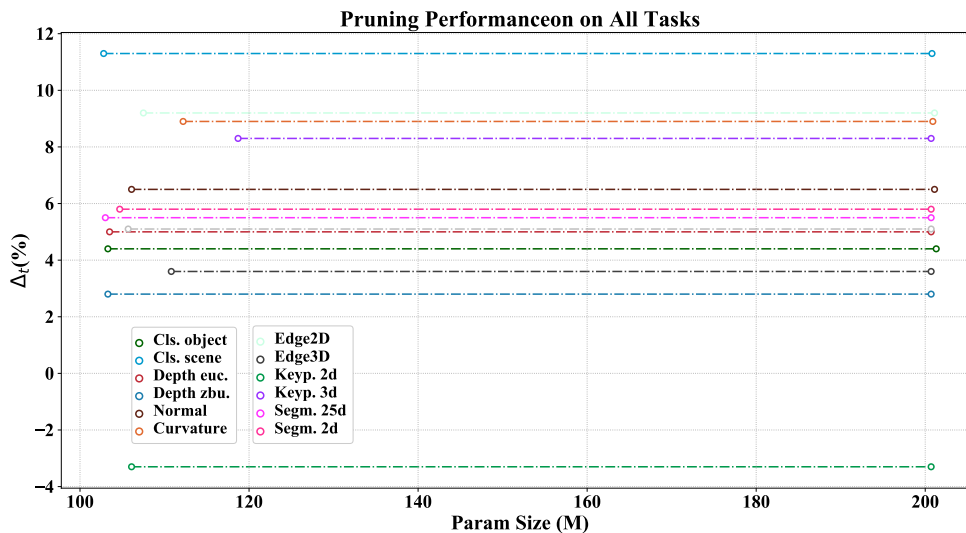# Supplementary Material for
# Mod-Squad: Designing Mixtures of Experts As Modular Multi-Task Learners

Table A1. Metric for more tasks on the taskonomy dataset. The experiment section in the paper demonstrates what metric we use for each task.

| Model | Depth Zbu. | | | | | | Edge2D | Keyp.3D | Segm. 2.5D |
|---|---|---|---|---|---|---|---|---|---|
| | $RMSE \downarrow$ | Error $\downarrow$ | | $\delta$, within $\uparrow$ | | | L1 dis. $\downarrow$ | L1 dis. $\downarrow$ | L1 dis. $\downarrow$ |
| | | Abs. | Rel. | 1.25 | $1.25^2$ | $1.25^3$ | | | |
| STL | 6.90 | 0.086 | 1.71 | 92.6 | 97.0 | 98.7 | 0.00500 | 0.072 | 0.155 |
| MTL | 6.75 | 0.083 | 1.30 | 93.1 | 96.8 | 98.9 | 0.00497 | 0.071 | 0.152 |
| $M^3ViT$ | 6.73 | 0.081 | 1.30 | **93.4** | 97.3 | 98.9 | 0.00490 | 0.070 | 0.150 |
| Mod-Squad | **6.70** | **0.080** | **1.28** | 93.3 | **97.5** | **99.1** | **0.00482** | **0.068** | **0.147** |

Figure A1. **Visualization of the pruning results on each task.** We set $\theta = 0.1\%$ and do the pruning. Every tasks keep the same performance while reducing the parameters size from over 200M to lower than 120M.



## A1. Metric for more task on the Taskonomy.

We provide more results with the task metric in Table. A1. Other tasks are shown in the paper already. Mod-Squad consistently have the best performance on these metrics.

## A2. Pruning on different tasks.

We show the pruning results for all tasks on the Taskonomy in Figure. A2. We set $\theta = 0.1\%$ and removing all experts have an activation frequency lower than $\theta$. Every tasks keep the same performance while reducing the parameters size from over 200M to lower than 120M. This demonstrate the effectiveness of Mod-Squad on all tasks.

Table A2. **Ablation study of Top-$K$ on MoE attention layer and MoE MLP layer.**

|  | FLOPs(G) | Params(M) | Hidden Dim | $\Delta_t$ |
|---|---|---|---|---|
| K=2 | 5.2 | 75 | 192 | -0.4 |
| K=4 | 5.2 | 75 | 96 | 1.3 |
| K=6 | 5.2 | 75 | 64 | **3.5** |
| K=9 | 5.2 | 75 | 42 | 3.1 |
| K=12 | 5.2 | 75 | 32 | 2.2 |

(a) **Ablation on Top-$K$ in the MoE attention layer.**

|  | FLOPs(G) | Params(M) | Hidden Dim | $\Delta_t$ |
|---|---|---|---|---|
| K=1 | 5.2 | 75 | 1536 | 3.3 |
| K=2 | 5.2 | 75 | 768 | 3.3 |
| K=4 | 5.2 | 75 | 384 | **3.5** |
| K=6 | 5.2 | 75 | 256 | 3.2 |
| K=8 | 5.2 | 75 | 192 | 3.1 |

(b) **Ablation on Top-$K$ in the MoE MLP layer.**

Table A3. **Ablation study of expert number E on MoE attention layer and MoE MLP layer.**

|  | Params(M) | $\Delta_t$ |
|---|---|---|
| E=6 | 66 | 2.0 |
| E=9 | 69 | 2.7 |
| E=12 | 72 | 3.2 |
| E=15 | 75 | **3.5** |
| E=18 | 78 | **3.5** |

(a) **Ablation on $E$ in the MoE attention layer.**

|  | Params(M) | $\Delta_t$ |
|---|---|---|
| E=4 | 33 | 2.2 |
| E=8 | 47 | 2.7 |
| E=12 | 61 | 3.3 |
| E=16 | 75 | **3.5** |
| E=20 | 90 | **3.5** |

(b) **Ablation on E in the MoE MLP layer.**

### A3. Ablation on Top-$K$.

As showin in Table. A2, we explore the effect on Top-$K$ for both MoE attention layer and MoE MLP layer. The experiment setting is the same as in § A4. To control the FLOPs to be the same for different Top-$K$, the hidden dimension of attention experts and the mlp experts is divided by $K$. All experiments have the same parameter size and the same FLOPs.

### A4. Ablation on number of experts.

As showin in Table. A3, we explore the effect on number of experts $E$ for both MoE attention layer and MoE MLP layer. For quick experiments, we use ViT-small as the backbone network. In default, we add MoE attention layer with 15 experts and Top-$K$ as 6 as well as MoE MLP layer with 16 experts and toop-k as 4. The MoE modules are added at every layer. We use $\Delta_t$ to evaluate the multi-task performance and report the parameters size. The $\Delta_t$ compare various version of models to the vanilla single task learning baseline with ViT-small. The experiments show that increasing experts number bring extra performance but the gain gradually diminish as the $E$ goes up. This conclusion hold true for both MoE attention layer and MoE MLP layer.

### A5. More visualization on Mod-Squad.

We visualize Mod-Squad based on ViT-small as shown in Figure. A2 and Figure. A3. To better understand the relation between experts and task in all layers, we insert MoE attention layer and MoE MLP layer on every transformer block. In Figure. A2, the activation frequencies of MoE attention modules are shown in all transformer blocks with 15 experts and Top-$K$ as 6. In Figure. A3, the activation frequencies of MoE MLP modules are shown in all transformer blocks with 16 experts and Top-$K$ as 4. Both visualization demonstrates the sparsity of Mod-Squad in all layers for all tasks.

### A6. Task relation from different layers of Mod-Squad.

In the paper, we define the similarity between tasks as the mean of the percentage of experts that they are sharing given the same input. We put all experts into the calculation of task similarity. However, we can also calculate the task similarity for each layer, by only put the experts from that specific layer into the calculation. As shown in Figure. A4, we visualize the task similarity for the first two layers, the middle two layers, and the last two layers of Mod-Squad. Mod-Squad is trained on the Taskonomy with ViT-base as the backbone and follow the same setting in our Taskonomy experiments in the paper. We notice that (1) all 3d tasks are very similar in the first layer; (2) all tasks are more similar to classification tasks in the middle two layers; (3) there is less similarity between tasks in the last two layers.

Figure A2. **Visualization of the frequency that experts being selected for each task in the MoE attention layer.** We visualize Mod-Squad based on ViT-small. The activation frequencies of MoE attention modules are shown in all transformer blocks with 15 experts and Top-$K$ as 6. The y-axis represents the tasks and the x-axis represents the experts. It demonstrates the sparsity of Mod-Squad in all layers for all tasks.
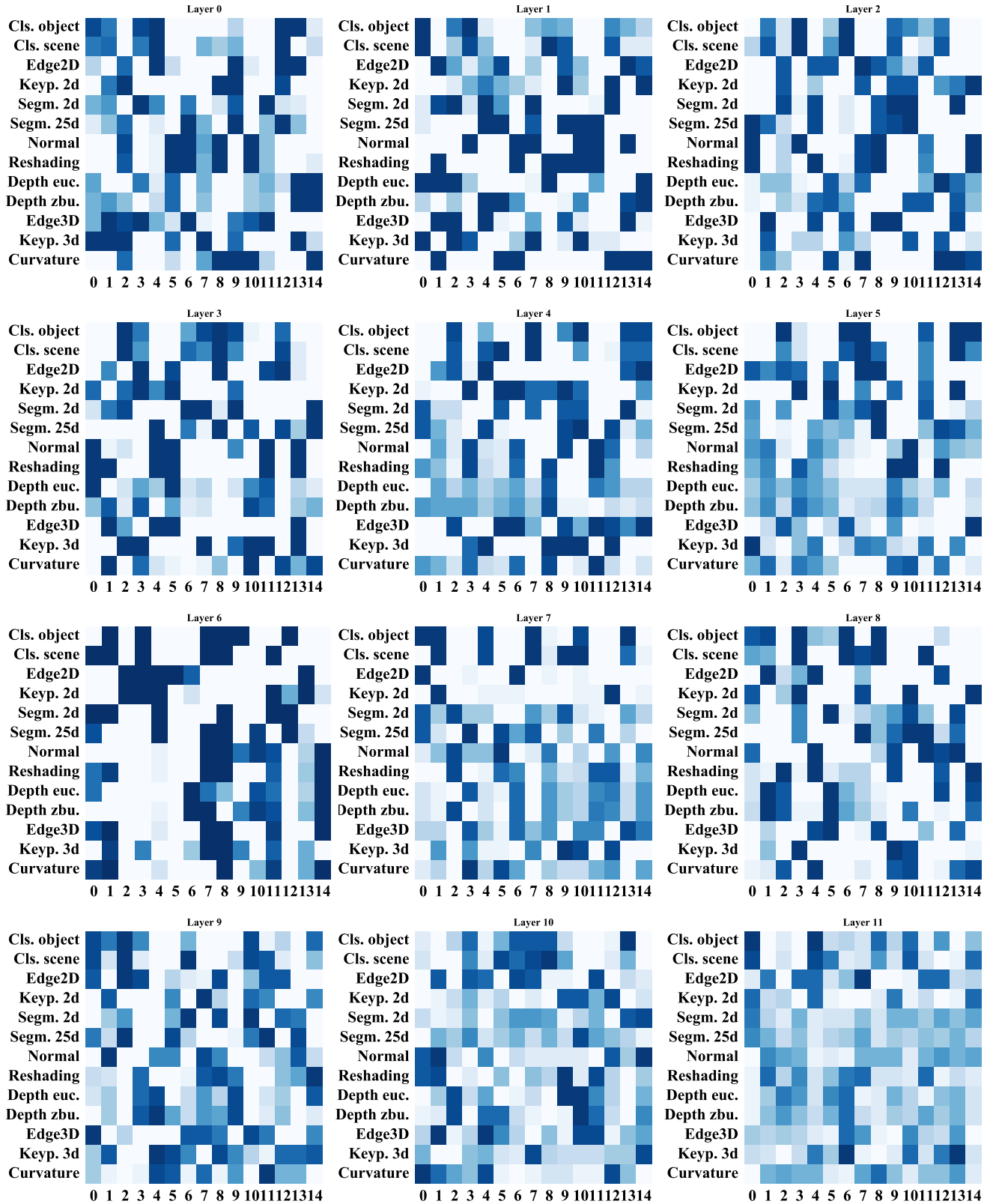
Figure A3. **Visualization of the frequency that experts being selected for each task in the MoE MLP layer.** We visualize Mod-Squad based on ViT-small. The activation frequencies of MoE MLP modules are shown in all transformer blocks with 16 experts and Top-$K$ as 4. The y-axis represents the tasks and the x-axis represents the experts. It demonstrates the sparsity of Mod-Squad in all layers for all tasks.

Figure A4. **Visualization of task similarity from the first, the middle, and the last layers.** For each layer $L_i$, we evaluate the similarity between a task pair as the mean of the percentage of experts in $L_i$ that task pairs are sharing with the same input.