

# Supplementary Material - CLR: Channel-wise Lightweight Reprogramming for Continual Learning

Yunhao Ge<sup>1†</sup>, Yuecheng Li<sup>1\*</sup>, Shuo Ni<sup>1\*</sup>, Jiaping Zhao<sup>2</sup> Ming-Hsuan Yang<sup>2</sup>, Laurent Itti<sup>1†</sup>

<sup>1</sup>University of Southern California <sup>2</sup>Google Research

\*Equal contribution as second author, †correspondence to {yunhaoge, itti}@usc.edu

## 1. Details of our 53-dataset for continual learning and method performance

Fig. 1 shows a summary of the 53 datasets we used as the continual learning benchmark in our main paper. The figure also shows the detailed per-task accuracy of our methods and baselines after learning all 53 tasks in the task incremental continual learning setting.

## 2. Channel-wise linear reprogramming ability

To further understand the performance of channel-wise lightweight reprogramming achieved by channel-wise linear transformation, we conduct qualitative experiments to explore the ability of CLR layer to transfer the feature map from a Pre-trained immutable parameter weight set (starting point) to a target parameter weight set (goal).

Usually, the Pre-trained weight is not good enough due to the domain gap between the Pre-trained dataset/learning paradigm and the target dataset. And a relatively good performance could be achieved by either finetuning the whole backbone on the target dataset (FINETUNE) or learning from scratch (randomly initialized backbone) on the target task dataset (SCRATCH). We will show that with the help of a very cheap CLR layer, a feature map after a pre-trained (non-optimal) model could be reprogrammed towards a "relatively optimal" feature map obtained by either finetuning the whole backbone (FINETUNE) or training from scratch (SCRATCH).

We choose two datasets: CLEVR dataset and the Kannada-MNIST dataset. Model performance on the CLEVR dataset reaches 46.09% with a Pre-trained ResNet-50 backbone + linear head, 97.66% with FINETUNE, and 91.41% with SCRATCH. In this scenario, pretrain has a large accuracy gap with FINETUNE and SCRATCH. It would be interesting to see if the CLR layer could reprogram a feature map obtained from pretrain towards a feature map obtained by FINETUNE and SCRATCH, which shows the ability of CLR layer to fill a large domain gap.

Model performance on the Kannada-MNIST dataset reaches 95.77% with a Pre-trained backbone + linear head,

99.62% with FINETUNE, and 100% with SCRATCH. Here, SCRATCH performs better than FINETUNE, which shows that the pretrained weights may have no benefit (or even harmful) for target task learning. Here we want to show that the CLR layer could reprogram a feature map obtained from pretrain towards the feature map obtained by SCRATCH. We use the feature map after the first convolutional layer in the different models (pretrain, FINETUNE, and SCRATCH). Taking the feature map after the pretrain model as input and the feature map after FINETUNE (or SCRATCH) as output, we utilize a CLR layer (3x3 2D depthwise convolutional kernels) to learn the mapping, i.e. the channel-wise linear transformation between them. The qualitative results are shown in Fig. 2. Specifically, in Fig. 2, we visualize the feature map that initially has a large initial gap between pretrain and FINETUNE (or SCRATCH). The results show that after channel-wise linear transformation, the feature after pretrain could be reprogrammed towards the goal feature (FINETUNE or SCRATCH)

## 3. Bootstrapping results

Fig.5 in the main paper shows the average accuracy as more tasks are learned. However, the gradient of the curve is also influenced by the order of the tasks (i.e., Hard tasks located in earlier sequence will cause average accuracy tends to increase, while easy tasks located in earlier sequence will cause average accuracy tends to decrease) which is entangled with the effect of catastrophic forgetting.

We use bootstrapping to show the tendency of average accuracy when more tasks are learned. Specifically, for any number of tasks ( $t \in (1, \dots, 53)$ ) that we want to conduct in one continual learning setting, we randomly sample  $t$  tasks from the 53 tasks 50,000 times with replacement and compute the mean accuracy (mean) and standard deviation (std). Fig. 3 shows the Bootstrapping statistic results, which show the change of mean and std when we increase the total number of tasks. The X-axis represents the task number  $t$  we want to conduct. For instance, if the continual learning task



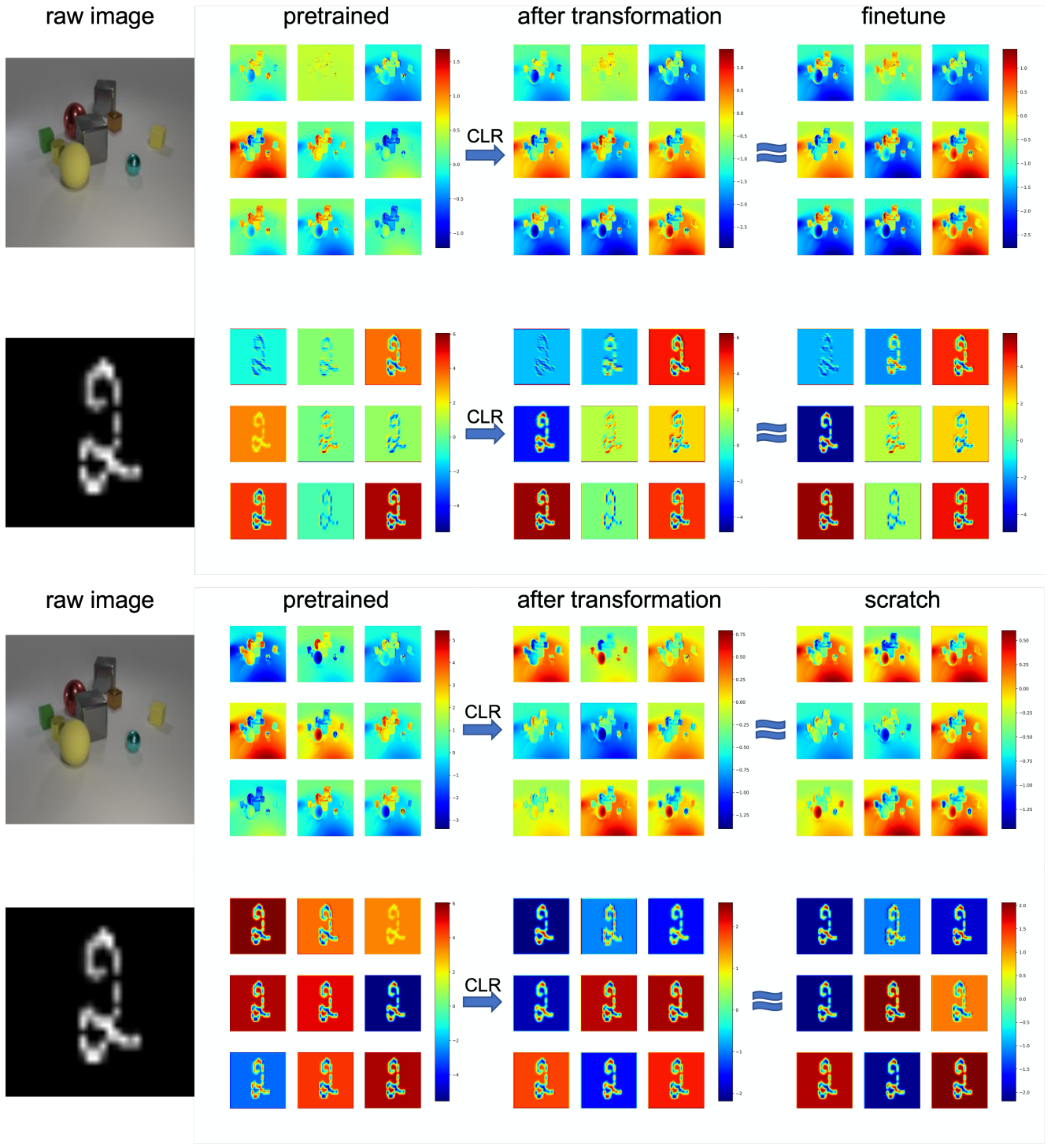


Figure 2. The Figure shows quantitative results of the CLR transformation ability on CLEVR and Kannada-MNIST datasets. We visualize the feature maps in the first residual group of ResNet-50 that initially has a large initial gap between pre-train and FINETUNE (or SCRATCH). The results show that after channel-wise linear transformation, the feature after pre-train could be reprogrammed towards the goal feature (FINETUNE or SCRATCH). Pretrained indicates the frozen Imagenet pretrained ResNet-50 backbone. Finetune is a finetuned ResNet-50 backbone with Imagenet pretrained initialization while Scratch is a trained ResNet-50 backbone with random initialization.

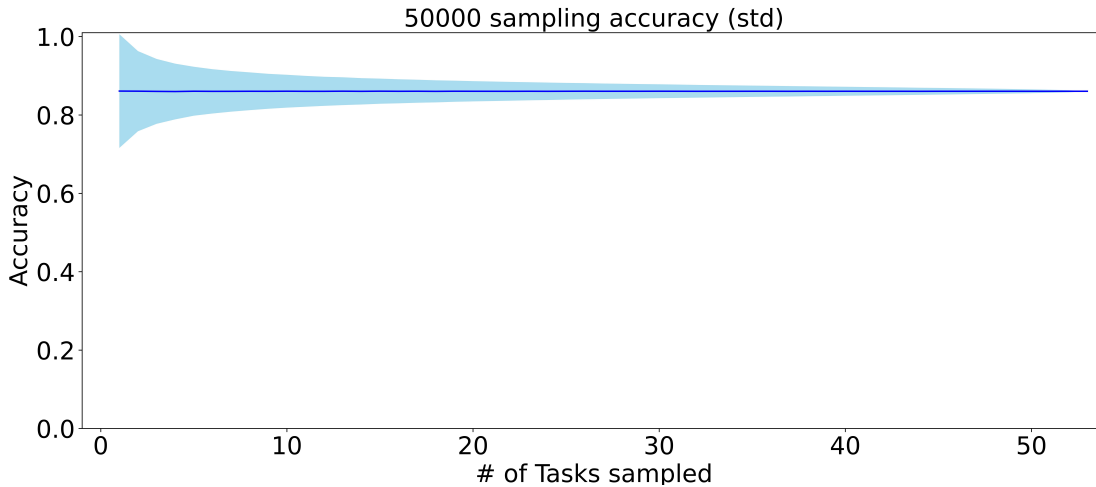


Figure 3. Bootstrapping statistic results. The X-axis represents the number of tasks  $t$  in a specific continual learning task. Y-axis shows the mean Accuracy (solid blue line) on the sampled tasks with replacement and std as the shaded light blue range.

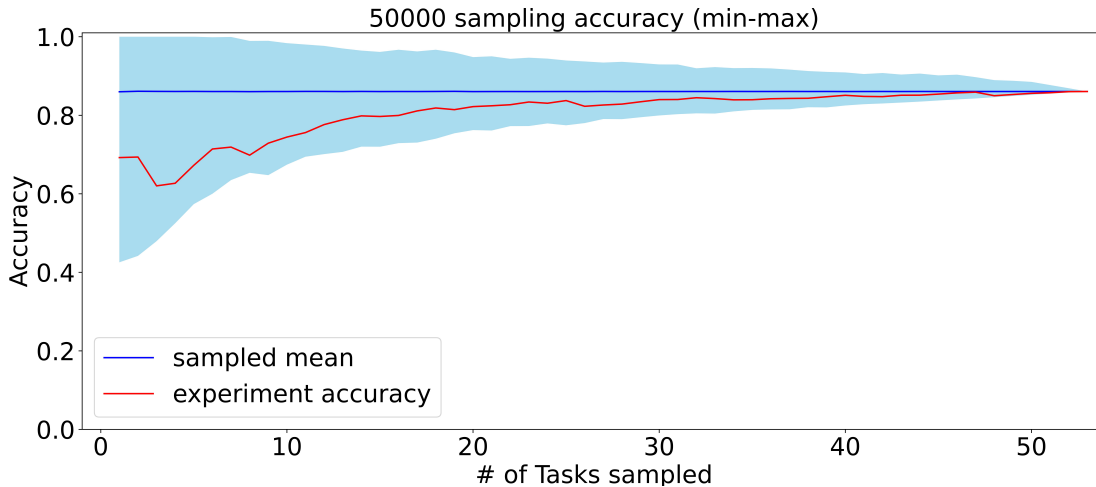


Figure 4. Bootstrapping statistic results with detailed accuracy log. The X-axis represents the number of tasks  $t$  in a specific continual learning task. Y-axis shows the mean Accuracy (solid blue line) on the sampled tasks with replacement. The shaded light blue range shows the min and max range in the given task number  $t$ . We use the solid red line to represent our reported results in the main paper Fig.5, which filled in the shaded light blue range.

CLR layer with  $1 \times 1$  2D reprogramming kernels after all  $1 \times 1$  original Conv kernels and normal CLR layer with  $3 \times 3$  2D reprogramming kernels after the rest CONV kernels. It reaches 85.7% in accuracy and costs  $1.08 \times$  parameters compared to our main method (CLR).

The CLR-mixed version learns a weighted combination of the original and our reprogrammed feature maps. The intuition is that we keep some proportion of the original feature and add the new features learned after reprogramming. Specifically, A trainable parameter  $A$  decides the weight of the summation of the reprogrammed feature and the original feature map. Equation 1 shows the details of the weighted

combination.

$$\hat{x}'_k = A * CLR_k(x'_k) + (1 - A) * x'_k \quad (1)$$

where  $x'_k$  is the  $k$ th channel of the feature map from the original Convolutional layer and  $CLR_k$  is the corresponding linear transformation. It reaches 86.25% in accuracy and costs  $1.79 \times$  parameters compared to our main method (CLR).

The results are shown in Fig. 5. Theoretically, more trainable parameters could lead to better performance, and it is true for CLR-mixed version, which achieves +0.2% than our main method. Interestingly, the CLR-full version achieves lower average accuracy than the main method,

Method	Average Acc
LwF	24%
iCARL	49%
RPS	57%
CCLL	85%
EWC	41%
SI	52%
CLR (Ours)	<b>94.2%</b>

Table 1. We applied our method on CIFAR-100 dataset with 10 tasks, each containing 10 classes with comparisons to baselines from CCLL, using ResNet-18 as the backbone.

while most of the per-task accuracy is higher (43 out of 53 tasks).

## 5. Transfer learning with CLR-based model

We apply our CLR method to the transfer learning problem, where we only care about the accuracy of the transferred dataset while do not need to maintain performance on previous datasets.

**Datasets.** We use the same 53-dataset to evaluate transfer learning performance. Specifically, we use the ImageNet pretrained ResNet-50 model as initialization and apply our method and 4 baseline transfer learning methods 53 times, on 53 different classification tasks.

**Baselines.** We have four baseline methods: 1) learn from scratch (SCRATCH), where the backbone ResNet-50 is randomly initialized with no prior knowledge, and then uses the training set of each task to train the whole network from scratch. 2) finetune the whole backbone and last layer (FINETUNE), 3) finetune only the last layer (LINEAR), 4) Head2Toe method [1] use the fixed backbone and need two steps: 1) feature selection: train the model by adding a large fully connection between all intermediate features and the last layer and select the important connection by adding regularization. 2) keep the important skip connection and re-train the added layers.

Fig. 6 shows the average accuracy on all 53 classification tasks and the details of each task, and Fig. 7 shows the detailed result for transfer learning. Our CLR achieves the best average accuracy on the 53-dataset compared with all baselines. Specifically, CLR achieves almost 5% average improvement on 53 datasets over Head2Toe, and even larger improvement over LINEAR, FINETUNE, and SCRATCH. This shows the effectiveness of the CLR-based model in transfer learning problems.

## 6. CIFAR-100 Result

We also show our method’s result on incremental CIFAR-100 dataset with other previous baselines in table

Task ID	Task Name	Ours	Different Versions of Ours			Task ID	Task Name	Ours	Different Versions of Ours		
		CLR	CLR-Full	CLR-Mixed	CLR-Reduced			CLR	CLR-Full	CLR-Mixed	CLR-Reduced
0	scenes	69.22%	70.90%	73.88%	71.27%	27	boat-types-recognition	89.15%	89.92%	89.92%	86.82%
1	birds	69.50%	70.50%	71.00%	69.50%	28	concrete-crack	100.00%	100.00%	100.00%	100.00%
2	wikiart	47.33%	49.79%	50.21%	44.86%	29	Malacca_Historical_Buildings	100.00%	100.00%	100.00%	100.00%
3	DescribableTextures	64.68%	68.51%	71.91%	66.38%	30	African_countries	84.77%	82.42%	82.81%	83.20%
4	OfficeHome_Clipart	85.54%	86.03%	84.56%	85.78%	31	SurgicalTools	98.51%	99.00%	100.00%	98.51%
5	OfficeHome_Product	92.16%	91.45%	91.45%	91.45%	32	Galaxy10	77.31%	75.38%	78.85%	75.00%
6	OfficeHome_Art	74.80%	74.80%	74.40%	72.80%	33	Stanford_Online_Products	73.41%	75.40%	76.59%	74.21%
7	Food-101	55.45%	57.62%	60.00%	56.63%	34	NWPU-RESISC45	84.81%	84.81%	87.41%	81.85%
8	EuroSAT	97.31%	97.31%	97.69%	98.08%	35	ilab80m-shell	92.94%	93.33%	91.76%	92.55%
9	HistAerial	88.42%	88.80%	86.10%	88.42%	36	CLEVR_v1.0	87.11%	67.58%	69.53%	86.33%
10	OriSet_classification	87.06%	87.06%	88.63%	85.88%	37	CelebA	85.88%	85.88%	87.06%	87.84%
11	Rice_Image_Dataset	100.00%	100.00%	100.00%	100.00%	38	Vegetable_images_Dataset	99.61%	100.00%	99.61%	99.61%
12	garbage_classification	93.65%	95.63%	96.43%	94.05%	39	Monkey_Species	99.23%	99.23%	100.00%	99.23%
13	PokemonData	92.67%	94.44%	94.44%	93.11%	40	aptos2019	73.63%	69.78%	69.23%	71.43%
14	Manga_Facial_Expressions	77.55%	77.55%	77.55%	75.51%	41	freiburg_groceries_dataset	82.40%	84.40%	84.80%	81.60%
15	oregon_wildlife	83.85%	84.23%	81.92%	81.92%	42	Weather_Type_Dataset	100.00%	99.11%	99.11%	99.11%
16	Blood_Cell_Dataset	99.61%	99.22%	99.22%	98.83%	43	Simpsons_Characters_Data	85.91%	89.09%	89.09%	86.36%
17	OCT2017	94.53%	95.70%	94.53%	92.58%	44	Hurricane_Damage_Dataset	97.27%	97.66%	98.05%	97.66%
18	Cataract_Dataset	73.77%	75.41%	72.13%	75.41%	45	Kannada-MNIST	99.62%	100.00%	100.00%	99.62%
19	Apparel_Images_Dataset	96.89%	99.22%	98.44%	97.67%	46	dragon-ball-super-saiyan-dataset	94.44%	72.22%	83.33%	88.89%
20	Zalando_Fashion_Dataset	86.70%	87.12%	85.41%	86.70%	47	ip02-dataset	42.55%	44.71%	45.49%	42.55%
21	House_Room_Images	88.24%	89.02%	89.41%	89.41%	48	planets-and-moons-dataset-ai-in-space	100.00%	100.00%	100.00%	100.00%
22	UIUC_Sports_Event_Dataset	98.75%	98.75%	98.12%	98.12%	49	polish-craft-beer-labels	100.00%	100.00%	100.00%	100.00%
23	Yoga-82	75.78%	75.16%	77.24%	74.95%	50	the-kvasircapsule-dataset	93.61%	95.43%	94.52%	94.06%
24	UMNIST	100.00%	100.00%	100.00%	100.00%	51	atari_shell	100.00%	100.00%	100.00%	100.00%
25	electronic-components	45.82%	47.41%	48.21%	46.61%	52	deepvp_shell	88.10%	85.71%	86.11%	88.49%
26	Breast_Ultrasound	91.14%	97.47%	94.94%	91.14%		Average	86.05%	85.85%	86.25%	85.70%

Figure 5. Per-task accuracy of our main method and other versions of our method after learning all 53 tasks in the continual learning setting.

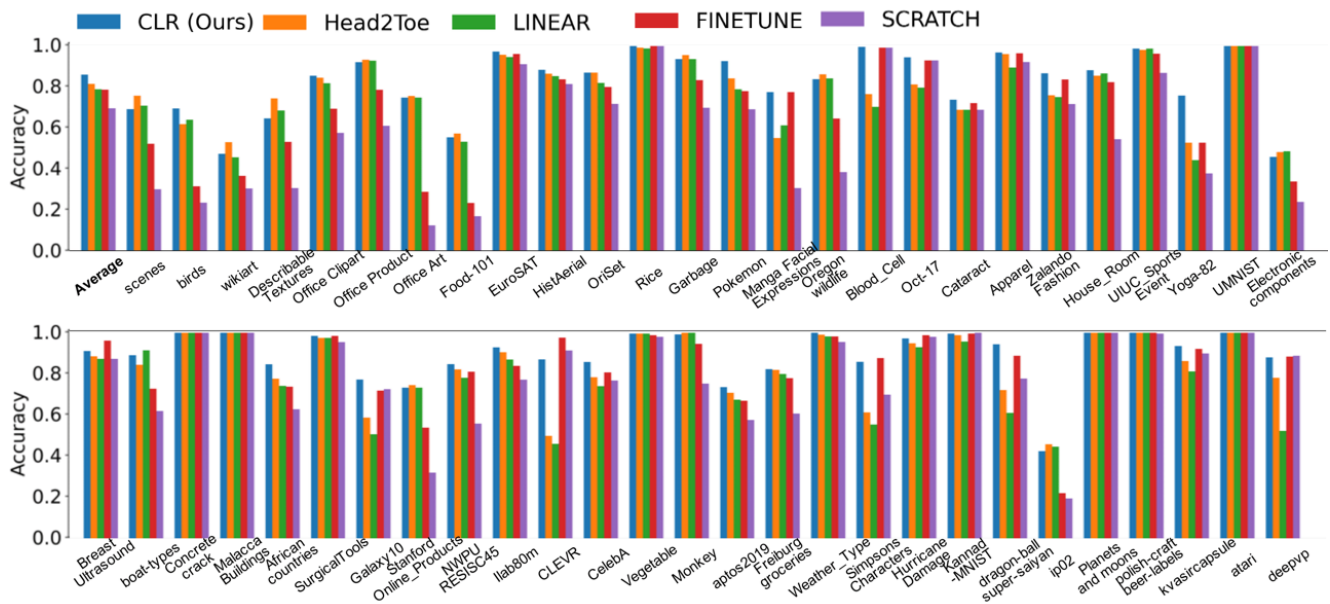


Figure 6. Bar plot of transfer learning performance on 53-dataset.

Task ID	Task Name	Ours	Transfer Learning				Task ID	Task Name	Ours	Transfer Learning			
		CLR	LINEAR	SCRATCH	FINETUNE	Head2Toe			CLR	LINEAR	SCRATCH	FINETUNE	Head2Toe
0	scenes	69.22%	70.90%	30.04%	52.24%	75.75%	27	boat-types-recognition	89.15%	91.47%	62.02%	72.87%	84.50%
1	birds	69.50%	64.00%	23.50%	31.50%	61.83%	28	concrete-crack	100.00%	100.00%	100.00%	100.00%	100.00%
2	wikiart	47.33%	45.68%	30.45%	36.63%	53.09%	29	Malacca_Historical_Buildings	100.00%	100.00%	100.00%	100.00%	100.00%
3	DescribableTextures	64.68%	68.51%	30.64%	53.19%	74.47%	30	African_countries	84.77%	74.22%	62.89%	73.83%	77.73%
4	OfficeHome_Clipart	85.54%	81.86%	57.60%	69.36%	84.56%	31	SurgicalTools	98.51%	97.51%	95.52%	98.51%	97.51%
5	OfficeHome_Product	92.16%	92.87%	61.05%	78.62%	93.35%	32	Galaxy10	77.31%	50.77%	72.69%	71.92%	58.85%
6	OfficeHome_Art	74.80%	74.80%	12.40%	28.80%	75.60%	33	Stanford_Online_Products	73.41%	73.41%	32.14%	53.97%	74.60%
7	Food-101	55.45%	53.27%	16.83%	23.37%	57.23%	34	NWPU-RESISC45	84.81%	78.15%	55.93%	81.11%	82.22%
8	EuroSAT	97.31%	94.62%	91.15%	96.15%	95.77%	35	ilab80m-shell	92.94%	87.06%	77.25%	83.92%	90.59%
9	HistAerial	88.42%	85.33%	81.47%	83.78%	86.49%	36	CLEVR_v1.0	87.11%	46.09%	91.41%	97.66%	50.00%
10	OriSet_classification	87.06%	81.96%	71.76%	80.00%	87.06%	37	CelebA	85.88%	74.12%	76.86%	80.78%	78.43%
11	Rice_Image_Dataset	100.00%	98.82%	100.00%	100.00%	99.22%	38	Vegetable_images_Dataset	99.61%	99.61%	98.04%	98.82%	99.61%
12	garbage_classification	93.65%	93.65%	69.84%	83.33%	95.63%	39	Monkey_Species	99.23%	100.00%	75.38%	94.62%	100.00%
13	PokemonData	92.67%	78.89%	69.11%	78.00%	84.22%	40	aptos2019	73.63%	67.58%	57.69%	67.03%	70.88%
14	Manga_Facial_Expressions	77.55%	61.22%	30.61%	77.55%	55.10%	41	freiburg_groceries_dataset	82.40%	80.00%	60.80%	78.00%	82.00%
15	oregon_wildlife	83.85%	84.23%	38.46%	64.62%	86.15%	42	Weather_Type_Dataset	100.00%	98.21%	95.54%	98.21%	99.11%
16	Blood_Cell_Dataset	99.61%	70.31%	99.22%	99.22%	76.56%	43	Simpsons_Characters_Data	85.91%	55.45%	70.00%	87.73%	61.36%
17	OCT2017	94.53%	79.69%	92.97%	92.97%	81.25%	44	Hurricane_Damage_Dataset	97.27%	92.97%	98.05%	98.83%	94.92%
18	Cataract_Dataset	73.77%	68.85%	68.85%	72.13%	68.85%	45	Kannada-MNIST	99.62%	95.77%	100.00%	99.62%	98.85%
19	Apparel_Images_Dataset	96.89%	89.49%	92.22%	96.50%	96.11%	46	dragon-ball-super-saiyan-dataset	94.44%	61.11%	77.78%	88.89%	72.22%
20	Zalando_Fashion_Dataset	86.70%	75.11%	71.67%	83.69%	75.97%	47	ip02-dataset	42.55%	44.71%	19.61%	22.16%	45.88%
21	House_Room_Images	88.24%	86.67%	54.51%	82.35%	85.49%	48	planets-and-moons-dataset-ai-in-space	100.00%	100.00%	100.00%	100.00%	100.00%
22	UIUC_Sports_Event_Dataset	98.75%	98.75%	86.87%	96.25%	98.13%	49	polish-craft-beer-labels	100.00%	100.00%	99.60%	100.00%	100.00%
23	Yoga-82	75.78%	44.26%	37.79%	52.82%	52.82%	50	the-kvasircapsule-dataset	93.61%	81.28%	89.95%	92.24%	86.30%
24	UMNIST	100.00%	100.00%	100.00%	100.00%	100.00%	51	atari_shell	100.00%	100.00%	100.00%	100.00%	100.00%
25	electronic-components	45.82%	48.61%	23.90%	33.86%	48.21%	52	deepvp_shell	88.10%	52.38%	88.89%	88.49%	78.17%
26	Breast_Ultrasound	91.14%	87.34%	87.34%	96.20%	88.61%		Average	86.05%	78.90%	69.59%	78.72%	81.53%

Figure 7. Transfer learning result on 53-dataset of our method and other baselines (LINEAR, SCRATCH, FINETUNE, and Head2Toe)

## References

- [1] Utku Evci, Vincent Dumoulin, Hugo Larochelle, and Michael C Mozer. Head2toe: Utilizing intermediate representations for better transfer learning. *arXiv preprint arXiv:2201.03529*, 2022. [5](#)