# A    Implementation Details of SNNs

## Spiking Neuron Model

For all SNNs, we use the Integrate-and-Fire (IF) model as the spiking neuron model, which acts as the activation layer in neural networks. As mentioned in (Fang et al. 2021b,a), $V_t$, $X_t$ and $S_t$ denote the state (membrane voltage), input (current) and output (spike) of the spiking neuron model respectively at time-step $t$, and the dynamics of the IF model can be described as follows:

$$H_t = V_{t-1} + X_t, \tag{1}$$
$$S_t = \Theta(H_t - V_{thresh}), \tag{2}$$
$$V_t = H_t(1 - S_t) + V_{reset}S_t. \tag{3}$$

While $V_t$ is the membrane voltage after the trigger of a spike, $H_t$ is also the membrane voltage, but after charging and before a spike firing. $\Theta(x)$ is the unit step function, so $S_t$ equals 1 when $H_t$ is greater than or equal to the threshold voltage $V_{thresh}$ and 0 otherwise. Meanwhile, when a spike fires, $V_t$ is reset to $V_{reset}$. Here, we set $V_{thresh} = 1$ and $V_{reset} = 0$.

In addition, because $\Theta(x)$ is non-differentiable at 0, the surrogate gradient method (Neftci, Mostafa, and Zenke 2019) is applied to approximate the derivative function during back-propagation. Here, we use the inverse tangent function as the surrogate gradient function

$$\sigma(x) = \frac{1}{\pi}\arctan(\pi x) + \frac{1}{2}, \tag{4}$$

and the derivative function is

$$\sigma'(x) = \frac{1}{1 + (\pi x)^2}. \tag{5}$$

## Architectures of SNNs

In our experiments on SNNs, we not only use SEW ResNet proposed by (Fang et al. 2021a), but also build several new SNNs. On the one hand, we improve the spike-element-wise block in SEW ResNet with new architectures referring to studies on ResNet (He et al. 2016; Zagoruyko and Komodakis 2016; Xie et al. 2017), as shown in Table 1. On the other hand, as the multi-branch structures in CNNs increase neural representation similarity to mouse visual cortex, we use depthwise separable convolutions and follow the overall architecture of MobileNetV2 (Howard et al. 2017; Sandler et al. 2018) to build the SpikingMobileNet, the basic block of which is shown in Figure 1.

Our implementation is based on SpikingJelly (Fang et al. 2020), an open-source framework of deep SNN.

## Hyper-Parameters in SNNs' Training

We use the ImageNet dataset to pre-train the new SNNs. Following the settings for training SEW ResNet (Fang et al. 2021a), we train the models for 320 epochs on 8 GPUs (NVIDIA V100), using SGD with a mini-batch size of 32. The momentum is 0.9 and the weight decay is 0. The initial learning rate is 0.1 and we decay it with a cosine annealing, where the maximum number of iterations is the same as the number of epochs. For all SNNs, we set the simulation duration $T = 4$.

# B    Overall model rankings

The results of model rankings are shown in Figure 2, 3 and 4. We also apply the Spearman's rank correlation to the overall model rankings of different metrics, which is shown in Figure 5.

# C    Score Comparisons among Model Groups

We conduct comparisons of similarity scores among CNNs, SNNs, and vision transformers. The results are shown in Figure 6.

# D    Overall CNN rankings

The results of CNN rankings are shown in Figure 7, 8 and 9.

# E    Correlations between the Model Sizes and the Similarity Scores

The results of linear regression to model sizes and the similarity scores are shown in Figure 10, 11 and 12.

# F    The ImageNet Accuracy and the Similarity Scores

The results are shown in Figure 13.

# References

Fang, W.; Chen, Y.; Ding, J.; Chen, D.; Yu, Z.; Zhou, H.; Masquelier, T.; Tian, Y.; and other contributors. 2020. SpikingJelly. https://github.com/fangwei123456/spikingjelly. Accessed: 2022-07-06.

Fang, W.; Yu, Z.; Chen, Y.; Huang, T.; Masquelier, T.; and Tian, Y. 2021a. Deep residual learning in spiking neural networks. In *Advances in Neural Information Processing Systems 34*, 21056–21069.

Fang, W.; Yu, Z.; Chen, Y.; Masquelier, T.; Huang, T.; and Tian, Y. 2021b. Incorporating learnable membrane time constant to enhance learning of spiking neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2661–2671.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.

Howard, A. G.; Zhu, M.; Chen, B.; Kalenichenko, D.; Wang, W.; Weyand, T.; Andreetto, M.; and Adam, H. 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv:1704.04861.

Neftci, E. O.; Mostafa, H.; and Zenke, F. 2019. Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks. *IEEE Signal Processing Magazine*, 36(6): 51–63.

Sandler, M.; Howard, A.; Zhu, M.; Zhmoginov, A.; and Chen, L.-C. 2018. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4510–4520.

Xie, S.; Girshick, R.; Dollár, P.; Tu, Z.; and He, K. 2017. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1492–1500.

Zagoruyko, S.; and Komodakis, N. 2016. Wide residual networks. In *Proceedings of the British Machine Vision Conference (BMVC)*.

| layers | output size | SEW ResNet8 | SEW ResNet14 | Wide SEW ResNet8 | Wide SEW ResNet14 | SEW ResNeXt11 | SEW ResNeXt20 |
|---|---|---|---|---|---|---|---|
| conv with bn | $112 \times 112$ | \multicolumn{6}{c}{$7 \times 7, 64$, stride 2} | | | | | |
| sn | $112 \times 112$ | | | | | | |
| max pool | $56 \times 56$ | \multicolumn{6}{c}{$3 \times 3$, stride 2} | | | | | |
| SEW block1 | $28 \times 28$ | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 1$ | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 128 \times 4 \\ 3 \times 3, 128 \times 4 \end{bmatrix} \times 1$ | $\begin{bmatrix} 3 \times 3, 128 \times 4 \\ 3 \times 3, 128 \times 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64, g = 32 \\ 1 \times 1, 128 \end{bmatrix} \times 1$ | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64, g = 32 \\ 1 \times 1, 128 \end{bmatrix} \times 2$ |
| SEW block2 | $14 \times 14$ | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 1$ | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 256 \times 4 \\ 3 \times 3, 256 \times 4 \end{bmatrix} \times 1$ | $\begin{bmatrix} 3 \times 3, 256 \times 4 \\ 3 \times 3, 256 \times 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128, g = 32 \\ 1 \times 1, 256 \end{bmatrix} \times 1$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128, g = 32 \\ 1 \times 1, 256 \end{bmatrix} \times 2$ |
| SEW block3 | $7 \times 7$ | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 1$ | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 512 \times 4 \\ 3 \times 3, 512 \times 4 \end{bmatrix} \times 1$ | $\begin{bmatrix} 3 \times 3, 512 \times 4 \\ 3 \times 3, 512 \times 4 \end{bmatrix} \times 2$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256, g = 32 \\ 1 \times 1, 512 \end{bmatrix} \times 1$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256, g = 32 \\ 1 \times 1, 512 \end{bmatrix} \times 2$ |
| global average pool | $1 \times 1$ | \multicolumn{6}{c}{$7 \times 7$} | | | | | |
| fc | 1000 | | | | | | |

Table 1: Architectures of SNNs. "sn" denotes the spiking neuron. "$g = 32$" denotes the grouped convolutions with 32 groups. The hyper-parameters of the spike-element-wise block are shown in the brackets with the number of stacked blocks outside.



Figure 1: The basic block of SpikingMobileNet. "PW CONV" is the pointwise convolution and "DW CONV" is the depthwise convolution. "SN" is the spiking neuron.

Figure 2: Overall model rankings of the similarity scores on Allen Brain mouse dataset. The similarity scores of CNNs, SNNs and vision transformers are shown by blue, green and orange bars, respectively.

## SVCCA

| Model | Similarity Score |
|---|---|
| Wide-SEW-ResNet14 | 0.5213 |
| Wide-SEW-ResNet8 | 0.5195 |
| ViT-B-32 | 0.5169 |
| ViT-L-32 | 0.5159 |
| SEW-ResNet101 | 0.5152 |
| LeViT-384 | 0.5143 |
| SEW-ResNet152 | 0.5134 |
| LeViT-256 | 0.5112 |
| LeViT-128S | 0.5105 |
| Inception-ResNet-V2 | 0.5094 |
| LeViT-192 | 0.5089 |
| SEW-ResNet50 | 0.5079 |
| PiT-S | 0.5077 |
| EfficientNet-B5 | 0.5076 |
| LeViT-128 | 0.5067 |
| EfficientNet-B1 | 0.5066 |
| PiT-B | 0.5060 |
| EfficientNet-B7 | 0.5047 |
| RegNetY-8GF | 0.5043 |
| ConvNeXt-Small | 0.5042 |
| SEW-ResNet34 | 0.5041 |
| SEW-ResNet14 | 0.5041 |
| ShuffleNet-V2-x1.0 | 0.5041 |
| PiT-XS | 0.5039 |
| ConvNeXt-Large | 0.5038 |
| SEW-ResNet8 | 0.5035 |
| SEW-ResNet18 | 0.5032 |
| MobileNet-V3-Large | 0.5032 |
| EfficientNet-B3 | 0.5030 |
| CORNet-Z | 0.5028 |
| ResNet101 | 0.5025 |
| PiT-Ti | 0.5023 |
| ConvNeXt-Base | 0.5022 |
| EfficientNet-B6 | 0.5020 |
| DeiT-Ti | 0.5017 |
| DeiT-S | 0.5016 |
| ResNet152 | 0.5011 |
| GoogLeNet | 0.5010 |
| RegNetX-8GF | 0.5007 |
| EfficientNet-B2 | 0.5006 |
| MNASNet-1.0 | 0.5001 |
| ResNeXt101-32x8d | 0.5001 |
| ResNet50 | 0.4997 |
| DeiT-B | 0.4988 |
| ConvNeXt-Tiny | 0.4986 |
| DenseNet201 | 0.4981 |
| RegNetX-800MF | 0.4980 |
| MobileNet-V3-Small | 0.4977 |
| ViT-B-16 | 0.4975 |
| EfficientNet-B4 | 0.4969 |
| Swin-Ti | 0.4966 |
| SEW-ResNeXt20 | 0.4964 |
| DenseNet121 | 0.4963 |
| EfficientNet-B0 | 0.4959 |
| DenseNet169 | 0.4959 |
| TNT-S | 0.4950 |
| Swin-S | 0.4950 |
| CaiT-XXS | 0.4947 |
| Swin-L | 0.4941 |
| RegNetY-800MF | 0.4938 |
| Inception-V3 | 0.4935 |
| Spiking-MobileNet | 0.4933 |
| ResNeXt50-32x4d | 0.4931 |
| SEW-ResNeXt11 | 0.4924 |
| Swin-B | 0.4923 |
| AlexNet | 0.4916 |
| ViT-L-16 | 0.4914 |
| MobileNet-V2 | 0.4906 |
| ResNet18 | 0.4887 |
| CaiT-S | 0.4881 |
| SqueezeNet1.1 | 0.4879 |
| VGG11 | 0.4873 |
| ResNet34 | 0.4861 |
| ConViT-S | 0.4848 |
| ConViT-B | 0.4842 |
| ConViT-Ti | 0.4836 |
| CORNet-S | 0.4816 |
| VGG13 | 0.4806 |
| VGG16 | 0.4784 |
| VGG19 | 0.4761 |
| CORNet-RT | 0.4741 |

## TSVD-Regression

| Model | Similarity Score |
|---|---|
| Wide-SEW-ResNet14 | 0.5027 |
| Wide-SEW-ResNet8 | 0.5013 |
| LeViT-384 | 0.4994 |
| ViT-L-32 | 0.4991 |
| ViT-B-32 | 0.4989 |
| SEW-ResNet101 | 0.4979 |
| SEW-ResNet152 | 0.4979 |
| Inception-ResNet-V2 | 0.4972 |
| LeViT-256 | 0.4968 |
| LeViT-128S | 0.4957 |
| PiT-S | 0.4949 |
| PiT-B | 0.4948 |
| LeViT-192 | 0.4942 |
| PiT-XS | 0.4938 |
| SEW-ResNet50 | 0.4929 |
| DeiT-Ti | 0.4918 |
| PiT-Ti | 0.4917 |
| EfficientNet-B1 | 0.4917 |
| DeiT-S | 0.4916 |
| ShuffleNet-V2-x1.0 | 0.4916 |
| EfficientNet-B5 | 0.4916 |
| SEW-ResNet14 | 0.4914 |
| LeViT-128 | 0.4911 |
| SEW-ResNet34 | 0.4909 |
| SEW-ResNet8 | 0.4908 |
| SEW-ResNet18 | 0.4904 |
| MobileNet-V3-Large | 0.4901 |
| EfficientNet-B7 | 0.4898 |
| EfficientNet-B2 | 0.4897 |
| CORNet-Z | 0.4897 |
| GoogLeNet | 0.4895 |
| ResNeXt101-32x8d | 0.4894 |
| MNASNet-1.0 | 0.4894 |
| RegNetY-8GF | 0.4893 |
| ConvNeXt-Small | 0.4892 |
| ResNet152 | 0.4886 |
| ConvNeXt-Large | 0.4886 |
| EfficientNet-B3 | 0.4886 |
| EfficientNet-B6 | 0.4885 |
| ConvNeXt-Base | 0.4884 |
| DeiT-B | 0.4882 |
| ResNet101 | 0.4881 |
| DenseNet169 | 0.4879 |
| EfficientNet-B4 | 0.4879 |
| TNT-S | 0.4878 |
| DenseNet201 | 0.4877 |
| ViT-B-16 | 0.4877 |
| RegNetX-8GF | 0.4876 |
| SEW-ResNeXt20 | 0.4873 |
| ResNet50 | 0.4873 |
| Swin-Ti | 0.4867 |
| DenseNet121 | 0.4858 |
| ViT-L-16 | 0.4857 |
| MobileNet-V3-Small | 0.4854 |
| CaiT-XXS | 0.4854 |
| ConvNeXt-Tiny | 0.4848 |
| SEW-ResNeXt11 | 0.4848 |
| RegNetX-800MF | 0.4846 |
| Swin-S | 0.4838 |
| ResNeXt50-32x4d | 0.4833 |
| EfficientNet-B0 | 0.4826 |
| Swin-L | 0.4825 |
| Inception-V3 | 0.4825 |
| Swin-B | 0.4819 |
| CaiT-S | 0.4813 |
| AlexNet | 0.4807 |
| RegNetY-800MF | 0.4807 |
| ResNet18 | 0.4806 |
| Spiking-MobileNet | 0.4800 |
| ConViT-S | 0.4798 |
| ConViT-Ti | 0.4785 |
| MobileNet-V2 | 0.4783 |
| ConViT-B | 0.4779 |
| SqueezeNet1.1 | 0.4776 |
| ResNet34 | 0.4773 |
| VGG11 | 0.4758 |
| VGG13 | 0.4725 |
| CORNet-S | 0.4713 |
| CORNet-RT | 0.4697 |
| VGG16 | 0.4685 |
| VGG19 | 0.4657 |

## RSA

| Model | Similarity Score |
|---|---|
| TNT-S | 0.4804 |
| Inception-ResNet-V2 | 0.4786 |
| CaiT-S | 0.4524 |
| EfficientNet-B3 | 0.4491 |
| EfficientNet-B1 | 0.4455 |
| LeViT-384 | 0.4436 |
| RegNetX-8GF | 0.4419 |
| ResNet101 | 0.4365 |
| PiT-S | 0.4349 |
| Wide-SEW-ResNet14 | 0.4310 |
| RegNetY-8GF | 0.4291 |
| LeViT-128 | 0.4291 |
| SEW-ResNet101 | 0.4283 |
| PiT-XS | 0.4281 |
| ResNet50 | 0.4277 |
| DenseNet121 | 0.4274 |
| SEW-ResNet152 | 0.4274 |
| ViT-L-32 | 0.4270 |
| MNASNet-1.0 | 0.4263 |
| SEW-ResNet18 | 0.4248 |
| SEW-ResNeXt11 | 0.4244 |
| ResNet152 | 0.4238 |
| SEW-ResNet50 | 0.4236 |
| ViT-B-32 | 0.4235 |
| DenseNet201 | 0.4234 |
| LeViT-256 | 0.4233 |
| RegNetX-800MF | 0.4231 |
| Spiking-MobileNet | 0.4224 |
| LeViT-192 | 0.4224 |
| EfficientNet-B4 | 0.4198 |
| MobileNet-V2 | 0.4192 |
| Swin-S | 0.4184 |
| PiT-Ti | 0.4177 |
| PiT-B | 0.4171 |
| ResNeXt101-32x8d | 0.4170 |
| SEW-ResNet14 | 0.4163 |
| ConvNeXt-Tiny | 0.4162 |
| Swin-B | 0.4153 |
| LeViT-128S | 0.4139 |
| Swin-Ti | 0.4136 |
| ShuffleNet-V2-x1.0 | 0.4125 |
| DenseNet169 | 0.4125 |
| EfficientNet-B6 | 0.4121 |
| MobileNet-V3-Large | 0.4112 |
| CaiT-XXS | 0.4099 |
| EfficientNet-B2 | 0.4094 |
| DeiT-B | 0.4092 |
| SEW-ResNet34 | 0.4070 |
| EfficientNet-B7 | 0.4064 |
| MobileNet-V3-Small | 0.4047 |
| SEW-ResNeXt20 | 0.4039 |
| Swin-L | 0.4036 |
| SEW-ResNet8 | 0.4025 |
| EfficientNet-B5 | 0.4005 |
| Wide-SEW-ResNet8 | 0.3984 |
| GoogLeNet | 0.3972 |
| Inception-V3 | 0.3963 |
| DeiT-Ti | 0.3960 |
| VGG16 | 0.3959 |
| EfficientNet-B0 | 0.3950 |
| RegNetY-800MF | 0.3934 |
| DeiT-S | 0.3923 |
| AlexNet | 0.3899 |
| ResNet34 | 0.3892 |
| VGG11 | 0.3886 |
| ConvNeXt-Small | 0.3884 |
| ConViT-S | 0.3881 |
| CORNet-Z | 0.3879 |
| ViT-L-16 | 0.3875 |
| ConvNeXt-Base | 0.3865 |
| ConvNeXt-Large | 0.3833 |
| CORNet-S | 0.3821 |
| ConViT-B | 0.3818 |
| ResNeXt50-32x4d | 0.3815 |
| VGG19 | 0.3784 |
| VGG13 | 0.3763 |
| ViT-B-16 | 0.3761 |
| CORNet-RT | 0.3728 |
| SqueezeNet1.1 | 0.3720 |
| ConViT-Ti | 0.3662 |
| ResNet18 | 0.3617 |

Figure 3: Overall model rankings of the similarity scores on Macaque-Face dataset.

Figure 4: Overall model rankings of the similarity scores on Macaque-Synthetic dataset.

Figure 5: The Spearman's rank correlation between the overall model rankings of different metrics. There is a strong correlation between SVCCA and TSVD-Reg, but RSA has weaker correlations with them.



Figure 6: For three datasets and three similarity metrics, each point indicates the final representation similarity score of a model. There are three types of models (CNNs, SNNs, and vision transformers). From top to bottom, each row respectively corresponds to Allen Brain dataset, Macaque-Face dataset, and Macaque-Synthetic dataset. The error bar is the 95% confidence interval across models from a group. The '*' on the horizontal brackets indicates how significant the differences are.

Figure 7: Overall CNN rankings of the similarity scores on Allen Brain mouse dataset. The top 20% models are above the red dotted lines. The networks that are in the top 20% for all three metrics are marked in red.

## SVCCA

| Model | Similarity Score |
|---|---|
| Inception-ResNet-V2 | 0.5094 |
| EfficientNet-B5 | 0.5076 |
| EfficientNet-B1 | 0.5066 |
| EfficientNet-B7 | 0.5047 |
| RegNetY-8GF | 0.5043 |
| ConvNeXt-Small | 0.5042 |
| ShuffleNet-V2-x1.0 | 0.5041 |
| ConvNeXt-Large | 0.5038 |
| MobileNet-V3-Large | 0.5032 |
| EfficientNet-B3 | 0.5030 |
| CORNet-Z | 0.5028 |
| ResNet101 | 0.5025 |
| ConvNeXt-Base | 0.5022 |
| EfficientNet-B6 | 0.5020 |
| ResNet152 | 0.5011 |
| GoogLeNet | 0.5010 |
| RegNetX-8GF | 0.5007 |
| EfficientNet-B2 | 0.5006 |
| MNASNet-1.0 | 0.5001 |
| ResNeXt101-32x8d | 0.5001 |
| ResNet50 | 0.4997 |
| ConvNeXt-Tiny | 0.4986 |
| DenseNet201 | 0.4981 |
| RegNetX-800MF | 0.4980 |
| MobileNet-V3-Small | 0.4977 |
| EfficientNet-B4 | 0.4969 |
| DenseNet121 | 0.4963 |
| EfficientNet-B0 | 0.4959 |
| DenseNet169 | 0.4959 |
| RegNetY-800MF | 0.4938 |
| Inception-V3 | 0.4935 |
| ResNeXt50-32x4d | 0.4931 |
| AlexNet | 0.4916 |
| MobileNet-V2 | 0.4906 |
| ResNet18 | 0.4887 |
| SqueezeNet1.1 | 0.4879 |
| VGG11 | 0.4873 |
| ResNet34 | 0.4861 |
| CORNet-S | 0.4816 |
| VGG13 | 0.4806 |
| VGG16 | 0.4784 |
| VGG19 | 0.4761 |
| CORNet-RT | 0.4741 |

## TSVD-Regression

| Model | Similarity Score |
|---|---|
| Inception-ResNet-V2 | 0.4972 |
| EfficientNet-B1 | 0.4917 |
| ShuffleNet-V2-x1.0 | 0.4916 |
| EfficientNet-B5 | 0.4916 |
| MobileNet-V3-Large | 0.4901 |
| EfficientNet-B7 | 0.4898 |
| EfficientNet-B2 | 0.4897 |
| CORNet-Z | 0.4897 |
| GoogLeNet | 0.4895 |
| ResNeXt101-32x8d | 0.4894 |
| MNASNet-1.0 | 0.4894 |
| RegNetY-8GF | 0.4893 |
| ConvNeXt-Small | 0.4892 |
| ResNet152 | 0.4886 |
| ConvNeXt-Large | 0.4886 |
| EfficientNet-B3 | 0.4886 |
| EfficientNet-B6 | 0.4885 |
| ConvNeXt-Base | 0.4884 |
| ResNet101 | 0.4881 |
| DenseNet169 | 0.4879 |
| EfficientNet-B4 | 0.4879 |
| DenseNet201 | 0.4877 |
| RegNetX-8GF | 0.4876 |
| ResNet50 | 0.4873 |
| DenseNet121 | 0.4858 |
| MobileNet-V3-Small | 0.4854 |
| ConvNeXt-Tiny | 0.4848 |
| RegNetX-800MF | 0.4846 |
| ResNeXt50-32x4d | 0.4833 |
| EfficientNet-B0 | 0.4826 |
| Inception-V3 | 0.4825 |
| AlexNet | 0.4807 |
| RegNetY-800MF | 0.4807 |
| ResNet18 | 0.4806 |
| MobileNet-V2 | 0.4783 |
| SqueezeNet1.1 | 0.4776 |
| ResNet34 | 0.4773 |
| VGG11 | 0.4758 |
| VGG13 | 0.4725 |
| CORNet-S | 0.4713 |
| CORNet-RT | 0.4697 |
| VGG16 | 0.4685 |
| VGG19 | 0.4657 |

## RSA

| Model | Similarity Score |
|---|---|
| Inception-ResNet-V2 | 0.4786 |
| EfficientNet-B3 | 0.4491 |
| EfficientNet-B1 | 0.4455 |
| RegNetX-8GF | 0.4419 |
| ResNet101 | 0.4365 |
| RegNetY-8GF | 0.4291 |
| ResNet50 | 0.4277 |
| DenseNet121 | 0.4274 |
| MNASNet-1.0 | 0.4263 |
| ResNet152 | 0.4238 |
| DenseNet201 | 0.4234 |
| RegNetX-800MF | 0.4231 |
| EfficientNet-B4 | 0.4198 |
| MobileNet-V2 | 0.4192 |
| ResNeXt101-32x8d | 0.4170 |
| ConvNeXt-Tiny | 0.4162 |
| ShuffleNet-V2-x1.0 | 0.4125 |
| DenseNet169 | 0.4125 |
| EfficientNet-B6 | 0.4121 |
| MobileNet-V3-Large | 0.4112 |
| EfficientNet-B2 | 0.4094 |
| EfficientNet-B7 | 0.4064 |
| MobileNet-V3-Small | 0.4047 |
| EfficientNet-B5 | 0.4005 |
| GoogLeNet | 0.3972 |
| Inception-V3 | 0.3963 |
| VGG16 | 0.3959 |
| EfficientNet-B0 | 0.3950 |
| RegNetY-800MF | 0.3934 |
| AlexNet | 0.3899 |
| ResNet34 | 0.3892 |
| VGG11 | 0.3886 |
| ConvNeXt-Small | 0.3884 |
| CORNet-Z | 0.3879 |
| ConvNeXt-Base | 0.3865 |
| ConvNeXt-Large | 0.3833 |
| CORNet-S | 0.3821 |
| ResNeXt50-32x4d | 0.3815 |
| VGG19 | 0.3784 |
| VGG13 | 0.3763 |
| CORNet-RT | 0.3728 |
| SqueezeNet1.1 | 0.3720 |
| ResNet18 | 0.3617 |

Figure 8: Overall CNN rankings of the similarity scores on Macaque-Face dataset.

## SVCCA

| Model | Similarity Score |
|---|---|
| VGG16 | 0.4716 |
| RegNetY-800MF | 0.4694 |
| ShuffleNet-V2-x1.0 | 0.4683 |
| ConvNeXt-Tiny | 0.4683 |
| EfficientNet-B1 | 0.4669 |
| VGG11 | 0.4664 |
| VGG13 | 0.4663 |
| RegNetY-8GF | 0.4655 |
| VGG19 | 0.4651 |
| EfficientNet-B0 | 0.4640 |
| MNASNet-1.0 | 0.4638 |
| EfficientNet-B5 | 0.4616 |
| ResNet101 | 0.4612 |
| DenseNet121 | 0.4605 |
| AlexNet | 0.4595 |
| MobileNet-V3-Large | 0.4593 |
| ResNet152 | 0.4590 |
| RegNetX-800MF | 0.4581 |
| EfficientNet-B7 | 0.4577 |
| ResNet34 | 0.4577 |
| RegNetX-8GF | 0.4575 |
| Inception-ResNet-V2 | 0.4572 |
| EfficientNet-B2 | 0.4568 |
| DenseNet201 | 0.4564 |
| EfficientNet-B6 | 0.4554 |
| SqueezeNet1.1 | 0.4541 |
| ResNeXt101-32x8d | 0.4541 |
| DenseNet169 | 0.4537 |
| Inception-V3 | 0.4534 |
| EfficientNet-B3 | 0.4521 |
| MobileNet-V2 | 0.4500 |
| ResNet18 | 0.4497 |
| CORNet-Z | 0.4480 |
| MobileNet-V3-Small | 0.4479 |
| CORNet-S | 0.4470 |
| EfficientNet-B4 | 0.4462 |
| ConvNeXt-Small | 0.4446 |
| ResNet50 | 0.4418 |
| GoogLeNet | 0.4416 |
| CORNet-RT | 0.4404 |
| ConvNeXt-Large | 0.4387 |
| ResNeXt50-32x4d | 0.4332 |
| ConvNeXt-Base | 0.4332 |

## TSVD-Regression

| Model | Similarity Score |
|---|---|
| Inception-ResNet-V2 | 0.5522 |
| RegNetY-8GF | 0.5513 |
| RegNetY-800MF | 0.5486 |
| VGG16 | 0.5472 |
| VGG19 | 0.5456 |
| AlexNet | 0.5452 |
| VGG11 | 0.5439 |
| ShuffleNet-V2-x1.0 | 0.5417 |
| VGG13 | 0.5416 |
| SqueezeNet1.1 | 0.5403 |
| GoogLeNet | 0.5379 |
| ResNet101 | 0.5375 |
| CORNet-Z | 0.5360 |
| DenseNet121 | 0.5358 |
| ResNet152 | 0.5351 |
| ResNet34 | 0.5344 |
| RegNetX-800MF | 0.5336 |
| DenseNet169 | 0.5333 |
| ResNeXt101-32x8d | 0.5300 |
| ResNet18 | 0.5296 |
| MNASNet-1.0 | 0.5291 |
| ConvNeXt-Tiny | 0.5287 |
| CORNet-S | 0.5281 |
| RegNetX-8GF | 0.5271 |
| ResNet50 | 0.5265 |
| Inception-V3 | 0.5261 |
| DenseNet201 | 0.5255 |
| MobileNet-V2 | 0.5237 |
| EfficientNet-B1 | 0.5221 |
| EfficientNet-B5 | 0.5217 |
| EfficientNet-B0 | 0.5214 |
| MobileNet-V3-Large | 0.5192 |
| ResNeXt50-32x4d | 0.5189 |
| ConvNeXt-Small | 0.5178 |
| EfficientNet-B3 | 0.5176 |
| EfficientNet-B2 | 0.5138 |
| ConvNeXt-Base | 0.5128 |
| MobileNet-V3-Small | 0.5125 |
| EfficientNet-B6 | 0.5120 |
| EfficientNet-B4 | 0.5103 |
| EfficientNet-B7 | 0.5101 |
| ConvNeXt-Large | 0.5095 |
| CORNet-RT | 0.5050 |

## RSA

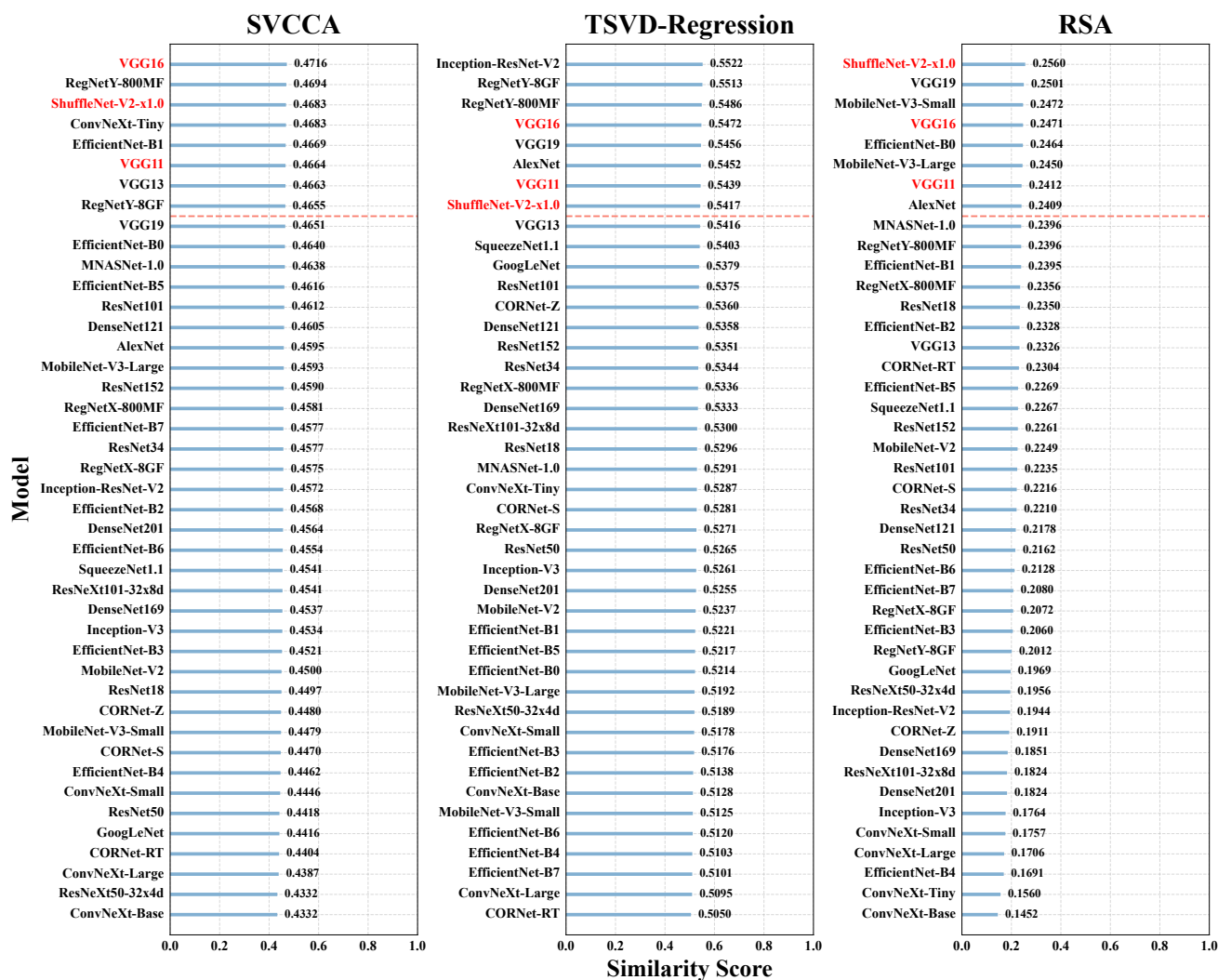| Model | Similarity Score |
|---|---|
| ShuffleNet-V2-x1.0 | 0.2560 |
| VGG19 | 0.2501 |
| MobileNet-V3-Small | 0.2472 |
| VGG16 | 0.2471 |
| EfficientNet-B0 | 0.2464 |
| MobileNet-V3-Large | 0.2450 |
| VGG11 | 0.2412 |
| AlexNet | 0.2409 |
| MNASNet-1.0 | 0.2396 |
| RegNetY-800MF | 0.2396 |
| EfficientNet-B1 | 0.2395 |
| RegNetX-800MF | 0.2356 |
| ResNet18 | 0.2350 |
| EfficientNet-B2 | 0.2328 |
| VGG13 | 0.2326 |
| CORNet-RT | 0.2304 |
| EfficientNet-B5 | 0.2269 |
| SqueezeNet1.1 | 0.2267 |
| ResNet152 | 0.2261 |
| MobileNet-V2 | 0.2249 |
| ResNet101 | 0.2235 |
| CORNet-S | 0.2216 |
| ResNet34 | 0.2210 |
| DenseNet121 | 0.2178 |
| ResNet50 | 0.2162 |
| EfficientNet-B6 | 0.2128 |
| EfficientNet-B7 | 0.2080 |
| RegNetX-8GF | 0.2072 |
| EfficientNet-B3 | 0.2060 |
| RegNetY-8GF | 0.2012 |
| GoogLeNet | 0.1969 |
| ResNeXt50-32x4d | 0.1956 |
| Inception-ResNet-V2 | 0.1944 |
| CORNet-Z | 0.1911 |
| DenseNet169 | 0.1851 |
| ResNeXt101-32x8d | 0.1824 |
| DenseNet201 | 0.1824 |
| Inception-V3 | 0.1764 |
| ConvNeXt-Small | 0.1757 |
| ConvNeXt-Large | 0.1706 |
| EfficientNet-B4 | 0.1691 |
| ConvNeXt-Tiny | 0.1560 |
| ConvNeXt-Base | 0.1452 |

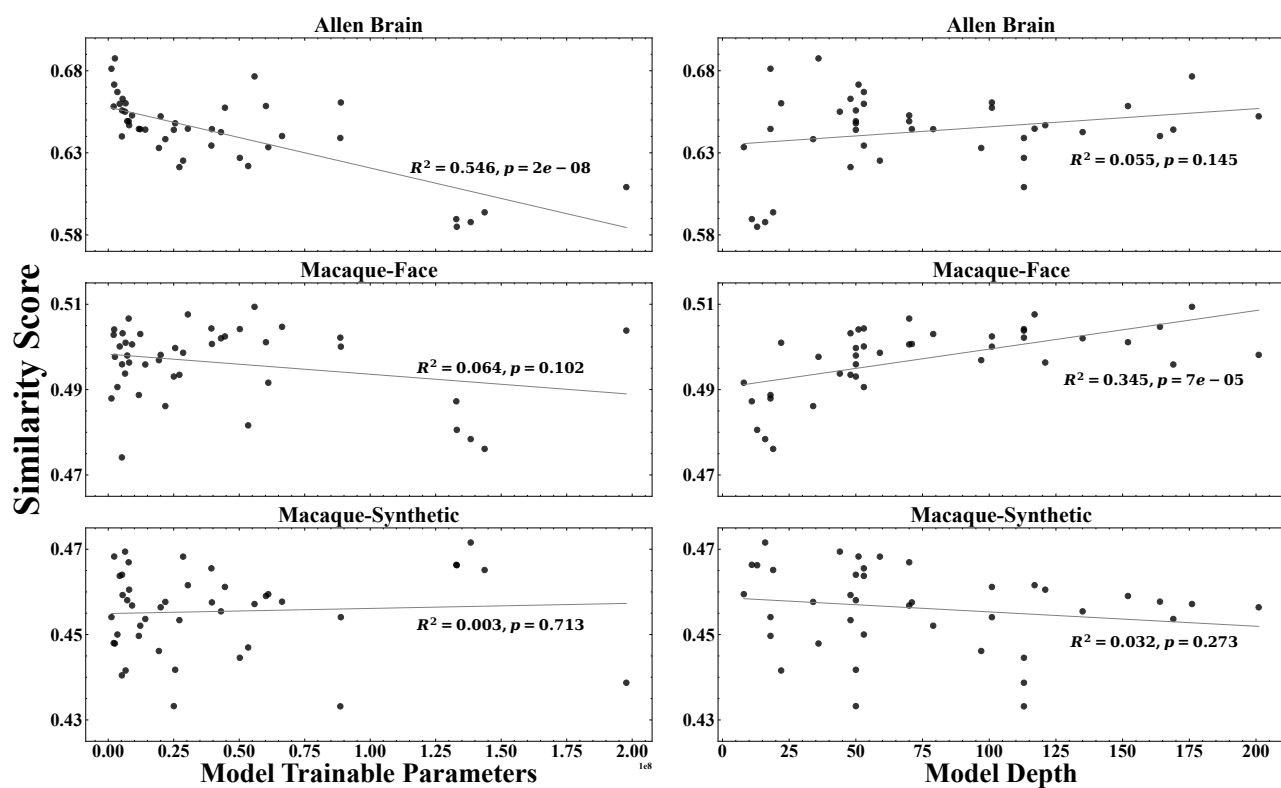Figure 9: Overall CNN rankings of the similarity scores on Macaque-Synthetic dataset.

Figure 10: The results of linear regression to model sizes and the similarity scores for SVCCA. Each point indicates a model.
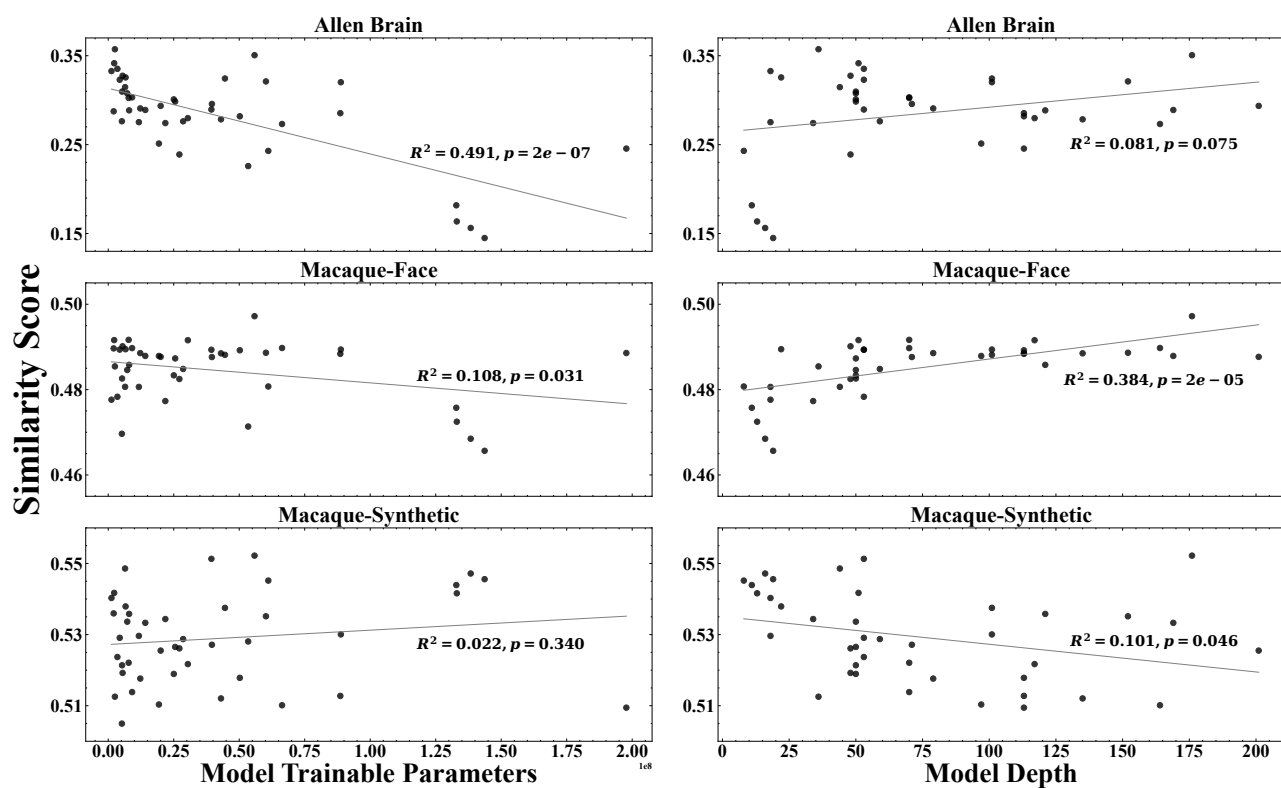
Figure 11: The results of linear regression to model sizes and the similarity scores for TSVD-Reg.
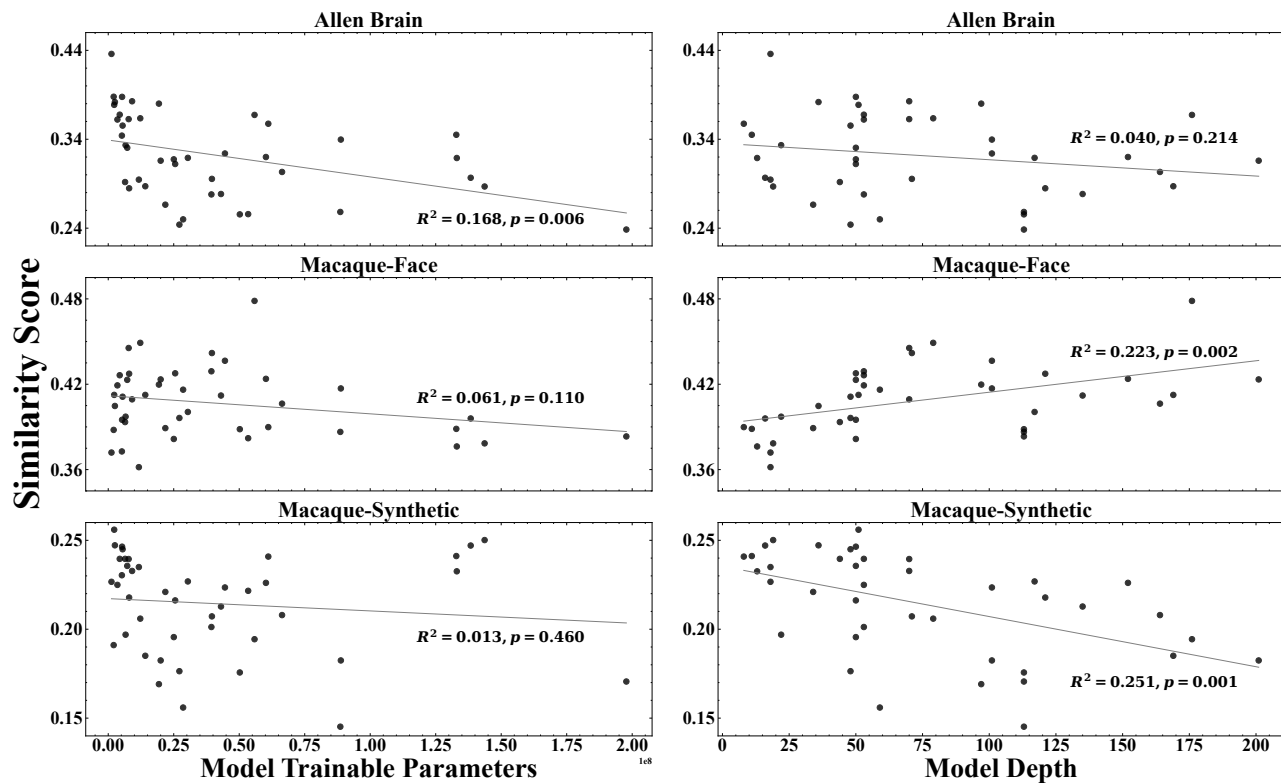
Figure 12: The results of linear regression to model sizes and the similarity scores for RSA.
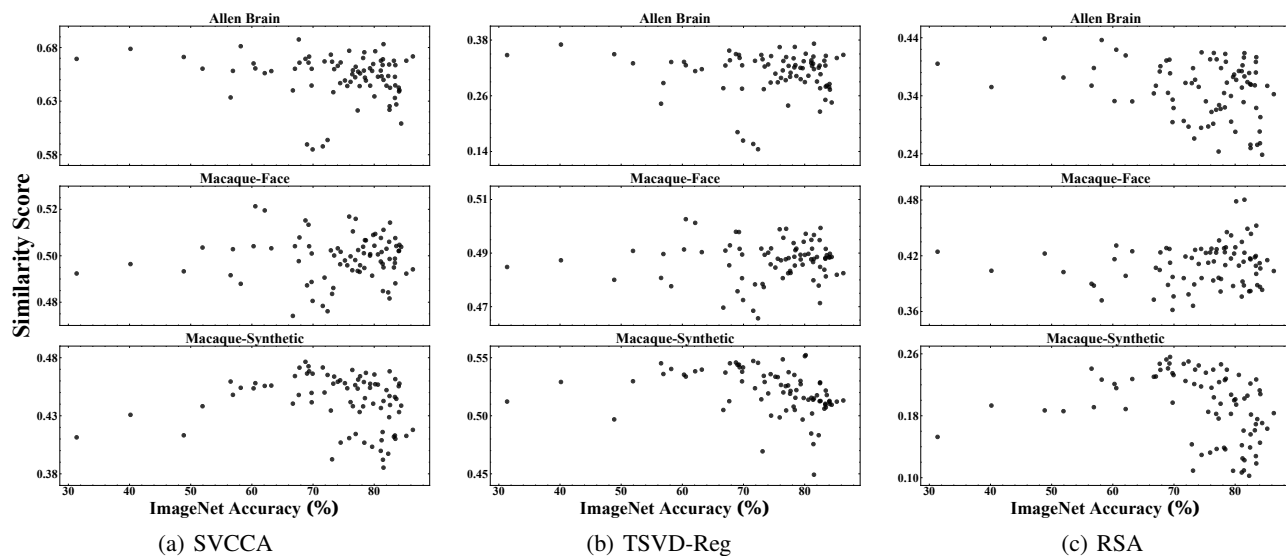


Figure 13: The ImageNet accuracy and the similarity scores for SVCCA, TSVD-Reg and RSA.