# Visualizing Transfer Learning

Róbert Szabó [1]   Dániel Katona [2]   Márton Csillag [1]   Adrián Csiszárik [3] [1]   Dániel Varga [3]

## Abstract

We provide visualizations of individual neurons of a deep image recognition network during the temporal process of transfer learning. These visualizations qualitatively demonstrate various novel properties of the transfer learning process regarding the speed and characteristics of adaptation, neuron reuse, spatial scale of the represented image features, and behavior of transfer learning to small data. We publish the large-scale dataset that we have created for the purposes of this analysis.

## 1. Introduction

Deep neural networks are still commonly conceptualized as black boxes, despite all the recent progress made in interpretability and feature visualization (Buhrmester et al., 2019; Olah et al., 2020; Hohman et al., 2019; Rathore et al., 2019; Bau et al., 2017; Selvaraju et al., 2019).

The current work is following in the footsteps of the Clarity research programme (Olah et al., 2017; 2018; Carter et al., 2019; Olah et al., 2020), both in the techniques employed, and in the qualitative flavour of the research: creating images of neurons, and trying to identify interesting patterns.

Our main focus is using feature visualization to get a better understanding of what happens during transfer learning, both by comparing neurons before-and-after transfer learning, and by observing what happens during the transfer learning process. We also present a channel visualization technique employing a learned prior (Nguyen et al., 2016) utilizing the StyleGAN2 generator (Karras et al., 2020).

Another output of the current work is a large-scale visualization of the transfer learning behavior of an InceptionV1 network. This mapping of the InceptionV1 network with its 57 convolutional layers and 7280 channels on four datasets resulted in approximately 30 000

[1]Eötvös Loránd University, Budapest, Hungary [2]Budapest University of Technology and Economics, Budapest, Hungary [3]Alfréd Rényi Institute of Mathematics, Budapest, Hungary. Correspondence to: Dániel Varga <daniel@renyi.hu>.

*Figure 1.* Visualizations for layer `Mixed_5c_Branch_3_b_1x1` channels. Columns correspond to channels, top row shows Lucid visual, bottom row shows StyleGAN2-based visual. Note the strong correspondence of facial features between the two kinds of visualizations.

visualization images. The produced dataset is presented in a browsable form at `https://bit.ly/visualizing-transfer-learning`.

## 2. Feature visualization via activation maximization

Gradient based methods of feature visualization strive to maximize the aggregated activation of a network layer, channel, or single neuron by computing the activation's gradient with respect to the input image, and doing gradient ascent (Zeiler & Fergus, 2014; Mahendran & Vedaldi, 2015; Simonyan et al., 2013; Mordvintsev et al., 2015; Olah et al., 2017). To achieve good results, some image parametrizations or priors must be added that guide the optimization process to an output interpretable to human observers.

Usually the goal is to add priors that bring the least amount of their own biases, but a particularly interesting exception is the use of generative models: feeding the output of an independently trained generator to the inspected model, and doing regularized gradient ascent in the latent space of the generator (Nguyen et al., 2016). We employ this technique on a CelebA classifier's top convolutional layer, using the StyleGAN2 generator (Karras et al., 2020). In effect, we can solve the highly nonlinear activation maximization of face recognition neurons, constrained to the manifold of face images.

More precisely, with an $F : \mathbb{R}^{w_1 \times h_1 \times 3} \to \mathbb{R}^m$ recognition network, and a $G : \mathbb{R}^d \to \mathbb{R}^{w_2 \times h_2 \times 3}$ synthesis network (where $w_1, h_1$ and $w_2, h_2$ are the spatial dimensions of the images for the two models, respectively), we solve the following soft-constrained optimization task using SGD:

$$w^* = \underset{w \in \mathbb{R}^d}{\arg\max} \, F_z(downscale(G(w))) - \lambda \|w - \hat{w}\|^2,$$

where $F_z$ is the average-pooled activation function of channel $z$, $\hat{w}$ is the center of gravity of StyleGAN2's intermediate $W$ space (Karras et al., 2020), $\lambda$ is a multiplier hyperparameter tuning the trade-off between more realistic/typical faces and faces activating neuron $z$. The optimization is started from $\hat{w}$. The visualization obtained is $G(w^*)$.

It is important to note that gradient-based feature visualization (even when combined with e.g. the Lucid framework's diversity feature (Olah et al., 2017)) only presents a facet of the functionality of a channel or neuron, and higher layer neurons are probably always multi-faceted (Szegedy et al., 2014; Olah et al., 2020). We will limit ourselves to observations that do not assume that an apparent functionality of a channel is its *only* functionality.

## 3. Setup

Unless otherwise noted, all our experiments use the Lucid framework for visualizing convolutional channels via activation maximization (Olah et al., 2017). We use Lucid's 2D FFT image representation with decorrelation. Lucid has the capability to optimize individual neural activations within a channel, and its authors had success with optimizing the spatially central neuron of each channel, but as we achieved better results on our tasks when optimizing the average-pooled activations of channels, this is what we use in all our presented visualizations.

The network that we analyze in all of our experiments is an InceptionV1 (Szegedy et al., 2015), with two dense layers at the top (these are: dense(1024, Relu), dense(number_of_classes, activation_func), where activation_func is softmax for Flowers17 and Animals, and coordinate-wise sigmoid for CelebA. Accordingly, the loss functions were categorical cross-entropy and mean coordinate-wise binary cross-entropy respectively.

The InceptionV1 network is a standard choice in the feature visualization community. After all the progress in classification performance, this network still appears to be the best option if the goal is easy to interpret gradient-based feature visualization. (We are not aware of any explanation of this phenomenon.) We however deviate from (Olah et al., 2017) in using a batch-normalized InceptionV1 variant. An unfortunate minor side effect of this choice is the lack of direct correspondence between our neurons and the Activation Atlas (Carter et al., 2019).

We employ three datasets for the transfer learning task: The CelebFaces Attributes Dataset (CelebA) (Liu et al., 2018) is a large-scale face attributes dataset with more than 200K face images, each annotated with 40 binary attributes. The Flowers17 dataset (Nilsback & Zisserman, 2006) contains 17 flower categories with 80 images for each class. The Animals dataset contains 31 dog and cat breed categories, 200 image of each.

In each case, we start from network weights trained on ImageNet.

No layers were frozen during transfer learning. The networks were trained for 200 epochs with the Adam optimizer with learning rate 0.001 and batch size of 32.

Validation accuracies were 0.94 for Animals, 0.99 for Flowers17, mean binary accuracy was 0.90 for CelebA. Feature visualizations were created based on these networks.

Visualization of the temporal process of transfer learning used a different setup. An InceptionV1 network was trained on the CelebA dataset for 3000 iterations with batch size 10. The low batch size was chosen to show a finer detail about the first few iterations that already result in large changes in the visualized features. The top only had a single dense layer.

A technical detail regarding the StyleGAN2 synthesis network is that it also has noise variables as input. In each gradient ascent step these are sampled, leading to nondeterministic results, as we will present in Figure 4.

## 4. Discussion

In this section we present qualitative observations we made while inspecting our visualizations. We highlight visual evidence for each, but as there is a thin line between highlighting and cherry-picking, we encourage the reader to browse the complete set of visualizations at `https://bit.ly/visualizing-transfer-learning` to verify the claims or make their own observations. Figure 7 also presents a completely unbiased sample of channels. Forming and validating quantitative hypotheses based on the observations is an important further step that we intend to make in follow-up work.

**Adaptation happens early in the transfer learning process.** While monitoring the progress of transfer learning, the feature visualizations show a surprisingly early emergence of features from the target domain. The visualizations of Figure 6 show that in the case of adapting e.g. from ImageNet to CelebA, several characteristic CelebA features appear after training on only 100-600 images from the target domain (2nd to 4th row in Figure 6). Note that while ImageNet data contains images of human faces, these are
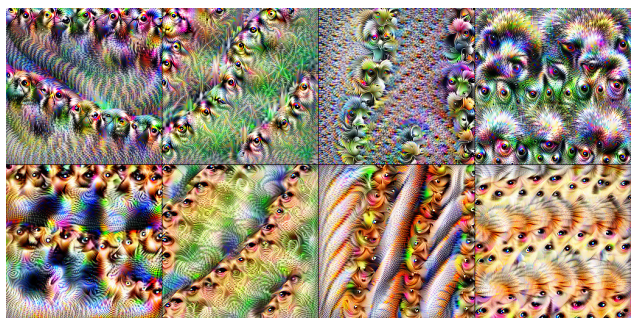
*Figure 2.* Top row shows pre transfer, bottom row shows post transfer (from ImageNet to CelebA) neuron visualizations. Layer: `Mixed_4c_Branch_2_b_3x3`.



*Figure 3.* Channels from a single layer. Each channel has a fundamental spatial scale, and it is stable during transfer learning. Layer: `Mixed_4e_Branch_2_b_3x3`.

on a very different spatial scale compared to CelebA data. We provide some before-transfer feature visualizations for comparison on Figure 7. They rarely show highly similar patterns at the higher layers, and this fact strongly suggests that the face features that emerge early are acquired during transfer learning rather than already latently present in the ImageNet model.

**Even middle layer detectors can be reused.** Contrasting with, but not contradicting the previous point, it does happen that CelebA facial feature detectors are direct descendants of ImageNet feature detectors. For example, ImageNet contains many animal faces at the right scale, and detectors for these face parts are often subtly adjusted for the target domain, as seen on Figure 2.

**Each channel has a fundamental spatial scale, and it is stable during transfer.** Our channel-wise optimized visualizations have a distinct organic repeating pattern, and the period can change from 4 pixels to circa 120 pixels. (The InceptionV1 input size is $224 \times 224$.) Each layer has a predominant period, and this value increases with the depth of the layer, as expected. But this period can still vary significantly within a single layer, and this demonstrates that neurons have an effective receptive field size which is not necessarily the same as their convolutional receptive field size. Remarkably, this effective receptive field size tends to be stable during transfer learning. Figure 3 and Figure 6 demonstrate the phenomenon.

**If the dataset is small, the channels solve the same tasks repeatedly.** When comparing the results of the large CelebA dataset with the small Animals and Flowers17 datasets, it is apparent that in the latter cases, the higher level neurons do not manage to find a diverse enough set of features, and this manifests itself in high redundancy across channels: very often it is impossible to distinguish the visuals of two channels. In principle, this could be an artifact of the optimization process used for

feature visualization, but its prevalence and consistency suggests otherwise. The paper only presents visualizations on CelebA, we refer the reader to our webpage `https://bit.ly/visualizing-transfer-learning` for the Animals and Flowers17 collection.

**The lower layers rarely adapt besides adjusting color space.** Below the `Mixed_3c` layer, the predominant form of adaptation is adjusting the color scheme the channel is most interested in, without altering the pattern. Moving further across layers, first subtle, then less subtle changes appear when comparing the pre- and post-transfer visuals. See the first two rows of Figure 7 for visualizations at lower depths.

**Image-scale structures can emerge.** A peculiar property of InceptionV1 feature visualizations is that the middle layers are the most interpretable to the human eye. We report that this is an artifact of the thematically diverse ImageNet dataset, and with our datasets the top layers are the most interpretable. In particular, the top convolutional layers of our CelebA network are apparently sensitive to complete face images, and the visualizations reproduce much of the complexity of human faces, see Figure 1 top row visualizing some top convolutional layer neurons.

### 4.1. Feature Visualization with a Generative Prior

Our generator based visualizations are mostly believable as real human faces, although they have some of the characteristics of caricatures, exaggerating facial features and overrepresenting unusual shapes and forms when utilizing smaller $\lambda$ values.

Repeated runs tend to converge to similar but not identical latent points and generated images, see Figure 4. Note that the only source of randomness during optimization is the noise input of the synthesis network, so the technique gives only a limited view into the true diversity of local optima.

*Figure 4.* Repeated runs result in similar but not identical images. Utilized layer for visualization: `Mixed_5c_Branch_3_b_1x1`.

Importantly, comparing the Lucid image and the corresponding StyleGAN2 image as done on Figure 1, we can see that they reinforce and refine each other's message. Reinforce in the sense that organic face-like patterns of Lucid and the photorealistic StyleGAN2 images often seem to picture the same (nonexistent) person, or at least they present unique facial features appearing on both images. Refine in the sense that some features that are quite robustly manifested even across many repeated runs, can turn out to be an artifact of the optimization process, for example the gender of the person. See the accompanying webpage https://bit.ly/visualizing-transfer-learning for more examples.

In a crude preliminary approximation of *circuit editing* (Olah et al., 2020) presented in Figure 5, we have modified the weights of a single top convolutional layer neuron, and visualized the result with the generative prior. The modification strategy was picking the top $k$ filter weights by magnitude and negate them. Most of the weight modifications do not meaningfully affect the visualization outcome, but some affect them significantly. In the last elements of the array we can see that at some point, damaging the filter makes it impossible for the optimizer to move away from the $\hat{w}$ prior.

## 5. Conclusion

In this work we provided visualizations of individual neurons of a deep image recognition network during the temporal process of transfer learning. These visualizations qualitatively demonstrate various novel properties of the transfer learning process regarding the speed and characteristics of adaptation, neuron reuse, spatial scale of the represented image features, and behavior of transfer learning to small data.

Even though we tried to limit our exposure to this effect by promoting observations that show an unambiguous signal, and can in principle be formalized and quantitatively veri-



*Figure 5.* Neural weight ablation: filter weights of a single channel of the top convolutional layer were ranked by magnitude. The first $k$ weights were negated, and the ablated channel was visualized with the generative prior. $k$ running from 0 to 39, presented in row major order on a grid.

fied, the main risk of the qualitative approach we pursue is that the human brain is prone to making up stories where information is ambiguous. Hence, utilizing human visual pattern matching must be paired with some quantitative follow-up analysis when the goal is to make claims about neural networks. This is our ongoing work.
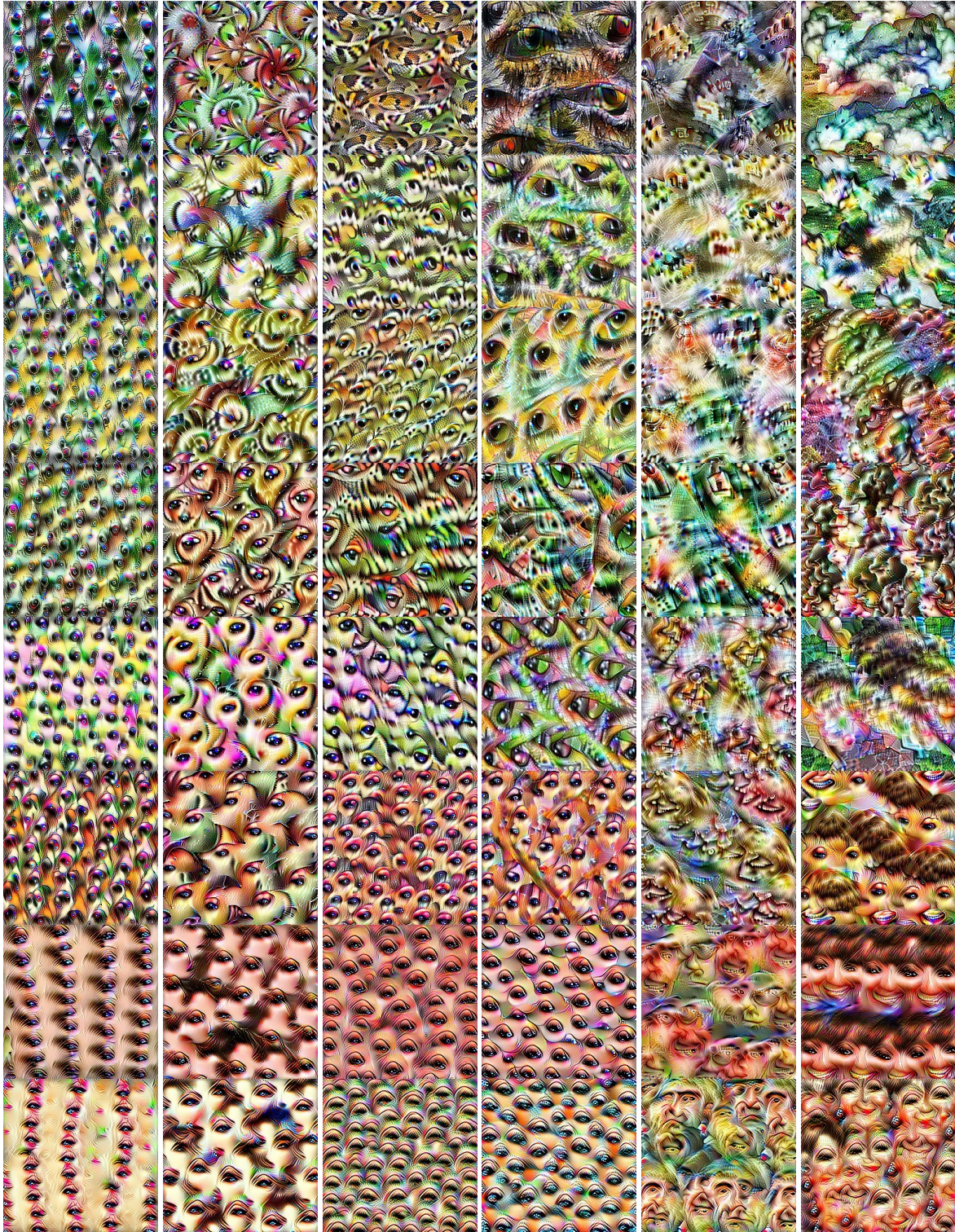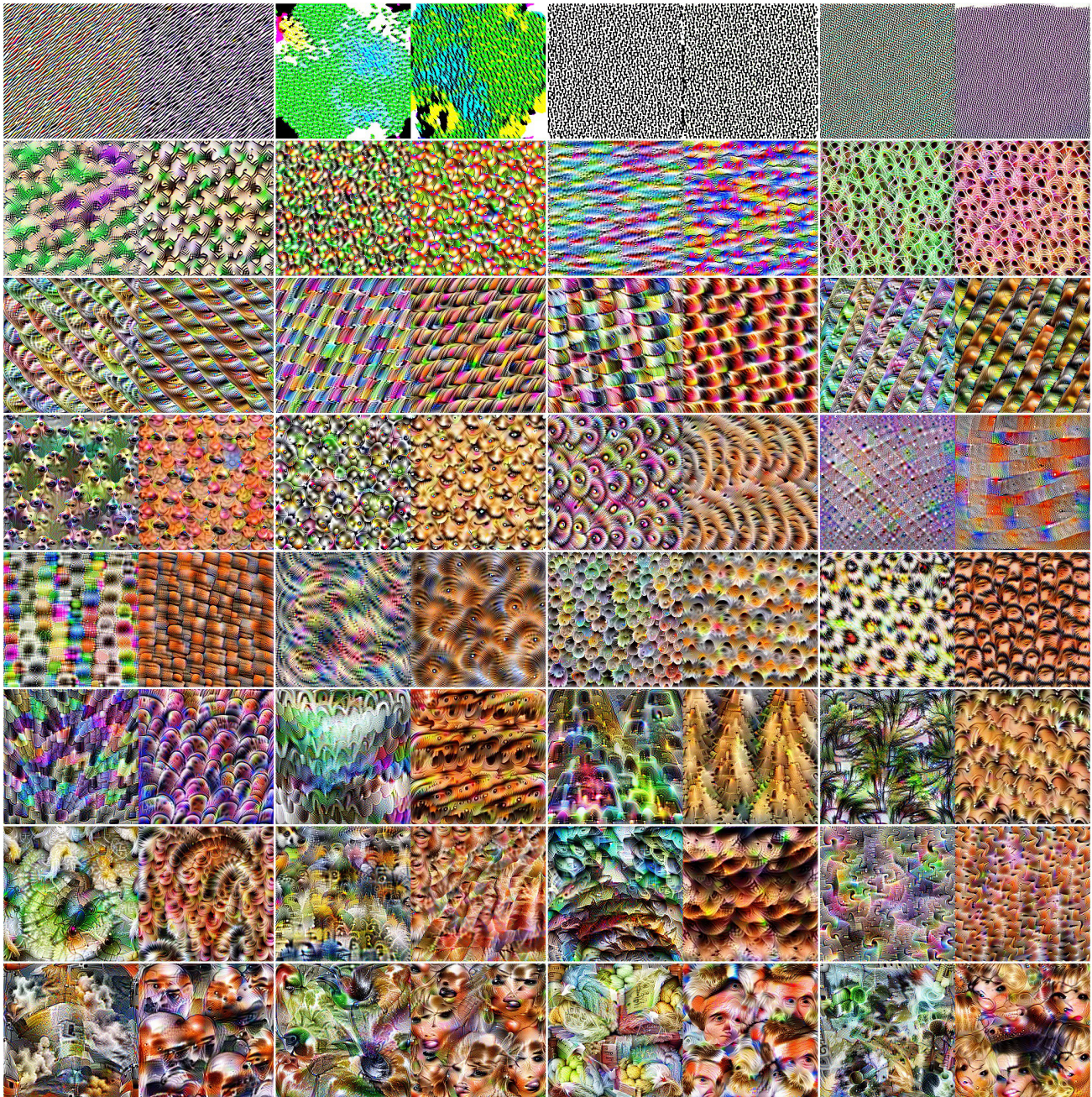
## Acknowledgements

*Figure 6.* Visualization of transfer learning (from ImageNet to CelebA) at different training iterations. Columns show distinct channels, rows show the following iterations, respectively: 0, 10, 20, 30, 60, 150, 1000, 3000. Note that for this experiment the batch size is 10 to visualize a finer grained detail of the transfer learning process.

(a) `4c_2_b_3x3`  (b) `4d_3_b_1x1`  (c) `4e_2_a_1x1`  (d) `4f_1_a_1x1`  (e) `5b_3_b_1x1`  (f) `5c_0_a_1x1`

*Figure 7.* Visualization of transfer learning from ImageNet to CelebA at different layer depths. Each row corresponds to a single layer, odd columns correspond to layer's first 4 channels pre-transfer, even columns are the same neurons post-transfer. The selected layers are every 7th layer ending with the deepest layer, namely `Conv2d_2b_1x1`, `Mixed_3b_Branch_3_b_1x1`, `Mixed_4b_Branch_1_a_1x1`, `Mixed_4c_Branch_2_a_1x1`, `Mixed_4d_Branch_0_a_1x1`, `Mixed_4e_Branch_1_b_3x3`, `Mixed_4f_Branch_2_b_3x3`, `Mixed_5b_Branch_3_b_1x1`. (The periodicity of the InceptionV1 layers is 6.) Layers are deeper from top to down.

## References

Bau, D., Zhou, B., Khosla, A., Oliva, A., and Torralba, A. Network dissection: Quantifying interpretability of deep visual representations, 2017.

Buhrmester, V., Münch, D., and Arens, M. Analysis of explainers of black box deep neural networks for computer vision: A survey, 2019.

Carter, S., Armstrong, Z., Schubert, L., Johnson, I., and Olah, C. Activation atlas. *Distill*, 4(3):e15, 2019.

Hohman, F., Park, H., Robinson, C., and Chau, D. H. Summit: Scaling deep learning interpretability by visualizing activation and attribution summarizations, 2019.

Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., and Aila, T. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8110–8119, 2020.

Liu, Z., Luo, P., Wang, X., and Tang, X. Large-scale celebfaces attributes (celeba) dataset (2018). *URL http://mmlab. ie. cuhk. edu. hk/projects/CelebA. html*, 2018.

Mahendran and Vedaldi. Understanding deep image representations by inverting them. *CVPR*, 2015.

Mordvintsev, A., Olah, C., and Tyka, M. Inceptionism: Going deeper into neural networks. 2015.

Nguyen, A., Dosovitskiy, A., Yosinski, J., Brox, T., and Clune, J. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 29*, pp. 3387–3395. Curran Associates, Inc., 2016.

Nilsback, M.-E. and Zisserman, A. A visual vocabulary for flower classification. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pp. 1447–1454. IEEE, 2006.

Olah, C., Mordvintsev, A., and Schubert, L. Feature visualization. *Distill*, 2017. doi: 10.23915/distill.00007. https://distill.pub/2017/feature-visualization.

Olah, C., Satyanarayan, A., Johnson, I., Carter, S., Schubert, L., Ye, K., and Mordvintsev, A. The building blocks of interpretability. *Distill*, 2018. doi: 10.23915/distill.00010. https://distill.pub/2018/building-blocks.

Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001. https://distill.pub/2020/circuits/zoom-in.

Rathore, A., Chalapathi, N., Palande, S., and Wang, B. Topoact: Exploring the shape of activations in deep learning, 2019.

Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2): 336–359, Oct 2019. ISSN 1573-1405. doi: 10.1007/s11263-019-01228-7. URL http://dx.doi.org/10.1007/s11263-019-01228-7.

Simonyan, K., Vedaldi, A., and Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014. URL http://arxiv.org/abs/1312.6199.

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.

Zeiler and Fergus. Visualizing and understanding convolutional networks. *ECCV*, 2014.