# Understanding Learning Dynamics of Neural Representations via Feature Visualization at Scale

Chandana Kuntala<sup>12</sup> Carlos R. Ponce<sup>1</sup> Deepak Kumar Sharma<sup>2</sup> Binxu Wang<sup>13\*</sup> <sup>1</sup> Department of Neurobiology, Harvard Medical School, Boston, MA <sup>2</sup> Department of Information Technology, Indira Gandhi Technical University for Women, New Delhi, India <sup>3</sup> Kempner Institute for the Study of Natural and Artificial Intelligence, Harvard University, Allston, MA {chandana007btit20, deepaksharma}@igdtuw.ac.in {carlos, binxu\_wang}@hms.harvard.edu

Editors: Marco Fumero, Emanuele Rodolà, Clementine Domine, Francesco Locatello, Gintare Karolina Dziugaite, Mathilde Caron

## Abstract

How does feature learning happen during the training of a neural network? We developed an accelerated pipeline to synthesize maximally activating images ("prototypes") for hidden units in a parallel fashion. Through this, we were able to perform feature visualization at scale and to track the emergence and development of visual features across the training of neural networks. Using this technique, we studied the "developmental" process of features in a convolutional neural network trained from scratch using SimCLR with or without color jittering augmentation. After creating over one million prototypes with our method, tracking and comparing these visual signatures showed that the training with color-jitter augmentation led to constantly diversifying high-level features, while no color-jittering led to more diverse low-level features but less development of high-level features. These results illustrate how feature visualization can be used to understand hidden learning dynamics under different training objectives and data distribution.



Figure 1: **A.** Dual perspectives on neural representation. Left, neural vector space view of representation; Right, image space tuning landscape view of representation. **B.** Visual summary of our approach.

Proceedings of the I edition of the Workshop on Unifying Representations in Neural Models (UniReps 2023).

<sup>\*</sup>Corresponding author: Binxu Wang (binxu\_wang@hms.harvard.edu)

# 1 Introduction

In the biological and artificial neural systems, the neural representation of images is often analyzed in a multi-dimensional vector space [4, 17, 7, 20], formed by the activation of neurons. For example, in this vector space, the representations of different object categories form "object manifolds" which allows classification and few-shot learning [4, 6, 28]. An alternative perspective is to think about neural representation in their domain i.e. on the image manifold [35, 31]. From this perspective, the tuning of each neuron is a function (i.e. landscapes) on the manifold, with peaks and troughs. The peaks of the landscape correspond to images that highly activate these neurons. Note that, the axes of the neural vector space are the tuning functions of these neurons, thus the highly activating images could be regarded as the meaning of these axes. Through this paper, we call the activation maximizing images for each neuron a "prototype" [27]. Thus, obtaining the prototypes for all units in a neural network could provide a full basis set for understanding the representation of this network.

Feature visualization has been a prominent technique for finding and synthesizing prototypes in deep artificial neural networks [22, 23, 10], and the biological brain [26, 34, 12]. But normally, these methods were applied to one unit at a time, hard for application at scale.

In this work, we developed an accelerated pipeline to extract "prototypes" in a parallel fashion. Through this, we were able to apply feature visualization on a large scale, tracking the emergence and change of "prototypes" across the whole training process of neural networks – creating a visual signature for each network checkpoint. We leveraged this method to study the "development" of features in a convolutional neural network trained from scratch via self-supervised learning. The preliminary results illustrate how different training objectives and data distribution led to different "development" dynamics inside a computational visual hierarchy.

layer1	140	j 🎽	4	*	•	٠	Ŧ	٠
layer2	0	Y	•		*	<b>N</b>	1	A.
layer3	in	C.				Wei -		
layer4								

(a) SimCLR with color jittering (clrjit)

(b) SimCLR without color jittering (keepclr)

Figure 2: Example prototypes for networks trained with color jittering (clrjit) and without (keepclr)

# 2 Methods

## 2.1 Feature Visualization at Scale

Our method is based on [21, 33, 30], where feature visualization is performed within the latent space of a pretrained generative adversarial network (GAN) [8]. This GAN can be regarded as the natural image prior or the regularizer for the optimization, which counteracts the adversarial artifacts [22]. For each target unit, we optimize its activation using a hybrid of Covariance Matrix Adaptation Evolutionary Strategy (CMA-ES) [15] and gradient optimization: we performed 10s of CMA steps to search for an initialization vector that evoked non-zero activation in the unit, then we performed 100s of gradient ascent steps to optimize the vector to a peak to visualize the features. We implemented both CMA and Adam optimization in a more paralleled fashion, which enables feature visualization for each and every channel in a layer in a single run. This method increased our overall throughput by 33-fold (details in Sec. 6.2).



Figure 3: Dynamics of prototypes diversity during training.

## 2.2 Experiment Setup for Self-Supervised Learning

SimCLR [2] is a popular self-supervised learning algorithm for learning visual representation from a data distribution. This algorithm trains a neural network to associate different augmented views of the same image as similar representations, and those of different images as dissimilar ones. One key component of this method is the augmentation pipeline, which determines what type of transformation should the neural network be invariant to. It's well-known that with the same SimCLR objective, different augmentation pipelines led to dramatically different final performances [2]. The network mechanism underlying this difference is not well understood. Thus, we treated this as an example problem and dissected the feature learning dynamics through the lens of the prototype distribution.

Here, we had two training conditions with different augmentation pipelines and tested their effect on the development process of prototypes. 1) **Color jitter** (abbreviated as *clrjit*), the default augmentation pipeline of SimCLR; 2) **Keep Color** (*keepclr*), the same pipeline with color jittering and random grayscale augmentation disabled, which keeps the original color of the image. As the two conditions exposed the neural networks to different image statistics and pushed them with different objectives, we'd like to see if we can understand these differences better through the lens of prototype distribution.

For all our experiments, we used ResNet18 [16] as our neural network architecture and trained it with SimCLR on STL10 [5] dataset for 100 epochs.

Specifically, we completed three training runs of ResNet18 from scratch with random seeds 1,2,3 with color jittering and keep color augmentations; resulting in 6 training sequences of 101 epochs neural network checkpoints. For each checkpoint, we performed prototype extraction twice for each channel of every major layer (details in Sec. 6.1). Thus, all these prototypes can be indexed by [training condition, run number, evolution repeat, epoch number, layer, channel].

We evaluated the quality of their representations using the linear probe protocol (see Sec.6.1), namely fitting a linear classifier to see how well it classifies the test set images. The models trained with color jittering augmentation have far higher classification accuracy ( $70.0 \pm 0.3\%$ ) than the models without ( $49.8 \pm 0.3\%$ ) (Fig.5). This is consistent with the original observations of the importance of color augmentation in SimCLR (Fig.5 in [2]). From this perspective, the *clrjit* models have better feature representations for object classification. The focus of our analysis is to dissect the difference in representation quality and link it back to the development of prototypes.

# **3** Results

## 3.1 Visual difference of the prototypes between conditions

How do the learned features differ between the two training conditions? We first visually inspected the distribution of prototypes in each layer for the two conditions (Fig.2).

For the color jittering condition (Fig.2a), in layer 1, the prototypes masked with their respective receptive fields primarily captured patterns like black stripes on a white background  $(a1-1,1-3^2)$ , and solid colors like Prussian blue (a1-2), white/off-white, black, and partial cyan, red, and green shades. In the second layer, more square-circle figures like squircles (a2-1,2-3) were observed along with intricate patterns like thick lines and irregular line figures that somewhat resembled a cracked earth texture or an abstract glass painting texture (a2-2,2-4). In layer 3, the features became finer, as

<sup>&</sup>lt;sup>2</sup>1-based row-col index in the grid a of Fig.2, same convention below

high-frequency textures were observed (a3-1), along with black and white squircles (a3-2), rectangles (a3-3), and grids (a3-4). In layer 4, high-frequency textures (a4-1) were observed as well as distorted grid-like structures (a4-2). Few prototypes showed a gradient of colors resembling fur-like (a4-3) and watercolor textures (a4-4).

In contrast, in the keep color condition (Fig.2 b), in layer 1, a vibrant array of colors were observed including magenta / pink (b1-1), green (b1-2), red, blues, cyan, and yellow along with colorful stripes (b1-4). In layer 2, high-frequency textures (b2-4) were present along with colored gemstone-like shapes embedded in high-frequency textures (b2-1,b2-2,b2-3). In higher layers (layer 3 and layer 4), there were a significant number of high-frequency textures present mainly in the warm hues like oranges and reds (row3,4). These texture patterns are perceptually more similar to each other than the ones in color-jittering conditions.



#### 3.2 Developmental process of prototypes during training

Figure 4: Development of prototypes through training for color jittering (clrjit) condition, layer 3. Columns denote 0,10,20,... to 90 epoch; Rows denote Units 1, 2, 12, and 98 (0-based index).

So how do the neural networks arrive at these features? We visualized the prototypes of each training epochs as a row. For instance, for these example units in layer3 of a clrjit network (Fig.4, see Fig.9 for keepclr condition), each of these units goes through an initial stage of rapid erratic change, and then settles down to a primitive version of the final feature at around epoch 30, then elaborate this primitive form until the end. The latter half of epochs have more similarities between each other for each unit than the initial epochs however each of these individual units keeps diversifying with respect to each other throughout the training process as seen in Fig. 3.

## 3.3 Distance structure between prototypes

Next, we quantified our perception by computing the distance structure between prototypes to understand their distribution and dynamics during training. We computed the Mean Squared Error (MSE) and Cosine distance in both pixel space and the embedding space of some pre-trained networks (detail in Sec.6.4). Further, we computed the prototype similarity with the images masked by their functional receptive field mask (Sec.6.3) to focus on the central feature.

We quantified the **diversity of prototypes** during training: for each epoch, we computed the pairwise distance matrix between prototypes of all channels, and then computed the mean distance between prototype pairs (Fig. 3). We found salient differences between the two training conditions (color jitter, keep color), and consistencies between repeated training runs and prototype evolutions. Here we showed results with ResNet50 as our embedding model and MSE as the distance metric. For layer 1, the diversity dropped drastically in the first few epochs, and then grew to a stable level. In the end, keepclr condition led to more diverse prototypes than clrjit. For layer 2, after the initial drop of diversity, the prototypes diversify again, and keepclr condition led to slightly higher diversity. However, for layers 3 and 4, the keepclr condition increased prototype diversity early on and then they plateaued; in contrast, the clrjit condition led to a constant increase in prototype diversity without plateau. When the cosine distance is used instead of MSE (Fig.6), a shift is observed in the dynamics, albeit the diverging trend between conditions remained similar to the MSE result. The consistency of the color-jittering networks being at the top tends to demonstrate how these networks develop more

diverse prototypes through their evolutions. Further observations about the rate of change and the stability of prototype across re-evolution are noted in Sec. 7.2

This observation is intriguing. We interpreted it as follows, the color jittering augmentation constantly drives the network to find higher-level visual features to solve the instance classification task; while without color jittering, slightly more diverse lower-level features (layer1,2) suffice to solve the task. Intuitively, when SimCLR training doesn't randomly augment the color (keepclr), one simple way to find views of the same image is to look for similar color palettes. Thus, it's intuitive that the keepclr network needs to be more sensitive to image colors (Fig. 2). In comparison, with color jittering, the network cannot rely on color matching as a reliable strategy, and it needs to discover higher-level form consistencies, which may drive the diversification and learning of deeper layer features.

#### 4 Related Work

**Understanding self-supervised representation** Self-supervised learning (SSL) has been a popular feature learning paradigm in vision. In this paradigm, pre-training uses different objectives to learn features, which are used in the downstream tasks, with minor fine-tunings. But what is a good feature representation? Usually, these features were evaluated based on the performance of the downstream task. One open question is to understand and evaluate representations by themselves without using a downstream task. Many works analyzed the representation similarity of SSL networks and supervised networks [14]. [32] understand the success of contrastive learning through the alignment of positive pairs and the uniformity of representation on the hypersphere. More recently, [19] and [13] found certain intrinsic measures of representation e.g. effective dimensionality, intrinsic dimensionality, and cluster learnability, predict the in-distribution task performance well. On the other front, [1] used generative models to understand the "interpretation" of the same image by pre-trained networks to show their different biases, creating a more visual and intuitive explanation.

#### 5 Discussion

It has been noticed that the randomly initialized neural networks have lower dimensional representations, i.e. the activations of hidden units are more correlated across populations; and supervised training increased the dimensionality of the representation, and the increase is more prominent in deeper layers (Fig. H.1A, [9]). In the other perspective, the units become less correlated to each other during training, which is consistent with our finding that the prototypes of the units become more and more diverse during training.

Prototype diversity seems like a promising proxy for the richness of neural representation, however, it may not be the full story, these prototypes need to be related to the training and testing distribution of images in a meaningful way to be useful. Thus, one deep and open question is to elucidate the relationship between these prototypes and the training distribution of the network. Classic work of efficient coding [11, 24, 25] relates Gabor-like V1 receptive field (i.e. prototypes) to the natural image distribution and the sparse coding objective. Future works of similar flavor may be able to illuminate the relationship between the collection of prototypes of deeper layers and the training distribution of the network.

#### Acknowledgments and Disclosure of Funding

B.W. was funded by Victoria Quan fellowship at Harvard Medical School, and Kempner research fellowship. C.K. and C.R.P. were funded by grants NIH-NEI 1DP2EY035176, NSF CAREER 2143077, and Packard Foundation. We thank O2 Cluster at Harvard Medical School and Kempner Cluster at Harvard FAS for providing computational resources for this work.

#### References

- [1] Florian Bordes, Randall Balestriero, and Pascal Vincent. High fidelity visualization of what your self-supervised representation knows about. *arXiv preprint arXiv:2112.09164*, 2021.
- [2] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

- [3] Xuefeng Chen, Xiabi Liu, and Yunde Jia. Combining evolution strategy and gradient descent method for discriminative learning of bayesian classifiers. In *Proceedings of the 11th Annual conference on Genetic and evolutionary computation*, pages 507–514, 2009.
- [4] SueYeon Chung, Daniel D Lee, and Haim Sompolinsky. Classification and geometry of general perceptual manifolds. *Physical Review X*, 8(3):031003, 2018.
- [5] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- [6] Uri Cohen, SueYeon Chung, Daniel D Lee, and Haim Sompolinsky. Separability and geometry of object manifolds in deep neural networks. *Nature communications*, 11(1):746, 2020.
- [7] James J DiCarlo and David D Cox. Untangling invariant object recognition. *Trends in cognitive sciences*, 11(8):333–341, 2007.
- [8] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. *Advances in neural information processing systems*, 29, 2016.
- [9] Eric Elmoznino and Michael F Bonner. High-performing neural network models of visual cortex benefit from high latent dimensionality. *bioRxiv*, pages 2022–07, 2022.
- [10] Dumitru Erhan, Y. Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. *Technical Report, Université de Montréal*, 01 2009.
- [11] David J Field. Relations between the statistics of natural images and the response properties of cortical cells. *Josa a*, 4(12):2379–2394, 1987.
- [12] Peter Földiák. Stimulus optimisation in primary visual cortex. *Neurocomputing*, 38-40:1217–1222, 2001. Computational Neuroscience: Trends in Research 2001.
- [13] Quentin Garrido, Randall Balestriero, Laurent Najman, and Yann Lecun. Rankme: Assessing the downstream performance of pretrained self-supervised representations by their rank. In *International Conference on Machine Learning*, pages 10929–10974. PMLR, 2023.
- [14] Tom George Grigg, Dan Busbridge, Jason Ramapuram, and Russ Webb. Do self-supervised and supervised methods learn similar visual representations? arXiv preprint arXiv:2110.00528, 2021.
- [15] Nikolaus Hansen, Sibylle D Müller, and Petros Koumoutsakos. Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evolutionary computation*, 11(1):1–18, 2003.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [17] Chou P Hung, Gabriel Kreiman, Tomaso Poggio, and James J DiCarlo. Fast readout of object identity from macaque inferior temporal cortex. *Science*, 310(5749):863–866, 2005.
- [18] Ilya Loshchilov. A computationally efficient limited memory CMA-ES for large scale optimization. GECCO 2014 - Proceedings of the 2014 Genetic and Evolutionary Computation Conference, pages 397–404, 2014.
- [19] Yuchen Lu, Zhen Liu, Aristide Baratin, Romain Laroche, Aaron Courville, and Alessandro Sordoni. Expressiveness and learnability: A unifying view for evaluating self-supervised learning. *arXiv preprint arXiv:2206.01251*, 2022.
- [20] Najib J Majaj, Ha Hong, Ethan A Solomon, and James J DiCarlo. Simple learned weighted sums of inferior temporal neuronal firing rates accurately predict human core object recognition performance. *Journal of Neuroscience*, 35(39):13402–13418, 2015.
- [21] Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks, 2016.

- [22] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2(11):e7, 2017.
- [23] Chris Olah, Arvind Satyanarayan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine Ye, and Alexander Mordvintsev. The building blocks of interpretability. *Distill*, 3(3):e10, 2018.
- [24] Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.
- [25] Bruno A Olshausen and David J Field. Natural image statistics and efficient coding. Network: computation in neural systems, 7(2):333, 1996.
- [26] Carlos R. Ponce, Will Xiao, Peter F. Schade, Till S. Hartmann, Gabriel Kreiman, and Margaret S. Livingstone. Evolving Images for Visual Neurons Using a Deep Generative Network Reveals Coding Principles and Neuronal Preferences. *Cell*, 177(4):999–1009.e10, may 2019.
- [27] Olivia Rose, James Johnson, Binxu Wang, and Carlos R Ponce. Visual prototypes in the ventral stream are attuned to complexity and gaze behavior. *Nature communications*, 12(1):1–16, 2021.
- [28] Ben Sorscher, Surya Ganguli, and Haim Sompolinsky. Neural representational geometry underlies few-shot concept learning. *Proceedings of the National Academy of Sciences*, 119(43):e2200800119, 2022.
- [29] Igor Susmelj, Matthias Heller, Philipp Wirth, Jeremy Prescott, and Malte Ebner et al. Lightly. *GitHub. Note: https://github.com/lightly-ai/lightly*, 2020.
- [30] Binxu Wang and Carlos R. Ponce. High-performance evolutionary algorithms for online neuron control. In *Proceedings of the Genetic and Evolutionary Computation Conference*, GECCO '22, page 1308–1316, New York, NY, USA, 2022. Association for Computing Machinery.
- [31] Binxu Wang and Carlos R. Ponce. Tuning landscapes of the ventral stream. *Cell Reports*, nov 2022.
- [32] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR, 2020.
- [33] Will Xiao and Gabriel Kreiman. XDream: Finding preferred stimuli for visual neurons using generative networks and gradient-free optimization. *PLOS Computational Biology*, 16(6):e1007973, jun 2020.
- [34] Yukako Yamane, Eric T Carlson, Katherine C Bowman, Zhihong Wang, and Charles E Connor. A neural code for three-dimensional object shape in macaque inferotemporal cortex. *Nature Neuroscience*, 11(11):1352–1360, nov 2008.
- [35] Jacob A Zavatone-Veth, Sheng Yang, Julian A Rubinfien, and Cengiz Pehlevan. Neural networks learn to magnify areas near decision boundaries. *arXiv preprint arXiv:2301.11375*, 2023.

#### 6 Extended Method

#### 6.1 Details for Self-Supervised Learning

We used a popular self-supervised learning pipeline for training neural networks, SimCLR [2]. We used the implementation in lightly-ai [29].

**Augmentations** We tested two augmentation conditions, 1) the default stochastic augmentation pipeline with color jittering, 2) the same pipeline with color jittering and random grayscaling disabled (cj\_prob=0.0, random\_gray\_scale=0.0)

**Model Architecture** For the model backbone, we used the ResNet18 model [16], with 128d projection head.

**Dataset** For computational feasibility, we experimented with the SimCLR algorithm on the STL10 dataset [5] with the 96-pixel resolution, a classic testbed for self-supervised learning. Where the unlabeled training set has 100000 images, the training set has 500 images for each of the 10 classes and the testing set has 800 images for each of the 10 classes. The 10 classes are airplane, bird, car, cat, deer, dog, horse, monkey, ship, and truck, where 4 of those are animate man-made vehicles, and 6 of those are animate species.

**Training hyperparameters** For all models, we trained 100 epochs with Stochastic Gradient Descent with Cosine Annealing learning rate  $(lr = 6 \times 10^{-2}, momentum = 0.9, weight\_decay = 5 \times 10^{-4})$ 

**Evaluation** For evaluation, we used the linear probe protocol: We fixed all parameters of the CNN and used it to map images from the training and test set to feature vectors. Here no image augmentation was used, only RGB value normalization. Then we fit a linear classifier based on the training set features and evaluated the classifiers on the test set features. We used three ways to fit the linear classifier: Logistic regression (LogisticRegression from sklearn), Linear Support Vector Classifier (LinearSVC from sklearn), and gradient descent (Adam) on Cross Entropy Loss.

#### 6.2 Scalable methods for synthesizing prototypes

This work requires us to synthesize a huge amount of highly activating images (prototypes), so we developed a more scalable way to synthesize them efficiently.

Specifically, we parallelized the hybrid of CMA-ES [15, 18] and gradient optimization [3] to optimize the images for each channel in a layer independently. Our rationale is as follows, gradient-based feature visualization is efficient and parallelizable: a batch of images can be sent into the network, with the activations extracted, then the optimization objective can be set as the sum of activations of the target unit for each image. Taking the gradient of this summed objective to the batch of samples is equivalent to optimizing each image independently.

$$\mathcal{L} = \sum_{i=1}^{B} f_{(i)}(G(z_i)) \tag{1}$$

$$\nabla_{z_i} \mathcal{L} = \nabla_{z_i} f_{(i)}(G(z_i)) \tag{2}$$

However, gradient-based optimization suffers from vanishing gradient problem, i.e. when the initial image didn't evoke any response in the target unit, then for ReLU activation, the gradient is zero, and optimization stopped for this unit. To mitigate this issue, we used Evolutionary algorithms i.e. CMA-ES to search for proper initialization for each unit. CMA-ES proposes a Gaussian distributed population of vectors in each iteration, and then adapts this distribution to optimize response. This algorithm is surprisingly good at navigating the landscape of activations and usually finds initializations with non-zero activation with 10s of iterations. Then we used the latent vector and image with maximal activation as the starting point to perform the gradient-based optimization.

As a concrete example, we need to synthesize,  $101 \times (64 + 128 + 256 + 512) = 96960$  prototypes for all channels in all major layers for each epoch of a training run in ResNet18. This will take around

269hrs on a single GPU using the previous non-parallelized CMA-ES algorithm per channel pipeline. Using our current method, it takes only 8hrs on a single GPU, which is a 33 times speed up. Using this method, we synthesized over 1 million prototypes in a reasonable time.

#### 6.3 Methods for computing receptive field of units

We used gradient-based receptive field mapping for hidden units. We denote the hidden unit as  $f : \mathbb{R}^{H \times W \times C} \to \mathbb{R}, \mathbf{x} \mapsto r$ . We sample random white noise patterns  $\mathbf{x}$  with image shape and then send the noise pattern through the neural network, and compute the gradient of f.

$$M_{raw} = \mathbb{E}_{\mathbf{x} \sim Unif[0,1]^{H \times W \times C}} \nabla_{\mathbf{x}} f(\mathbf{x})$$
(3)

We averaged this gradient across 200 samples of x and then took the sum of squares over the channel C dimension as a  $H \times W$  spatial mask. Finally, we fit this mask with a 2D Gaussian function, and the fitted mask is denoted as  $M_{fit}$ . This Gaussian mask was called the receptive field mask and was used to mask the prototypes and highlight the central features.

#### 6.4 Methods for comparing image similarity

To compare the image similarity between prototypes generated by various networks, we computed similarity metric (MSE and Cosine) in the pixel space and the embedding space. For both spaces, we first masked the image with the Gaussian receptive field mask for the target unit. The method to calculate the receptive field size and mask is mentioned in subsection 6.3. In the pixel space, the pixels from the masked images were vectorized and used to compute the distance matrices. In the embedding space, we mapped the masked prototypes to vector embeddings with ten pre-trained neural networks and computed similarity metrics there. Specifically, we sent the masked prototypes into the ten pre-trained convolutional neural networks (CNNs): ResNet50, ResNet101, ResNet152, InceptionV3, VGG16, VGG19, DenseNet121, DenseNet169, DenseNet201 and MobileNetV2. The activations from the last fully connected layer of the network, like the 'avgpool' layer from ResNet50 were chosen for each network to extract the activations. These were then used to compute the distance matrices.

Note that, we used the receptive field mask to mask the prototype image, before computing their similarity. Because of this, the similarity or distance value might not be comparable across layers, as the units in different layers have different sizes of receptive field masks, thus the different masked prototypes will have different amounts of black backgrounds.

## 7 Extended Results

#### 7.1 Linear Probe Evaluation of Learned Features



Figure 5: Linear feature evaluation of representation during SimCLR training. A. Test set accuracy during SimCLR training, Each panel corresponds to one way of fitting the linear readout layer: Logistic regression, Linear Support Vector Classifier, and gradient optimization (Adam) on Cross Entropy. B. Final test set accuracy for the self-supervised models, separated by whether they used color jittering augmentation.

#### 7.2 Additional observations on the distance structure of prototypes

Aside from prototype diversity, we also examined the **rate of prototype change** during training: we computed the distances between the prototype of the same unit c during neighboring epochs Ep and Ep + 1. We averaged the distance for channels in each layer and showed it across epochs and networks (Fig.7). We can see, that all the networks experienced a transient peak at the first step, showing the drastic change of representation between randomly initialized network to network after one training epoch. Here we also saw clrjit networks experienced a higher rate of change of prototypes in layers 3 and 4, which might be the cause of their higher diversity.

Finally, we examined the **consistency of prototype across repeated Evolution**. This is related to the overall geometry of the landscape, i.e. how multimodal the tuning of the unit is. We computed the distances between the prototype of the same channel *c* for the two extractions. We averaged the distance for channels in each layer and showed it across epochs and networks (Fig.8). Generally, the distance between prototypes of the same channel (repeated evolution) is smaller than the distance between prototypes of different channels (cf. Fig.3). Intriguingly this distance between re-evolved prototypes is increasing through the training process, especially for the deeper layers of the models trained with color jittering. This highlights that during training, for the same unit, repeated evolution led to increasingly different prototypes — namely the tuning functions of units are becoming more and more multimodal. This increase in multimodality may also benefit the final quality of representation.



Figure 6: **Dynamics of prototypes diversity during training**. Cosine distance metric, resnet50 embedding



Figure 8: Similarity of prototype between repeated evolution.

7.3 Developmental process of prototypes during training (keepclr)



Figure 9: Development of prototypes through training without color jittering (keepclr) condition, layer 3. Columns correspond to 0,10,20,... to 90 epoch; Rows correspond to Units 71, 93, 211, and, 248 (0-based index)