

ConvNeXt V2: Co-designing and Scaling ConvNets with Masked Autoencoders

Sanghyun Woo^{1*} Shoubhik Debnath² Ronghang Hu²
 Xinlei Chen² Zhuang Liu² In So Kweon¹ Saining Xie^{3†}
¹KAIST ² Meta AI, FAIR ³New York University

Code: <https://github.com/facebookresearch/ConvNeXt-V2>

Abstract

Driven by improved architectures and better representation learning frameworks, the field of visual recognition has enjoyed rapid modernization and performance boost in the early 2020s. For example, modern ConvNets, represented by ConvNeXt [33], have demonstrated strong performance in various scenarios. While these models were originally designed for supervised learning with ImageNet labels, they can also potentially benefit from self-supervised learning techniques such as masked autoencoders (MAE) [14]. However, we found that simply combining these two approaches leads to subpar performance. In this paper, we propose a fully convolutional masked autoencoder framework and a new Global Response Normalization (GRN) layer that can be added to the ConvNeXt architecture to enhance inter-channel feature competition. This co-design of self-supervised learning techniques and architectural improvement results in a new model family called ConvNeXt V2, which significantly improves the performance of pure ConvNets on various recognition benchmarks, including ImageNet classification, COCO detection, and ADE20K segmentation. We also provide pre-trained ConvNeXt V2 models of various sizes, ranging from an efficient 3.7M-parameter Atto model with 76.7% top-1 accuracy on ImageNet, to a 650M Huge model that achieves a state-of-the-art 88.9% accuracy using only public training data.

1. Introduction

Building on research breakthroughs in earlier decades [16, 25, 28, 37, 44], the field of visual recognition has ushered in a new era of large-scale visual representation learning. Pre-trained, large-scale vision models have become essential tools for feature learning and enabling a wide range of vision applications. The performance of a visual representation learning system is largely influenced by three main factors: the neural network architecture chosen, the method

* Work done during an internship at FAIR.

† Corresponding author.

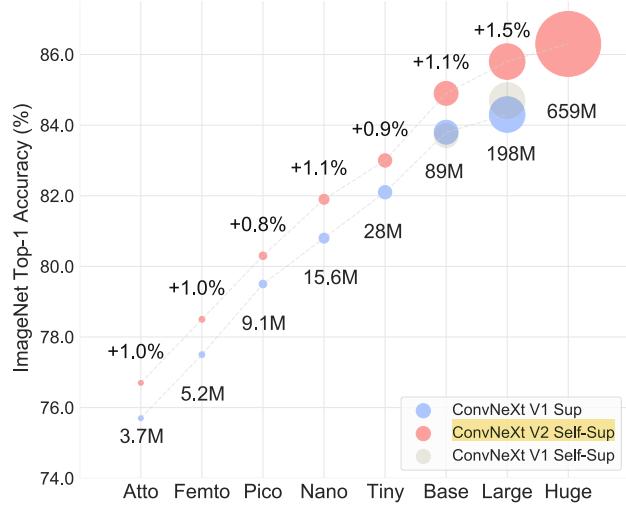


Figure 1. **ConvNeXt V2 model scaling.** The ConvNeXt V2 model, which has been pre-trained using our fully convolutional masked autoencoder framework, performs significantly better than the previous version across a wide range of model sizes.

used for training the network, and the data used for training. In the field of visual recognition, progress in each of these areas contributes to overall improvements in performance.

Innovation in neural network architecture design has consistently played a major role in the field of representation learning. Convolutional neural network architectures (ConvNets) [16, 25, 28] have had a significant impact on computer vision research by allowing for the use of generic feature learning methods for a variety of visual recognition tasks [10, 15], rather than relying on manual feature engineering. In recent years, the transformer architecture [44], originally developed for natural language processing, has also gained popularity due to its strong scaling behavior with respect to model and dataset size [7]. More recently, ConvNeXt [33] architecture has modernized traditional ConvNets and demonstrated that pure convolutional models could also be scalable architectures. However, the most common method for exploring the design space for neural network architectures is still through benchmarking *supervised learning* performance on ImageNet.

In a separate line of research, the focus of visual representation learning has been shifting from supervised learning with labels to self-supervised pre-training with pre-text objectives. Among many different self-supervised algorithms, masked autoencoders (MAE) [14] have recently brought success in masked language modeling to the vision domain and quickly become a popular approach for visual representation learning. However, a common practice in self-supervised learning is to use a *predetermined* architecture designed for supervised learning, and assume the design is fixed. For instance, MAE was developed using the vision transformer [7] architecture.

It is possible to combine the design elements of architectures and self-supervised learning frameworks, but doing so may present challenges when using ConvNeXt with masked autoencoders. One issue is that MAE has a specific encode-decoder design that is optimized for the sequence processing capabilities of transformers, which allows the compute-heavy encoder to focus on visible patches and thus reduce the pre-training cost. This design may not be compatible with standard ConvNets, which use dense sliding windows. Additionally, if the relationship between the architecture and the training objective is not taken into consideration, it may be unclear whether optimal performance can be achieved. In fact, previous research has shown that training ConvNets with mask-based self-supervised learning can be difficult [24], and empirical evidence suggests that transformers and ConvNets may have different feature learning behaviors that can affect representation quality.

To this end, we propose to *co-design* the network architecture and the masked autoencoder under the same framework, with the aim of making mask-based self-supervised learning effective for ConvNeXt models and achieving results similar to those obtained using transformers.

In designing the masked autoencoder, we treat the masked input as a set of sparse patches and use sparse convolutions [12] to process only the visible parts. The idea is inspired by the use of sparse convolutions in processing large-scale 3D point clouds [5, 52]. In practice, we can implement ConvNeXt with sparse convolutions, and at fine-tuning, the weights are converted back to standard, dense layers without requiring special handling. To further improve the pre-training efficiency, we replace the transformer decoder with a single ConvNeXt block, making the entire design fully convolutional. We have observed mixed results with these changes: the learned features are useful and improve upon the baseline results, but the fine-tuning performance is still not as good as the transformer-based model.

We then conduct a feature space analysis of different training configurations for ConvNeXt. We identify a potential issue of feature collapse at the MLP layer when training ConvNeXt directly on masked input. To address this issue, we propose adding a Global Response Normalization layer

to enhance inter-channel feature competition. This change is most effective when the model is pre-trained with masked autoencoders, suggesting that reusing a fixed architecture design from supervised learning may be suboptimal.

In summary, we introduce ConvNeXt V2 which demonstrates improved performance when used in conjunction with masked autoencoders. We have found that this model significantly improves the performance of pure ConvNets across various downstream tasks, including ImageNet classification [37], COCO object detection [30] and ADE20K segmentation [55]. The ConvNeXt V2 models can be used in a variety of compute regimes and includes models of varying complexity: from an efficient 3.7M-parameter *Atto* model that achieves 76.7% top-1 accuracy on ImageNet, to a 650M *Huge* model that reaches a state-of-the-art 88.9% accuracy when using IN-22K labels.

2. Related Work

ConvNets. The design of ConvNets, which were first introduced in the 1980s [27] and trained using back-propagation, has undergone numerous improvements in terms of optimization, accuracy, and efficiency over the years [17, 18, 21, 25, 35, 38, 39, 51]. These innovations have mainly been discovered through the use of supervised training on the ImageNet dataset. In recent years, some efforts have been made to perform architecture search using self-supervised pre-text tasks such as rotation prediction and colorization, as in the case of UnNAS [31]. Recently, ConvNeXt [33] conducted a comprehensive review of the design space and demonstrated pure ConvNets can be as scalable as the vision transformers [7, 32], which have become the dominant architecture in many applications. ConvNeXt has particularly excelled in scenarios requiring lower complexity [4, 46, 47]. Our ConvNeXt V2 model, which is powered by self-supervised learning, provides a simple way to upgrade existing models and achieve a significant boost in performance across a wide range of use cases.

Masked Autoencoders. Masked image modeling, represented by masked autoencoders [14], is one of the latest self-supervised learning strategies. As a neural network pre-training framework, masked autoencoders have shown a broad impact on visual recognition. However, original masked autoencoders are not directly applicable to ConvNets due to their asymmetric encoder-decoder design. Alternative frameworks such as [2, 53] have attempted to adapt the approach for use with ConvNets, but with mixed results. MCMAE [9] uses a few convolutional blocks as input tokenizers. To the best of our knowledge, there are no pre-trained models that show self-supervised learning can improve upon the best ConvNeXt supervised results.

3. Fully Convolutional Masked Autoencoder

Our approach is conceptually simple and runs in a fully convolutional manner. The learning signals are generated by randomly *masking* the raw input visuals with a high masking ratio and letting the model *predict* the missing parts given the remaining context. Our framework is illustrated in Figure 2, and we will now describe its main components in more detail.

Masking. We use a random masking strategy with a masking ratio of 0.6. As the convolutional model has a hierarchical design, where the features are downsampled in different stages, the mask is generated in the last stage and upsampled recursively up to the finest resolution. To implement this in practice, we randomly remove 60% of the 32×32 patches from the original input image. We use minimal data augmentation, only including random resized cropping.

Encoder design. We use ConvNeXt [33] model as the encoder in our approach. One challenge in making masked image modeling effective is preventing the model from learning shortcuts that allow it to copy and paste information from the masked regions. This is relatively easy to prevent in transformer-based models, which can leave the visible patches as the only input to the encoder. However, it is more difficult to achieve this with ConvNets, as the 2D image structure must be preserved. While naive solutions involve introducing learnable masked tokens in the input side [2, 53], these approaches decrease the efficiency of pre-training and result in a train and test time inconsistency, as there are no mask tokens at test time. This becomes especially problematic when the masking ratio is high.

To tackle this issue, our new insight is to view the masked image from a “sparse data perspective”, which was inspired by learning on sparse point clouds in 3D tasks [5, 52]. Our key observation is that the masked image can be represented as a 2D sparse array of pixels. Based on this insight, it is natural to incorporate sparse convolution into our framework to facilitate pre-training of the masked autoencoder. In practice, during pre-training, we propose to convert the standard convolution layer in the encoder with the submanifold sparse convolution, which enables the model to operate *only* on the visible data points [5, 11, 12]. We note that the sparse convolution layers can be converted back to standard convolution at the fine-tuning stage without requiring additional handling. As an alternative, it is also possible to apply a binary masking operation before and after the dense convolution operation. This operation has numerically the same effect as sparse convolutions, is theoretically more computationally intensive, but can be more friendly on AI accelerators like TPU.

Decoder design. We use a lightweight, plain ConvNeXt block as the decoder. This forms an asymmetric encoder-decoder architecture overall, as the encoder is heavier and

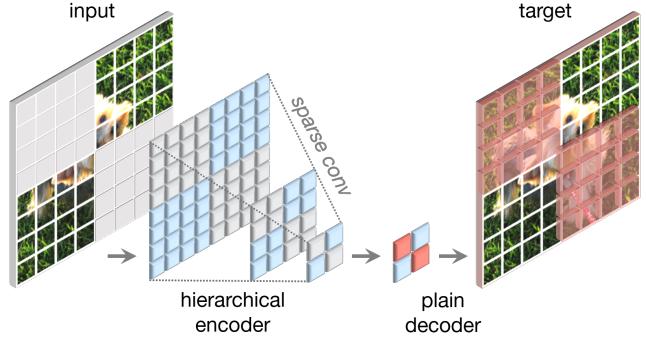


Figure 2. **Our FCMAE framework.** We introduce a fully convolutional masked autoencoder (FCMAE). It consists of a sparse convolution-based ConvNeXt encoder and a lightweight ConvNeXt block decoder. Overall, the architecture of our autoencoder is asymmetric. The encoder processes only the visible pixels, and the decoder reconstructs the image using the encoded pixels and mask tokens. The loss is calculated only on the masked region.

has a hierarchy. We also considered more complex decoders such as hierarchical decoders [29, 36] or transformers [7, 14], but the simpler single ConvNeXt block decoder performed well in terms of fine-tuning accuracy and reduced pre-training time considerably, demonstrated in Table 1. We set the dimension of the decoder to 512.

Reconstruction target. We compute the mean squared error (MSE) between the reconstructed and target images. Similar to MAE [14], the target is a patch-wise normalized image of the original input, and the loss is applied only on the masked patches.

FCMAE. We now present a Fully Convolutional Masked AutoEncoder (FCMAE) by combining the proposals described above. To evaluate the effectiveness of this framework, we use the ConvNeXt-Base model as the encoder and conduct a series of ablation studies. Throughout the paper, we focus on the end-to-end fine-tuning performance because of its practical relevance in transfer learning, and use that to assess the quality of the learned representation.

We pre-train and fine-tune using the ImageNet-1K (IN-1K) dataset for 800 and 100 epochs, respectively, and report the top-1 IN-1K validation accuracy for a single 224×224 center crop. Additional details about the experimental setup can be found in the appendix.

To understand the impact of using sparse convolution in our FCMAE framework, we first investigate how it affects the quality of the learned representation during masked image pre-training. Our empirical findings show that it is essential to prevent information leakage from the masked region in order to achieve good results.

	w/o Sparse conv.	w/ Sparse conv.
	79.3	83.7



dec. type	ft	hours	speedup
UNet w/ skip	83.7	12.9	-
UNet w/o skip	83.5	12.9	-
Transformer [14]	83.4	8.5	1.5x
ConvNeXt block	83.7	7.7	1.7x

(a) **Decoder design.** A simple convolutional block outperforms more complex decoder designs.

blocks	ft
1	83.7
2	83.5
4	83.7
8	83.6
12	83.3

(b) **Decoder depth.** A single block yields competitive fine-tuning performance.

dim	ft
128	83.5
256	83.7
512	83.7
768	83.6
1024	83.5

(c) **Decoder width.** A decoder width of 256 or 512 achieves the best performance.

Table 1. **MAE decoder ablation experiments** with ConvNeXt-Base on ImageNet-1K. We report fine-tuning (ft) accuracy (%). The pre-training schedule is 800 epochs. In the decoder design exploration, the wall-clock time is benchmarked on a 256-core TPU-v3 pod using JAX. The speedup is relative to the UNet decoder baseline. Our final design choices employed in the paper are marked in gray.

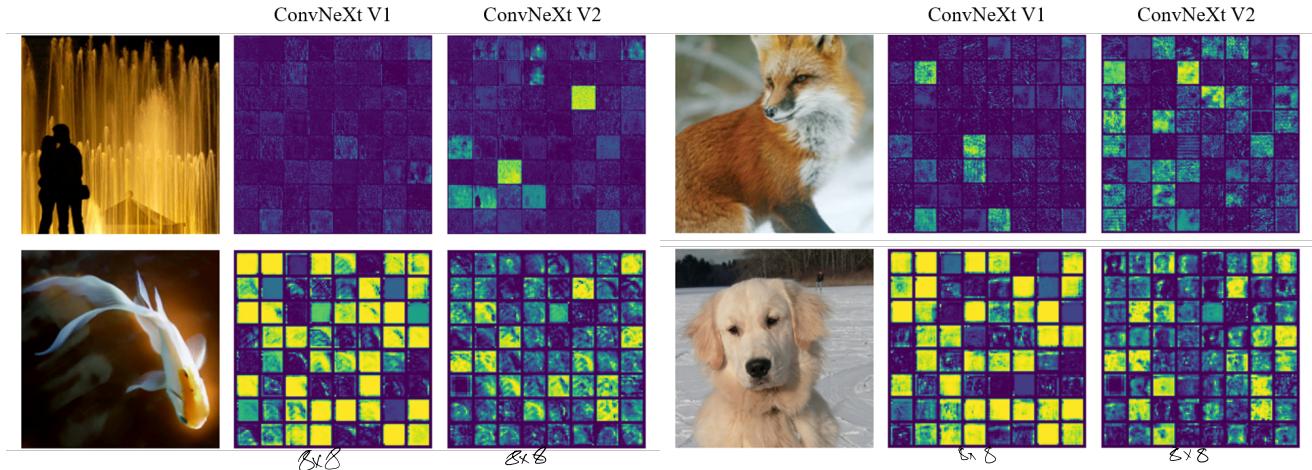


Figure 3. **Feature activation visualization.** We visualize the activation map for each feature channel in small squares. For clarity, we display 64 channels in each visualization. The ConvNeXt V1 model suffers from a feature collapse issue, which is characterized by the presence of redundant activations (dead or saturated neurons) across channels. To fix this problem, we introduce a new method to promote feature diversity during training: the global response normalization (GRN) layer. This technique is applied to high-dimensional features in every block, leading to the development of the ConvNeXt V2 architecture.

Next, we compare our self-supervised approach to supervised learning. Specifically, we obtain two baseline experimental results: the supervised 100 epoch baseline using the same recipe and the 300 epoch supervised training baseline provided in the original ConvNeXt paper [33]. We find that our FCMAE pre-training provides better initialization than the random baseline (*i.e.*, $82.7 \rightarrow 83.7$), but it still needs to catch up to the best performance obtained in the original supervised setup.

Sup, 100ep	Sup, 300ep. [33]	FCMAE
82.7	83.8	83.7

This is in contrast to the recent success of masked image modeling using transformer-based models [2, 14, 53], where the pre-trained models significantly outperform the supervised counterparts. This motivates us to investigate the unique challenges faced by the ConvNeXt encoder during masked autoencoder pre-training, which we discuss next.

4. Global Response Normalization

In this section, we introduce a new Global Response Normalization (GRN) technique to make FCMAE pre-training more effective in conjunction with the ConvNeXt architecture. We first motivate our approach through both qualitative and quantitative feature analyses.

Feature collapse. To gain more insight into the learning behavior, we first perform qualitative analysis in the feature space. We visualize the activations of a FCMAE pre-trained ConvNeXt-Base model and notice an intriguing “feature collapse” phenomenon: there are many dead or saturated feature maps and the activation becomes redundant across channels. We show some of the visualizations in Figure 3. This behavior was mainly observed in the dimension-expansion MLP layers in a ConvNeXt block [33].

Feature cosine distance analysis. To further validate our observation quantitatively, we perform a feature cosine distance analysis. Given an activation tensor $X \in R^{H \times W \times C}$,

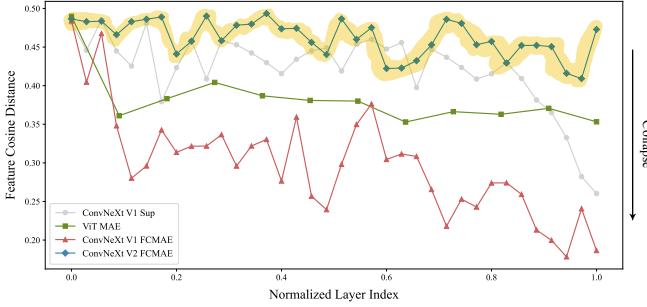


Figure 4. Feature cosine distance analysis. As the number of total layers varies for different architectures, we plot the distance values against the normalized layer indexes. We observe that the ConvNeXt V1 FCMAE pre-trained model exhibits severe feature collapse behavior. The supervised model also shows a reduction in feature diversity, but only in the final layers. This decrease in diversity in the supervised model is likely due to the use of the cross-entropy loss, which encourages the model to focus on class-discriminative features while suppressing the others.

$X_i \in \mathcal{R}^{H \times W}$ is the feature map of the i -th channel. We reshape it as a HW dimensional vector and compute the average pair-wise cosine distance across the channels by $\frac{1}{C^2} \sum_i^C \sum_j^C \frac{1 - \cos(X_i, X_j)}{2}$. A higher distance value indicates more diverse features, while a lower value indicates feature redundancy.

To perform this analysis, we randomly select 1,000 images from different classes in the ImageNet-1K validation set and extract the high-dimensional features from each layer of different models, including the FCMAE models, the ConvNeXt supervised model [33] and the MAE pre-trained ViT model [14]. We then compute the distance per layer for each image and average the values across all images. The results are plotted in Figure 4. The FCMAE pre-trained ConvNeXt model exhibits a clear tendency towards feature collapse, consistent with our observations from the previous activation visualizations. This motivates us to consider ways to diversify the features during learning and prevent feature collapse.

Approach. There are many mechanisms in the brain that promote neuron diversity. For example, lateral inhibition [3, 13] can help to sharpen the response of the activated neuron and increase the contrast and selectivity of individual neurons to the stimulus while also increasing the diversity of responses across the population of neurons. In deep learning, this form of lateral inhibition can be implemented by response normalization [26]. In this work, we introduce a new response normalization layer called global response normalization (GRN), which aims to increase the contrast and selectivity of channels. Given an input feature, $X \in \mathcal{R}^{H \times W \times C}$, the proposed GRN unit consists of three steps: 1) global feature aggregation, 2) feature normalization, and 3) feature calibration.

Algorithm 1 Pseudocode of GRN in a PyTorch-like style.

```
# gamma, beta: learnable affine transform parameters
# X: input of shape (N,H,W,C)

gx = torch.norm(X, p=2, dim=(1,2), keepdim=True)
nx = gx / (gx.mean(dim=-1, keepdim=True)+1e-6)
return gamma * (X * nx) + beta + X
```

First, we aggregate a spatial feature map X_i into a vector gx with a global function $\mathcal{G}(\cdot)$:

$$\mathcal{G}(X) := X \in \mathcal{R}^{H \times W \times C} \rightarrow gx \in \mathcal{R}^C. \quad (1)$$

This can be viewed as a simple pooling layer. We experimented with different functions in Table 2a. Interestingly, global average pooling, a widely used feature aggregator [19, 48], did not perform well in our case. Instead, we found that using norm-based feature aggregation, specifically, using L2-norm, resulted in better performance. This gives us a set of aggregated values $\mathcal{G}(X) = gx = \{\|X_1\|, \|X_2\|, \dots, \|X_C\|\} \in \mathcal{R}^C$ where $\mathcal{G}(X)_i = \|X_i\|$ is a scalar that aggregates the statistics of the i -th channel.

Next, we apply a response normalization function $\mathcal{N}(\cdot)$ to the aggregated values. Concretely, we use a standard divisive normalization as follows,

$$\mathcal{N}(\|X_i\|) := \|X_i\| \in \mathcal{R} \rightarrow \frac{\|X_i\|}{\sum_{j=1, \dots, C} \|X_j\|} \in \mathcal{R}, \quad (2)$$

where $\|X_i\|$ is the L2-norm of the i -th channel.¹ Intuitively, for the i -th channel, Eqn. 2 computes its *relative importance* compared to all the other channels. Similar to other forms of normalization [23, 26, 44], this step creates a *feature competition* across channels by mutual inhibition. In Table 2b, we also examine the use of other normalization functions and find that the simple divisive normalization works best, though standardization $(\|X_i\| - \mu)/\sigma$ yields similar results when applied to the same L2-norm aggregated values.

Finally, we calibrate the original input responses using the computed feature normalization scores:

$$X_i = X_i * \mathcal{N}(\mathcal{G}(X)_i) \in \mathcal{R}^{H \times W} \quad (3)$$

The core GRN unit is very easy to implement, requiring only three lines of code, and has no learnable parameters. The pseudo-code for the GRN unit is in Algorithm 1.

To ease optimization, we add two additional learnable parameters, γ and β , and initialize them to zero. We also add a residual connection between the input and output of the GRN layer. The resulting final GRN block is $X_i = \gamma * X_i * \mathcal{N}(\mathcal{G}(X)_i) + \beta + X_i$. This setup allows a GRN layer

¹To account for the increased number of channels at deeper layers, in practice, we also scale the normalized value by the channel count C .

<table border="1"> <thead> <tr> <th>case</th> <th>ft</th> </tr> </thead> <tbody> <tr> <td>g.avg.</td> <td>83.7</td> </tr> <tr> <td>L1</td> <td>84.3</td> </tr> <tr> <td>L2</td> <td>84.6</td> </tr> </tbody> </table>	case	ft	g.avg.	83.7	L1	84.3	L2	84.6	<table border="1"> <thead> <tr> <th>case</th> <th>ft</th> </tr> </thead> <tbody> <tr> <td>$(\ X_i\ - \mu)/\sigma$</td> <td>84.5</td> </tr> <tr> <td>$1/\sum \ X_i\$</td> <td>83.8</td> </tr> <tr> <td>$\ X_i\ /\sum \ X_i\$</td> <td>84.6</td> </tr> </tbody> </table>	case	ft	$(\ X_i\ - \mu)/\sigma$	84.5	$1/\sum \ X_i\ $	83.8	$\ X_i\ /\sum \ X_i\ $	84.6	<table border="1"> <thead> <tr> <th>case</th> <th>ft</th> </tr> </thead> <tbody> <tr> <td>w/o skip</td> <td>84.0</td> </tr> <tr> <td>w/ skip</td> <td>84.6</td> </tr> </tbody> </table>	case	ft	w/o skip	84.0	w/ skip	84.6															
case	ft																																						
g.avg.	83.7																																						
L1	84.3																																						
L2	84.6																																						
case	ft																																						
$(\ X_i\ - \mu)/\sigma$	84.5																																						
$1/\sum \ X_i\ $	83.8																																						
$\ X_i\ /\sum \ X_i\ $	84.6																																						
case	ft																																						
w/o skip	84.0																																						
w/ skip	84.6																																						
(a) Global aggregation $G(\cdot)$. L2 Norm-based aggregation function produces the best result.																																							
(b) Normalization operator , $N(\cdot)$. Divisive normalization is an effective channel importance calibrator.																																							
(c) Residual connection helps with GRN optimization and leads to better performance.																																							
<table border="1"> <thead> <tr> <th>case</th> <th>ft</th> </tr> </thead> <tbody> <tr> <td>Baseline</td> <td>83.7</td> </tr> <tr> <td>LRN [26]</td> <td>83.2</td> </tr> <tr> <td>BN [22]</td> <td>80.5</td> </tr> <tr> <td>LN [1]</td> <td>83.8</td> </tr> <tr> <td>GRN</td> <td>84.6</td> </tr> </tbody> </table>	case	ft	Baseline	83.7	LRN [26]	83.2	BN [22]	80.5	LN [1]	83.8	GRN	84.6	<table border="1"> <thead> <tr> <th>case</th> <th>ft</th> <th>#param</th> </tr> </thead> <tbody> <tr> <td>Baseline</td> <td>83.7</td> <td>89M</td> </tr> <tr> <td>SE [19]</td> <td>84.4</td> <td>109M</td> </tr> <tr> <td>CBAM [48]</td> <td>84.5</td> <td>109M</td> </tr> <tr> <td>GRN</td> <td>84.6</td> <td>89M</td> </tr> </tbody> </table>	case	ft	#param	Baseline	83.7	89M	SE [19]	84.4	109M	CBAM [48]	84.5	109M	GRN	84.6	89M	<table border="1"> <thead> <tr> <th>case</th> <th>ft</th> </tr> </thead> <tbody> <tr> <td>Baseline</td> <td>83.7</td> </tr> <tr> <td>drop at ft.</td> <td>78.8</td> </tr> <tr> <td>add at ft.</td> <td>80.6</td> </tr> <tr> <td>both</td> <td>84.6</td> </tr> </tbody> </table>	case	ft	Baseline	83.7	drop at ft.	78.8	add at ft.	80.6	both	84.6
case	ft																																						
Baseline	83.7																																						
LRN [26]	83.2																																						
BN [22]	80.5																																						
LN [1]	83.8																																						
GRN	84.6																																						
case	ft	#param																																					
Baseline	83.7	89M																																					
SE [19]	84.4	109M																																					
CBAM [48]	84.5	109M																																					
GRN	84.6	89M																																					
case	ft																																						
Baseline	83.7																																						
drop at ft.	78.8																																						
add at ft.	80.6																																						
both	84.6																																						
(d) Feature normalization . GRN outperforms other normalizations through global contrasting.																																							
(e) Feature re-weighting . GRN does effective and efficient feature re-weighting without parameter overhead.																																							
(f) GRN in pre-training/fine-tuning . To be effective, GRN should be used in both stages.																																							

Table 2. **GRN ablations** with ConvNeXt-Base. We report fine-tuning accuracy on ImageNet-1K. Our final proposal is marked in gray .

to initially perform an identity function and gradually adapt during training. The importance of residual connection is demonstrated in Table 2c.

ConvNeXt V2. We incorporate the GRN layer into the original ConvNeXt block, as illustrated in Figure 5. We empirically found that LayerScale [41] becomes unnecessary when GRN is applied and can be removed. Using this new block design, we create various models with varying efficiency and capacity, which we refer to as the ConvNeXt V2 model family. These models range from lightweight (e.g. Atto [46]) to compute-intensive (e.g. Huge) ones. Detailed model configurations can be found in the appendix.

Impact of GRN. We now pre-train ConvNeXt V2 using the FCMAE framework and evaluate the impact of GRN. From visualization in Figure 3 and cosine distance analysis in Figure 4, we can observe that ConvNeXt V2 effectively mitigates the feature collapse issue. The cosine distance values are consistently high, indicating that feature diversity is maintained across layers. This behavior is similar to that of the MAE pre-trained ViT model [14]. Overall, this suggests that ConvNeXt V2 learning behavior can resemble ViT, under a similar masked image pre-training framework.

Next, we evaluate the fine-tuning performance.

V1 + Sup, 300ep.	V1 + FCMAE	V2 + FCMAE
83.8	83.7	84.6

When equipped with GRN, the FCMAE pre-trained model can significantly outperform the 300 epoch supervised counterpart. GRN improves the representation quality by enhancing the feature diversity, which was absent in the V1 model but has proven crucial for masked-based pre-training. Note this improvement is achieved *without* adding additional parameter overhead or increased FLOPS.²

Relation to feature normalization methods. Can other normalization layers [1, 22, 26, 43, 49] perform as well as the

²The additional affine parameters γ/β are negligible.

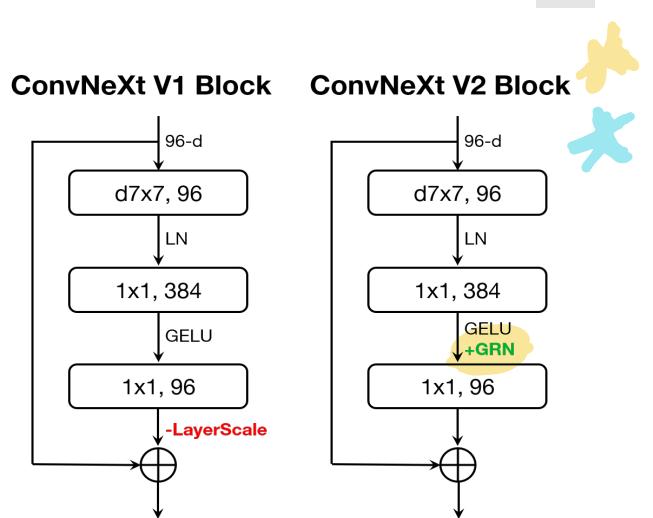


Figure 5. **ConvNeXt Block Designs.** In ConvNeXt V2, we add the GRN layer after the dimension-expansion MLP layer and drop LayerScale [41] as it becomes redundant.

global response normalization (GRN) layer? In Table 2d, we compare GRN with the three widely used normalization layers: Local Response Normalization (LRN) [26], Batch Normalization (BN) [22], and Layer Normalization (LN) [1]. We observe that only GRN can significantly outperform the supervised baseline. LRN lacks global context as it only contrasts channels within nearby neighbors. BN normalizes spatially along the batch axis, which is unsuitable for masked inputs. LN implicitly encourages feature competition through global mean and variance standardization but does not work as well as GRN.

Relation to feature gating methods. Another way to enhance competition across neurons is to use dynamic feature gating methods [19, 34, 45, 48, 54]. In Table 2e, we compare our GRN with two classic gating layers: squeeze-and-excite (SE) [19] and convolutional block attention module (CBAM) [48]. SE focuses on channel gating, while CBAM focuses on spatial gating. Both modules can increase the

Backbone	Method	#param	FLOPs	Val acc.
ConvNeXt V1-B	Supervised	89M	15.4G	83.8
ConvNeXt V1-B	FCMAE	89M	15.4G	83.7
ConvNeXt V2-B	Supervised	89M	15.4G	84.3 (+0.5)
ConvNeXt V2-B	FCMAE	89M	15.4G	84.6 (+0.8)
ConvNeXt V1-L	Supervised	198M	34.4G	84.3
ConvNeXt V1-L	FCMAE	198M	34.4G	84.4
ConvNeXt V2-L	Supervised	198M	34.4G	84.5 (+0.2)
ConvNeXt V2-L	FCMAE	198M	34.4G	85.6 (+1.3)

Table 3. Co-design matters. When the architecture and the learning framework are co-designed and used together, masked image pre-training becomes effective for ConvNeXt. We report the fine-tuning performance from 800 epoch FCMAE pre-trained models. The relative improvement is bigger with a larger model.

contrast of individual channels, similar to what GRN does. GRN is much simpler and more efficient as it does not require additional parameter layers (such as MLPs).

The role of GRN in pre-training/fine-tuning. Finally, we examine the importance of GRN in pre-training and fine-tuning. We present results in Table 2f where we either remove GRN from fine-tuning or add newly initialized GRN only at the time of fine-tuning. Either way, we observe a significant performance degradation, suggesting that keeping GRN in both pre-training and fine-tuning is important.

5. ImageNet Experiments

In this section, we present and analyze two key proposals, the FCMAE *pre-training framework* and ConvNeXt V2 *architecture*, which are co-designed to make masked-based self-supervised pre-training successful. We show these designs synergize well and provide a strong foundation for scaling the model to various sizes. Additionally, we compare our approach to previous masked image modeling approaches through experiments. Furthermore, we show that our largest ConvNeXt V2 Huge model, which has been pre-trained using the FCMAE framework and fine-tuned on the ImageNet-22K dataset, can achieve a new state-of-the-art of 88.9% top-1 accuracy on the ImageNet-1K dataset, using only publicly available data.

Co-design matters. In this paper, we conduct a unique study that involves *co-designing* both the self-supervised learning framework (FCMAE) and the model architecture improvement (GRN layer), through an empirical study of their learning behavior. The results presented in Table 3 demonstrate the importance of this approach.

We found that using the FCMAE framework without modifying the model architecture has a limited impact on representation learning quality. Similarly, the new GRN layer has a rather small effect on performance under the supervised setup. However, the combination of the two results in a significant improvement in fine-tuning perfor-

Backbone	Method	#param	PT epoch	FT acc.
ViT-B	BEiT	88M	800	83.2
ViT-B	MAE	88M	1600	83.6
Swin-B	SimMIM	88M	800	84.0
ConvNeXt V2-B	FCMAE	89M	800	84.6
ConvNeXt V2-B	FCMAE	89M	1600	84.9
ViT-L	BEiT	307M	800	85.2
ViT-L	MAE	307M	1600	85.9
Swin-L	SimMIM	197M	800	85.4
ConvNeXt V2-L	FCMAE	198M	800	85.6
ConvNeXt V2-L	FCMAE	198M	1600	85.8
ViT-H	MAE	632M	1600	86.9
Swin V2-H	SimMIM	658M	800	85.7
ConvNeXt V2-H	FCMAE	659M	800	85.8
ConvNeXt V2-H	FCMAE	659M	1600	86.3

Table 4. Comparisons with previous masked image modeling approaches. The pre-training data is the IN-1K training set. All self-supervised methods are benchmarked by the end-to-end fine-tuning performance with an image size of 224. We underline the highest accuracy for each model size and bold our best results.

mance. This supports the idea that both the model and learning framework should be considered together, particularly when it comes to self-supervised learning.

Model scaling. In this study, we evaluated a range of 8 models with different sizes, from a low-capacity 3.7M Atto model to a high-capacity 650M Huge model. We pre-trained these models using the proposed FCMAE framework and compared the fine-tuning results to the fully supervised counterparts.

The results, shown in Figure 1, demonstrate strong model scaling behavior, with consistently improved performance over the supervised baseline across all model sizes. This is the first time the benefit of masked image modeling has been demonstrated in such a broad model spectrum, both in terms of effectiveness and efficiency. The complete tabulated results can be found in the appendix.

Comparisons with previous methods. We compare our approach to previous masked auto-encoder methods [2, 14, 53], which were all designed for transformer-based models. The results are summarized in Table 4. Our framework outperforms the Swin transformer pre-trained with SimMIM [53] across all model sizes. Compared to the plain ViT pre-trained with MAE [14], our approach performs similarly up to the Large model regime, despite using much fewer parameters (198M vs 307M). However, in the huge model regime, our approach slightly lagged behind. This might be because a huge ViT model can benefit more from self-supervised pre-training. As we will see next, the gap might be closed with additional intermediate fine-tuning.

ImageNet-22K intermediate fine-tuning. We also present ImageNet-22K intermediate fine-tuning results [2]. The training process involves three steps: 1) FCMAE pre-

Type	Backbone	size	#param	FLOPS	Val acc.
Conv	Efficient V2-XL	480 ²	208M	94.0G	87.3
	ConvNeXt V1-XL	384 ²	350M	179.0G	87.8
Hybrid	CoAtNet-4	512 ²	275M	360.9G	88.1
	MaxViT-XL	384 ²	475M	293.7G	88.5
Trans	MaxViT-XL	512 ²	475M	535.2G	88.7
	MViTV2-H	384 ²	667M	388.5G	88.6
Conv	MViTV2-H	512 ²	667M	763.5G	88.8
	ConvNeXt V2-H	384 ²	659M	337.9G	88.7
Conv	ConvNeXt V2-H	512 ²	659M	600.7G	88.9

Table 5. **ImageNet-1K fine-tuning results using IN-21K labels.**

The ConvNeXt V2 Huge model equipped with the FCMAE pre-training outperforms other architectures and sets a new state-of-the-art accuracy of 88.9% among methods using public data only.

training, 2) ImageNet-22K fine-tuning, and 3) ImageNet-1K fine-tuning. We use 384² resolution images for pre-training and fine-tuning [20]. We compare our results to the state-of-the-art architecture designs, including convolution-based [33, 40], transformer-based [8], and hybrid designs [6, 42]. All these results were trained with ImageNet-22K supervised labels. The results are summarized in Table 5. Our method, using a convolution-based architecture, sets a new state-of-the-art accuracy using publicly available data only (*i.e.* ImageNet-1K and ImageNet-22K).

6. Transfer Learning Experiments

We now benchmark the transfer learning performance. First, we evaluate the impact of our co-design, *i.e.* comparing ConvNeXt V1 + supervised vs. ConvNeXt V2 + FCMAE. We also directly compare our approach with Swin transformer models pre-trained with SimMIM [53]. The training and testing details are provided in the appendix.

Object detection and segmentation on COCO. We fine-tune Mask R-CNN [15] on the COCO dataset [30] and report the detection mAP^{box} and the segmentation mAP^{mask} on the COCO val2017 set. The results are shown in Table 6. We see a gradual improvement as our proposals are applied. From V1 to V2, the GRN layer is newly introduced and enhances performance. Upon this, the model further benefits from better initialization when moving from supervised to FCMAE-based self-supervised learning. The best performances are achieved when both are applied together. Additionally, our final proposal, ConvNeXt V2 pre-trained on FCMAE, outperforms the Swin transformer counterparts across all model sizes, with the largest gap achieved in the huge model regime.

Semantic segmentation on ADE20K. To summarize, we conduct experiments on the ADE20K [56] semantic segmentation task using the UperNet framework [50]. Our results show a similar trend to the object detection experiments, and our final model significantly improves over the V1 supervised counterparts. It also performs on par with

Backbone	Method	FLOPS	AP ^{box}	AP ^{box} ₅₀	AP ^{box} ₇₅	AP ^{mask}	AP ^{mask} ₅₀	AP ^{mask} ₇₅
ConvNeXt V1-B	Supervised	486G	50.3	71.6	56.1	44.9	68.5	48.8
ConvNeXt V2-B	Supervised	486G	51.0	72.4	56.6	45.6	69.5	49.7
Swin-B	SimMIM	497G	52.3	—	—	—	—	—
ConvNeXt V2-B	FCMAE	486G	52.9	72.6	58.9	46.6	70.0	51.1
ConvNeXt V1-L	Supervised	875G	50.6	71.5	56.3	45.1	68.7	49.2
ConvNeXt V2-L	Supervised	875G	51.5	72.5	57.3	45.8	69.4	49.9
Swin-L	SimMIM	904G	53.8	—	—	—	—	—
ConvNeXt V2-L	FCMAE	875G	54.4	73.9	60.4	47.7	71.4	52.3
Swin V2-H	SimMIM	—	54.4	—	—	—	—	—
ConvNeXt V2-H	FCMAE	2525G	55.7	75.2	61.8	48.9	72.8	53.6

Table 6. **COCO object detection and instance segmentation results** using Mask-RCNN. FLOPS are calculated with image size (1280, 800). Swins’ results are from [53]. All COCO fine-tuning experiments rely on ImageNet-1K pre-trained models.

Backbone	Method	input	mIoU	#param	FLOPS
ConvNeXt V1-B	Supervised	512 ²	49.9	122M	1170G
ConvNeXt V2-B	Supervised	512 ²	50.5	122M	1170G
Swin-B	SimMIM	512 ²	52.8	121M	1181G
ConvNeXt V2-B	FCMAE	512 ²	52.1	122M	1170G
ConvNeXt V1-L	Supervised	512 ²	50.5	235M	1573G
ConvNeXt V2-L	Supervised	512 ²	51.6	235M	1573G
Swin-L	SimMIM	512 ²	53.5	234M	1601G
ConvNeXt V2-L	FCMAE	512 ²	53.7	235M	1573G
Swin V2-H	SimMIM	512 ²	54.2	—	—
ConvNeXt V2-H	FCMAE	512 ²	55.0	707M	3272G
ConvNeXt V2-H	FCMAE, 22K ft	640 ²	57.0	707M	5113G

Table 7. **ADE20K semantic segmentation results** using UPerNet. Swins’ results are from [53]. FLOPS are based on input sizes of (2048, 512) or (2560, 640). All ADE20K fine-tuning experiments rely on ImageNet-1K pre-trained model except FCMAE, 22K ft, in which case the ImageNet-1K pre-training is followed by ImageNet-22K supervised fine-tuning.

the Swin transformer in the base and large model regimes but outperforms Swin in the huge model regime.

7. Conclusion

In this paper, we introduce a new ConvNet model family called ConvNeXt V2 that covers a broader range of complexity. While the architecture has minimal changes, it is specifically designed to be more suitable for self-supervised learning. Using our fully convolutional masked autoencoder pre-training, we can significantly improve the performance of pure ConvNets across various downstream tasks, including ImageNet classification, COCO object detection, and ADE20K segmentation.

Acknowledgments. We thank Ross Wightman for the initial design of the small-compute ConvNeXt model variants and the associated training recipe. We also appreciate the helpful discussions and feedback provided by Kaiming He.

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [2] Hangbo Bao, Li Dong, and Furu Wei. BEiT: BERT pre-training of image transformers. In *ICLR*, 2022.
- [3] Fergus W Campbell and John G Robson. Application of fourier analysis to the visibility of gratings. *The Journal of physiology*, 1968.
- [4] Thomas Capelle. Finding the New Resnet18. https://wandb.ai/fastai/fine_tune_timm/reports/Finding-the-New-Resnet18--VmlldzoyMDI0MjU3, 2022.
- [5] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *CVPR*, 2019.
- [6] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. In *NeurIPS*, 2021.
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [8] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *CVPR*, 2021.
- [9] Peng Gao, Teli Ma, Hongsheng Li, Jifeng Dai, and Yu Qiao. MCMAE: Masked Convolution Meets Masked Autoencoders. In *NeurIPS*, 2022.
- [10] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- [11] Benjamin Graham, Martin Engelcke, and Laurens Van Der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, 2018.
- [12] Benjamin Graham and Laurens van der Maaten. Submanifold sparse convolutional networks. *arXiv preprint arXiv:1706.01307*, 2017.
- [13] H K Hartline, Henry G Wagner, and Floyd Ratliff. Inhibition in the eye of limulus. *The Journal of general physiology*, 1956.
- [14] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022.
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [18] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco An-
dreetto, and Hartwig Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv:1704.04861*, 2017.
- [19] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018.
- [20] Ronghang Hu, Shoubhik Debnath, Saining Xie, and Xinlei Chen. Exploring long-sequence masked autoencoders. *arXiv preprint arXiv:2210.07224*, 2022.
- [21] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017.
- [22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [23] Kevin Jarrett, Koray Kavukcuoglu, Marc'Aurelio Ranzato, and Yann LeCun. What is the best multi-stage architecture for object recognition? In *ICCV*, 2009.
- [24] Li Jing, Jiachen Zhu, and Yann LeCun. Masked siamese convnets. *arXiv preprint arXiv:2206.07700*, 2022.
- [25] Alex Krizhevsky, Ilya Sutskever, and Geoff Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012.
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 2017.
- [27] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1989.
- [28] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [29] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017.
- [30] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [31] Chenxi Liu, Piotr Dollár, Kaiming He, Ross Girshick, Alan Yuille, and Saining Xie. Are labels necessary for neural architecture search? In *ECCV*, 2020.
- [32] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- [33] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022.
- [34] Jongchan Park, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Bam: Bottleneck attention module. In *BMVC*, 2018.
- [35] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *CVPR*, 2020.
- [36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.

- [37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015.
- [38] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [39] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.
- [40] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models and faster training. In *ICML*, 2021.
- [41] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *ICCV*, 2021.
- [42] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. In *ECCV*, 2022.
- [43] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv:1607.08022*, 2016.
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [45] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *CVPR*, 2020.
- [46] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- [47] Ross Wightman and Jeremy Howard. Which image models are best? <https://www.kaggle.com/code/jhoward/which-image-models-are-best>, 2022.
- [48] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, 2018.
- [49] Yuxin Wu and Kaiming He. Group normalization. In *ECCV*, 2018.
- [50] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, 2018.
- [51] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017.
- [52] Saining Xie, Jiatao Gu, Demi Guo, Charles R Qi, Leonidas Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In *ECCV*, 2020.
- [53] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhiliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *CVPR*, 2022.
- [54] Zongxin Yang, Linchao Zhu, Yu Wu, and Yi Yang. Gated channel transformation for visual recognition. In *CVPR*, 2020.
- [55] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using Places database. In *NeurIPS*, 2014.
- [56] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ADE20K dataset. In *IJCV*, 2019.