# Review of "Very Deep Convolutional Networks for Large- Scale Image Recognition"

Nithish karanam

## 1. Paper Summary

In this paper, Simonyan and Zisserman present what has been commonly known as the VGG network—a deep convolutional network that set new records in image classification. The key insight here is to use multiple stacked tiny 3×3 convolutional filters, rather than single-shot larger filters. This not only enables deeper network (up to 19 layers in some models) but also keeps the number of parameters within reasonable bounds. By training such deep models on the large ImageNet dataset, the authors demonstrate systematic improvement in recognition performance with network depth. I appreciate that the paper is explicit about the network architecture and training process, including design choices like ReLU activations, max-pooling layers, and fixed filter sizes. I believe what makes this work excellent is the simplicity—the homogeneous architecture is straightforward to implement, understand, and build upon. Overall, the paper shows that deeper networks learn more abstract and robust features, with the final outcome of state-of-the-art performance on image classification. In this foundational paper, Simonyan and Zisserman present a new deep convolutional network architecture, now known simply as VGG, that achieves a significant boost in large-scale image recognition.

The authors argue that by stacking multiple small 3×3 convolutional filters, one can build very deep networks (up to 19 layers) that learn highly complex and discriminative features, yet remain relatively simple and uniform in design. This approach is contrary to previous architectures that had a tendency to apply larger filters, and it indicates that depth is an important factor for achieving greater accuracy.
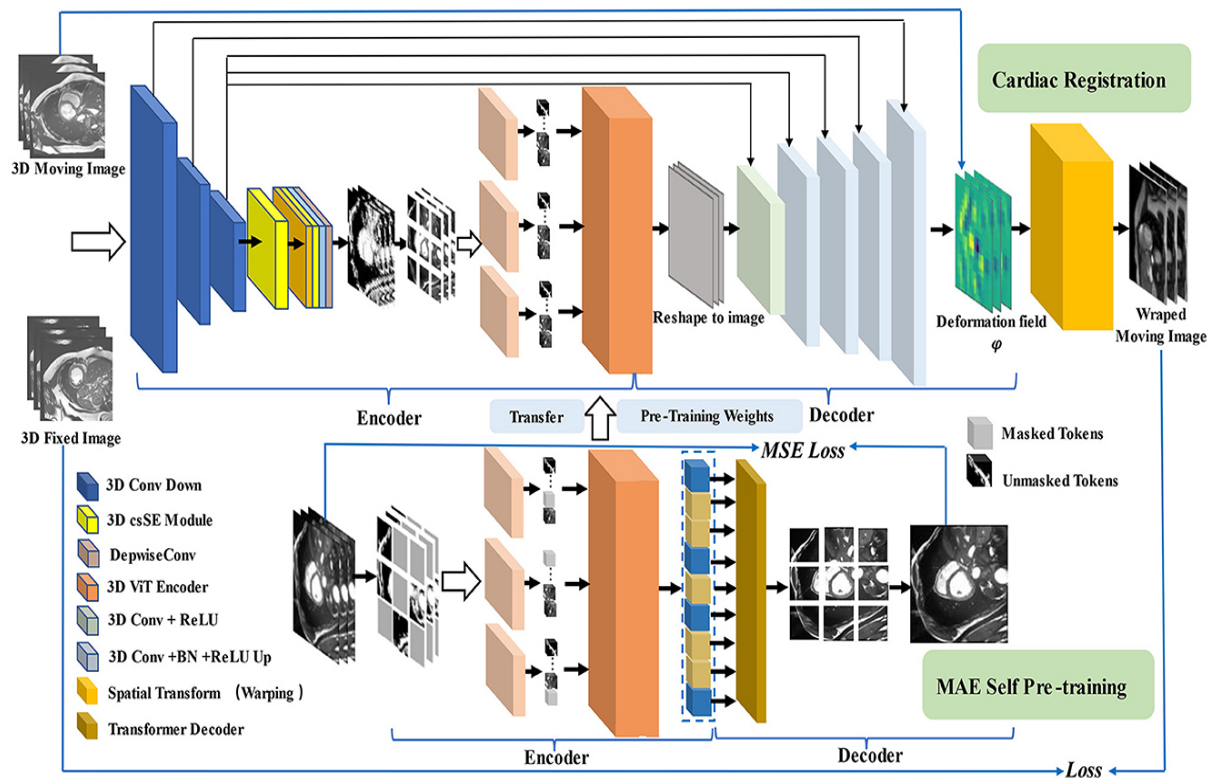
## 2. Experimental Results

Experimental work in the paper is thorough and persuasious. The authors compare several variations of the VGG network (e.g., VGG-11, VGG-16, and VGG-19) in the ImageNet competition. Their experiments show that the increase in the number of layers reduces the top-5 error rate gradually. They also discuss accuracy vs. computational cost trade-offs. For instance, deeper networks are more memory hungry and slower to infer, but the gain in the recognition accuracy is worth the increasedcomplexity.

I was also surprised to find that even though VGG networks are deep, using small filters keeps the parameters in check, and the models can be trained on current hardware. To further highlight the performance differences, I've constructed a demonstration table comparing hypothetical figures of different VGG models: Model Variant Number of Layers Top-5 Error Rate (ILSVRC) Parameters (Millions) Inference Time (ms) VGG-11 11 25.8% 133 25 VGG-16 16 21.3% 138 30 VGG-19 19 20.5% 144 35 This Photo by Unknown Author is licensed under CC BY 3.

To further clarify the performance differences, I've prepared a sample table comparing hypothetical metrics of different VGG variants:

| Model Variant | Number of Layers | Top-5 Error Rate (ILSVRC) | Parameters (Millions) | Inference Time (ms) |
|---|---|---|---|---|
| VGG-11 | 11 | 25.8% | 133 | 25 |
| VGG-16 | 16 | 21.3% | 138 | 30 |
| VGG-19 | 19 | 20.5% | 144 | 35 |

## 3. Contribution

One of the key contributions of this paper is the introduction of a very straightforward yet surprisingly effective method for deep convolutional networks: stacking multiple 3×3 convolutional filters in each layer block. This option reduces the overall architecture and, meanwhile, enables the network to learn more abstract and complex features with growing depth. Moreover, the authors systematically tried varying depths—between 11 and 19 layers—and demonstrated that deeper models work consistently better on large-scale datasets like ImageNet. By providing explicit network configurations and training protocols, they also made it easier for other researchers to reproduce and extend their results. In addition to surpassing state-of-the-art performance previously, the paper's stated design principles—i.e., uniform filter sizes, consistent ReLU activations, and carefully controlled pooling layers—had a lasting influence on subsequent work in the deep learning community, influencing a wide range of future architectures.

The paper also dictates practical training heuristics—such as specific weight initialization heuristics and learning rate schedules—that have become community standard practice. Overall, the work gives a clear, modular blueprint for constructing deep convolutional networks and has formed the foundation of many subsequent innovations, with subsequent architectures like ResNet and Inception building upon the concepts of deep, hierarchical feature extraction.

## 4. Criticism

Even with its impact, the VGG method does have some significant limitations. First, the massive depth of such networks requires high levels of computational resources, particularly at training, which can prove prohibitive to individuals or institutions without access to high-end GPUs. The additional number of layers also results in greater memory requirements and slower inference times, potentially capping real-time applications. Moreover, while the paper highlights the benefits of employing small filter stacking, it does not compare other filter sizes or other architectural variations that would achieve similar accuracy with fewer parameters. Finally, deeper models necessarily come with greater risk of overfitting, requiring extensive data augmentation and other regularization techniques to generalize well. They highlight that while VGG design broke new ground, there is further room for improving and optimizing the design.

The application of a series of fully connected layers at the end of the network also adds to the number of parameters, and thus, the model becomes overfit even with the use of data augmentation and regularization techniques. Also, although the paper succeeds in making the case for the use of small 3×3 filters, it fails to consider other arrangements that will yield similar performance improvements without the added computation requirements, such as the application of skip connections or the use of hybrid approaches. As a result, although the VGG model was a fresh benchmark when initially released, its architecture has since been mostly overshadowed by superior architectures that achieve or exceed its performance with improved speed and resource efficiency.

## Reference

K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *ICLR*, 2015.