

# 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  $3\times 3$  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  $3\times 3$  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.

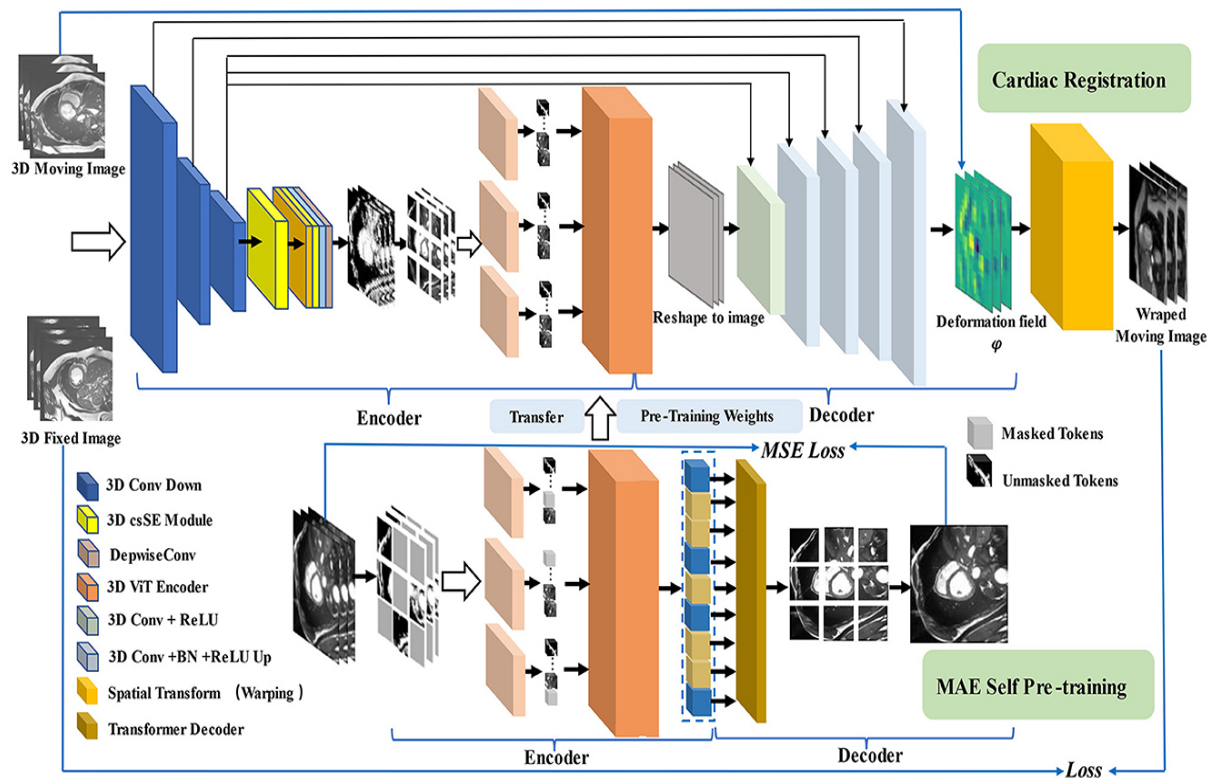
## 2. Experimental Results

Experimental work in the paper is thorough and persuasive. 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 increased complexity.

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.



[This Photo](#) by Unknown Author is licensed under [CC BY](#)

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

It makes it simpler to design and understand:

- Empirical Evidence for Depth: By controlled experimentation with networks at different depths, the paper develops strong evidence of why depth boost is beneficial to image recognition.
- Simplicity and Replicability: It is a light design, something that can reduce it immensely such that other scientists can easily reproduce and expand upon the research. The design is influential in more contemporary deep learning.
- Benchmark Performance: VGG network demonstrated new state-of-the-art performances on the ImageNet data set, and they demonstrated that careful but humble design can drastically improve classification accuracy.
- Computational Resources: The networks are computationally very intensive with deep structures, both during training and memory requirements. This can become a bottleneck for researchers with limited access to high-end GPUs.

### 4. Criticism

While the paper is highly influential, there are a few aspects that could possibly be criticized:

- Computational Resources: The networks are extremely computationally heavy with deep structures, both in the training process and in terms of memory consumption.
- This can prove to be a bottleneck for researchers with limited access to high-end GPUs.  
Inference Speed: With higher depth, the inference time also increases, which can impose a limit on the practical application of VGG networks in real-time scenarios
- Limited Investigation of Alternatives: The paper is almost entirely focused on the use of  $3 \times 3$  convolutions. It would have this, other regularization techniques might be explored.
- Overfitting: The deeper the models are, more prone to overfitting they are, and while the paper uses data augmentation and other techniques to prevent this, other regularization techniques might be explored

### Reference

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