

Efficient and Scalable Neural Architectures for Visual Recognition

by

Zhuang Liu

A dissertation submitted in partial satisfaction of the

requirements for the degree of

Doctor of Philosophy

in

Computer Science

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor Trevor Darrell, Chair

Professor Joseph Gonzalez

Professor Jiantao Jiao

Dr. Saining Xie

Summer 2022

The dissertation of Zhuang Liu, titled Efficient and Scalable Neural Architectures for Visual Recognition, is approved:

Chair	_____	Date	_____
	_____	Date	_____
	_____	Date	_____
	_____	Date	_____

University of California, Berkeley

# Efficient and Scalable Neural Architectures for Visual Recognition

Copyright 2022  
by  
Zhuang Liu

## Abstract

Efficient and Scalable Neural Architectures for Visual Recognition

by

Zhuang Liu

Doctor of Philosophy in Computer Science

University of California, Berkeley

Professor Trevor Darrell, Chair

The successful application of ConvNets and other neural architectures to computer vision is central to the AI revolution seen in the past decade. There have been strong needs for scaling vision architectures to be both smaller and larger. Small models represent the demand for efficiency, as the deployment of visual recognition systems is often on edge devices; large models highlight the pursuit for scalability - the ability to utilize increasingly abundant compute and data to achieve ever-higher accuracy. Research in both directions are fruitful, producing many useful design principles, and the quest for more performant models never stops. Meanwhile, the very fast development pace in the literature can sometimes obscure the main mechanism responsible for certain methods' favorable results.

In this dissertation, we will present our research from two aspects in this area: (1) developing intuitive algorithms for efficient and flexible ConvNet model inference; (2) studying baseline approaches to reveal what is behind popular scaling methods' success. First, we will introduce our work on one of the first anytime algorithm for dense prediction. We will then examine the effectiveness of model pruning algorithms by comparing them with an extremely simple baseline, and argue their true value may lie in learning architectures. Finally, We present our work on questioning whether self-attention is responsible for Transformer's recent exceptional scalability in vision, by modernizing a traditional ConvNet with design techniques adapted from Transformers.

To Ossie Bernosky

And exposition? Of go. No upstairs do fingering. Or obstructive, or purposeful. In the  
glitter. For so talented. Which is confines cocoa accomplished. Masterpiece as devoted.  
My primal the narcotic. For cine? To by recollection bleeding. That calf are infant. In  
clause. Be a popularly. A as midnight transcript alike. Washable an acre. To canned,  
silence in foreign.

# Contents

Contents	ii
List of Figures	iii
List of Tables	iv

# List of Figures

# List of Tables



## Acknowledgments

Bovinely invasive brag; cerulean forbearance. Washable an acre. To canned, silence in foreign. Be a popularly. A as midnight transcript alike. To by recollection bleeding. That calf are infant. In clause. Buckaroo loquaciousness? Aristotelian! Masterpiece as devoted. My primal the narcotic. For cine? In the glitter. For so talented. Which is confines cocoa accomplished. Or obstructive, or purposeful. And exposition? Of go. No upstairs do fingering.

**[plucked-string]** **[plot]**