

Data-dependent Resource Allocation and Routing in Multi-path Neural Networks for Vision-based Learning

M H G Dumindu Tissera

(188013F)

Thesis submitted in partial fulfillment of the requirements for the degree Doctor
of Philosophy

Department Electronic and Telecommunication Engineering

University of Moratuwa

Sri Lanka

September 2023

DECLARATION

I declare that this is my own work and this dissertation does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature:  ***UOM Verified Signature***

Date: 2023/09/10

The above candidate has carried out research for the PhD thesis under my supervision.

Signature of the Supervisor(s):

Date: 10 September 2023

UOM Verified Signature

Dr. Ranga Rodrigo

Abstract

As the depth of a neural network increases, the non-linearity and more parameters allow it to learn more complex functions. While network deepening has been proven effective, there is still an opportunity for efficient feature extraction within a layer that will improve the overall performance for the complexity of the network. Widening networks by adding more filters to each layer is the naive approach towards strengthening layer-wise feature extraction. It is an inefficient scaling option, considering the number of parameters being quadratic with the number of filters employed per layer. In contrast, parallel extractors in each layer provide an efficient scaling option. However, without context-dependent input allocation among these processes, such parallel computations tend to learn similar features, collapsing to a single computation.

Thus, it is vital to study the parallel stacking of computations layer-wise and design a routing method that allocates incoming feature maps to these computations. The expected outcome is to group homogeneous feature maps in parallel layers and employ exclusive filter sets to each of the groups (paths) so that the filter sets of each path can specialize in extracting features exclusive to each group.

To allow the network input to be routed end-to-end over such parallel paths, we propose data-dependent parallel resource allocation methods layer-wise. Given a layer of parallel tensors, we first employ sub-networks that produce gating coefficients to weigh cross-connections to the next layer of parallel tensors. Then, the next layer’s parallel tensors are constructed by getting summations of the current layer’s tensors, each weighted by the corresponding gating coefficient. We demonstrate that our multi-path networks outperform previous widening and adaptive feature extraction, ensembles, and deeper networks with comparable complexity using image recognition challenges.

To further regularize gating sub-networks, we think of a gating network’s path allocation as a soft clustering of its input feature maps. Thus, we propose a neural mixture model-based clustering objective to use as a regularization loss for the gating networks, which we first study as a standalone neural network-based clustering approach. The proposed clustering framework uses a neural network to learn cluster distributions in mixture modeling instead of tuning human-defined distributions. We adopt the Expectation-Maximization (EM) algorithm to train the network and perform batch-wise EM iterations where the forward pass acts as the E-step and the backward pass as the M-step. For image clustering, we use the mixture-based EM objective as the clustering objective, along with consistency optimization. Our networks outperform traditional and single-stage deep clustering methods that still depend on k-means.

Finally, we propose using this clustering objective to regulate gating networks to get distributed gating activation patterns. We show that the skewed gating patterns can be improved with such regularization loss as a local regularization. We further present the need for a global regularization method that takes the end task performance into account. We also suggest extending research towards sparse resource allocation, along with gating networks to handle more diversity.

Keywords- Multi-Path Networks, Data-Dependent Routing, Deep Clustering, Neural Mixture Models

DEDICATION

To the resilient people of Sri Lanka,
who have faced the challenges of the country's financial crisis with unwavering
determination and faith,
who continue to endure and strive for a better future for the nation.
This dissertation is a tribute to your strength and commitment to a brighter
tomorrow.

ACKNOWLEDGMENTS

I would like to express my deep appreciation and gratitude to the following individuals for their invaluable contributions to the completion of my Ph.D. thesis.

First and foremost, I would like to thank my adviser, Dr. Ranga Rodrigo, for his unwavering support, guidance, and encouragement throughout my Ph.D. studies. His expertise, insight, and mentorship have been invaluable to my academic and personal growth.

I also extend my gratitude to Dr. Subha Fernando for her technical support, which has contributed significantly to the success of this research. Additionally, I am grateful to Prof. Sanath Jayasena and Dr. Jayathu Samarawickrama for their constructive suggestions and insightful comments, which have helped refine and improve this thesis's technical content.

I would like to acknowledge the generous support of CodeGen International (PVT) Ltd and its CEO Dr. Harsha Subashinghe, who provided funding for this research. I am grateful for their support, which has allowed me to pursue my academic goals.

I would also like to thank Prof. Dileeka Dias, Prof. Sanath Jayasena, Dr. Ranga Rodrigo, Dr. Subha Fernando, Dr. Jayathu Samarawickrama, Dr. Harsha Subhasingshe, and CodeGen International (PVT) Ltd for their role in the establishment of QBITS Lab at the University of Moratuwa and the Ph.D. program.

I extend my sincere appreciation to Prof. Ruwan Udayanga and Dr. Pratha-

pasinghe Dharmawansa, the Research Coordinators, for coordinating the program and providing feedback on technical content. I am also indebted to Dr. Charith Chiththaranjan, the chair of the progress review panel, for reviewing and providing valuable feedback on my work.

I thank my colleagues and collaborators, Rukshan Wijesinghe, Kasun Vithanage, Alex Xavier, and Pubudu Ekanayake, for their outstanding technical contributions and unwavering support throughout this research.

Finally, I thank the thesis examiners Prof Salman Khan, Dr. Sadeep Jayasumana and Dr. Shehan Perera, and PhD examination chair Prof. Gihan Dias for providing valuable suggestions.

Thank you for your invaluable contributions, guidance, and support in making this achievement possible.

TABLE OF CONTENTS

Declaration	i
Abstract	ii
Dedication	iii
Acknowledgments	iv
Table of Contents	ix
List of Figures	xi
List of Tables	xii
1 Introduction	1
1.1 Towards Improving Layer-wise Feature Extraction	1
1.2 Usage of Parallel Stacks of Feature Maps (Paths)	2
1.3 Need for Layer-wise Routing Layers	3
1.4 Contributions	5
1.5 Summary	6

2	Background	7
2.1	Neural Networks	7
2.2	Deepening Neural Networks	8
2.3	Enriching Layer-wise Feature Extraction—Widening	9
2.4	Multi-path Networks with Cross-connections	10
2.5	Adaptive Feature Extraction Methods	11
2.6	Mixture of Experts	11
2.7	Summary	12
3	Data-Dependent End-to-End Routing	13
3.1	Cross-Prediction Based Routing	13
3.2	Cross-Connection Based Routing	17
3.3	Backpropagation through a Cross-Connection Layer	20
3.4	Experiments	23
3.4.1	Conventional Convolutional Neural Networks with Parallel Paths	23
3.4.2	Residual Networks with Parallel Paths	27
3.4.3	Multi-path ResNets on ILSVRC2012	30
3.5	Visualizations	33
3.5.1	Routing Visualization	34
3.5.2	Gate Maximization Patterns	35

3.5.3	Class-Wise Gating Patterns	37
3.5.4	Parallel Computation Weights	38
3.6	Conclusion	39
4	Neural Mixture Models for Clustering	41
4.1	Introduction	41
4.2	Related Work	44
4.3	EM Algorithm in Mixture Modeling	46
4.4	Formulating Mixture-EM on a Neural Network	48
4.4.1	Approximate Cluster Distributions	49
4.4.2	Deploying EM Batch-Wise	53
4.4.3	Image Clustering with Consistency Optimization	56
4.5	Experiments	58
4.5.1	Two-Dimensional Space	58
4.5.2	Image Clustering	61
4.5.3	Visualizations	64
4.6	Conclusion	67
5	Regularize Routing with Clustering Loss	68
5.1	Regularization of Cross-connection based Routing	68
5.2	Towards Global Regularization	69

5.3 Towards Sparse Multi-path Networks	70
6 Conclusions	72
List of Publications	75
References	90

LIST OF FIGURES

1.1	Need for routing throughout a network with parallel resources . . .	4
3.1	Cross-prediction-based Routing layer of m inputs and n outputs. . .	14
3.2	Insertion of cross-prediction-based routing layers to a two-path CNN in image classification task.	16
3.3	Cross-connecting two layers that have two parallel tensors on each layer.	18
3.4	CNN with two paths and adaptive cross-connections placed at spe- cific locations.	19
3.5	Backpropagation through the simplified cross-connecting layer be- tween two successive layers each containing two parallel tensors. . .	21
3.6	ResNet variant performance (accuracy) in CIFAR with the number of parameters (millions).	31
3.7	Route visualizations of VGG13-2-CC.	34
3.8	Features that maximize gates.	36
3.9	Gate activation histograms.	38
3.10	Parallel operations' weight histograms.	39
4.1	Overview of mixture-EM formulation with a neural network. . . .	42

4.2	Sigmoid function and standard normal variable.	52
4.3	Deploying Mixture-EM method for end-to-end training of a neural network for clustering.	54
4.4	2-Dimensional clustering space.	58
4.5	Effect of the relevance score normalization.	59
4.6	Learning curves in clustering STL10.	64
4.7	Two-dimensional mapping of clustering network response for STL10 subset.	65
4.8	Image cluster visualization.	66
4.9	Clustering network. Convolutional filter visualization	67
5.1	Regularization of gating with clustering loss.	70

LIST OF TABLES

3.1	Notations and details of the compared convolutional neural networks.	24
3.2	CIFAR10 CNN ablation study.	26
3.3	Comparison of ResNet variants in CIFAR.	29
3.4	ILSVRC 2012 Dataset: Single-crop and 10-crop validation error (%).	33
4.1	Average cluster distribution outputs with a batch of 128 samples.	60
4.2	Network architectures used for clustering image datasets.	61
4.3	Clustering accuracy (%) comparison.	62