

Real-time segmentation on smartphone

AXEL DEMBORG

Master in Computer Science

Date: April 13, 2018

Supervisor: Hossein Azizpour

Examiner: Danica Kragic

Swedish title: Realtids segmentering på smartphone

School of Electrical Engineering and Computer Science

Abstract

Sammanfattning

Contents

- 1 Introduction 1**
 - 1.1 Research Question 2
- 2 Background 3**
 - 2.1 Related work 3
 - 2.1.1 Quantization of weights 6
 - 2.1.2 Student-teacher learning 7
 - 2.1.3 Architectural optimizations 10
- 3 Methods 12**
- Bibliography 13**
- A Unnecessary Appended Material 18**

Chapter 1

Introduction

1.1 Research Question

Chapter 2

Background

2.1 Related work

Convolutional Neural Networks (CNN) were first introduced in 1998 [28] and since then larger and larger CNNs have slowly become the state of the art method for most areas of computer vision. Notably *AlexNet* [27] in 2012 proved that deep CNNs could be used for high resolution image classification by beating the previous state of the art [35] on the *ImageNet* classification challenge [9]. To do this *AlexNet* used 60 million parameters and 650,000 neurons and training of the network was only made feasible by the use of multiple graphical processing units (GPUs) [27].

In the areas of object detection and semantic segmentation it was *Regions with CNN features (R-CNN)* [12] that in 2014 first showed that CNNs could be successfully be applied to these fields by significantly improving over the previous state of the art in object detection [33] and after minor modifications matching the performance of the state of the art in semantic segmentation [5] with a system not specifically built for the task. For object detection *R-CNN* works as a hybrid system with *selective search* [40] producing proposals for object regions and a CNN, pre-trained on *ImageNet* [9] and fine-tuned for region classification, generating fixed length features for each region, finally classifying each region by running class specific *support vector machines* [3] on these features. Some issues with *R-CNN* are that it requires a multistage training, training the CNN to give good features and training the SVMs for classification, and that it is slow, for training but most notably at inference where one image is processed in 47s. These prob-

lems are addressed with further work resulting in *Fast R-CNN* [11] where the CNN isn't run once per proposed region but instead once for the entire network generating a convolutional feature map that is then pooled with a region of interest (RoI) pooling layer to produce a feature vector for each region. These feature vectors are then feed into a fully connected neural network with two sibling output layers that perform both classification and bounding box refinement in parallel. With these improvements *Fast R-CNN* achieves faster inference and higher accuracy than its predecessors and does so with a arguably much more elegant design. Even though *Fast R-CNN* improved speed significantly it was nowhere near real-time performance, further performance improvements were introduced with *Faster R-CNN* [32] where *selective search* for region proposals is replaced with Region Proposal Networks, fully convolutional neural networks that take as input the convolutional feature maps as described from *Fast R-CNN* and outputs region proposals. Since this approach for region proposals shares most of its computation with the classification network the region proposals are practically free and frame rates of 5fps are achievable. The region proposal networks not only speed up computation but also prove to give better accuracy region proposals and thus raise over all accuracy in the system as well [32]. Even further improvements to this framework was achieved with the introduction of *Mask R-CNN* [18] which expands upon *Faster R-CNN* by adding a third branch for a segmentation mask besides the branches for bounding box refinement and classification making the system able to predict not only the general bounding box of items in the image but also which exact pixels belong to the object. Since segmentation is a pixel-by-pixel prediction problem *Mask R-CNN* replaces the spatially quantizing RoIPool operation from *Fast R-CNN* with a quantization-free layer called RoIAlign.

Some parallel work on semantic segmentation of images resulted in *SegNet* [2], a fully convolutional encoder-decoder network. Here the encoder network encodes the input image down into a lower dimensional feature space while keeping storing the indices of the max pooling operations. The low dimensional representations are then run through a decoder network which is architecturally a mirror image of the encoder network but where the max pooling operations have been replaced with upsampling layers that use the stored indices from the corresponding pooling layers to maintain the granularity of the im-

ages. The final layer in the network is a softmax and hence the outputs are the probabilities of each pixel belong to each class. Continued work on segmentation utilizes blocks of so called *DenseNets* [21], CNNs where every layer is connected to every layer after it enabling the training of really deep network architectures by alleviating the vanishing gradient problem and promoting feature reuse between the layers. By using these *DenseNets* in a very deep encoder-decoder structure where skip connections restore image granularity during upsampling the state of the art in image segmentation has been pushed even further [25] while still reducing the amount of parameters required for the models by a factor 10 as compared to the previous state of art.

Despite their impressive performance on a wide range of problems neural networks are still prohibited from running locally on mobile devices with slow processors, limited power envelopes or limited memory due to their large size and big computational load. For example modern neural networks can't fit on the on-chip SRAM cache and instead they have to reside in the much more power hungry off-chip DRAM memory making applications up to 100 times more power consuming [17]. Regarding inference speed the most modern networks for object segmentation [18] run at 5fps but that is on high performance GPUs meaning that mobile performance is far from real-time. Due to these limitations applications of neural networks for mobile use cases are either to forced give up on state of the art performance or to be run on off-site servers which requires steady network connections and incurs delays, both of which may be intolerable for real-time mobile applications, self driving cars and robotics [26]. However work on understanding the structure of the learned weights in neural networks [10] has showed that there is significant redundancy in the parameterization of several deep learning models and that up to 95% of weights in networks can be predicted from the remaining 5% without any drop in accuracy. This indicates that models could be made much smaller while still maintaining performance and several such approaches for squeezing high performance networks into small memory footprints and computational loads have been proposed. The most prominent approaches will be presented below.

2.1.1 Quantization of weights

Modern neural networks are usually based on 32-bit floating point representations of parameters. It has been shown however that networks are quite resilient to noise and even that some noise can improve training [30]. Since reduced precision variables can be modeled as noise this means that networks can be compressed by changing to a less accurate format without any loss in performance. This can be done either by reducing the bit accuracy after training [43] or by doing the entire training in reduced accuracy [22] [14]. The benefits of using a reduced format like this for representation is not only that the models take less space but also that the individual multiplications become cheaper and hence the networks run faster.

Weight sharing

One of the most direct approaches for removing the redundancy in parametrization from neural networks is by forcing the networks to share weights between different connections. This is precisely what *HashedNets* [7] does by fixing the amount of weights K^l that are to be used in each layer making the weights $\vec{w}^l \in \mathbb{R}^{K^l}$ and using hashing functions to map each element in the virtual weight matrices V_{ij}^l to one of these weights $V_{ij} = w_{h(i,j)}$ with $h()$ being a hashing function. With the weight matrices defined in this fashion *HashedNets* can be trained like normal networks with the gradients with respect to the weights calculated from the gradients with respect to the virtual matrices as

$$\frac{\partial \mathcal{L}}{\partial w_k^l} = \sum_{ij} \frac{\partial \mathcal{L}}{\partial V_{ij}^l} \frac{\partial V_{ij}^l}{\partial w_k^l}$$

This method gave a compression of about 20 times before any notable loss in accuracy was introduced during tests on variations of the MNIST dataset which seems to agree very well with the results from [10].

Other notable work focuses on the use of k-means clustering to cluster the weights in networks after training [13], this proves to work very well and manages to compress the models with a factor 16 with no more than a 0.5% drop in classification accuracy on the ImageNet dataset. Further work in this area explores the effects of pruning away low-weight connections and iterative retraining of the pruned networks

[17]. This lets the authors compress models with a factor 9 - 13 without any loss in performance while getting sparser weight matrices that could potentially speed up calculations. These two lines of research, clustering and pruning, were merged into a single framework called *deep compression* [15] where a three stage approach is taken to model compression. First low-weight connections are pruned away and the network is retrained to compensate for this, in the second stage k-means clustering is performed on the weights and again the network is retrained to make the clusters take the most useful values, finally Huffman coding [42] is used to reduce the storage required for the weights. This process allows *deep compression* to compress networks with a factor 35 without any loss in accuracy. Despite these very impressive results however *deep compression* comes with a major drawback, it can't be run efficiently in its compressed form and the full weight matrices have to be rebuilt at inference time to use the models on commodity hardware. To alleviate these problems hardware has been designed that could perform prediction directly from the compressed models. This so called *efficient inference engine* [16] would enable inference 13 times faster than GPU while being 3400 times more energy efficient.

2.1.2 Student-teacher learning

Student-teacher learning is a type of model compression where a smaller and/or faster to compute *student* network is trained by making it learn the representations learned by a larger *teacher* network. This idea was first introduced for compressing ensemble models produced by *Ensemble Selection* [6] which consist of hundreds of models of many different kinds, support vector machines, neural networks, memory based models, and decision trees into a single neural network [4]. This work leverages the neural networks property of being universal approximators [8], meaning that given sufficiently much training data and a big enough hidden layer a neural network can learn to approximate any function with arbitrary precision, by not directly training the student network on the relatively limited labeled training data available but instead on large amounts of pseudo random data that has been given labels by first being passed through the large teacher ensemble. This compression technique yielded student networks up to 1000 times smaller and 1000 times faster to compute than their teachers with a negligible drop in accuracy on some test problems.

Further work on student-teacher learning experiments with why deep neural networks usually perform better than shallow ones, even when they have the same amount of parameters. This was done by training shallow student models to mimic deep teachers [1]. The work introduces two major modifications that make training of these student models feasible, firstly the student model isn't tasked with just recreating the same label as the teacher but also the same distribution which is achieved by regressing the student to the logits, log probability, values of the teacher as they were before softmax. Getting predictions from the student is then achieved by adding a softmax layer to the end of it after training. Secondly a bottleneck linear layer is added to the network to speed up training. With these modifications they are able to train flat neural networks for both the TMIT and CIFAR-10 datasets with performance closely matching that of single deep networks. Continued analysis of flat networks however shows that depth and convolutions are critical for getting good performance on image classification datasets [41]. Empirically this claim is supported by training state of the art, deep, convolutional models for classification on the CIFAR-10 dataset and then building an ensemble of such models using that as a teacher for shallow students. The student models were then compared to deep convolutional benchmarks that were not trained in a student-teacher fashion. To make sure that the networks were all performing to the best of their abilities and thus making the comparison fair Bayesian hyperparameter optimization [38] was used. Through this thorough analysis it was shown that shallow networks are unable to mimic the performance of deep networks if the number of parameters is held constant between them, these findings are also in agreement with the theoretical results that the representational efficiency of neural networks grows exponentially with the number of layers [29].

Improvements to the student-teacher learning method have been proposed where the student is tasked with minimizing the weighted average of the cross-entropy between its own output and the teacher output when the last layer is softmax with increased temperature, yielding softer labels, and the cross-entropy between the student output and the correct labels when they are available. This framework is called *Distillation* [19] and proves to work very well for transferring of information from teacher to student. The framework is demonstrated by training a student model with only 13.2% test error on the MNIST

dataset despite only having seen 7s and 8s during its own training. These results mean that distillation manages to transfer knowledge about how a 6 looks from the teacher to the student by only telling it to what degree different 7s and 8s don't look like 6s.

Continued work lead to the creation of *FitNets* [34] which goes in the opposite direction to previous attempts at student architectures and instead proposes very deep but thin students. To enable learning in these deep student networks a stage-wise training procedure is used. In the first stage intermediate layers in the teacher and student networks are selected, these are called *hint* and *guided* layers respectively. The guided layer in the student is then tasked to mimic the hint layer in the teacher through a convolutional regressor that compensates for the difference in number of outputs between the networks, this procedure gives a good initialization for the first layers in the student and allows for it to learn the internal representations of the data from the teacher. The second stage of training is then distillation as described above but with the small addition that the weight of the loss against the teacher is slowly annealed during training. This annealing allows for the student to lean heavily on the teacher for support in early stages of training and learn samples which the even the teacher struggles with towards the end of its training. Using this approach the *FitNets* manage to produce predictions at the same level or in some cases even better than models with 10 times more parameters.

Some more recent work [36] builds upon the ideas from *FitNets* with not only letting the students mimic the output of teachers but also some intermediary representations. Unlike the way it is done *FitNets* however the student is not tasked with reconstructing the exact activations of the teacher in the intermediate layers but instead the attention maps, regions in the image that the teacher uses to make its predictions, and thus teaches the student where to look. A few different methods for calculating these attention maps are proposed in the paper but notable is that they are all non parametric meaning that no extra layers of convolution have to be learned to make the student attention maps comparable to the ones from the teacher. This attention transferring approach is proves to give good results on a number of difficult datasets including *ImageNet* and is also shown to work well together with distillation.

2.1.3 Architectural optimizations

Another orthogonal approach for compression is to optimize the convolutional layers themselves making them require less parameters or less computation to perform their tasks but still keep as much as possible of their representational power. One of the simplest things that can be done here is to replace single layers of $N \times N$ convolutional filters with two layers with $N \times 1$ and $1 \times N$ filters respectively, this reduces the amount of parameters that have to be stored per channel from N^2 to $2N$ and the amount of multiplications that have to be made scale in the same way. These asymmetrical convolutions have seen successful use in inception models [39]. Other variations on the convolutional operator that help compress the networks are dilated convolutions [44] where an exponentially expanding receptive field is achieved without the need for any extra parameters. There have also been some promising results from *depthwise separable convolutions* where the convolution is factored into a depthwise convolution followed by a pointwise 1×1 convolution reducing the computational load with a factor 8 to 9 for 3×3 convolutional kernels [20]. This scheme was introduced in [37] and has since been successfully used in *Inception* models [24].

SqueezeNet [23] presents a different take on how to get smaller models in that it rather optimizes the architecture of the network than any of the constituent parts, this approach gives a network with *AlexNet* performance but with 50 times fewer parameters than normal *AlexNet*. This is done by focusing on the usage of 1×1 convolutional filters, reducing the amount of channels that go in to the larger filters and by holding out on downsampling so that feature maps are kept large through the network. It was also proven that these results were orthogonal from compression by running the SqueezeNet through the *deep compression* framework [15] and getting further 10 times compression with out accuracy loss.

MobileNets [20] combine these two approaches, utilizing both *depthwise separable convolutions* and a heavily optimized architecture to build networks specially suited for mobile vision applications. In doing so *MobileNets* also introduce two hyper-parameters, *width-multiplier* and *resolution-multiplier* that help design models with a optimal trade off between latency and precision given the limitations of the available hardware.

Another network specially designed for real-time segmentation on

mobile devices is *ENet* [31]. Here dilated convolutions are used together with asymmetrical convolutions to give a large receptive field without introducing that many parameters. The network is built as an encoder-decoder network but with a much smaller decoder, the argument behind this being that decoder should simply upsample the output while fine-tuning the details which should be a simpler task than the information processing and extraction that the encoder is performing. Attention has also been paid to quickly downsampling the feature maps which saves on computation but then not downsampling so aggressively after that keeping much of the spatial information in the images. Together these improvements give a network that performs on par with *SegNet* but that requires 79 times fewer parameters and is 18 times faster at inference.

Chapter 3

Methods

Bibliography

- [1] Jimmy Ba and Rich Caruana. “Do deep nets really need to be deep?” In: *Advances in neural information processing systems*. 2014, pp. 2654–2662.
- [2] Vijay Badrinarayanan, Ankur Handa, and Roberto Cipolla. “Segnet: A deep convolutional encoder-decoder architecture for robust semantic pixel-wise labelling”. In: *arXiv preprint arXiv:1505.07293* (2015).
- [3] Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. “A training algorithm for optimal margin classifiers”. In: *Proceedings of the fifth annual workshop on Computational learning theory*. ACM. 1992, pp. 144–152.
- [4] Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. “Model compression”. In: *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM. 2006, pp. 535–541.
- [5] Joao Carreira et al. “Semantic segmentation with second-order pooling”. In: *European Conference on Computer Vision*. Springer. 2012, pp. 430–443.
- [6] Rich Caruana et al. “Ensemble selection from libraries of models”. In: *Proceedings of the twenty-first international conference on Machine learning*. ACM. 2004, p. 18.
- [7] Wenlin Chen et al. “Compressing neural networks with the hashing trick”. In: *International Conference on Machine Learning*. 2015, pp. 2285–2294.
- [8] George Cybenko. “Approximation by superpositions of a sigmoidal function”. In: *Mathematics of control, signals and systems* 2.4 (1989), pp. 303–314.

- [9] Jia Deng et al. "Imagenet: A large-scale hierarchical image database". In: *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE. 2009, pp. 248–255.
- [10] Misha Denil et al. "Predicting parameters in deep learning". In: *Advances in neural information processing systems*. 2013, pp. 2148–2156.
- [11] Ross Girshick. "Fast r-cnn". In: *arXiv preprint arXiv:1504.08083* (2015).
- [12] Ross Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014, pp. 580–587.
- [13] Yunchao Gong et al. "Compressing deep convolutional networks using vector quantization". In: *arXiv preprint arXiv:1412.6115* (2014).
- [14] Suyog Gupta et al. "Deep learning with limited numerical precision". In: *International Conference on Machine Learning*. 2015, pp. 1737–1746.
- [15] Song Han, Huizi Mao, and William J Dally. "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding". In: *arXiv preprint arXiv:1510.00149* (2015).
- [16] Song Han et al. "EIE: efficient inference engine on compressed deep neural network". In: *Computer Architecture (ISCA), 2016 ACM/IEEE 43rd Annual International Symposium on*. IEEE. 2016, pp. 243–254.
- [17] Song Han et al. "Learning both weights and connections for efficient neural network". In: *Advances in neural information processing systems*. 2015, pp. 1135–1143.
- [18] Kaiming He et al. "Mask r-cnn". In: *Computer Vision (ICCV), 2017 IEEE International Conference on*. IEEE. 2017, pp. 2980–2988.
- [19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network". In: *arXiv preprint arXiv:1503.02531* (2015).
- [20] Andrew G Howard et al. "Mobilenets: Efficient convolutional neural networks for mobile vision applications". In: *arXiv preprint arXiv:1704.04861* (2017).

- [21] Gao Huang et al. "Densely connected convolutional networks". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. Vol. 1. 2. 2017, p. 3.
- [22] Itay Hubara et al. "Quantized neural networks: Training neural networks with low precision weights and activations". In: *arXiv preprint arXiv:1609.07061* (2016).
- [23] Forrest N Iandola et al. "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size". In: *arXiv preprint arXiv:1602.07360* (2016).
- [24] Sergey Ioffe and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift". In: *International conference on machine learning*. 2015, pp. 448–456.
- [25] Simon Jégou et al. "The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation". In: *Computer Vision and Pattern Recognition Workshops (CVPRW), 2017 IEEE Conference on*. IEEE. 2017, pp. 1175–1183.
- [26] Jonghoon Jin, Aysegul Dundar, and Eugenio Culurciello. "Flattened convolutional neural networks for feedforward acceleration". In: *arXiv preprint arXiv:1412.5474* (2014).
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems*. 2012, pp. 1097–1105.
- [28] Yann LeCun et al. "Gradient-based learning applied to document recognition". In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.
- [29] Shiyu Liang and R Srikant. "Why deep neural networks for function approximation?" In: *arXiv preprint arXiv:1610.04161* (2016).
- [30] Alan F Murray and Peter J Edwards. "Enhanced MLP performance and fault tolerance resulting from synaptic weight noise during training". In: *IEEE Transactions on neural networks* 5.5 (1994), pp. 792–802.
- [31] Adam Paszke et al. "Enet: A deep neural network architecture for real-time semantic segmentation". In: *arXiv preprint arXiv:1606.02147* (2016).

- [32] Shaoqing Ren et al. "Faster r-cnn: Towards real-time object detection with region proposal networks". In: *Advances in neural information processing systems*. 2015, pp. 91–99.
- [33] Xiaofeng Ren and Deva Ramanan. "Histograms of sparse codes for object detection". In: *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. IEEE. 2013, pp. 3246–3253.
- [34] Adriana Romero et al. "Fitnets: Hints for thin deep nets". In: *arXiv preprint arXiv:1412.6550* (2014).
- [35] Jorge Sánchez and Florent Perronnin. "High-dimensional signature compression for large-scale image classification". In: *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE. 2011, pp. 1665–1672.
- [36] Nikos Komodakis Sergey Zagoruyko. "Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer". In: *International Conference on Learning Representations* (2017). URL: https://openreview.net/forum?id=Sks9_ajax.
- [37] Laurent Sifre and PS Mallat. "Rigid-motion scattering for image classification". PhD thesis. Citeseer, 2014.
- [38] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. "Practical bayesian optimization of machine learning algorithms". In: *Advances in neural information processing systems*. 2012, pp. 2951–2959.
- [39] Christian Szegedy et al. "Rethinking the inception architecture for computer vision". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, pp. 2818–2826.
- [40] Jasper RR Uijlings et al. "Selective search for object recognition". In: *International journal of computer vision* 104.2 (2013), pp. 154–171.
- [41] Gregor Urban et al. "Do deep convolutional nets really need to be deep and convolutional?" In: *arXiv preprint arXiv:1603.05691* (2016).
- [42] Jan Van Leeuwen. "On the Construction of Huffman Trees." In: *ICALP*. 1976, pp. 382–410.

- [43] Vincent Vanhoucke, Andrew Senior, and Mark Z Mao. "Improving the speed of neural networks on CPUs". In: *Proc. Deep Learning and Unsupervised Feature Learning NIPS Workshop*. Vol. 1. Cite-seer. 2011, p. 4.
- [44] Fisher Yu and Vladlen Koltun. "Multi-scale context aggregation by dilated convolutions". In: *arXiv preprint arXiv:1511.07122* (2015).

Appendix A

Unnecessary Appended Material