

A multitask deep learning model for real-time deployment in embedded systems

Miquel Martí*[†] and Atsuto Maki[†]

*Universitat Politècnica de Catalunya, Barcelona, Catalonia/Spain

[†]KTH Royal Institute of Technology, Stockholm, Sweden

Email: {miquelmr,atsuto}@kth.se

Abstract—We propose an approach to Multitask Learning (MTL) to make deep learning models faster and lighter for applications in which multiple tasks need to be solved simultaneously, which is particularly useful in embedded, real-time systems. We develop a multitask model for both Object Detection and Semantic Segmentation and analyze the challenges that appear during its training. Our multitask network is 1.6x faster, lighter and uses less memory than deploying the single-task models in parallel. We conclude that MTL has the potential to give superior performance in exchange of a more complex training process that introduces challenges not present in single-task models.

I. INTRODUCTION

Deep learning models (and machine learning models in general) focus on solving one task at a time. However, applications often require more than one task to be performed simultaneously and the naive solution is to deploy in parallel one model for each task.

Applications imposing real-time constraints require small inference times and prohibit off-board computation, forcing the deployment of deep learning models in embedded systems, in which not only storage and memory available are limited but also computing power. This presents challenges in terms of resources, accentuated when deploying multiple models: weights storage and forward pass memory usage and computational complexity.

Multitask Learning [1] learns to solve tasks in parallel using a shared representation of a common input, improving the generalization capabilities of the models. Multitask networks share a base trunk and a number of task-specific branches emerging from it. Only the task-specific layers are computed separately.

In this work, we propose using multitask models to get benefits in terms of speed, memory usage and storage during deployment. For studying our approach, we train and evaluate a multitask model for both semantic segmentation and object detection. We highlight the challenges imposed by applying MTL, explain how they affect the performance of our model and show that it compares positively in terms of inference time, memory usage and model size against deploying one model per task.

II. RELATED WORK

A. Multitask Learning

MTL has been successfully used in different domains, including CV [2], [3]. Some challenges appear when applying it [1]: *learning speed* differences between tasks

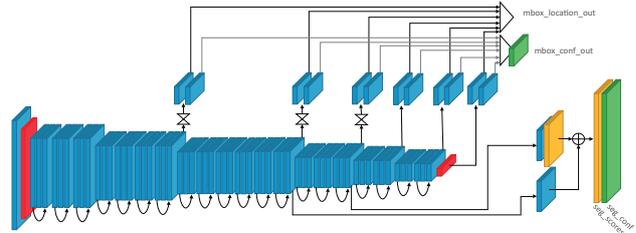


Fig. 1. Our multitask network architecture for object detection and semantic segmentation.

and deciding *what to share* according to the *relatedness* between tasks in the multitask architecture [4], [5].

B. Semantic Segmentation

Semantic segmentation aims at partitioning parts of images belonging to the same semantic class, typically via pixel-wise classification. Fully convolutional networks (FCN) [6] have improved both accuracy and speed for dense prediction problems by using only convolutional layers. Upsampling layers allow a segmentation output size equal to the input and skip connections add finer details. Other approaches add post-processing steps [7], learnable *deconvolution* layers [8] or global context [9].

C. Object Detection

Object detection aims at finding in an image all instances of objects and classifying them in a number of classes. Faster R-CNN [10] was the first to give close to real-time performance. YOLO [11] avoids the generation of region proposals for increased speed. SSD [12] avoids fully-connected layers for speed and takes features at different levels for improved accuracy.

III. METHODS

A. Model details

Extended details on our new model architecture, datasets and training strategies can be found in [13]. We select a ResNet-50 [14] as our base network, FCN with a skip connection for the semantic segmentation task and SSD for objection detection due to their speed vs. accuracy trade-offs. Fig. 1 shows our multitask model.

B. Datasets

For the Pascal VOC dataset [15], we only use the subset of samples that have ground-truth annotations for both tasks from VOC07, VOC12 and SBD[16].

TABLE I
PASCAL VOC RESULTS.

	References			Ours			
	[18]	[12]	[19]	SSD	FCN	Multi	Color
mIoU (%)	-	-	62.5	-	56.4	54.4	55.2
mAP (%)	74.3	70.4	-	51.3	-	51.8	52.6



Fig. 2. Pascal VOC detection and segmentation samples.

As there is no aerial view dataset with annotations for both tasks, we interleave images from Stanford Drone Dataset [17] for object detection and the much smaller Okutama/Swiss [18] datasets for semantic segmentation in each mini-batch.

IV. RESULTS

A. Pascal VOC

We train single-task baselines with the limitations imposed by using MTL and compare those with the multitask model and reference models from the literature. Table I shows the results. Some example output images can be seen in Fig. 2. We add color distortion for data augmentation.

B. Aerial view

We train single-task baselines with no limitations and compare them to our multitask model. Table II shows the results in terms of accuracy and resources used when deployed on an NVIDIA GTX Titan X for the single-task models, their combination and our multitask model. Sample output images can be seen in Fig. 3.

C. Analysis

The compromises made due to the particularities of MTL and especially the lack of a strong data augmentation caused the final accuracy of our multitask model to lag behind that of the single-task ones trained without these although they improved in terms of speed and usage of resources, being 1.6x faster, lighter and consuming less memory than the naive solution. Compared to the single-task models trained with the same limitations, the multitask models matched or outperformed their accuracy for one of the tasks. Find a more detailed analysis in [13].

V. CONCLUSION

We conclude that MTL has the potential to give superior performance in the accuracy vs. speed over using multiple single-task models simultaneously as far as the two analyzed are concerned. This is in exchange of a training process that is more complex as it requires extra choices to be made, a new tuning of the parameters that jointly works well for each task and overcoming new challenges that were not present in the training of single-task models.

TABLE II
AERIAL VIEW RESULTS.

	Base	FCN	SSD	Multitask	Naive
mIoU (%)	-	70.9	-	65.3	70.9
mAP (%)	-	-	54.3	28.2	54.3
Inf. Time (ms)	19	24	27	30	49
Memory (MB)	1203	1233	1525	1552	2511
Size (MB)	95	95	140	140	235



Fig. 3. Aerial view detection and segmentation samples.

ACKNOWLEDGMENT

The work was partially conducted while the first author was at Prendering Lab, participating in the NII International Internship Program in Tokyo, Japan.

REFERENCES

- [1] R. Caruana, "Multitask learning," *Mach. Learn.*, vol. 28, no. 1, pp. 41–75, Jul. 1997.
- [2] I. Kokkinos, "Urbnet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory," *CoRR*, vol. abs/1609.02132, 2016.
- [3] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," *CoRR*, vol. abs/1703.06870, 2017.
- [4] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert, "Cross-stitch networks for multi-task learning," 2016, pp. 3994–4003.
- [5] Y. Lu, A. Kumar, S. Zhai, Y. Cheng, T. Javidi, and R. Feris, "Fully-adaptive feature sharing in multi-task networks with applications in person attribute classification," *arXiv:1611.05377*, 2016.
- [6] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," June 2015.
- [7] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected crfs," *arXiv:1412.7062*, 2014.
- [8] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," December 2015.
- [9] W. Liu, A. Rabinovich, and A. C. Berg, "ParseNet: Looking wider to see better," *CoRR*, vol. abs/1506.04579, 2015.
- [10] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," *CoRR*, vol. abs/1506.01497, 2015.
- [11] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," 2016, pp. 779–788.
- [12] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg, "SSD: single shot multibox detector," *CoRR*, vol. abs/1512.02325, 2015.
- [13] M. Martí, "Multitask deep learning models for real-time deployment in embedded systems," Master's thesis, KTH, Sweden, 2017.
- [14] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *arXiv:1512.03385*, 2015.
- [15] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes challenge: A retrospective," *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, 2015.
- [16] B. Hariharan, P. Arbelaez, L. Bourdev, S. Maji, and J. Malik, "Semantic contours from inverse detectors," 2011.
- [17] A. Robicquet, A. Sadeghian, A. Alahi, and S. Savarese, "Learning social etiquette: Human trajectory understanding in crowded scenes." Springer, 2016, pp. 549–565.
- [18] J. Laurmaa, "A deep learning model for scene segmentation of images captured by drones," Master's thesis, EPFL, Switzerland, 2016.
- [19] J. Mahadeokar, "pynetbuilder," <https://github.com/jay-mahadeokar/pynetbuilder>, 2016.