

A Mobile Cloud Framework for Deep Learning and Its Application to Smart Car Camera

Chien-Hung Chen, Che-Rung Lee^(✉), and Walter Chen-Hua Lu

National Tsing Hua University, Hsinchu 30013, Taiwan
cherung@cs.nthu.edu.tw

Abstract. Deep learning has become a powerful technology in image recognition, gaming, information retrieval, and many other areas that need intelligent data processing. However, huge amount of data and complex computations prevent deep learning from being practical in mobile applications. In this paper, we proposed a mobile cloud computing framework for deep learning. The architecture puts the training process and model repository in cloud platforms, and the recognition process and data gathering in mobile devices. The communication is carried out via Git protocol to ensure the success of data transmission in unstable network environments. We used smart car camera that can detect objects in recorded videos during driving as an example application, and implemented the system on NVIDIA Jetson TK1. Experimental results show that detection rate can achieve four frame-per-second with Faster R-CNN and ZF model, and the system can work well even when the network connection is unstable. We also compared the performance of system with and without GPU, and found that GPU still plays a critical role in the recognition side for deep learning.

1 Introduction

Deep learning has become the most promising machine learning approach. From LeCun's LeNet [1] to Khrizevsky AlexNet [2] and recent GoogleNet [3], deep learning has shown its capability of solving difficult computer vision problems. Its detection performance surpasses other artificial classifiers which rely on the hand-crafted features. With the success in vision work, deep learning also attracts attentions from other fields, such as sentiment analysis [4], language processing [5], region of interest localization and description [6], and medical use [7].

The success of deep learning is brought off by three factors: the advance of numerical optimization methods, growth of data volume, and fast computational hardware. New numerical methods solved the convergence problems, which are more and more critical when the number of layers goes deeper and deeper. Large enough training data sets make the extracted features by deep learning sufficiently representative. These required data can be continuously collected by the sensors equipped in modern embedded systems and mobile devices.

Last, the accelerators, such as Graphic Processing Units (GPUs), provide strong computing power to support the training of deep learning model.

In the era of Internet of Things (IOT), the deployment of deep models to mobile devices is a nature request. On the one hand, the mobile devices can be smarter with the deep learning ability. Moreover, they also help the data gathering from various sources. On the other hand, the limited power, storage, and computational resources of mobile devices prevent the complex computation, like deep learning, from being practical. Therefore, the mobile cloud computing model, which combines the mobility of mobile devices and the computational resources of cloud platforms via ubiquitously accessible network, becomes a practical solution for mobile applications that utilize deep learning.

Mobile cloud computing has three types of models to coordinate the works between mobile devices and cloud platforms. The first one is to off-load all the work to the cloud platforms, and the mobile devices take care of data input and output only. However, such model usually require frequent data communication, which is not suitable for heavy data transmission. In addition, when the network is unstable, mobile devices cannot work alone. The second model relies on mobile devices to handle most of the works. Such model is only limited to light weighted works that can be processed on mobile devices. The last one splits the works to mobile devices and cloud platforms. This model could balance the computation and communication, but requires good synergy from both sides.

Most of the current solution of mobile deep learning utilizes the first solution since the computation of deep learning is heavy. In this paper, we utilized the third model for deep learning, which puts the training and model repository on the cloud platform, and recognition and data gathering on the mobile devices. Such architecture cuts the storage and computation requirement of mobile devices and the data transmission between cloud and end users. In addition, we employed the Git protocol for data communication data so that the transmissions can be resumed even the network connection is not available sometimes.

We used the smart car camera as an example application to demonstrate how the framework works. The smart car camera can select suitable deep learning models for video recognition, decide which part of video clips contain important information and send them to cloud platform, and update the deep learning models when necessary. We have implemented the system on NVIDIA Jetson TK1 embedded development board. Experimental results show that detection rate can achieve 2.8 to 4 FPS (frame-per-second) with Faster R-CNN and ZF model, and the system can work well even when the network connection is unstable. We also compared the performance of system with and without GPU, and found that GPU still plays a critical role in the recognition side for deep learning.

The rest of this paper is organized as follows. Section 2 introduces background knowledge and related works. Section 3 presents the framework and the implementation details of the smart car camera. Section 4 shows the experimental results of system performance. The conclusion and future directions are given in the last section.

2 Background

2.1 Mobile Cloud Computing

In the last decade, mobile devices have been widely used to provide various services, such as gaming, video streaming, health-care, and position-based services. However, mobile devices are still limited by storage capacity, battery life and computing power. As a result, Mobile Cloud computing (MCC) came up not long after cloud computing. Mobile cloud computing can be considered as a combination of mobile services and cloud computing which communicate through wireless network. The key concept of both cloud computing and MCC is offloading heavy tasks to cloud platforms so that users can access a variety of services without complicate computation on their devices.

Mobile Cloud Computing (MCC) architectures possess many desired advantages, as well as many critical issues. In [8], authors listed three merits of MCC: extending battery life, improving data storage capacity and processing power, and improving reliability. On the other hand, MCC has issues to overcome. Wireless network is less reliable and has lower bandwidth compared to wired network. Other than bandwidth issue, network for mobile devices is usually unstable, which could result service unavailable sometimes. Discussion of issues includes offloading, security, heterogeneity are remained in [8].

2.2 Deep Learning

Deep learning, or deep neural network, refers to a neural network that consists of one input layer, one output layer and multiple hidden layers. The earliest appearance of deep neural network can be traced to late 1960s in the work published by Alexey Grigorevich and V.G. Lapa. However, deep learning grows at a slow pace due to immature training scheme and architecture in the next few decades. In the 1990s, LeCun trained LeNet, a deep neural net with convolution and pooling layers, with back-propagation algorithm for digit recognition [1]. Stochastic gradient descent was invented in the 1970s to train artificial neural networks.

LeNet is the earliest work to take deep learning into recognition task. In 2006, layer-wise unsupervised pre-training was proposed for parameter initialization [9]. The new training method has two phases. In the pre-training phase, every two consecutive hidden layers are considered a restricted Boltzmann machine and weights are adjusted with an unsupervised learning manner. After pre-training, the whole model is fine-tuned with an additional supervised layer in the second phase. Pre-training makes layer parameters get better values compared to random initialization. As a result, trained models reach more sensible local minimum and quality gets more stable. In 2012, Krizhevsky et al. developed the eight-layer AlexNet and won the ILSVRC 2012 prize with about 85% classification and 66% location accuracy [2]. AlexNet won their competitors who used linear classifier over ten percentage. And the winners of ILSVRC classification

task in the following years all used deep neural networks. Instead of using hand-crafted classifiers, deep neural network learns high level features during training time. And the ILSVRC challenge results prove its high performance.

With big success in image classification, deep learning attracts focus from other fields. Beside from classification, deep learning is also used in object detection [10–12], speech recognition [13, 14], language processing [5], sentiment analysis [4], and many other works. Recent deep learning bloom can be ascribed to new optimization methods, appearance of more powerful processor and rapid growing data volume. With more powerful CPU and multi-core processor, especially general purpose GPU, training time can be cut down from months to days. With growing of various data, under-fitting can be prevented and makes deep learning be applied to solve different types of problem. Latest publications not only improve accuracy but boost performance. PReLU activation function pushes classification accuracy over 95% and saves time for deriving gradient [15]. Dropout prevents training from over-fitting [16]. Also, new initialization schemes make pre-training phase unnecessary [17].

2.3 Git Version Control System

Git is an open source version control system. The designated work space to be managed is called repository. The system would manage files added to tracked list in the repository. Users can copy files from remote directory through clone command. Cloned repository would be managed by Git system and repository status be kept consistent with remote if git pull command is called. Push command would be used in the situation that a user makes some modification and wants to update file in original repository. Modification of tracked file would not be send back to remote after a commit is created. And according to the official report, Git is faster than many other version control systems.

3 Framework and Implementation

3.1 Mobile Cloud Framework for Deep Learning

Our design of mobile cloud computing framework for deep learning is based on two facts. The first one is about deep learning processes, and the second one is for mobile network. Deep learning approach has two phases: training and recognition. In the recognition phases, deep learning application only needs to go through forward propagation with input data. After getting output data, some post-processes are required to retrieve desired information. On the contrary, the training phase needs both forward and backward propagation and a large amount of training data. As a result, training process usually takes days or weeks.

For mobile network, according to 4GTest [18], LTE network has a median of 5.64 Mbps upload bandwidth. Take image process for consideration, common car cameras have 1080p resolution, which also indicates that a raw image can have size overs 6 megabytes. Retrieving results after uploading and processing

an image on the cloud becomes infeasible. Even we compress files first, unstable wireless network may prevents users some services.

Our design offloads the training phase to the cloud platforms and keeps trained models to the mobile devices for recognition. In addition, the mobile devices collect received data and feed back to cloud. Various data collected from users can help to improve accuracy of deep learning models. Also, the server side can evaluate how the model performs by analyzing prediction result.

3.2 Smart Car Camera System

We used car camera object detection as an example application. Many cars have equipped cameras to record videos during driving. However, the capacity of device storage is limited, and many parts of the recorded videos do not have interested objects. With the object detection ability, the car camera can find out the video segments that possibly have objects of interests. Those interested video segments can be uploaded to cloud platforms for further analysis, and the unwanted video segments can be deleted to save storage space. We designed a deep learning object detection system which not only provides detection service but handles model maintenance.

The designed system is presented in Fig. 1. When the system starts, version checker is executed first. If the client does not have model files in local disk storage, it will connect to server and get the desired model. Otherwise the system would check whether an updated version is available. In the second phase, the system has two threads: update thread and detection thread. The update thread keeps checking new model released periodically. This thread is only responsible for model checking and downloading. Models would not be replaced with the updated version until system restarts for safety concern. The detection thread handles data collection. Detection thread would collect network inputs and outputs if possible object appears in the input images. Meanwhile, it keeps a counter to control collected file size. Once the counter hits a designated threshold, the detection thread creates another thread for uploading data. The collected data after upload would be destroyed to save local storage space.

3.3 Detection Task

We have surveyed some object detection works, and chose Faster R-CNN as our detection approach. We have re-produced the C version Faster R-CNN using Caffe [19] and openCV. Figure 2 introduces the workflow of Faster R-CNN. Although Faster R-CNN can take any size of input, we cropped and resized every input image to a fixes size: 224×224 pixels to make it faster. The following section describes Faster R-CNN implementation in details.

Faster R-CNN trains RPN by sliding a small network over the feature map from the last shared convolutional layer. At each location on the feature map, RPN predicts 3 scales and 3 aspect ratios (the ratio of width to height). Centers and sizes of each box is encoded based on a reference box, called anchor box.

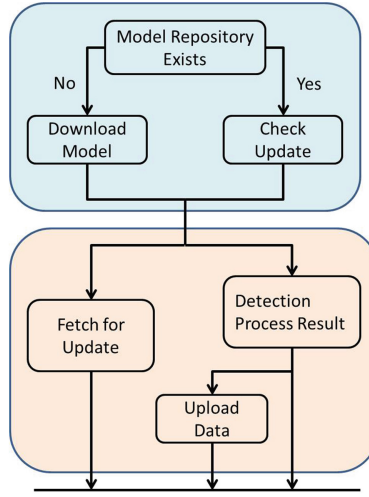


Fig. 1. Designed system

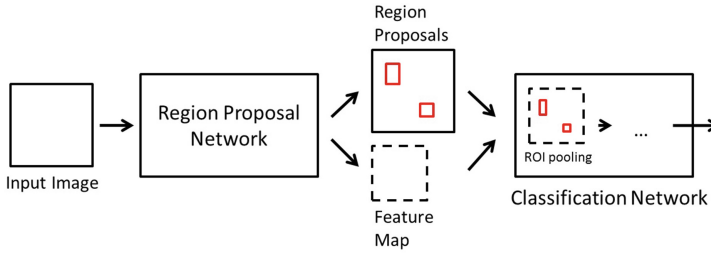


Fig. 2. Faster R-CNN workflow

At detection time, region proposal network (RPN) outputs two data blobs, predicted boxes and box scores, after a series of convolutional and pooling processes. Prediction box is represented by four elements (N, C, H, W) , where N represents the batch size, C stores proposed boxes, and H and W correspond to the height and the width of feature map. After extracting proposal regions, we drop boxes that are too small and crop the remaining boxes for ensuring that the box size does not exceed input image size. Next, we sort the boxes by confidence c from score blob and apply non-maximal suppression to filter out highly overlapped boxes with lower confidence score.

Classification net takes n filtered regions and feature map generated by the last convolutional layer of region proposal network (RPN) as inputs, and outputs two blobs. Predicted boxes are stored in the same fashion as RPN's output in the first blob. The second blob contains prediction scores of background and desired object classes. If the number of region proposals from RPN exceeds the batch size N , the classification net will run more than one time. The batch size N is

adjustable. In Sect. 4.2, we will discuss how to tune the batch size to optimize performance.

Besides, we trained the models using py-faster-rcnn with VOC2007 dataset. We trained ZF model with provided definition files and using the alternating optimization [20]. Different from previously described Faster R-CNN, the model takes over intermediate output, which does not require coordinate transformation between RPN and the following layers. Also, the output predicted boxes information are not related with anchor boxes. Instead, it retrieves coordinates of box vertices from the output blob. As to box regression, refined information can be calculated in the same manner with original Faster R-CNN.

3.4 Model Maintenance

We chose Git version control protocol for model transmission. By using Git, we save the effort to manage modifications. Also, Git provides the basic authentication service, by which users need account and password to access remote repository.

First time download and follow-up model update can be achieved by Git clone and Git pull functions. And the feedback part can be implemented by Git push. We created two Git repositories on our device: one for storing model files and the other for feedback data. When starts, the system checks whether the model repository exists. If files in need have not been downloaded, Git clone would be triggered to get the latest files from server. Otherwise, the system calls update function to keep version be consistent with the server. In the detection phase, update-thread calls Git fetch constantly and push-thread starts to push feedback repository once collected data reaches a threshold.

4 Performance Evaluation

We have implemented the smart car camera system using Nvidia Tegra K1, the processor and GPU designed for mobile device. Tegra K1 has four plus one core ARM Cortex-A15 CPU with upto 2.3 GHz frequency and a Kepler GPU. The experimental board has a 2 GB memory, and is connected to wired network. Data can be transmitted through WAN. We chose relatively light weight ZF model provided by origin work to do the detection task. ZF net has 5 convolutional layers and 3 fully connected layers. The trained model can be obtained from [20].

We evaluated our work on the trained dataset PASCAL VOC 2007 test data. In addition, we used pedestrian detection database from University of Pennsylvania and the clips downloaded from Youtube, which are videos recorded by car camera on highway.

We have four sets of experiments. The first set of experiments compares the performance with and without GPU. The second one is to tune the batch size for the best performance of the model. The third one compares the system performance for different data sets and different batches. The last one uses simulation to demonstrate how our system works under unstable network environments.

4.1 CPU and GPU Performance Comparison

We evaluated CPU and GPU performance for recognition process. Figure 3 shows average CPU and GPU forward pass time. As can be seen, GPU accelerates the forward propagation and performs over 20 times faster than CPU only version, both for RPN and classification net. The significant difference indicates the importance of accelerator. For the rest of our experiments, we only showed the performance results with GPU.

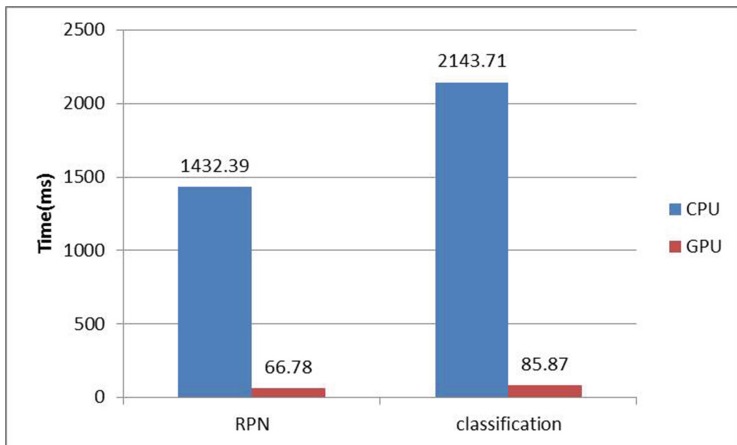


Fig. 3. CPU and GPU performance

4.2 Batch Size Tuning

In the recognition phase, the batch size, which is the number of Region Of Interests (ROI) to be processed at the same time, influences the performance significantly. In this experiment, we evaluated the time (seconds) of forward pass of classification net for different batch size. As Fig. 4 shows, execution time of one iteration grows as the batch size gets larger. This is not a surprised results since larger batch size means more work to do. However, if we normalized the execution time for 180 frames, the time decreases first and then curves up later. This is because when the batch size grows, the number of iterations to process the ROIs of 180 frames decreases, but the time of each iteration increases. The experimental results show the best performance falls around batch size 60, which is the parameter we used in the later experiments.

4.3 System Performance

We evaluated the system performance for different number of batches to be processed for the classification net. Since we resized and cropped input images,

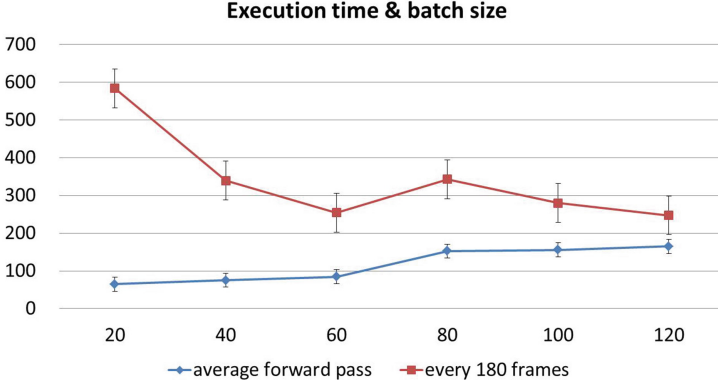


Fig. 4. Performance and batch size

execution time only relates to the number of candidate Region of Interests (ROI) per image. We can accelerate the system performance by dropping the ROIs with lower scores. The influence of dropping is that some of the objects may not be detected. But experiments show the dropping of low score ROIs does not alter the results much (less than 1%). Each frame usually has 60 to 180 ROIs. Since we let batch size be 60, the number of batches per frame is ranged from 1 to 3. Figure 5 shows the performance result, measured by frame-per-second, for different number of batches. The 1 batch means that we only keep at most 60 ROIs per frame; the 2 batch means each frame can have up to 120 ROIs; and the 3 batch means 180 ROIs are detected for each frame. The results show that the detection task could process over 4 frames per second for 1 batch. When the number of batches grows, the performance drops to around 3 frames per second. The results of two test data have similar behaviors.

4.4 Network Transmission

Our smart car camera system can dynamically download and update the required models from cloud platform. According to [18], the downlink of current LTE network is much faster than that of the uplink. However, the model files of deep learning usually are usually very large, around hundred megabytes, even with compression. In unstable network environments, transmission of such files could cause a problem, because file transmission needs to restart after re-connection. One solution to reduce the duplicate transmission is to split whole model into chunks with smaller file size. With this approach, duplicate transmitted data size would be limited to the chunk size. On the other side, it needs longer preparation time for file splitting, but that can be done in advance on cloud platform.

We used the Markov chain model, proposed in [21], to simulate the package loss to evaluate the performance of the file splitting strategy. The maximum number of retransmission in the simulation is 15, after which the file chunk is required to retransmit. The bandwidth of uplink in simulation is 15 Mbps and

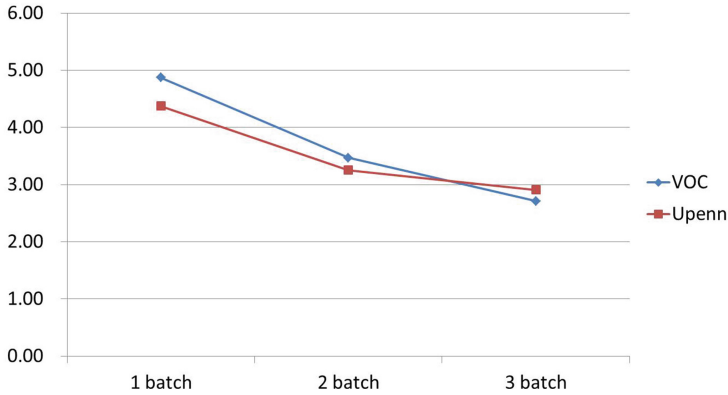


Fig. 5. Performance measured by frames per second

the model file size is 224MB, which is a real data size of ZF model for 20 categories. The latency of each split file transmission is 0.1 s. The performance of three chunk sizes (number of chunks) are compared: 1 (no slitting), 20, and 100. The evaluated package drop rates are 0.1, 0.4, 0.5, and 0.6.

Table 1. The transmission time (seconds) of different chunk sizes for different package drop rates.

Drop rate\No. chunks	1	20	100
0	120	121	130
0.1	133	134	143
0.4	208	202	210
0.5	4382	271	256
0.6	(unable)	2877	451

Table 1 summarizes the transmission time of each method in seconds for different package drop rate. The result is the average of 10 times of simulation. It shows when the package drop rate is small, which means the network environment is stable, the splitting method takes longer time owing to the overhead of file transmission latency. However, for larger package drop rate, the time of un-split file grows rapidly. For drop rate 0.6, the un-split file cannot be completed in hours. On the other hand, for chunk number 100, the time grows much slower even when the network is extremely unstable.

5 Conclusion

In this paper, we proposed a mobile cloud computing framework for deep learning. The architecture leaves training on the cloud and performs recognition in

the client side with some post processing. The advantage is that the mobile devices can still work under unstable network environment without huge storage and computing capacity. Experiments show that GPU is still critical for the recognition of deep learning, since it can accelerate the computation over 20 times.

Limited storage space and computing power remain a challenge for embedded systems using deep learning. S. Han et al. proposed DeepCompression [22], which reduces storage space by learning only important connections, quantize weights and apply Hoffman encoding. They reduced size of AlexNet from 240 MB to 6.9 MB without accuracy loss. With deeper models applied on mobile devices. For future work, we will integrated their work to the system, and investigate further optimization of storage and computation for deep learning.

Acknowledgment. This study is conducted under the The Core Technologies of Smart Handheld Devices (3/4) of the Institute for Information Industry; which is subsidized by the Ministry of Economy Affairs, Taiwan. The authors thank the Institute for Information Industry for the financial support under grant number 105-EC-17-A-24-0691.

References

1. LeCun, Y., Jackel, L., Bottou, L., Brunot, A., Cortes, C., Denker, J., Drucker, H., Guyon, I., Muller, U., Sackinger, E., et al.: Comparison of learning algorithms for handwritten digit recognition. In: International Conference on Artificial Neural Networks, vol. 60, pp. 53–60 (1995)
2. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems, pp. 1097–1105 (2012)
3. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1–9 (2015)
4. Socher, R., Perelygin, A., Wu, J.Y., Chuang, J., Manning, C.D., Ng, A.Y., Potts, C.: Recursive deep models for semantic compositionality over a sentiment treebank. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), vol. 1631, p. 1642. Citeseer (2013)
5. Collobert, R., Weston, J.: A unified architecture for natural language processing: deep neural networks with multitask learning. In: Proceedings of the 25th International Conference on Machine Learning, pp. 160–167. ACM (2008)
6. Johnson, J., Karpathy, A., Fei-Fei, L.: Denscap: fully convolutional localization networks for dense captioning. arXiv preprint [arXiv:1511.07571](https://arxiv.org/abs/1511.07571) (2015)
7. Fakoor, R., Ladhak, F., Nazi, A., Huber, M.: Using deep learning to enhance cancer diagnosis and classification. In: Proceedings of the International Conference on Machine Learning (2013)
8. Dinh, H.T., Lee, C., Niyato, D., Wang, P.: A survey of mobile cloud computing: architecture, applications, and approaches. *Wirel. Commun. Mob. Comput.* **13**, 1587–1611 (2013)

9. Bengio, Y., Lamblin, P., Popovici, D., Larochelle, H., et al.: Greedy layer-wise training of deep networks. *Adv. Neural Inform. Process. Syst.* **19**, 153 (2007)
10. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 580–587 (2014)
11. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: unified, real-time object detection. *arXiv preprint [arXiv:1506.02640](https://arxiv.org/abs/1506.02640)* (2015)
12. Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., LeCun, Y.: Overfeat: integrated recognition, localization and detection using convolutional networks. *arXiv preprint [arXiv:1312.6229](https://arxiv.org/abs/1312.6229)* (2013)
13. Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A.R., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T.N., et al.: Deep neural networks for acoustic modeling in speech recognition: the shared views of four research groups. *IEEE Sig. Process. Mag.* **29**, 82–97 (2012)
14. Dahl, G.E., Yu, D., Deng, L., Acero, A.: Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition. *IEEE Trans. Audio Speech Lang. Process.* **20**, 30–42 (2012)
15. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1026–1034 (2015)
16. Srivastava, N., Hinton, G.E., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958 (2014)
17. Sutskever, I., Martens, J., Dahl, G.E., Hinton, G.E.: On the importance of initialization and momentum in deep learning. *ICML* **28**(3), 1139–1147 (2013)
18. Huang, J., Qian, F., Gerber, A., Mao, Z.M., Sen, S., Spatscheck, O.: A close examination of performance and power characteristics of 4G LTE networks. In: *Proceedings of the 10th International Conference on Mobile Systems, Applications, and Services*, pp. 225–238. *ACM* (2012)
19. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: convolutional architecture for fast feature embedding. *arXiv preprint [arXiv:1408.5093](https://arxiv.org/abs/1408.5093)* (2014)
20. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: *Advances in Neural Information Processing Systems*, pp. 91–99 (2015)
21. Sanneck, H.A., Carle, G.: Framework model for packet loss metrics based on loss runlengths. In: *Proceedings of the SPIE, Multimedia Computing and Networking 2000*, vol. 3969 (1999)
22. Han, S., Mao, H., Dally, W.J.: Deep compression: compressing deep neural network with pruning, trained quantization and huffman coding. *CoRR*, abs/1510.00149 **2** (2015)

Internet of Vehicles – Technologies and Services
Third International Conference, IOV 2016, Nadi, Fiji,
December 7–10, 2016, Proceedings
Hsu, C.-H.; Wang, S.; Zhou, A.; Shawkat, A. (Eds.)
2016, XV, 250 p. 119 illus., Softcover
ISBN: 978-3-319-51968-5