



ARL-TN-0781 • SEP 2016



# Mission Driven Scene Understanding: Candidate Model Training and Validation

by Arnold Tunick

Approved for public release; distribution unlimited.

## **NOTICES**

### **Disclaimers**

The findings in this report are not to be construed as an official Department of the Army position unless so designated by other authorized documents.

Citation of manufacturer's or trade names does not constitute an official endorsement or approval of the use thereof.

Destroy this report when it is no longer needed. Do not return it to the originator.



# **Mission Driven Scene Understanding: Candidate Model Training and Validation**

**by Arnold Tunick**

*Computational and Information Sciences Directorate, ARL*

**REPORT DOCUMENTATION PAGE**

*Form Approved  
OMB No. 0704-0188*

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing the burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

**PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.**

<b>1. REPORT DATE (DD-MM-YYYY)</b> September 2016		<b>2. REPORT TYPE</b> Technical Note		<b>3. DATES COVERED (From - To)</b> 10/2015 – 08/2016	
<b>4. TITLE AND SUBTITLE</b> Mission Driven Scene Understanding: Candidate Model Training and Validation				<b>5a. CONTRACT NUMBER</b>	
				<b>5b. GRANT NUMBER</b>	
				<b>5c. PROGRAM ELEMENT NUMBER</b>	
<b>6. AUTHOR(S)</b> Arnold Tunick				<b>5d. PROJECT NUMBER</b>	
				<b>5e. TASK NUMBER</b>	
				<b>5f. WORK UNIT NUMBER</b>	
<b>7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)</b> US Army Research Laboratory ATTN: RDRL-CII-A 2800 Powder Mill Road Adelphi, MD 20783-1138				<b>8. PERFORMING ORGANIZATION REPORT NUMBER</b>  ARL-TN-0781	
<b>9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)</b>				<b>10. SPONSOR/MONITOR'S ACRONYM(S)</b>	
				<b>11. SPONSOR/MONITOR'S REPORT NUMBER(S)</b>	
<b>12. DISTRIBUTION/AVAILABILITY STATEMENT</b> Approved for public release; distribution unlimited.					
<b>13. SUPPLEMENTARY NOTES</b>					
<b>14. ABSTRACT</b> Army missions take place in dynamic environments, where changing illumination, precipitation, and vegetation can modify saliency and context of an outdoor scene, obscure features, and degrade object recognition. For Army missions, scene understanding tools need to account for dynamic environments that change as a function of space and time and should be tested in mission simulating conditions. In addition, the impact of dynamic environments should be included in the scene understanding approach. At this stage, we are evaluating different computational frameworks that may be useful to incorporate dynamic environments into mission driven scene understanding. One of the candidate engines that we are evaluating is a convolutional neural network (CNN) program installed on a Windows 10 notebook computer. In this report, we present progress toward the proof-of-principle testing of the candidate model to examine the impact of dynamic environments on scene understanding model results.					
<b>15. SUBJECT TERMS</b> computer vision, context, saliency, convolutional neural network, dynamic environments					
<b>16. SECURITY CLASSIFICATION OF:</b>			<b>17. LIMITATION OF ABSTRACT</b>  UU	<b>18. NUMBER OF PAGES</b>  18	<b>19a. NAME OF RESPONSIBLE PERSON</b> Arnold Tunick
<b>a. REPORT</b> Unclassified	<b>b. ABSTRACT</b> Unclassified	<b>c. THIS PAGE</b> Unclassified			<b>19b. TELEPHONE NUMBER (Include area code)</b> (301) 394-1233

## **Contents**

---

---

<b>List of Figures</b>	<b>iv</b>
<b>List of Tables</b>	<b>iv</b>
<b>Acknowledgments</b>	<b>v</b>
<b>1. Introduction</b>	<b>1</b>
<b>2. CNN Deep Learning Libraries and Open Source Codes</b>	<b>2</b>
<b>3. Candidate CNN Model: Training and Validation</b>	<b>3</b>
<b>4. Candidate CNN Model: Results</b>	<b>4</b>
<b>5. Summary and Conclusions</b>	<b>6</b>
<b>6. References</b>	<b>7</b>
<b>Distribution List</b>	<b>9</b>

## List of Figures

---

---

Fig. 1	For Army mission activities, the impact of dynamic environments should be included in the scene understanding approach (e.g., space-time coordinates, weather conditions and trends, visibility, terrain, scene descriptors, anomalies, and other salient features) (data from the ImageNet Large Scale Visual Recognition Challenge 2012 [ILSVRC2012]) .....2
Fig. 2	Candidate CNN training and validation. Top-1 training accuracy (red line). Top-1 validation accuracy (blue diamonds). .....4
Fig. 3	Candidate CNN results showing the top-5 most likely classification labels and corresponding top-5 confidence levels .....5

## List of Tables

---

---

Table 1	CNN deep learning libraries: open source frameworks .....3
Table 2	CNN open-source codes .....3

## **Acknowledgments**

---

I thank RE Meyers, G Warnell, P David, B Byrne, C Karan, and S Gutstein for helpful discussions. This research was supported by the US Army Research Laboratory.

INTENTIONALLY LEFT BLANK.



## 1. Introduction

---

Rapid and robust scene understanding is a critically important goal for the development of Army autonomous intelligent systems to support the Army mission.<sup>1</sup> Army missions take place in dynamic environments, where changing illumination, precipitation, and vegetation can modify saliency and context of an outdoor scene, obscure features, and degrade object recognition. For Army missions, scene understanding tools need to account for dynamic environments that change as a function of space and time and should be tested in mission simulating conditions. In addition, the impact of dynamic environments should be included in the scene understanding approach.<sup>2</sup> Image features that can potentially help the mission are relevant. For example, important image features may be related to space-time coordinates, weather conditions and trends, visibility, terrain, scene descriptors, anomalies, and other salient features.<sup>3-5</sup>

To explore the impact of dynamic environments on scene understanding, we need a computational engine for scene exploration of new images. At this stage, we are evaluating different computational frameworks that may be useful to incorporate dynamic environments into mission driven scene understanding. One of the candidate engines that we are evaluating is a convolutional neural network (CNN) program (i.e., Theano-AlexNet<sup>6,7</sup>) installed on a Windows 10 notebook computer. To the best of our knowledge, an implementation of the open-source, Python-based AlexNet CNN on a Windows notebook computer has not been previously reported.

In this report, we present progress toward the proof-of-principle testing of the candidate CNN model to examine the impact of dynamic environments on scene understanding model results. While we found previously<sup>5</sup> that the CNN was able to determine the correct class labels for images taken from the 2,560 image training data set, the validation process did not appear to provide optimal results for images not previously seen. As a result, we performed additional trials and analysis using the larger ImageNet<sup>8</sup> data set containing approximately 1.2 million images (Fig. 1). In Section 3, we show that the CNN achieved 79.7% validation accuracy for the top-5 class labels, which is in close agreement with results published by its developers.

We start our discussion by presenting an overview of representative deep learning libraries (i.e., available open-source computational engines/frameworks) as well as a summary of several current CNN open-source codes.



**Fig. 1** For Army mission activities, the impact of dynamic environments should be included in the scene understanding approach (e.g., space-time coordinates, weather conditions and trends, visibility, terrain, scene descriptors, anomalies, and other salient features) (data from the ImageNet<sup>8</sup> Large Scale Visual Recognition Challenge 2012 [ILSVRC2012])

## **2. CNN Deep Learning Libraries and Open Source Codes**

---

CNN deep learning methods have influenced and advanced many applications in computer vision, especially those related to image classification.<sup>6,8</sup> A recent paper by Bahrampour et al.<sup>9</sup> presented a comparative study of 5 current deep learning software frameworks with regard to their capability to incorporate different types of CNN architectures, their hardware usage (central processing unit [CPU] and graphical processing unit [GPU]), and an evaluation of their training/testing speed. We present a summary of these open-source libraries, as well as 3 additional frameworks, in Table 1, to include a listing of the principal software developers, the primary programming language used, the Internet location of the open-source codes, the Internet location of installation and user's guide documentation, and key reference citations. Similarly, Table 2 presents a summary of representative CNN open-source codes to include the candidate CNN program<sup>7</sup> that we trained and validated. Note that in Table 2, some CNN codes achieve better validation accuracy or train at greater speeds than AlexNet<sup>6,7</sup> in Theano<sup>10</sup>, particularly those associated with the Caffe<sup>11</sup> and Computational Network Toolkit (CNTK)<sup>14,15</sup> frameworks. Nevertheless, Bahrampour et al.<sup>9</sup> commented that the Theano-based libraries and

codes benefit from the flexibility and ease in development using the Python language. In contrast, the primary programming language for Caffe and CNTK is C<sup>++</sup>.

**Table 1 CNN deep learning libraries: open source frameworks**

Name	Developer	Language	Availability	Documentation	Computation	Key Reference
Caffe	Berkeley Vision and Learning Center	C++ with Python/MATLAB wrappers	github.com/BVLC/caffe	Tutorial and installation guide: <ul style="list-style-type: none"> <li>caffe.berkeleyvision.org/tutorial/</li> <li>caffe.berkeleyvision.org/installation.html</li> </ul>	Support for CPU, GPU	Jia et al. <sup>11</sup>
Torch	R Collobert C Farabet K Kavukcuoglu S. Chintala	Lua	github.com/torch/torch7	Tutorials, demos, examples, developer guide: <ul style="list-style-type: none"> <li>www.torch.ch</li> </ul>	Support for CPU, GPU	Collobert et al. <sup>12</sup>
Theano	The Theano Development Team	Python	github.com/Theano/Theano	Tutorial and installation guide: <ul style="list-style-type: none"> <li>deeplearning.net/software/theano/</li> </ul>	Support for CPU, GPU	Al-Rfou et al. <sup>10</sup>
TensorFlow	Google	C++, Python	github.com/tensorflow/tensorflow	Installation and user's guide: <ul style="list-style-type: none"> <li>tensorflow.org/</li> </ul>	Support for CPU, GPU	Abadi et al. <sup>13</sup>
CNTK	Microsoft (Computational Network Toolkit)	C++	github.com/Microsoft/CNTK	Tutorial, setup, examples: <ul style="list-style-type: none"> <li>www.cntk.ai/</li> <li>github.com/Microsoft/CNTK/wiki</li> </ul>	Support for CPU, GPU	Yu et al. <sup>14,15</sup>
Neon	Nervana Systems	Python	github.com/nervanasystems/neon	Tutorial: <ul style="list-style-type: none"> <li>neon.nervanasys.com/docs/latest/tutorials.html</li> </ul>	Support for CPU, GPU	Neon DLL <sup>16</sup>
Deeplearning4j	SkyMind	Java, Scala	github.com/deeplearning4j/deeplearning4j	Tutorial and user's guide: <ul style="list-style-type: none"> <li>deeplearning4j.org/</li> </ul>	Support for CPU, GPU	Deeplearning4j <sup>17</sup>
VLFeat	A Vedaldi B Fulkerson	C with Matlab interface	www.vlfeat.org/install-matlab.html	User's manual: <ul style="list-style-type: none"> <li>http://www.vlfeat.org/matconvnet/matconvnet-manual.pdf</li> </ul>	Support for CPU, GPU	Vedaldi and Fulkerson <sup>18</sup>

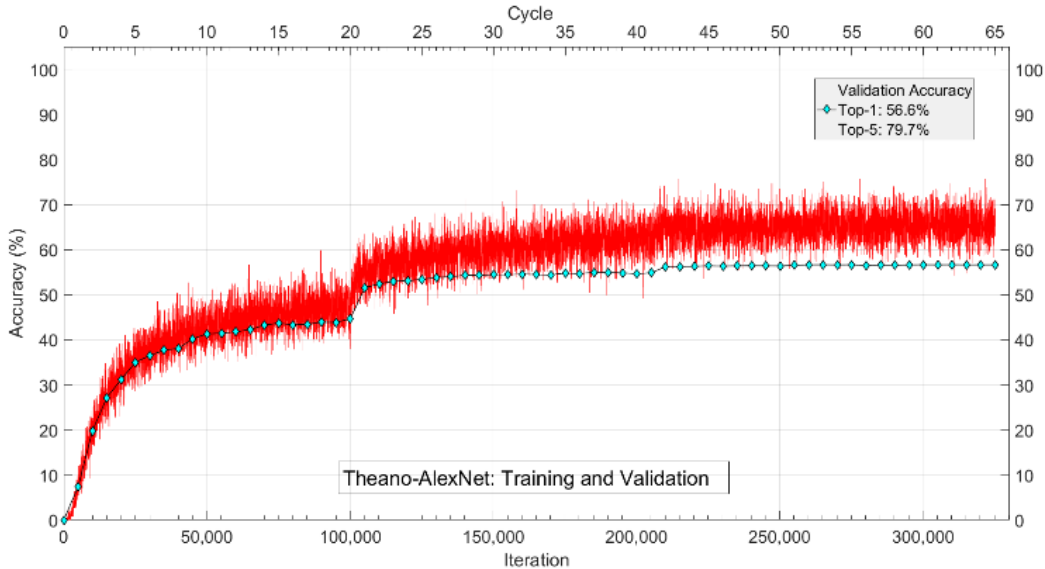
**Table 2 CNN open-source codes**

Name	Developer	Language	Availability	Top-5 Accuracy	Reference
Cuda-Convnet (in Caffe)	A Krizhevsky I Sutskever GE Hinton	Python	github.com/akrizhevsky/cuda-convnet2	81.8 % (2 GPUs)	Krizhevsky et al. <sup>6</sup>
AlexNet (in Theano)	W Ding, R Wang, F Mao, G Taylor	Python	github.com/uoguelph-mlrg/theano_alexnet	80.1% (2 GPUs)	Ding et al. <sup>7</sup>
GoogLeNet (in Caffe)	Google	C++ with Python/MATLAB wrappers	github.com/BVLC/caffe/tree/master/models/bvlc_googlenet	93.3 % (CPU)	Szegedy et al. <sup>19</sup>
VGG (in Caffe)	K Simonyan A Zisserman	C++	www.robots.ox.ac.uk/~vgg/research/very_deep/	92.5% (4 GPUs)	Simonyan and Zisserman <sup>20</sup>
OverFeat (in Torch)	NYU Computational Intelligence, Learning, Vision, and Robotics Lab	C++ with Python/Lua wrappers	github.com/semanet/OverFeat civvr.nyu.edu/doku.php?id=software:overfeat:start	86.4% (1 GPU)	Sermanet et al. <sup>21</sup>
Matconvnet (in VLFeat)	A Vedaldi K Lenc	Matlab	github.com/vlfeat/matconvnet	~80% (1 GPU)	Vedaldi and Karel <sup>22</sup>

### 3. Candidate CNN Model: Training and Validation

In this section, we present the candidate CNN model training and validation results that were achieved implementing the program code on a Windows 10 notebook computer using a single GPU. A description of the installed software and dependencies was given in in a previous report<sup>5</sup> and is therefore not repeated here.

The CNN was executed for 65 epoch (i.e., cycles), wherein 5,004 mini-batches of 256 images were processed for each training cycle. Here, we used image data from the ImageNet<sup>8</sup> Large Scale Visual Recognition Challenge 2012 (ILSVRC2012). A few example images from the training data set are shown in Fig. 1. On average, training on 20 mini-batches (or iterations) took approximately 172 s. As a result, the entire 65 cycles took approximately 32 days to complete.\* Even though the training time was long, the CNN achieved 56.6% validation accuracy for the top-1 class labels and 79.7% accuracy for the top-5 class labels (Fig. 2). These results are in close agreement with those reported by Krizhevsky et al.<sup>6</sup> (i.e., a top-5 accuracy of 81.8%) and Ding et al.<sup>7</sup> (i.e., a top-5 accuracy of 80.1%). Thus, in the next section, we show our initial top-5 class label results achieved from testing the CNN with 4 single images gleaned from the training data set.



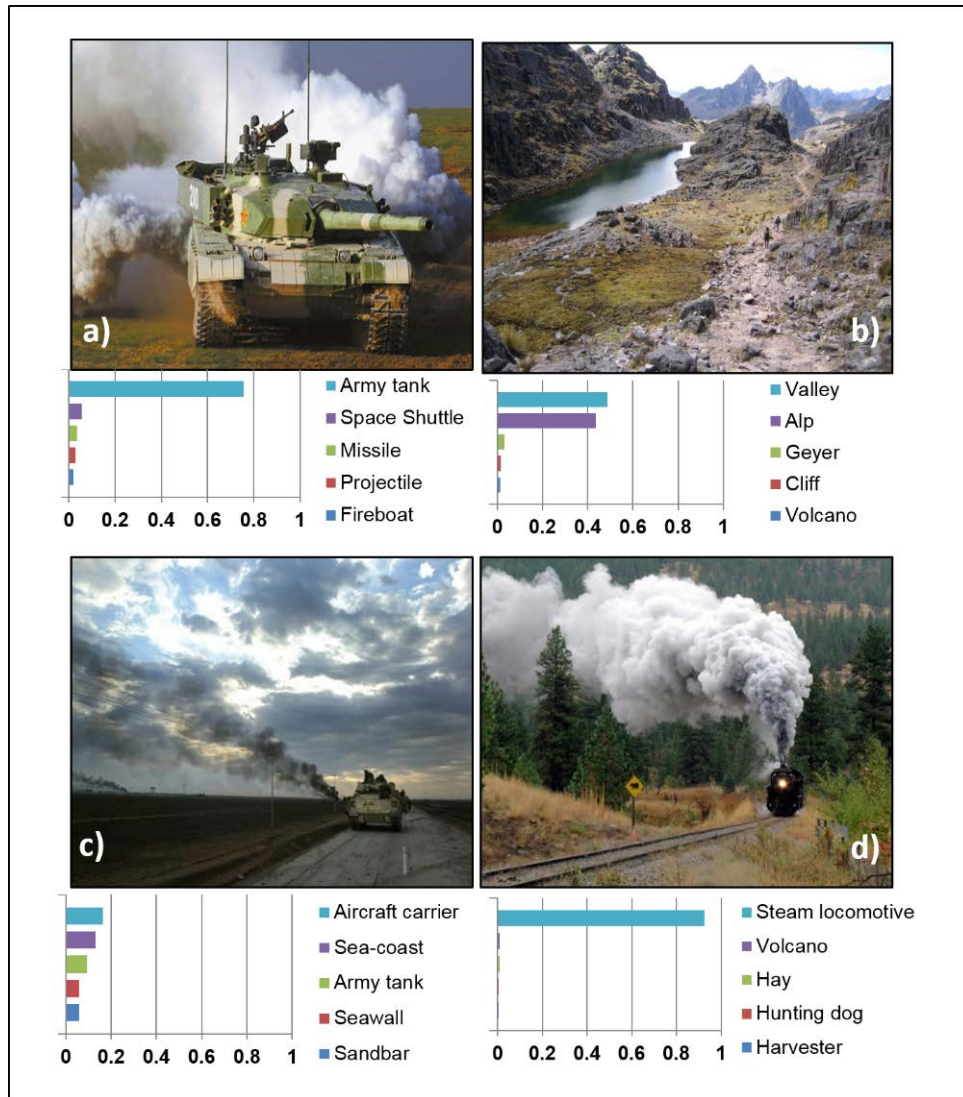
**Fig. 2** Candidate CNN training and validation. Top-1 training accuracy (red line). Top-1 validation accuracy (blue diamonds).

## 4. Candidate CNN Model: Results

In this section, we test the candidate model to examine the impact of dynamic environments on scene understanding model results. For the example shown in Fig. 3, we had the model output the top-5 most likely classification labels and corresponding confidence levels (i.e., top-5 probabilities) for the 4 images shown in Fig. 1. To do this, we modified the model code and incorporated an inference calculation<sup>23</sup> to extract the desired results. We found that the CNN predicted the

\* For comparison, Ding et al.<sup>7</sup> reported training times of about 40–49 s per 20 iterations for 1 GPU (e.g., approximately 9 days to complete 65 cycles) and 24–29 s per 20 iterations for 2 GPUs.

correct class label for the principal object(s) shown in the test images, generally with high confidence. Nevertheless, a person viewing these images (e.g., a Soldier-in-the-loop) would likely see several additional features, such those related to the environment that were not identified (e.g., clouds, haze, smoke plumes, sandy soil, rocky terrain, mountains, river water, trees, and forests). More importantly though, we noticed that, in Fig. 3c, low light and visibility conditions negatively affected the candidate model results (i.e., much lower probabilities). Hence, it is this kind of adverse impact on scene understanding model results that require further testing (e.g., with sets of new images that contain similar objects, but depict a wide variety of relevant dynamic environment features).



**Fig. 3** Candidate CNN results showing the top-5 most likely classification labels and corresponding top-5 confidence levels



## 5. Summary and Conclusions

---

Two key aspects of scene understanding modeling are readily apparent from our research so far:

- 1) Scene understanding tools need to account for dynamic environments to better support Army missions performed by autonomous intelligent systems, and
- 2) Images depicting adverse dynamic environment features (e.g., low visibility and illumination) tend to negatively impact the scene understanding model results.

We can conduct further testing of candidate models to quantify these aspects in more detail. Nevertheless, it is clear that improved or retrained models are needed to better address the impact of dynamic environments on mission driven scene understanding.

## 6. References

---

---

1. Army Research Laboratory (US). US Army Research Laboratory S&T Campaign Plans 2015–2035. Adelphi (MD): Army Research Laboratory (US); 2014. System Intelligence & Intelligent Systems; p. 98–99.
2. Meyers RE. Army Research Laboratory (US), Adelphi, MD. Personal communication, 2014.
3. Tunick A. Space and time scale characterization of image data in varying environmental conditions for better scene understanding. Adelphi (MD): Army Research Laboratory (US); 2015 Sep. Report No.: ARL-TR-7419.
4. Tunick A. Space–time environmental image information for scene understanding. Adelphi (MD): Army Research Laboratory (US); 2016 Apr. Report No.: ARL-TR-7652.
5. Tunick A. Mission driven scene understanding: Dynamic environment. Adelphi (MD): Army Research Laboratory (US); 2016 Jun. Report No.: ARL-TN-0764.
6. Krizhevsky A, Sutskever I, Hinton GE. ImageNet classification with deep convolutional neural networks. Proc. Advances in Neural Info Processing Sys. 2012;25.
7. Ding W, Wang R, Mao F, Taylor G. Theano-based large-scale visual recognition with multiple GPUs. arXiv:1412.2302v4; 2015.
8. Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg AC, Fei-Fei L. ImageNet large scale visual recognition challenge. arXiv:1409.0575; 2014.
9. Bahrapour S, Ramakrishnan N, Schott L, Shah M. Comparative study of deep learning software frameworks. arXiv:1511.06435v3; 2016.
10. Al-Rfou R. et al. Theano: a Python framework for fast computation of mathematical expressions. arXiv:1605.02688v1; 2016.
11. Jia Y, Shelhamer E, Donahue J, Karayev S, Long J, Girshick R, Guadarrama S, Darrell T. Caffe: convolutional architecture for fast feature embedding. arXiv:1408.5093v1; 2014.
12. Collobert R, Kavukcuoglu K, Farabet C. Torch7: A MATLAB-like environment for machine learning. Proc Advances in Neural Info Processing Sys. 2011.

13. Abadi M et al. TensorFlow: large-scale machine learning on heterogeneous distributed systems. arXiv:1603.04467v2; 2016.
14. Yu D, Yao K, Zhang Y. The computational network toolkit. IEEE Signal Processing Magazine. 2015:123–126.
15. Yu D et al. An introduction to computational networks and the computational network toolkit. Redmond (WA): Microsoft; 2014. Report No.: MSR-TR-2014–112.
16. Neon deep learning library (DLL). San Diego (CA): Nervana Systems; 2016 Jul. 15. [accessed 2016 Jul]. <http://neon.nervanasys.com/docs/latest/index.html>.
17. Deeplearning4j: Open-source distributed deep learning for the JVM. San Francisco (CA): Skymind; 2016 [accessed 2016 Jul]. <http://deeplearning4j.org>.
18. Vedaldi A, Fulkerson B. VLFeat: An open and portable library of computer vision algorithms. Proceedings of ACM International Conference on Multimedia; 2010; Firenze, Italy.
19. Szegedy C, Liu W, Yangqing J, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A. Going deeper with convolutions. In: CVPR 2015. Proceedings of IEEE Conference on Computer Vision and Pattern Recognition; 2015; Boston, MA.
20. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. arXiv:1409.1556v6; 2015.
21. Sermanet P, Eigen D, Zhang X, Mathieu M, Fergus R, LeCun T. Overfeat: Integrated recognition, localization and detection using convolutional networks. arXiv:1312.6229v4; 2014.
22. Vedaldi A, Karel L. Matconvnet: convolutional neural networks for MATLAB. Proceedings of 23rd ACM International Conference on Multimedia; 2015; Brisbane, Australia.
23. Ma H. University of Guelph, Ontario. Personal communication, 2016.



1 DEFENSE TECH INFO CTR  
(PDF) DTIC OCA

2 US ARMY RSRCH LAB  
(PDF) IMAL HRA MAIL & RECORDS MGMT  
RDRL CIO L TECHL LIB

1 GOVT PRNTG OFC  
(PDF) A MALHOTRA

9 US ARMY RESEARCH LABORATORY  
(PDF) RDRL CII  
B BROOME  
M THOMAS  
RDRL CII A  
S YOUNG  
A TUNICK  
D BARAN  
G WARNELL  
P DAVID  
RDRL CIN T  
R MEYERS  
K DEACON

INTENTIONALLY LEFT BLANK.