
[Re] Synbols : Probing Learning Algorithms with Synthetic Datasets

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 Scope of Reproducibility

3 This report focuses on the reproduction of some results presented of the above-mentioned paper [3]. Authors introduced
4 a new data generator called Synbols allowing fast generation of low-resolution images rich in latent features. Researchers
5 explored the capabilities of the tool by training popular machine learning algorithms in various M.L paradigms with
6 their synthetically generated datasets. The tool is also trying to address some broader issues relevant to the whole field
7 (i.e. faster iteration cycles for the training, less reliance on expensive hardware, etc.).

8 To assess the features of Synbols and its capacity to explore well known neural network architectures, we decided to
9 reproduce the results of the Supervised Learning classification task and the Unsupervised Representation Learning
10 experiments. We then generated some datasets with the same attributes to assure the results were consistent. Additionally,
11 we tried to get further insights for the unsupervised task by modifying classifier downstream. The final code used for
12 implementing the replicated results can be found here: [Re] Synbols Repository

13 Methodology

14 Regarding our methodology, we predominantly followed authors instructions and their publicly available code. Modifi-
15 cations to the original code made in order to further explore some findings will be discussed later in the corresponding
16 section.

17 Results

18 We manage to reproduce the original results falling within a 2% margin of the reported values. We were pleasantly
19 surprised given the number of models and datasets tested. And thus conclude that Synbols is a well designed tool for
20 rapidly generating a wide variety of low resolution images of UTF-8 characters and strings.

21 What was easy / What was difficult

22 We applaud authors reproducibility efforts and their availability whenever we had questions. A repository specifically
23 made in order to facilitate the reproduction was available and an up-to-date docker image was also at our disposal to
24 help generate more datasets with the tool. No hidden/forgotten assumptions were needed to reproduce their results.
25 Originally, for the two paradigms tested, twelve different models were trained. Although important hyper-parameters
26 and architecture choices were always mentioned or referenced, we sometimes struggled to check their implementation
27 to see if everything was performed as reported.

28 Communication with original authors

29 We actively reached out the original authors through e-mail and meeting sessions. The authors always made time
30 to answer our questions. Hereby, we sincerely thank the authors for providing us adequate supports during the
31 reproducibility.

32 1 Introduction

33 The original paper [3] introduces Synbols, a dataset generator with a rich latent feature space. It generates low resolution
34 images to support quick iteration times. More than 1000 artistic fonts over 14 different languages were collected. The
35 diversity of background and foreground can also vary from solid, gradient, camouflage and natural. Occlusion can also
36 be added to the foreground. In each symbol or character, one can modify the inherent attributes of the image or the
37 character itself. This includes *translation, scale, rotation, shear, bold, and italic*. The authors used this versatile tool to
38 probe the limits of existing algorithms in different machine learning paradigms relevant in the field of computer vision.
39 The motivation behind designing a low-resolution dataset generator is that, usually in order to obtain state-of-the-art
40 performance, the model is expected to train on large-scale dataset, especially when the model complexity is high. But it
41 comes at the cost of slow iteration cycles, taking sometimes weeks of training before obtaining the expected results.
42 On the other hand, applying small-scale datasets to train new SOTA models would limit the capability of testing their
43 generalization capacity but also prevent meaningful model comparison. Still, relying on very large datasets creates
44 a high barrier to entry for many organizations and researchers wanting to get into the Deep Learning Revolution [4].
45 Finally, current research is biased towards fast methods leveraging big datasets instead of considering a more qualitative
46 approach. Synbols aims at solving those issues. Our team is confident that this field of research is of importance for the
47 future and hope that the following reproducibility report will help assess with more confidence the presented claims to
48 allow more research to be conducted on this topic.

49 Our report is articulated around three key questions ;

- 50 • Are the original results reproducible ?
- 51 • Were there any hidden assumptions in order to obtain the same results ?
- 52 • Can we quickly generate similar datasets ?

53 2 Methodology

54 In the original article, authors probed six machine learning paradigms in order to test their synthetically generated
55 datasets. Researchers aim was to further investigate strengths and weaknesses of popular machine learning models by
56 exposing them to a wide range of challenging datasets generated by Synbols. We focused our efforts on replicating the
57 supervised learning and the unsupervised representation learning experiments.

58 In order to facilitate the reproducibility of the experiments and the results presented in the paper, authors made
59 the code used for the benchmarks publicly available. The repository contained the model architectures, the training/testing/validation in HDF5 format storing the images but also the corresponding attributes used in the generation.
60 Each dataset was generated three times using different pseudo-random seed in order to test more thoroughly each
61 dataset. For the more computationally demanding models we ran the experiment using only one seed. We additionally
62 decided to generate the camouflage dataset using the same attributes and seed. The two datasets were identical and
63 provided consistent results. To gain further insights on the unsupervised representation task we edited the source code.
64 More specifically, we modified the classifier downstream on the pipeline by tweaking the original MLP and then trying
65 with a linear regression. We tried implementing a different classifier (EfficientNet) but it did not provide any meaningful
66 insight to understanding the low performance in the unsupervised task.
67

68 3 Reproducibility resources

69 The computational resources required to reproduce the experiment were very accessible. Authors originally used Tesla
70 V100 (TDP of 300W) type hardware for a cumulative 23916 hours of computation needed for the whole paper (this
71 includes debugging, failed experiments and hyperparameter search). By focusing on two experiments and reducing the
72 number of seed tested, we were able to reproduce their results in approximately 194 hours using a Tesla K80 (TDP of
73 300W) type GPU with 12GB of GDDR5 memory available on Google Cloud Platform. Total emissions are estimated to
74 be 1.16 kgCO₂eq. [2]. All models were implemented using Pytorch.

75 3.1 Datasets

76 3.1.1 Supervised Learning

77 The Synbols default dataset will serve as baseline for other dataset and it consists of samples of English characters with
78 a font uniformly selected from the font collection and the attributes are selected to have high variance. Respectively for

79 the Camouflage and Natural datasets, the according feature was added to the default dataset. The Less Variations dataset
 80 removes the italic and bold attributes and reduces the variations of other attributes. Finally, the Korean dataset consists
 81 of a uniformly selected Hangul characters (reduced to the first 1000 symbols). The width and height and channels of all
 82 of the images is 32x32x3 and the dataset size was 100k¹. The authors also decided to confront those synthetic datasets
 83 to popular benchmark datasets, namely MNIST and SVHN. We did not reproduce the results for those standard datasets
 84 instead choosing to focus our efforts on the synthetic datasets generated by the tool.

85 3.1.2 Unsupervised Representation Learning

86 In the Unsupervised Representation Learning, the paper leverages three variants of datasets, namely, solid, camouflage,
 87 and shades. In these datasets the bold attribute was kept on while a low variance was applied on the scale. The first
 88 variant, the solid dataset used black and white contrast while a smooth gradient was applied on the shade variant. In the
 89 camouflage dataset the corresponding attribute was added. The width and height and channels of all of the images is
 90 32x32x3. Moreover, due to limited resources we only used one of the three variant of each dataset ².

91 4 Model Architecture

92 All the models were trained using adaptive learning rate optimization algorithm [1]. Also, the results were obtained
 93 using a partition size of (60%, 20%, 20%) for the training, validation and testing sets and the learning rate was selected
 94 using the validation set. Models were trained using Mixed precision, a NVIDIA extension enabling distributed training
 95 for Pytorch. Tables containing information about the architectures in a more condensed manner can be found in App. B.

96 5 Reproduction Results

97 In this section, we present our reproduction results for the Supervised and Unsupervised experiments. We followed
 98 as closely the ideas presented by the authors but as previously mentioned the default dataset of size 1 million nor the
 99 standard deviation on some results (where we only reproduced one seed) were reported. Because the standard deviation
 100 was relatively small and the default dataset followed the same data distribution, we believe our overall conclusion on
 101 the reproducibility still holds.

102 5.1 Supervised Learning

103 The results of supervised learning experiment were used as the baseline for all the other experiments presented in the
 104 article. For this reason it seemed imperative for us to start by reproducing those results. Here are the results we obtained,
 105 see table 1.

Dataset	Synbols Default	Camouflage	Korean	Less Variation
Size	100k	100k	100k	100k
MLP	14.56 +0.27	3.98 +0.10	0.11 +0.1	0.06 +0.05
Conv-4-Flat	68.47 +0.04	34.62 -2.27	2.07 -0.45	0.22 -0.01
Conv-4-GAP	70.83 -0.69	28.90 +0.70	33.96 -0.38	3.53 -0.37
ResNet-12	95.58 -0.15	90.44 -0.30	96.92 +0.16	38.51 +0.9
ResNet-12+	97.24 -0.08	94.39 -0.04	98.58 -0.04	57.63 -0.21
WRN-28-4	93.74 -0.17	86.64 +0.30	96.47 -0.68	22.18 +0.92
WRN-28-4+	97.38 -0.03	95.54 +0.01	99.27 -0.13	67.02 +1.4

Table 1: **Reproduction of Supervised Learning Results:** Accuracy of various models on supervised classification tasks. Deviation from original results are in gray.

106 We can see that the results are very similar to the results reported in the paper [3] confirming the assessment of the
 107 authors on the versatility of the synthetic data generated. While all models were able to achieve +98% accuracy on
 108 MNIST dataset, only the state-of-the-art models were able to achieve high accuracy on more sophisticated datasets
 109 generated by the tool. We can also assert that Synbols can be used to provide meaningful data augmentation ³, increasing
 110 by a factor of three the accuracy achieved on the hardest dataset (i.e Less Variations). In addition, we trained a second
 111 time the MLP using the same datasets generated on our own and obtained very similar results.

¹Due to limited resources we were not able to run the larger Default variant dataset.

²Originally generated using three different pseudo random seed to replicate the results

³Here, data augmentation consists of uniformly sampled affine deformations in the attributes.

112 **5.2 Unsupervised Representation Learning**

113 The reproduction results are reported in the following table.

	Character Accuracy			Font Accuracy		
	Solid Pattern	Shades	Camouflage	Solid Pattern	Shades	Camouflage
Deep InfoMax	82.69 +1.18	6.15 +0.37	5.48 -0.63	15.37 +1.07	0.23 +0.08	0.25 +0.03
VAE	60.73 +2.75	22.17 +0.26	2.98 +0.87	2.11 +0.57	0.27 +0.09	0.11 +0.07
HVAE	68.92 -2.2	28.32 -0.54	3.79 +0.12	1.9 +0.81	0.29 +0.1	0.16 +0.01

Table 2: **Reproduction of Unsupervised Representation Learning Results:** Accuracy of a MLP classifier downstream. Deviation from original results are in gray.

114 Again, we observe the reproduced results are aligned with the ones reported in the paper. Although all models perform
 115 well in character classification on the solid pattern dataset, we observe the same significant drop on the Shades and
 116 Camouflage variants. Those results are very different from the ones reported in the Supervised experiment. In Sec. 5.2
 117 we mention some of our hypothesis regarding this issue.

118 **6 Discussion of findings**

119 **6.1 Supervised Learning**

120 The table shown in In Sec. 5.1 report the test set loss from our reproduction. However, it is still interesting to mention
 121 how fast different supervised learning models reduce the validation loss to the optimum through iterations of epoch.
 122 This can reflect the ability of models tackling the synthetic datasets. Specifically, except on the Less Variation dataset we
 123 noted that WRN performed really well on the classification task. We believe this model, thanks to its wider convolutional
 124 layers, benefits from the rich composition of latent features generated by Symbols. What is also impressive is the speed
 125 at which it achieves high accuracy and robustness (i.e Generalization). Even on the hardest dataset tested, the optimal
 126 training and validation losses were reached at the 25th epoch as shown in Figure 1. ⁴

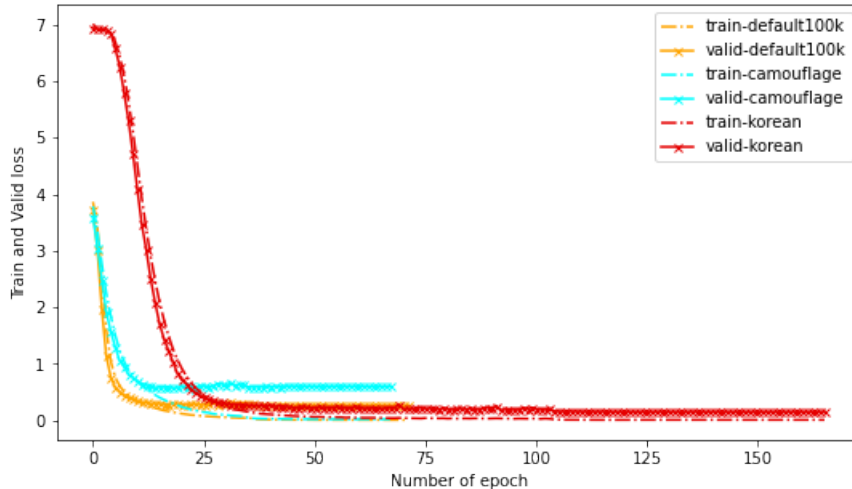


Figure 1: Cross Entropy loss of WRN on various datasets

127 We have noticed a difference in the choice of channels between what was reported in the paper and their code for
 128 Conv-4. From inspecting their code, we found that instead of the 64 channels for all layers claimed in the article, the 4
 129 layers had 32, 64, 128, 256 channels respectively.

⁴The final test loss for each model is reported in Tab. 2.

130 6.2 Unsupervised Representation Learning

131 Although all the models were able to learn meaningful representations on the Solid dataset, a major drop was observed
132 when adding Camouflage. The best performing model, Deep InfoMax for Solid and Camouflage Pattern was the least
133 performing on Shades. It seems that due to the global structure of the gradient pattern Deep InfoMax, the model
134 struggles to capture meaningful latent features in the limited size representation. Intuitively speaking, we believe that
135 the local feature in the gradient pattern can be very different from the global feature of the original image ⁵ and this is
136 why Deep InfoMax did not capture meaningful representations for Shades dataset.
137 We tried to increase the accuracy by performing a grid search on the MLP classifier downstream and also tried with a
138 linear regression model, both methods lead to similar performance (5% margin).

139 7 Conclusion

140 Despite a couple of points that were different in the code from what was reported in the paper, we applaud authors
141 reproducibility efforts and their availability when we had questions. We were able to reproduce the original results
142 without major drawbacks. We thus conclude by answering the three key questions as followed;

- 143 • Are the original results reproducible? *Yes*.
- 144 • Were there any hidden assumptions in order to obtain the same results? *No*.
- 145 • Can we quickly generate similar datasets? *Yes*.

146 Synbols is a very versatile tool for rapidly generating rich composition of latent features in low resolution images
147 effectively probing a wide range of machine learning algorithms. We also observe that it can help identify latent
148 properties and increase the robustness of a model on smaller datasets.

149 Although its limited generation capabilities (i.e : UTF-8 symbols only), authors are planning to add more features to
150 the current generator and also extend the concept to video generation/visual question answering support. We are very
151 excited to see its impact on the computer vision field and hopefully on the whole field of deep learning.

152 8 Discussion

153 This report focuses on the reproduction of some results presented of the above-mentioned paper [3]. Authors introduced
154 a new data generator called Synbols allowing fast generation of low-resolution images rich in latent features. Researchers
155 explored the capabilities of the tool by training popular machine learning algorithms in various M.L paradigms with
156 their synthetically generated datasets. The tool is also trying to address some broader issues relevant to the whole
157 field (i.e.faster iteration cycles for the training, less reliance on expensive hardware, etc.). In this report, we follow the
158 replication instructions and the published code provided by the authors in order to verify some of those claims. The
159 final code used for implementing the replicated results can be found here: [Re] Synbols Repository.

160 8.1 What was easy

161 We applaud authors reproducibility efforts and their availability whenever we had questions. A repository specifically
162 made in order to facilitate the reproduction was available and an up-to-date docker image was also at our disposal to
163 help generate more datasets with the tool. No hidden/forgotten assumptions were needed to reproduce their results.
164 Thanks to those all those efforts our task was significantly simplified.

165 8.2 What was difficult

166 Originally, for the two paradigms tested, twelve different models were trained. Although the important hyper-parameters
167 were always mentioned or referenced, we sometimes struggled to check their implementation to see if everything was
168 performed as reported. But authors always made time to explain implementation details that were more difficult to
169 understand at first glance.

⁵The model is more likely to confuse gradient changes with important symbol information.

170 **8.3 Communication with original authors**

171 We actively reached out the original authors through e-mail and meeting sessions. The authors always made time
 172 to answer our questions. Hereby, we sincerely thank the authors for providing us adequate supports during the
 173 reproducibility.

174 **References**

175 [1] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

176 [2] A. Lacoste, A. Luccioni, V. Schmidt, and T. Dandres. Quantifying the carbon emissions of machine learning. *arXiv*
 177 *preprint arXiv:1910.09700*, 2019.

178 [3] A. Lacoste, P. Rodríguez López, F. Branchaud-Charron, P. Atighehchian, M. Caccia, I. H. Laradji, A. Drouin,
 179 M. Craddock, L. Charlin, and D. Vázquez. Symbols: Probing learning algorithms with synthetic datasets. *Advances*
 180 *in Neural Information Processing Systems*, 33, 2020.

181 [4] T. J. Sejnowski. *The Deep Learning Revolution*. The MIT Press, 10 2018.

182 **Appendix**

183 **A. Supervised Learning**

184

MLP Parameters	Value
Layers	3
Hidden size	256
Activation	Leaky ReLU non-linearities
Learned parameters	72k (fully connected)

185

Conv-4-GAP Parameters	Value
Convolution layers	4
Channels per layer	64
Pooling	Global average
Learned parameters	112k

186

Resnet-12 Parameters	Value
Residual Layers	12
Residual blocks	4
Channel/Output per block	{64,128,256,512}
CNN per block	3
CNN structure	3x3
Activation	ReLU non-linearities
Pooling (at the end of each block)	Max
Dropout(first& second convolution at each block)	0.1
Learned parameters	8M

187

WRN-28-4 Parameters	Value
Residual Layers	28
Residual blocks	{16,4,4,4}
Output per block	{16,32,64,128}*4
CNN structure	3x3
Activation	ReLU non-linearities
Pooling	Global average
Dropout	0.1
Batch size	128
Learned parameters	5.8M

188 **B. Unsupervised Supervised Learning**

Deep InfoMax hyperparameters		Value
Seed		2
Dropout		0.3
Activation Function		ReLU
Kernel		3
Stride		1
Padding		1
Feature Vector Size		64
Global Discriminator Number of Convolutional Layers		2
alpha		0.5
Local Discriminator Number of Convolutional Layers		3
Beta		1.0
Prior Discriminator Number of Fully-Connected Layers		3
Gamma		0.1

189

Variational Auto-Encoder hyperparameters		Value
Dropout		0.3
Activation Function		leaky ReLU
Kernel		3
Stride		1
Padding		1
Pooling		2x2
Beta		0.01
Feature Vector Size		64
Hierarchical = True for HVAE		False

190