Towards Biologically Plausible Learning by Stacking Circular Autoencoders

Anonymous Full Paper Submission 14

OD1 Abstract

Training deep neural networks in biological systems 002 is faced with major challenges such as scarce labeled 003 data and obstacles for propagating error signals in 004 005 the absence of symmetric connections. We introduce 006 Tourbillon, a new architecture that uses circular au-007 to encoders trained with various recirculation algorithms in a self-supervised mode, with an optional 008 top layer for classification or regression. Tourbillon 009 010 is designed to address biological learning constraints rather than enhance existing engineering applica-011 tions. Preliminary experiments on small benchmark 012 datasets (MNIST, Fashion MNIST, CIFAR-10) show 013 that Tourbillon performs comparably to models 014 trained with backpropagation and may outperform 015 other biologically plausible approaches. The code 016 and models are available at https://anonymous. 017 4open.science/r/Circular-Learning-4E1F. 018

1 Introduction

Decades of machine learning have taught us that 020 gradient descent is the sole effective optimization 021 method in high-dimensional spaces. Other strategies, 022 like random search, are bound to fail. Backpropa-023 024 gation, the algorithm behind gradient computation in artificial neural networks, has been incredibly 025 successful. It powers advancements in Artificial 026 Intelligence, from protein folding (e.g., AlphaFold 027 [1]) to natural language understanding and gener-028 ation (e.g., GPT-4 [2, 3]). Backpropagation effi-029 ciently computes the gradient in a network with W030 weights using O(W) operations. Considering that 031 at least O(W) operations are necessary to adjust 032 W synapses, backpropagation demonstrates optimal 033 efficiency. Consequently, if learning is viewed as an 034 optimization problem in a high-dimensional space of 035 synaptic weights, this suggests that the brain likely 036 employs learning algorithms based on gradient com-037 putation, either exact or approximate. Yet there 038 are several well-known reasons in the literature why 039 backpropagation is implausible in biological systems 040 [4, 5, 6, 7, 8, 9]. Thus, in short, we hypothesize that 041 biological systems must strive to approximate gradi-042 ent descent methods without being able to compute 043 exact gradients by backpropagation. Here, we set 044 out to propose a plausible strategy for achieving this 045 goal. 046

047 Let us first briefly enumerate some of the ma-

jor reasons why backpropagation is not plausible 048 in biological neural networks: 1) Symmetry of 049 **Connections** (weight transport): Backpropagation 050 requires precisely symmetric connections between 051 the forward and backward passes. This constraint 052 cannot be satisfied in a biological neural system and 053 might be hard to realize in some physical neural 054 systems. 2) Forward Nonlinearities (F prime): 055 Backpropagation relies on the exact memory of for-056 ward pass nonlinearities (e.g., activation functions) 057 to compute weight updates. This is not supported 058 in biological or physical neural systems. 3) Local-059 ity: In a biological neural system, the learning rule 060 for adjusting synaptic weights must be local, i.e. it 061 must rely solely on variables available locally, both in 062 space (spatial locality) and time (temporal locality), 063 at each synapse. 4) Clocked Computation: In 064 backpropagation, the forward and backward passes 065 are manually clocked to compute activations and 066 update weights. In contrast, in a biological sys-067 tem, neurons communicate stochastically, lacking 068 the precise clocking mechanism observed in back-069 propagation. 5) Labeling: Backpropagation relies 070 on large labeled datasets, unlike biological systems, 071 which lack access to such data. 6) Spike: Bio-072 logical neurons use noisy spikes for communication, 073 while artificial neurons typically use deterministic 074 analog values. 7) **Distances**: Backpropagation 075 necessitates propagating signals over considerable 076 neural distances in deep models, which can result 077 in signal dilution and lead to distorted or unsta-078 ble gradients. 8) Developmental Modularity: 079 Backpropagation in general, requires having a com-080 plete architecture in place before training can begin, 081 which may not be realistic for biological systems 082 undergoing development and other changes. 083

Several solutions have been suggested to try to 084 address these problems, in isolation or small combi-085 nations, but no approach addresses all of them at 086 once. Here we propose a neural architecture called 087 Tourbillon and its training algorithms to address 088 all the implausibility discussed above by combining 089 different ideas, including stacked autoencoders, recir-090 culation, and asynchronous training. We emphasize 091 that the primary goal here is to address the obsta-092 cles listed above for biological (or neuromorphic) 093 neural systems and not to derive a new architecture 094 or algorithm that is practically useful for digital 095 applications of deep learning. 096

⁰⁹⁷ 2 Biological Plausibility

Several approaches have been proposed to address 098 the biological implausibilities enumerated above. 099 The most notable ones include: Feedback Alignment 100 (FA), Difference Target Propagation (DTP), Stacked 101 Autoencoders, and the Forward-Forward (FF) algo-102 rithm. However, each of these methods addresses 103 only a limited subset of the biological implausibilities 104 (Section A.1 and Table 1). Self-supervised learning, 105 in particular stacked autoencoders, provides one way 106 of addressing the data labeling issue. However, stan-107 dard autoencoders suffer from several other issues 108 which we now address. 109

Circular Autoencoders. In a standard feed-110 forward autoencoder (AE), the data itself provides 111 the targets (self-supervised learning). The data 112 and hence the targets are available in the input 113 layer. However, they are not available in the output 114 layer, in the sense that they are not physically local 115 (spatial-locality) to the output layer. This problem is 116 addressed in circular autoencoders (CAE) [8] where 117 the output layer is physically equal (or physically 118 adjacent) to the input layer (Figure 1). With the 119 circular layout, targets and errors can be computed 120 at the level of the input/output layer. 121

Recirculation Algorithms. Standard back-122 propagation, or even FA, of these targets, would 123 require a channel (wires) running backward from 124 the output layer to the hidden layer. However, be-125 cause of the circular layout, it is possible to use the 126 forward connections to propagate target and error 127 information during learning. This is the fundamen-128 tal idea behind recirculation, a family of algorithms 129 for training CAEs that do not require backward 130 connections [10, 11, 8]. 131

Consider a CAE with layers numbered from 0 to 132 L, where 0 corresponds to the input layer. We use 133 the index t to denote different cyclic passes through 134 the autoencoder, with the first pass indexed by t =135 0. After the first pass, one can *locally* compute 136 the error $T - H_L^0$, where T is the target located 137 at the input layer. This error could be used to 138 train the top layer of the CAE by gradient descent, 139 and then train the other layers by using a form of 140 random backpropagation where the error signal is 141 obtained by propagating the error $T - H_L^0$ using the 142 forward weights of the CAE. This however requires 143 propagating two different kinds of signals, activities, 144 and errors, through the CAE. Thus rather than 145 recirculating the error, a more uniform approach can 146 be obtained by recirculating activities. If H_i^t denotes 147 the activation of layer i during the forward pass 148 indexed by t, the main idea behind the recirculation 149 family of algorithms is to use H_i^t as the target for 150 the output $H_i^{t'}$ taken at a later time t' to produce 151 the post-synaptic term for the weight update. The 152 intuition is that the data may become increasingly 153

corrupted as it is being recycled, thus earlier pass serve as targets for later passes. Different variations can be obtained, by varying, for instance, the postand pre-synaptic terms. In addition to the original learning rule of recirculation [8] Equation 1, we propose two variations of the learning rule, all shown in Equation 9 (b) and (c) in the Appendix.

$$\Delta W_i = \eta (H_i^t - H_i^{t'})^{post} (H_{i-1}^t)^{pre}$$
(1) 16:

These rules follow a Hebbian-product form, resem-162 bling backpropagation but with a postsynaptic re-163 circulation error, denoted as $[H_i^0 - H_i^1]^{post}$. This 164 error term is both spatially and temporally local, 165 assuming that consecutive passes through the circu-166 lar autoencoder fall within the proper time window. 167 In the input layer, the vector H_0^0 represents the in-168 put data, including the targets for an autoencoder. 169 The presynaptic term can be computed at different 170 times (t or t') or even as the difference between the 171 activities at two different times. Using the first form 172 or presynaptic terms (Equation 1), the recircula-173 tion learning equation for the top layer of weights 174 is identical to backpropagation. Although in this 175 work we are not using spiking neurons, all learning 176 rules described in Equation 9 are closely related to 177 the concept of spike time-dependent synaptic plas-178 ticity (STDP) [12]. STDP Hebbian or anti-Hebbian 179 learning rules have been proposed using the tem-180 poral derivative of the activity of the post-synaptic 181 neuron [13] to encode error derivatives. 182

3 Tourbillon: A CAE Stack 183

We propose the Tourbillon architecture as a stack of 184 circular autoencoders, capped by a classification or 185 regression layer connecting the hidden representa-186 tion of the top circular autoencoder and the output 187 layer. Each circular autoencoder has an encoder 188 and decoder components. The hidden layer that is 189 shared by the encoding and decoding components is 190 called the hinge layer. In the stack, the hinge layer 191 of the *i*th circular autoencoder becomes the input 192 layer of the i + 1th circular autoencoder (Figure 2). 193 The Tourbillon architecture addresses the issues of 194 target labels and spatial locality. With the recircula-195 tion algorithms, it also addresses the issues of weight 196 transport, forward non-linearities, temporal locality, 197 and distances. Using a novel training algorithm, we 198 set out to address issues of clocking and modularity. 199

Sequential Training. In sequential training, the 200 CAEs are trained separately. Training of the i-th 201 CAE in the stack must be completed before training 202 of the i + 1-th CAE can begin. The input data to 203 the i + 1-th CAE is provided by the hidden repre-204 sentations produced by the hinge layer of the i-th 205 CAE. Finally, we can stack N trained CAE, each 206 is trained to further compress the hinge layer of its 207



Figure 1. From left to right: Recirculation, forward forward, difference target propagation, (direct) feedback alignment. The learning rule for each model is written at the top of the architecture schematics.

Table 1. A comparison of physical plausibility between different neural architectures from a biological standpoint. $\mathbf{X}, \mathbf{\Theta}$, and $\mathbf{\checkmark}$ correspond to no plausibility, partial plausibility, and full plausibility, respectively.

	W Transport	F prime	Locality	Clocked	Labeling	Spike	Distance	Modular
Backpropagation	×	×	×	×	×	×	×	×
Feedback Alignment (FA)	✓	✓	×	×	×	×	×	×
Direct Feedback Alignment (DFA)	✓	✓	۲	×	×	×	✓	×
Difference Target Propagation (DTP)	۲	×	×	×	۲	×	✓	×
Stacked Autoencoders	✓	✓	۲	×	✓	×	✓	۲
Forward Forward (FF)	✓	✓	~	×	۲	×	✓	×
Tourbillon	~	✓	✓	✓	۲	×	~	✓

(2)

(3)

predecessor CAE. These trained CAEs as the build-208 ing blocks of tourbillon must be trained from the 209 bottom to the top. So each autoencoder would be 210 able to generate the training data for its successor 211 building blocks. This is explained in Equation 3 212 where E and D represent the encoder and decoder 213 characterized by a deep neural network trained with 214 recirculation. x is the original training data (e.g. 215 MNIST) and H shows the hidden representation of 216 the circular autoencoder. 217

218
$$H_1 = E_1(x_1), \hat{x}_1 = D_1(H_1),$$

219

$$H_2 = E_2(H_1), \hat{H}_1 = D_2(H_2)$$

Asynchronous Training. The sequential train-220 ing algorithm is the standard method for training 221 a stack of autoencoders by fully training the first 222 AE, then the second one, and so on. This requires 223 a high degree of orchestration. Moreover, sequen-224 tial training does not fully support modularity, as 225 training each CAE depends on the prior training of 226 all preceding CAEs. To remove the need for clocked 227 orchestration and increase biological plausibility via 228 modularity, we introduce the asynchronous training 229 algorithm where all the CAEs are trained simulta-230 neously and asynchronously. The main idea behind 231 asynchronous training is to train the weights of one 232 CAE of the stack that is randomly chosen. 233

In this case, each CAE can be viewed as a "spinning wheel" and these wheels can spin independently of each other. At any random time, a CAE may elect to recirculate whatever happens to be in its 237 input layer (either random data, older data, or the 238 data at the hinge layer activated by previous CAEs) 239 and adapt its synapses accordingly. Further details 240 and a precise algorithm on asynchronous training 241 are given in the Appendix. 242



Figure 2. Tourbillon architecture with a stack of three circular autoencoders (CAE) trained by recirculation.

4 Experiments and Results 243

We begin by training CAEs to investigate the effects of various parameters, including the number 245 of cycles (between t and t'), and the CAE size (i.e., 246

305

the number of hidden layers in the CAE except the 247 input and output layers). We use the learning rule in 248 Equation 1 in training CAEs, however a discussion 249 on the effect of different learning rules (Equation 9) 250 is presented in Section A.2.1. Then, using the best 251 set of these parameters, we develop and train several 252 Tourbillon architectures, both with and without a 253 top classifier layer, using different stack depths and 254 training algorithms. The goal of these experiments 255 is not to outperform existing deep learning mod-256 els but to show that Tourbillon architectures can 257 learn complex tasks while satisfying the plausibility 258 constraints. Similar to recently proposed plausible 259 architectures [14], we use relatively small datasets 260 and models, leaving the scaling up to future stud-261 ies. Details on the hyperparameters, hardware, and 262 CAE implementation can be found in the Appendix 263 and the GitHub repository. 264

NLDL

#14

265 4.1 Training Tourbillon CAEs

We train CAEs using the learning rule in Equation 266 1. We optimize each architecture using a mean-267 squared reconstruction loss. In all experiments, we 268 use symmetric CAEs, where the number of hidden 269 layers in the encoder and decoder are equal. To 270 satisfy distance plausibility, we use CAEs with a 271 small number of hidden layers. Additionally, to 272 maintain the temporal locality of the variables, we 273 limit the number of cycles (difference between t and 274 t') to one, two, and three. 275

We train CAEs with fully connected layers for the 276 MNIST and Fashion MNIST datasets and convolu-277 tional layers for the CIFAR-10 dataset. Table 2 dis-278 plays the reconstruction losses on the test datasets. 279 Notably, using a CAE size of one and one cycle 280 281 (t = 0 and t' = 1) yields the lowest testing loss. This also corresponds to the highest level of spatial and 282 temporal locality of variables. 283

Using the best values for the CAE size and number 284 of cycles, we further show the viability of training 285 CAEs with the learning rule above. We compare 286 the mean-squared loss of the trained CAE with 287 the same autoencoder trained with backpropagation 288 (BP) and feedback alignment (FA). Figure 3 shows 289 the training and test error curves for the MNIST 290 and Fashion MNIST datasets. Figure 5 (first row) 291 shows the trajectories of training and test error 292 for CAEs with convolutional layers and similar au-293 to encoders trained with BP and FA. Additionally, 294 randomly sampled MNIST images are reconstructed 295 using three models shown in Figure A.3 with the 296 average reconstruction loss. The results presented 297 in these figures show that recirculation outperforms 298 FA across different datasets and architectures both 299 in training and reconstruction loss. Recirculation 300 also shows a reconstruction loss very close to that of 301 BP in both fully-connected and convolutional archi-302 tectures (0.001% relative error). Detailed training 303

0.05 Train Test 0.040 Test 0.040

parameters for each CAE are given in the Appendix. 304



Figure 3. Train and test loss of three equivalent autoencoders. The first row shows the performance of the models for the MNIST dataset, and the second row shows the performance for the Fashion MNIST dataset.

4.2 Stacking CAEs with Various 306 Stack Depths and Training Algo- 307 rithms 308

We construct stacks of two, three, and four compres-309 sive CAEs and train them using sequential and asyn-310 chronous training algorithms. Table 3 demonstrates 311 the reconstruction errors of the stacks, indicating 312 the effect of stack depth and training algorithm ap-313 plied to each CAE. Detailed parameters, such as 314 CAE dimensions, batch size, and learning rate, are 315 given in the Appendix. 316

During training, we found the asynchronous train-317 ing algorithm to be sensitive to learning rate changes. 318 Experimenting with different schedules for each CAE 319 in the stack reveals that decreasing the learning rates 320 from the bottom to the top layers is crucial. Despite 321 the algorithm's inherent randomness, it achieves an 322 identical test reconstruction error to that of Tourbil-323 lon trained sequentially. Thus, asynchronous train-324 ing increases the model's biological plausibility while 325 maintaining its performance. 326

After training the stacks, we add a top classifier 327 layer for MNIST, Fashion MNIST, and CIFAR-10 328 datasets. We conduct classification experiments to 329 compare the performance of different Tourbillons 330 with neural networks of similar, but sequential ar-331 chitecture trained using BP, FA, and a stack of 332 autoencoders (SAEs) with identical parameters to 333 CAEs but trained with backpropagation. 334

Three different Tourbillons are considered in this sate experiment: (1) Tourbillon with the sequential learning algorithm (Tourbillon-seq); (2) Tourbillon with sate the asynchronous learning algorithm (Tourbillonasynch); (3) and Tourbillon with the asynchronous sate learning algorithm that only uses 10% of the labeled sate

NLDL

#14

Table 2. The test mean-squared reconstruction loss of CAEs with different cycles and CAE sizes trained on MNIST, Fashion MNIST, and CIFAR-10. Each number is the mean of five distinct runs. The top results are in boldface.

		MNIST		Fas	hion MNI	ST		CIFAR-10	
CAE size		Cycles			Cycles			Cycles	
	1	2	3	1	2	3	1	2	3
1	0.0090	0.0093	0.0099	0.0123	0.0124	0.0132	0.0013	0.0012	0.0014
3	0.0151	0.0154	0.0165	0.0204	0.0198	0.0204	0.0324	0.0323	0.0381

Table 3. The test mean-squared reconstruction loss of Tourbillon with different stack depths and training algorithms on MNIST, Fashion MNIST, and CIFAR-10. The top results are in boldface.

	MNIST			Fashion MNIST Stack Depth			CIFAR-10	
Training Algorithm	Stack Depth		Stack Depth					
	2	3	4	2	3	4	2	3
Sequential	0.009	0.026	0.027	0.013	0.031	0.032	0.023	0.046
Asynchronous	0.009	0.019	0.025	0.014	0.029	0.042	0.034	0.074

data to train the top classifier (Tourbillon-10%). All
three models achieve a similar level of good classification accuracy. The Tourbillon-10% model is the
most biologically plausible and performs on par with
SAE which is less biologically plausible.

Although Tourbillon's performance falls short of 346 BP in fully connected architectures (Figure 4), its 347 viability is demonstrated by its trainability and 348 competitive accuracy, achieving 92% test accuracy 349 compared to 99% for BP. Additionally, Tourbillon 350 matches the performance of SAE, which is trained 351 with BP but is less biologically plausible. When 352 using convolutional layers, Tourbillon performs com-353 parably to BP and surpasses other less plausible 354 models like FA and SAEs (Figure 5). 355

Additionally, we assess Tourbillon's reconstruction 356 capability. Initially, we employ the trained stack of 357 circular autoencoders to generate compressed rep-358 resentations of the input data. These compressed 359 representations are then mapped onto a 2D space us-360 ing the t-SNE method [15]. Figure 7 shows these 2D 361 maps, highlighting Tourbillon's unsupervised abil-362 ity to cluster elements in each class. Moreover, we 363 utilize the trained stack of circular autoencoders 364 with feed-forward fully connected layers to recon-365 struct the original input images by leveraging the 366 decoder weights of the circular autoencoders. Figure 367 7 demonstrates that the Tourbillon models success-368 fully capture the crucial information of the input 369 images. During testing, the uncorrupted input is 370 passed through the encoder channel to the top of the 371 model and then back through the separate decoder 372 channel to produce the reconstructed sample. In 373 this scenario, each CAE except the top one decodes 374 a representation in their hinge layer that comes from 375 the CAE immediately above, and not from the CAE 376 itself. This may explain why the reconstruction er-377 ror seems to increase with the depth of the stack 378 (Table 3). However, this does not imply that stacks 379

with multiple CAEs are ineffective as we are interested in using the hidden representation at the top of the stack for classification or regression and not for reconstruction.



Figure 4. Train and test accuracy of six classifiers trained on the MNIST dataset BP, FA, Tourbillon Sequential (seq), Tourbillon Asynchronous (asynch), Tourbillon with 10% of the labeled data, and SAE. Each line corresponds to the mean of five distinct runs with the standard deviation shown as the shaded area.

4.3 Tourbillon Requires Less Labeled 384 Data 385

A key advantage of Tourbillon is its ability to lever-386 age unlabeled data for unsupervised training of the 387 stack, allowing the top classifier to be trained with 388 significantly less labeled data. To evaluate this, we 389 use a stack of CAEs, identical to the one in the 390 previous section, with the asynchronous training 391 algorithm. Once the stack is trained in an unsu-392 pervised manner, we train the top classifier in a 393 supervised manner using 10%, 25%, 50%, and 100%394 of the labeled data. The results for both training 395 and testing are presented in Figure 6. 396

Remarkably, Tourbillon's performance remains 397 unchanged even with only 10% of the data. This is 398 because the stack of CAEs produces an abstract and 399 informative representation of the input, enabling the 400 small top-layer classifier to learn the classification 401 task with a minimal amount of labeled data. 402

430



Figure 5. Top: The mean-squared reconstruction loss for three autoencoders trained with BP, FA, and CAE on CIFAR-10. Bottom: Train/test accuracy for six classifiers trained with BP, FA, Tourbillon-seq, Tourbillonasynch, SAE, and Tourbillon-10% on CIFAR-10. Each trajectory is the mean of five runs, with the shaded area representing the standard deviation.

403 4.4 Conversion to Tourbillon

Finally, we introduce an algorithm designed to convert a traditional neural network with a sequential
layout into its Tourbillon version.

By applying this algorithm to a layered feed-407 forward architecture with no skip connections, we 408 can create a more biologically plausible counterpart 409 410 that preserves the original model's functionality and core structure. The Tourbillon version has roughly 411 twice as many parameters but requires only minimal 412 additional memory. Importantly, during inference, 413 414 both the computational resources and parameter count of the Tourbillon version are identical to those 415 of the original model. 416



Figure 6. The train and test accuracy of four asynchronous Tourbillons with the use of 10%, 25%, 50%, and 100% of the labeled data. Each line corresponds to the mean of five distinct runs with the standard deviation shown as the shaded area.

The key idea is to substitute each sequential layer 417 with a small CAE, enabling the encoding of inter-418 mediate data in a manner that facilitates its trans-419 mission to the higher layers of the network. This 420 conversion algorithm is represented in Algorithm 421 A.2. We use Algorithm A.2 to convert a U-Net [16] 422 and a feed-forward architecture into their Tourbillon 423 424 versions. The details of the training are given in the



Figure 7. tSNE plots for the 2D visualization of 1000 random samples from the MNIST, Fashion MNIST, and CIFAR-10 datasets. Examples of reconstructed images from the MNIST and the Fashion MNIST datasets using Tourbillon are shown on the right.

Appendix. Table A.3 illustrates that the process of 425 converting an architecture to its more biologically 426 plausible Tourbillon version incurs a relatively small 427 cost in terms of accuracy degradation (94% vs 99%). 428

Table 4. Comparison of networks trained with BP vs their Tourbillon version: U-Net numbers show meansquared error, feed-forward numbers show classification error.

U-Net				
	MNIST	CIFAR-10		
BP Train BP Test Tourbillon Train Tourbillon Test	$\begin{array}{c} 0.0031 \; (5.8e\text{-}05) \\ 0.0031 \; (5.8e\text{-}05) \\ 0.0111 \; (e\text{-}04) \\ 0.0113 \; (2e\text{-}04) \end{array}$	$\begin{array}{c} 0.0024 \ (5.8e{-}05) \\ 0.0026 \ (e{-}04) \\ 0.0733 \ (3.2e{-}04) \\ 0.0783 \ (3.2e{-}04) \end{array}$		
Fully Connected				
	MNIST	CIFAR-10		
BP Train BP Test Tourbillon Train	$\begin{array}{c} 1.04 \ (0.20) \\ 1.18 \ (0.20) \\ 5.21 \ (0.40) \end{array}$	-		
Tourbillon Test	5.21(0.40) 5.82(0.40)	-		

5 Conclusion

Tourbillon represents a systematic approach towards 431 addressing the problems of deep learning in biolog-432 ical or neuromorphic networks. In essence, it is a 433 stack of circular autoencoders trained by recircu-434 lation, hence the name Tourbillon associated with 435 turbulence in French. Moreover, in horology, a tour-436 billon is an addition to the mechanics of a watch 437 escapement to increase its accuracy. While we do not 438 claim to have increased accuracy, we have shown that 439 the Tourbillon approach shows similar performance 440 to backpropagation, at least on MNIST, Fashion 441 MNIST, and CIFAR-10. 442

Several issues have been identified that will require 443 additional research. The first one is to study the 444 scaling of the Tourbillon architecture so it can be 445 trained on real-world data sets such as ImageNet [17]. 446 The second one is to study whether local learning 447 algorithms are necessary to also fine-tune the stack 448 during regression or classification. For instance, 449 using decoder weights as the random matrices of FA 450 may be a possibility. The third one is the issue of 451 convolutions that was only incrementally addressed 452 here by showing that Tourbillon with convolutions 453 does better than FA, but lags behind BP. And last, 454 there is the study of Tourbillon with spiking neurons. 455

NLDL

#14

456 References

- [1] John Jumper, Richard Evans, Alexander
 Pritzel, Tim Green, Michael Figurnov, Olaf
 Ronneberger, Kathryn Tunyasuvunakool, Russ
 Bates, Augustin Žídek, Anna Potapenko, et al.
 Highly accurate protein structure prediction
 with alphafold. *Nature*, 596(7873):583–589,
 2021.
- 464 [2] Alec Radford, Karthik Narasimhan, Tim Sali465 mans, Ilya Sutskever, et al. Improving language
 466 understanding by generative pre-training. 2018.
- 467 [3] Alec Radford, Jeffrey Wu, Rewon Child, David
 468 Luan, Dario Amodei, Ilya Sutskever, et al. Lan469 guage models are unsupervised multitask learn470 ers. OpenAI blog, 1(8):9, 2019.
- [4] Tobias Bonhoeffer, Volker Staiger, and AMHJ
 Aertsen. Synaptic plasticity in rat hippocampal
 slice cultures: local" hebbian" conjunction of
 pre-and postsynaptic stimulation leads to distributed synaptic enhancement. *Proceedings of*the National Academy of Sciences, 86(20):8113–
 8117, 1989.

478 [5] Yoshua Bengio, Dong-Hyun Lee, Jorg Bornschein, Thomas Mesnard, and Zhouhan Lin.
480 Towards biologically plausible deep learning.
481 arXiv preprint arXiv:1502.04156, 2015.

- [6] Pierre Baldi and Peter Sadowski. A theory of
 local learning, the learning channel, and the optimality of backpropagation. *Neural Networks*,
 83:61-74, 2016.
- [7] James CR Whittington and Rafal Bogacz. An approximation of the error backpropagation algorithm in a predictive coding network with
 local hebbian synaptic plasticity. *Neural com- putation*, 29(5):1229–1262, 2017.
- [8] P. Baldi and P. Sadowski. Learning in the
 machine: Recirculation is random backpropagation. Neural Networks, 108:479–494, 2018.
- [9] Timothy P Lillicrap, Adam Santoro, Luke Marris, Colin J Akerman, and Geoffrey Hinton.
 Backpropagation and the brain. Nature Reviews Neuroscience, 21(6):335–346, 2020.
- [10] Geoffrey E Hinton and James L McClelland.
 Learning representations by recirculation. In *Neural information processing systems*, pages
 358–366. New York: American Institute of
 Physics, 1988.
- [11] Randall C O'Reilly. Biologically plausible errordriven learning using local activation differences: The generalized recirculation algorithm. *Neural computation*, 8(5):895–938, 1996.

- [12] Henry Markram, Joachim Lübke, Michael 507
 Frotscher, and Bert Sakmann. Regulation of 508
 synaptic efficacy by coincidence of postsynaptic aps and epsps. *Science*, 275(5297):213–215, 510
 1997. 511
- [13] Xiaohui Xie and H. Sebastian Seung. Spike-512
 based learning rules and stabilization of per-513
 sistent neural activity. In S. A. Solla, T. K. 514
 Leen, and K. Müller, editors, Advances in Neu-515
 ral Information Processing Systems 12, pages 516
 199–208. MIT Press, 2000. 517
- [14] Geoffrey Hinton. The forward-forward algorithm: Some preliminary investigations. arXiv 519 preprint arXiv:2212.13345, 2022. 520
- [15] Laurens Van der Maaten and Geoffrey Hinton. 521
 Visualizing data using t-sne. Journal of machine learning research, 9(11), 2008. 523
- [16] Olaf Ronneberger, Philipp Fischer, and Thomas 524
 Brox. U-net: Convolutional networks for 525
 biomedical image segmentation. In Interna- 526
 tional Conference on Medical image computing 527
 and computer-assisted intervention, pages 234- 528
 241. Springer, 2015. 529
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, 530
 Kai Li, and Li Fei-Fei. Imagenet: A large-scale 531
 hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009. 534
- [18] Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. Random
 synaptic feedback weights support error backpropagation for deep learning. Nature communications, 7(1):1–10, 2016.
- [19] P. Baldi, Z. Lu, and P. Sadowski. Learning in 540 the machine: Random backpropagation and the 541 deep learning channel. Artificial Intelligence, 542 260:1–35, 2018. Also: arXiv:1612.02734. 543
- [20] Arild Nokland. Direct feedback alignment provides learning in deep neural networks. arXiv 545 preprint arXiv:1609.01596, 2016. 546
- Maria Refinetti, Stéphane d'Ascoli, Ruben 547
 Ohana, and Sebastian Goldt. The dynamics 548
 of learning with feedback alignment. arXiv 549
 preprint arXiv:2011.12428, 2020. 550
- [22] P. Baldi, Z. Lu, and P. Sadowski. Learning 551
 in the machine: the symmetries of the deep 552
 learning channel. Neural Networks, 95:110–133, 553
 2017. 554
- [23] Roman Pogodin, Yash Mehta, Timothy Lilli crap, and Peter E Latham. Towards biolog ically plausible convolutional networks. Ad vances in Neural Information Processing Sys tems, 34:13924–13936, 2021.

598

605

621

[24] Qianli Liao, Joel Leibo, and Tomaso Poggio.
How important is weight symmetry in backpropagation? In *Proceedings of the AAAI Con- ference on Artificial Intelligence*, volume 30,
2016.

- [25] Sergey Bartunov, Adam Santoro, Blake A
 Richards, Luke Marris, Geoffrey E Hinton,
 and Timothy Lillicrap. Assessing the scalability of biologically-motivated deep learning
 algorithms and architectures. arXiv preprint
 arXiv:1807.04587, 2018.
- 571 [26] Theodore H Moskovitz, Ashok Litwin-Kumar,
 572 and LF Abbott. Feedback alignment in
 573 deep convolutional networks. arXiv preprint
 574 arXiv:1812.06488, 2018.
- J. Ott, E. Linstead, N. LaHaye, and P. Baldi.
 Learning in the machine: To share or not
 to share. Neural Networks, 2020. In
 press. Available online March 25, 2020.
 https://doi.org/10.1016/j.neunet.2020.03.016.
- [28] G.E. Hinton, S. Osindero, and Y.W. Teh. A fast
 learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.
- [29] Pierre Baldi. Autoencoders, unsupervised learning, and deep architectures. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pages 37–49. JMLR Workshop and Conference Proceedings, 2012.
- [30] Michael Tschannen, Olivier Bachem, and
 Mario Lucic. Recent advances in autoencoderbased representation learning. arXiv preprint
 arXiv:1812.05069, 2018.
- [31] Ila R Fiete, Walter Senn, Claude ZH Wang, and
 Richard HR Hahnloser. Spike-time-dependent
 plasticity and heterosynaptic competition organize networks to produce long scale-free sequences of neural activity. *Neuron*, 65(4):563–
 576, 2010.

A Appendix

In this appendix, we provide additional details regarding the algorithms and the experiments. Each section in this appendix corresponds to the section with the same title in the main article. All the experiments are conducted using a single NVidia Titan X GPU. 604

A.1 Biological Plausibility

We describe the learning equation (weight update) GOGusing post- and pre-synaptic terms. For a forward GOGweight W_i at layer *i*, the backpropagation learning GOGequation can be written as: GOG

$$\Delta W_i = \eta B_i^{\text{post}} H_{i-1}^{\text{pre}}, \quad B_L^{\text{post}} = T - H_L, \quad (4) \quad \text{610}$$

$$B_{i-1}^{\text{post}} = F_{i-1}' \circ W_i^T B_i^{\text{post}}$$
⁶¹¹

where η denotes the learning rate, B_i^{post} denotes 612 the postsynaptic backpropagated error at layer i, 613 and H_{i-1}^{pre} denotes the pre-synaptic activity. $F'_{i-1} \circ$ 614 denotes the component wise multiplication by the 615 vector of activation function derivatives in layer i-1. 616 T and H_L are the targets and the activation at the 617 top layer L. Note that in order to avoid further clut-618 tering the notation, we omit the transpose sign for 619 all the presynatpic terms throughout this document. 620

A.1.1 Feedback Alignment

Feedback Alignment (FA) or Random Backpropagation (RBP) refers to a family of algorithms [18, 19, 20, 21] that address the weights transport problem by using non-symmetric, and usually random, weights in the backward pass as follows: [22]

$$\Delta W_i = \eta B_i^{post} H_{i-1}^{pre}, \quad B_L^{post} = T - H_L, \quad (5) \quad 627$$
$$B_{i-1}^{post} = F'_{i-1} \circ B_i B_i^{post} \quad 628$$

$$B_{i-1}^{post} = F_{i-1}' \circ R_i B_i^{post}$$
⁶²⁸

where R_i denotes the random fixed matrices (random backward channels) to fix the issue of weights framework.

While FA and its variants address the weight 632 transport problem, by themselves they do not ad-633 dress the other problems. Among several flavors 634 of FA [19, 22], Direct Feedback Alignment (DFA) 635 [20], backpropagates the error signal obtained at the 636 top layer to each of the lower layers independently 637 using direct fixed random matrices. By this means, 638 DFA can address the issues of spatial locality and 639 distance implausibilities. The learning equation of 640 DFA can be written as follows: 641

$$\Delta W_i = \eta B_i^{post} H_{i-1}^{pre}, \quad B_i^{post} = R_i (T - H_L) \quad (6) \quad 642$$

A version with component-wise multiplication by 643 the derivatives of the corresponding activation functions is also possible. Experiments reported in the 645 (7)

697

707

literature suggest that FA and its variants do not
work well with convolutional layers [23, 24, 25]. A
few methods have been proposed to address this
apparent weakness of FA algorithms, however, most
of them introduce more constraints and dependence
on the forward pass which may make them less biologically plausible [24, 26, 27].

653 A.1.2 Difference Target Propagation

Difference Target Propagation (DTP) [5] as another biologically plausible model, trains the weights using a local target at each layer \hat{Y}_i that is propagated from the original target Y to each of the lower layers using learnable weights G_i .

659
$$\Delta W_i = \eta \hat{Y}_i^{post} H_{i-1}^{pre}, \quad \hat{Y}_i = \hat{Y}_{i+1} G_{i+1}$$

The G_i s are trained in the forward pass to approxi-660 661 mate the inverse of the forward operation at each layer. Propagating the target using G_i s at the top 662 two layers is dependent on the backpropagation and 663 weight transport [25]. Also, the forward and back-664 ward passes through the network are completely 665 clocked to learn G_i s. However, since the information 666 flows through layers independently, the variables are 667 local in space, thus, this architecture can address 668 space-locality and distance implausibilities. 669

670 A.1.3 Stacked Autoencoders

A well known approach to address the labeling issue 671 is using a stack of autoencoders [28, 29], where each 672 autoencoder learns to reproduce the hidden repre-673 sentation of the previous one in a self-supervised 674 675 manner, allowing the stack to learn increasingly abstract representations without labels. Labels are 676 only used to train the top layer in a supervised way, 677 with the option to fine-tune all layers via backprop-678 agation [30]. 679

This approach also addresses the distance and 680 developmental modularity issues since backpropaga-681 tion within each autoencoder limits error gradients 682 to short distances and allows training to begin be-683 fore the entire architecture is complete. However, 684 stacked autoencoders do not solve the locality and 685 weight transport issues. Each autoencoder, being 686 deep, requires backpropagation across at least two 687 adaptive layers, necessitating a learning channel for 688 error signals and symmetric weights to implement 689 backpropagation. 690

691 A.1.4 Forward Forward

The recently introduced Forward Forward algorithm (FF) [14], attempts to address the implausibility through a contrastive learning framework. Positive and negative data are fed through the network. Then the weights can be updated using a local target defined at each layer as follows:

$$\Delta W_i = \eta (||H_i^+ - H_i^-||^2)^{post} H_{i-1}^{pre} \tag{8}$$

FF uses variables that are local in space and can be assumed to be local in time (due to the short room neural distance). Given the contrastive learning room framework, it can be trained in a self-supervised room neural however, the computation remains heavily room clocked for feeding positive and negative data one room at a time. 705

A.2 Training Tourbillon CAEs 706

A.2.1 Recurculation Learning Rules

We start by introducing two new learning rules (in 708 addition to the main rule of recirculation which can 709 be found in Equation 1). 710

$$\Delta W_{i} = \begin{cases} \eta (H_{i}^{t} - H_{i}^{t'})^{post} (H_{i-1}^{t})^{pre}, & (a) \\ \eta (H_{i}^{t} - H_{i}^{t'})^{post} (H_{i-1}^{t'})^{pre}, & (b) & (9) \\ \eta (H_{i}^{t} - H_{i}^{t'})^{post} (H_{i}^{t} - H_{i-1}^{t'})^{pre} (c) \end{cases}$$

Rule (a) is the main learning rule of recirculation. (b) 712 is a slightly different version where the pre-synaptic 713 term is computed at the later cycle (t'). In our 714 experiments we observe that learning rule (b) provides trainability, however, it leads to a higher reconstruction error in a CAE across all datasets and 717 architectures. 718

Given that rules (c) and (a) have been identified 719 as the best-performing rules, our focus is on comparing these two rules. Specifically, we examine the characteristics of these rules in terms of their impact 722 on training dynamics. 723

In Figure A.1, we observe the training process 724 of a circular autoencoder using rule (c), which in-725 corporates symmetric terms for both pre- and post-726 synaptic activation differences. The graph illustrates 727 that the use of the symmetric learning rule leads to 728 the convergence of the loss function during training. 729 Moreover, compared to rule (a), this rule exhibits 730 smoother training dynamics. The smoother training 731 dynamics associated with rule (c) or the symmetric 732 learning rule suggest that it promotes more stable 733 and consistent updates to the model parameters, 734 leading to improved training performance. 735

Due to the inherent cyclic structure of circular 736 autoencoders, the data circulates through the model 737 in multiple time steps where neurons produce dif-738 ferent activations. Assuming that the time frame 739 between two consecutive cycles is short, the differ-740 ence between two consecutive activations of a neuron 741 $[H_i^{t+1} - H_i^t]$ can be interpreted as the activation rate 742 of that neuron at time t. Therefore, rule (c) be also 743 seen as the multiplication of the post-synaptic and 744 pre-synaptic neurons' activation rates. This implies 745 that the connection between two neurons (Δw_{ij}) 746



Figure A.1. The different behavior of the two recirculation learning rules (a) and (c). The REC learning rule corresponds to rule (a) and the symmetric learning rule corresponds to rule (c). While both demonstrate similar performance in training models using MNIST (first row) and Fashion MNIST (second row), the symmetric learning rule has smoother trajectories. All trajectories correspond to the best run out of three. The noise was removed for clarity.

must be strengthened (or weakened) if neurons' activation is correlated. This behavior closely matches
the concept of spike time-dependent synaptic plasticity Hebbian (or anti-Hebbian) learning rules which
was proposed using the temporal derivative of the activity of the post-synaptic neuron [13, 31] to encode
error derivatives.

However, during training a circular autoencoder 754 using rule (c), the model tends to be trapped into 755 a mode-collapse state where the reconstructed im-756 ages are the mean of the entire dataset. This mode 757 collapse state can be seen in Figure A.2. Therefore, 758 despite the interesting interpretation and intuition 759 behind rule (c) with its symmetric terms, the prob-760 lem of mode collapse confined our studies to the use 761 of the main recirculation learning associated with 762 rule (a). 763

764 A.2.2 Training CAEs

Table A.1 summarizes the parameters used for train-ing the CAEs in Section 4.1 of the main article.

Additionally, all models were trained for 100 epochs 767 with a batch size of 64. To optimize the activa-768 tion function and learning rates, a grid search was 769 conducted, resulting in the use of *tanh* activation 770 function and a learning rate of 0.01 for the initial 771 layers. Subsequently, a smaller learning rate of 0.001 772 was employed for the remaining fully connected lay-773 ers across all architectures. For the models that 774 incorporated convolutional layers and were trained 775 using CIFAR-10, a learning rate of 0.001 was used 776 for the initial layer, while a learning rate of 0.0001 777 778 was applied to the subsequent layers.



Figure A.2. Examples of reconstructed images from the MNIST and Fashion MINST datasets using the circular autoencoder trained with rule (c). There is a mode-collapse effect and the model reconstructs the mean of the data.

In the case of the convolutional models, zero 779 padding is crucial to maintain the size of the input, 780 ensuring that the spatial dimensions are preserved 781 during the convolutions. The plots in the first row 782 of Figure 5 display the training and test error curves 783 for the convolutional circular autoencoder trained 784 using the CIFAR-10 dataset. Notably, these exper-785 iments demonstrate that the use of recirculation 786 leads to comparable reconstruction errors, and in 787 some cases even superior, to those achieved with 788 traditional backpropagation and random backprop-789 agation techniques. 790



Figure A.3. Samples of reconstructed images from MNIST using autoencoders trained with BP, FA, and a CAE trained with recirculation.

Table A.1. CAE size refers to the number of hidden layers except for the input and output layers. For CIFAR-10, we only show the kernel size of the encoder part of the CAEs. The decoder of the CAEs used a symmetric kernel size.

CAE size	MNIST and Fashion MNIST	CIFAR-10	
	Hidden Layers Dim	Kernel	Stride
1	784-256-784	$(3 \times 5 \times 5 \times 3)$	(1×1)
3	784-256-64-256-784	$(3 \times 5 \times 5 \times 3)$ - $(3 \times 5 \times 5 \times 6)$	(1×1)

A.3 Stacking CAEs With Various Depth and Training Algorithms

Here we explain the details of the experiments conducted in Section 4.2. Specifically, Table A.2 summarizes the parameters used for stacking and training
the CAEs.

The size of the CAE architectures is explained in Table A.2. We use the same architectures for both sequential and asynchronous training algorithms. For the sequential training, each CAE is trained for 100 epochs with a batch size of 64. We use the learning rates explained in the previous section.

For the asynchronous training algorithm, we use a 803 batch size of 64 and we train the entire stack for 3000 804 iterations. According to Algorithm A.1, an iteration 805 refers to feeding one batch of data through the stack 806 and updating the weights of one CAE within the 807 808 stack. During our observations, we noted that when employing the same learning rate as the sequential 809 training algorithm, the reconstruction loss exhibited 810 a noisy trajectory with a limited convergence rate. 811 812 To address this issue, we conducted several grid searches to identify an optimal approach. Our find-813 ings indicated that utilizing lower learning rates for 814 the upper CAEs in the stack during the initial stages 815 of training was crucial. By gradually increasing the 816 learning rate of the upper CAEs, we observed a de-817 crease in the reconstruction loss, eventually aligning 818 all the CAEs in the stack to use the same learning 819 rate. This adjusted learning rate strategy enabled 820 more stable and efficient training, facilitating im-821 proved convergence of the reconstruction loss. We 822

CAE Image: Cae</td

Figure A.4. Random phases of the asynchronous training. Each time, one CAE is selected randomly and trained by recirculating the information.

provide a pseudocode of the asynchronous training in Algorithm A.1. To further enhance clarity, we have depicted the schematics of the asynchronous training algorithm in Figure A.4.

Algorithm A.1 Asynchronous train

Input: T: A stack of m sequential circular autoencoders $T = CAE_m \circ ... \circ CAE_1$, $CAE_i = D_i \circ E_i(datasample)$, data: training data, S: steps, E_i and D_i are the encoder and decoder of CAE_i for i = 1 to S do $1 \leq j = random \leq m$ $h = E_{j-1} \circ ... \circ E_1$ circulation(CAE_j, h)

end for

To evaluate the performance of the Tourbillon 827 architecture when adding the top classification layer, 828 we conducted a comparison with similar architec-829 tures trained using backpropagation, FA, and DFA. 830 Figure 4 presents the results of this experiment 831 specifically for the Tourbillons, trained sequentially 832 with a stack of three fully connected CAEs using the 833 Fashion MNIST dataset. Additionally, the second 834 row of Figure 5 showcases similar results obtained for 835 the Tourbillons, trained sequentially with a stack of 836 two convolutional CAEs using the CIFAR-10 dataset. 837 All the parameters of the CAEs are explained above. 838 In all cases, we can observe that Tourbillon can 839 achieve comparable performance to models trained 840 using BP or FA, showing the viability of Tourbillon 841 while being the most biologically plausible model. 842

A.4 Conversion to Tourbillon 843

As described in Section 4.4, we follow the Algorithm 844 A.2 to build the biologically plausible version of the 845 architecture for a U-Net autoencoder model that 846 addresses all the problems mentioned in the Intro-847 duction. In addition, we also convert a feed-forward 848 fully connected network architecture to its biologi-849 cally plausible version. Results obtained on MNIST 850 and Fashion MNIST are very similar, therefore we 851 reported the results obtained on MNIST and CIFAR-852 10. The U-Net architecture for the MNIST dataset 853 comprises a two-layer encoder and a two-layer de-854 coder. The encoder layers compress the data from 855

NLDL

kernel size of the encoder part of the CAEs. The decoder of the CAEs used a symmetric kernel size.

 Depth
 MNIST and Fashion MNIST
 CIFAR-10

 Input and Hinge Layers Dim
 Kernel
 Stride

Table A.2. Depth refers to the number of CAEs used to construct the stack. For CIFAR-10, we only show the

Depth	MINIST and Fashion MINIST	CIFAR-1	0
- •F •	Input and Hinge Layers Dim	Kernel	Stride
2	(784,256)- $(256,128)$	$(3 \times 5 \times 5 \times 3)$	(1×1)
3	(784,256)-(256,128)-(128,64)	$(3 \times 5 \times 5 \times 3)$	(1×1)
4	(784,256)- $(256,128)$ - $(128,64)$ - $(64,64)$		

Algorithm A.2 Conversion to Tourbillon

Input: M: A network with m sequential layers $L = [L_1, ..., L_m], data:$ training data. **Output:** M': The biological plausible version of the architecture. Initialize L' = []. for i = 1 to m do if L_i has learnable parameters then $l = train_cae([L_1, ..., L_i])$ L'.append(l)else $L'.append(L_i)$ end if end for M' = [data]for i = 1 to m do $M' = L'_i(M')$ end for function $train_cae([L_1, ..., L_i])$: $input = L_{i-1}(...(L_1(data)))$ circular_ae = $G_i(L_i(input)), \quad G_i:$ decoder *recirculation*(circular_ae) return L_i endfunction

784 to 128, and then from 128 to 64. The decoder 856 layers expand the data from 64 to 128, and then 857 from 128 to 784. The U-Net architecture for the 858 CIFAR-10 dataset comprises two 2D convolutional 859 layers with kernels of size (5×5) . For both the 860 MNIST and CIFAR-10 U-Nets we use the same 861 batch size, learning rate, number of epochs, and 862 activation function used in the previous section. For 863 the feed-forward architecture, we use a three-layer 864 network with 256, 64, and 10 hidden units. For this 865 architecture and its Tourbillon twin, we also use the 866 same batch size, learning rate, number of epochs, 867 and activation function as described in the previous 868 section. 869

Table A.3. Train/Test time and number of parameters used in the models. Models for MINST and Fashion MNIST are of identical sizes. Train times are the average time of one epoch. All times are in seconds.

Model Name	Train/Test Time	Parameters
	Section 4.1	
BP (MNIST)	3.95/411	403488
FA (MNIST)	4.05/420	403488
CAÈ (MNIST)	4.21/422	403488
BP (CIFAR-10)	2.40/231	244
FA (CIFAR-10)	2.44/231	244
CAE (CIFAR-10)	2.59/245	462
	Section 4.2	
BP (MNIST)	5.90/595	49220
FA (MNIST)	6.17/605	49220
SAE (MNIST)	6.20/620	981250
Tourb (MNIST/seq)	6.29/628	981250
Tourb (MNIST/asynch)	7.01/628	981250
BP (CIFAR-10)	3.66/350	720
FA (CIFAR-10)	3.66/350	720
SAE (CIFAR-10)	3.81/379	1380
Tourb (CIFAR-10/seq)	3.94/396	1380
Tourb (CIFAR-10/asynch)	4.19/420	1380
	Section 4.3	
Tourb-10% (MNIST)	7.01/628	981250
Tourb-25% (MNIST)	7.01'/628	981250
Tourb-50% (MNIST)	7.01/628	981250
Tourb-100% (MNIST)	7.01/628	981250
Tourb-10% (CIFAR-10)	4.19/420	1380
Tourb-25% (CIFAR-10)	4.19/420	1380
Tourb-50% (CIFAR-10)	4.19/420	1380
Tourb-100% (CIFAR-10)	4.19/420	1380
	Section 4.4	
U-Net (MNIST)	5.12/520	327100
Tourb-Ù-Net (MNIST)	6.71/590	654200
U-Net (CIFAR-10)	3.35/377	630
Tourb-U-Net (CIFAR-10)	3.88/390	1260
FC (MNIST)	5.88/600	49220
Tourb-FC (MNIST)	7.01/628	981250