

# Center for Brains, Minds & Machines

arXiv:1510.05067v3 [cs.LG] 2 Dec 2015

CBMM Memo No. 36

November 29, 2015

## How Important is Weight Symmetry in Backpropagation?

by

**Qianli Liao, Joel Z. Leibo, Tomaso Poggio**

Center for Brains, Minds and Machines, McGovern Institute, MIT

**Abstract:** Gradient backpropagation (BP) requires symmetric feedforward and feedback connections—the same weights must be used for forward and backward passes. This “weight transport problem” [1] is thought to be one of the main reasons of BP’s biological implausibility. Using 15 different classification datasets, we systematically study to what extent BP really depends on weight symmetry. In a study that turned out to be surprisingly similar in spirit to Lillicrap et al.’s demonstration [2] but orthogonal in its results, our experiments indicate that: (1) the magnitudes of feedback weights do not matter to performance (2) the signs of feedback weights do matter—the more concordant signs between feedforward and their corresponding feedback connections, the better (3) with feedback weights having random magnitudes and 100% concordant signs, we were able to achieve the same or even better performance than SGD. (4) some normalizations/stabilizations are indispensable for such asymmetric BP to work, namely Batch Normalization (BN) [3] and/or a “Batch Manhattan” (BM) update rule.



This work was supported by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC award CCF - 1231216.

# 1 Introduction

Deep Neural Networks (DNNs) have achieved remarkable performance in many domains [4, 5, 6, 7, 8, 9]. The simple gradient backpropagation (BP) algorithm has been the essential “learning engine” powering most of this work.

Deep neural networks are universal function approximators [10]. Thus it is not surprising that solutions to real-world problems exist within their configuration space. Rather, the real surprise is that such configurations can actually be discovered by gradient backpropagation.

The human brain may also be some form of DNN. Since BP is the most effective known method of adapting DNN parameters to large datasets, it becomes a priority to answer: could the brain somehow be implementing BP? Or some approximation to it?

For most of the past three decades since the invention of BP, it was generally believed that it could not be implemented by the brain [11, 12, 13, 14, 15]. BP seems to have three biologically implausible requirements: (1) feedback weights must be the same as feedforward weights (2) forward and backward passes require different computations, and (3) error gradients must somehow be stored separately from activations.

One biologically plausible way to satisfy requirements (2) and (3) is to posit a distinct “error network” with the same topology as the main (forward) network but used only for backpropagation of error signals. The main problem with such a model is that it makes requirement (1) implausible. There is no known biological way for the error network to know precisely the weights of the original network. This is known as the “weight transport problem” [1]. In this work we call it the “weight symmetry problem”. It is arguably the crux of BP’s biological implausibility.

In this report, we systematically relax BP’s weight symmetry requirement by manipulating the feedback weights. We find that some natural and biologically plausible schemes along these lines lead to exploding or vanishing gradients and render learning impossible. However, useful learning is restored if a simple and indeed *more* biologically plausible rule called Batch Manhattan (BM) is used to compute the weight updates. Another technique, called Batch Normalization (BN) [3], is also shown effective. When combined together, these two techniques seem complementary and significantly improve the performance of our asymmetric version of backpropagation.

The results are somewhat surprising: if the aforementioned BM and/or BN operations are applied, the magnitudes of feedback weights turn out not to be important. A much-relaxed *sign-concordance* property is all that is needed to attain comparable performance to mini-batch SGD on a large number of tasks.

Furthermore, we tried going beyond sign concordant feedback. We systematically reduced the probability of feedforward and feedback weights having the same sign (the *sign concordance probability*). We found that the effectiveness of backpropagation is strongly dependent on high sign concordance probability. That said, completely random and fixed feedback still outperforms chance e.g., as in the recent work of Lillicrap et al. [2].

Our results demonstrate that the perfect forward-backward weight symmetry requirement of backpropagation can be significantly relaxed and strong performance can still be achieved. To summarize, we have the following conclusions:

**(I) The magnitudes of feedback weights do not matter to performance.** This surprising result suggests that our theoretical understanding of why backpropagation works is probably far from complete.

**(II) Magnitudes of the weight updates also do not matter.**

**(III) Normalization / stabilization methods such as Batch Normalization and Batch Manhattan are necessary for these asymmetric backpropagation algorithms to work.** Note that this result was missed by previous work on random feedback weights [2].

**(IV) Asymmetric backpropagation algorithms evade the weight transport problem.** Thus it is plausible that the brain could implement them.

**(V) These results indicate that sign-concordance is very important for achieving strong performance.** However, even fixed random feedback weights with Batch Normalization significantly outperforms chance. This is intriguing and motivates further research.

**(VI) Additionally, we find Batch Manhattan to be a very simple but useful technique in general.** When used with Batch Normalization, it often improves the performance. This is especially true for smaller training sets.

## 2 Asymmetric Backpropagations

A schematic representation of backpropagation is shown in Fig. 1. Let  $E$  be the objective function. Let  $W$  and  $V$  denote the feedforward and feedback weight matrices respectively. Let  $X$  denote the inputs and  $Y$  the outputs.  $W_{ij}$  and  $V_{ij}$  are the feedforward and feedback connections between the  $j$ -th output  $Y_j$  and the  $i$ -th input  $X_i$ , respectively.  $f(\cdot)$  and  $f'(\cdot)$  are the transfer function and its derivative. Let the derivative of the  $i$ -th input with respect to the objective function be  $\frac{\partial E}{\partial X_i}$ , the formulations of forward and back propagation are as follows:

$$Y_j = f(N_j), \text{ where } N_j = \sum_i W_{ij} X_i \quad (1)$$

$$\frac{\partial E}{\partial X_i} = \sum_j V_{ij} f'(N_j) \frac{\partial E}{\partial Y_j} \quad (2)$$

The standard BP algorithm requires  $V = W$ . We call that case *symmetric backpropagation*. In this work we systematically explore the case of *asymmetric backpropagation* where  $V \neq W$ .

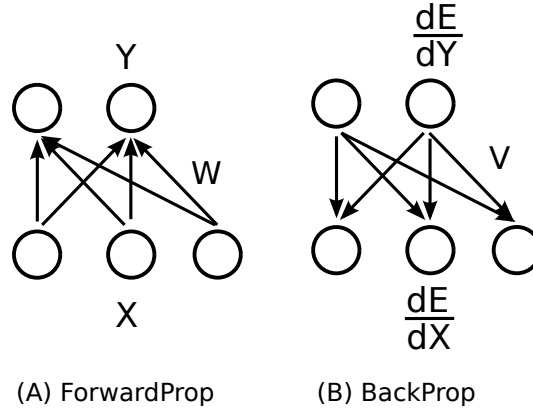


Figure 1: A simple illustration of backpropagation

By varying  $V$ , one can test various asymmetric BPs. Let  $sign(\cdot)$  denote the function that takes the sign (-1 or 1) of each element. Let  $\circ$  indicate element-wise multiplication.  $M, S$  are matrices of the same size as  $W$ .  $M$  is a matrix of uniform random numbers  $\in [0, 1]$  and  $S_p$  is a matrix where each element is either 1 with probability  $1 - p$  or  $-1$  with probability  $p$ . We explored the following choices of feedback weights  $V$  in this paper:

### 1. Uniform Sign-concordant Feedbacks (uSF):

$$V = sign(W)$$

### 2. Batchwise Random Magnitude Sign-concordant Feedbacks (brSF):

$V = M \circ sign(W)$ , where  $M$  is redrawn after each update of  $W$  (i.e., each mini-batch).

### 3. Fixed Random Magnitude Sign-concordant Feedbacks (frSF):

$V = M \circ sign(W)$ , where  $M$  is initialized once and fixed throughout each experiment.

### 4. Batchwise Random Magnitude p-percent-sign-concordant Feedbacks (brSF-p):

$V = M \circ sign(W) \circ S_p$ , where  $M$  and  $S_p$  is redrawn after each update of  $W$  (i.e., each mini-batch).

### 5. Fixed Random Magnitude p-percent-sign-concordant Feedbacks (frSF-p):

$V = M \circ sign(W) \circ S_p$ , where  $M$  and  $S_p$  is initialized once and fixed throughout each experiment.

### 6. Fixed Random Feedbacks (RndF):

Each feedback weight is drawn from a zero-mean gaussian distribution and fixed throughout each experiment:  $V \sim \mathcal{N}(0, \sigma^2)$ , where  $\sigma$  was chosen to be 0.05 in all experiments.

The results are summarized in the Section 5. The performances of 1, 2 and 3, which we call **strict sign-concordance** cases, are shown in Experiment A. The performances of 4 and 5 with different choices of  $p$ , which we call **partial sign-concordance** cases, are shown in Experiment B. The performances and control experiments about setting 6, which we call **no concordance** cases, are shown in Experiments C1 and C2.

### 3 Normalizations/stabilizations are necessary for “asymmetric” backpropagations

#### Batch Normalization (BN)

Batch Normalization (BN) is a recent technique proposed by [3] to reduce “internal covariate shift” [3]. The technique consists of element-wise normalization to zero mean and unit standard deviation. Means and standard deviations are separately computed for each batch. Note that in [3], the authors proposed the use of additional learnable parameters after the whitening. We found the effect of this operation to be negligible in most cases. Except for the “BN” and “BN+BM” entries (e.g., in Table 2), we did not use the learnable parameters of BN. Note that batch normalization may be related to the homeostatic plasticity mechanisms (e.g., Synaptic Scaling) in the brain [16, 17, 18].

#### Batch Manhattan (BM)

We were first motivated by looking at how BP could tolerate noisy operations that could be seen as more easily implementable by the brain. We tried relaxing the weight updates by discarding the magnitudes of the gradients. Let the weight at time  $t$  be  $w(t)$ , the update rule is:

$$w(t+1) = w(t) + \eta * \tau(t) \quad (3)$$

where  $\eta$  is the learning rate.

We tested several settings of  $\tau(t)$  as follows:

**Setting 0 (SGD):**  $\tau(t) = -\sum_b \frac{\partial E}{\partial w} + m * \tau(t-1) - d * w(t)$

**Setting 1:**  $\tau(t) = -\text{sign}(\sum_b \frac{\partial E}{\partial w}) + m * \tau(t-1) - d * w(t)$

**Setting 2:**  $\tau(t) = \text{sign}(-\text{sign}(\sum_b \frac{\partial E}{\partial w}) + m * \tau(t-1) - d * w(t))$

**Setting 3:**  $\tau(t) = \text{sign}(\kappa(t))$

where  $\kappa(t) = -\text{sign}(\sum_b \frac{\partial E}{\partial w}) + m * \kappa(t-1) - d * w(t)$

where  $m$  and  $d$  are momentum and weight decay rates respectively.  $\text{sign}()$  means taking the sign (-1 or 1),  $E$  is the objective function, and  $b$  denotes the indices of samples in the mini-batch. Setting 0 is the SGD algorithm (by “SGD” in this paper, we always refer to the mini-batch version with momentum and weight decay). Setting 1 is same as 0 but rounding the accumulated gradients in a batch to its sign. Setting 2 takes an extra final sign after adding the gradient term with momentum and weight decay terms. Setting 3 is something in between 1 and 2, where an final sign is taken, but not accumulated in the momentum term.

We found these techniques to be surprisingly powerful in the sense that they did not lower performance in most cases (as long as learning rates were reasonable). In fact, sometimes they improved performance. This was especially true for smaller training sets. Recall that asymmetric BPs tend to have exploding/vanishing gradients, these techniques are immune to such problems since the magnitudes of gradients are discarded.

We also found that the performance of this technique was significantly influenced by batch size on some experiments. In the cases of very small batch sizes, discarding the magnitudes of the weight updates was detrimental to performance (see Appendix).

This class of update rule is very similar to a technique called the Manhattan update rule [19, 20]. We suggest calling it “Batch Manhattan” (BM) to distinguish it from the stochastic version [21]. By default, we used setting 1 for BM throughout the Experiments A, B, C1 and C2. The “miscellaneous experiment” at the end of the Section 5.3 demonstrates that settings 1, 2 and 3 give similar performances, so the conclusions we draw broadly apply to all of them.

### 4 Related Work

Since the invention of backpropagation (BP) [22], its biological plausibility has been a long-standing controversy. Several authors have argued that BP is not biologically plausible [11, 12, 13, 14, 15]. Various biologically plausible modifications have been proposed. Most involve bidirectional connections e.g. Restricted Boltzmann Machines [23, 24] and so-called recirculation algorithms [25, 13] which despite their name provided, in the case of an autoencoder, an elegant early demonstration

that adaptive backwards weight can work without being identical to the forward ones. Recently, there have also been BP-free auto-encoders [26] based on “target propagation” [27].

The most relevant work to ours is a recent paper by Lillicrap et al. [2] of which we became aware after most of this work was done. Lillicrap et al. showed that fixed random feedback weights can support the learning of good representations for several simple tasks: (i) approximating a linear function, (ii) digit recognition on MNIST and (iii) approximating the outputs of a random 3 or 4 layer nonlinear neural network. Our work is very similar in spirit but rather different and perhaps complementary in its results, since we conclude that signs must be concordant between feedforward and corresponding feedback connections for consistent good performance, whereas the magnitudes do not matter, unlike Lillicrap et al. experiments in which both signs and magnitudes were random (but fixed). To explain the difference in our conclusions it is useful to consider the following points:

1. We systematically explored performance of the algorithms using 15 different datasets because simple tasks like MNIST by themselves do not always reveal differences between algorithms.
2. We tested deeper networks, since the accuracy of asymmetric BP’s credit assignment may critically attenuate with depth (for task (i) and (ii) Lillicrap et al. used a 3-layer (1 hidden layer) fully-connected network, and for task (iii) they used a 3 or 4 layer fully-connected network, whereas in most of our experiments, we use deeper and larger CNNs as shown in Table 1).
3. We found that local normalizations/stabilizations is critical for making asymmetric BP algorithms work. As shown by our results in Table 4, the random feedbacks scheme (i.e. the “RndF” column) suggested by Lillicrap et al. seem to work well only on one or two tasks, performing close to chance on most of them. Only when combined with Batch Normalization (“RndF+BN” or “RndF+BN+BM” in Table 4), it appears to become competitive.
4. We investigated several variants of asymmetric BPs such as sign-concordance (Table 2 and 3), batchwise-random vs. fixed-random feedbacks (Table 3) and learning with clamped layers (Table 4 Exp. C2).

## 5 Experiments

We were interested in relative differences between algorithms, not absolute performance. Thus we used common values for most parameters across all datasets to facilitate comparison.

**Method:** Key to our approach was the development of software allowing us to easily evaluate the “cartesian product” of models (experimental conditions) and datasets (tasks). Each experiment was a {model,dataset} pair, which was run 5 times using different learning rates (reporting the best performance).

### 5.1 Datasets

We extensively test our algorithms on 15 datasets of 5 Categories, including popular datasets of machine learning, object recognition, fine-grained recognition, face recognition, scene recognition and speech recognition. No data augmentation (e.g., cropping, flip, etc.) is used in any of the experiments.

The datasets are briefly described below:

#### 5.1.1 Machine Learning

1. MNIST[28]: handwritten digits recognition; 10 classes of digits 0 to 9; Inputs are 28x28 binary images; 60000 training and 10000 testing samples in total.
2. CIFAR-10[29]: classification of 10 object classes; Tiny 32x32 color natural images; 5000 training and 100 testing samples per class.
3. CIFAR-100[29]: similar to CIFAR-10, but with 100 classes. Training and testing samples per class decrease by 10 fold.
4. SVHN[30]: Street View House Numbers Recognition; Classification of 10 classes, one for each digit; 32x32 color images; Closely cropped version is used; 73257 digits for training, 26032 digits for testing; Extra data NOT used.

5. STL10[31]: classification of 10 object classes: airplane, bird, car, cat, deer, dog, horse, monkey, ship, truck; 96x96 color natural images; 500 training and 800 testing samples per class; There are, in addition, 100,000 unlabeled images, but NOT used in our experiments.

### 5.1.2 Fine-grained Recognition

1. Flowers17[32]: Classification of 17 categories of flowers. 60 training and 20 testing samples per class.
2. Flowers102[33]: Classification of 102 categories of flowers. Each class consists of between 40 and 258 images.

### 5.1.3 Basic-level Recognition

1. Caltech101[34]: classification of 102 classes (101 objects + background class). We use 30 training and 10 testing samples per class.
2. Caltech256-101[35]: to achieve a reasonable performance on this dataset (and also interesting to compare with Caltech101), we train/test on a subset of randomly sampled 102 classes. 30 training and 10 testing samples are used per class.
3. iCubWorld dataset[36]: We follow the “categorization” protocol; 10 class categorization task; 600 training and 200 testing samples per category.

### 5.1.4 Face Identification

We perform face identification tasks on three face datasets — Pubfig83, SUFR-W and Labeled Faces in the Wild. The testing protocol is somewhat unconventional, since the above datasets are usually used for face verification (predicting if two images depict the same person). In our experiments, we randomly sample 80 individuals from each dataset and perform a 80-way classification task.

1. Pubfig83-ID[37]: 80-individual identification task. 85 training and 15 testing samples per individual
2. SUFR-W-ID[38]: 80-individual identification task. 10 training and 5 testing samples per individual
3. LFW-ID[39]: 80-individual identification task. 10 training and 5 testing samples per individual

### 5.1.5 Scene Recognition

We tested our algorithms on one scene dataset.

1. MIT-indoor67[40]: classification of 67 indoor categories; 80 training and 20 testing samples per category

### 5.1.6 Non-visual tasks

For completeness, we also tested running a fully-connected neural network on a non-visual task: phoneme recognition on TIMIT dataset.

1. TIMIT-80[41]: the development set of TIMIT is used. For each frame, 123 dimensional MFCC features are extracted. We stack each frame with 7 frames before and after it, forming a 15x123 matrix, which is flattened to 1845x1 and fed into a 3-layer feed-forward neural network. We train/test by classifying each frame into 80 subphonetic structures (instead of the standard set of 183). There are 400 training and 100 testing samples per class.

	All others	MNIST	CIFAR10&100	SVHN	TIMIT-80
InputSize	119x119x3	28x28x1	32x32x3	32x32x3	1845x1x1
1	Conv 9x9x48/2	Conv 5x5x20/1	Conv 5x5x32/1	Conv 5x5x20/1	FC 512
2	Max-Pool 2x2/2	Max-Pool 2x2/2	Max-Pool 3x3/2	Max-Pool 2x2/2	FC 256
3	Conv 5x5x128/1	Conv 5x5x50/1	Conv 5x5x64/1	Conv 7x7x512/1	FC 80
4	Max-Pool 2x2/2	Max-Pool 2x2/2	Avg-Pool 3x3/2	Max-Pool 2x2/2	
5	FC max(256,#Class*3)	FC 500	Conv 5x5x64/1	FC 40	
6	FC #Class*2	FC 10	Avg-Pool 3x3/2	FC 10	
7	FC #Class		FC 128		
8			FC 10/100		

Table 1: Network architectures used in the experiments:  $AxBxC/D$  means  $C$  feature maps of size  $AxB$ , with stride  $D$ . The CIFAR10&100 architecture has a 2 units zero-padding for every convolution layer and 1 unit right-bottom zero-padding for every pooling layer. The other models do not have paddings. “FC  $X$ ” denotes Fully Connected layer of  $X$  feature maps. In the first model, the number of hidden units in FC layers are chosen according to the number of classes (denoted by “#Class”) in the classification task. “max(256,#Class\*3)” denotes 256 or #Class\*3, whichever is larger. Rectified linear units (ReLU) are used as nonlinearities for all models.

## 5.2 Training Details

The network architectures for various experiments are listed in Table 1. The input sizes of networks are shown in the second row of the table. All images are resized to fit the network if necessary.

Momentum was used with hyperparameter 0.9 (a conventional setting). All experiments were run for 65 epochs. The base learning rate: 1 to 50 epochs  $5 * 10^{-4}$ , 51 to 60 epochs  $5 * 10^{-5}$ , and 61 to 65 epochs  $5 * 10^{-6}$ . All models were run 5 times on each dataset with base learning rate multiplied by 100, 10, 1, 0.1, 0.01 respectively. This is because different learning algorithms favor different magnitudes of learning rates. The best validation error among all epochs of 5 runs was recorded as each model’s final performance. The batch sizes were all set to 100 unless stated otherwise. All experiments used a softmax for classification and the cross-entropy loss function.

## 5.3 Results

### Experiment A: sign-concordant Feedback

In this experiment, we show the performances of setting 1, 2 and 3 in Section 2, which we call **strict sign-concordance** cases: while keeping the signs of feedbacks the same as feedforward ones, the magnitudes of feedbacks are either randomized or set to uniform. The results are shown in Table 2 : Plus sign (+) denotes combination of methods. For example, uSF+BM means Batch Manhattan with uniform sign-concordant feedback. **SGD**: Stochastic gradient descent, the baseline algorithm. **BM**: SGD + Batch Manhattan. **BN**: SGD + Batch Normalization. **BN+BM**: SGD + Batch Normalization + Batch Manhattan. **uSF**: Uniform sign-concordant feedback. This condition often had exploding gradients. **NuSF**: same as uSF, but with feedback weights normalized by dividing the number of inputs of the feedforward filter (filter width \* filter height \* input feature number). This scheme avoids the exploding gradients but still suffers from vanishing gradients. **uSF+BM**: this setting is somewhat unstable for small batch sizes. Two training procedures were explored: (1) batchsize 100 for all epochs (2) batchsize 100 for 3 epochs and then batchsize 500 for the remaining epochs. The best performance was reported. While this gives a little advantage to this model since more settings were tried, we believe it is informative to isolate the stability issue and show what can be achieved if the model converges well. Note that uSF+BM is the only entry with slightly different training procedures. All other models share exactly the same training procedure & parameters. **uSF+BN**, **uSF+BN+BM**, **brSF+BN+BM**, **frSF+BN+BM**: These are some combinations of uSF, brSF, frSF, BN and BM. The last three are the most robust, well-performing ,and biologically-plausible algorithms.

### Experiment B: Violating Sign-Concordance with probability $p$

In this experiment, we test the effect of **partial sign-concordance**. That is, we test settings 4 and 5 as described in Section 2. In these cases, the feedback weight magnitudes were all random. Strict sign-concordance was relaxed by manipulating the

Experiment A	SGD	BM	BN	BN+BM	uSF	NuSF	uSF +BM	uSF +BN	uSF +BN +BM	brSF +BN +BM	frSF +BN +BM
MNIST	0.67	0.99	0.52	0.73	<b>88.65</b>	0.60	0.95	0.55	0.83	0.80	0.91
CIFAR	22.73	23.98	<b>16.75</b>	<b>17.94</b>	<b>90.00</b>	<b>40.60</b>	<b>26.25</b>	<b>19.48</b>	<b>19.29</b>	<b>18.44</b>	<b>19.02</b>
CIFAR100	55.15	<b>58.44</b>	<b>49.44</b>	<b>51.45</b>	<b>99.00</b>	<b>71.51</b>	<b>65.28</b>	57.19	53.12	<b>50.74</b>	52.25
SVHN	9.06	10.77	7.50	9.88	<b>80.41</b>	<b>14.55</b>	9.78	8.73	9.67	9.95	10.16
STL10	48.01	<b>44.14</b>	45.19	<b>43.19</b>	<b>90.00</b>	<b>56.53</b>	46.41	48.49	<b>41.55</b>	<b>42.74</b>	<b>42.68</b>
Cal101	74.08	<b>66.70</b>	<b>66.07</b>	<b>61.75</b>	<b>98.95</b>	<b>70.50</b>	75.24	<b>63.33</b>	<b>60.70</b>	<b>59.54</b>	<b>60.27</b>
Cal256-101	87.06	<b>83.43</b>	<b>82.94</b>	<b>81.96</b>	<b>99.02</b>	85.98	86.37	<b>82.16</b>	<b>80.78</b>	<b>78.92</b>	<b>80.59</b>
iCub	57.62	55.57	<b>46.43</b>	<b>37.08</b>	<b>89.96</b>	<b>66.57</b>	<b>70.61</b>	<b>61.37</b>	<b>48.38</b>	<b>47.33</b>	<b>46.08</b>
Flowers17	35.29	<b>31.76</b>	36.76	32.35	<b>94.12</b>	<b>42.65</b>	38.24	35.29	32.65	<b>29.41</b>	<b>31.47</b>
Flowers102	77.30	77.57	75.78	74.92	<b>99.67</b>	77.92	79.25	<b>71.74</b>	<b>73.20</b>	<b>73.31</b>	<b>73.57</b>
PubFig83-ID	63.25	<b>54.42</b>	<b>51.08</b>	<b>41.33</b>	<b>98.75</b>	<b>78.58</b>	65.83	<b>54.58</b>	<b>40.67</b>	<b>42.67</b>	<b>40.33</b>
SUFR-W-ID	80.00	<b>74.25</b>	<b>75.00</b>	<b>65.00</b>	<b>98.75</b>	<b>83.50</b>	79.50	<b>72.00</b>	<b>65.75</b>	<b>66.25</b>	<b>66.50</b>
LFW-ID	79.25	<b>74.25</b>	<b>73.75</b>	<b>55.75</b>	<b>98.75</b>	<b>85.75</b>	80.75	<b>73.75</b>	<b>56.25</b>	<b>57.25</b>	<b>55.75</b>
Scene67	87.16	85.37	86.04	<b>82.46</b>	<b>98.51</b>	88.21	87.09	87.09	<b>81.87</b>	<b>82.31</b>	<b>81.79</b>
TIMIT80	23.04	25.92	23.92	24.40	23.60	<b>29.28</b>	25.84	25.04	25.12	25.24	24.92

Table 2: **Experiment A:** The magnitudes of feedbacks do not matter. Sign concordant feedbacks can produce strong performance. Numbers are error rates (%). **Yellow:** performances worse than baseline(SGD) by 3% or more. **Blue:** performances better than baseline(SGD) by 3% or more.

probability  $p$  of concordance between feedforward and feedback weight signs. Feedback weights  $V = M \circ \text{sign}(W) \circ S_p$  depend on the matrix  $S_p$  as defined in Section 2. Table 3 Part 1 reports results from setting 4, the case where a new  $M$  and  $S_p$  is sampled for each batch. Table 3 Part 2 reports results of setting 5, the case where  $M$  and  $S_p$  are fixed.

### Experiment C1: Fixed Random Feedback

Results are shown in Table 4: **RndF:** fixed random feedbacks. **RndF+BN, RndF+BN+BM:** some combinations of RndF, BN and BM. **uSF+BN+BM:** one of our best algorithms, for reference. The “RndF” setting is the same as the one proposed by [2]. Apparently it does not perform well on most datasets. However, combining it with Batch Normalization makes it significantly better. Although it remains somewhat worse than its sign concordant counterpart. Another observation is that random feedback does not work well with BM alone (but better with BN+BM).

### Experiment C2: Control experiments for Experiment C1

The fact that random feedback weights can learn meaningful representations is somewhat surprising. We explore this phenomenon by running some control experiments, where we run two control models for each model of interest: **1. (.)Bottom:** The model’s last layer is initialized randomly and clamped/frozen. All learning happens in the layers before the last layer. **2. (.)Top:** The model’s layers before the last layer are initialized randomly and clamped/frozen. All learning happens in the last layer.

There are some observations: **(i)** When only the last layer was allowed to adapt, all models behaved similarly. This was expected since the models only differed in the way they backpropagate errors. **(ii)** With the last layer is clamped, random feedback cannot learn meaningful representations. **(iii)** In contrast, the models with sign concordant feedback can learn surprisingly good representations with the last layer locked—almost as if the last layer was not locked. We can draw the following conclusions from such observations: sign concordant feedback ensures that meaningful error signals reach lower layers by itself, while random feedback is not sufficient. If all layers can learn, random feedback seems to work via a “mysterious co-adaptation” between the last layer and the layers before it. This “mysterious co-adaptation” was first observed by [2], where it was called “feedback alignment” and some analyses were given. Note that our Experiment C shows that the effect is more significant if Batch Normalization is applied.

**Miscellaneous Experiment: different settings of Batch Manhattan** We show that the 3 settings of BM (as described in Section 3) produce similar performances. This is the case for both symmetric and asymmetric BPs. Results are in Table 5.



Experiment B	Baseline	Part 1: Random $M$ and $S$ every batch					Part 2: Fixed random $M$ and $S$				
	SGD	100%	75%	50%	25%	Same Sign (brSF +BN+BM)	100%	75%	50%	25%	Same Sign (frSF +BN+BM)
MNIST	0.67	<b>7.87</b>	<b>8.34</b>	<b>7.54</b>	1.01	0.80	<b>4.25</b>	3.56	1.37	0.84	0.91
CIFAR	22.73	<b>71.18</b>	<b>75.87</b>	<b>70.60</b>	19.98	<b>18.44</b>	<b>71.41</b>	<b>68.49</b>	<b>31.54</b>	20.85	<b>19.02</b>
CIFAR100	55.15	<b>93.72</b>	<b>96.23</b>	<b>94.58</b>	<b>68.98</b>	<b>50.74</b>	<b>96.26</b>	<b>95.22</b>	<b>71.98</b>	56.02	52.25
SVHN	9.06	<b>75.25</b>	<b>77.91</b>	<b>74.64</b>	10.94	9.95	<b>37.78</b>	<b>32.72</b>	<b>15.02</b>	11.50	10.16
STL10	48.01	<b>69.65</b>	<b>72.50</b>	<b>72.10</b>	<b>44.82</b>	<b>42.74</b>	<b>72.46</b>	<b>68.96</b>	50.54	<b>43.85</b>	<b>42.68</b>
Cal101	74.08	<b>91.46</b>	<b>93.99</b>	<b>91.36</b>	<b>67.65</b>	<b>59.54</b>	<b>94.20</b>	<b>87.57</b>	<b>68.60</b>	<b>63.12</b>	<b>60.27</b>
Cal256-101	87.06	<b>93.24</b>	<b>94.02</b>	<b>92.75</b>	<b>82.75</b>	<b>78.92</b>	<b>95.98</b>	<b>90.88</b>	84.22	<b>83.04</b>	<b>80.59</b>
iCub	57.62	<b>75.56</b>	<b>79.36</b>	<b>83.61</b>	<b>52.42</b>	<b>47.33</b>	<b>68.97</b>	<b>63.82</b>	<b>64.62</b>	<b>51.22</b>	<b>46.08</b>
Flowers17	35.29	<b>70.29</b>	<b>82.06</b>	<b>73.82</b>	35.00	<b>29.41</b>	<b>79.12</b>	<b>61.76</b>	<b>44.12</b>	35.59	<b>31.47</b>
Flowers102	77.30	<b>90.54</b>	<b>92.58</b>	<b>89.56</b>	<b>73.18</b>	<b>73.31</b>	<b>93.69</b>	<b>93.77</b>	77.62	77.85	<b>73.57</b>
PubFig83-ID	63.25	<b>95.33</b>	<b>95.67</b>	<b>94.17</b>	61.25	<b>42.67</b>	<b>94.25</b>	<b>89.25</b>	64.33	<b>47.17</b>	<b>40.33</b>
SUFR-W-ID	80.00	<b>95.00</b>	<b>95.75</b>	<b>92.50</b>	<b>71.00</b>	<b>66.25</b>	<b>95.50</b>	<b>95.50</b>	77.25	<b>68.00</b>	<b>66.50</b>
LFW-ID	79.25	<b>95.25</b>	<b>96.00</b>	<b>94.50</b>	<b>64.50</b>	<b>57.25</b>	<b>95.75</b>	<b>95.25</b>	<b>75.00</b>	<b>59.75</b>	<b>55.75</b>
Scene67	87.16	<b>94.48</b>	<b>94.78</b>	<b>93.43</b>	<b>83.13</b>	<b>82.31</b>	<b>95.37</b>	<b>92.69</b>	88.36	84.40	<b>81.79</b>
TIMIT80	23.04	<b>55.32</b>	<b>61.36</b>	<b>55.28</b>	<b>26.48</b>	25.24	<b>62.60</b>	<b>33.48</b>	<b>27.56</b>	<b>26.08</b>	24.92

Table 3: **Experiment B Part 1 (left):** Feedbacks have random magnitudes, varying probability of having different signs (percentages in second row, column 3-7) from the feedforward ones. The  $M$  and  $S$  redrawn in each mini-batch. Numbers are error rates (%). **Yellow:** performances worse than baseline(SGD) by 3% or more. **Blue:** performances better than baseline(SGD) by 3% or more. **Experiment B Part 2 (right):** Same as part 1, but The  $M$  and  $S$  were fixed throughout each experiment.

## 6 Discussion

This work aims to establish that there exist variants of the gradient backpropagation algorithm that could plausibly be implemented in the brain. To that end we considered the question: how important is weight symmetry to backpropagation? Through a series of experiments with increasingly asymmetric backpropagation algorithms, our work complements a recent demonstration[2] that perfect weight symmetry can be significantly relaxed while still retaining strong performance.

These results show that Batch Normalization and/or Batch Manhattan are crucial for asymmetric backpropagation to work. Furthermore, they are complementary operations that are better used together than individually. These results highlight the importance of sign-concordance to asymmetric backpropagation by systematically exploring how performance declines with its relaxation.

Finally, let us return to our original motivation. How does all this relate to the brain? Based on current neuroscientific understanding of cortical feedback, we cannot make any claim about whether such asymmetric BP algorithms are actually implemented by the brain. Nevertheless, this work shows that asymmetric BPs, while having less constraints, are not computationally inferior to standard BP. So if the brain were to implement one of them, it could obtain most of the benefits of the standard algorithm.

This work suggests a hypothesis that could be checked by empirical neuroscience research: if the brain does indeed implement an asymmetric BP algorithm, then there is likely to be a high degree of sign-concordance in cortical forward-backward connections.

These empirical observations concerning Batch Manhattan updating may shed light on the general issue of how synaptic plasticity may implement learning algorithms. They show that changes of synaptic strength can be rather noisy. That is, the *sign* of a long term accumulation of synaptic potentiation or depression, rather than precise magnitude, is what is important. This scheme seems biologically implementable.

**Acknowledgements** We thank G. Hinton for useful comments. This work was supported by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC award CCF 1231216.

Experiment C1	SGD	RndF	NuSF	BN	RndF+BN	RndF+BM	RndF+BN+BM	uSF+BN+BM
MNIST	0.67	1.81	0.60	0.52	0.83	1.89	1.07	0.83
CIFAR	22.73	<b>42.69</b>	<b>40.60</b>	<b>16.75</b>	24.35	<b>62.71</b>	<b>25.75</b>	<b>19.29</b>
CIFAR100	55.15	<b>90.88</b>	<b>71.51</b>	<b>49.44</b>	<b>60.83</b>	<b>97.11</b>	<b>64.69</b>	53.12
SVHN	9.06	<b>12.35</b>	<b>14.55</b>	7.50	<b>12.63</b>	<b>13.63</b>	<b>12.79</b>	9.67
STL10	48.01	<b>57.80</b>	<b>56.53</b>	45.19	<b>51.60</b>	<b>80.39</b>	47.39	<b>41.55</b>
Cal101	74.08	<b>88.51</b>	<b>70.50</b>	<b>66.07</b>	72.81	<b>98.42</b>	<b>67.12</b>	<b>60.70</b>
Cal256-101	87.06	<b>94.12</b>	85.98	<b>82.94</b>	85.49	<b>98.14</b>	<b>83.63</b>	<b>80.78</b>
iCub	57.62	<b>67.87</b>	<b>66.57</b>	<b>46.43</b>	58.82	<b>84.41</b>	59.02	<b>48.38</b>
Flowers17	35.29	<b>69.41</b>	<b>42.65</b>	36.76	<b>43.53</b>	<b>91.18</b>	38.24	32.65
Flowers102	77.30	<b>92.31</b>	77.92	75.78	<b>81.22</b>	<b>96.13</b>	78.99	<b>73.20</b>
PubFig83-ID	63.25	<b>95.42</b>	<b>78.58</b>	<b>51.08</b>	<b>67.00</b>	<b>97.67</b>	<b>55.25</b>	<b>40.67</b>
SUFR-W-ID	80.00	<b>94.75</b>	<b>83.50</b>	<b>75.00</b>	77.75	<b>97.25</b>	<b>69.00</b>	<b>65.75</b>
LFW-ID	79.25	<b>95.75</b>	<b>85.75</b>	<b>73.75</b>	79.25	<b>97.75</b>	<b>67.50</b>	<b>56.25</b>
Scene67	87.16	<b>95.75</b>	88.21	86.04	87.84	<b>97.69</b>	87.09	<b>81.87</b>
TIMIT80	23.04	<b>26.76</b>	<b>29.28</b>	23.92	<b>26.52</b>	<b>33.12</b>	<b>26.32</b>	25.12

Experiment C2	SGD	SGD	RndF	RndF	RndF+BN+BM	RndF+BN+BM	uSF+BN+BM	uSF+BN+BM
	<i>Bottom</i>	<i>Top</i>	<i>Bottom</i>	<i>Top</i>	<i>Bottom</i>	<i>Top</i>	<i>Bottom</i>	<i>Top</i>
MNIST	0.65	<b>3.85</b>	<b>85.50</b>	<b>3.81</b>	<b>86.74</b>	<b>3.81</b>	0.66	<b>3.85</b>
CIFAR	23.12	<b>56.80</b>	<b>89.71</b>	<b>56.77</b>	<b>78.90</b>	<b>58.54</b>	<b>16.72</b>	<b>57.84</b>
CIFAR100	<b>59.49</b>	<b>80.71</b>	<b>98.79</b>	<b>80.65</b>	<b>98.69</b>	<b>84.34</b>	<b>61.61</b>	<b>84.10</b>
SVHN	8.31	<b>75.22</b>	<b>82.86</b>	<b>75.12</b>	<b>84.84</b>	<b>69.99</b>	10.96	<b>71.89</b>
STL10	49.96	<b>74.69</b>	<b>88.36</b>	<b>72.44</b>	<b>81.31</b>	<b>76.11</b>	<b>52.18</b>	<b>76.09</b>
Cal101	71.97	<b>82.72</b>	<b>98.63</b>	<b>79.14</b>	<b>98.21</b>	<b>80.40</b>	<b>63.86</b>	<b>79.98</b>
Cal256-101	86.08	89.71	<b>98.43</b>	88.92	<b>98.14</b>	89.02	<b>82.84</b>	89.12
iCub	<b>46.73</b>	<b>83.96</b>	<b>87.56</b>	<b>83.26</b>	<b>80.36</b>	<b>84.31</b>	<b>49.33</b>	<b>84.16</b>
Flowers17	38.24	<b>70.59</b>	<b>92.35</b>	<b>70.00</b>	<b>87.35</b>	<b>77.06</b>	<b>45.00</b>	<b>77.06</b>
Flowers102	78.99	<b>86.57</b>	<b>97.89</b>	<b>86.84</b>	<b>98.11</b>	<b>84.24</b>	78.09	<b>84.57</b>
PubFig83-ID	<b>66.75</b>	<b>89.58</b>	<b>97.67</b>	<b>89.58</b>	<b>97.67</b>	<b>89.67</b>	<b>43.83</b>	<b>89.50</b>
SUFR-W-ID	80.50	<b>90.50</b>	<b>97.50</b>	<b>90.50</b>	<b>97.50</b>	<b>89.50</b>	<b>71.50</b>	<b>89.50</b>
LFW-ID	79.75	<b>92.50</b>	<b>98.25</b>	<b>93.00</b>	<b>97.00</b>	<b>89.50</b>	<b>65.00</b>	<b>89.50</b>
Scene67	88.73	<b>91.57</b>	<b>97.84</b>	<b>91.49</b>	<b>97.54</b>	<b>91.19</b>	85.97	<b>91.04</b>
TIMIT80	23.52	<b>46.20</b>	<b>95.00</b>	<b>46.20</b>	<b>93.00</b>	<b>39.76</b>	24.96	<b>40.16</b>

Table 4: **Experiment C1:** fixed random feedbacks. **Experiment C2:** (*Bottom*): The model’s last layer is initialized randomly and clamped/frozen. All learning happens in the layers before the last layer. (*Top*): The model’s layers before the last layer are initialized randomly and clamped/frozen. All learning happens in the last layer. Numbers are error rates (%). **Yellow:** performances worse than baseline(SGD) by 3% or more. **Blue:** performances better than baseline(SGD) by 3% or more.

	SGD	BM1	BM2	BM3	uSF+BN	uSF+BN+BM1	uSF+BN+BM2	uSF+BN+BM3
MNIST	0.67	0.99	1.30	1.09	0.55	0.83	0.72	0.61
CIFAR	22.73	23.98	23.09	20.47	<b>19.48</b>	<b>19.29</b>	<b>18.87</b>	<b>18.38</b>
CIFAR100	55.15	<b>58.44</b>	<b>58.81</b>	52.82	57.19	53.12	52.38	54.68
SVHN	9.06	10.77	<b>12.31</b>	<b>12.23</b>	8.73	9.67	10.54	9.20
STL10	48.01	<b>44.14</b>	<b>44.74</b>	45.23	48.49	<b>41.55</b>	47.71	46.45
Cal101	74.08	<b>66.70</b>	<b>65.96</b>	<b>70.28</b>	<b>63.33</b>	<b>60.70</b>	<b>64.38</b>	<b>62.59</b>

Table 5: Different settings of Batch Manhattan (as described in Section 3) seem to give similar performances. SGD: setting 0, BM1: setting 1, BM2: setting 2, BM3: setting 3. The interaction of BM with sign concordant feedback weights (uSF) and Batch Normalization are shown in “uSF+BN(+)” entries. Numbers are error rates (%). **Yellow:** performances worse than baseline (SGD) by 3% or more. **Blue:** performances better than baseline(SGD) by 3% or more.

## References

- [1] S. Grossberg, “Competitive learning: From interactive activation to adaptive resonance,” *Cognitive science*, vol. 11, no. 1, pp. 23–63, 1987.
- [2] T. P. Lillicrap, D. Cownden, D. B. Tweed, and C. J. Akerman, “Random feedback weights support learning in deep neural networks,” *arXiv preprint arXiv:1411.0247*, 2014.
- [3] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015.
- [4] A. Krizhevsky, I. Sutskever, and G. Hinton, “ImageNet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, (Lake Tahoe, CA), 2012.
- [5] O. Abdel-Hamid, A. Mohamed, H. Jiang, and G. Penn, “Applying convolutional neural networks concepts to hybrid NN-HMM model for speech recognition,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4277–4280, 2012.
- [6] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, *et al.*, “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups,” *Signal Processing Magazine, IEEE*, vol. 29, no. 6, pp. 82–97, 2012.
- [7] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” in *Advances in Neural Information Processing Systems (NIPS)*, pp. 3111–3119, 2013.
- [8] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, “Deepface: Closing the gap to human-level performance in face verification,” in *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pp. 1701–1708, IEEE, 2014.
- [9] A. Graves, G. Wayne, and I. Danihelka, “Neural Turing machines,” *arXiv preprint arXiv:1410.5401*, 2014.
- [10] K. Hornik, M. Stinchcombe, and H. White, “Multilayer feedforward networks are universal approximators,” *Neural networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [11] F. Crick, “The recent excitement about neural networks.,” *Nature*, vol. 337, no. 6203, pp. 129–132, 1989.
- [12] P. Mazzoni, R. A. Andersen, and M. I. Jordan, “A more biologically plausible learning rule for neural networks.,” *Proceedings of the National Academy of Sciences*, vol. 88, no. 10, pp. 4433–4437, 1991.
- [13] R. C. O’Reilly, “Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm,” *Neural computation*, vol. 8, no. 5, pp. 895–938, 1996.
- [14] L. V. Chintala and D. B. Tweed, “Adaptive optimal control without weight transport,” *Neural computation*, vol. 24, no. 6, pp. 1487–1518, 2012.
- [15] Y. Bengio, D.-H. Lee, J. Bornschein, and Z. Lin, “Towards biologically plausible deep learning,” *arXiv preprint arXiv:1502.04156*, 2015.
- [16] G. G. Turrigiano and S. B. Nelson, “Homeostatic plasticity in the developing nervous system,” *Nature Reviews Neuroscience*, 2004.
- [17] D. Stellwagen and R. C. Malenka, “Synaptic scaling mediated by glial TNF- $\alpha$ ,” *Nature*, vol. 440, no. 7087, pp. 1054–1059, 2006.
- [18] G. G. Turrigiano, “The self-tuning neuron: synaptic scaling of excitatory synapses,” *Cell*, vol. 135, no. 3, pp. 422–435, 2008.
- [19] J. Heaton, “Manhattan propagation, encog.” <http://www.heatonresearch.com/javadoc/encog-3.0/org/encog/neural/networks/training/propagation/manhattan/ManhattanPropagation.html>.
- [20] M. Riedmiller and H. Braun, “A direct adaptive method for faster backpropagation learning: The rprop algorithm,” in *Neural Networks, 1993., IEEE International Conference on*, pp. 586–591, IEEE, 1993.
- [21] E. Zamanidoost, F. M. Bayat, D. Strukov, and I. Kataeva, “Manhattan rule training for memristive crossbar circuit pattern classifiers,” 2015.
- [22] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Cognitive modeling*, 1988.
- [23] G. E. Hinton and R. R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [24] P. Smolensky, “Information processing in dynamical systems: Foundations of harmony theory,” 1986.
- [25] G. E. Hinton and J. L. McClelland, “Learning representations by recirculation,” in *Neural information processing systems*, pp. 358–366, New York: American Institute of Physics, 1988.
- [26] Y. Bengio, “How auto-encoders could provide credit assignment in deep networks via target propagation,” *arXiv preprint arXiv:1407.7906*, 2014.
- [27] Y. Le Cun, “Learning process in an asymmetric threshold network,” in *Disordered systems and biological organization*, pp. 233–240, Springer, 1986.
- [28] Y. LeCun, C. Cortes, and C. J. Burges, “The mnist database: <http://www.hollywoodreporter.com/news/earthquake-twitter-users-learned-tremors-226481>.”
- [29] A. Krizhevsky, “Learning multiple layers of features from tiny images,” 2009.

- [30] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, “Reading digits in natural images with unsupervised feature learning,” in *NIPS workshop on deep learning and unsupervised feature learning*, vol. 2011, p. 5, 2011.
- [31] A. Coates, A. Y. Ng, and H. Lee, “An analysis of single-layer networks in unsupervised feature learning,” in *International conference on artificial intelligence and statistics*, pp. 215–223, 2011.
- [32] M.-E. Nilsback and A. Zisserman, “A visual vocabulary for flower classification,” in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, IEEE, 2006.
- [33] M.-E. Nilsback and A. Zisserman, “Automated flower classification over a large number of classes,” in *Computer Vision, Graphics & Image Processing, 2008. ICVGIP’08. Sixth Indian Conference on*, IEEE, 2008.
- [34] L. Fei-Fei, R. Fergus, and P. Perona, “Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories,” *Computer Vision and Image Understanding*, vol. 106, no. 1, pp. 59–70, 2007.
- [35] G. Griffin, A. Holub, and P. Perona, “Caltech-256 object category dataset,” 2007.
- [36] S. R. Fanello, C. Ciliberto, M. Santoro, L. Natale, G. Metta, L. Rosasco, and F. Odone, “icub world: Friendly robots help building good vision data-sets,” in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2013 IEEE Conference on*, pp. 700–705, IEEE, 2013.
- [37] N. Pinto, Z. Stone, T. Zickler, and D. Cox, “Scaling up biologically-inspired computer vision: A case study in unconstrained face recognition on facebook,” in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Computer Society Conference on*, pp. 35–42, IEEE, 2011.
- [38] J. Z. Leibo, Q. Liao, and T. Poggio, “Subtasks of Unconstrained Face Recognition,” in *International Joint Conference on Computer Vision, Imaging and Computer Graphics, VISIGRAPP*, (Lisbon, Portugal), 2014.
- [39] G. B. Huang, M. Mattar, T. Berg, and E. Learned-Miller, “Labeled faces in the wild: A database for studying face recognition in unconstrained environments,” in *Workshop on faces in real-life images: Detection, alignment and recognition (ECCV)*, (Marseille, Fr), 2008.
- [40] A. Quattoni and A. Torralba, “Recognizing indoor scenes,” in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pp. 413–420, IEEE, 2009.
- [41] J. Garofolo, L. Lamel, W. Fisher, J. Fiscus, D. Pallett, N. Dahlgren, and V. Zue, “Timit acoustic-phonetic continuous speech corpus. <https://catalog.ldc.upenn.edu/ldc93s1>.”

## A The influence of batch size on Batch Manhattan

We show the influence of batch size on Batch Manhattan (setting 1) in Table 6.

	SGD-100	SGD-30	SGD-10	SGD-5	BM-100	BM-30	BM-10	BM-5
MNIST	0.69	0.66	0.52	0.67	1.00	1.74	2.72	2.66
CIFAR	22.74	22.09	22.28	22.69	23.56	24.53	<b>31.64</b>	<b>35.53</b>
CIFAR100	54.94	55.81	55.16	56.56	57.81	56.17	<b>78.49</b>	<b>83.14</b>
SVHN	9.02	8.45	8.54	8.47	<b>12.83</b>	<b>17.06</b>	<b>30.60</b>	<b>51.19</b>
STL10	47.86	46.86	47.16	46.09	<b>44.00</b>	<b>44.34</b>	<b>44.85</b>	<b>42.75</b>

Table 6: SGD-x denotes Stochastic Gradient Descent with batchsize x. BM-x denotes Batch Manhattan (setting 1) with batchsize x. Performance on some datasets declines as the batch size decreases. Numbers are error rates (%). **Yellow**: performances worse than baseline(SGD-100) by 3% or more. **Blue**: performances better than baseline(SGD-100) by 3% or more. The performance of Batch Manhattan is proportional to the batch size on some datasets. This effect seems to not present in SGD.