
Block-local learning with probabilistic latent representations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The ubiquitous backpropagation algorithm requires sequential updates across
2 blocks of a network, introducing a locking problem. Moreover, backpropaga-
3 tion relies on the transpose of weight matrices to calculate updates, introducing a
4 weight transport problem across blocks. Both these issues prevent efficient par-
5 allelisation and horizontal scaling of models across devices. We propose a new
6 method that introduces a twin network that propagates information backwards from
7 the targets to the input to provide auxiliary local losses. Forward and backward
8 propagation can work in parallel and with different sets of weights, addressing the
9 problems of weight transport and locking. Our approach derives from a statistical
10 interpretation of end-to-end training which treats activations of network layers as
11 parameters of probability distributions. The resulting learning framework uses
12 these parameters locally to assess the matching between forward and backward
13 information. Error backpropagation is then performed locally within each block,
14 leading to “block-local” learning. Several previously proposed alternatives to error
15 backpropagation emerge as special cases of our model. We present results on vari-
16 ous tasks and architectures, including transformers, demonstrating state-of-the-art
17 performance using block-local learning. These results provide a new principled
18 framework to train very large networks in a distributed setting and can also be
19 applied in neuromorphic systems.

20 1 Introduction

21 Recent developments in machine learning have seen deep neural network architectures scaling to
22 billions of parameters [Touvron et al., 2023, Brown et al., 2020]. This development has boosted
23 the capabilities of these models to unprecedented levels but simultaneously pushed the computing
24 hardware on which large network models are running to its limits. It is therefore becoming increas-
25 ingly important to distribute learning algorithms over a large number of independent compute nodes.
26 However, today’s machine learning algorithms are ill-suited for distributed computing. The error
27 backpropagation (backprop) algorithm requires an alternation of inter-dependent forward and back-
28 ward phases, introducing a locking problem (the two phases have to wait for each other) [Jaderberg
29 et al., 2016a]. Furthermore, the two phases rely on the same weight matrices to calculate updates,
30 introducing a weight transport problem across blocks [Grossberg, 1987, Lillicrap et al., 2014a]. These
31 two issues make efficient parallelisation and horizontal scaling of large machine learning models
32 across compute nodes extremely difficult.

33 We propose a new method to address these problems by distributing a globally defined optimisation
34 algorithm across a large network of nodes that use only local learning. Our approach uses a message-
35 passing approach that uses results from probabilistic models and communicates uncertainty messages
36 forward and backwards between compute nodes in parallel. To do so, we augment a network

37 architecture with a twin network that propagates information backwards from the targets to the
38 input to provide uncertainty measures and auxiliary targets for local losses. Forward and backward
39 messages comprise information about extracted features and feature uncertainties and are matched
40 against each other using local probabilistic losses. Importantly, forward and backward propagation can
41 work in parallel, reducing the locking problem. Inside each block, conventional error backpropagation
42 is performed locally (“block-local”). These local updates can be used in the forward network and its
43 backward twin for adapting parameters during training. The developed theoretical learning provides a
44 new principled method to distribute very large networks over multiple compute nodes. The solutions
45 emerging from this framework show striking similarities to earlier models that used random feedback
46 weights as local targets [Lillicrap et al., 2020, Frenkel et al., 2021] but also provide a principled way
47 to train these feedback weights.

48 In summary, the contribution of this paper is threefold:

- 49 1. We provide a theoretical framework on how interpreting the representations of deep neural
50 networks as probability distributions provides a principled approach for block-local training
51 of these networks. This can be used to distribute learning and inference over many interacting
52 neural network blocks for various neural network architectures.
- 53 2. We demonstrate an instance of this probabilistic learning model on several benchmark
54 classification tasks, where classifiers are split into multiple blocks and trained without
55 end-to-end gradient computation.
- 56 3. We demonstrate how this framework can be used to allow deep networks to produce
57 uncertainty estimates over their predictions. This principle is showcased on an autoencoder
58 network that automatically predicts uncertainties alongside pixel intensity values after
59 training.

60 2 Related work

61 A number of methods for using local learning in DNNs had been introduced previously. Lomnitz et al.
62 [2022] introduced Target Projection Stochastic Gradient Descent (tpSGD), which uses layer-wise
63 SGD and local targets generated via random projections of the labels, but does not adapt the backward
64 weights. LocoProp [Amid et al., 2022] uses a layer-wise loss that consists of a target term and a
65 regularizer, which is used however to enable 2nd order learning and does not focus on distributing
66 the gradient optimization. Jimenez Rezende et al. [2016] used a generative model and a KL-loss for
67 local unsupervised learning of 3D structures.

68 Some previous methods are based on probabilistic or energy-based cost functions and use a contrastive
69 approach with positive and negative data samples. Contrastive learning Chen et al. [2020], Oord
70 et al. [2019] can be used to construct block-local losses Xiong et al. [2020], Illing et al. [2021].
71 Equilibrium propagation replaces target clamping with a target nudging phase [Scellier and Bengio,
72 2017]. Another interesting contrastive approach was recently introduced [Hinton, 2022, Ororbia and
73 Mali, 2023, Zhao et al., 2023]. However, it needs task-specific negative examples. [Han et al., 2018]
74 uses a local predictive loss to improve recurrent networks’ performance. In contrast to these methods,
75 our approach does not need separate positive and negative data samples and focuses on block-local
76 learning.

77 Feedback alignment [Lillicrap et al., 2020, Sanfiz and Akrouf, 2021] uses random projections to
78 propagate gradient information backwards. Jaderberg et al. [2016b] used pseudo-reward functions
79 which are optimized simultaneously by reinforcement learning to improve performance. Random
80 feedback alignment [Amid et al., 2022, Refinetti et al., 2021] and related approaches [Clark et al., 2021,
81 Nøklund, 2016, Launay et al., 2020], use fixed random feedback weights to back-propagate errors.
82 [Jaderberg et al., 2017] used decoupled synthetic gradients for local training. Target propagation
83 demonstrates non-trivial performance with random projections for target labels instead of errors
84 [Frenkel et al., 2021]. In contrast to these methods, we provide a principled way to adapt feedback
85 weights.

86 Other methods [Belilovsky et al., 2019, Löwe et al., 2019] used greedy local, block- or layer-wise
87 optimization. Notably, Nøklund and Eidnes [2019] achieved good results by combining a matching
88 and a local cross-entropy loss. [Siddiqui et al., 2023] recently used block-local learning based on a
89 cross-correlation metric over feature embeddings [Zbontar et al., 2021], demonstrating promising

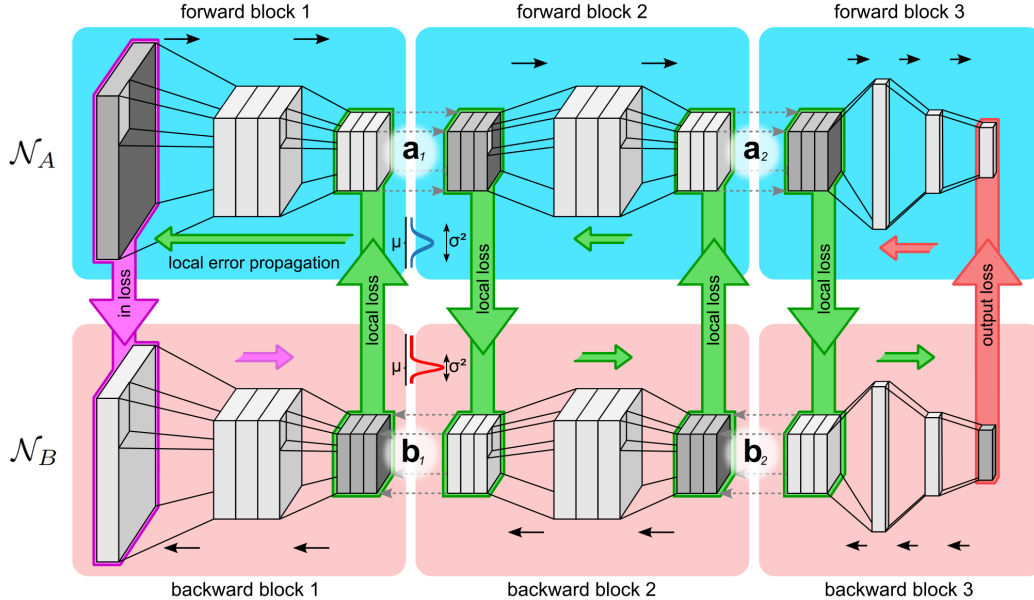


Figure 1: Illustration of use of block-local representations as learning signals on intermediate network layers. A deep neural network architecture \mathcal{N}_A is split into multiple blocks (forward blocks) and trained on an auxiliary local loss. Targets for local losses are provided by a twin backward network \mathcal{N}_B .

90 performance. [Wu et al., 2021] used greedy layer-wise optimization of hierarchical autoencoders for
 91 video prediction. [Wu et al., 2022] used an encoder-decoder stage for pretraining. In contrast to these
 92 methods, we do not rely solely on local greedy optimization but provide a principled way to combine
 93 local losses with feedback information without locking and weight transport across blocks.

94 3 A probabilistic formulation of distributed learning

95 At a high level, our method interprets the activations of a neural network as the parameters of
 96 probability distributions of latent variables. We use these intermediate representations at each block
 97 to derive block local losses. These latent variables over multiple blocks implicitly define a Markov
 98 chain, which allows us to tractably minimize the block’s local loss. We show that the derived block
 99 local losses and the resulting block local learning (BLL) are a general form of various existing local
 100 losses and provide an upper bound to a global loss.

101 3.1 Using latent representations to construct probabilistic block-local losses

102 Learning in deep neural networks can be formulated probabilistically [Ghahramani, 2015] in
 103 terms of maximum likelihood, i.e. the problem is to minimize the negative log-likelihood
 104 $\mathcal{L} = -\log p(\mathbf{x}, \mathbf{y}) = -\log p(\mathbf{y} | \mathbf{x}) - \log p(\mathbf{x})$ with respect to the network parameters θ .
 105 For many practical cases where we may not be interested in the prior distribution $p(\mathbf{x})$, we would
 106 like to directly minimize $\mathcal{L} = -\log p(\mathbf{y} | \mathbf{x})$.

107 This probabilistic interpretation of deep learning can be used to define block-local losses and distribute
 108 the learning over multiple blocks of networks by introducing intermediate latent representations. The
 109 idea is illustrated in Fig. 1. A neural network that computes the distribution $\log p(\mathbf{y} | \mathbf{x})$ takes \mathbf{x} as
 110 input and outputs the statistical parameters to the conditional distribution. The deep neural network
 111 is split at an intermediate layer k (in Fig. 1 we used $k \in (1, 2)$) and end-to-end estimation of the
 112 gradient is replaced by two estimators that optimize the sub-networks $\mathbf{x} \rightarrow \mathbf{z}_k$ and $\mathbf{z}_k \rightarrow \mathbf{y}$ separately.
 113 To do this, consider the gradient of the log-likelihood loss function

$$-\frac{\partial}{\partial \theta} \mathcal{L} = \frac{\partial}{\partial \theta} \log p(\mathbf{y} | \mathbf{x}) . \quad (1)$$

114 For any deep network, it is possible to choose any intermediate activation at layer k as latent
 115 representations \mathbf{z}_k , such that $\log p(\mathbf{y} | \mathbf{x}) = \left\langle p(\mathbf{y} | \mathbf{z}_k) p(\mathbf{z}_k | \mathbf{x}) \right\rangle_{p(\mathbf{z}_k | \mathbf{x}, \mathbf{y})}$, where $\left\langle \cdot \right\rangle_p$ denotes
 116 expectation with respect to p . Therefore, the representations of \mathbf{y} depend on \mathbf{x} only through \mathbf{z}_k as
 117 expected for a feed-forward network. Using this conditional independence property, the log-likelihood
 118 (1) expands to

$$-\frac{\partial}{\partial \theta} \mathcal{L} = \frac{\partial}{\partial \theta} \log p(\mathbf{y} | \mathbf{x}) = \left\langle \frac{\partial}{\partial \theta} \log p(\mathbf{y} | \mathbf{z}_k) + \frac{\partial}{\partial \theta} \log p(\mathbf{z}_k | \mathbf{x}) \right\rangle_{p(\mathbf{z}_k | \mathbf{x}, \mathbf{y})}. \quad (2)$$

119 This well-known result is the foundation of the Expectation-Maximization (EM) algorithm [Dempster
 120 et al., 1977]. Computing the marginal with respect to $p(\mathbf{z}_k | \mathbf{x}, \mathbf{y})$ corresponds to the E-step and
 121 calculating the gradients corresponds to the M-step. The sum inside the expectation separates the
 122 gradient estimators into two parts: $\mathbf{x} \rightarrow \mathbf{z}_k$ and $\mathbf{z}_k \rightarrow \mathbf{y}$.

123 However, the E-step is impractical to compute for most interesting applications because of the
 124 combinatorial explosion in the state space of \mathbf{z}_k . To get around this, we use a variational lower bound
 125 to EM, based on the ELBO loss $\mathcal{L}_V = -\log p(\mathbf{y} | \mathbf{x}) + \mathcal{D}_{KL}(q | p)$ [Mnih and Gregor, 2014] and
 126 demonstrate that this yields a practical solution to split gradients in a similar fashion to Eq. (2). In the
 127 next section, we describe how we construct the variational distribution q .

128 3.2 Auxiliary latent representations

129 As described earlier, the output of any layer of a DNN can be interpreted as parameters to a distribution
 130 over latent random variable \mathbf{z}_k . The sequence of blocks across a network therefore implicitly defines
 131 a Markov chain $\mathbf{x} \rightarrow \mathbf{z}_1 \rightarrow \mathbf{z}_2 \rightarrow \dots$ (see Fig. 2A). This probabilistic interpretation of hidden layer
 132 activity is valid under relatively mild assumptions, studied in more detail in the Supplement. It is
 133 important to note that the network at no point produces samples from the implicit random variables
 134 \mathbf{z}_k , but they are introduced here only to conceptualize the mathematical framework. Instead the
 135 network outputs the parameters to $\alpha_k(\mathbf{z}_k)$ which is the probability distribution over \mathbf{z}_k (e.g. means
 136 and variances if α_k is Gaussian). The network thus translates $\alpha_{k-1} \rightarrow \alpha_k \rightarrow \dots$ by outputting the
 137 statistical parameters of the conditional distribution $\alpha_k(z_k)$ and taking $\alpha_k(z_{k-1})$ parameters as input.
 138 More precisely, the network implicitly computes a marginal distribution

$$\alpha_k(\mathbf{z}_k) = p(\mathbf{z}_k | \mathbf{x}) = \left\langle p_k(\mathbf{z}_k | \mathbf{z}_{k-1}) \right\rangle_{p(\mathbf{z}_{k-1} | \mathbf{x})} = \left\langle p_k(\mathbf{z}_k | \mathbf{z}_{k-1}) \right\rangle_{\alpha_{k-1}(\mathbf{z}_{k-1})}, \quad (3)$$

139 where $\left\langle \cdot \right\rangle_p$ denotes expectation with respect to the probability distribution p . Consequently, the
 140 network realizes a conditional probability distribution $p(\mathbf{y} | \mathbf{x})$ (where \mathbf{x} and \mathbf{y} are network inputs
 141 and outputs, respectively). And by the universal approximator property of deep neural networks,
 142 an accurate representation of this distribution can be learnt in the network weights through error
 143 back-propagation (as demonstrated for the example in Fig. 2). Eq. (3) is an instance of the belief
 144 propagation algorithm to efficiently compute conditional probability distributions.

145 To construct the variational distribution q we introduce the backward network \mathcal{N}_B that propagates
 146 messages β_k backwards according to Eq. 4 (see Fig. 1 for an illustration). Inference over the posterior
 147 distribution $p(\mathbf{z}_k | \mathbf{x}, \mathbf{y})$ for any latent variable \mathbf{z}_k can be made using the belief propagation algorithm,
 148 propagating messages $\alpha_k(\mathbf{z}_k)$ forward through the network using Eq. (3). In addition messages
 149 $\beta_k(\mathbf{z}_k)$ need to be propagated backward according to

$$\beta_k(\mathbf{z}_k) = p(\mathbf{y} | \mathbf{z}_k) = \left\langle p(\mathbf{y} | \mathbf{z}_{k+1}) \right\rangle_{p_k(\mathbf{z}_{k+1} | \mathbf{z}_k)} = \left\langle \beta_{k+1}(\mathbf{z}_{k+1}) \right\rangle_{p_k(\mathbf{z}_{k+1} | \mathbf{z}_k)}, \quad (4)$$

150 such that the posterior $p(\mathbf{z}_k | \mathbf{x}, \mathbf{y})$ can be computed up to normalization

$$\rho_k(\mathbf{z}_k) = p(\mathbf{z}_k | \mathbf{x}, \mathbf{y}) \propto p(\mathbf{z}_k | \mathbf{x}) p(\mathbf{y} | \mathbf{z}_k) = \alpha_k(\mathbf{z}_k) \beta_k(\mathbf{z}_k). \quad (5)$$

151 We make use of the fact that, through Eq. (3), the parameters of a probability distribution $p(\mathbf{z}_k | \mathbf{x})$
 152 are a function of the parameters to $p(\mathbf{z}_i | \mathbf{x})$, for $0 < i < k$, e.g. if α is assumed to be Gaussian
 153 we have $(\mu(\alpha_k), \sigma^2(\alpha_k)) = f(\mu(\alpha_i), \sigma^2(\alpha_i))$, where $\mu(\cdot)$ and $\sigma^2(\cdot)$ are the mean and vari-
 154 ance of the distribution respectively. Thus, if a network outputs $(\mu(\alpha_i), \sigma^2(\alpha_i))$ on layer i and
 155 $(\mu(\alpha_k), \sigma^2(\alpha_k))$ on layer k , a suitable probabilistic loss function will allow the network to learn

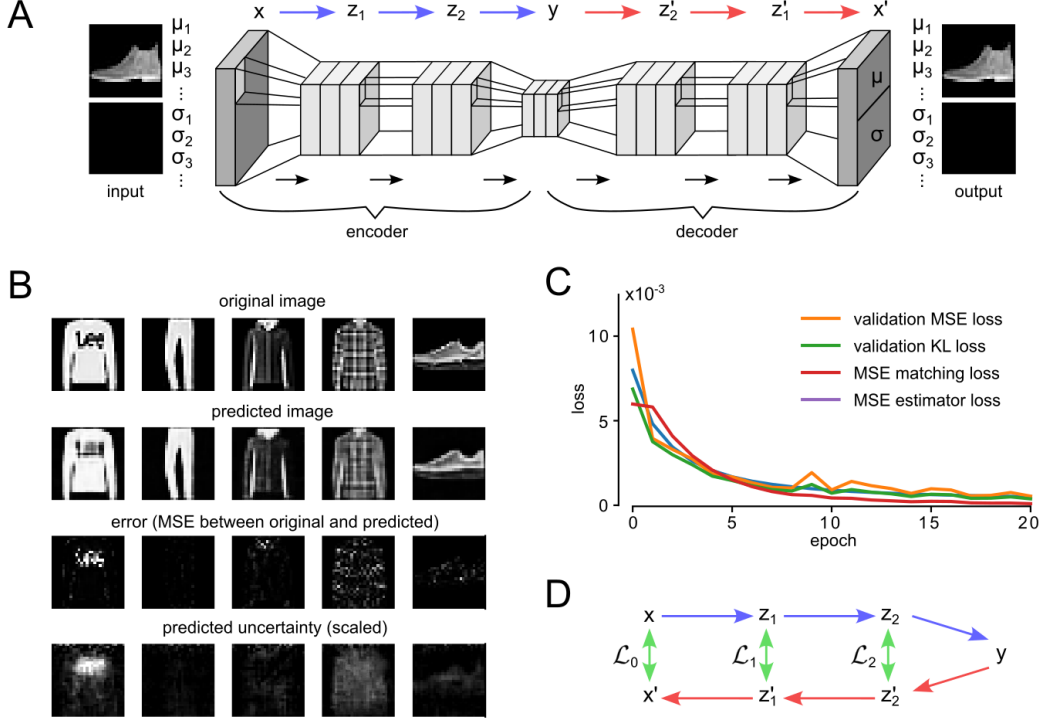


Figure 2: Zero shot learning of predicted uncertainties. **A:** Gaussian convolutional autoencoder network. Variance inputs and outputs are set to a constant during the whole training process. The network implements an implicit Markov chain. **B:** Example images showing self-prediction of uncertainties. **C:** Uncertainty mismatch metrics throughout learning. **D:** The network in (A) can be ‘folded’ to provide targets for local losses $\mathcal{L}_0, \mathcal{L}_1, \dots$

156 f from examples. Therefore, the conditional distributions $p_k(\mathbf{z}_k | \mathbf{z}_{k-1})$ and the expectation in
 157 Eq. (3) are only implicitly encoded in the network weights. We will study the exponential family of
 158 probability distributions for which this observation can be formalized more thoroughly.

159 **Exponential family distributions:** To derive concrete losses and update rules for the forward
 160 and backward networks, we assume that α_k are from the exponential family (EF) of probability
 161 distributions, given by

$$\alpha_k(\mathbf{z}_k) = \prod_j \alpha_{kj}(z_{kj}) = \prod_j h(z_{kj}) \exp(T(z_{kj}) \phi_{kj} - A(\phi_{kj})) , \quad (6)$$

162 with base measure h , sufficient statistics T , log-partition function A , and natural parameters ϕ_{kj} .
 163 This rich class contains the most common distributions, such as Gaussian, Poisson or Bernoulli, as
 164 special cases. For the example of a Bernoulli random variable we have $z_{kj} \in \{0, 1\}$, $T(z_{kj}) = z_{kj}$
 165 and $A(\phi_{kj}) = \log(1 + e^{\phi_{kj}})$ [Koller and Friedman, 2009]. A network directly implements an EF
 166 distribution if the activations a_{kj} encode the natural parameters, $a_{kj} = \phi_{kj}$. Using this result, a
 167 feed-forward DNN $\mathcal{N}_A: \mathbf{x} \rightarrow \mathbf{y}$, can be split into N blocks by introducing implicit latent variables
 168 $\mathbf{z}_k: \mathbf{x} \rightarrow \mathbf{z}_k \rightarrow \mathbf{y}$, and generating the respective natural parameters. In principle, blocks can be
 169 separated after any arbitrary layer, but some splits may turn out more natural for a particular network
 170 architecture.

171 Conveniently, if both α_{kj} and β_{kj} are members of the EF with natural parameters a_{kj} and b_{kj} , then
 172 ρ_{kj} is also EF with parameters $a_{kj} + b_{kj}$. We will use this property to deconstruct a single global
 173 loss into multiple block-local losses.

174 **3.3 Illustrative example: forward-backward networks as an autoencoder**

175 Probability representations in DNNs are useful since they provide a principled way to represent
 176 uncertainties in the network. Before we establish our main result to show how a DNN can be
 177 deconstructed into local blocks, we first demonstrate how representations of Bayesian uncertainty
 178 can emerge in DNNs by using appropriate probabilistic losses. We consider the autoencoder network
 179 illustrated in Fig. 2A and use it to learn representations for the Fashion-MNIST dataset [Xiao
 180 et al., 2017]. The CNN comprises a bottleneck layer y that implicitly splits the architecture into a
 181 decoder and encoder part (Fig. 2A). It is well known that such a network is able to learn compact
 182 representations and features that allow it to reconstruct the gray scale pixel intensities of a given
 183 input [Kingma and Welling, 2013]. Here we demonstrate that autoencoders are also able to learn
 184 representations of uncertainties, i.e. to automatically output high uncertainties for pixel values that
 185 are poorly represented in the learnt features.

186 To show this, we augmented the pixel representations on the inputs and outputs with additional
 187 channels that represented the logarithms of the variances of a Gaussian distribution (see Supplement
 188 for details). The input and outputs now represent the parameters of probability distributions, where
 189 the variances are proxies for the uncertainties. An appropriate loss function for this architecture is
 190 one that measures the distance between probability distributions. We used the Kullback-Leibler (KL)
 191 divergence between Gaussian distributions. This augmentation to conventional deep auto-encoders
 192 requires us to also provide uncertainty values for training data samples. Since the Fashion-MNIST
 193 dataset does not contain this information, we set the variances of pixels for all training samples to the
 194 same small constant values, reflecting high confidence (low variance) in the training set. Thus, during
 195 training, the network has only seen the same constant inputs (and outputs) for the variance channels.

196 Fig. 2B shows representative sample outputs for the test dataset after training. As expected, the
 197 network is able to represent the means of gray scale values in the dataset well and generalize to
 198 new images. Interestingly, the network also learned meaningful representations of the variances.
 199 Although the network has only seen constant values for the variances during training, it is able to
 200 infer information about its own uncertainty during testing. The true MSE errors between inputs and
 201 predictions qualitatively match the pixel-level variance predictions across a wide variety of inputs.
 202 For example, the network poorly represents the logo on the shirt (leftmost example) and predicts
 203 high variance in the output for these pixels. Other samples like the trousers (second from left) that
 204 are well represented correctly predict low variance. To further quantify this result, we developed
 205 additional metrics that measure the mismatch between estimated and true prediction errors (Fig. 2C,
 206 see Supplement for details). These metrics consistently decrease throughout training even though
 207 they were not directly minimized. These results suggest that DNNs are able to represent uncertainties
 208 well enough that they show zero-shot generalizations to unseen data from very limited training data.

209 **3.4 Modularized learning using local variational losses**

210 The autoencoder example described in Section 3.3 shows that DNNs can represent probability
 211 distributions well in principle, and also provides an idea of how probabilistic losses could be
 212 constructed locally at any layer. By ‘folding’ the network along the bottleneck layer y we are able
 213 to construct a sequence of pairs of auxiliary targets $(\mathbf{z}_1, \mathbf{z}'_1), (\mathbf{z}_2, \mathbf{z}'_2), \dots$ (see Fig. 2D). Finally, by
 214 introducing suitable loss functions $\mathcal{L}_0, \mathcal{L}_1, \dots$, the mismatch between the encoder and decoder parts
 215 of the network can be minimized on a per-layer basis.

216 The forward and backward networks \mathcal{N}_A and \mathcal{N}_B can be used to construct local loss functions $\mathcal{L}_V^{(k)}$ at
 217 blocks k . In the Supplement, we show in detail that minimizing $\mathcal{L}_V^{(k)}$ locally and in parallel optimizes
 218 a lower bound to the log-likelihood loss \mathcal{L} (Eq. 1), without propagating gradients end-to-end. To
 219 arrive at this result, we take the forward α_k and posterior messages ρ_k to be given by EF distributions
 220 with natural parameters ϕ_{kj} and γ_{kj} . Using this we show in the Supplement that the local loss can be
 221 optimized using the modularized gradient estimator

$$-\frac{\partial}{\partial \theta} \mathcal{L}_V^{(k)} = \sum_j \underbrace{(\mu(\rho_{kj}) - \mu(\alpha_{kj}))}_{\text{forward weight}} \frac{\partial}{\partial \theta} \phi_{kj} + \underbrace{\sigma^2(\rho_{kj}) (\phi_{kj} - \gamma_{kj})}_{\text{posterior weight}} \frac{\partial}{\partial \theta} \gamma_{kj}, \quad (7)$$

222 where $\mu(\cdot)$ and $\sigma^2(\cdot)$ are means and variances of EF distribution. Note that the gradients of the natural
 223 parameters ϕ_{kj} and γ_{kj} are computed independently and modulated by the *forward* and *posterior*
 224 *weight*, respectively.

225 The result in Eq. (7) holds for general EF distributions. For the special case of Bernoulli random
 226 variables we get

$$-\frac{\partial}{\partial \theta} \mathcal{L}_V^{(k)} = \sum_{k,j} (\rho_{kj} - \alpha_{kj}) \frac{\partial}{\partial \theta} a_{kj} - \rho_{kj} (1 - \rho_{kj}) b_{kj} \left(\frac{\partial}{\partial \theta} a_{kj} + \frac{\partial}{\partial \theta} b_{kj} \right), \quad (8)$$

227 where $a_{kj} = f_j(\mathbf{a}_{k-1})$ and $b_{kj} = g_j(\mathbf{b}_{k+1})$, are the outputs of the forward and backward network at
 228 block k ,

$$\rho_{kj} = S(a_{kj} + m b_{kj}) \quad \text{and} \quad \alpha_{kj} = S(a_{kj}), \quad (9)$$

229 where m is a mixing parameter described below and $S(x) = 1/(1 + e^{-x})$ is the sigmoid/logistic
 230 function.

231 The Bernoulli solution in Eq. (8) is convenient because it is a single parameter distribution (mean
 232 and variance share one parameter) such that all channels in \mathbf{z} can be treated independently. Also the
 233 structure of Eq. 9 is well suited for a DNN implementation. In our experiments, we focus on this
 234 Bernoulli variant of the general result in Eq. (7). In the Supplement, we study a number of other
 235 relevant members of the EF. Furthermore, it is interesting to study the structure of Eq. (8) more
 236 carefully. The first term minimizes the mismatch between the forward and the posterior distribution
 237 with respect to the forward blocks. The second term is the uncertainty-weighted backward activation
 238 b_{kj} which modulates local gradients (see Supplement). Therefore, the backward activations b_{kj}
 239 act directly as learning signals for local updates. The BLL method is therefore related to feedback
 240 alignment [Lillicrap et al., 2020] and target propagation [Frenkel et al., 2021] where backward
 241 information is provided through random weights. However, since the gradients of the backward
 242 blocks appear in the second term, our model also provides a principled way to optimize the backward
 243 flow of information from the targets.

244 **Data mixing schedule:** The equation for the posterior distribution Eq. 9 contains a data mixing
 245 parameter m , with $0 \leq m \leq 1$, that scales the influence of the backward messages in the posterior
 246 distribution. This parameter serves two important functions, (1) It scales the balance between forward
 247 and backward messages in the posterior distribution ρ and (2) it scales the first term in the parameter
 248 updates Eq. 8. We found that an annealing schedule for this parameter that decreases m slowly
 249 during learning works well in practice. If not stated otherwise, we used $m = (1 + \tau M)^{-1}$ in
 250 our experiments, where M is the index of the current epoch and τ is a scaling parameter (see the
 251 Supplement for further details).

252 4 Experimental results

253 We evaluated the BLL model on a number of vision and sequence learning tasks. All models used the
 254 Bernoulli BLL gradients described in Eq. (8) for local optimization. Additional details of the network
 255 models can be found in the Supplement.

256 4.1 Block-local learning of vision benchmark tasks

257 We compare the performance of our block local learning (BLL) algorithm with that of end-to-end
 258 backprop (BP) and Feedback Alignment (FA) Lillicrap et al. [2014b]. Three datasets are considered:
 259 MNIST, Fashion MNIST and CIFAR10 together with two residual network architectures [He et al.,
 260 2016]: ResNet-18 and ResNet-50, each trained with one of the three methods (BP, FA, BLL).

261 The BLL architectures were split into 4 blocks that were trained locally using the Bernoulli loss
 262 in Eq. (8). Splits were introduced after residual layers of the ResNet architecture by grouping
 263 subsequent layers into blocks. Group sizes were (4,5,4,5) for ResNet-18 and (12,13,12,13) for
 264 ResNet-50. Backward twin networks were here constructed simply by using the same network
 265 architecture (ResNet-18 or ResNet-50) in reverse order, introducing appropriate splits to provide
 266 intermediate targets. For CIFAR-10 gradients were propagated between two neighboring blocks
 267 (see Supplement for details and a comparison with purely local gradients). The kernels of ResNet-
 268 18/ResNet-50 + FA architectures used during backpropagation are fixed and uniformly initialised
 269 following the Kaiming He et al. [2015] initialisation method. The bias is set to one.

270 The results are summarized in Table 1. Test top-1, top-3 and train top-1 accuracies are shown. Top-3
 271 accuracies count the number of test samples for which the correct class was among the network's

	MNIST			Fashion-MNIST			CIFAR-10		
	test-1	test-3	train-1	test-1	test-3	train-1	test-1	test-3	train-1
ResNet-18 + BP	99.5	100	99.7	92.7	99.3	96.0	95.2	99.3	100
ResNet-50 + BP	99.5	99.9	100	89.0	98.9	92.7	94.0	99.2	99.8
ResNet-18 + FA	99.0	99.9	100	87.9	98.6	92.1	70.4	92.5	80.9
ResNet-50 + FA	98.9	99.9	100	83.1	97.9	83.7	70.3	92.0	79.3
ResNet-18 + BLL	99.4	100	99.6	91.2	98.8	91.0	72.2	93.0	98.8
ResNet-50 + BLL	99.4	99.8	99.2	88.7	99.0	85.9	73.4	92.7	99.7

Table 1: Classification accuracy (% correct) on vision tasks. BP: end-to-end backprop, FA: feedback alignment, BLL: block local learning. Test-1, test-3 and train-1 represent the top-1, top-3 test accuracy and top-1 training accuracy respectively.

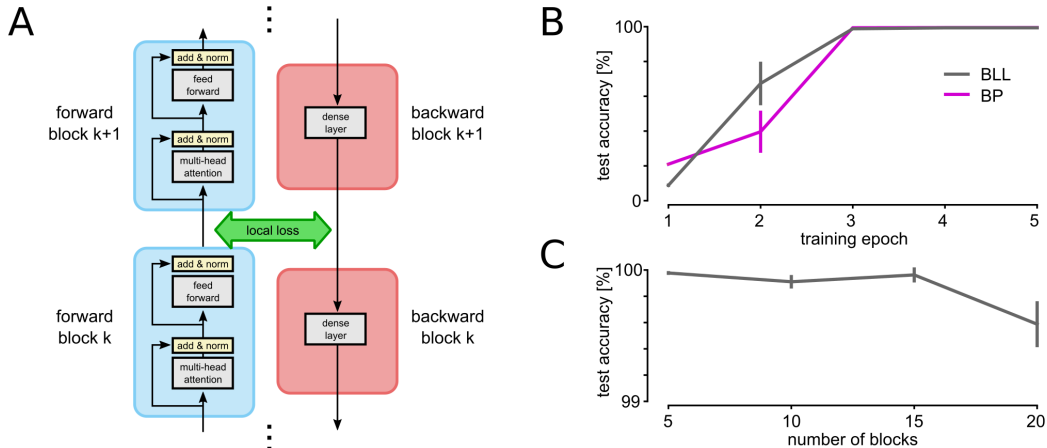


Figure 3: Block local learning of transformer architecture. **A:** Illustration of the transformer twin network. **B:** Learning curves of block local (BLL) and backprop (BP) training. **C:** Test accuracy vs. number of blocks in the transformer model. Error bars show standard deviations over 5 independent runs.

272 3 highest output activations. See Supplement for results over multiple runs. BLL achieved good
 273 performance on MNIST and Fashion-MNIST, closely matching end-to-end training and outperforming
 274 FA networks. Note that in contrast to FA and BP, BLL does not need to compute error gradients at
 275 the output but can work directly with the target labels. Performance on CIFAR-10 was significantly
 276 lower than BP but outperformed FA. Interestingly the performance on the training set was close to
 277 perfect for ResNet-50 suggesting over-fitting the task.

278 4.2 Block-local transformer architecture for sequence-to-sequence learning

279 Transformer architectures are in principle well suited for distributed computing due to their modular
 280 network structure that comprises a repetition of homogeneous blocks. We demonstrate a proof-of-
 281 concept result on training a transformer with BLL. We used a transformer model with 20 self-attention
 282 blocks with a single attention head each. Block local losses were added after each layer and blocks
 283 were trained locally. A backward twin network was constructed by projecting targets through dense
 284 layers and used the Bernoulli loss Eq. (8) for local training (see Fig. 3 A for an illustration). The
 285 transformer was trained on a sequence-to-sequence task, where a random permutation of numbers
 286 0..9 was presented on the input and had to be re-generated at the output in reverse order. We trained
 287 the network for 5 epochs.

288 BLL achieves convergence speed that is comparable to that of end-to-end BP on this task. Fig. 3 B
 289 shows learning curves of BLL and BP. Both algorithms converge after around 3 epochs to nearly
 290 perfect performance. BLL also achieved good performance for a wide range of network depths.
 291 Fig. 3 C shows the performance after 5 epochs for different transformer architectures. Using only 5
 292 transformer blocks yields performance of around 99.9% (average over five independent runs). The

293 test accuracy on this task for the 20 block transformer was 99.6%. These results suggest that the BLL
294 method is equally applicable to transformer architectures.

295 **5 Discussion**

296 In this work, we have demonstrated a general purpose probabilistic framework for rigorously defining
297 block-local losses for deep architectures. This not only provides a novel way of performing distributed
298 training of large models but also hints at new paradigms of self-supervised training that are biologically
299 plausible. We have also shown that our block-local training approach outperforms existing local
300 training approaches while still getting around the locking and weight transport problems. Our method
301 introduces a twin network that propagates information backwards from the targets to the input
302 and automatically estimates uncertainties on intermediate layers. This is achieved by representing
303 probability distributions in the network activations. The forward network and its backward twin can
304 work in parallel and with different sets of weights.

305 The proposed method may also help further blur the boundary between deep learning and probabilistic
306 models. A number of previous models have shown that DNNs are capable of representing probability
307 distribution [Abdar et al., 2021, Pawlowski et al., 2017, Tran et al., 2019, Malinin and Gales, 2019].
308 Unlike these previous methods, our method does not require Monte Carlo sampling or contrastive
309 training, but instead exploits the log-linear structure of exponential family distributions to efficiently
310 propagate uncertainty-aware messages through a network using a belief-propagation strategy. We
311 have demonstrated that implicit uncertainty messages can be learnt from sparse data and accurately
312 represent the network’s performance.

313 Greedy block-local learning has recently shown compelling performance on a number of tasks
314 [Nøkland and Eidnes, 2019, Siddiqui et al., 2023]. These methods use local losses with an information-
315 theoretic motivation but are agnostic to global back-propagating information. In future work, it may
316 be interesting to combine these approaches with the proposed model to get the best of both worlds.
317 Being able to produce block-level uncertainty predictions can also be useful for enhancing the sparsity
318 of the network and using optimal amount of compute for predictions. The uncertainty predictions
319 can also be used to handle missing labels, and for evaluating the model’s confidence about its
320 predictions. Since the framework is flexible enough to apply to self-supervised training, it can be
321 used on unlabelled and multi-modal datasets as well. Due to the local nature of the training process,
322 our method is particularly attractive for application on neuromorphic systems that co-locate memory
323 and compute and use orders of magnitude less energy if the computation is local.

324 This work addresses potential problems of modern ML: The estimation of uncertainties in neural
325 networks is an important open problem and understanding the underlying mechanisms better will
326 likely help to make ML models safer and more reliable. Also the main focus of this work, which is on
327 distributing large ML models over many compute nodes may make these model more energy efficient
328 in the future. The energy consumption and resulting carbon footprint of ML is a major concern and
329 the proposed model may provide a new direction to approach this problem. This method may enable
330 training of larger models which also come with associated risks in terms of biases and inappropriate
331 use in the real world. It is also not known what biases using this method itself and extensions with
332 sparsity may introduce in the models predictions.

333 **References**

- 334 Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad
335 Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A
336 review of uncertainty quantification in deep learning: Techniques, applications and challenges.
337 *Information Fusion*, 76:243–297, 2021.
- 338 Ehsan Amid, Rohan Anil, and Manfred Warmuth. LocoProp: Enhancing BackProp via local loss
339 optimization. In *Proceedings of The 25th International Conference on Artificial Intelligence and
340 Statistics*, pages 9626–9642. PMLR, 2022. URL [https://proceedings.mlr.press/v151/
341 amid22a.html](https://proceedings.mlr.press/v151/amid22a.html).
- 342 Eugene Belilovsky, Michael Eickenberg, and Edouard Oyallon. Greedy layerwise learning can scale
343 to ImageNet. In *Proceedings of the 36th International Conference on Machine Learning*, pages 583–
344 593. PMLR, 2019. URL <https://proceedings.mlr.press/v97/belilovsky19a.html>.

- 345 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
346 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
347 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
348 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
349 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
350 Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. *arXiv:2005.14165 [cs]*,
351 July 2020. URL <http://arxiv.org/abs/2005.14165>.
- 352 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
353 contrastive learning of visual representations. In *International conference on machine learning*,
354 pages 1597–1607. PMLR, 2020.
- 355 David Clark, L F Abbott, and Sueyeon Chung. Credit assignment through broadcasting a global error
356 vector. In *Advances in Neural Information Processing Systems*, volume 34, pages 10053–10066.
357 Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/
358 532b81fa223a1b1ec74139a5b8151d12-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/532b81fa223a1b1ec74139a5b8151d12-Abstract.html).
- 359 A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the *EM*
360 algorithm. 39(1):1–22, 1977. ISSN 00359246. doi: 10.1111/j.2517-6161.1977.tb01600.x. URL
361 <https://onlinelibrary.wiley.com/doi/10.1111/j.2517-6161.1977.tb01600.x>.
- 362 Charlotte Frenkel, Martin Lefebvre, and David Bol. Learning without feedback: Fixed random
363 learning signals allow for feedforward training of deep neural networks. 15, 2021. ISSN 1662-
364 453X. URL <https://www.frontiersin.org/articles/10.3389/fnins.2021.629892>.
- 365 Zoubin Ghahramani. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):
366 452–459, 2015.
- 367 Stephen Grossberg. Competitive learning: From interactive activation to adaptive resonance. *Cogni-
368 tive science*, 11(1):23–63, 1987.
- 369 Kuan Han, Haiguang Wen, Yizhen Zhang, Di Fu, Eugenio Culurciello, and Zhongming Liu. Deep
370 predictive coding network with local recurrent processing for object recognition, 2018. URL
371 <http://arxiv.org/abs/1805.07526>.
- 372 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing
373 human-level performance on imagenet classification, 2015.
- 374 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
375 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
376 pages 770–778, 2016.
- 377 Geoffrey Hinton. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint
378 arXiv:2212.13345*, 2022.
- 379 Bernd Illing, Jean Ventura, Guillaume Bellec, and Wulfram Gerstner. Local plasticity
380 rules can learn deep representations using self-supervised contrastive predictions. In *Ad-
381 vances in Neural Information Processing Systems*, volume 34, pages 30365–30379. Cur-
382 ran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/
383 feade1d2047977cd0cefdafc40175a99-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/feade1d2047977cd0cefdafc40175a99-Abstract.html).
- 384 Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, David
385 Silver, and Koray Kavukcuoglu. Decoupled Neural Interfaces using Synthetic Gradients.
386 *arXiv:1608.05343 [cs]*, August 2016a. URL <http://arxiv.org/abs/1608.05343>.
- 387 Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z. Leibo, David
388 Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks, 2016b.
389 URL <http://arxiv.org/abs/1611.05397>.
- 390 Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, David
391 Silver, and Koray Kavukcuoglu. Decoupled neural interfaces using synthetic gradients. In
392 *Proceedings of the 34th International Conference on Machine Learning*, pages 1627–1635. PMLR,
393 2017. URL <https://proceedings.mlr.press/v70/jaderberg17a.html>.

- 394 Danilo Jimenez Rezende, S. M. Ali Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and
395 Nicolas Heess. Unsupervised learning of 3d structure from images. In *Advances in Neural Informa-*
396 *tion Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2016/hash/1d94108e907bb8311d8802b48fd54b4a-Abstract.html)
397 [neurips.cc/paper/2016/hash/1d94108e907bb8311d8802b48fd54b4a-Abstract.html](https://proceedings.neurips.cc/paper/2016/hash/1d94108e907bb8311d8802b48fd54b4a-Abstract.html).
- 398 Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. *arXiv:1312.6114 [cs, stat]*,
399 December 2013. URL <http://arxiv.org/abs/1312.6114>.
- 400 Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT
401 press, 2009.
- 402 Julien Launay, Iacopo Poli, François Boniface, and Florent Krzakala. Direct feedback alignment scales
403 to modern deep learning tasks and architectures. In *Advances in Neural Information Processing Sys-*
404 *tems*, volume 33, pages 9346–9360. Curran Associates, Inc., 2020. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2020/hash/69d1fc78dbda242c43ad6590368912d4-Abstract.html)
405 [neurips.cc/paper/2020/hash/69d1fc78dbda242c43ad6590368912d4-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/69d1fc78dbda242c43ad6590368912d4-Abstract.html).
- 406 Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Random feedback
407 weights support learning in deep neural networks. *arXiv:1411.0247 [cs, q-bio]*, November 2014a.
408 URL <http://arxiv.org/abs/1411.0247>.
- 409 Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Random feedback
410 weights support learning in deep neural networks, 2014b.
- 411 Timothy P. Lillicrap, Adam Santoro, Luke Marris, Colin J. Akerman, and Geoffrey Hinton. Backprop-
412 agation and the brain. 21(6):335–346, 2020. ISSN 1471-0048. doi: 10.1038/s41583-020-0277-3.
413 URL <https://www.nature.com/articles/s41583-020-0277-3>.
- 414 Michael Lomnitz, Zachary Daniels, David Zhang, and Michael Piacentino. Learning with local
415 gradients at the edge, 2022. URL <http://arxiv.org/abs/2208.08503>.
- 416 Sindy Löwe, Peter O’ Connor, and Bastiaan Veeling. Putting an end to end-to-end: Gradient-isolated
417 learning of representations. In *Advances in Neural Information Processing Systems*, volume 32.
418 Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/851300ee84c2b80ed40f51ed26d866fc-Abstract.html)
419 [851300ee84c2b80ed40f51ed26d866fc-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/851300ee84c2b80ed40f51ed26d866fc-Abstract.html).
- 420 Andrey Malinin and Mark Gales. Reverse kl-divergence training of prior networks: Improved
421 uncertainty and adversarial robustness. *Advances in Neural Information Processing Systems*, 32,
422 2019.
- 423 Andriy Mnih and Karol Gregor. Neural variational inference and learning in belief networks. In
424 *International Conference on Machine Learning*, pages 1791–1799. PMLR, 2014.
- 425 Arild Nøkland and Lars Hiller Eidnes. Training neural networks with local error signals. In
426 *International conference on machine learning*, pages 4839–4850. PMLR, 2019.
- 427 Arild Nøkland. Direct feedback alignment provides learning in deep neural net-
428 works. In *Advances in Neural Information Processing Systems*, volume 29. Curran
429 Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper/2016/hash/](https://proceedings.neurips.cc/paper/2016/hash/d490d7b4576290fa60eb31b5fc917ad1-Abstract.html)
430 [d490d7b4576290fa60eb31b5fc917ad1-Abstract.html](https://proceedings.neurips.cc/paper/2016/hash/d490d7b4576290fa60eb31b5fc917ad1-Abstract.html).
- 431 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive
432 coding, 2019. URL <http://arxiv.org/abs/1807.03748>.
- 433 Alexander Ororbia and Ankur Mali. The predictive forward-forward algorithm. *arXiv preprint*
434 *arXiv:2301.01452*, 2023.
- 435 Nick Pawłowski, Andrew Brock, Matthew CH Lee, Martin Rajchl, and Ben Glocker. Implicit weight
436 uncertainty in neural networks. *arXiv preprint arXiv:1711.01297*, 2017.
- 437 Maria Refinetti, Stéphane d’Ascoli, Ruben Ohana, and Sebastian Goldt. Align, then memorise:
438 the dynamics of learning with feedback alignment, 2021. URL [http://arxiv.org/abs/2011.](http://arxiv.org/abs/2011.12428)
439 [12428](http://arxiv.org/abs/2011.12428).
- 440 Albert Jiménez Sanfiz and Mohamed Akrouf. Benchmarking the accuracy and robustness of feedback
441 alignment algorithms, 2021. URL <http://arxiv.org/abs/2108.13446>.

- 442 Benjamin Scellier and Yoshua Bengio. Equilibrium propagation: Bridging the gap between energy-
443 based models and backpropagation. *Frontiers in computational neuroscience*, 11:24, 2017.
- 444 Shoaib Ahmed Siddiqui, David Krueger, Yann LeCun, and Stéphane Deny. Blockwise self-supervised
445 learning at scale, 2023. URL <http://arxiv.org/abs/2302.01647>.
- 446 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
447 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
448 Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation Language
449 Models, February 2023. URL <http://arxiv.org/abs/2302.13971>.
- 450 Dustin Tran, Mike Dusenberry, Mark Van Der Wilk, and Danijar Hafner. Bayesian layers: A module
451 for neural network uncertainty. *Advances in neural information processing systems*, 32, 2019.
- 452 Bohan Wu, Suraj Nair, Roberto Martin-Martin, Li Fei-Fei, and Chelsea Finn. Greedy hierarchi-
453 cal variational autoencoders for large-scale video prediction. In *Proceedings of the IEEE/CVF*
454 *Conference on Computer Vision and Pattern Recognition*, pages 2318–2328, 2021.
- 455 Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan. TinyViT:
456 Fast pretraining distillation for small vision transformers, 2022. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2207.10666)
457 [2207.10666](http://arxiv.org/abs/2207.10666).
- 458 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
459 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 460 Yuwen Xiong, Mengye Ren, and Raquel Urtasun. LoCo: Local contrastive representation learning.
461 In *Advances in Neural Information Processing Systems*, volume 33, pages 11142–11153. Cur-
462 ran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/7fa215c9efebb3811a7ef58409907899-Abstract.html)
463 [7fa215c9efebb3811a7ef58409907899-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/7fa215c9efebb3811a7ef58409907899-Abstract.html).
- 464 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised
465 learning via redundancy reduction, 2021. URL <http://arxiv.org/abs/2103.03230>.
- 466 Gongpei Zhao, Tao Wang, Yidong Li, Yi Jin, Congyan Lang, and Haibin Ling. The cascaded forward
467 algorithm for neural network training. *arXiv preprint arXiv:2303.09728*, 2023.