

---

# Stacked unsupervised learning with a network architecture found by supervised meta-learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Stacked unsupervised learning (SUL) seems more biologically plausible than  
2 backpropagation, because learning is local to each layer. But SUL has fallen far  
3 short of backpropagation in practical applications, undermining the idea that SUL  
4 can explain how brains learn. Here we show an SUL algorithm that can perform  
5 completely unsupervised clustering of MNIST digits with comparable accuracy  
6 relative to unsupervised algorithms based on backpropagation. Our algorithm is  
7 exceeded only by self-supervised methods requiring training data augmentation by  
8 geometric distortions. The only prior knowledge in our unsupervised algorithm is  
9 implicit in the network architecture. Multiple convolutional “energy layers” contain  
10 a sum-of-squares nonlinearity, inspired by “energy models” of primary visual  
11 cortex. Convolutional kernels are learned with a fast minibatch implementation of  
12 the K-Subspaces algorithm. High accuracy requires preprocessing with an initial  
13 whitening layer, representations that are less sparse during inference than learning,  
14 and rescaling for gain control. The hyperparameters of the network architecture  
15 are found by supervised meta-learning, which optimizes unsupervised clustering  
16 accuracy. We regard such dependence of unsupervised learning on prior knowledge  
17 implicit in network architecture as biologically plausible, and analogous to the  
18 dependence of brain architecture on evolutionary history.

## 19 1 Introduction

20 Recently there has been renewed interest in the hypothesis that the brain learns through some version  
21 of the backpropagation algorithm [31]. This hypothesis runs counter to the neuroscience textbook  
22 account that local learning mechanisms, such as Hebbian synaptic plasticity, are the basis for learning  
23 by real brains. The concept of local learning has fallen out of favor because it has been far eclipsed  
24 by backpropagation in practical applications. This was not always the case. Historically, a popular  
25 approach to visual object recognition was to repeatedly stack a single-layer unsupervised learning  
26 module to generate a multilayer network, as exemplified by Fukushima’s pioneering Neocognitron  
27 [11]. Stacked unsupervised learning (SUL) avoids the need for the backward pass of backpropagation,  
28 because learning is local to each layer.

29 In the 2000s, SUL was quite popular. There were attempts to stack diverse kinds of unsupervised  
30 learning modules, such as sparse coding [21], restricted Boltzmann machines [14, 29], denoising  
31 autoencoders [46], K-Means [7], and independent subspace analysis [25].

32 SUL managed to generate impressive-looking feature hierarchies that are reminiscent of the hierarchy  
33 of visual cortical areas. Stacking restricted Boltzmann machines yielded features that were sensitive  
34 to oriented edges in the first layer, eyes and noses in the second, and entire faces in the third layer  
35 [29]. Stacking three sparse coding layers yielded an intuitive feature hierarchy where higher layers  
36 were more selective to whole MNIST digits and lower layers were selective to small strokes [43].  
37 Although these feature hierarchies are pleasing to the eye, they have not been shown to be effective

38 for visual object recognition, in spite of recent efforts to revive SUL using sparse coding [43, 6] and  
39 similarity matching [39].

40 Here we show an SUL algorithm that can perform unsupervised clustering of MNIST digits with high  
41 accuracy (2% error). The clustering accuracy is as good as the best unsupervised learning algorithms  
42 based on backpropagation. As far as we know, our accuracy is only exceeded by self-supervised  
43 methods that require training data augmentation or architectures with hand-designed geometric  
44 transformations. Such methods use explicit prior knowledge in the form of geometric distortions to  
45 aid learning.

46 Our network contains three convolutional energy layers inspired by energy models of primary visual  
47 cortex [1], which contain a sum-of-squares nonlinearity. Our energy layer was previously used by  
48 [16] in their independent subspace analysis (ISA) algorithm for learning complex cells. The kernels  
49 of our convolutional energy layers are trained by K-Subspaces clustering [45] rather than ISA. We  
50 also provide a novel minibatch algorithm for K-Subspaces learning.

51 After training, the first energy layer contains neurons that are selective for simple features but invariant  
52 to local distortions. These are analogous to the complex cells in energy models of the primary visual  
53 cortex [1]. The invariances are learned here rather than hand-designed, similar to previous work  
54 [16, 15]. We go further by stacking multiple energy layers. The second and third energy layers learn  
55 more sophisticated kinds of invariant feature selectivity. As mentioned above, feature hierarchies  
56 have previously been demonstrated for SUL. The novelty here is the learning of a feature hierarchy  
57 that is shown to be useful for pattern recognition.

58 In the special case that the sum-of-squares contains a single term, or equivalently the subspaces are  
59 restricted to be rank one, our convolutional energy layer reduces to a conventional convolutional layer.  
60 Accuracy worsens considerably, consistent with the idea that the energy layers are important for  
61 learning invariances. The energy layers contain an adaptive thresholding that allows representations  
62 to be less sparse for inference than for learning. This is also shown to be important for attaining high  
63 accuracy, as has been reported for other SUL algorithms [8, 24]. Representations are rescaled for  
64 gain control, and the energy layers are preceded by a convolutional whitening layer. These aspects of  
65 the network are also important for high accuracy.

66 The detailed architecture of our unsupervised network depends on subspace number and rank, kernel  
67 size, and sparsity. We use automatic tuning software [2] to systematically search for a hyperparameter  
68 configuration that optimizes the clustering accuracy of the unsupervised network. Evaluating the  
69 clustering accuracy requires labeled examples, so the meta-learning is supervised, while the learning  
70 is unsupervised. A conceptually similar meta-learning approach has previously been applied to search  
71 for biologically plausible unsupervised learning algorithms [37].

72 For each iteration of the meta-learning, the weights of our network are initialized randomly before  
73 running the SUL algorithm. Therefore the only prior knowledge available to SUL resides in the  
74 architectural hyperparameters; no weights are retained from previous networks. We regard this  
75 implicit encoding of prior knowledge in network architecture as biologically plausible, because brain  
76 architecture also contains prior knowledge gained during evolutionary history, and meta-learning is  
77 analogous to biological evolution. In its own “lifetime” our network is able to learn with no labels at  
78 all. This is possible because the network is “born” with an architecture inherited from networks that  
79 “lived” previously.

## 80 2 Related work

81 **Independent subspace analysis** Our method is related to past works on independent subspace  
82 analysis [16, 17, 26], mixtures of principal components analyzers [13], and subspace clustering  
83 [45, 47]. A core idea behind these works is that invariances can be represented via subspaces. The  
84 most similar of these works to ours is [26] who stacked 2 layers of subspace features learned with  
85 independent subspace analysis to action recognition datasets.

86 **K-Means based features** Mathematically the work of [7, 9] is similar to ours. They use a variant of  
87 K-Means to learn patch features. Their learning algorithm (Algorithm 1) is a special case of our alg. 1  
88 where they use 1D subspaces and full batch updates. Their inference procedure is also very similar,  
89 in that they use a dynamic threshold to sparsify patch vector representations. They use spatial pooling  
90 layers for invariance, whereas our pooling is learned by the energy layers. Their primary mode of

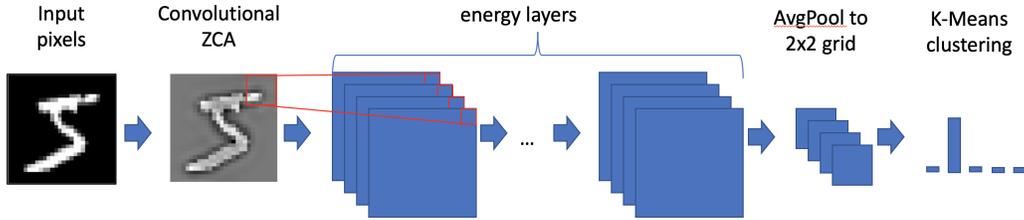


Figure 1: Our multilayer convolutional energy network. The last step of K-Means clustering can be regarded as part of the evaluation rather than the network itself.

91 evaluation was linear evaluation with the full labeled dataset, whereas we will seek to produce an  
 92 unsupervised learning algorithm which clusters inputs without labels. They additionally employ a  
 93 receptive field selection method.

94 **Capsule networks** Capsules are vector representations where the vector’s length represents the  
 95 probability of entity existence and the direction represents properties of that entity [41]. In our  
 96 networks, the  $r$ -dimensional subspace vectors  $\mathbf{V}_x$  can be interpreted as learned “pose” vectors.  
 97 However our goal is to learn invariance, so we only propagate the norm  $\|\mathbf{V}_x\|$ , thus suppressing the  
 98 pose details at each layer. Compared to the unsupervised capsule networks [23], our networks do not  
 99 use any backpropagation of gradients, and do not rely on hand-designed affine transformations to  
 100 generate representations.

101 **Meta-learning of unsupervised learning rules** Our work will rely on using a label-based clustering  
 102 objective to evaluate and tune an unsupervised learning rule. In other words there are two levels of  
 103 learning, an unsupervised inner loop and a supervised outer loop. This is the domain of meta-learning.  
 104 The more common scenario for meta-learning is to focus on optimizing over a distribution of tasks,  
 105 but for this work we will focus on one task. The work of [36] is perhaps the most closely related to  
 106 ours. They use a few-shot supervised learning rule to tune an unsupervised learning algorithm that is  
 107 a form of randomized backward propagation [30]. We wish to go further and remove any form of  
 108 gradient feedback from a higher layer  $L$  to a lower layer  $L - 1$ .

### 109 3 Network architecture

110 The overall network architecture is shown in Figure 1, and consists of a whitening layer followed by  
 111 multiple convolutional energy layers and a final average pooling layer. The energy layers include  
 112 adaptive thresholding to control sparsity of activity, as well as normalization of activity by rescaling.

113 **Convolutional ZCA** We define ZCA whitening for image patches in terms of the eigenvalues of  
 114 the pixel-pixel correlation matrix. Our ZCA filter attenuates the top  $k - 1$  eigenvalues, setting them  
 115 equal to the  $k$ th largest eigenvalue. The smaller eigenvalues pass through unchanged. Our definition  
 116 is slightly different from [10], which zeros out the smaller eigenvalues completely.

117 The first layer of our network is a convolutional variant of ZCA, in which each pixel of the output  
 118 image is computed by applying ZCA whitening to the input patch centered at that location, and  
 119 discarding all but the center pixel of the whitened output patch (p. 118 of [10]). The kernel has a single-  
 120 pixel center with a diffuse surround (see Appendix). Patch size and number of whitened eigenvalues  
 121 are specified in Table 1. Reflection padding is used to preserve the output size. Unsupervised learning  
 122 algorithms such as sparse coding [40] and ICA [5] are often preceded by whitening when applied to  
 123 images.

124 **Convolutional energy layer** We define the following modification of a convolutional layer, which  
 125 computes  $k$  output feature maps given  $m$  input feature maps. We define the “feature vector” at a  
 126 location to consist of the values of a set of feature maps at a given location.

- 127 1. Convolve the  $m$  input feature maps with kernels to produce  $kr$  feature maps (“S-maps”) in  
 128 the standard way, except with no bias or threshold.
- 129 2. Divide the  $kr$  feature maps into  $k$  groups of  $r$ . For each group, compute the Euclidean  
 130 norm of the  $r$ -dimensional feature vector at each location. The result is  $k$  feature maps  
 131 (“C-maps”).

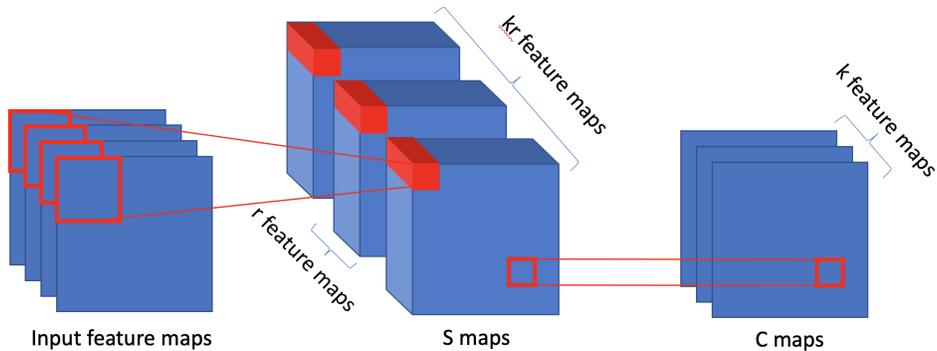


Figure 2: Diagram of a single energy layer. The outputs of this layer are the C feature maps, while S are just intermediate feature maps. We produce  $kr$  feature maps (S) with a standard convolution. We produce (C) feature maps by square root the sum of squares inside each group of S feature maps, followed by an adaptive thresholding and rescaling.

132 Individual values in the S- and C-maps of Steps 1 and 2 will be called S-cells and C-cells, respectively.  
 133 Each C-cell computes the square root of the sum-of-squares (Euclidean norm) of  $r$  S-cells. The terms  
 134 S-cells and C-cells are used in homage to [11]. They are reminiscent of energy models of primary  
 135 visual cortex, in which complex cells compute the sum-of-squares of simple cell outputs [1].

136 The sum-of-squares can be regarded as a kind of pooling operation, but applied to S-cells at the same  
 137 location in different feature maps. Pooling is typically performed on neighboring locations in the  
 138 same feature map, yielding invariance to small translations by design. We will see later on that the  
 139 sum-of-squares in the energy layer acquires invariances due to the learned kernels.

140 **Adaptive thresholding and normalization** It turns out to be important to postprocess the output of a  
 141 convolutional energy layer as follows.

- 142 3. For each  $k$ -dimensional feature vector, adaptively threshold so that there are  $W$  winners  
 143 active.
- 144 4. Normalize the feature vector to unit Euclidean length, and rescale by multiplying with the  
 145 Euclidean norm of the input patch at the same location.

146 Define  $\mathbf{f}$  as the  $k$ -dimensional vector which is the  $k$  feature map values at a location  $u$  in the  $C$  feature  
 147 maps. In Step 3, the adaptive thresholding of  $\mathbf{f}$  takes the form  $\max\{0, f_i - \tau\}$  for  $i = 1$  to  $k$  where  $\tau$   
 148 is the  $W + 1$ st largest element of  $f$ . Note that  $\tau$  is set adaptively for each location. Such adaptive  
 149 thresholding was used by [44] and is a version of the well-known  $W$ -winners-take-all concept [33].  
 150 After thresholding, C-cells are sparsely active, with sparseness controlled by the hyperparameter  $W$ .  
 151 S-cells, on the other hand, will typically be densely active, since they are linear.

152 Step 4 normalizes the  $k$ -dimensional feature vector, and also multiplies by the Euclidean norm of the  
 153 input patch, to prevent the normalized output from being large even if the input is vanishingly small.  
 154 The kernels in the layer are size  $p \times p$ , so that the input patch at any location contains  $mp^2$  values  
 155 where  $m$  is the number of input feature maps.

156 **Final average pooling layer** To reduce the output dimensionality, we average pool each output  
 157 feature map of the last energy layer to a  $2 \times 2$  grid, exactly as in [8]. We have generally avoided  
 158 pooling because we want to learn invariances rather than design them in. However, a final average  
 159 pooling will turn out to be advantageous later on for speeding up meta-learning. In Table 2 we show  
 160 that this pooling layer has only a modest impact on the final clustering accuracy.

## 161 4 Stacked unsupervised learning

162 The outputs of the ZCA layer are used as inputs to a convolutional energy layer. We train the kernels  
 163 of this layer, and then we freeze the kernels. The outputs of the first convolutional energy layer are  
 164 used as inputs to a second convolutional energy layer, and the kernels in this layer are trained. We  
 165 repeat this procedure to stack a total of three convolutional energy layers, and then conclude with a  
 166 final pooling layer.

LAYER	# SUBSPACES ( $k$ )	SUBSPACE RANK ( $r$ )	# WINNERS ( $w$ )	KERNEL SIZE	PADDING
L1	37	2	9	8	2
L2	9	3	8	5	1
L3	58	16	2	21	2

Table 1: Detailed architecture of our three energy layer net. The first energy layer is preceded by convolutional ZCA whitening of the input image, where the kernel size is 9 and the number of whitened eigenvalues is 9. These parameters are found with automated hyperparameter tuning.

REPRESENTATION	PIXELS	ZCA	LAYER 1	LAYER 2	LAYER 3	2X2 POOL
CLUSTERING ERROR (%)	46.2	50.0	22.7	22.2	2.3	2.1

Table 2: Clustering error after every layer of our network with three energy layers.

167 **K-Subspaces clustering** Consider an energy layer with  $k$  C-cells and  $kr$  S-cells at each location. The  
168 S-cells at one location are linearly related to the input patch at that location by the set of  $k$  matrices  
169  $\mathbf{V}_j$  for  $j = 1$  to  $k$ , each of size  $r \times mp^2$ . Here  $mp^2$  is the size of the flattened input patch, where  $m$   
170 is the number of input feature maps and  $p$  is the kernel size.

171 We can think of these matrices as defining a set of  $k$  linear subspaces of rank  $r$  embedded in  $\mathbb{R}^{mp^2}$ .  
172 We learn these matrices with a convolutional extension of the K-Subspace learning algorithm [45].  
173 Let  $\mathbf{x}_n$  be the previous layer’s  $m$ -dimensional feature maps for pattern  $n$ . Define  $\mathbf{x}_{n,i}$  as the  $mp^2$ -  
174 dimensional feature vector created by flattening a  $p \times p$  patch centered around location  $i$ . K-Subspaces  
175 aims to learn  $k$   $r$ -dimensional subspaces  $\mathbf{V}_k \in \mathbb{R}^{r \times mp^2}$  such that every patch is well modeled by  
176 one of these subspaces. This is formalized with the following optimization:

$$\min_{\mathbf{V}} \min_{\mathbf{C}} \sum_{n,i} \sum_k c_{nik} \|\mathbf{x}_{ni} - \mathbf{V}_k^\top \mathbf{V}_k \mathbf{x}_{ni}\|^2 \quad (1)$$

177 such that  $c_{nik} \in \{0, 1\}$  and  $\sum_k c_{nik} = 1$ . We provide a novel minibatch algorithm for this opti-  
178 mization in Appendix A. During the learning, each matrix  $\mathbf{V}_j$  is constrained so that its rows are  
179 orthonormal. Therefore at each location the C-cells contain the Euclidean norms of the projection of  
180 the input patch onto each of the subspaces, before the adaptive thresholding and rescaling. This is  
181 why the K-Subspaces algorithm is naturally well-suited for learning the convolutional energy layer.

182 During K-Subspaces learning, an input patch is assigned to a single subspace, which means there  
183 is winner-take-all competition between C-cells at a given location. During inference, on the other  
184 hand, there can be many C-cells active at a given location (depending on the "number of winners"  
185 hyperparameter  $w$ ).

186 This idea of making representations less sparse for inference than for learning has been exploited by  
187 a number of authors [7, 24]. This may be advantageous because overly sparse representations can  
188 suffer from sensitivity to distortions [32].

189 **Experiments with MNIST digits** We train a network with the architecture of Fig. 1. The details  
190 of the ZCA layer and the three convolutional energy layers are specified by Table 1. Each energy  
191 layer is trained using a single pass through the 60,000 MNIST [27] training examples (without using  
192 the labels), with a minibatch size of 512. Training on a single Nvidia Titan 1080-Ti GPU takes 110  
193 seconds.

194 **Evaluation of learned representations** We adopt the intuitive notion that the output representation  
195 vectors of a “good” network should easily cluster into the underlying object classes as defined by  
196 image labels. This is quantified by applying the trained network to 10,000 MNIST test examples.  
197 The resulting output representation vectors are clustered with the *scikit-learn* implementation of  
198 K-Means. This uses full batch EM-style updates and additionally returns the lowest mean squared  
199 error clustering found by running the algorithm using 10 different initializations.

200 We set the number of clusters to be 10, to match the number of MNIST digit classes. We compute the  
201 disagreement between the cluster assignments and image labels, and minimize over permutations of  
202 the clusters. The minimal disagreement is the final evaluation, which we will call the clustering error  
203 [49].

REPRESENTATION	CLUSTERING ERROR (%)
PIXELS	46.2
NMF <sup>‡</sup> [28]	44.0
STACKED DENOISING AUTOENCODERS <sup>†</sup> [46]	18.8
UMAP [35]	17.9
GENERATIVE ADVERSARIAL NETWORKS <sup>†</sup> [12]	17.2
VARIATIONAL AUTOENCODER <sup>†</sup> [22]	16.8
DEEP EMBEDDED CLUSTERING <sup>†</sup> [48]	15.7
VaDE [19]	5.5
CLUSTERGAN <sup>‡</sup> [38]	5.0
UMAP + GMM [35]	3.6
N2D [34]	2.1
<b>OURS (THREE LAYER NET)</b>	<b>2.1</b>
INVARIANT INFORMATION CLUSTERING (AVG SUB-HEAD) <sup>†</sup> [18]	1.6
STACKED CAPSULE AUTOENCODERS [23]	1.3
INVARIANT INFORMATION CLUSTERING (BEST SUB-HEAD) <sup>†</sup> [18]	0.8

Table 3: Comparison of clustering accuracy for other unsupervised learning algorithms. For methods which do not generate clusters,  $k$ -means is used to cluster representations into 10 clusters, with the exception of UMAP + GMM in which case we use a Gaussian Mixture Model to cluster. The errors for methods with a dagger <sup>†</sup> are all taken from [18], methods with double dagger <sup>‡</sup> are taken from [38]. UMAP uses “out-of-the-box” parameter settings.

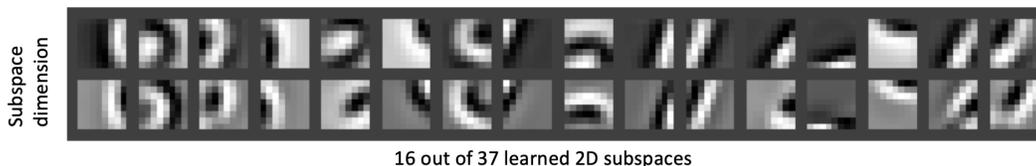


Figure 3: Layer 1 subspaces learned with our algorithm. Each  $8 \times 8$  image corresponds to a kernel.

204 Table 2 quantifies the accuracy of each layer. The accuracy of the ZCA representation is actually  
205 worse than that of the raw pixels. However, the accuracy of the representations improves with each  
206 additional convolutional energy layer, until the final error is just 2.1%.

207 Comparisons with other algorithms are shown in Table 3. It is helpful to distinguish between algo-  
208 rithms that require training data augmentation, and those that do not. Many well-known unsupervised  
209 algorithms that do not make use of training data augmentation, such as GANs, variational autoen-  
210 coders, and stacked denoising autoencoders, yield clustering errors of 15 to 19%. Methods such as  
211 VaDE and ClusterGAN encourage clustered latent representations and these give rise to much lower  
212 clustering error. Clustering 2D UMAP representations with K-Means gives suprisingly high error,  
213 and this is likely the returned clusters are not spherical. Using a Gaussian Mixture Model instead  
214 gives much lower error. See the appendix for more discussion.

215 Self-supervised algorithms require training data augmentation, using prior knowledge to create same-  
216 class image pairs. For MNIST clustering one of the highest performance is Invariant Information  
217 Clustering, with a clustering error of 1.6 - 0.8% [18]. Our approach delivers roughly 2 % error and is  
218 noticeably better than the other algorithms that do not require training data augmentation.

219 Stacked capsule autoencoders [23] also achieve high accuracy. However, this algorithm incorporates  
220 a model of geometric distortions. Furthermore, despite the modifier “stacked,” the algorithm does not  
221 conform to the original idea of repeatedly stacking the same learning module. The architecture uses  
222 several distinct types of layers and still backpropagates gradients.

223 **Learned kernels** Figure 3 shows that the kernels in the first energy layer look like bars or edges,  
224 more often curved than straight. Since the subspaces are of rank 2 (Table 1), the kernels come in  
225 pairs. The kernels in a pair look quite similar to each other, and typically appear to be related by  
226 small distortions. Therefore the two S-cells in a pair should prefer slightly different versions of the  
227 same feature. By computing the square root of the sum-of-squares of the S-cells, the C-cell should  
228 detect the same feature with more invariance to distortions.

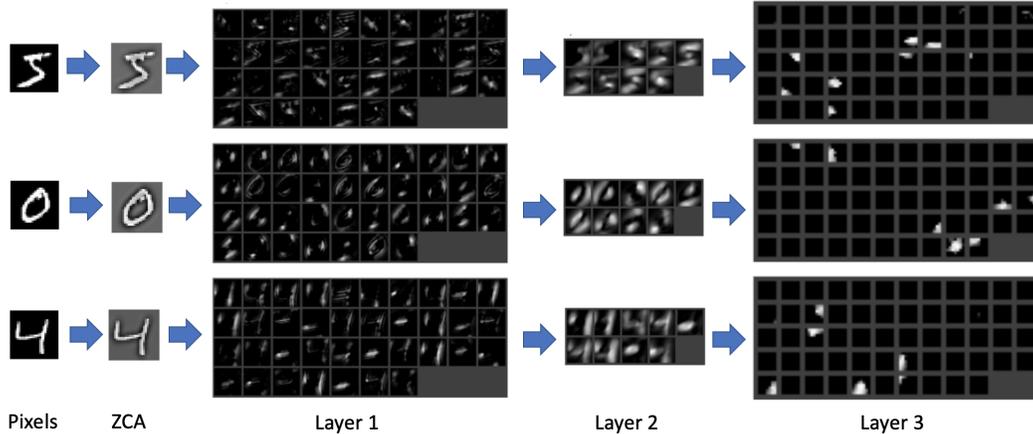


Figure 4: Output "C" feature maps for 3 input images.

229 This behavior is reminiscent of energy models of simple and complex cells in primary visual cortex  
 230 [1, 42]. A complex cell computes the sum-of-squares of simple cells, which are quadrature pairs of  
 231 Gabor filters. The sum-of-squares is invariant to phase shifts, much as the sum-of-squares of sine  
 232 and cosine is constant. Our model also contains a sum-of-squares, but the filters are learned rather  
 233 than hand-designed. Learning of complex cells was previously demonstrated by [16] for an energy  
 234 model and by [20] for a similar model. [15] showed how to substitute rectification for quadratic  
 235 nonlinearity. Our contribution is to stack multiple energy layers, and investigate the accuracy of the  
 236 resulting network at a clustering task.

237 **Feature maps** Figure 4 shows the C-maps for the first 3 MNIST digits. The first energy layer exhibits  
 238 an intermediate degree of sparsity, (approximately 20 % of the maps are active at central locations).  
 239 The second energy layer exhibits dense representations (nearly all maps are active at central locations).  
 240 The final energy layer shows quite sparse representations (approximately 5% of the maps are active).  
 241 The feature maps appear to be broad spots.

## 242 5 Meta-learning of network architecture for unsupervised learning

243 The detailed architecture of the network is specified in Table 1, and one might ask where it came from.  
 244 For example, the kernel size is 21 in the third energy layers, but sizes are only single digit integers in  
 245 other layers. The subspaces are rank 2 in the first energy layer, but rank 3 and 16 in the subsequent  
 246 energy layers. The second energy layer is highly dense (all but one neuron is active at each location),  
 247 while the third layer is highly sparse (only two neurons are active at each location). There is a total  
 248 of 17 numbers in Table 1, and we can regard them as hyperparameters of the unsupervised learning  
 249 algorithm.

250 In the initial stages of our research, we set hyperparameters by hand, guided by intuitive criteria such  
 251 as increasing the subspace rank with layer (meaning more invariances learned). Later on, we resorted  
 252 to automated hyperparameter tuning. For this purpose, we employed the Optuna software package,  
 253 which implements a Bayesian method [2]. We found that we were able to find hyperparameter  
 254 configurations with considerably better performance than our hand-designed configurations.

255 In particular, the configuration of Table 1 was obtained by automated search through 2000 hyper-  
 256 parameter configurations, which took approximately 20 hours using 8 Nvidia GTX 1080 Ti GPUS.  
 257 Each hyperparameter configuration was used to generate an unsupervised clustering of the training  
 258 set, and its accuracy with respect to all 60,000 training labels was the objective function of the search.

259 For the ZCA layer we tune the kernel size  $k_{zca} \in [1, 11]$  and number of whitened eigenvalues  
 260  $n_{zca} \in [0, k_{zca}^2]$ . For each subspace layer we tune number of subspaces  $k \in [2, 64]$ , subspace  
 261 dimension  $r \in [1, 16]$ , number of winners  $w \in [1, k]$ , kernel size  $k_s \in [1, \text{input\_size}]$ , and padding  
 262  $p_s \in [0, \text{floor}(k_s/2)]$ . The stride is fixed at 1.

263 In some respects, the optimized configuration of Table 1 ended up conforming to our qualitative  
 264 expectations. Sparsity and kernel size increased with layer, consistent with the idea of a feature

NUMBER OF LABELS IN TRAINING SET	10	30	50	100	500	5000	60000
% MISCLASSIFIED (TEST)	30.0	9.2	10.7	5.3	4.6	2.9	2.1
% MISCLASSIFIED (TRAIN)	0.0	0.0	0.0	3.0	2.8	2.9	2.5

Table 4: Label efficiency of our stacked learning algorithm. In each case, the algorithm has access to all 60K unlabeled images from the training set. What varies is the number of labels we use to evaluate each setting of learning parameters.

EXPERIMENT	% ERROR (TEST)	% ERROR (TRAIN)
4 ENERGY LAYERS	3.1	3.6
3 ENERGY LAYERS	<b>2.1</b>	<b>2.5</b>
2 ENERGY LAYERS	2.9	3.3
1 ENERGY LAYERS	20.2	20.7
NO ZCA - 3 ENERGY LAYERS	4.0	4.8
NO RESCALING - 3 ENERGY LAYERS	4.1	4.6
NO ZCA & NO RESCALING - 3 ENERGY LAYERS	9.7	10.1
1D SUBSPACES - 3 ENERGY LAYERS	16.0	17.3
RANDOM SUBSPACES - 3 ENERGY LAYERS	29.4	31.1

Table 5: Systematic studies with our multilayer energy model. The learning hyperparameters are tuned using all 60K labels from the training set.

265 hierarchy with progressively greater selectivity and invariance. However, the number of subspaces  
 266 behaved nonmonotonically with layer, which was unexpected.

267 We can think of the hyperparameter tuning as an outer loop surrounding the unsupervised learning  
 268 algorithm. We will refer to this outer loop as meta-learning. Given an architecture, the unsupervised  
 269 learning algorithm requires no labels at all. However, the outer loop searches for the optimal network  
 270 architecture by using training labels. Therefore, while the learning is unsupervised, the meta-learning  
 271 is supervised. Alternatively, the outer loop can use only a fraction of the training labels, in which  
 272 case the meta-learning is semi-supervised.

273 It is interesting to vary the number of training labels used for hyperparameter search. The results  
 274 are shown in Table 4. The best accuracy is obtained when all 60,000 training set labels are used.  
 275 Accuracy degrades slightly for 5000 labels, and more severely for fewer labels than that. The test  
 276 error can be lower than the training error. This is not a mistake, and it appears to result from a  
 277 non-random ordering of patterns in the MNIST train and test sets.

278 To be clear about the use of data, we note that neither test images nor labels are used during learning  
 279 or meta-learning. Training images but not labels are used during learning. Training labels are used  
 280 by meta-learning. With each iteration of meta-learning, the weights of the network are randomly  
 281 initialized.

## 282 6 Experiments

283 The hyperparameter search explores the space of networks defined by Fig. 1. We can widen the  
 284 space of exploration by performing ablation studies, with results given in Table 5. In all experiments,  
 285 we completely retune the hyperparameters using the full 60K training labels to evaluate clustering  
 286 accuracy. With the exception of the experiment where we vary the number of layers, we use the 3  
 287 energy layer network in this section.

288 **Vary number of energy layers** One can vary the depth of the network by adding or removing energy  
 289 layers. 2 and 4 energy layers yield similar accuracy, and are roughly 1% absolute error (50% relative  
 290 error) worse than for 3 energy layers. A 1 energy layer net is dramatically worse (>20% train/test  
 291 error), suggesting the stacking is critical for performance.

292 The hyperparameters for each of these optimal architectures are provided in the Appendix. The  
 293 optimal 2 energy layer net resembles the 3 energy layer net with its 2nd energy layer removed, while  
 294 simultaneously making the 1st energy layer less sparse.

295 **Remove ZCA whitening** Removing the ZCA layer increases both train and test error of the three  
296 energy layer net by roughly  $2\times$ . This might seem surprising, as Table 2 shows that whitening by  
297 itself decreases accuracy if we directly cluster whitened pixels instead of raw pixels. Apparently, it  
298 is helpful to “take a step backward, before taking 3 steps forward” in the case of our networks with  
299 three energy layers. We have no theoretical explanation for this interesting empirical fact.

300 Preprocessing images by whitening has a long history. The retina has been said to perform a whitening  
301 operation, which is held to be an “efficient coding” that reduces redundancy in an information theoretic  
302 sense [4, 3]. Our experiment suggests that whitening is useful because it improves the accuracy  
303 of subsequent representations produced by stacked unsupervised learning. This seems at least  
304 superficially different from efficient coding theory, because the invariant feature detectors in our  
305 networks appear to discard some information (Figure 4).

306 **Remove rescaling** We now remove the rescaling operation from our energy layers. We observe  
307 a roughly  $2\times$  increase in both test and train error. It may be surprising that such a seemingly  
308 innocuous change can double the resulting train/test errors. This is likely because we only have 17  
309 hyperparameters to optimize; with a limited set of parameters to tune, small changes in architecture  
310 can lead to dramatic performance changes as we only have limited flexibility to tune hyperparameters.

311 **Remove rescaling and ZCA whitening** Removing both ZCA and the per-layer rescaling operations  
312 causes a  $4\times$  increase in train and test error. Apparently the damage to the resulting performance is  
313 multiplicative: removing ZCA alone doubles the train/test error, removing rescaling alone doubles  
314 train/test error, and removing both quadruples the train/test error.

315 **1D subspaces** We rerun the automated hyperparameter tuning experiments from the previous section,  
316 this time restricting our subspaces to be 1D. The subspace norm can be thought of as a conventional  
317 dot product followed by absolute value nonlinearity  $\|\mathbf{V}\mathbf{x}\| = |\mathbf{v} \cdot \mathbf{x}|$  where  $\mathbf{v}$  is the one row of the  
318 subspace matrix  $\mathbf{V}$ . When the subspaces are 1D, our algorithm in fact reduces to that of [9].

319 **Random subspaces** Finally we ask the question: does learning actually help or is it simply the  
320 architecture that matters? To do so, we run tuning experiments with random orthogonal subspaces.  
321 We observe that performance is almost completely erased by using random subspaces, telling us that  
322 the unsupervised learning component is indeed critical for clustering performance.

## 323 7 Discussion

324 The elements of our SUL algorithm were already known in the 2000s: ZCA whitening, energy layers,  
325 K-Subspaces, and adaptive thresholding and rescaling of activities during inference. To achieve  
326 state-of-the-art unsupervised clustering accuracy on MNIST, we employed one more trick, automated  
327 tuning of hyperparameters. Such meta-learning is more feasible now than it was in the 2000s, because  
328 computational power has increased since then.

329 Given that meta-learning optimizes a supervised criterion, is our SUL algorithm really unsupervised?  
330 It is true that the complete system is supervised. However, even if the outer loop (meta-learning) is  
331 supervised, it is accurate to say that the inner loop (learning) is unsupervised. For any hyperparameter  
332 configuration, the network starts from randomly initialized weights, and proceeds to learn in a purely  
333 label-free unsupervised manner.

334 Although MNIST is an easy dataset by today’s standards, we think that the success of our SUL  
335 algorithm at unsupervised clustering is still surprising. We are only tuning 17 hyperparameters, and it  
336 is not clear a priori that our approach would be flexible enough to succeed even for MNIST.

337 In future work, it will be important to investigate more complex datasets or tasks. Following the  
338 more common scenario for meta-learning, future work should train an unsupervised algorithm on a  
339 distribution of tasks, and test transfer to some held-out distribution. We have done some preliminary  
340 experiments with the CIFAR-10 dataset. The meta-learning is more time-consuming because larger  
341 network architectures must be explored. The research is still in progress, but we speculate that the  
342 winner-take-all behavior of K-Subspaces learning may turn out to be a limitation. If so, it will be  
343 important to relax the winner-take-all condition in Equation (1).

344 Evolution created brains through eons of trial-and-error. For us to discover how the brain learns, it  
345 will be important to exploit our computational resources, and use meta-learning to empirically search  
346 the space of biologically plausible learning algorithms.

347 **References**

- 348 [1] EH Adelson and JR Bergen. Spatiotemporal energy models for the perception of motion.  
349 *Journal of the Optical Society of America. A, Optics and Image Science*, 2(2):284–299, 1985.
- 350 [2] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:  
351 A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM*  
352 *SIGKDD international conference on knowledge discovery & data mining*, pages 2623–2631,  
353 2019.
- 354 [3] Joseph J Atick and A Norman Redlich. What does the retina know about natural scenes? *Neural*  
355 *computation*, 4(2):196–210, 1992.
- 356 [4] Horace B Barlow. Unsupervised learning. *Neural computation*, 1(3):295–311, 1989.
- 357 [5] Anthony J Bell and Terrence J Sejnowski. The “independent components” of natural scenes are  
358 edge filters. *Vision research*, 37(23):3327–3338, 1997.
- 359 [6] Yubei Chen, Dylan Paiton, and Bruno Olshausen. The sparse manifold transform. *Advances in*  
360 *neural information processing systems*, 31, 2018.
- 361 [7] Adam Coates and Andrew Ng. Selecting receptive fields in deep networks. *Advances in neural*  
362 *information processing systems*, 24, 2011.
- 363 [8] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsuper-  
364 vised feature learning. In *Proceedings of the fourteenth international conference on artificial*  
365 *intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 366 [9] Adam Coates and Andrew Y Ng. Learning feature representations with k-means. In *Neural*  
367 *networks: Tricks of the trade*, pages 561–580. Springer, 2012.
- 368 [10] D. Eigen. Predicting images using convolutional networks: Visual scene understanding with  
369 pixel maps, 2015.
- 370 [11] K. Fukushima. Neocognitron: a self organizing neural network model for a mechanism of  
371 pattern recognition unaffected by shift in position. *Biol Cybern*, 36(4):193–202, 1980.
- 372 [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
373 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural*  
374 *information processing systems*, 27, 2014.
- 375 [13] Geoffrey E Hinton, Michael Revow, and Peter Dayan. Recognizing handwritten digits using  
376 mixtures of linear models. *Advances in neural information processing systems*, pages 1015–  
377 1022, 1995.
- 378 [14] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with  
379 neural networks. *science*, 313(5786):504–507, 2006.
- 380 [15] Haruo Hosoya and Aapo Hyvärinen. Learning visual spatial pooling by strong pca dimension  
381 reduction. *Neural computation*, 28(7):1249–1264, 2016.
- 382 [16] Aapo Hyvärinen and Patrik Hoyer. Emergence of phase-and shift-invariant features by decom-  
383 position of natural images into independent feature subspaces. *Neural computation*, 12(7):1705–  
384 1720, 2000.
- 385 [17] Aapo Hyvärinen and Urs Köster. Fastisa: A fast fixed-point algorithm for independent subspace  
386 analysis. In *EsANN*, pages 371–376. Citeseer, 2006.
- 387 [18] Xu Ji, Joao F Henriques, and Andrea Vedaldi. Invariant information clustering for unsuper-  
388 vised image classification and segmentation. In *Proceedings of the IEEE/CVF International*  
389 *Conference on Computer Vision*, pages 9865–9874, 2019.
- 390 [19] Zhuxi Jiang, Yin Zheng, Huachun Tan, Bangsheng Tang, and Hanning Zhou. Variational  
391 deep embedding: An unsupervised and generative approach to clustering. *arXiv preprint*  
392 *arXiv:1611.05148*, 2016.

- 393 [20] Yan Karklin and Michael S Lewicki. Emergence of complex cell properties by learning to  
394 generalize in natural scenes. *Nature*, 457(7225):83–86, 2009.
- 395 [21] Koray Kavukcuoglu, Pierre Sermanet, Y-Lan Boureau, Karol Gregor, Michaël Mathieu, Yann  
396 Cun, et al. Learning convolutional feature hierarchies for visual recognition. *Advances in neural  
397 information processing systems*, 23, 2010.
- 398 [22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint  
399 arXiv:1312.6114*, 2013.
- 400 [23] Adam Kosiorek, Sara Sabour, Yee Whye Teh, and Geoffrey E Hinton. Stacked capsule autoen-  
401 coders. *Advances in neural information processing systems*, 32, 2019.
- 402 [24] Dmitry Krotov and John J Hopfield. Unsupervised learning by competing hidden units. *Pro-  
403 ceedings of the National Academy of Sciences*, 116(16):7723–7731, 2019.
- 404 [25] Quoc V Le, Will Y Zou, Serena Y Yeung, and Andrew Y Ng. Learning hierarchical invariant  
405 spatio-temporal features for action recognition with independent subspace analysis. In *CVPR  
406 2011*, pages 3361–3368. IEEE, 2011.
- 407 [26] Quoc V Le, Will Y Zou, Serena Y Yeung, and Andrew Y Ng. Learning hierarchical invariant  
408 spatio-temporal features for action recognition with independent subspace analysis. In *CVPR  
409 2011*, pages 3361–3368. IEEE, 2011.
- 410 [27] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>,  
411 1998.
- 412 [28] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix  
413 factorization. *Nature*, 401(6755):788–791, 1999.
- 414 [29] Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Y Ng. Convolutional deep belief  
415 networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of  
416 the 26th annual international conference on machine learning*, pages 609–616, 2009.
- 417 [30] Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. Random synap-  
418 tic feedback weights support error backpropagation for deep learning. *Nature communications*,  
419 7(1):1–10, 2016.
- 420 [31] Timothy P Lillicrap, Adam Santoro, Luke Marris, Colin J Akerman, and Geoffrey Hinton.  
421 Backpropagation and the brain. *Nature Reviews Neuroscience*, 21(6):335–346, 2020.
- 422 [32] Kyle Luther and H Sebastian Seung. Sensitivity of sparse codes to image distortions. *arXiv  
423 preprint arXiv:2204.07466*, 2022.
- 424 [33] Alireza Makhzani and Brendan Frey. K-sparse autoencoders. *arXiv preprint arXiv:1312.5663*,  
425 2013.
- 426 [34] Ryan McConville, Raul Santos-Rodriguez, Robert J Piechocki, and Ian Craddock. N2d:(not  
427 too) deep clustering via clustering the local manifold of an autoencoded embedding. In *2020  
428 25th International Conference on Pattern Recognition (ICPR)*, pages 5145–5152. IEEE, 2021.
- 429 [35] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation  
430 and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 431 [36] Luke Metz, Niru Maheswaranathan, Brian Cheung, and Jascha Sohl-Dickstein. Meta-learning  
432 update rules for unsupervised representation learning. *arXiv preprint arXiv:1804.00222*, 2018.
- 433 [37] Luke Metz, Niru Maheswaranathan, Brian Cheung, and Jascha Sohl-Dickstein. Meta-learning  
434 update rules for unsupervised representation learning. In *7th International Conference on  
435 Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net,  
436 2019.
- 437 [38] Sudipto Mukherjee, Himanshu Asnani, Eugene Lin, and Sreeram Kannan. Clustergan: Latent  
438 space clustering in generative adversarial networks. In *Proceedings of the AAAI conference on  
439 artificial intelligence*, volume 33, pages 4610–4617, 2019.

- 440 [39] Dina Obeid, Hugo Ramambason, and Cengiz Pehlevan. Structured and deep similarity matching  
441 via structured and deep hebbian networks. *Advances in neural information processing systems*,  
442 32, 2019.
- 443 [40] Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by  
444 learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.
- 445 [41] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. *arXiv*  
446 *preprint arXiv:1710.09829*, 2017.
- 447 [42] Ko Sakai and Shigeru Tanaka. Spatial pooling in the second-order spatial structure of cortical  
448 complex cells. *Vision Research*, 40(7):855–871, 2000.
- 449 [43] Jeremias Sulam, Vardan Papayan, Yaniv Romano, and Michael Elad. Multilayer convolutional  
450 sparse modeling: Pursuit and dictionary learning. *IEEE Transactions on Signal Processing*,  
451 66(15):4090–4104, 2018.
- 452 [44] Louis Thiry, Michael Arbel, Eugene Belilovsky, and Edouard Oyallon. The unreasonable effec-  
453 tiveness of patches in deep convolutional kernels methods. *arXiv preprint arXiv:2101.07528*,  
454 2021.
- 455 [45] René Vidal. Subspace clustering. *IEEE Signal Processing Magazine*, 28(2):52–68, 2011.
- 456 [46] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol,  
457 and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep  
458 network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010.
- 459 [47] Dingding Wang, Chris Ding, and Tao Li. K-subspace clustering. In *Joint European Conference*  
460 *on Machine Learning and Knowledge Discovery in Databases*, pages 506–521. Springer, 2009.
- 461 [48] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering  
462 analysis. In *International conference on machine learning*, pages 478–487. PMLR, 2016.
- 463 [49] Yi Yang, Dong Xu, Feiping Nie, Shuicheng Yan, and Yueting Zhuang. Image clustering using  
464 local discriminant models and global integration. *IEEE Transactions on Image Processing*,  
465 19(10):2761–2773, 2010.

## 466 Checklist

- 467 1. For all authors...
- 468 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
469 contributions and scope? [Yes] In the abstract we claim 2% clustering error with a  
470 stacked-unsupervised method, and we show this in Table 2, and provide the method in  
471 section 3,4,5
- 472 (b) Did you describe the limitations of your work? [Yes] See discussion
- 473 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 474 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
475 them? [Yes]
- 476 2. If you are including theoretical results...
- 477 (a) Did you state the full set of assumptions of all theoretical results? [Yes] In the ap-  
478 pendix, we state our assumptions for our theoretical result (that our power updates for  
479 optimizing the K-Subspace objective converges in the full batch setting)
- 480 (b) Did you include complete proofs of all theoretical results? [Yes] In the appendix, we  
481 state prove our theoretical result that the algorithm for optimizing the K-Subspace  
482 objective converges in the full batch setting
- 483 3. If you ran experiments...
- 484 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
485 mental results (either in the supplemental material or as a URL)? [Yes] We will include  
486 code in supplemental material

- 487 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
488 were chosen)? [Yes] In section 5 we provide these details
- 489 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
490 ments multiple times)? [No] Unfortunately we did not have adequate time to complete  
491 these experiments, as each meta-learning experiment required training and evaluating  
492 1000s of different architectures
- 493 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
494 of GPUs, internal cluster, or cloud provider)? [Yes] In section 5 on hyperparameter  
495 tuning, we state that each 3 layer meta-learning experiment takes approximately 20  
496 hours on 8 GPUs
- 497 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 498 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite mnist set  
499 creators.
- 500 (b) Did you mention the license of the assets? [No]
- 501 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 502 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
503 using/curating? [No] Dataset only includes handwritten digits.
- 504 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
505 information or offensive content? [No] Dataset only includes handwritten digits.
- 506 5. If you used crowdsourcing or conducted research with human subjects...
- 507 (a) Did you include the full text of instructions given to participants and screenshots, if  
508 applicable? [N/A]
- 509 (b) Did you describe any potential participant risks, with links to Institutional Review  
510 Board (IRB) approvals, if applicable? [N/A]
- 511 (c) Did you include the estimated hourly wage paid to participants and the total amount  
512 spent on participant compensation? [N/A]

## 513 A Online minimization of Eq. 1

---

### Algorithm 1: K-Subspace clustering via minibatch power iterations

---

Initialize subspaces  $\{\mathbf{V}_k \in \mathbb{R}^{r \times mp^2}\}$   
**for**  $i = 1, 2, 3$  **do**  
    Sample patches from a minibatch of images  $\{\mathbf{x}_u \in \mathbb{R}^{mp^2}\}$   
    Cluster patches  
        
$$\mathbf{X}_k := \{\mathbf{x}_u : k = \underset{q}{\operatorname{argmin}} \|\mathbf{x}_u - \mathbf{V}_q^\top \mathbf{V}_q \mathbf{x}_u\|\} \quad \text{for } k = 1, 2, 3, \dots$$
  
    Apply one orthogonal power iteration to every subspace  
        
$$\mathbf{USV}^\top := \mathbf{X}_k^\top \mathbf{X}_k \mathbf{V}_k^\top \quad \text{SVD decomposition}$$
  
        
$$\mathbf{V}_k := \mathbf{U}^\top \quad \text{for } k = 1, 2, 3, \dots$$
  
**end**

---

514 A standard algorithm for minimizing Eq. 1 is a full batch EM algorithm that alternates between  
515 cluster assignment (E-step) and using PCA to set  $\mathbf{V}_k$  to the top  $r$  principle components of the patches  
516 assigned to each cluster  $k$  (M-step) [45]. The full batch requirement makes this algorithm rather slow.

517 The algorithm we present and use in this paper is described in Algorithm 1. We make two core  
518 changes to the standard EM algorithm. One, we use minibatch updates instead of full batch updates.  
519 Two, we perform a single step of power-iteration after each cluster assignment step, instead of  
520 performing a full PCA.

521 Proper initialization can impact on the quality of learned subspaces. We initialize subspaces by  
522 setting the first dimension to a randomly chosen patch, and the other dimensions to white Gaussian

523 noise with  $\mu = 0, \sigma = 0.01$ . We then perform a few “warmup” iterations of the main loop in alg. 1,  
 524 except we only cluster using the first subspace dimension. We use *warmup\_iter* = 10 in all our  
 525 experiments. During these warmup iterations, the power updates are still performed on the full rank- $r$   
 526 subspaces. Intuitively, this warmup procedure generates clusters with 1D subspace clustering, and  
 527 initializes subspaces be the top  $r$  components within these clusters.

528 Convergence of our algorithm is discussed in the Appendix. We prove that the clustering+power  
 529 iteration step ensures the loss computed over the minibatch is non-decreasing. We show empirically  
 530 that the loss computed over the whole dataset decreases with iteration for a reasonable setting of  
 531 parameters. Because we only use a single pass through the data to train each layer, our algorithm is  
 532 much faster than applying the full batch EM-style algorithm described in [45] to the whole dataset.

## 533 B Convergence analysis of Algorithm 1

### 534 B.1 Theory: full-batch convergence

535 We will not be able to provide a full proof that Algorithm 1 converges in the minibatch setting.  
 536 However we can at least show that in the full batch setting, every iteration of Algorithm 1 decreases  
 537 the energy in Equation (1). Of course the cluster assignment portion of the algorithm decreases the  
 538 energy. What we will show here is that the power iteration step also decreases the energy at every  
 539 iteration.

540 We recall the notation from Algorithm 1. We define the matrix  $\mathbf{X}_k$  whose rows are the input patches  
 541 assigned to cluster  $k$ :

$$\mathbf{X}_k := \{\mathbf{x}_u : k = \operatorname{argmin}_q \|\mathbf{x}_u - \mathbf{V}_q^\top \mathbf{V}_q \mathbf{x}_u\|\} \quad (2)$$

542 The energy in Equation (1) can be written as a sum of energies for each cluster  $e = \sum_k e_k$  where:

$$e_k = \|\mathbf{X}_k - \mathbf{X}_k \mathbf{V}_k^\top \mathbf{V}_k\|_F^2 \quad (3)$$

543 We will show that a power update for cluster  $k$  now gives a non-decreasing energy  $e_k$ . To avoid  
 544 notational clutter, we will drop the subspace index  $k$  and assume we are working with a single fixed  
 545 cluster for now. It will actually be easier to show that a more general class of updates than the SVD  
 546 update in Algorithm 1 cause the energy to remain or decrease. Suppose we have any subspace update  
 547 defined by:

$$\mathbf{V}' := \mathbf{Q}(\mathbf{V}\mathbf{C}\mathbf{V})^\dagger \mathbf{C}\mathbf{V} \quad (4)$$

548 where  $\mathbf{Q}$  is any orthogonal  $r \times r$  matrix (that can also be a function of  $\mathbf{V}, \mathbf{C}$ ). The power update  
 549 in Algorithm 1 is an example of such an update. We will show that the energy for every subspace  
 550 is non-decreasing with this update:  $e' \leq e$ . We do so by relating the update in Equation (4) to one  
 551 sequence of steps in an alternating least squares problem. Define:

$$h(\mathbf{A}, \mathbf{B}) := \|\mathbf{X} - \mathbf{A}\mathbf{B}^\top\|^2 \quad (5)$$

552 Define  $\mathbf{A} = \mathbf{X}\mathbf{V}^\top$  and  $\mathbf{B} = \mathbf{V}^\top$ . Then  $e = h(\mathbf{A}, \mathbf{B})$ . Define the sequence:

$$\begin{aligned} \mathbf{B}' &= \operatorname{argmin}_{\mathbf{U}} h(\mathbf{A}, \mathbf{U}) = \mathbf{X}\mathbf{A}(\mathbf{A}^\top \mathbf{A})^\dagger \\ \mathbf{A}' &= \operatorname{argmin}_{\mathbf{U}} h(\mathbf{U}, \mathbf{B}') = \mathbf{X}\mathbf{B}'((\mathbf{B}')^\top \mathbf{B}')^\dagger \end{aligned} \quad (6)$$

553 We can multiply  $\mathbf{A}'(\mathbf{B}')^\top$  and substitute the above equations to get:

$$\mathbf{A}'(\mathbf{B}')^\top = \mathbf{X}(\mathbf{V}')^\top \mathbf{V}' \quad (7)$$

554 We therefore have that  $e' = h(\mathbf{A}', \mathbf{B}') \leq h(\mathbf{A}, \mathbf{B}) = e$ , so the subspace energy is non-increasing  
 555 under the updates in Algorithm 1.

### 556 B.2 Experiment: learning curves

557 We show learning curves using Algorithm 1 applied to MNIST digits in Figure 5. We run this  
 558 algorithm for a single pass through the 60k training set patterns. Inputs are first whitened with

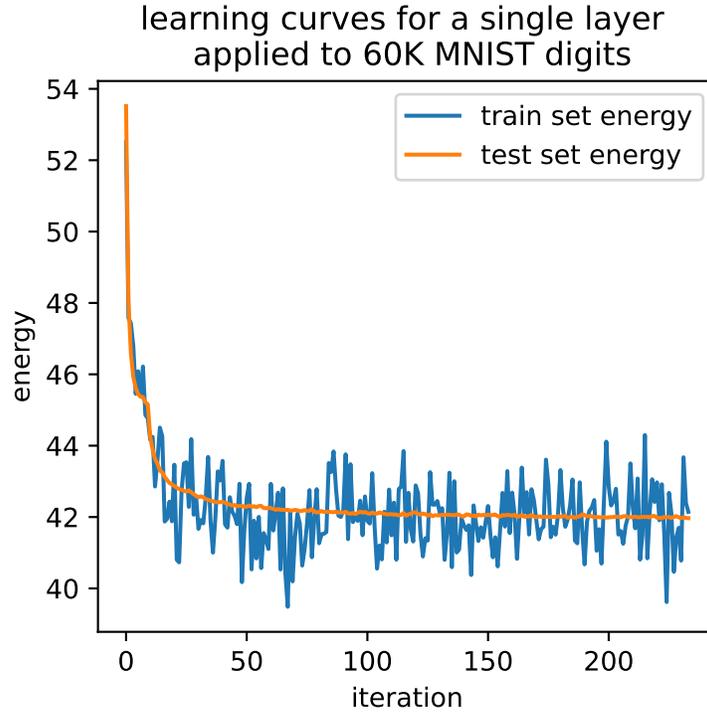


Figure 5: Learning algorithms for Algorithm 1 applied to MNIST digits.

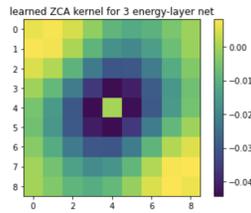


Figure 6: Learned ZCA kernel for 3 energy layer net. We have zero-ed out the central pixel, as its value is 1.17 and thus would make it harder to see the surround structure. This filter resembles an oblique center surround filter.

559 ConvZCA using a kernel size of 11 and num component of 16. We learn 64 subspaces, each of  
 560 dimension 5, and kernel size of 9. We use a minibatch size of 256.

561 The training set energy is computed over a minibatch of 256 inputs at each iteration. The test set  
 562 energy is computed over all 10K test set patterns, which is why it is less noisy. Empirically we  
 563 see that for this setting of parameters at least, our minibatch K-Subspace algorithm does lead to a  
 564 decreasing energy computed over unseen patterns.

565 We apply 10 warmup iterations (only using the first subspace dimension to cluster for the first 10  
 566 iterations), which is why we observe the sharp drop-off in energy after 10 iterations.

## 567 C Hyperparameters for 1,2,3,4 energy layer nets

568 In Table 6 we show the learned architectures for the 1,2,3,4 energy layer networks.

Table 6: Network architectures for 1,2,3,4 energy layer nets from Table 5.

HYPERPARAMETER	1 LAYER	2 LAYER	3 LAYER	4 LAYER
CONVZCA KERNEL SIZE	5	9	9	11
CONVZCA N COMPONENTS	0	18	9	18
LAYER1 SUBSPACE NUMBER ( $k$ )	59	20	37	15
LAYER1 SUBSPACE RANK ( $r$ )	2	5	2	2
LAYER1 ACTIVE FEATURES ( $w$ )	1	16	9	10
LAYER1 KERNEL SIZE	10	10	8	7
LAYER1 PADDING	4	2	2	3
LAYER2 SUBSPACE NUMBER ( $k$ )	-	55	9	63
LAYER2 SUBSPACE RANK ( $r$ )	-	16	3	12
LAYER2 ACTIVE FEATURES ( $w$ )	-	1	8	1
LAYER2 KERNEL SIZE	-	19	5	17
LAYER2 PADDING	-	1	1	1
LAYER3 SUBSPACE NUMBER ( $k$ )	-	-	58	57
LAYER3 SUBSPACE RANK ( $r$ )	-	-	16	7
LAYER3 ACTIVE FEATURES ( $w$ )	-	-	2	6
LAYER3 KERNEL SIZE	-	-	21	8
LAYER3 PADDING	-	-	2	2
LAYER4 SUBSPACE NUMBER ( $k$ )	-	-	-	22
LAYER4 SUBSPACE RANK ( $r$ )	-	-	-	1
LAYER4 ACTIVE FEATURES ( $w$ )	-	-	-	1
LAYER4 KERNEL SIZE	-	-	-	3
LAYER4 PADDING	-	-	-	1

## 569 D Learned ZCA filter

570 In Figure 6 we show the learned ConvZCA kernel for the 3 energy layer network. It resembles an  
 571 oblique center surround filter. It is interesting to calculate the relative weight of the negative surround  
 572 vs positive center term. Specifically we calculate:

$$f = \frac{I[4, 4] - \sum_{u,v} \max\{0, -I_{u,v}\}}{I[4, 4]} = 0.10 \quad (8)$$

573 where  $I$  is the 9x9 kernel and  $I[4, 4]$  is the central pixel of that kernel. In other words, there is a small  
 574 DC component (the central pixel is not perfectly cancelled out by the negative surround).

## 575 E Clustering UMAP embeddings

576 In the main text we reported a surprisingly low performance for K-Means applied to UMAP em-  
 577 bedding, and that using a Gaussian mixture model instead can significantly improve the accuracy.  
 578 Here we show the UMAP embeddings, and corresponding clusterings generated via K-Means and  
 579 Gaussian Mixture Models.

580 Note that our result is not inconsistent with the result given in Table 1 of [34] who reported 82.5%  
 581 clustering accuracy (compared to our result of 96.4% in Table 5) when using a GMM to cluster the  
 582 embedding vectors return by UMAP. This is because we are using UMAP to embed to 2D while they  
 583 used UMAP to embed to 10 dimennsions. For this task it seems that the extremely low dimensional  
 584 embeddings are actually more suitable for downstream clustering.

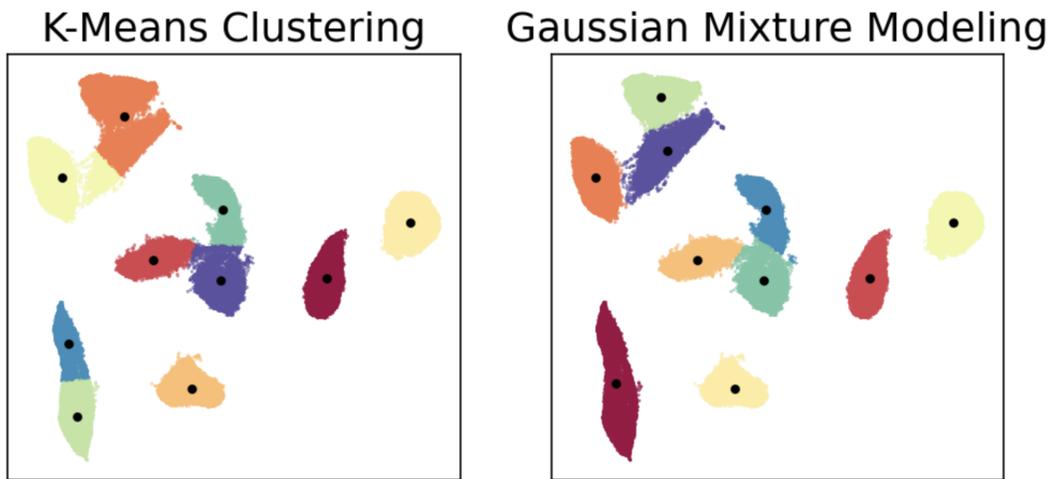


Figure 7: K-Means vs. Gaussian Mixture-based clustering applied to UMAP with "out-of-the-box" parameters applied to MNIST handwritten digits. Colors are assigned to each point based of the clustering (not the ground truth labels). K-Means erroneously merges and splits some of the clusters, while Gaussian Mixture models give a much more intuitive clustering (and ultimately a much lower clustering error as shown in Table 5)