

---

# Deep Networks on Graph-Structured Data

---

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

### 1 Introduction

In recent times, deep learning models have proven extremely successful on a wide variety of tasks, from computer vision and acoustic modeling to natural language processing [8]. At the core of their success lies an important assumption on the statistical properties of the data, namely the *stationarity* and the *compositionality* through local statistics, which is present in natural images, video, and speech. These properties are exploited efficiently by ConvNets [7, 6], which are designed to extract local features that are shared across the signal domain. Thanks to this, they are able to greatly reduce the number of parameters in the network with respect to generic deep architectures, without sacrificing the capacity to extract informative statistics from the data. Similarly, Recurrent Neural Nets (RNNs) trained on temporal data implicitly assume a stationary distribution.

One can think of such data examples as being signals defined on a low-dimensional grid. In this case stationarity is well defined via the natural translation operator on the grid, locality is defined via the metric of the grid, and compositionality is obtained from downsampling, or equivalently thanks to the multi-resolution property of the grid. However, there exist many examples of data that lack the underlying low-dimensional grid structure. For example, text documents represented as bags of words can be thought of as signals defined on a graph whose nodes are vocabulary terms and whose weights represent some similarity measure between terms, such as co-occurrence statistics. In medicine, a patient's gene expression data can be viewed as a signal defined on the graph imposed by the regulatory network. In fact, computer vision and audio, which are the main focus of research efforts in deep learning, only represent a special case of data defined on an extremely simple low-dimensional graph. These complex graphs might be of higher dimension, and the statistical properties of data defined on such graphs might not satisfy the stationarity, locality and compositionality assumptions previously described. For such type of data of dimension  $N$ , deep learning strategies are reduced to learning with fully-connected layers, which have  $O(N^2)$  parameters, and regularization is carried out via weight decay and dropout [15].

When the graph structure of the input is known, [2] introduced a model to generalize ConvNets using low learning complexity similar to that of a ConvNet, and which was demonstrated on simple low-dimensional graphs. In this work, we are interested in generalizing ConvNets to high-dimensional, general datasets, and, most importantly, to the setting where the graph structure is not known a priori. In this context, learning the graph structure amounts to estimating the similarity matrix, which has complexity  $O(N^2)$ . One may therefore wonder whether the graph estimation followed by graph convolutions offers advantages with respect to learning directly from the data with fully connected layers. We attempt to answer this question experimentally and to establish baselines for future reference.

We explore these approaches in two areas of application for which it has not been possible to apply convolutional networks before: text categorization and bioinformatics. Our results show that

054 our method is capable of matching or outperforming large, fully-connected networks trained with  
055 dropout using fewer parameters. Our main contributions can be summarized as follows:  
056

- 057 • We extend the ideas from [2] to large-scale classification problems, specifically object  
058 Recognition, text categorization and bioinformatics.
- 059 • We consider the most general setting where no prior information on the graph structure is  
060 available, and propose unsupervised and supervised graph estimation strategies.
- 061 • Finally, we introduce an alternative formulation for efficient learning in graph-structured  
062 domains which works directly in the feature domain.  
063

064 The rest of the paper is structured as follows. Section 2 reviews similar works in the literature. Sec-  
065 tion 3 discusses generalizations of convolutions on graphs, and Section 4 addresses the question of  
066 graph estimation. Finally, Section 5 shows numerical experiments on large scale object recogniton,  
067 text categorization and bioinformatics.

## 068 2 Related Work 069

070 There have been several works which have explored architectures using the so-called local receptive  
071 fields [5, 4, 12], mostly with applications to image recognition. In particular, [4] proposes a scheme  
072 to learn how to group together features based upon a measure of similarity that is obtained in an  
073 unsupervised fashion. However, it does not attempt to exploit any weight-sharing strategy.  
074

075 Recently, [2] proposed a generalization of convolutions to graphs via the Graph Laplacian. By  
076 identifying a linear, translation-invariant operator in the grid (the Laplacian operator), with its coun-  
077 terpart in a general graph (the Graph Laplacian), one can view convolutions as the family of linear  
078 transforms commuting with the Laplacian. By combining this commutation property with a rule  
079 to find localized filters, the model requires only  $O(1)$  parameters per “feature map”. However,  
080 this construction requires prior knowledge of the graph structure, and was shown only on simple,  
081 low-dimensional graphs. More recently, [11] introduced Shapenet, another generalization of con-  
082 volutions on non-Euclidean domains based on geodesic polar coordinates, which was successfully  
083 applied to shape analysis, and allows comparison across different manifolds. However, it also re-  
084 quires prior knowledge of the manifolds.  
085

086 The graph or similarity estimation aspects have also been extensively studied in the past. For in-  
087 stance, [13] studies the estimation of the graph from a statistical point of view, through the identifi-  
088 cation of a certain graphical model using  $\ell_1$ -penalized logistic regression. Also, [3] considers the  
089 problem of learning a deep architecture through a series of Haar contractions, which are learnt using  
an unsupervised pairing criteria over the features.

## 090 3 Generalizing Convolutions in Graphs 091

### 092 3.1 Spectral Networks 093

094 Our work builds upon [2] which introduced spectral networks. We recall the definition here and its  
095 main properties. A spectral network generalizes a convolutional network through the Graph Fourier  
096 Transform, which is in turn defined via a generalization of the Laplacian operator on the grid to the  
097 graph Laplacian. An input vector  $x \in \mathbb{R}^N$  is seen as a signal defined on a graph  $G$  with  $N$  nodes.  
098

099 **Definition 1.** Let  $W$  be a  $N \times N$  similarity matrix representing an undirected graph  $G$ , and let  
100  $L = D - W$  be its graph Laplacian with  $D = W \cdot \mathbf{1}$  eigenvectors  $U = (u_1, \dots, u_N)$ . Then a graph  
101 convolution of input signals  $x$  with filters  $g$  on  $G$  is defined by  $x *_G g = U^T (Ux \odot Ug)$ , where  $\odot$   
represents a point-wise product.  
102

103 Here, the unitary matrix  $U$  plays the role of the Fourier Transform in  $\mathbb{R}^d$ . There are several ways  
104 of computing the graph Laplacian  $L$  [1]. In this paper, we choose the normalized version  $L =$   
105  $I - D^{-1/2}WD^{-1/2}$ , where  $D$  is a diagonal matrix with entries  $D_{ii} = \sum_j W_{ij}$ . Note that in the case  
106 where  $W$  represents the lattice, from the definition of  $L$  we recover the discrete Laplacian operator  
107  $\Delta$ . Also note that the Laplacian commutes with the translation operator, which is diagonalized in  
the Fourier basis. It follows that the eigenvectors of  $\Delta$  are given by the Discrete Fourier Transform

(DFT) matrix. We then recover a classical convolution operator by noting that convolutions are by definition linear operators that diagonalize in the Fourier domain (also known as the Convolution Theorem [10]).

Learning filters on a graph thus amounts to learning spectral multipliers  $w_g = (w_1, \dots, w_N)$

$$x *_G g := U^T(\text{diag}(w_g)Ux) .$$

Extending the convolution to inputs  $x$  with multiple input channels is straightforward. If  $x$  is a signal with  $M$  input channels and  $N$  locations, we apply the transformation  $U$  on each channel, and then use multipliers  $w_g = (w_{i,j} ; i \leq N, j \leq M)$ .

However, for each feature map  $g$  we need convolutional kernels are typically restricted to have small spatial support, independent of the number of input pixels  $N$ , which enables the model to learn a number of parameters independent of  $N$ . In order to recover a similar learning complexity in the spectral domain, it is thus necessary to restrict the class of spectral multipliers to those corresponding to localized filters.

For that purpose, we seek to express spatial localization of filters in terms of their spectral multipliers. In the grid, smoothness in the frequency domain corresponds to the spatial decay, since

$$\left| \frac{\partial^k \hat{x}(\xi)}{\partial \xi^k} \right| \leq C \int u^k |x(u)| du ,$$

where  $\hat{x}(\xi)$  is the Fourier transform of  $x$ . In [2] it was suggested to use the same principle in a general graph, by considering a smoothing kernel  $\mathcal{K} \in \mathbb{R}^{N \times N_0}$ , such as splines, and searching for spectral multipliers of the form

$$w_g = \mathcal{K} \tilde{w}_g .$$

The algorithm which implements the graph convolution is described in 1.

---

#### Algorithm 1 Train Graph Convolution Layer

---

- 1: Given GFT matrix  $U$ , interpolation kernel  $\mathcal{K}$ , weights  $w$ .
  - 2: **Forward Pass:**
  - 3: Fetch input batch  $x$  and gradients w.r.t outputs  $\nabla y$ .
  - 4: Compute interpolated weights:  $w_{f'f} = \mathcal{K} \tilde{w}_{f'f}$ .
  - 5: Compute output:  $y_{sf'} = U^T \left( \sum_f U x_{sf} \odot w_{f'f} \right)$ .
  - 6: **Backward Pass:**
  - 7: Compute gradient w.r.t input:  $\nabla x_{sf} = U^T \left( \sum_{f'} \nabla y_{sf'} \odot w_{f'f} \right)$
  - 8: Compute gradient w.r.t interpolated weights:  $\nabla w_{f'f} = U^T \left( \sum_s \nabla y_{sf'} \odot x_{sf} \right)$
  - 9: Compute gradient w.r.t weights  $\nabla \tilde{w}_{f'f} = \mathcal{K}^T \nabla w_{f'f}$ .
- 

### 3.2 Pooling with Hierarchical Graph Clustering

In image and speech applications, and in order to reduce the complexity of the model, it is often useful to trade off spatial resolution for feature resolution as the representation becomes deeper. For that purpose, pooling layers compute statistics in local neighborhoods, such as the average amplitude, energy or maximum activation.

The same layers can be defined in a graph by providing the equivalent notion of neighborhood. In this work, we construct such neighborhoods at different scales using multi-resolution spectral clustering [17], and consider both average and max-pooling as in standard convolutional network architectures.

## 4 Graph Construction

Whereas some recognition tasks in non-Euclidean domains, such as those considered in [2] or [11], might have a prior knowledge of the graph structure of the input data, many other real-world applications do not have such knowledge. It is thus necessary to estimate a similarity matrix  $W$  from

162 the data before constructing the spectral network. In this paper we consider two possible graph  
 163 constructions, one unsupervised by measuring joint feature statistics, and another one supervised using  
 164 an initial network as a proxy for the estimation.

#### 166 4.1 Unsupervised Graph Estimation

168 Given data  $X \in \mathbb{R}^{L \times N}$ , where  $L$  is the number of samples and  $N$  the number of features, the  
 169 simplest approach to estimating a graph structure from the data is to consider a distance between  
 170 features  $i$  and  $j$  given by

$$171 \quad d(i, j) = \|X_i - X_j\|^2,$$

172 where  $X_i$  is the  $i$ -th column of  $X$ . While correlations are typically sufficient to reveal the intrinsic  
 173 geometrical structure of images [14], the effects of higher-order statistics might be non-negligible in  
 174 other contexts, especially in presence of sparsity. Indeed, in many situations the pairwise Euclidean  
 175 distances might suffer from unnormalized measurements. Several strategies and variants exist to  
 176 gain some robustness, for instance replacing the Euclidean distance by the  $Z$ -score (thus renormalizing  
 177 each feature by its standard deviation), the “square-correlation” (computing the correlation of  
 178 squares of previously whitened features), or the mutual information.

179 This distance is then used to build a Gaussian diffusion Kernel [1]

$$180 \quad 181 \quad \omega(i, j) = \exp^{-\frac{d(i, j)}{\sigma^2}}. \quad (1)$$

182 In our experiments, consider the variant of self-tuning diffusion kernel [18]

$$183 \quad 184 \quad \omega(i, j) = \exp^{-\frac{d(i, j)}{\sigma_i \sigma_j}},$$

186 where  $\sigma_i$  is computed as the distance  $d(i, i_k)$  corresponding to the  $k$ -th nearest neighbor  $i_k$  of feature  
 187  $i$ .

188 The main advantage of (1) is that it does not require labeled data. Therefore, it is possible to estimate  
 189 the similarity using several datasets that share the same features, for example in text classification.

#### 191 4.2 Supervised Graph Estimation

193 As discussed in the previous section, the notion of feature similarity is not well defined, as it depends  
 194 on our choice of kernel and criteria. Therefore, in the context of supervised learning, the relevant  
 195 statistics from the input signals might not correspond to our imposed similarity criteria. It may thus  
 196 be interesting to ask for the feature similarity that best suits a particular classification task.

197 A particularly simple approach is to use a fully-connected network to determine the feature similarity.  
 198 Given a training set with normalized <sup>1</sup> features  $X \in \mathbb{R}^{L \times N}$  and labels  $y \in \{1, \dots, C\}^L$ , we  
 199 initially train a fully connected network  $\phi$  with  $K$  layers of weights  $W_1, \dots, W_K$ , using standard  
 200 ReLU activations and dropout. We then extract the first layer features  $W_1 \in \mathbb{R}^{N \times M_1}$ , where  $M_1$  is  
 201 the number of first-layer hidden features, and consider the distance

$$202 \quad 203 \quad d_{sup}(i, j) = \|W_{1,i} - W_{1,j}\|^2, \quad (2)$$

204 that is then fed into the Gaussian kernel as in (1). The interpretation is that the supervised criterion  
 205 will extract through  $W_1$  a collection of linear measurements that best serve the classification  
 206 task. Thus two features are similar if the network decides to extract similar collections of linear  
 207 measurements.

208 Another supervised strategy that further exploits the dependencies learnt at the deeper layers of the  
 209 network consists in using the gradients of the supervised loss function with respect to the inputs as  
 210 discriminative information. More concretely, if  $\mathcal{L}(X_l, y_l, (W_k)_{k \leq K})$  is the loss function associated  
 211 to training example  $l$  and  $Z_{l,i} = \nabla_{x_i} \mathcal{L}(X_l, y_l, (W_k)_{k \leq K})$  is the input gradient corresponding to that  
 212 example with respect to the  $i$ -th feature, then we consider the distance

$$213 \quad 214 \quad d_{gsup}(i, j) = \|Z_{.,i} - Z_{.,j}\|^2, \quad (3)$$

---

215 <sup>1</sup>In our experiments we simply normalized each feature by its standard deviation, but one could also whiten  
 216 completely the data.

216 that is, we measure the average similarity between the gradients for features  $i$  and  $j$ . That is, the  
217 similarity of two features depends on their respective influence (averaged over the training set) to  
218 the outcome of the network.

219 These constructions can be seen as distilling the information learnt by a first network into a kernel.  
220 In the general case where no assumptions are made on the dimension of the graph, it amounts  
221 to extracting  $N^2$  parameters from the first learning stage (which typically involves a much larger  
222 number of parameters). If, moreover, we assume a low-dimensional graph structure of dimension  
223  $m$ , then  $mN$  parameters are extracted by projecting the resulting kernel into its leading  $m$  directions.  
224

225 Finally, observe that one could simply replace the eigen-basis  $U$  obtained by diagonalizing the graph  
226 Laplacian by an arbitrary unitary matrix, which is then optimized by back-propagation together with  
227 the rest of the parameters of the model. We do not report results on this strategy, although we point  
228 out that it was the same learning complexity as the Fully Connected network (requiring  $O(KN^2)$   
229 parameters, where  $K$  is the number of layers and  $N$  is the input dimension).

## 234 5 Experiments

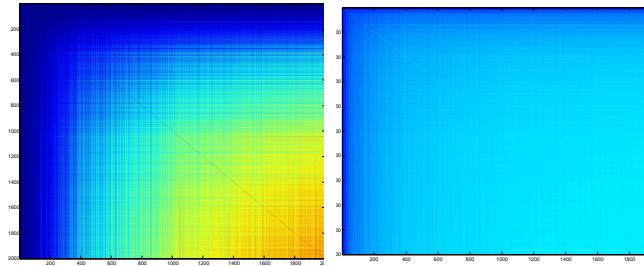
235 In order to measure the performance of spectral networks on real-world data and to explore the  
236 effect of the graph estimation procedure, we conducted experiments on three datasets from text  
237 categorization, computational biology and computer vision. All experiments were done using the  
238 Torch machine learning environment with a custom CUDA backend, which is publicly available at  
239 URL.

240 In this section we adopt the following notation to describe network architectures:  $\text{GC}_k$  denotes a  
241 graph convolution layer with  $k$  feature maps,  $\text{Pk}$  denotes a graph pooling layer with stride  $k$  and  
242 pool size  $2k$ , and  $\text{FC}_k$  denotes a fully connected layer with  $k$  hidden units.  
243

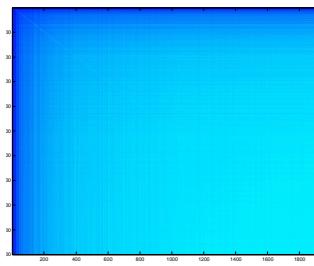
### 252 5.1 Reuters

253  
254  
255 We used the Reuters dataset described in [16], which consists of training and test sets each con-  
256 taining 201,369 documents from 50 mutually exclusive classes. Each document is represented as a  
257 log-normalized bag of words for 2000 common non-stop words. As a baseline we used the fully-  
258 connected network of [16] with two hidden layers consisting of 2000 and 1000 hidden units regu-  
259 larized with dropout.

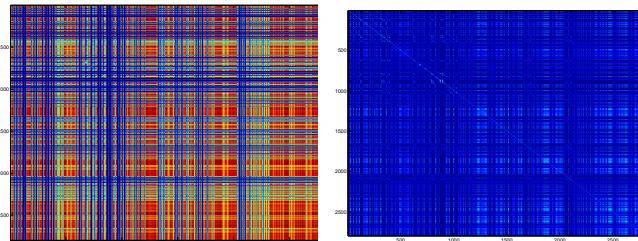
260 We based the spectral network architecture on that of a classical convolutional network, namely by  
261 interleaving graph convolution and graph pooling layers and ending with a fully connected layer.  
262 Performing pooling at the beginning of the network was especially important to reduce the dimen-  
263 sionality in the graph domain and mitigate the cost of the expensive Graph Fourier Transform op-  
264 eration. We chose hyperparameters by performing initial experiments on a validation set consisting  
265 of one-tenth of the training data. Specifically, we set the number of subsampled weights to  $k = 60$ ,  
266 learning rate to 0.01 and used max pooling rather than average pooling. We also found that using  
267 AdaGrad [?] made training faster. All architectures were then trained using the same hyperparam-  
268 eters. Since the experiments were computationally expensive, we did not train all models until full  
269 convergence. This enabled us to explore more model architectures and obtain a clearer understand-  
ing of the effects of graph construction.



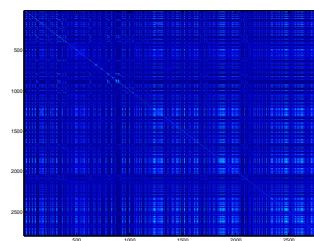
283 (a) Global scaling..



283 (b) Local scaling.



298 (c) Global scaling.



298 (d) Local scaling.

299 Figure 1: Similarity graphs for the Reuters (top) and Merck DPP4 (bottom) datasets.

301 Table 1: Results for Reuters dataset

Graph	Architecture	Parameters	Acc. (200)	Acc. (500)
-	FC-2000-1000-50	$6 \cdot 10^6$	70.2	70.2
Supervised	GC4-P4-FC1000	$2 \cdot 10^6$	69.41	70.09
Supervised	GC8-P8-FC1000	$2 \cdot 10^6$	69.15	-
Supervised low rank	GC4-P4-FC1000	$2 \cdot 10^6$	69.25	-
Supervised low rank	GC8-P8-FC1000	$2 \cdot 10^6$	running cims	-
Supervised	GC16-P4-GC16-P4-FC1000	$2 \cdot 10^6$	69.04	-
Supervised	GC64-P8-GC64-P8-FC1000	$2 \cdot 10^6$	running cims	-
RBF kernel	GC4-P4-FC1000	$2 \cdot 10^6$	67.85	-
RBF kernel	GC8-P8-FC1000	$2 \cdot 10^6$	66.95	-
RBF kernel	GC16-P4-GC16-P4-FC1000	$2 \cdot 10^6$	67.16	-
RBF kernel	GC64-P8-GC64-P8-FC1000	$2 \cdot 10^6$	running cims	-
RBF kernel (local)	GC4-P4-FC1000	$2 \cdot 10^6$	68.56	-
RBF kernel (local)	GC8-P8-FC1000	$2 \cdot 10^6$	67.66	-

## 5.2 Merck Molecular Activity Challenge

The Merck Molecular Activity Challenge is a computational biology benchmark where the task is to predict activity levels for various molecules based on the distances in bonds between different atoms.

324 For our experiments we used the DPP4 dataset which has 8193 samples and 2796 features. We chose  
 325 this dataset because it was one of the more challenging and was of relatively low dimensionality  
 326 which made the spectral networks tractable. As a baseline architecture, we used the network of [9]  
 327 which has 4 hidden layers and is regularized using dropout and weight decay.

328 As before, we log-normalized the data and used one-tenth of the training set to tune hyperparameters  
 329 of the network. For this task we found that  $k = 40$  subsampled weights worked best, and that average  
 330 pooling performed better than max pooling. Since the task is to predict a continuous variable, all  
 331 networks were trained by minimizing the Root Mean-Squared Error loss.  
 332

333 Table 2: Results for Merck DPP4 dataset  
 334

Graph	Architecture	Parameters	$R^2$
-	FC-4000-2000-1000-1000-50	$22.1 \cdot 10^6$	0.2728
Supervised	GC64-P8-GC64-P8-1000-1000-1	$3.8 \cdot 10^6$	0.2629
RBF Kernel	GC64-P8-GC64-P8-1000-1000-1	$3.8 \cdot 10^6$	0.1992
RBF Kernel (local)	GC64-P8-GC64-P8-1000-1000-1	$3.8 \cdot 10^6$	0.1487

### 342 5.3 ImageNet 343

344  
 345  
 346 In the experiments above our graph construction relied on an approximate estimation from the data.  
 347 To measure the influence of the graph construction compared to the filter learning in the frequency  
 348 domain, we performed the same experiments on the ImageNet dataset for which the graph is already  
 349 known, namely it is the 2-D grid. The spectral network was thus a convolutional network whose  
 350 weights were defined in the frequency domain. Training was performed exactly as in Figure 1,  
 351 except that the linear transformation was a Fast Fourier Transform.

352 Our network consisted of 4 convolution/ReLU/max pooling layers with 48, 128, 256 and 256 feature  
 353 maps, followed by 3 fully-connected layers each with 4096 hidden units regularized with dropout.  
 354 We trained two versions of the network: one classical convolutional network and one as a spectral  
 355 network where the weights were defined in the frequency domain only and were interpolated using  
 356 a spline kernel. Both networks were trained for 40 epochs over the ImageNet dataset where input  
 357 images were scaled down to  $128 \times 128$  to accelerate training.

359 Table 3: ImageNet results  
 360

Graph	Architecture	Test Accuracy (Top 5)	Test Accuracy (Top 1)
2-D Grid	Convolutional Network	71.854	46.24
2-D Grid	Spectral Network	71.998	46.71

364  
 365  
 366 We see that both models yield nearly identical performance. Interestingly, the spectral network learns  
 367 faster than the ConvNet during the first part of training, although both networks converge around the  
 368 same time. This requires further investigation.

369 Word datasets. Mention that Using Word embeddings is also a possibility for future work.  
 370

## 371 6 Discussion 372

373 The graph construction using the Supervised Graph Estimation. Related to discovering latent graph-  
 375 ical models

376 Limitations: When is weight sharing appropriate/ not.  
 377

When is locality appropriate/not appropriate.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

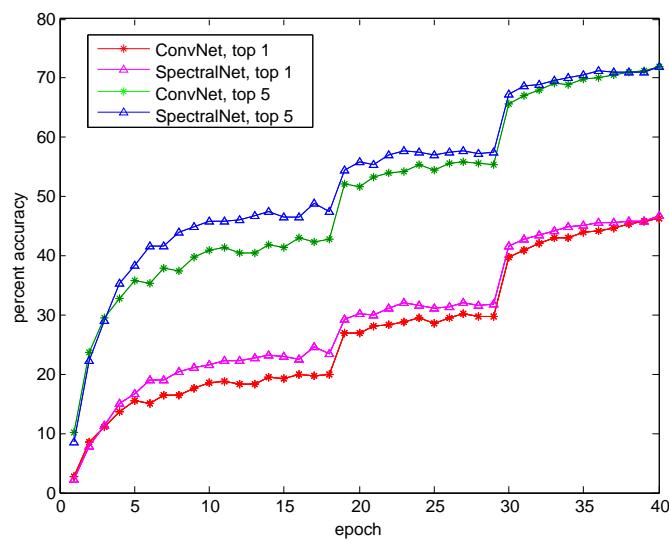


Figure 2: ConvNet vs. SpectralNet on ImageNet.

432 Complexity.

433  
434 Alternatives we are

435  
436 **References**

- 438 [1] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embed-  
439 ding and clustering. In *NIPS*, volume 14, pages 585–591, 2001.
- 440 [2] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and deep  
441 locally connected networks on graphs. In *Proceedings of the 2nd International Conference on*  
442 *Learning Representations*, 2013.
- 443 [3] Xu Chen, Xiuyuan Cheng, and Stéphane Mallat. Unsupervised deep haar scattering on graphs.  
444 In *Advances in Neural Information Processing Systems*, pages 1709–1717, 2014.
- 445 [4] Adam Coates and Andrew Y Ng. Selecting receptive fields in deep networks. In *Advances in*  
446 *Neural Information Processing Systems*, pages 2528–2536, 2011.
- 447 [5] Karol Gregor and Yann LeCun. Emergence of complex-like cells in a temporal product net-  
448 work with local receptive fields. *arXiv preprint arXiv:1006.0448*, 2010.
- 449 [6] Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel rahman Mohamed, Navdeep Jaitly,  
450 Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath, and Brian Kingsbury. Deep  
451 neural networks for acoustic modeling in speech recognition. *Signal Processing Magazine*,  
452 2012.
- 453 [7] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep  
454 convolutional neural networks. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger,  
455 editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran  
456 Associates, Inc., 2012.
- 457 [8] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–  
458 444, 05 2015.
- 459 [9] Junshui Ma, Robert P. Sheridan, Andy Liaw, George E. Dahl, and Vladimir Svetnik. Deep neu-  
460 ral networks as a method for quantitative structure-activity relationships. *Journal of Chemical*  
461 *Information and Modeling*, 2015.
- 462 [10] Stéphane Mallat. *A wavelet tour of signal processing*. Academic press, 1999.
- 463 [11] Jonathan Masci, Davide Boscaini, Michael M. Bronstein, and Pierre Vandergheynst. Shapenet:  
464 Convolutional neural networks on non-euclidean manifolds. *CoRR*, abs/1501.06297, 2015.
- 465 [12] Jiquan Ngiam, Zhenghao Chen, Daniel Chia, Pang W Koh, Quoc V Le, and Andrew Y Ng.  
466 Tiled convolutional neural networks. In *Advances in Neural Information Processing Systems*,  
467 pages 1279–1287, 2010.
- 468 [13] Pradeep Ravikumar, Martin J Wainwright, John D Lafferty, et al. High-dimensional ising  
469 model selection using ?1-regularized logistic regression. *The Annals of Statistics*, 38(3):1287–  
470 1319, 2010.
- 471 [14] Nicolas L Roux, Yoshua Bengio, Pascal Lamblin, Marc Joliveau, and Balázs Kégl. Learning  
472 the 2-d topology of images. In *Advances in Neural Information Processing Systems*, pages  
473 841–848, 2008.
- 474 [15] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.  
475 Dropout: A simple way to prevent neural networks from overfitting. *The Journal of*  
476 *Machine Learning Research*, 15(1):1929–1958, 2014.
- 477 [16] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.  
478 Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine*  
479 *Learning Research*, 15:1929–1958, 2014.
- 480 [17] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–  
481 416, 2007.
- 482 [18] Lihi Zelnik-Manor and Pietro Perona. Self-tuning spectral clustering. In *Advances in neural*  
483 *information processing systems*, pages 1601–1608, 2004.