

Evaluating ‘Graphical Perception’ with CNNs

Daniel Haehn, James Tompkin, and Hanspeter Pfister

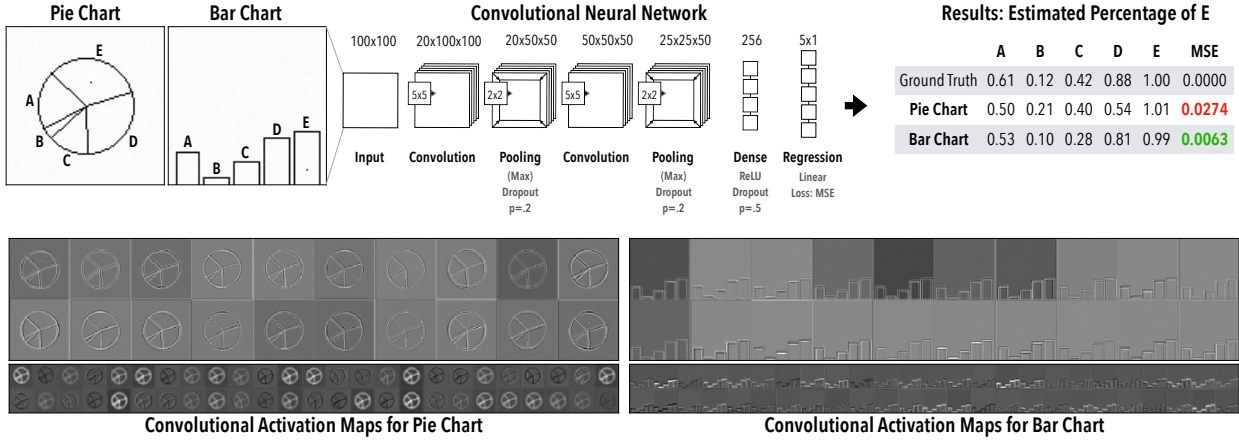


Fig. 1. **Computing Cleveland and McGill's Position-Angle Experiment using Convolutional Neural Networks.** We replicate the original experiment by asking visual cortex inspired machine learning classifiers to assess the relationship between values encoded in pie charts and bar charts. Similar to the findings of Cleveland and McGill [6], our experiments show that CNNs read quantities more accurately from bar charts (mean squared error, MSE in green).

Abstract—Convolutional neural networks can successfully perform many computer vision tasks on images, and their learned representations are often said to mimic the early layers of the visual cortex. But can CNNs understand graphical perception for visualization? We investigate this question by reproducing Cleveland and McGill's seminal 1984 experiments, which measured human perception efficiency of different visual encodings and defined elementary perceptual tasks for visualization. We measure the graphical perceptual capabilities of four classifiers on a) elementary perceptual tasks with increasing parametric complexity, b) the position-angle experiment that compares pie charts to bar charts, c) the position-length experiment that compares grouped and divided bar charts, and d) the bars and framed rectangles experiment where visual cues aid perception. We also study how feed-forward neural networks obey Weber's law, which defines the proportional relation between perceivable information and distribution density. We present the results of these experiments to foster the understanding of how CNN classifiers succeed and fail when applied to data visualizations.

Index Terms—Machine Perception, Deep Learning

1 INTRODUCTION

Deep multilayer neural networks are being successfully applied in a wide range of applications that are regularly outperforming humans in object recognition [16, 26, 27]. Originally inspired by neuroscientific discoveries, the recent advances in deep learning have been the direct results of engineering efforts, more specifically in convolutional neural networks (CNNs). Our current knowledge of biological vision suggests that CNNs indeed mimic feature learning similar to the early layers of the visual cortex but without incorporating the many details of biological neural networks [12, 20, 31]. While there has been significant advancement, fully understanding how CNNs work is a still ongoing process which does not detract from their success [9, 25] and recent works target the systematic evaluation of perceptual limits of CNNs [14, 23].

We are interested in applying neural networks to graphs, charts, and visual encodings. A first step in this direction is to evaluate if and how convolutional neural networks perceive low-level graphical elements, the building blocks of information visualizations. This is a very different task than applying such networks to natural images (like the majority of computer vision research) since such scenes are usually more complex and identifiers such as edges are not as prominent. More related towards our goal are the *graphical perception* experiments of Cleveland and McGill [6]. In this paper, we explore whether these experiments which were performed with humans can be reproduced with CNNs. Cleveland and McGill's work in the 1980s has led to many insights for modern information visualization research such as the identification of elementary perceptual tasks or that bar charts are easier to understand than pie charts.

In order to perform this evaluation, we parametrize different visual encodings such as the elementary perceptual tasks suggested by Cleveland and McGill [6]. We then replicate the original experimental design by defining linear regression tasks for continuous variables. However, we select three modern feature generators based on convolutional neural networks and combine them with a multilayer perceptron (MLP) to include non-linearities. We include the LeNet-5 network [19], the VGG19 network [26], and the Xception classifiers [5]. While we train LeNet-5 from scratch, for VGG19 and Xception we use pre-trained imagenet [17] weights to further mimic the human visual system. By also using the MLP directly without convolutional feature detection as baseline, we test four different classifiers as part of the following

- Daniel Haehn, and Hanspeter Pfister are with the Paulson School of Engineering and Applied Sciences at Harvard University.
E-mail: {haehn,pfister}@seas.harvard.edu.
- James Tompkin is with the Thomas J. Watson Sr. Center for Information Technology at Brown University.
E-mail: james_tompkin@brown.edu.

Manuscript received xx xxx. 201x; accepted xx xxx. 201x. Date of Publication xx xxx. 201x; date of current version xx xxx. 201x. For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org.
Digital Object Identifier: xx.xxx/TVCG.201x.xxxxxxx

scenarios: a) elementary perceptual tasks with increasing parameteric complexity, b) position-angle experiment comparing pie charts to bar charts, c) position-length experiment that compares grouped and divided bar charts, and d) the bars and framed rectangles experiments where visual cues aid perception. We also investigate other properties of CNNs such as whether they obey Weber’s law which defines a proportional dependency between an initial distribution and perceivable change. Our experimental setup includes repetitions, randomizations, and regularizations to prevent any bias.

Our motivation for replicating Cleveland and McGill’s experiments stems from the thought that computational perception seems to be closely related to biological vision. If human perception yields certain results, maybe we can replicate these results with machine vision. As our first contribution, we study the elementary perceptual tasks by Cleveland and McGill and then systematically parametrize and evaluate them for computational perception with our four classifiers. This setup includes also cross-network evaluations which give insight into the generalizability of the classifiers. It also yields our second contribution: a ranking of Cleveland and McGill’s elementary perceptual tasks for our tested CNN architectures. Further, we replicate the *position-angle* and the *position-length* experiments which contributes to the general knowledge of superior perceivability of bar charts to pie charts in certain conditions. We then reproduce the *bars-and-framed-rectangles* experiment with our classifiers. Here, we also include an additional experiment to test the Weber-Fechner’s law for CNNs. Our experiments yield a TODO. Finally, we discuss our findings and derive recommendations for allowing CNNs to perceive visualizations.

We accompany this paper and detailed supplemental material with open source code¹, data, and results to enable a framework for the development and evaluation of new network architectures for graphical perception.

2 PREVIOUS WORK

Graphical Perception. Cleveland and McGill [6] introduce the fundamental concept of *graphical perception* and investigate how different visual attributes and encodings are perceivable by humans. They define *elementary perceptual tasks* as mental-visual stimuli to understand encodings in visualizations. Based on these definitions, the authors propose and perform different experiments such as the *position-angle* experiment which compares bar charts and pie charts, the *position-length* experiment where users judge relations between encoded values in grouped and divided bar charts, and the *bars-and-framed-rectangles* experiment to evaluate Weber’s law. Heer and Bostock later reproduced the Cleveland-McGill experiments crowd-sourced on Mechanical Turk [13] which lead to follow-up work from Harrison *et al.* [10] who replicated the experiments while observing emotional states. Both papers report similar results to Cleveland and McGill which increased our motivation to mimick their pioneering work. Our experimental setup replicates the original setup of Cleveland and McGill - just instead of humans, we use convolutional neural networks due to the connection with the human visual system. While we focus on Cleveland and McGill’s work from 1984, many other excellent articles from the last decades target low-level visual encoding [1, 4, 7, 21, 28–30].

Interesting are also the rankings of correlation visualization using Weber’s law [11]. This law defines the proportional relation between the initial distribution density and perceivable change. In this paper, we investigate with a simple experiment whether this holds for convolutional neural networks.

Computational Visualization Understanding.

Pineo *et al.* [22] create computational model of human vision based on neural networks and their experiments show that understanding visualization triggers neural activity in high-level areas of cognition. The authors suspect that this level of understanding is produced by low-level neurons performing elementary perceptual tasks. We

¹Code, data, results and more are available at: <http://rhoana.org/perception>

are further investigating this suspicion. Other work tries to parse infographics by finding higher-level saliency models [3], or by extracting text or key visual elements [2, 15, 24]. However, none of these works focus on computational understanding of lower-level building blocks of visualizations such as curvature, lengths, or position.

Visual Cortex Inspired Machine Learning. The human visual cortex is an extremely powerful system which allows the ability, and seemingly without effort, to recognize an enormous amount of distinct objects in the world. This visual system is organized in layers and has inspired the theory of computational classifiers based on multilayer neural networks. Fukushima and Miyake developed the Neocognitron quantitative model [8] that ultimately led to the important work of Hinton, Bengio, and LeCun: *deep neural networks* [18], visual cortex inspired machine learning. Nowadays, such classifiers exist with many different architectures. For this paper, we select the traditional *LeNet-5* [19] which was designed to recognize hand-written digits, the VGG19 [26] classifier with 16 convolutional layers, and the Xception [5] classifier with 36 convolutional layers. Selecting these specific networks allows us to compare architectures with different depths.

3 EXPERIMENTAL SETUP

The experiments shown in this paper are either supervised regression or classification tasks. We formulate any estimation of quantities (e.g. angles, positions, lengths etc.) as a regression problem between 0 and 1. The output indicates the percentage in regards to the degrees of freedom of the individual experiment. If the experiment involves a choice, we formulate it as a classification problem.

3.1 Measures

Accuracy. We use the same metric as Cleveland and McGill to measure accuracy.

$$\log_2(|\text{predicted percent} - \text{true percent}| + .125) \quad (1)$$

Confidence Intervals. We follow the notion of Cleveland and McGill to compute the confidence intervals.

Efficiency. We use the convergence rate based on the decrease of loss per training epoch as an indicator for the efficiency of the classifier in combination with a visual encoding. For regression tasks the loss is defined as mean squared error (MSE) and for classification tasks the loss is categorical cross-entropy.

3.2 Classifiers

Our classifiers are built upon a multilayer perceptron (MLP) which is a feedforward artificial neural network. We combine this MLP with different convolutional neural networks (CNNs) for preprocessing and feature generation. These include the traditional LeNet trained from scratch, as well as VGG19 and Xception trained using ImageNet.

Multilayer Perceptron. The multilayer perceptron in this paper has 256 neurons which are activated as rectified linear units (Fig. 2). We then add a dropout layer to prevent overfitting and compute linear regression or classification (softmax).

Convolutional Neural Networks. We use CNNs to generate additional features as input to the MLP. We train the *LeNet* classifier with tune it specifically towards each visualization. For *VGG19* and *Xception*, we generate features using previously trained weights on ImageNet.

Optimization. All networks are optimized using stochastic gradient descent with Nesterov momentum using fixed parameters (Table 1). We train for 1000 epochs but stop early if the loss does not decrease for ten epochs.

Table 1. We use different feature generators as input to a multilayer perceptron which performs linear regression or the classification task. This yields different sets of trainable parameters. We also train the MLP directly on the visualizations without any additional feature generation.

Classifier	Trainable Parameters	Optimization
MLP	2,560,513	SGD (Nesterov momentum)
LeNet + MLP	8,026,083	Learning rate: 0.0001
VGG19 + MLP	21,204,545	Momentum: 0.9
Xception + MLP	25,580,585	Batchsize: 32
		Epochs: 1000 (Early Stopping)

Feature Generation	Multilayer Perceptron
<p>LeNet</p>	<p>Dense ReLU Dropout (p=.5) Regression Linear, Loss: MSE</p>
<p>VGG19</p>	
<p>Xception</p>	

Fig. 2. The multilayer perceptron (MLP) in our experiments has 256 neurons which are activated as rectified linear units (ReLU). We use Dropout regularization to prevent overfitting. We learn categorical and unordered dependent variables using the softmax function and perform linear regression for continuous variables. The MLP can learn the visualizations directly but we also learn features generated by LeNet (2 conv. layers, filter size 5), VGG19 trained on ImageNet (16 conv. layers, filter size 3×3), or Xception trained on ImageNet (36 conv. layers, filter size 3×3) to increase the number of trainable parameters.

Environment. We run all experiments on an NVIDIA DGX1 machine with Tesla V100 graphical processing units. We use the KERAS framework with tensorflow.

3.3 Data

We create all visualizations as parametrized rasterized images without interpolation. The number of parameters differs per experiment as summarized in Table 2 and section 4.1. We add subtle random noise (0.05) to each pixel to introduce additional variation.

Training/Validation/Test Splits. We specify the size of each split set as follows: 60,000 training images, 20,000 validation images, and 20,000 test images. We then randomly add parameterized visualizations to the sets while guaranteeing that each set is disjunct from each other in terms of encoded variables. This eliminates leakage during training and evaluation. We also scale each set independently: images to the range of $-.5$ to $.5$ and labels to the range of 0.0 to 1.0 .

Cross-classifier variability. We also evaluate classifiers previously trained with one visualization on the same type of visualizations with different parameters by decreasing and increasing the variability of the generated images.

4 ELEMENTARY PERCEPTUAL TASKS

Cleveland and McGill describe the mapping of graphical elements to quantitative variables as *elementary perceptual tasks* and introduce a list of ten different encodings in their paper [6]. We create visualizations of these tasks as rasterized images (Fig. 3).

4.1 Parametrizations

We generate multiple parameterizations for each elementary perceptual task (Fig. 3) and sequentially increase the number of parameters. For instance, for *position non-aligned scale* we first only vary the origin of the coordinate system which yields just 10 different parameters. We then include translation along the y-axis with a significant increase in variability. We then also add x-movement and a variable spot size. This

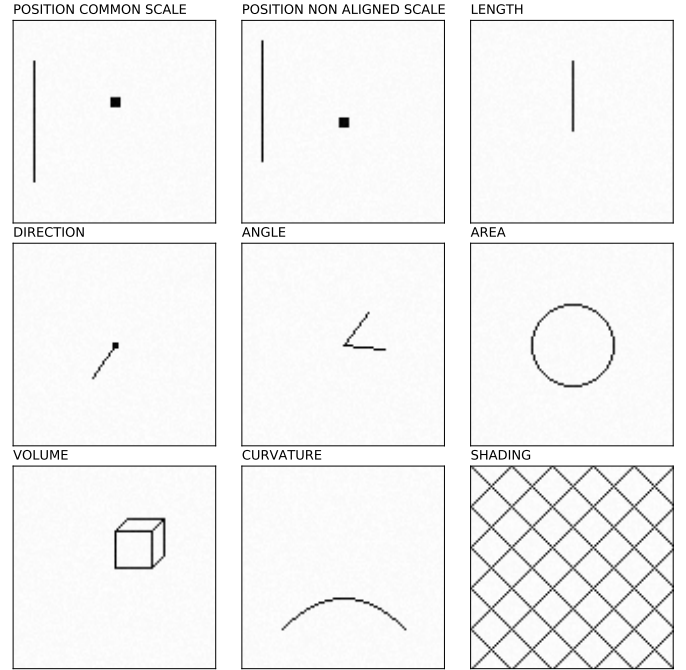


Fig. 3. **Elementary Perceptual Tasks.** Rasterized visualizations of the elementary perceptual tasks as defined by Cleveland and McGill [6] (color saturation excluded). We vary the parameters of each perceptual task and then assess the interpretability of feed-forward neural networks.

results in more complex datasets depending on the variability setting. Table 2 shows the different settings. It is important to consider this variability when evaluating different classifiers with individual trainable parameters (Table 1).

4.2 Hypotheses

We proposed four hypotheses entering the elementary perceptual task experiment:

- **H1.1 Visual cortex inspired classifiers are able to connect graphical elements to their quantitative variables.** While much simpler models than their biological pendant, convolutional neural networks are heavily influenced by our biological knowledge of the visual system. Such classifiers therefore follow the same principles as human perception.
- **H1.2 Computed perceptual performance is dependent on classifier complexity.** We evaluate multiple classifiers with different numbers of trainable parameters. A more complex classifier (with higher number of parameters) will perform better on elementary perceptual tasks.
- **H1.3 Some visual encodings are better than others for computations.** Cleveland and McGill order the elementary perceptual tasks by accuracy. We investigate whether this order is also relevant for computing graphical perception.
- **H1.4 Classifiers trained on perceptual tasks can generalize to more or less complex variations of the same task.** Recent research suggests that convolutional neural networks generalize extremely well. While the underlying reasons are mainly yet unknown, this property allows them to perform on variations of a similar perceptual task.

5 POSITION-ANGLE EXPERIMENT

The position-angle experiment was originally performed by Cleveland and McGill to measure whether humans can better perceive quantities encoded as positions or as angles [6]. The actual experiment then

Table 2. **Variability of Elementary Perceptual Tasks.** We sequentially increase the number of parameters for every visual encoding of the elementary perceptual tasks. This introduces variability and increasingly more complex datasets.

Elementary Perceptual Task	Variability	Parameters
<i>Position Common Scale</i>	Position Y	60
	+ Position X	3600
	+ Spot Size	21600
<i>Position Non-Aligned Scale</i>	Position Y	600
	+ Position X	36000
	+ Spot Size	216000
<i>Length</i>	Length	60
	+ Position Y	2400
	+ Position X	144000
	+ Width	864000
<i>Direction</i>	Angle	360
	+ Position Y	21600
	+ Position X	1296000
<i>Angle</i>	Angle	90
	+ Position Y	5400
	+ Position X	324000
<i>Area</i>	Radius	40
	+ Position Y	800
	+ Position X	16000
<i>Volume</i>	Cube Sidelength	20
	+ Position Y	400
	+ Position X	8000
<i>Curvature</i>	Midpoint Curvature	80
	+ Position Y	1600
	+ Position X	64000
<i>Shading</i>	Density	100
	+ Position Y	2000
	+ Position X	40000

compares pie charts versus bar charts since these map down to elementary position and angle judgement. We create rasterized images mimicking Cleveland and McGill’s proposed encoding and investigate computational perception of our four classifiers.

5.1 Hypotheses

We proposed four hypotheses entering the elementary perceptual task experiment:

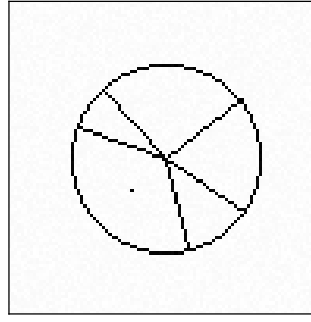
- **H2.1 Computed perceptual performance is better using bar charts than pie charts.** Cleveland and McGill report that position judgements are almost twice as accurate as angle judgements. This renders bar charts superior to pie charts and should also be the case for convolutional neural networks.
- **H2.2 Classifiers can learn position faster than angles.** We assume that understanding bar charts is easier than understanding pie charts. We suspect that our classifiers learn encodings of positions faster than of angles resulting in more efficient training and faster convergence.

6 POSITION-LENGTH EXPERIMENT

This is the one where we estimate two selected bars compared to the longest one - very similar to the previous one but, not yet done. We basically test divided versus grouped bar chart and we estimate one relation between two marked quantities: what percent the smaller is of the larger.

There are five types: type 1-3 this is a position judgement along a common scale. (btw all classifiers seem to do that extremely well in the

PIE CHART



BAR CHART

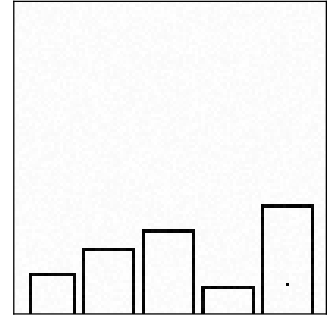


Fig. 4. **Position-Angle Experiment.** We create rasterized visualizations of pie charts and bar charts to follow Cleveland and McGill’s position-angle experiment. The experimental task involves the judgement of different encoded values in comparison to the largest encoded values. The pie chart (left) and the bar chart (right) visualize the same data point. In their paper, Cleveland and McGill report less errors using bar charts.

elementary tasks so we assume this will work well here too). Types 4-5 are length judgements and we know that the classifiers struggle with that quite a bit.

The setup from Cleveland McGill is first a classification task: which one is smaller? and then a regression task: how much smaller. So we have to see how to encode this.

6.1 Hypotheses

We proposed two hypotheses entering the elementary perceptual task experiment:

- **H3.1 Grouped bar charts are better computational perceivable than divided bar charts.** A grouped bar chart involves judging a position while a divided bar chart most likely (if not the bottom is looked at) requires length judgements. Classifiers are better at judging position than at judging length so grouped bar charts are easier to grasp in terms of computational perception.
- **H3.2 not yet** Any ideas?

6.2 Discussion

JT: Look at the relative difficulty of the tasks. In Cleveland and McGill, types 1-5 were post-ordered by their log error such that type 1 was easiest and type 5 was hardest. Is this still the case with our CNNs?

7 BARS AND FRAMED RECTANGLES EXPERIMENT

Visual cues can help converting graphical elements back to their real world variables. Cleveland and McGill introduced the bars and framed rectangles experiment which judges the elementary perceptual task of position along non-aligned scales [6].

7.1 Hypotheses

We proposed two hypotheses entering the elementary perceptual task experiment:

- **H4.1 Classifiers can leverage additional visual cues.** The original bar and framed rectangle experiment shows how visual cues aid humans in mapping graphical elements to quantitative variables. This should be the same for feed-forward neural networks since they are inspired by the visual system.
- **H4.2 Weber’s law can be transferred to computational perception.** Cleveland and McGill confirmed Weber’s law based on the bar and framed rectangle experiment. For humans, the ability to perceive change within a distribution is proportional to the size of the initial distribution.

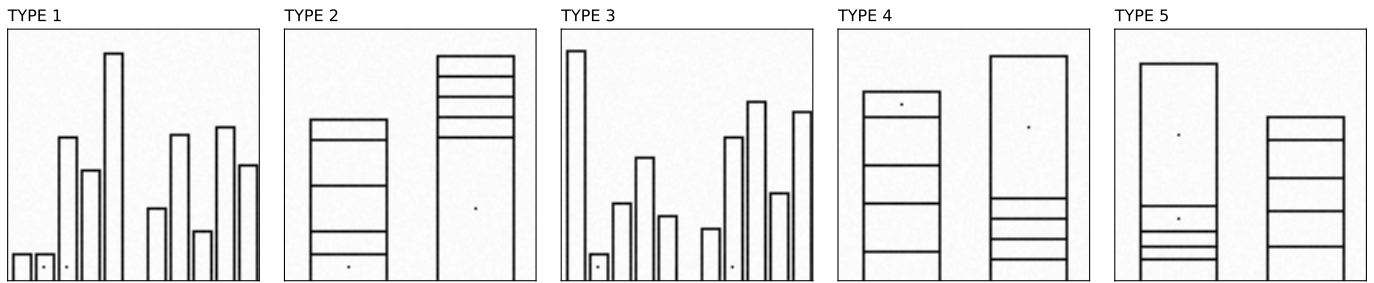


Fig. 5. **Position-Length Experiment.** (Not yet) Rasterized versions of the graphs of Cleveland and McGill's position-length experiment. The perceptual task involves comparing the two dot-marked quantities across five different visual encodings of either grouped or divided bar charts. We evaluate which type of bar chart performs better with our neural networks as a combined classification and regression problem. The first task is to select which of the marked quantities is smaller (classification) and the second task is to specify how much smaller it is (regression).

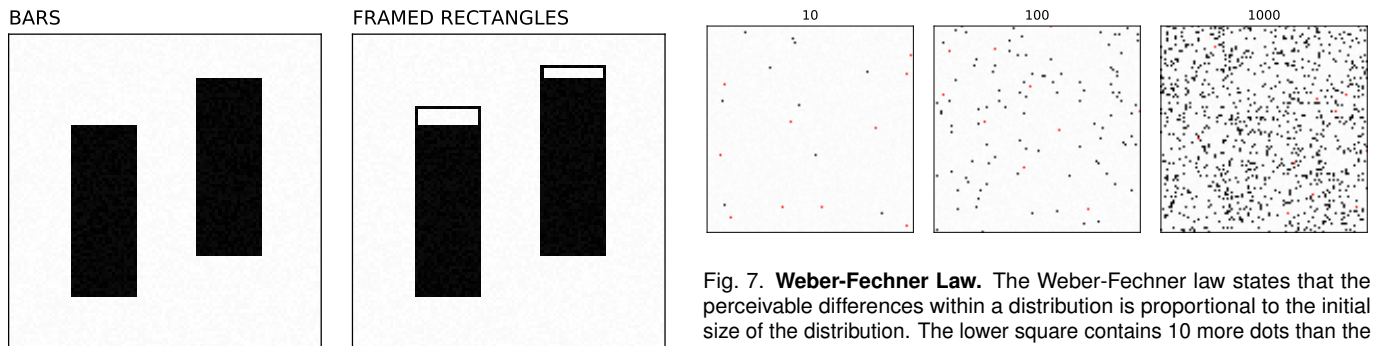


Fig. 6. **Bars and Framed Rectangles Experiment.** Cleveland and McGill introduce the bars and framed rectangles experiment which measures the perceptual task of judging position along non-aligned scales. For humans, it is easier to decide which of two bars represent a larger height if a scale is introduced by adding framed rectangles (right). In this case, the right bar is higher as visible with less free space when adding the frame. We evaluate whether such a visual aid also helps machines to perceive visually encoded quantities.

7.2 Weber-Fechner's Law

As identified by Cleveland and McGill, the bar and framed rectangle experiment is closely related to Weber's law. This psychophysics law states that perceivable difference within a distribution is proportional to the initial size of the distribution. Weber's law goes hand-in-hand with Fechner's law. We conduct an additional experiment based on the original illustrations of the Weber-Fechner law to investigate whether this law can be applied to computational perception of our classifiers (Fig. 7).

8 RESULTS AND DISCUSSION

8.1 Elementary Perceptual Tasks

some are good and some are bad.. why?

Computational Perception Ranking.

Cleveland McGills Ranking - can we observe something similar?

1. Position along a common scale e.g. scatter plot
2. Position on identical but nonaligned scales e.g. multiple scatter plots
3. Length e.g. bar chart
4. Angle & Slope (tie) e.g. pie chart
5. Area e.g. bubbles
6. Volume, density, and color saturation (tie) e.g. heatmap

Fig. 7. **Weber-Fechner Law.** The Weber-Fechner law states that the perceivable differences within a distribution is proportional to the initial size of the distribution. The lower square contains 10 more dots than the upper one on both sides. However, the difference is easily perceivable on the left while the squares on the right almost look the same. We generate rasterized visualizations similar to this setup and evaluate our classifiers.

7. Color hue e.g. newsmag

Cross-classifier variability.

Can a neural network generalize on simple perceptual tasks?

8.2 Position-Angle Experiment

Bar charts are more accurate (Fig. 11) and networks converge faster (Fig. 10). This is great.

8.3 Position-Length Experiment

8.4 Bars and Framed Rectangles Experiment

First run indicates that framed rectangles perform better but we dont really know it yet.

9 CONCLUSIONS

Future work: allow insights for infovis for machines

REFERENCES

- [1] J. Bertin and M. Barbut. *Semiologie graphique: les diagrammes, les reseaux, les cartes*. Mouton, 1967.
- [2] Z. Bylinskii, S. Alsheikh, S. Madan, A. Recasens, K. Zhong, H. Pfister, F. Durand, and A. Oliva. Understanding infographics through textual and visual tag prediction. *CoRR*, abs/1709.09215, 2017.
- [3] Z. Bylinskii, A. Recasens, A. Borji, A. Oliva, A. Torralba, and F. Durand. Where should saliency models look next? In *European Conference on Computer Vision*, pp. 809–824. Springer, 2016.
- [4] M. Carpendale. Considering visual variables as a basis for information visualisation. 2003.
- [5] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, pp. 1800–1807. IEEE Computer Society, 2017.
- [6] W. S. Cleveland and R. McGill. Graphical perception: Theory, experimentation, and application to the development of graphical methods. *Journal of the American statistical association*, 79(387):531–554, 1984.
- [7] W. S. Cleveland and R. McGill. Graphical perception and graphical methods for analyzing scientific data. *Science*, 229(4716):828–833, 1985.

- [8] K. Fukushima and S. Miyake. Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition. In *Competition and cooperation in neural nets*, pp. 267–285. Springer, 1982.
- [9] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [10] L. Harrison, D. Skau, S. Franconeri, A. Lu, and R. Chang. Influencing visual judgment through affective priming. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 2949–2958. ACM, 2013.
- [11] L. Harrison, F. Yang, S. Franconeri, and R. Chang. Ranking visualizations of correlation using weber’s law. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):1943–1952, Dec 2014. doi: 10.1109/TVCG.2014.2346979
- [12] D. Hassabis, D. Kumaran, C. Summerfield, and M. Botvinick. Neuroscience-inspired artificial intelligence. *Neuron*, 95(2):245–258, 2017.
- [13] J. Heer and M. Bostock. Crowdsourcing graphical perception: using mechanical turk to assess visualization design. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pp. 203–212. ACM, 2010.
- [14] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. L. Zitnick, and R. B. Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. *CoRR*, abs/1612.06890, 2016.
- [15] A. Kembhavi, M. Salvato, E. Kolve, M. Seo, H. Hajishirzi, and A. Farhadi. A diagram is worth a dozen images. In *European Conference on Computer Vision*, pp. 235–251. Springer, 2016.
- [16] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, eds., *Advances in Neural Information Processing Systems 25*, pp. 1097–1105. Curran Associates, Inc., 2012.
- [18] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [19] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998. doi: 10.1109/5.726791
- [20] D. Linsley, S. Eberhardt, T. Sharma, P. Gupta, and T. Serre. What are the visual features underlying human versus machine vision? *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pp. 2706–2714, 2017.
- [21] T. Munzner. *Visualization Analysis and Design*. AK Peters Visualization Series. CRC Press, 2015.
- [22] D. Pineo and C. Ware. Data visualization optimization via computational modeling of perception. *IEEE Transactions on Visualization and Computer Graphics*, 18(2):309–320, Feb 2012. doi: 10.1109/TVCG.2011.52
- [23] M. Ricci, J. Kim, and T. Serre. Not-So-CLEVR: Visual Relations Strain Feedforward Neural Networks. *ArXiv e-prints*, Feb. 2018.
- [24] M. J. Seo, H. Hajishirzi, A. Farhadi, and O. Etzioni. Diagram understanding in geometry questions. In *AAAI*, pp. 2831–2838. AAAI Press, 2014.
- [25] R. Shwartz-Ziv and N. Tishby. Opening the black box of deep neural networks via information. *CoRR*, abs/1703.00810, 2017.
- [26] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [27] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015.
- [28] A. Treisman and S. Gormican. Feature analysis in early vision: evidence from search asymmetries. *Psychological review*, 95(1):15, 1988.
- [29] D. Wigdor, C. Shen, C. Forlines, and R. Balakrishnan. Perception of elementary graphical elements in tabletop and multi-surface environments. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI ’07, pp. 473–482. ACM, New York, NY, USA, 2007. doi: 10.1145/1240624.1240701
- [30] L. Wilkinson. *The grammar of graphics*. Springer Science & Business Media, 2006.
- [31] D. L. Yamins and J. J. DiCarlo. Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*, 19(3):356, 2016.

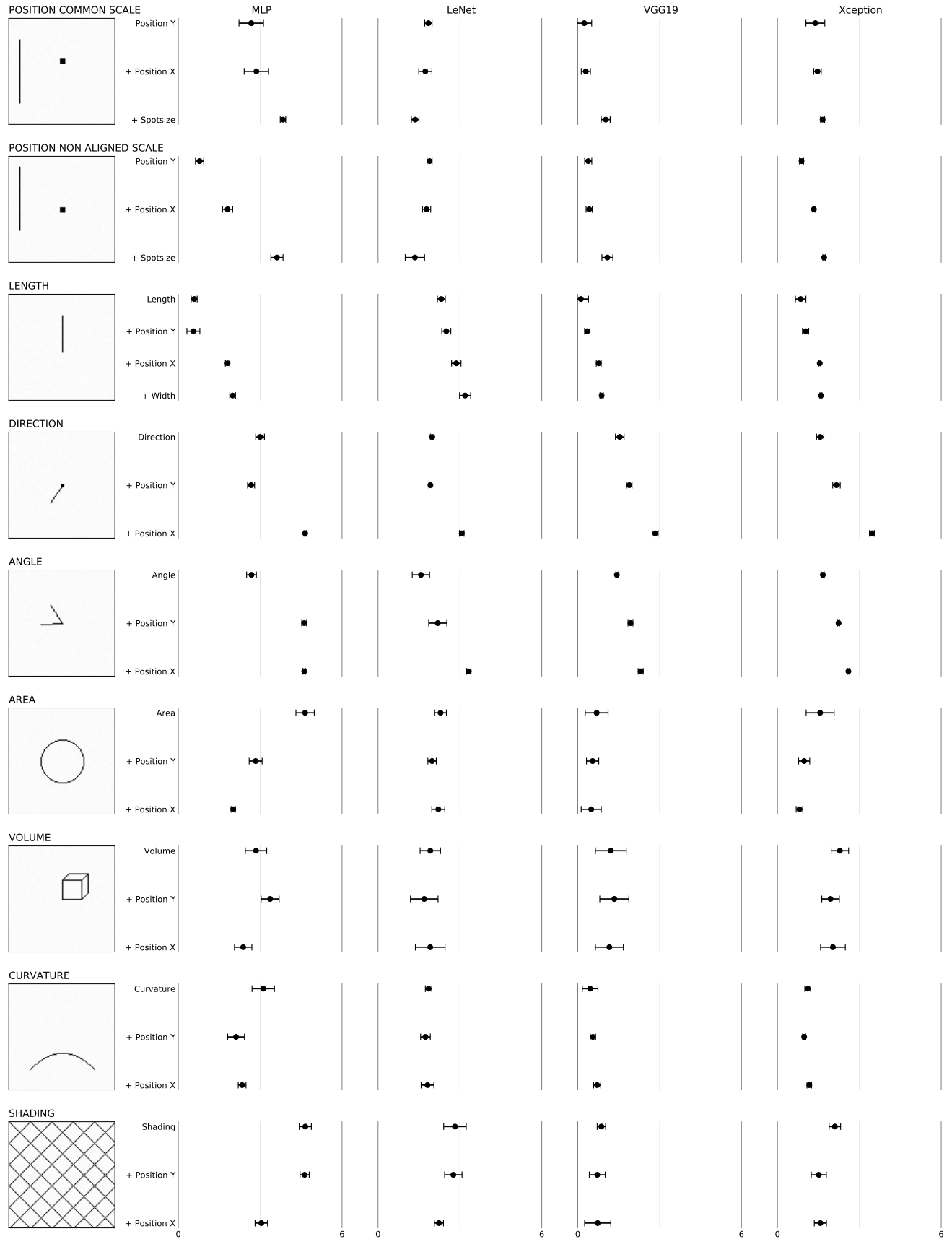


Fig. 8. **Computational results of Elementary Perceptual Tasks experiment.** Log absolute error means and 95% confidence intervals for computed perception of different classifiers on the *elementary perceptual tasks* introduced by Cleveland and McGill 1984 [6]. We test the performance of a Multi-layer Perceptron (MLP), the LeNet Convolutional Neural Network, as well as feature generation using the VGG19 and Xception networks trained on ImageNet.

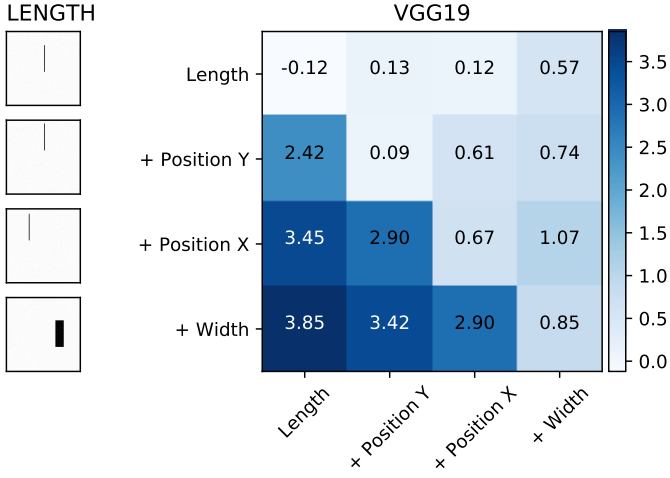


Fig. 9. **Cross-classifier variability for the perceptual task of measuring length.** We use predictions of LeNet classifiers trained on different parametrizations of the *curvature* elementary perceptual task and measure the mean logistic absolute error (MLAE). The lower score, the better. Classifiers trained on curves with variable position can generalize even if the axis of translation varies. However, classifiers trained on fixed positions of curves are not able to measure translated curves.

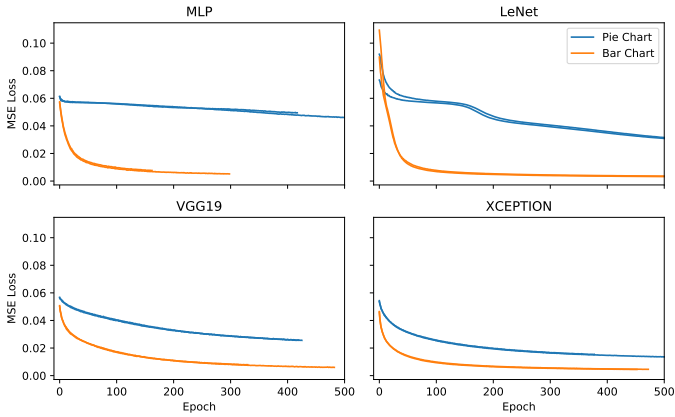


Fig. 10. **Classifier Efficiency of the Position-Angle experiment.** Mean Square Error (MSE) loss for the *position-angle experiment* as described by Cleveland and McGill [6] which compares the visualization of pie charts and bar charts. We report the MSE measure for both encodings of four different classifier on previously unseen validation data.

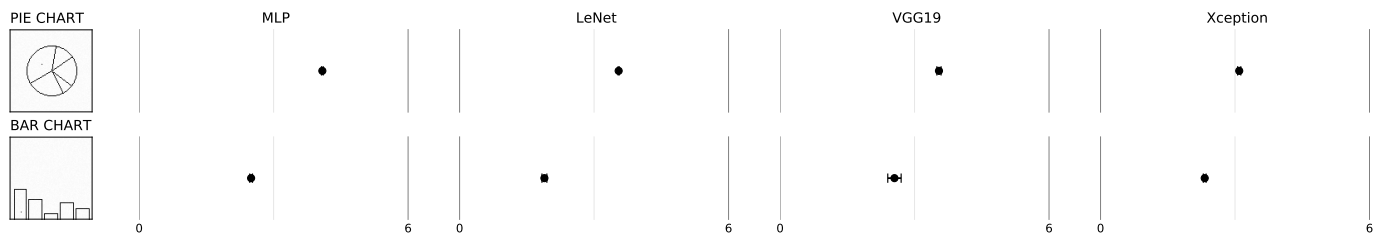


Fig. 11. **Computational results of the Position-Angle experiment.** Log absolute error means and 95% confidence intervals for the *position-angle experiment* as described by Cleveland and McGill [6]. We test the performance of a Multi-layer Perceptron (MLP), the LeNet Convolutional Neural Network, as well as feature generation using the VGG19 and Xception networks trained on ImageNet.

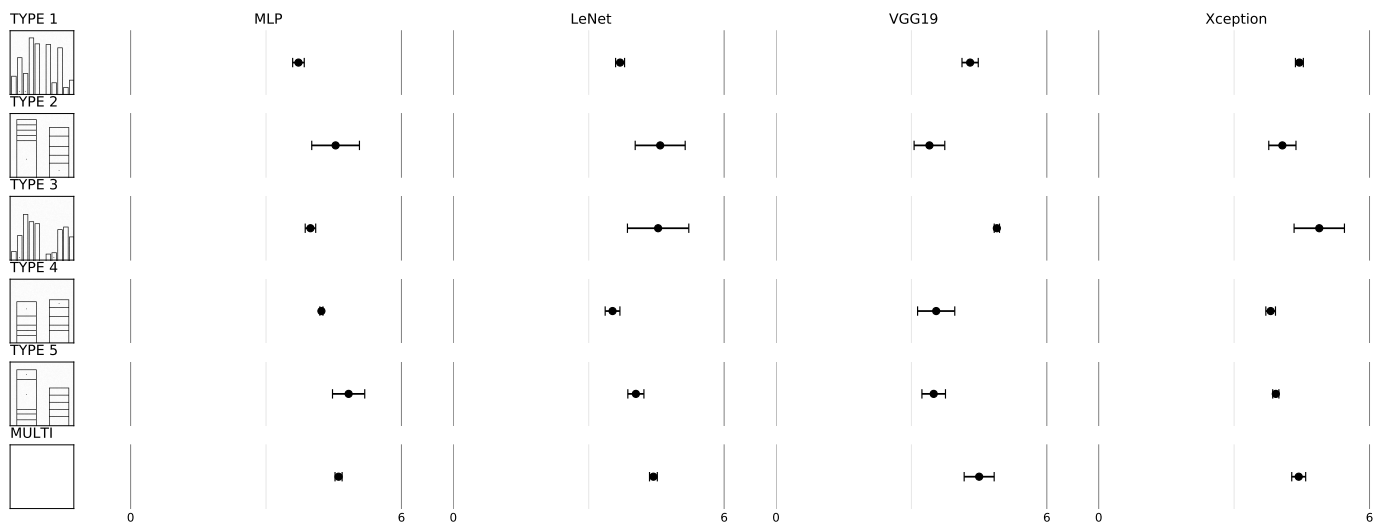


Fig. 12. **Computational results of the Position-Length experiment.** Log absolute error means and 95% confidence intervals for the *position-length experiment* as described by Cleveland and McGill [6]. We test the performance of a Multi-layer Perceptron (MLP), the LeNet Convolutional Neural Network, as well as feature generation using the VGG19 and Xception networks trained on ImageNet.

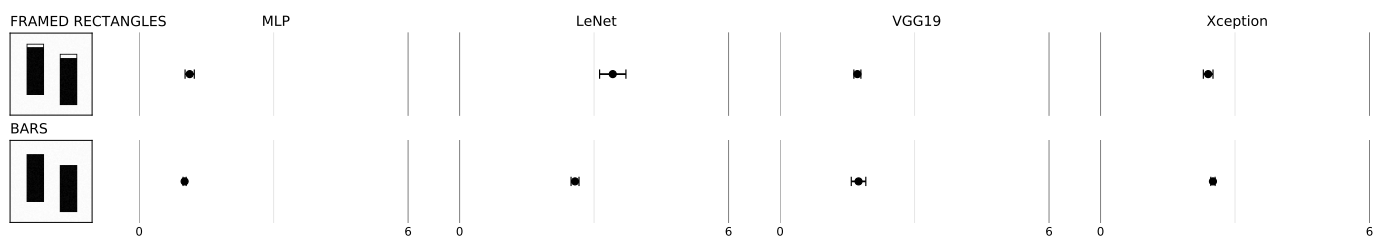


Fig. 13. **Computational results of the Bars-and-Framed-Rectangles experiment.** Log absolute error means and 95% confidence intervals for the *bars-and-framed-rectangles experiment* as described by Cleveland and McGill [6]. We test the performance of a Multi-layer Perceptron (MLP), the LeNet Convolutional Neural Network, as well as feature generation using the VGG19 and Xception networks trained on ImageNet.