
Naive Few-Shot Learning: Uncovering the fluid intelligence of machines

Tomer Barak¹ Yonatan Loewenstein^{1,2}

Abstract

In this paper, we aimed to help bridge the gap between human fluid intelligence - the ability to solve novel tasks without prior training - and the performance of deep neural networks, which typically require extensive prior training. An essential cognitive component for solving intelligence tests, which in humans are used to measure fluid intelligence, is the ability to identify regularities in sequences. This motivated us to construct a benchmark task, which we term *sequence consistency evaluation* (SCE), whose solution requires the ability to identify regularities in sequences. Given the proven capabilities of deep networks, their ability to solve such tasks after extensive training is expected. Surprisingly, however, we show that naive (randomly initialized) deep learning models that are trained on a *single* SCE with a *single* optimization step can still solve non-trivial versions of the task relatively well. We extend our findings to solve, without any prior training, real-world anomaly detection tasks in the visual and auditory modalities. These results demonstrate the fluid-intelligent computational capabilities of deep networks. We discuss the implications of our work for constructing fluid-intelligent machines.

1. Introduction

It has been demonstrated that deep learning models can successfully solve intelligence tests (Santoro et al., 2017; Barrett et al., 2018; Zhuo & Kankanhalli, 2020; Kim et al., 2020; Webb et al., 2022). However, they required extensive prior training to achieve this goal. Are deep learning models capable of exhibiting “fluid intelligence” that does not rely

¹The Edmond and Lily Safra Center for Brain Sciences, The Hebrew University, Jerusalem ²Department of Cognitive Sciences, The Federmann Center for the Study of Rationality, The Alexander Silberman Institute of Life Sciences, The Hebrew University, Jerusalem. Correspondence to: Tomer Barak <tomer.barak@mail.huji.ac.il>.

on prior training?

To address this question, we focused on a specific sub-task shared among presumably different intelligence tests, the extraction of simple rules from sequences, (Sternberg, 1977; 1983; Lohman, 2000; Siebers et al., 2015). For example, when solving a Raven’s Progression Matrix (Raven et al., 1998), humans extract the rules governing the change in the matrix’s rows and columns, and then use these rules to select the most consistent answer (Carpenter et al., 1990). Our focus here is the extent to which *naive* models can solve this computational task, extracting simple rules from sequences of inputs *without prior training*. This approach is an unusual setting for deep learning models, in which pretraining is considered crucial even in the context of few-shot learning (Chollet, 2019; Vogelstein et al., 2022).

We will start by describing the SCE, a task designed to study rule extraction. Then we will introduce Contrastive Predictive Coding (CPC) and Relation Network (RN) models and compare their performances on the SCE task. Our main result is that CPC can successfully solve non-trivial versions of the task by a single optimization step (starting with random weights) over a single sequence of 5 images. We conclude by demonstrating the applicability of our approach to real-world problems using two anomaly detection tasks.

2. Methods

2.1. Sequence consistency evaluation (SCE) tests

Each SCE test¹ is a sequence of $K = 5$ gray-scale images \mathbf{x}_j and $n = 4$ optional-choice images (Fig. 1). Each image includes 1-9 identical objects arranged on a 3×3 grid. An image is characterized by a low-dimensional vector of features, \mathbf{f}_j , where f_j^i denotes the value of feature i in image j . We use the following five features: the number of objects in an image (possible values: 1 to 9), their color (6 linearly distributed gray-scale values), the shapes (circle, triangle, square, star, hexagon), their size (6 linearly distributed values for the shapes’ enclosing circle circumference), and positions (a vector of grid positions that was used to place

¹Code is available in <https://github.com/Tomer-Barak/Naive-Few-Shot-Learning>

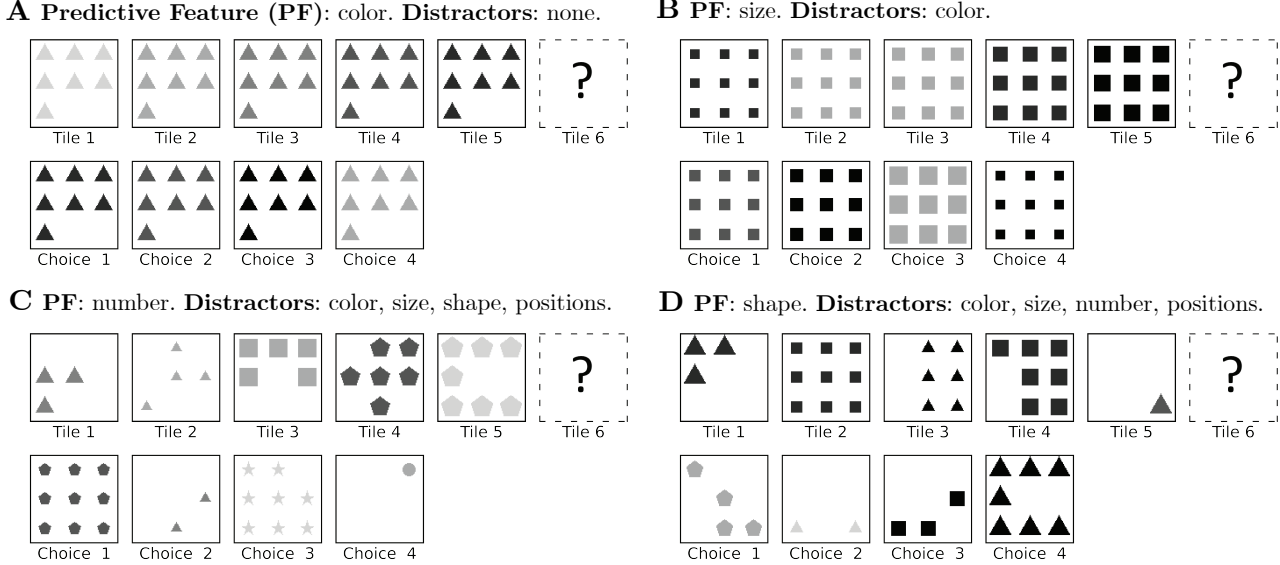


Figure 1. SCE Tests. The predictive features can be the (A) Color, (B) Size, (C) Number, which increases monotonically, or the (D) Shape of the objects, which alternates between a triangle and a square. Given a predictive feature, the rest of the features are either constant or random. We refer to the random features as *distractors*, and their number determines the test difficulty. *The correct choice in all of the tests above is 3.*

the shapes in order). An image \mathbf{x}_j is constructed according to its characterizing features by a non-linear and complex generative function $\mathbf{x}_j = G(\mathbf{f}_j)$.

One of the features f^p predictably changes along the sequence according to a simple deterministic rule $f_{j+1}^p = U(f_j^p)$ while the other features are either constant over the images or change randomly (values are i.i.d). We refer to the randomly-changing features as *distractors*, and their number is considered a measure of the difficulty of the test. Given a sequence of K images, an agent’s task is to select the correct $K + 1^{\text{th}}$ image from the set of n optional choice images that are generated using the same generative function G from the feature space. In the correct choice, f^p follows the deterministic rule $f_{K+1}^p = U(f_K^p)$, whereas in the incorrect choices it does not follow that rule and is instead randomly chosen from the remaining possible values. The features that are constant or randomly changing in the sequence are also constant or change randomly in all optional choice images.

2.2. Abstract relations models

Images in the SCE task are related to each other by abstract relations. These abstract relations are between low-dimensional latent representations of the images. For example, two images \mathbf{x}_i and \mathbf{x}_j with the same number of shapes are related by their latent variables $Z(\mathbf{x}_i)$ and $Z(\mathbf{x}_j)$, which encode their number of shapes. Specifically, in their case,

$Z(\mathbf{x}_i) = Z(\mathbf{x}_j)$. To account for relations more complex than equality and allow for these relations to be learned empirically, we used a learnable abstract relation function between images \mathbf{x}_i and \mathbf{x}_j ,

$$R_\theta(Z_\phi(\mathbf{x}_i), Z_\phi(\mathbf{x}_j)),$$

where θ and ϕ are the parameters of artificial neural networks R and Z , respectively.

We chose to construct a network that can identify the relation between consecutive sequence images in the sense that it can correctly identify the $K + 1$ image out of the n choice images. We compared two main candidates of deep learning models: Markov Contrastive Predictive Coding (Markov-CPC) (Oord et al., 2018) and Relation Network (RN) (Sung et al., 2018).

2.2.1. MARKOV CONTRASTIVE PREDICTIVE CODING

The Markov-CPC model (Fig. 2 left) has an inductive bias that assumes a causal predictive structure between inputs. It, therefore, uses a predictor function T_θ to predict a latent variable $Z_\phi(\mathbf{x}_j)$ given another latent variable $Z_\phi(\mathbf{x}_i)$. This manifests in the prediction error between the two variables, defined as,

$$\epsilon_{i,j}(Z_\phi, T_\theta) = \left(T_\theta(Z_\phi(\mathbf{x}_i)) - Z_\phi(\mathbf{x}_j) \right)^2.$$

To solve an SCE test, we used the model to reduce the predic-

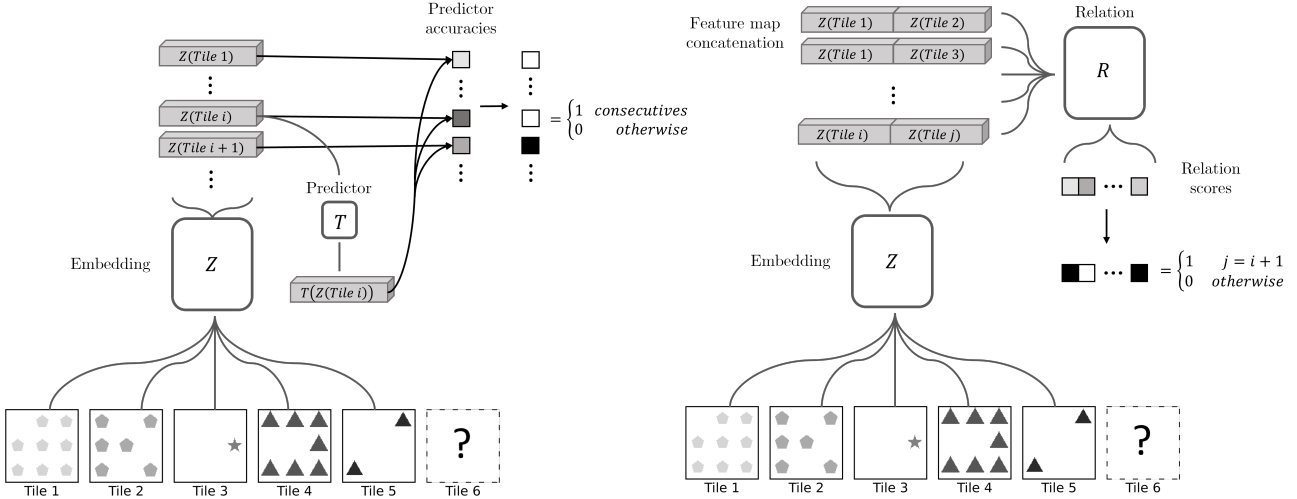


Figure 2. Comparison between Markov-CPC (left) and RN (right). The main difference between the models is that Markov-CPC imposes a *causal structure* between consecutive latent representations, while RN can learn a *general relation* R between two latent representations.

tion error between the latent variables of consecutive inputs. By construction, for an encoder and predictor functions that match the generative function and the deterministic rule of the SCE test, namely $Z^* = G^{-1}$ and $T^* = U$, the prediction error $\epsilon_{i,j}(Z^*, T^*) = 0$ if i and j are two consecutive images ($j = i + 1$), and $\epsilon_{i,j}(Z^*, T^*) > 0$ otherwise.

The challenge is that Z^* and T^* are unknown. However, given a sequence of K ordered images, we can approximate Z^* and T^* by finding parameters ϕ and θ that minimize the prediction error for consecutive images and maximize it for the non-consecutive ones. For that, we defined a contrastive infoNCE loss based on those prediction errors,

$$\mathcal{L}_{M-CPC} = -\frac{1}{K-1} \sum_{i=1}^{K-1} \log \frac{e^{-\epsilon_{i,i+1}}}{\sum_{i'=1}^K e^{-\epsilon_{i,i'}}$$

and find ϕ and θ that minimize it.

2.2.2. RELATION NETWORK

Unlike the Markov-CPC model, which is inductively biased for finding causal relations between images, the RN model (Fig. 2 right) can learn *general* abstract relations. It achieves that by learning a general function $R_\theta(Z_\phi(x_i) \oplus Z_\phi(x_j))$ where \oplus is a concatenation operator.

Specifically for our tests, we used RN to learn the relation between consecutive pairs of inputs. Following (Sung et al., 2018), we did that by using an MSE classification loss that classified consecutive and non-consecutive images. Given a sequence of K images, we used the following RN loss

function,

$$\mathcal{L}_{RN} = \frac{1}{(K-1)^2} \sum_{i,j=1}^{K-1} \left(R_\theta(Z_\phi(x_i) \oplus Z_\phi(x_j)) - \delta_{j,i+1} \right)^2$$

in which the relation between consecutive images is assigned the label 1, and between non-consecutive images - the label 0.

3. Results

3.1. Naive few-shot learning

We applied the Markov-CPC and RN models to an SCE test in the following way: First, we randomly initialized the models' networks (see architectures in Appendix A1). We then updated these networks' weights with a *single* optimization step in the direction that minimizes the RN or Markov-CPC loss functions, given the $K = 5$ sequence images (optimizer details are also in Appendix A1). After the single optimization step, we evaluated the consistency of each choice image with the sequence based on the resulting Markov-CPC or RN loss function, when these choices were applied as the sixth image. We selected the most consistent choice image, out of the $n = 4$ choices, as the answer.

We found that both models performed much better than chance (0.25) in almost all tasks and all levels of difficulty (Fig. 3). However, their performance decreased with the number of distractors, indicating that this number is a good measure of the task's difficulty. We also found that some rules, Color and Size, seem easier than other rules. For example, Markov-CPC perfectly performed the easier Size

Naive Few-Shot Learning

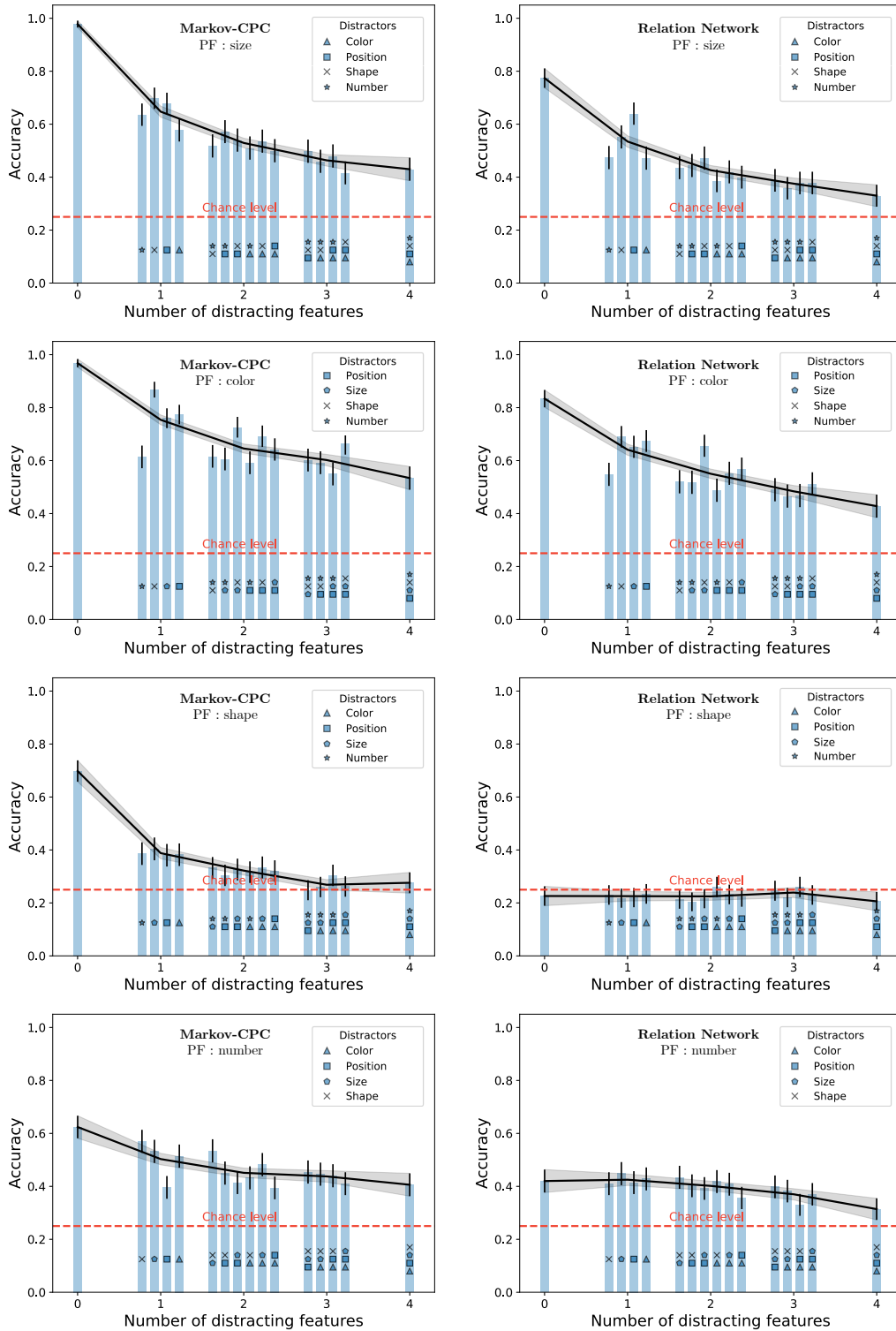


Figure 3. Performance of the Markov-CPC (left column) and RN (right column) models on SCE tests with four different predictive features (rows). For each predictive feature, we tested the networks over 16 test conditions where the rest of the features were either distractors (marked according to the legend) or constant (not marked). Each test condition included 500 randomly generated intelligence tests. Error bars are 95% confidence intervals. The black line and its shade are the average accuracy per difficulty and the standard deviation. The dashed line denotes the chance level ($\frac{1}{n}$).

and Color tasks. On the other hand, the Number rule seemed more difficult, which indicates that numerosity is not encoded accurately by naive networks. Notably, a similar claim has been made in the psychology literature about the encoding of numbers in “naive” humans (Leibovich et al., 2017; Siemann & Petermann, 2018). Overall, the Markov-CPC performed better than the RN. This was particularly pronounced when considering the Shape rule, in which the performance of RN, unlike Markov-CPC, did not exceed the chance level. Importantly, naive networks that did not train on the sequence images (did not perform the single gradient step) did not solve the task (see Appendix A2).

To study the determinants of the model’s performance, we tested various variations of Markov-CPC and RN (see Appendix A4). The emerging picture is that a match between the inductive bias of the model and the task crucially affects its performance. First, the Markov CPC that posits a causal relation between the consecutive images does better than the RN model, which allows for a more general relation between the images. Second, in our task, which is Markovian, the performance of a non-Markovian CPC, which allows a more complex relation between the latent variables, is slightly worse than that of the more restrictive Markov-CPC. Third, using the particular residual predictor T_θ ($T_\theta(Z) = Z + \Delta T_\theta(Z)$) resulted in higher performance than a non-residual predictor T_θ . Finally, increasing the dimensionality of the latent variable was detrimental to per-

formance (albeit the effect was small).

Additionally, adding complexity to the model did not improve its performance: both for the Markov CPC and the RN models, the relatively shallow network we used for the encoder Z was better than deeper and more complicated ones. This result may seem contradictory to the general trend of preferring ever-deeper networks. However, achieving higher performance in these more complex networks requires more data, while our networks were trained on a minimal number of examples.

3.2. Prior training

3.2.1. EXPRESSIVITY

While Markov-CPC performed substantially better than chance, it still failed to select the correct choice image in some tests. To see whether this resulted from a limited expressivity of the model, we pretrained Markov-CPC with SCE tests as training episodes. Specifically, we performed one optimization step for each of these training-episode tests and then tested the model on novel tests. We found that with 1000 training episodes, the performance of Markov-CPC on tests in which Size and Color are the predictive features can exceed 90%, even in the hardest trials (Table 1, best performance for each testing rule). These results indicate that, at least for these predictive features, the performance of the Markov-CPC model is not limited by its expressivity. By

Training \ Testing	Size (easy)	Size (hard)	Color (easy)	Color (hard)	Number (easy)	Number (hard)	Shape (easy)	Shape (hard)
Size (easy)	0.99 ±0.01	0.97 ±0.01	0.65 ±0.03	0.68 ±0.03	0.52 ±0.05	0.5 ±0.03	0.72 ±0.09	0.42 ±0.14
Size (hard)	0.9 ±0.0	0.88 ±0.01	0.24 ±0.0	0.25 ±0.01	0.24 ±0.01	0.24 ±0.01	0.32 ±0.04	0.31 ±0.04
Color (easy)	0.75 ±0.03	0.78 ±0.01	1.0 ±0.0	1.0 ±0.0	0.79 ±0.05	0.81 ±0.05	0.61 ±0.13	0.71 ±0.1
Color (hard)	0.26 ±0.01	0.27 ±0.01	0.93 ±0.01	0.95 ±0.01	0.27 ±0.01	0.26 ±0.01	0.53 ±0.03	0.34 ±0.05
Number (easy)	0.45 ±0.01	0.44 ±0.01	0.41 ±0.01	0.39 ±0.02	0.67 ±0.03	0.6 ±0.03	0.33 ±0.06	0.32 ±0.07
Number (hard)	0.26 ±0.0	0.26 ±0.01	0.28 ±0.01	0.25 ±0.01	0.7 ±0.02	0.67 ±0.03	0.33 ±0.03	0.26 ±0.04
Shape (easy)	0.18 ±0.02	0.25 ±0.03	0.29 ±0.02	0.33 ±0.03	0.2 ±0.02	0.26 ±0.03	0.37 ±0.06	0.31 ±0.05
Shape (hard)	0.24 ±0.01	0.25 ±0.01	0.26 ±0.01	0.26 ±0.0	0.24 ±0.01	0.25 ±0.01	0.27 ±0.01	0.24 ±0.02
Naive	0.97 ±0.0	0.42 ±0.01	0.96 ±0.0	0.55 ±0.01	0.60 ±0.01	0.40 ±0.01	0.71 ±0.0	0.26 ±0.01

Table 1. Each cell in the table is the average accuracies of 10 networks, each trained on 1000 episodes of one rule and tested using 500 tests on a different (or the same) rule. The errors are 95% confidence intervals.

contrast, when the Number or Shape were the predictive features, performance level improved with training but could only reach 70%–80%, leaving the question of expressivity of the model open.

3.2.2. SPECIALIZATION

The challenge of fluid intelligence is identifying regularities in a domain on which the agent has not been trained. One possibility to address this problem is to train on a large number of different examples, with the hope that they will generalize to the new domain. To test whether our networks can generalize, we tested whether pretraining in one type of test improves (or impairs) the model’s performance in another type of test. We considered 8 different test types for training episodes (four different predictive features, either in the easiest setting of no distractors or in the hardest setting of four distractors). We tested the performance on those 8 different tests, yielding an 8×8 performance matrix (Table 1). We found that typically, pretraining within a domain (same predictive feature) improved performance in that domain relative to the naive network. The exception was training on Shape, which, surprisingly, was detrimental to performance on Shape tests. Interestingly, we noted that within a domain, training using easy episodes was typically more effective than training using hard ones.

However, training in one domain was typically detrimental

to performance in other domains (relative to naive networks). This result highlights the potential advantage of “fluid” models over trained ones when the test domain is unknown.

3.3. Anomaly detection

To test naive few-shot learning models in natural settings and compare the performance of Markov-CPC and RN, we tested the abilities of these models to detect unlabeled anomalies in two datasets.

The first dataset, UMN (Mehran et al., 2009), consists of 11 security camera videos that A) start with humans walking normally; B) towards the end of the videos, they start running; C) a label of “Abnormal Crowd Activity” appears shortly after² (Fig. 4 A-C). We assigned each video frame an anomaly score based on its consistency with its five preceding frames (see details in appendix A3). We found that both the Markov-CPC and the RN models assigned larger anomaly scores to the frames in which the humans started to run (Fig. 4). Detecting the time in which running begins is not a difficult task for trained networks, and previous anomaly detection models achieved nearly perfect scores in this dataset (Pang et al., 2020). However, these previous models all relied on prior training. By contrast, our objective

²In the tests we removed the top 30 pixels of all images to exclude this label.

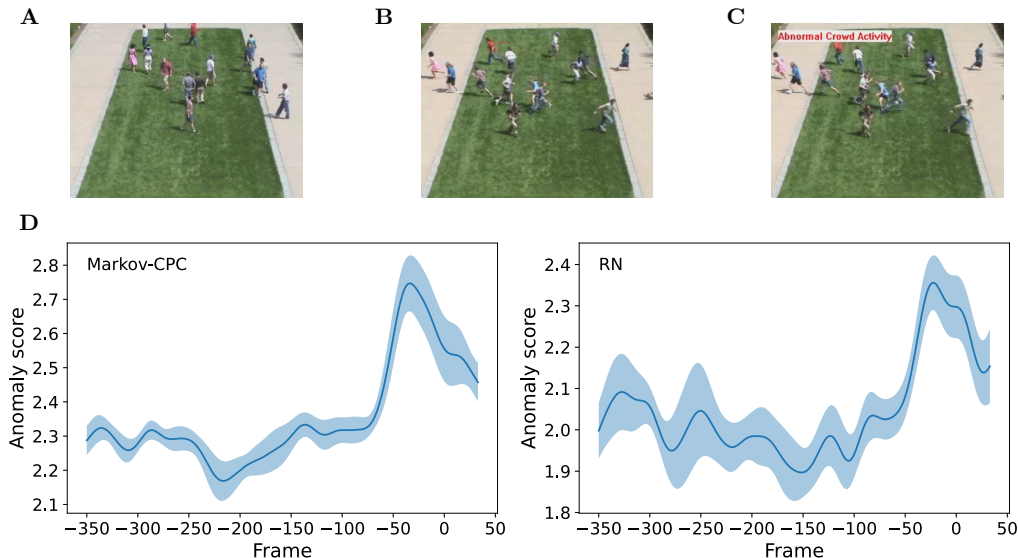


Figure 4. UMN Anomaly scores. (A-C) Example frames from the first UMN video: (A) normal walking human behavior; (B) humans start to run; (C) “Abnormal Crowd Activity” label appears. (D) Anomaly scores of the frames. The frames are numbered with respect to the first appearance of the “Abnormal Crowd Activity” label. We ran a model 5 times for each frame. We averaged the anomaly scores over the 5 runs per frame and smoothed the scores with a Gaussian kernel with a standard deviation of 10 frames. Finally, we averaged the resulting scores over the 11 videos.

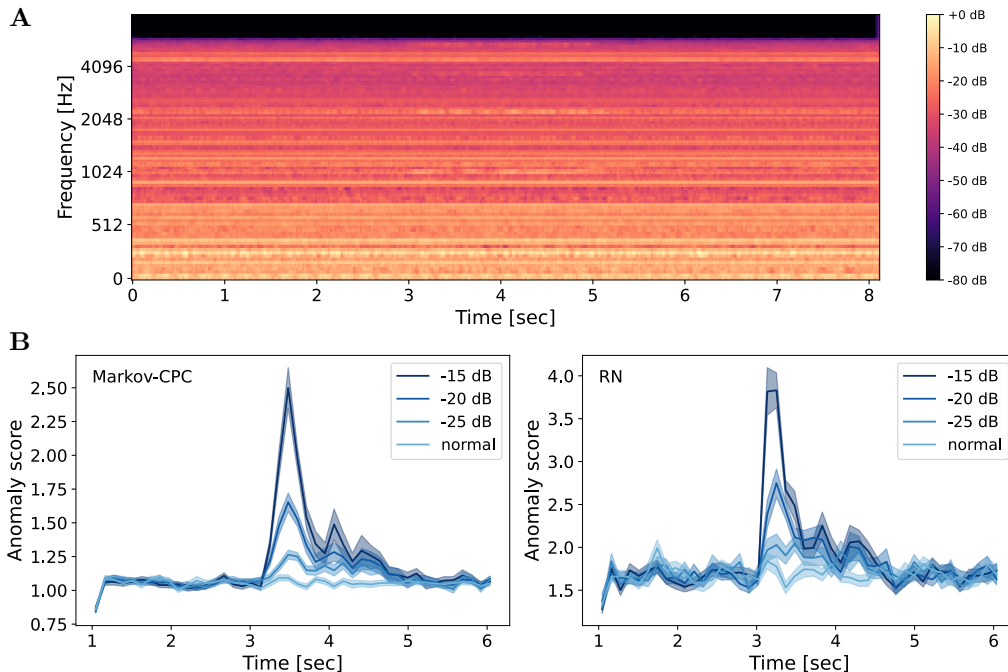


Figure 5. **ADS anomaly scores.** (A) Example spectrogram of a phone that rings over natural background noise. The anomaly-to-noise (ANR) ratio in this example is -20 dB. (B) Anomaly scores of Markov-CPC and RN. We averaged the anomaly scores over five iterations per frame of each sound snippet and the 140 sound snippets in each ANR.

was to identify these anomalies without prior training.

The second dataset, Anomaly Detection of Sound (ADS) (Koizumi et al., 2019), consists of 140 short sound snippets. Each snippet consists of natural background noise, such as the noise of an air conditioner recorded in a natural environment, interrupted by various types of anomalous sounds, such as keys falling or a drawer closing. The detection difficulty was determined by the anomaly-to-noise power ratios (ANRs), which were set to -15 dB, -20 dB, or -25 dB. This dataset is more challenging than the UMN dataset; even trained networks fail to detect anomalous sounds in some examples. To test our models on this dataset, we converted the sound snippets into a video in which each frame was a Mel-spectrogram window of ~ 0.5 seconds of the snippet, hopped in ~ 0.1 second steps. Then, similar to the analysis of the UMN dataset, we computed an anomaly score for each image based on the previous five frames using either the Markov-CPC or the RN models. We found that both models assigned relatively higher anomaly scores to anomalous events (Fig. 5).

To compare our results with other models, we used the standard AUC measure (area under the ROC curve) (Pedregosa et al., 2011) based on the maximal anomaly score of each snippet. For Markov-CPC, The AUC scores were around 0.96, 0.84, and 0.67 for the -15, -20, and -25 dB ANRs respectively, which are not much worse than those of previous

models. This is despite the fact that our model was naive, whereas previous models required substantial prior training (Fig. 6). The RN model exhibited substantially poorer performance in this task than the Markov-CPC model. The finding that the Markov-CPC model performs better than the RN model is consistent with its relative success in the SCE tests.

4. Discussion

Our results showed that deep learning tools could be used to solve non-trivial tasks without any prior training. The Markov-CPC model successfully solved SCE tests and detected anomalies in data streams, starting with random weights. This approach has practical implications. Today, deep learning models lack “fluid” abilities, impairing their performance in uncertain and shifting environments such as driving in bad weather conditions (Zang et al., 2019). Markov-CPC’s ability to detect anomalies in data streams without prior training is useful in those uncertain environments, where prior training might even be detrimental (Section 3.2.2).

Moreover, our work suggests a method to develop fluid intelligent models. A complex task can be separated into sub-components. Our findings suggest that deep learning models can solve some of these sub-components without

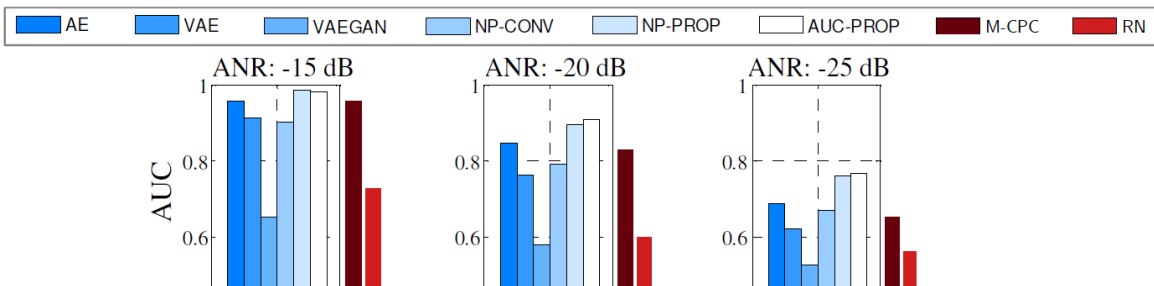


Figure 6. **ADS model comparison.** AUC scores for anomaly detection of different models. Bars in the boxes depict the AUC scores of six models that require prior training (adapted from (Koizumi et al., 2019)). Darker and lighter red depict the AUC scores of the Markov-CPC and RN models, respectively.

relying on prior knowledge. Like our anomaly detection model, these fluid-intelligence models will be flexible to changing environments. The main challenge is finding these sub-components for which there is an appropriate inductive bias. We encourage the finding of such sub-components, and the further improvement of naive models’ ability to solve the SCE task.

Regarding the relationship between our approach and standard deep learning methodology, in which models are extensively trained to solve complex problems. While we focused on naive untrained models, it is, of course, possible to pretrain fluid intelligent models. Indeed, today’s standard few-shot learning approach is to pretrain deep learning models over large relevant, diverse datasets and then fine-tune them for the new tasks (Brown et al., 2020; Reed et al., 2022). This approach has limited generalization ability to datasets that are different from the datasets they were trained on (Li et al., 2017; Nalisnick et al., 2019; Yin et al., 2020; Rajendran et al., 2020). However, we expect the combination of pretrained weights with high fluid intelligence to work best in relatively stable environments.

In a field that is dominated by extensive training with an extensive number of examples, our work demonstrates the potential power of “fluid” models. Pursuing this line of research, making artificial fluid-intelligent models, will make machines considerably more flexible.

References

- Barrett, D. G. T., Hill, F., Santoro, A., Morcos, A. S., and Lillicrap, T. Measuring abstract reasoning in neural networks. 2018. ISSN 17740746. doi: 10.1051/agro/2009059. URL <http://arxiv.org/abs/1807.04225>. ISBN: 1807.04225v1.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Aspell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language Models are Few-Shot Learners. Technical Report arXiv:2005.14165, arXiv, July 2020. URL <http://arxiv.org/abs/2005.14165>. arXiv:2005.14165 [cs] type: article.
- Carpenter, P. A., Just, M. A., and Shell, P. What One Intelligence Test Measures: A Theoretical Account of the Processing in the Raven Progressive Matrices Test. (3): 28, 1990.
- Chollet, F. On the Measure of Intelligence. *arXiv:1911.01547 [cs]*, November 2019. URL <http://arxiv.org/abs/1911.01547>. arXiv:1911.01547.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Henaff, O. Data-Efficient Image Recognition with Contrastive Predictive Coding. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 4182–4192. PMLR, November 2020. URL <https://proceedings.mlr.press/v119/henaff20a.html>. ISSN: 2640-3498.
- Hochreiter, S. and Schmidhuber, J. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. Conference Name: Neural Computation.
- Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam, H. Mobilenets:

- Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Kim, Y., Shin, J., Yang, E., and Hwang, S. J. Few-shot Visual Reasoning with Meta-analogical Contrastive Learning. *arXiv:2007.12020 [cs, stat]*, July 2020. URL <http://arxiv.org/abs/2007.12020>. arXiv:2007.12020.
- Koizumi, Y., Saito, S., Uematsu, H., Kawachi, Y., and Harada, N. Unsupervised Detection of Anomalous Sound Based on Deep Learning and the Neyman–Pearson Lemma. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(1):212–224, January 2019. ISSN 2329-9304. doi: 10.1109/TASLP.2018.2877258. Conference Name: IEEE/ACM Transactions on Audio, Speech, and Language Processing.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. ImageNet classification with deep convolutional neural networks. In Pereira, F., Burges, C., Bottou, L., and Weinberger, K. (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. URL <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>.
- Leibovich, T., Katzin, N., Harel, M., and Henik, A. From “sense of number” to “sense of magnitude”: The role of continuous magnitudes in numerical cognition. *Behavioral and Brain Sciences*, 40:e164, 2017. doi: 10.1017/S0140525X16000960.
- Li, D., Yang, Y., Song, Y.-Z., and Hospedales, T. M. Learning to Generalize: Meta-Learning for Domain Generalization. *arXiv:1710.03463 [cs]*, October 2017. URL <http://arxiv.org/abs/1710.03463>. arXiv:1710.03463.
- Lohman, D. F. Complex Information Processing and Intelligence. In Sternberg, R. J. (ed.), *Handbook of Intelligence*, pp. 285–340. Cambridge University Press, Cambridge, 2000. ISBN 978-0-521-59648-0. doi: 10.1017/CBO9780511807947.015.
- Mehran, R., Oyama, A., and Shah, M. Abnormal crowd behavior detection using social force model. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 935–942, June 2009. doi: 10.1109/CVPR.2009.5206641. URL http://mha.cs.umn.edu/proj_events.shtml#crowd. ISSN: 1063-6919.
- Nalisnick, E., Matsukawa, A., Teh, Y. W., Gorur, D., and Lakshminarayanan, B. Do Deep Generative Models Know What They Don’t Know? Technical Report arXiv:1810.09136, arXiv, February 2019. URL <http://arxiv.org/abs/1810.09136>. arXiv:1810.09136 [cs, stat] type: article.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation Learning with Contrastive Predictive Coding. 2018. URL <http://arxiv.org/abs/1807.03748>.
- Pang, G., Yan, C., Shen, C., Hengel, A. v. d., and Bai, X. Self-Trained Deep Ordinal Regression for End-to-End Video Anomaly Detection. pp. 12173–12182, 2020. URL https://openaccess.thecvf.com/content_CVPR_2020/html/Pang_Self-Trained_Deep_Ordinal_Regression_for_End-to-End_Video_Anomaly_Detection_CVPR_2020_paper.html.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830, 2011.
- Rajendran, J., Irpan, A., and Jang, E. Meta-Learning Requires Meta-Augmentation. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 5705–5715. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/3e5190eeb51ebe6c5bbc54ee8950c548-Paper.pdf>.
- Raven, J., Raven, J. C., and Court, J. H. *Manual for Raven’s progressive matrices and vocabulary scales*. Pearson, San Antonio, TX, 1998. ISBN 978-0-15-868643-1 978-0-15-468622-0 978-0-15-468623-7 978-1-85639-022-4. OCLC: 697438611.
- Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S. G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J. T., Eccles, T., Bruce, J., Razavi, A., Edwards, A., Heess, N., Chen, Y., Hadsell, R., Vinyals, O., Bordbar, M., and de Freitas, N. A Generalist Agent. Technical Report arXiv:2205.06175, arXiv, May 2022. URL <http://arxiv.org/abs/2205.06175>. arXiv:2205.06175 [cs] type: article.
- Santoro, A., Raposo, D., Barrett, D. G. T., Malinowski, M., Pascanu, R., Battaglia, P., and Lillicrap, T. A simple neural network module for relational reasoning. *Advances in Neural Information Processing Systems*, 2017-Decem:4968–4977, 2017. ISSN 10495258.

- Siebers, M., Dowe, D. L., Schmid, U., Hernández-Orallo, J., and Martínez-Plumed, F. Computer models solving intelligence test problems: Progress and implications. *Artificial Intelligence*, 230:74–107, 2015. ISSN 00043702. doi: 10.1016/j.artint.2015.09.011. URL <http://dx.doi.org/10.1016/j.artint.2015.09.011>. ISBN: 9780999241103 Publisher: Elsevier B.V.
- Siemann, J. and Petermann, F. Innate or acquired? – disentangling number sense and early number competencies. *Frontiers in Psychology*, 9, 2018. ISSN 1664-1078. doi: 10.3389/fpsyg.2018.00571. URL <https://www.frontiersin.org/articles/10.3389/fpsyg.2018.00571>.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Sternberg, R. J. Component processes in analogical reasoning. *Psychological Review*, 84(4):353–378, 1977. ISSN 0033295X. doi: 10.1037/0033-295X.84.4.353.
- Sternberg, R. J. Components of human intelligence. *Cognition*, 15(1-3):1–48, 1983. ISSN 00100277. doi: 10.1016/0010-0277(83)90032-X.
- Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P. H. S., and Hospedales, T. M. Learning to Compare: Relation Network for Few-Shot Learning. pp. 1199–1208, 2018. URL https://openaccess.thecvf.com/content_cvpr_2018/html/Sung_Learning_to_Compare_CVPR_2018_paper.html.
- Vogelstein, J. T., Verstynen, T., Kording, K. P., Isik, L., Krakauer, J. W., Etienne-Cummings, R., Ogburn, E. L., Priebe, C. E., Burns, R., Kuttan, K., Knierim, J. J., Potash, J. B., Hartung, T., Smirnova, L., Worley, P., Savonenko, A., Phillips, I., Miller, M. I., Vidal, R., Sulam, J., Charles, A., Cowan, N. J., Bichuch, M., Venkataraman, A., Li, C., Thakor, N., Kebschull, J. M., Albert, M., Xu, J., Shuler, M. H., Caffo, B., Ratnanather, T., Geisa, A., Roh, S.-E., Yezerets, E., Madhyastha, M., How, J. J., Tomita, T. M., Dey, J., Ningyuan, Huang, Shin, J. M., Kinfu, K. A., Chaudhari, P., Baker, B., Schapiro, A., Jayaraman, D., Eaton, E., Platt, M., Ungar, L., Wehbe, L., Kepecs, A., Christensen, A., Osuagwu, O., Brunton, B., Mensh, B., Muotri, A. R., Silva, G., Puppo, F., Engert, F., Hillman, E., Brown, J., White, C., and Yang, W. Prospective Learning: Back to the Future. *arXiv:2201.07372 [cs]*, January 2022. URL <http://arxiv.org/abs/2201.07372>. arXiv: 2201.07372.
- Webb, T. W., Holyoak, K. J., and Lu, H. Emergent analogical reasoning in large language models. *ArXiv*, abs/2212.09196, 2022.
- Yin, M., Tucker, G., Zhou, M., Levine, S., and Finn, C. Meta-Learning without Memorization. *arXiv:1912.03820 [cs, stat]*, April 2020. URL <http://arxiv.org/abs/1912.03820>. arXiv: 1912.03820.
- Zang, S., Ding, M., Smith, D., Tyler, P., Rakotoarivelo, T., and Kaafar, M. A. The impact of adverse weather conditions on autonomous vehicles: How rain, snow, fog, and hail affect the performance of a self-driving car. *IEEE Vehicular Technology Magazine*, 14(2):103–111, 2019. doi: 10.1109/MVT.2019.2892497.
- Zhuo, T. and Kankanhalli, M. Solving Raven’s Progressive Matrices with Neural Networks. *arXiv:2002.01646 [cs]*, February 2020. URL <http://arxiv.org/abs/2002.01646>. arXiv: 2002.01646.

Appendix

A1. Models details

A1.1. Markov-CPC

Markov-CPC’s encoder Z and predictor T were implemented by deep neural networks. We used a relatively shallow convolutional neural network for the encoder Z . For the predictor T , we used a residual network, such that

$$T(Z(\mathbf{x})) = Z(\mathbf{x}) + \Delta T(Z(\mathbf{x}))$$

where ΔT is a fully connected neural network.

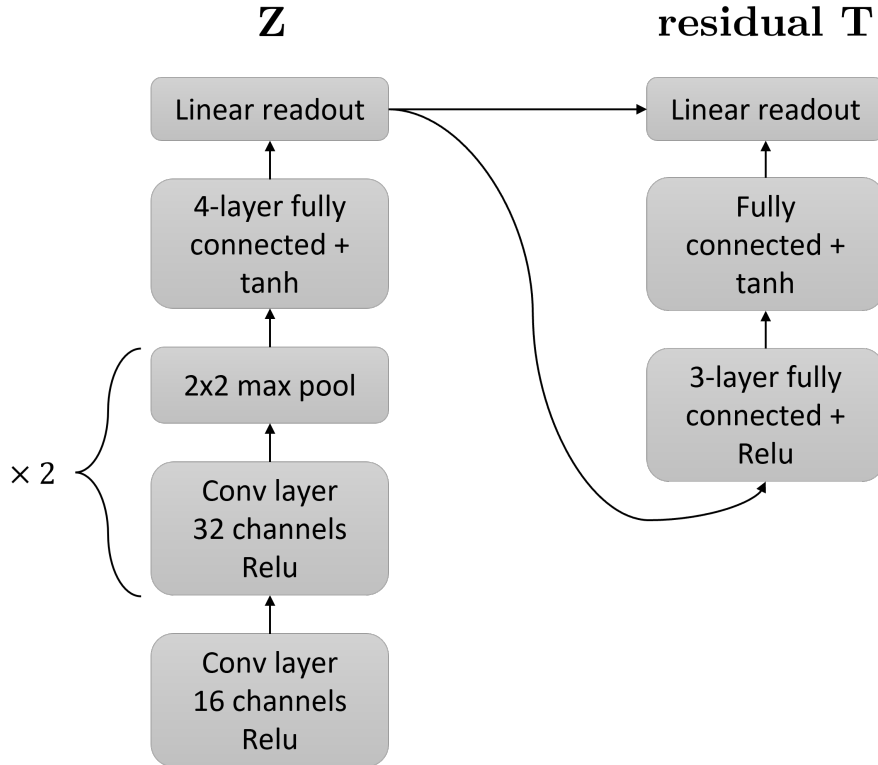


Figure f1. M-CPC architecture.

To update the weights of the networks, we used the RMSprop optimizer with a learning rate $\eta = 4 \cdot 10^{-4}$. The rest of the optimizer hyperparameters were set to PyTorch defaults.

A1.2. Relation Network

We used the same Relation Network loss function as in the original paper (Sung et al., 2018). For the networks, we either used the networks from (Sung et al., 2018) or used the same Z and R as Markov-CPC’s Z and T . In both cases, we used the RMSprop optimizer. For the original RN networks from (Sung et al., 2018) we used a learning of $\eta = 4 \cdot 10^{-6}$. When we used the Z of Markov-CPC, we used the learning rate $\eta = 4 \cdot 10^{-4}$. The rest of the optimizer parameters were set to PyTorch default.

A2. Control experiment

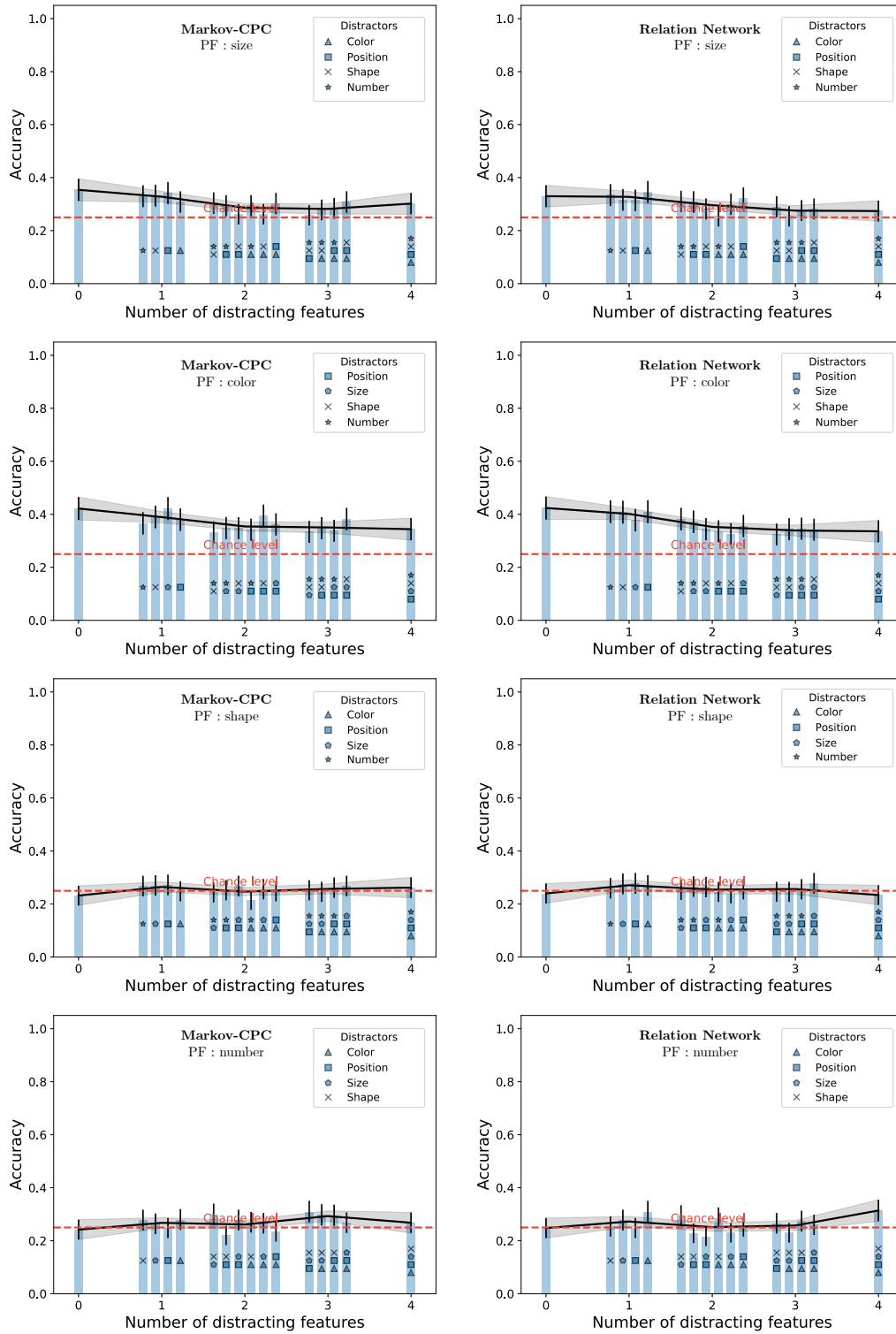


Figure f2. Control experiment. The accuracies that Markov-CPC and the RN models achieved when they solved tests without performing a gradient step.

A3. Anomaly score based on SCE

To assign an anomaly score to a candidate image \mathbf{x}_c based on its preceding images, we optimized a naive Markov-CPC or RN model, with a single optimization step, on the $K = 5$ preceding images. To determine the candidate’s anomaly score s_c with the Markov-CPC model, we checked by how much the prediction error of the candidate image given the K^{th} image, $\epsilon_{K,c}$, deviates from the prediction errors of the $K - 1$ preceding consecutive pair images. Mathematically, the score is defined by,

$$s_c = \frac{\epsilon_{K,c} - \langle \epsilon_p \rangle}{\text{std}(\epsilon_p)},$$

where $\langle \epsilon_p \rangle$ is the average prediction error of the $K - 1$ preceding consecutive pairs and $\text{std}(\epsilon_p)$ is their standard deviation.

To determine s_c with the RN model, we checked by how much the candidate’s classification error,

$$\epsilon_{K,c} = \left(R_\theta(Z_\phi(x_K) \oplus Z_\phi(x_c)) - 1 \right)^2 \quad (1)$$

deviates from the classification errors of the $K - 1$ preceding consecutive pair images. We used the same mathematical score defined above, where the classification errors are used to define as in Eq (1).

A4. Model variations

Markov-CPC achieved an average accuracy of 0.52 ± 0.02 over all the test conditions, while RN’s average accuracy was 0.42 ± 0.02 . We tested how this average accuracy changes for various model variations.

Markov-CPC was implemented with a residual T , such that $T(Z(\mathbf{x})) = Z(\mathbf{x}) + \Delta T(Z(\mathbf{x}))$ where ΔT is a neural network (see appendix A1). we also tested a variant of Markov-CPC in which T is non-residual. We found that the performance of a Markov-CPC with a non-residual T is substantially worse than that with a residual T (Table t1). This indicates that a residual T is a good inductive bias for finding these rules.

Model variant	Total accuracy
Residual T	0.52 ± 0.02
Non-residual T	0.29 ± 0.02

Table t1. Residual versus non-residual T .

CPC models are usually non-Markovian (Oord et al., 2018; Henaff, 2020). To test the effect of memory on the performance of the CPC model, we used latent variables Z whose values were also dependent on previous images in the sequence, either via a regular recurrent neural network (RNN) or an LSTM (Hochreiter & Schmidhuber, 1997). Both the recurrent connections and the LSTM impaired performance, while taking twice the time to compute (Table t2). This indicates that adding recurrent weights to the latent variables, when the data can be explained by Markov latent variables, can impair the ability of the network to extract the rule.

Model variant	Total accuracy	Tests per second
Markov-CPC	0.52 ± 0.02	13 ± 2 [Hz]
LSTM-CPC	0.48 ± 0.02	6.6 ± 0.4 [Hz]
RNN-CPC	0.47 ± 0.02	6.6 ± 0.4 [Hz]

Table t2. Markov versus recurrent networks. The running times were measured on a basic laptop with Nvidia RTX 2070 GPU.

Throughout the paper, we used Markov-CPC with a 1-dimensional latent space. We tried changing the latent space dimension of Markov-CPC. We found that even a 1-dimensional latent space is enough and that the performance does not strongly depend on this value (Table t3).

Naive Few-Shot Learning

Model variant	Total accuracy
Markov-CPC-1D	0.52 ± 0.02
Markov-CPC-10D	0.53 ± 0.02
Markov-CPC-100D	0.51 ± 0.02
Markov-CPC-1000D	0.50 ± 0.02

Table t3. Different Z latent dimensions.

We also measured the performance of the model without a contrastive loss, minimizing the prediction errors (Eq. 2.2.1) of consecutive inputs only (Table t4).

Model variant	Total accuracy
Contrastive loss	0.52 ± 0.02
No contrast	0.45 ± 0.02

Table t4. No contrast.

We trained Markov-CPC throughout the paper with an RMSprop optimizer. We also measured the performance with a standard SGD optimizer (Table t5). The optimal learning rate was found to be $\eta = 40$.

Model variant	Total accuracy
RMSprop, lr = $4 \cdot 10^{-4}$	0.52 ± 0.02
SGD, lr = 40	0.5 ± 0.02

Table t5. SGD.

We also compared the performance of the original RN networks (Sung et al., 2018) to the more shallow networks that we used for Markov-CPC (Fig. f2). The shallower networks were better (Table t6), indicating that in the RN models, simpler networks are better in this task than deeper networks.

Model variant	Total accuracy
Shallow RN (Fig. f2)	0.42 ± 0.02
Deep RN (Sung et al., 2018)	0.27 ± 0.02

Table t6. Deep versus shallow RN networks.

Both in the Markov-CPC and RN models, we used a relatively shallow network for the encoder Z (Fig. f2). We tested other, more complex and deep, networks from the literature as candidate encoder backbones of the Markov-CPC model. The shallow encoder was better than various complicated and deep networks from the literature (Table t7).

Model variant	Total accuracy
Paper's network (Fig. f2)	0.52 ± 0.02
VGG11 (Simonyan & Zisserman, 2014)	$0.45 \pm 0.02^*$
DenseNet121 (Huang et al., 2017)	0.33 ± 0.02
AlexNet (Krizhevsky et al., 2012)	0.31 ± 0.02
ResNet18 (He et al., 2016)	$0.29 \pm 0.02^*$
MobileNet v3 small (Howard et al., 2017)	0.29 ± 0.02

Table t7. Encoder Z networks comparison. *ResNet18 and VGG11 achieved relatively high accuracy on the color rule (0.47 ± 0.04 and 0.78 ± 0.04 respectively).