# A Study of Compositional Generalization in Neural Models

**Tim Klinger**[*]
Thomas J. Watson Research Center
IBM Research AI
Yorktown, NY, USA
`tklinger@us.ibm.com`

**Dhaval Adjodah**[*]
MIT Media Lab
Cambridge, MA
`dval@mit.edu`

**Vincent Marois**
Thomas J. Watson Research Center
IBM Research AI
Yorktown, NY, USA
`vincent.marois@ibm.com`

**Josh Joseph**
MIT Media Lab
Cambridge, MA
`jmjoseph@mit.edu`

**Matthew Riemer**
Thomas J. Watson Research Center
IBM Research AI
Yorktown, NY, USA
`mdriemer@us.ibm.com`

**Alex 'Sandy' Pentland**
MIT Media Lab
Cambridge, MA
`pentland@mit.edu`

**Murray Campbell**
Thomas J. Watson Research Center
IBM Research AI
Yorktown, NY, USA
`mcam@us.ibm.com`

## Abstract

Compositional and relational learning is a hallmark of human intelligence, but one which presents challenges for neural models. One difficulty in the development of such models is the lack of benchmarks with clear compositional and relational task structure on which to systematically evaluate them. In this paper, we introduce an environment called ConceptWorld, which enables the generation of images from compositional and relational concepts, defined using a logical domain specific language. We use it to generate images for a variety of compositional structures: 2x2 squares, pentominoes, sequences, scenes involving these objects, and other more complex concepts. We perform experiments to test the ability of standard neural architectures to generalize on relations with compositional arguments as the compositional depth of those arguments increases and under substitution. We compare standard neural networks such as MLP, CNN and ResNet, as well as state-of-the-art relational networks including WReN and PrediNet in a multi-class image classification setting. For simple problems, all models generalize well to close concepts but struggle with longer compositional chains. For more complex tests involving substitutivity, all models struggle, even with short chains. In high-lighting these difficulties and providing an environment for further experimentation, we hope to encourage the development of models which are able to generalize effectively in compositional, relational domains.

## 1 Introduction

Humans have a relational and compositional view which enables them to better manage the complexity of the world [Nav77, FP⁺88, SK07], but extracting this view from images and text is a challenge

---

[*]equal contribution.

for neural networks [LB18, Lou18]. Recent work has begun to address this challenge through the development of relational datasets and frameworks - for example [JHvdM+17, HM19b] for Visual Question Answering. These have in turn spurred research for novel relational models. We believe there is a similar need for datasets of compositional, relational concepts, however we know of no work to enable formal specification and image generation for them. In this paper, we describe such an environment, which we call ConceptWorld.

Concepts in ConceptWorld are specified hierarchically in a logical language. We see several benefits in this approach: (1) it makes it easier to define concepts whose structure is clear (to the author and others) because it is logical and declarative rather than procedural (it specifies what it is, not how to compute it); (2) it allows rapid prototyping and experimentation because it reduces the amount of code that needs to be written; (3) it supports easy sharing between domains as concepts are hierarchical and lower level ones can be re-used. As a test of ConceptWorld's ability to represent concepts of interest, we use it to recreate the key-lock task of Box-World [ZRS+18] in our setting.

We consider two types of (zero-shot) compositional generalization [HDMB20]. In the *productivity* experiments, we examine compositional generalization of a concept relation to greater compositional depths than have been seen in training (for example, from a length 2 sequence of squares to a length 3 sequence). In the *substitutivity* experiments, we maintain the same depth of composition from training to test but change the object being composed. For example, given training on (1) a concept which consists of red squares and f-pentominoes [2], and (2) sequences of red squares, can the model generalize to same-length sequences of f-pentominoes? Although the concepts we use in these experiments are artificial, the compositional patterns we discuss occur frequently in natural and man-made images (the patterned fabric of a dress for example) and the ability to recognize them in a way which compositionally generalizes is one which we believe is central to more efficient learning and effective generalization.

We conduct experiments to evaluate compositional generalization on four example domains (specified using ConceptWorld) and a variety of standard models (MLP, CNN, ResNet [HZRS16]) and newer ones designed specifically to extract relational representations (WReN [Ada18], PrediNet [SNC+19]). None of the models we evaluate are explicitly biased to encourage compositional generalization.

To summarize, our paper makes the following contributions:

- ConceptWorld: A concept specification language and generator for compositional relational concepts.

- Four tasks each with their own domain: two which test compositional *productivity* (experiment 1: pure and mixed sequences; and experiment 2: Box-World sequences); and two which test compositional *substitutivity* (experiment 3: 2x2->4x4 patterned squares; and experiment 4: sequence substitutions).

- An evaluation of standard and relational models on these domains. Our results demonstrate that the evaluated models struggle in these settings, including recently proposed relational models. This suggests that new research is needed to encourage compositional generalization with neural models and ConceptWorld provides a much needed test-bed in this direction.

## 2 Related Work

We focus here on work related to compositionality of neural models, compositional generalization, relational learning in neural models, and reasoning tests.

[HDMB20] discusses a variety of interpretations of the term "compositional generalization". Their taxonomy includes "productivity" (generalization to deeper compositional depths) and "substitutivity" (generalization to objects of the same type as seen in training), which we evaluate here. In [LB18], they develop a compositional dataset called SCAN, and test compositional generalization on recurrent neural network models. This is extended in [Lou18]. A framework for measuring compositional generalization is proposed in [Ada18] using Raven-style matrices. It is based on two principles: (1) minimize the difference between train and test distributions for "atomic" concepts (the base concepts

---

[2]A pentomino is a set of 5 adjacent points. Pentominoes come in different shapes which are labeled with letters (like "f") which they resemble. See Figure 1b.
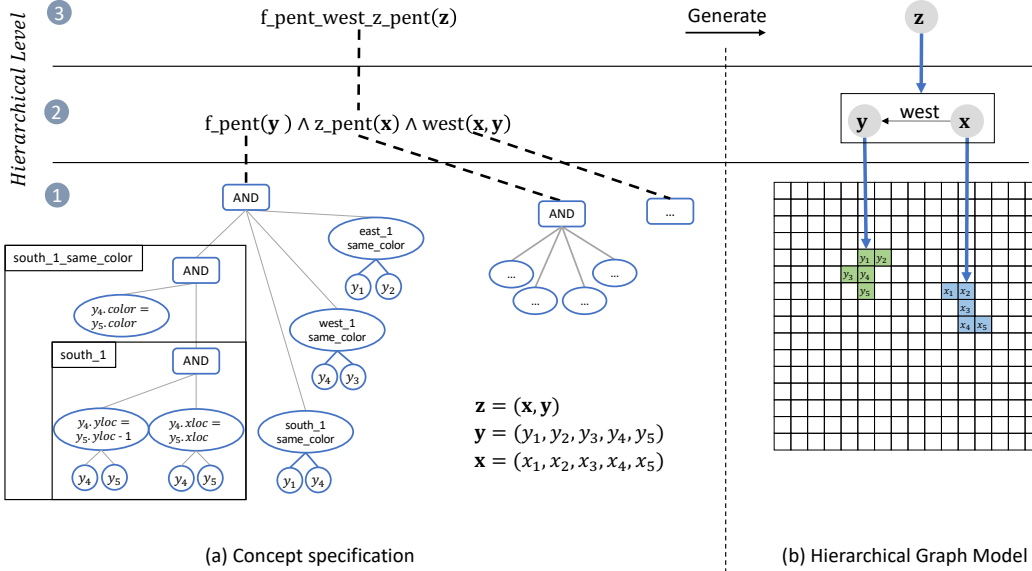
Figure 1: Representation of the `f-pent-west-z-pent` concept. (a): View of the concept specification. The compositional levels are labelled 1-3. (b) A hierarchical scene graph with accompanying generated image (32 x 32 pixel).

used in compositions); and (2) maximize the difference between train and test distributions for "compound" (composite) concepts and uses a procedural language with logical operators. We have attempted to follow these principles in the construction of our experiments. Similarly, [SNC+19] introduces the Relations Game as a set of binary image classification tasks, where the aim is to label whether a specified relation holds. They use visual shapes such as pentominoes and hexominoes to test substitutivity, when the unique relation is given. We use the WReN and PrediNet models introduced in these papers as baselines in our experiments.

Compositional generalization is considered in [CGLG18] where they introduce a model called Compositional Recursive Learner (CRL) for multi-task learning. The emphasis there is on transformations such as language translation rather than classification, but the CRL model or a Routing Network [RKR18] might provide a good starting point for an architecture capable of generalizing more systematically than the ones considered in our experiments. Unfortunately, there are several known challenges for stabilizing the learning of these models [RCRK19].

In relational reasoning, the Visual Question Answering (VQA) [AAL+15] setting is used frequently. [JHvdM+17] introduces the CLEVR dataset, linking templatized natural language questions to elementary visual scenes. Critics of such artificial language datasets [Man19, ZGSS+16, JJVDM16] have pointed to the linguistic and semantic simplicity of the questions, as well as tendencies in the answers distributions as circumventing the need for "true" visual reasoning. [HM19b] introduced the GQA dataset, to remediate some of these shortcomings.

VQA research has spawned the development of several relevant models. In [Dre18], they show that iterative, attention-based reasoning leads to more data-efficient learning. [HM19a, Jia19] draw on the strengths of neural and symbolic AI to perform visual reasoning. [SRB+17] proposes a simple neural module to reason about entities (visual or textual) and their relations. The Neuro-Symbolic Concept Learner (NSCL) [Jia19] is a multi-modal model which can learn visual concepts from training on images paired with textual questions and answers. In [HSM+18], they propose a generative model called SCAN (unrelated to the dataset SCAN), based on a $\beta$-VAE which can learn grounded hierarchical concepts from a small number symbol-image pairs.

3

# 3 ConceptWorld

## 3.1 Definitions

A *concept* is a unary or binary relation over objects. Objects can be simple points or vectors of objects, themselves possibly vectors. Concepts whose parameters are points are called *primitive* concepts; while those with parameters which are vectors satisfying lower-level concepts are called *higher-order*. For example, the unary relation $\texttt{red}(x_1)$ (The object $x_1$ is red) and the binary relation $\texttt{west}(x_1, x_2)$ ($x_2$ is west of $x_1$), for $x_1$ and $x_2$ grid points are both primitive concepts. See Table 1 for an example of a higher order unary concept whose argument is a vector $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5)$ constrained to satisfy the definition of an f-pentomino.

Concepts are specified in a logical language which supports composite terms (variables and constants; denoted here in bold). For example, given the concepts of an f-pentomino and a z-pentomino, we can define the composite concept of an f-pentomino-west-of-a-z-pentomino. This is specified in a higher order relation: $\texttt{f-pent-west-z-pent}(\mathbf{z}) \equiv \texttt{west}(\mathbf{x}, \mathbf{y}) \wedge \texttt{f-pent}(\mathbf{y}) \wedge \texttt{z-pent}(\mathbf{x})$. $\texttt{west}$ is taken to mean that all points $y_i$ are west of all points $x_i$. We define it as $\texttt{west}(\mathbf{x}, \mathbf{y}) \equiv [\bigwedge_{i,j} \texttt{west\_point}(x_i, y_j)] \wedge \mathbf{x} == (x_i) \wedge \mathbf{y} == (y_j)$ using the primitive concept $\texttt{west\_point}(x, y)$ which says that point $y$ is one grid point west of point $x$. Fig. 1a illustrates the decomposition.

| Concept | Definition |
|---|---|
| red | $\texttt{red}(x) \equiv x.color = \texttt{RED}$ |
| point | $\texttt{point}(x) \equiv (\texttt{red}(x) \vee \texttt{blue}(x) \vee \cdots \vee \texttt{yellow}(x)) \wedge$ $(x.x\_loc == 0 \vee \cdots \vee x.x\_loc == \texttt{GRID\_SIZE} - 1) \wedge$ $(x.y\_loc == 0 \vee \cdots \vee x.y\_loc == \texttt{GRID\_SIZE} - 1)$ |
| 2x2 square of points | $\texttt{2x2\_square}(\mathbf{x}) \equiv \mathbf{x} == (x_1, x_2, x_3, x_4) \wedge \texttt{east}_1(x_1, x_2) \wedge$ $\texttt{south}_1(x_2, x_3) \wedge \texttt{west}_1(x_3, x_4) \wedge \texttt{north}_1(x_3, x_1) \wedge$ $\texttt{point}(x_1) \wedge \texttt{point}(x_2) \wedge \texttt{point}(x_3) \wedge \texttt{point}(x_4)$ |
| south 1 point and same color | $\texttt{south}_1\_\texttt{same\_color}(x_1, x_2) \equiv \texttt{south}_1(x_1, x_2) \wedge x_1.color == x_2.color$ |
| f-pentomino | $\texttt{f-pent}(\mathbf{x}) \equiv \texttt{east}_1\_\texttt{same\_color}(x_1, x_2) \wedge \texttt{south}_1\_\texttt{same\_color}(x_4, x_1) \wedge$ $\texttt{west}_1\_\texttt{same\_color}(x_4, x_3) \wedge \texttt{south}_1\_\texttt{same\_color}(x_4, x_5) \wedge \mathbf{x} == (x_1, x_2, x_3, x_4, x_5)$ $\wedge \texttt{point}(x_1) \wedge \texttt{point}(x_2) \wedge \texttt{point}(x_3) \wedge \texttt{point}(x_4) \wedge \texttt{point}(x_5)$ |

Table 1: Examples of concepts used in our study.

## 3.2 Generation: Concepts to Images

For each concept, we can generate a hierarchical scene graph [JKS+15] whose lowest level corresponds directly to the image we want to generate. The details of the generation process can be found in Appendix 1. Here, we give a sketch of the procedure for the example shown in Figure 1. We don't discuss disjunction in this example but the algorithm can handle it by converting the concept definition to Disjunctive Normal Form (DNF) and applying this procedure to each conjunctive clause, until one is found which successfully generates.

For a unary concept like $\texttt{f\_pent\_west\_z\_pent}(\mathbf{z})$, we generate a single node for its object $\mathbf{z} = (\mathbf{x}, \mathbf{y})$ and recursively generate the constituent objects $\mathbf{x}$ and $\mathbf{y}$ according to the definition. This requires first building the scene graph for that definition, with nodes $\mathbf{x}$ and $\mathbf{y}$ and an edge between them labeled with the concept they satisfy, $\texttt{west}$. This graph is then traversed in breadth-first order starting from the first argument of the relation (here $\mathbf{x}$). First, the $\mathbf{x}$ node is recursively generated, followed by the $\texttt{west}$ relation, which recursively generates $\mathbf{y}$ so that it is west of $\mathbf{x}$. The $\texttt{f\_pent}$ is itself a composite on the points $y_1, y_2, y_3, y_4, y_5$, which constrains them using primitive concepts such as $\texttt{south\_1\_same\_color}$, which in turn is defined using the relational concepts ($==, \neq$, etc) on points. The scene graph corresponding to the primitive concept level is an abstract representation of the grid and its points which can be directly rendered as an image for image classification experiments.

If the concept specification is inconsistent, the generator will fail and report an error. If the concept is disjunctive, another disjunct is chosen and the process is repeated until a successful graph can be generated. This generation procedure is therefore sensitive to disjunctions and may fail to terminate if there are too many. We timeout if that is the case, but all concepts in this paper generate quickly.

A full list of concepts, including examples of more complex concepts not used in this paper, is included in Appendix 1.

## 4 Experimental Setup

The task for all experiments is multi-class prediction on images generated from a small set of concepts (3-5). Our goal is to better understand how standard as well as explicitly relational models generalize compositionally, for images with clear relational and compositional structure. We focus in particular on two specific types of compositional generalization: *productivity* and *substitutivity* [HDMB20]. The experiments on productivity (composition length generalization) test how well these models are able to learn recursive concepts if trained on a small number of compositions. In the substitutivity experiments, we test whether a model can learn to generalize correctly to other objects of the same type as its argument(s).[3]

We choose an image size of 32x32 similar to CIFAR10 [K$^+$09] and five colors: blue, red, green, yellow and white. We compare an MLP; a 2-layers CNN; ResNet18 [HZRS16]; PrediNet [SNC$^+$19], which uses multi-head attention to extract relational structure; and WReN [Ada18], which uses relation network modules [SRB$^+$17]. These baselines were chosen to provide a balance between well-known architectures and recent ones with relational inductive biases. All models have approximately the same number of parameters for fair comparison, and we performed hyper-parameter optimization on each model, starting with published values for Predinet, ResNet and WReN.

When performing multiple tests with the same trained model, such as in Sec 5.1.2 where we test on sequence length 3 and 5, we keep the same test data for concepts which do not change. In Sec 5.1.2, the "2x2 colored square" concept does not depend on sequence length: we use the same samples for both tests. The statistics for such classes therefore do not change between these variations. We report F1 score rather than accuracy as it better reflects performance on false positives and negatives. All numerical results shown in this section were averaged over 10 random seeds. We refer the reader to the appendix for more information (number of training samples etc.).

## 5 Evaluation

### 5.1 Productivity Experiments: Generalization to Longer Compositions

#### 5.1.1 Pure and Mixed Sequences of 2x2 Squares

This experiment evaluates the ability of a model to learn a recursive concept: $\texttt{east}_1(x_1, x_2)$, which requires at least one point in $x_2$ to be 1 grid point east of one point in $x_1$. We define 2x2 red or blue squares (see Table 1) and use them to construct horizontal sequences, by composing the $\texttt{east}_1$ concept on its first argument. We consider three sequence concepts: *all red*, *all blue*, and *mixed red and blue*. We train on sequences of length 1 (a single red or blue square) and 2, then test on lengths 3, 5, and 7. Aggregated results are shown in Fig. 2b; per-concept results are available in Appendix 2.

ResNet performs best overall but experiences a limited but noticeable degradation in generalization over longer sequences. Pretraining on ImageNet [DDS$^+$09] appears to be beneficial, particularly for longer sequences, presumably because pretraining helps the model learn more robust visual features which are helpful for the more complex test sequences. We use only the pretrained version of ResNet in subsequent experiments. The CNN, while architecturally simpler than ResNet, shows good performance but struggles on longer sequences. WReN, which uses a CNN to extract features, does better, particularly on mixed sequences. PrediNet performs similarly to the CNN, but suffers from high variance (we observed that PrediNet requires a significant amount of hyper-parameter tuning to reduce variance over different seeds). The MLP appears to have learned to average pixel values to classify the image. As the mixed sequences contain both red and blue squares, the MLP cannot discriminate them and labels them as all red or all blue.

The sequences discussed here are spatial, constructed by composing the $\texttt{east}_1$ relation. To classify them correctly, models need to take the whole sequence into consideration, which is challenging as the length increases. While the WReN relational model performs reasonably well, PrediNet does

---

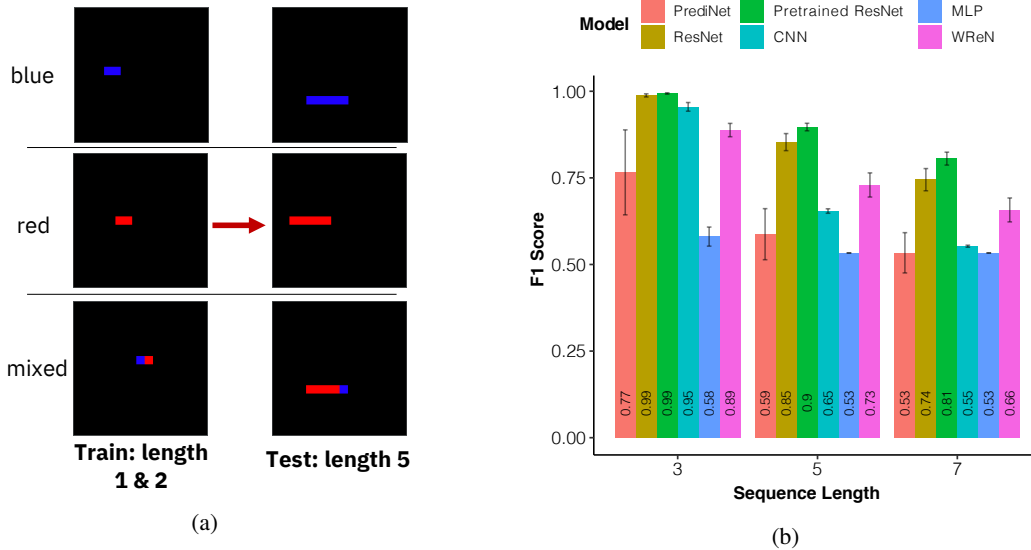[3]We will release the dataset upon publication.

Figure 2: (a) Some train and test samples for experiment 1 (pure and mixed sequences, Sec. 5.1.1). (b) F1 scores, aggregated over all concepts (sequence types). Error bars are 95% CI over 10 seeds.

poorly. The task here involves composing a single spatial relation and the images contain only a single object. In the next experiment, we try compositions involving both a spatial and color relation and multiple objects.

### 5.1.2 Box-World

In this experiment, we recreate a simplified version of the Box-World environment of [BHB+18], which is a grid-world domain with keys, locks and a gem. In the original reinforcement learning version of the game, the agent needs to find an initial (free) key and perform a series of "unlock" / "lock" steps, to find the gem. In our version, the task is to distinguish images with valid sequences from those with invalid, distractor, sequences, which can't be solved because the chain is broken somewhere with a key that doesn't open any locks. The problem is rendered as an image by representing keys, locks, and the gem as 2x2 squares (pixels); representing a paired lock and key by a common color; and a locked object (key or gem) as a square 1 pixel west of the lock. See Appendix 3 for an example.

This experiment tests a more complex sequence concept than experiment 1, as it requires learning to recurse on a conjunction of a spatial relation $\texttt{west}_1$ (for "locks") and non-spatial $\texttt{same\_color}$ (for "unlocks").

The concepts to be distinguished here are $\texttt{solution}$ when a path leads to a gem, $\texttt{distractor}$ when a path does not lead to a gem, and $\texttt{2x2 square}$, the "pixel" or base element of the paths. We train the models with paths of length 1 and 2 and test on lengths 3 and 5. The results are shown in Fig. 3a. Full results are available in Appendix 3.

All models, except PrediNet, retain the ability to identify the 2x2 square concept they were trained on. However, most models show poor generalization ability on learning longer $\texttt{solution}$ and $\texttt{distractor}$ sequences, even for length 3. WReN performs best (although not far from random), beating the pretrained ResNet, suggesting that it makes better use of the relational biases in its architecture. Nevertheless, it has difficulty on the $\texttt{distractor}$ concept. This is clearly visible in its confusion matrix (Fig. 3b), where it can be seen that it mistakes most of the distractor paths for valid ones and mistakes fewer valid ones for distractor ones.

### 5.2 Substitutivity Experiments: Scaling up 2x2 Patterned Squares

In contrast to the previous two experiments where we test if models can learn to generalize to longer sequences of the *same* object, the two experiments below test whether a concept generalizes to
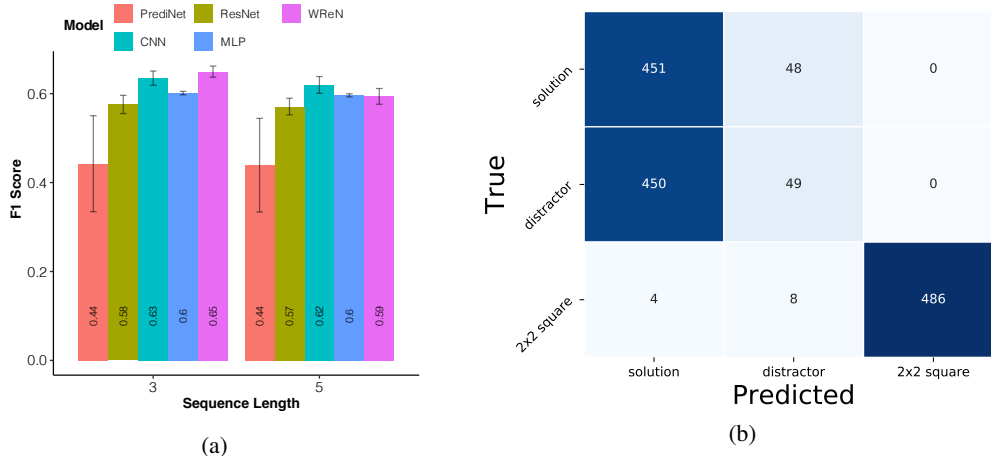
Figure 3: Results for Experiment 2 (Box-World sequences). (a): Generalization (zero-shot transfer) over longer sequences of keys and locks. (b) Confusion Matrix of WReN on sequences of length 5 (averaged over 10 runs).

different objects (still satisfying the same concept) than seen in training. We consider variations on the number of arguments to substitute and the type of arguments to substitute.

### 5.2.1 Generalizing from 2x2 to 4x4 patterned squares

In this experiment, we create concepts which are 2x2 squares (see Table 1) of four classes: *all blue*, *all red*, *vertical alternating red/blue stripe*s, and a *checkerboard pattern of red/blue*. The concept we want to test is that of a 2x2 square composed of 4 identical smaller squares – one for each quadrant. For the training concepts, the "squares" being composed are 4 points (we consider these 1x1 squares here). We test whether a model can learn to generalize this relation by substituting 2x2 squares for the points creating a 2x2 square of 2x2 squares, which is a 4x4 square of points. The 2x2 squares to be substituted are the solid red and blue ones given in training. This substitution corresponds visually to "scaling up". We assign these scaled up 4x4 squares the same concept ids as the corresponding 2x2 versions. Figure 4a shows some examples.

We observe (table of F1 scores available in Appendix 4) that no model generalizes to the striped and checkered 4x4 squares. Among the models, ResNet, CNN and MLP are able to recognize the 4x4 all blue and all red squares after being trained on corresponding 2x2 ones. However, they aren't able to use the 2x2 squares compositionally to generalize to the 4x4 stripes or checkerboard. As an example, Fig. 4b shows the MLP's confusion matrix, which suggests that it has learned a strategy that generalizes to some degree but not systematically.

### 5.2.2 Generalizing sequences under element substitutions

In this experiment, we test whether a model trained on a concept `c` and sequences of instances $e \in$ `c`, can generalize to sequences of elements $e' \in$ `c` for $e \neq e'$. Figure 5a shows examples where the element class is concept `type 1` which consists of a 2x2 red square and an f-pentomino; or `type 2` which consists of a blue square and a z-pentomino. We evaluate the ability of the model to discriminate three types of sequences of these elements: pure `type 1`, pure `type 2`, and `mixed 1+2`. The model receives training on `Type 1`, `Type 2` and sequences involving squares from these classes. It must generalize to classify sequences involving *both* pentominoes and squares. There are five classes in total including the element types and sequences. For the pentominoes, color is ignored and only type matters.

We test generalization of the sequences as follows. A `type 1` sequence is created with both a red square and an f-pentomino; a `type 2` sequence is created with both a blue square and a z-pentomino; and a mixed sequence with a blue square and f-pentomino or a red square and z-pentomino. This tests the ability of the model to perform one substitution in the sequence. We also evaluate 2 substitutions,
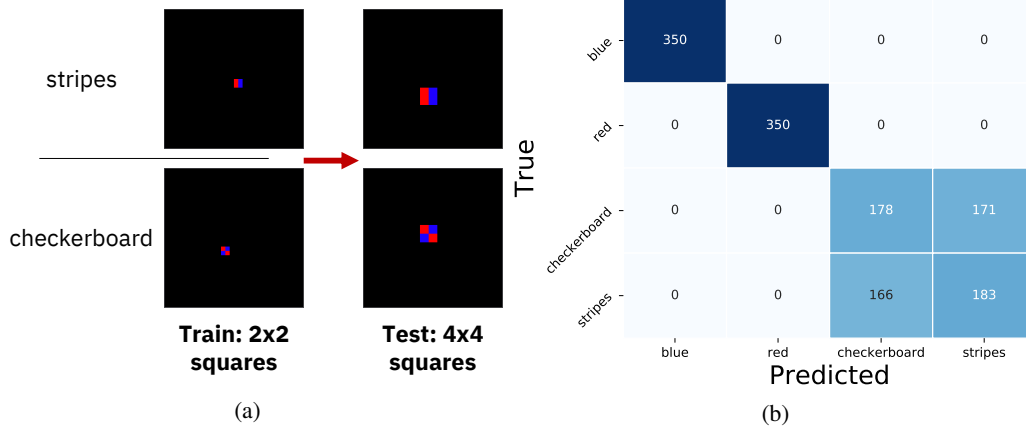
7

Figure 4: (a) Train and test examples of the patterned squares concepts in experiment 3 (Sec. 5.2.1). (b) Confusion matrix over the 4x4 squares concepts for the MLP model.

by changing both elements of the sequence. A `type 1` sequence is an f-pent pair; a `type 2` sequence is a z-pent pair; and a `mixed 1+2` sequence is an f-z pair.

To solve this problem robustly, the models need to learn to associate visually different objects as a common class and generalize a concept trained on some elements of this class to correctly classify new objects of the class not seen in training.



Figure 5: Experiment 4 (sequence substitutions, Sec. 5.2.2): (a) Some train and test samples of the types and sequences. (b) Average confusion matrix over 1-substitute pairs for ResNet: it recognizes true `type 1` or `type 2` examples, but also labels as such the pairs.

The models generally learn to recognize the `type 1` and `type 2` classes but with low precision, as they confuse the pairs for the single arguments. This is particularly visible with ResNet (Fig. 5b). One issue could be that learning the `type 1` and `type 2` concepts concurrently with the higher-order concepts may be too difficult. We tested a curriculum variant, where we trained until convergence on `type 1` and `type 2`, then added the remaining 3 sequence concepts. Except for ResNet, which improved slightly, the models all degraded. We hypothesize that a curriculum is not helpful without a composition mechanism to properly make use of it. Interestingly, while the relational models PrediNet and WReN do not surpass other models on the `type 1` and `type 2` classes, they perform relatively better on the sequences. Absolute performance remains poor.

### 5.3 Discussion

All productivity experiments show degradation in performance of all models as composition length increases, indicating they have not learned a recursive generalization. For the Box-World experiment,

8

no models do well (F1 scores are around 0.6 on sequence length 3 and performance degrades slightly for length 5). We note however that the pretrained ResNet performs well, only to be surpassed by WReN on the Box-World solution paths. For the relational models, WReN performs reasonably well but PrediNet does poorly. This indicates to us that their relational bias is perhaps not helpful without an additional compositional bias to employ the learned relations recursively. For the substitutivity experiments, while all models fail to properly generalize, WReN and PrediNet perform relatively better with the elements substitutions. These tasks are qualitatively harder and seem to require a more explicit compositional bias to generalize at all.

## 6  Conclusion

In this paper, we investigated the problem of compositional, relational generalization as a multi-class image recognition problem. We introduced a concept specification language, which allows a description of hierarchical concepts on a grid of colored points, and sketched the generation process that renders them into images. We found that having a declarative definition of the concepts facilitated more agile experimentation and concept sharing. Our experiments on compositional productivity and substitutivity provide evidence that, without specific biases for both relation representation as well as composition, neural networks do not generalize well in this setting and suffer from the same degradation on composition length as seen in experiments on text and multi-modal data.

# References

[AAL+15]     Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. VQA: Visual Question Answering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2425–2433, 2015.

[Ada18]      Adam Santoro and Felix Hill and David G. T. Barrett and Ari S. Morcos and Timothy P. Lillicrap. Measuring Abstract Reasoning in Neural Networks. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 4477–4486. PMLR, 2018.

[BHB+18]     Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational Inductive Biases, Deep Learning, and Graph Networks. *arXiv preprint arXiv:1806.01261*, 2018.

[CGLG18]     Michael B Chang, Abhishek Gupta, Sergey Levine, and Thomas L Griffiths. Automatically Composing Representation Transformations as a means for Generalization. *Proceedings of the International Conference on Learning Representations*, 2018.

[DDS+09]     Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A Large-Scale Hierarchical Image Database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. Ieee, 2009.

[Dre18]      Drew A. Hudson and Christopher D. Manning. Compositional Attention Networks for Machine Reasoning. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

[FP+88]      Jerry A Fodor, Zenon W Pylyshyn, et al. Connectionism and Cognitive Architecture: A Critical Analysis. *Cognition*, 28(1-2):3–71, 1988.

[HDMB20]     Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. Compositionality Decomposed: How do Neural Networks Generalise? *Journal of Artificial Intelligence Research*, 67:757–795, 2020.

[HM19a]      Drew Hudson and Christopher D Manning. Learning by Abstraction: The Neural State Machine. In *Advances in Neural Information Processing Systems 32*, pages 5903–5916. Curran Associates, Inc., 2019.

[HM19b]      Drew A Hudson and Christopher D Manning. GQA: A New Dataset for Real-World Visual Reasoning and Compositional Question Answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6700–6709, 2019.

[HSM+18]     Irina Higgins, Nicolas Sonnerat, Loic Matthey, Arka Pal, Christopher P. Burgess, Matko Bosnjak, Murray Shanahan, Matthew Botvinick, Demis Hassabis, and Alexander Lerchner. SCAN: Learning Hierarchical Compositional Visual Concepts. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

[HZRS16]     Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.

[JHvdM+17]   Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2901–2910, 2017.

[Jia19]      Jiayuan Mao and Chuang Gan and Pushmeet Kohli and Joshua B. Tenenbaum and Jiajun Wu. The Neuro-Symbolic Concept Learner: Interpreting Scenes, Words, and Sentences From Natural Supervision. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[JJVDM16]     Allan Jabri, Armand Joulin, and Laurens Van Der Maaten. Revisiting Visual Question Answering Baselines. In *European Conference on Computer Vision*, pages 727–739. Springer, 2016.

[JKS+15]      J. Johnson, R. Krishna, M. Stark, L. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei. Image Retrieval using Scene Graphs. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3668–3678, 2015.

[K+09]        Alex Krizhevsky et al. Learning Multiple Layers of Features from Tiny Images. *University of Toronto*, 2009.

[KB15]        Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.

[LB18]        Brenden M. Lake and Marco Baroni. Generalization without Systematicity: On the Compositional Skills of Sequence-to-Sequence Recurrent Networks. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 2879–2888. PMLR, 2018.

[Lou18]       Loula, João and Baroni, Marco and Lake, Brenden. Rearranging the Familiar: Testing Compositional Generalization in Recurrent Networks. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 108–114. Association for Computational Linguistics, 2018.

[Man19]       Manjunatha, Varun and Saini, Nirat and Davis, Larry S. Explicit Bias Discovery in Visual Question Answering Models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9562–9571, 2019.

[Nav77]       David Navon. Forest Before Trees: The Precedence of Global Features in Visual Perception. *Cognitive Psychology*, 353:383, 1977.

[PGM+19]      Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

[RCRK19]      Clemens Rosenbaum, Ignacio Cases, Matthew Riemer, and Tim Klinger. Routing networks and the challenges of modular and compositional computation. *arXiv preprint arXiv:1904.12774*, 2019.

[RKR18]       Clemens Rosenbaum, Tim Klinger, and Matthew Riemer. Routing Networks: Adaptive Selection of Non-Linear Functions for Multi-Task Learning. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

[SK07]        Elizabeth S Spelke and Katherine D Kinzler. Core Knowledge. *Developmental science*, 10(1):89–96, 2007.

[SNC+19]      Murray Shanahan, Kyriacos Nikiforou, Antonia Creswell, Christos Kaplanis, David Barrett, and Marta Garnelo. An Explicitly Relational Neural Network Architecture. *arXiv preprint arXiv:1905.10307*, 2019.

[SRB+17]      Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. A Simple Neural Network Module for Relational Reasoning. In *Advances in Neural Information Processing Systems*, pages 4967–4976, 2017.

[ZGSS+16]     Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Yin and Yang: Balancing and Answering Binary Visual Questions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5014–5022, 2016.

[ZRS+18]   Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, et al. Deep Reinforcement Learning with Relational Inductive Biases. *Proceedings of the International Conference on Learning Representations*, 2018.

# 1 List of Concepts and Details of the Generation Algorithm

## 1.1 List of concepts

Table 2 contains the definition of the concepts we used in our experiments. For vector objects, we use the notation $x :: \mathbf{y}$ to mean the function which splits the vector into its first element ($x$) and the rest of the elements ($\mathbf{y}$). More generally, we can split off $k$ elements from the front of the sequence by writing $(x_1, x_2, \cdots, x_k) :: \mathbf{y}$. We write singleton vectors $\mathbf{x}$ as $\mathbf{x} == (x)$. Quantification over vectors is interpreted to mean over the set of points of the vector, forgetting the structure. We allow two reduction operations on integer-valued properties of a vector: $\arg\min$ and $\arg\max$. For example, $\arg\min_{x \in \mathbf{x}}(x.x\_loc)$ produces the the element $x$ in $\mathbf{x}$ for which $x.x\_loc$ is minimal. We have defined $\texttt{east}_1$ to use these functions rather than conjunction (as in the main paper) for increased efficiency in generation.

| Concept | Definition |
|---|---|
| point | $\mathtt{point}(x) \equiv (\mathtt{red}(x) \lor \mathtt{blue}(x) \lor \cdots \lor \mathtt{yellow}(x)) \land$<br>$(x.x\_loc == 0 \lor \cdots \lor x.x\_loc == \mathtt{GRID\_SIZE} - 1) \land$<br>$(x.y\_loc == 0 \lor \cdots \lor x.y\_loc == \mathtt{GRID\_SIZE} - 1)$ |
| red (point) | $\mathtt{red\_point}(x) \equiv \mathtt{point}(x) \land x.color == \mathtt{RED}$ |
| red (object) | $\mathtt{red}(\mathbf{x}) \equiv \forall x \in \mathbf{x}\ \mathtt{red\_point}(x)$ |
| object with red or blue points | $\mathtt{red\_or\_blue}(\mathbf{x}) \equiv \forall x \in \mathbf{x}\ [\mathtt{red}(x) \lor \mathtt{blue}(x)]$ |
| same color (point) | $\mathtt{same\_color\_point}(x_1, x_2) \equiv x_1.color == x_2.color \land \mathtt{point}(x_1) \land \mathtt{point}(x_2)$ |
| east 1 grid point (point) | $\mathtt{east\_point}_1(x_1, x_2) \equiv x_2.x\_loc == x_1.x\_loc + 1 \land$<br>$x_2.y\_loc == x_1.y\_loc \land \mathtt{point}(x_1) \land \mathtt{point}(x_2)$ |
| east (point) | $\mathtt{east\_point}(x_1, x_2) \equiv x_2.x\_loc > x_1.x\_loc \land x_2.y\_loc == x_1.y\_loc$ |
| east (object) | $\mathtt{east}(\mathbf{x}_1, \mathbf{x}_2) \equiv x_1 == \arg\max_{x \in \mathbf{x}_1}\{x.x\_loc\} \land x_2 == \arg\min_{y \in \mathbf{x}_2}\{y.x\_loc\} \land \mathtt{east\_point}(x_1, x_2)$ |
| east 1 grid point (object) | $\mathtt{east}_1(\mathbf{x}_1, \mathbf{x}_2) \equiv \mathtt{east}(\mathbf{x}_1, \mathbf{x}_2) \land \exists x \in \mathbf{x}_1, y \in \mathbf{x}_2\ \mathtt{east\_point}_1(x, y)$ |
| east 1 sequence | $\mathtt{east\_seq}(\mathbf{x}) \equiv [\mathbf{x} == (x)] \lor [\mathbf{x} == x :: \mathtt{xs} \land \mathtt{xs} == y :: \mathtt{ys} \land \mathtt{east\_point}_1(x, y) \land \mathtt{east\_seq}(\mathtt{xs})]$ |
| red east sequence (Exp 1) | $\mathtt{red\_east\_seq}(\mathbf{x}) \equiv \mathtt{east\_seq}(\mathbf{x}) \land \mathtt{red}(\mathbf{x})$ |
| red or blue east sequence (Exp 1) | $\mathtt{red\_or\_blue\_east\_seq}(\mathbf{x}) \equiv \mathtt{east\_seq}(\mathbf{x}) \land \mathtt{red\_or\_blue}(\mathbf{x})$ |
| north 1 grid point (point) | $\mathtt{north\_point}_1(x_1, x_2) \equiv x_2.y\_loc == x_1.y\_loc - 1 \land x_2.x\_loc == x_1.x\_loc$<br>$\land \mathtt{point}(x_1) \land \mathtt{point}(x_2)$ |
| north 1 grid point (object) | $\mathtt{north}_1(\mathbf{x}, \mathbf{y}) \equiv \forall x \in \mathbf{x}, y \in \mathbf{y}\ \mathtt{north\_point}_1(x, y)$ |
| 2x2 square of points (Exp 1,2,3,4) | $\mathtt{2x2\_square\_point}(\mathbf{x}) \equiv \mathtt{east}_1(x_1, x_2) \land \mathtt{south}_1(x_2, x_3) \land \mathtt{west}_1(x_3, x_4) \land \mathtt{north}_1(x_4, x_1) \land$<br>$\mathbf{x} == (x_1, x_2, x_3, x_4) \land \mathtt{point}(x_1) \land \mathtt{point}(x_2) \land \mathtt{point}(x_3) \land \mathtt{point}(x_4)$ |
| 2x2 checkerboard (points) (Exp 3) | $\mathtt{2x2\_checkerboard}(\mathbf{x}) \equiv \mathtt{2x2\_square\_point}(\mathbf{x}) \land \mathbf{x} == (x_1, x_2, x_3, x_4) \land \mathtt{red}(x_1) \land$<br>$\mathtt{blue}(x_2) \land \mathtt{red}(x_3) \land \mathtt{blue}(x_4)$ |
| 2x2 vertical stripes (points) (Exp 3) | $\mathtt{2x2\_vert\_stripe}(\mathbf{x}) \equiv \mathtt{2x2\_square\_point}(\mathbf{x}) \land \mathbf{x} = (x_1, x_2, x_3, x_4)$<br>$\land \mathtt{red}(x_1) \land \mathtt{blue}(x_2) \land \mathtt{blue}(x_3) \land \mathtt{red}(x_4)$ |
| square shape | $\mathtt{square}(\mathbf{x}) \equiv \mathbf{x} == (x_1, x_2, x_3, x_4) \land \mathtt{point}(x_1) \land \mathtt{point}(x_2) \land \mathtt{point}(x_3) \land \mathtt{point}(x_4) \land$<br>$\mathtt{east}(x_1, x_2) \land \mathtt{south}(x_2, x_3) \land \mathtt{west}(x_3, x_4) \land \mathtt{north}(x_4, x_1)$ |
| 2x2 square of squares (Exp 3) | $\mathtt{2x2\_square\_of\_squares}(\mathbf{x}) \equiv \mathbf{x} = (\mathbf{x_1}, \mathbf{x_2}, \mathbf{x_3}, \mathbf{x_4}) \land$<br>$\mathtt{square}(\mathbf{x_1}) \land \mathtt{square}(\mathbf{x_2}) \land \mathtt{square}(\mathbf{x_3}) \land \mathtt{square}(\mathbf{x_4}) \land$<br>$\mathbf{x_1} == (ul_1, ur_1, lr_1, ll_1) \land \mathbf{x_2} == (ul_2, ur_2, lr_2, ll_2) \land$<br>$\mathbf{x_3} == (ul_3, ur_3, lr_3, ll_3) \land \mathbf{x_4} == (ul_4, ur_4, lr_4, ll_4) \land$<br>$\mathtt{east}_1(ur_1, ul_2) \land \mathtt{east}_1(lr_1, ll_2) \land \mathtt{south}_1(ll_2, ul_3) \land$<br>$\mathtt{south}_1(lr_2, ur_3) \land \mathtt{west}_1(ul_3, ur_4) \land \mathtt{west}_1(ll_3, lr_4)$ |
| 2x2 square of squares checkerboard (Exp 3) | $\mathtt{2x2\_square\_of\_squares\_checkerboard}(\mathbf{x}) \equiv \mathtt{2x2\_square\_of\_squares}(\mathbf{x}) \land$<br>$\mathbf{x} == (\mathbf{x_1}, \mathbf{x_2}, \mathbf{x_3}, \mathbf{x_4}) \land \mathtt{red}(\mathbf{x_1}) \land \mathtt{blue}(\mathbf{x_2}) \land \mathtt{red}(\mathbf{x_3}) \land \mathtt{blue}(\mathbf{x_4})$ |
| adjacency | $\mathtt{adj}(x_1, x_2) \equiv \mathtt{north}_1(x_1, x_2) \lor \mathtt{east}_1(x_1, x_2) \lor \mathtt{south}_1(x_1, x_2) \lor \mathtt{west}_1(x_1, x_2)$ |
| pentomino | $\mathtt{pentomino}(\mathbf{x}) \equiv \mathtt{adj\_same\_color}(x_1, x_2) \land \mathtt{adj\_same\_color}(x_2, x_3) \land$<br>$\mathtt{adj\_same\_color}(x_3, x_4) \land \mathtt{adj\_same\_color}(x_4, x_5) \land \mathbf{x} == (x_1, x_2, x_3, x_4, x_5)$ |
| "f"-pentomino (Exp 4) | $\mathtt{f\text{-}pent}(\mathbf{x}) \equiv \mathtt{east}_1\_\mathtt{same\_color}(x_1, x_2) \land \mathtt{south}_1\_\mathtt{same\_color}(x_4, x_1) \land$<br>$\mathtt{west}_1\_\mathtt{same\_color}(x_4, x_3) \land \mathtt{south}_1\_\mathtt{same\_color}(x_4, x_5) \land \mathbf{x} == (x_1, x_2, x_3, x_4, x_5)$<br>$\land \mathtt{point}(x_1) \land \mathtt{point}(x_2) \land \mathtt{point}(x_3) \land \mathtt{point}(x_4) \land \mathtt{point}(x_5)$ |
| "z" pentomino (Exp 4) | $\mathtt{z\text{-}pent}(\mathbf{x}) \equiv \mathtt{east}_1\_\mathtt{same\_color}(x_1, x_2) \land \mathtt{south}_1\_\mathtt{same\_color}(x_2, x_3) \land$<br>$\mathtt{south}_1\_\mathtt{same\_color}(x_3, x_4) \land \mathtt{east}_1\_\mathtt{same\_color}(x_4, x_5) \land \mathbf{x} == (x_1, x_2, x_3, x_4, x_5)$<br>$\land \mathtt{point}(x_1) \land \mathtt{point}(x_2) \land \mathtt{point}(x_3) \land \mathtt{point}(x_5)$ |
| key<br>lock<br>gem | $\mathtt{key}(\mathbf{x}) \equiv \mathtt{2x2\_square\_point}(\mathbf{x});$<br>$\mathtt{lock}(\mathbf{x}) \equiv \mathtt{2x2\_square\_point}(\mathbf{x});$<br>$\mathtt{gem}(\mathbf{x}) \equiv \mathtt{2x2\_square\_point}(\mathbf{x})$ |
| key or gem | $\mathtt{key\_or\_gem}(\mathbf{x}) \equiv \mathtt{key}(\mathbf{x}) \lor \mathtt{gem}(\mathbf{x})$ |
| locked object | $\mathtt{locks}(\mathbf{l}, \mathbf{o}) \equiv \mathtt{key\_or\_gem}(\mathbf{o}) \land \mathtt{lock}(\mathbf{l}) \land \mathtt{east}_1(\mathbf{o}, \mathbf{l})$ |
| key unlocks lock | $\mathtt{unlocks\_lock}(\mathbf{k}, \mathbf{l}) \equiv \mathtt{key}(\mathbf{k}) \land \mathtt{lock}(\mathbf{k}) \land \mathtt{same\_color}(\mathbf{k}, \mathbf{l})$ |
| solution (Exp 2) | $\mathtt{solution}(\mathbf{s}) \equiv [\mathbf{s} == (\mathbf{g}) \land \mathtt{gem}(\mathbf{g})] \lor$<br>$\mathbf{s} == \mathbf{k} :: \mathbf{s}_1 \land \mathbf{s}_1 == \mathbf{l} :: \mathbf{s}_2 \land \mathtt{unlocks\_lock}(\mathbf{k}, \mathbf{l}) \land \mathtt{solution}(\mathbf{s}_2)$ |
| distractor position 2 (Exp 2) | $\mathtt{distractor}(\mathbf{s}) \equiv \mathbf{s} == (\mathbf{k}_1, \mathbf{l}_1) :: \mathbf{s}_1 \land$<br>$\mathtt{unlocks\_lock}(\mathbf{k}_1, \mathbf{l}_1) \land \mathbf{s}_1 == (\mathbf{k}_2, \mathbf{l}_2) :: \mathbf{s}_2 \land$<br>$\mathtt{locks}(\mathbf{l}_1, \mathbf{k}_2) \land \mathtt{not\_unlocks\_lock}(\mathbf{k}_2, \mathbf{l}_2) \land \mathtt{solution}(\mathbf{s}_2)$ |
| composite x pentomino (i,u,f,z,x) Appx Fig 1.(c) | $\mathtt{composite\_x\_pent}(\mathbf{x}) \equiv \mathtt{x\text{-}pent}(\mathbf{x}) \land \mathbf{x} == (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5) \land$<br>$\mathtt{i\text{-}pent}(\mathbf{x}_1) \land \mathtt{u\text{-}pent}(\mathbf{x}_2) \land \mathtt{f\text{-}pent}(\mathbf{x}_3) \land \mathtt{z\text{-}pent}(\mathbf{x}_4) \land \mathtt{x\text{-}pent}(\mathbf{x}_5) \land$<br>$\mathtt{south}_1(\mathbf{x}_1, \mathbf{x}_3) \land \mathtt{west}_1(\mathbf{x}_3, \mathbf{x}_2) \land \mathtt{east}_1(\mathbf{x}_3, \mathbf{x}_4) \land \mathtt{south}_1(\mathbf{x}_3, \mathbf{x}_5)$ |

Table 2: The definition of the concepts used in our experiments. The left column contains a short textual description and in which experiment(s) was the concept used.

## 1.2 Details of the Generation Algorithm

For clarity, we have separated the description of the algorithms in several parts, from high-level structure (Fig. 6) to specific functions (Fig. 8).

---

**Algorithm 1:** `generate_concept`

---

**Input:** c : Concept

```
/* A Concept has:  (1) a name; (2) an argument; (3) a definition.
we assume a single arg (unary) because if a concept has more we can always
wrap them up as a vector at a higher level.  Example:
c(x) := x=(x1,x2) and red(x_1) and blue(x_2) here the name is "c"
the argument is "x" and the definition is "red(x_1) and blue(x_2).

The constituent elements of x are x1 and 2.  */


/* create a map of variables to their bindings,

which are vectors of their constituent variables or point objects */
```
bindings = Map()

```
/* convert concept to disjunctive normal form (DNF) with an OR of conjunctive concepts */
```
conj_concepts = convert_to_dnf(c)

```
/* generate each conjunctive concept.  if we fail; try another */
```
**forall** conj_concept ∈ conj_concepts **do**
   **try**
      └ generate_conjunctive(conj_concepts, bindings)
   **catch**
      └ continue

---

Figure 6: Generation Algorithm: Iterate over the different conjunctive clauses of a concept until one can be successfully generated.

**Algorithm 2:** `generate_conjunctive_concept`

**Input:** c : Concept; bindings : Map

```
/* create a graph from the conjunctive concept definition:
variables become vertices; unary relations become types associated with the vertex;
binary relations become edges between vertices.
We assume connectivity here but it's not a hard requirement.
For simplicity, we assume one unary relation on any node and at most one edge between nodes.
If there are, for example, 2 unary relations on a node, say c(x), d(x), then these can be
grouped as e(x) := c(x) ∧ d(x).  Similarly for binary relations.  So this is w.l.o.g.
 */
```
$g = \text{create\_graph}(c)$

```
/* pick a variable to act as the root of the search,
chosen randomly but bound variables are prioritized.  */
```
$\text{root\_variable} = \text{pick\_root}(g)$

```
/* get the root variable concept and the vector of variables used
in the root variable concept definition */
```
$\text{root\_concept} = \text{get\_unary\_concept}(\text{root\_variable})$
$\text{root\_variable\_elements} = \text{get\_elements}(\text{root\_concept})$

```
/* if no constituent elements (i.e.  primitive concept), then create a point, set its properties and bind
    it.  */
```
**if** root_variable_elements.is_empty*()* **then**

> `/* function defined below */`
> $\text{generate\_primitive\_conjunctive\_unary\_relation}(c, \text{bindings})$
> **return** bindings

```
/* composite concept:  add binding of the root variable to its elements in the binding map */
```
$\text{bindings}[\text{root\_variable}] = \text{root\_variable\_elements}$

```
/* create a queue for the BFS and add the root variable to it */
```
$q = \text{Queue}()$
$\text{q.push}(\text{root\_variable})$

```
/* perform BFS on the graph starting from the root variable.
Mark nodes to avoid revisiting.  Fail if there is an inconsistency */
```
$\text{visited} = \text{Set}()$
**while** q **do**

> `/* get next variable */`
> $v = \text{q.pop}()$
>
> `/* mark it as visited */`
> $\text{visited.add}(v)$
>
> `/* get concept for current variable */`
> $\text{concept} = \text{get\_unary\_concept}(v)$
>
> `/* generate the concept definition (recursively) */`
> $\text{generate\_conjunctive}(\text{concept}, \text{bindings})$
>
> `/* generate any binary concepts r */`
> **forall** edges $(v, w, \text{binary\_concept}) \in g \wedge w \notin \text{visited}$ **do**
>> $\text{binary\_concept\_definition} = \text{get\_definition}(\text{binary\_concept})$
>> $\text{generate\_binary\_concept}(\text{binary\_concept}, \text{bindings})$
>> $\text{q.push}(w)$

**return** bindings

Figure 7: Generation Algorithm for a given conjunctive concept.

---

**Algorithm 3:** `generate_primitive_conjunctive_unary_relation`

---

**Input:** c : PrimitiveConjunctiveClause, bindings : Map

`/* generates a point object with properties determined by the supplied concept (e.g.`
   $x \equiv \text{red}(x) \land x.x\_loc == 0 \land x.y\_loc == 16)$                            `*/`

$p = \text{Point}()$ ;
$\text{bindings}[\mathbf{x}] = p$ ;
**forall** unary $\in$ c **do**

   |  `/* each primitive unary concept has its own generator implementation which sets the properties of the`
   |    `point.`                                 `*/`
   |  generate_primitive_unary_relation(unary, bindings)

---

---

**Algorithm 4:** `generate_binary_concept`

---

`/* a binary concept` $\mathbf{r(x, y)}$ `must be defined as a primitive relation using a reduction operator`
`on both` $\mathbf{x}$ `and` $\mathbf{y}$ `to reduce them to points.`
`for now these reduction operators are limited to argmin and argmax.`
`Example:  see east.  There, the first reduction operator is` $\arg\max$ `; the second is` $\arg\min$
`and the primitive relation on them is east_point */`


$\text{reduction}_1 = \text{get\_reduction}_1(c)$
$\text{reduction}_2 = \text{get\_reduction}_2(c)$
$\text{prim\_concept} = \text{get\_prim\_concept}(c)$
$x_1 = \text{reduction}_1(\mathbf{x_1})$
**forall** pairs $(x_1, y)$, y $\in$ **y do**

   |  **try**
   |    |  `/* the primitive generators must be defined individually for each primitive concept */`
   |    |  generate_binary_concept_primitive(prim_concept, $x_1$, $y$, bindings)
   |    |  `/* check that the reduction over y is satisfied by the generated bindings */`
   |    |  check_reduction($\mathbf{y}$, $\text{reduction}_2$, bindings)
   |  **catch**
   |    |  `continue`

---

Figure 8: Specific functions, used in Alg. 7.

## 1.3 Additional concept examples

Fig. 9 presents some examples of concepts which we generated with ConceptWorld but did not include in our experiments.
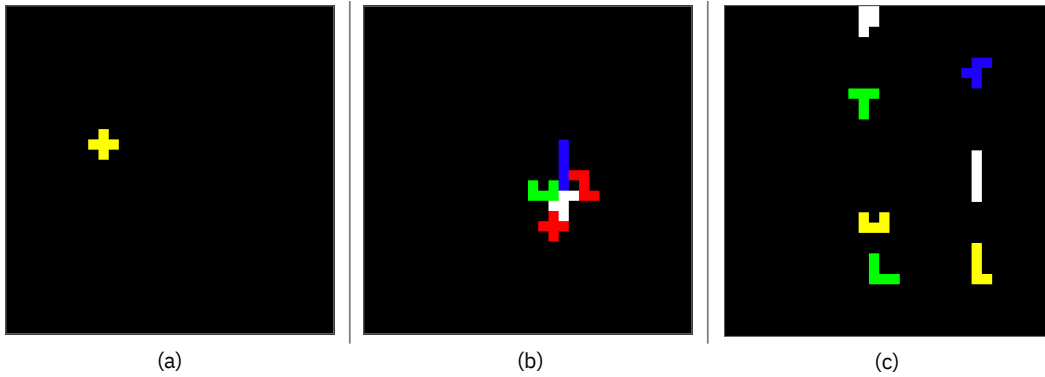


Figure 9: Some examples of additional concepts not included in our experiments. (a): the "X" pentomino. (b): An "X" pentomino shape ("+") made of "I", "U", "F", "Z", "X" pentominoes. (c): two loose stacks of pentominoes, one east of the other.

# 2 Experiment 1: Complete set of results and additional details

## 2.1 Training details and hyper-parameters search

We developed the codebase using PyTorch [PGM$^+$19]. All experiments being multi-class image classification tasks, we used Adam [KB15] and Cross-Entropy loss. We fixed the number of epochs to 100, and early-stopped, based on the validation loss. The default learning rate was $1^{-3}$ and the batch size 128. We used different GPU-accelerated platforms, with NVIDIA's GeForce GTX TITAN X and Tesla K80 GPUs.

For the CNN, MLP and ResNet [HZRS16], we started with published or common hyper-parameter values:

- CNN: 2 convolution layers, each with output size 128, kernel size 5, stride 1, padding 0; 2 maxpool layers of kernel size 2.
- MLP: 1 hidden layer of size 128, along with input and output layers.
- ResNet: we used ResNet18, for which we replaced the last 2 layers (layer 3 and 4) by identity functions. This was done to reduce the number of parameters, making it comparable to the other baselines.

For the 3 models above, we performed random hyper-parameter search, using Experiment 1 (pure and mixed sequences of 2x2 squares) over the learning rate, hidden size and number of layers, optimizing for validation accuracy. We found that the original hyper-parameters worked well for all experiments, and therefore used them for all reported experiments.

In contrast, we observed that WReN and PrediNet are sensitive to hyper-parameters. Starting with the published values, we performed random search over Experiment 1 – verifying that performance was equally good in other experiments – over the learning rate, number and size of the input convolutional layers, as well as the key size in the case of PrediNet.

We obtained, selecting based on validation accuracy, the following sets of hyperparameters for PrediNet and WReN:

PrediNet:

- Learning rate: $1^{-5}$,
- Training batch size: 64,
- Key size: 16,
- Number of attention heads: 2,
- Number of relations: 2,
- 1 input convolutional layer (with bias and batch normalization) with output size 32, kernel size 12, stride 6.

For WReN:

- Learning rate: $1^{-5}$,
- Training batch size: 128,
- Key size: 16,
- Number of attention heads: 2,
- 3 input convolutional layers (with batch normalization) with output size 64, kernel size 2, stride 2.

With these parameter values, all models have comparable size, between 400k and 600k trainable parameters.

## 2.2 Additional examples

Fig. 10 contains some examples of the sequences we considered in this experiment.
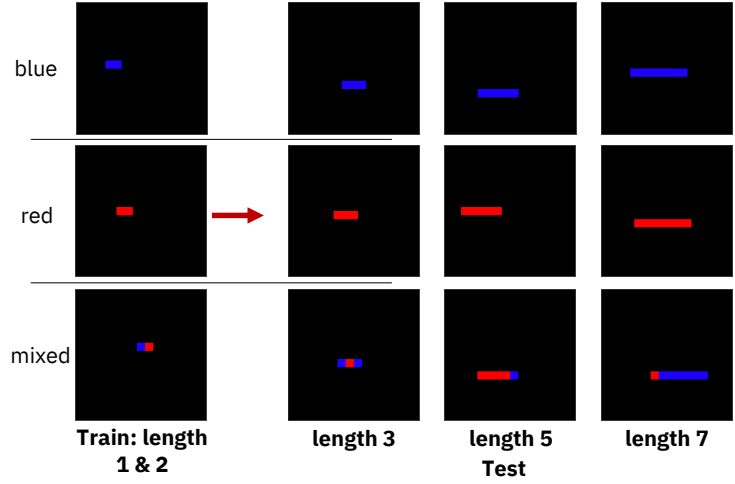
Figure 10: Some examples of the *all-red*, *all-blue* and *mixed* sequences considered in Experiment 1, for the training and test sets (*best viewed in color*).

## 2.3  Results

The number of samples per concept for the training and test sets is available in Table 3. The training set is an union over sequences of length 1 and 2, except for the *mixed* sequences, which can only be defined from length 2 and up. The test sets, one per sequence length, are independent from one another.

| Concept | Train | | Test | | |
|---|---|---|---|---|---|
| | Seq len 1 | Seq len 2 | Seq len 3 | Seq len 5 | Seq len 7 |
| red | 800 | 800 | 600 | 600 | 500 |
| blue | 800 | 800 | 600 | 600 | 500 |
| mixed | - | 1600 | 600 | 600 | 500 |

Table 3: Number of samples per concept and sequence length for Experiment 1. Due to the constant image size (32 x 32), the longer the sequence, the smaller the space of corresponding samples.

Table 4 shows F1 scores per class, sorted by model and sequence length. These scores were obtained by averaging over 10 runs, each with a different random seed.

| Model | Concept | F1 score on test set | | |
| --- | --- | --- | --- | --- |
| | | Seq len 3 | Seq len 5 | Seq len 7 |
| ResNet | Pretrained | | | |
| | Blue | 0.993 | 0.935 | 0.890 |
| | Red | 0.997 | 0.929 | 0.881 |
| | Mixed | 0.990 | 0.827 | 0.646 |
| | Non-Pretrained | | | |
| | Blue | 0.991 | 0.913 | 0.867 |
| | Red | 0.992 | 0.905 | 0.855 |
| | Mixed | 0.981 | 0.742 | 0.512 |
| WReN | Blue | 0.921 | 0.842 | 0.814 |
| | Red | 0.915 | 0.848 | 0.809 |
| | Mixed | 0.828 | 0.499 | 0.350 |
| SimpleConvNet | Blue | 0.964 | 0.829 | 0.803 |
| | Red | 0.945 | 0.844 | 0.822 |
| | Mixed | 0.929 | 0.304 | 0.05 |
| PrediNet | Blue | 0.795 | 0.709 | 0.689 |
| | Red | 0.798 | 0.733 | 0.717 |
| | Mixed | 0.705 | 0.321 | 0.195 |
| SimpleFeedForward | Blue | 0.805 | 0.800 | 0.800 |
| | Red | 0.818 | 0.800 | 0.800 |
| | Mixed | 0.119 | 0 | 0 |

Table 4: F1 scores per concept and test sequence length for all 5 baselines in Experiment 1.

# 3 Experiment 2: Complete set of results and additional details

Fig. 11 illustrates some of the valid and invalid $key - lock$ paths of different lengths we considered for Experiment 2 (Box-World).
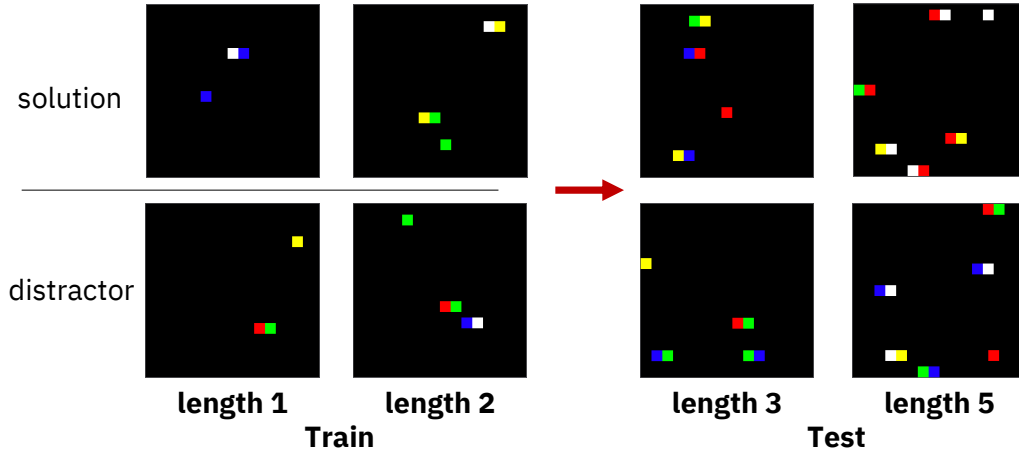


Figure 11: Some examples of the `solution` and `distractor` concepts in Experiment 2 (*best viewed in color*).

The number of samples per concept for the training and test sets is available in Table 5. The training set is an union over paths of length 1 (1 $key - lock$ pair) and 2 (2 $key - lock$ pairs). The test sets, one per sequence length, are independent from one another.

| Concept | Train | | Test | |
|---|---|---|---|---|
| | Seq len 1 | Seq len 2 | Seq len 3 | Seq len 5 |
| solution | 500 | 500 | 500 | 500 |
| distractor | 500 | 500 | 500 | 500 |
| 2x2 square | 1000 | | 500 | |

Table 5: Number of samples per concept and sequence length for Experiment 2. The `2x2 square` concept is independent of sequence length, and the same samples are used for both test sets.

Table 6 contains the F1 scores over the test sets, for the 5 baselines. These scores were also obtained by averaging over 10 runs.

| Model | Concept | F1 score on test set | |
| --- | --- | --- | --- |
| | | Seq len 3 | Seq len 5 |
| ResNet (Pretrained) | Solution | 0.301 | 0.303 |
| | Distractor | 0.427 | 0.411 |
| | 2x2 Square | 0.999 | 0.999 |
| WReN | Solution | 0.584 | 0.642 |
| | Distractor | 0.379 | 0.154 |
| | 2x2 Square | 0.986 | 0.986 |
| SimpleConvNet | Solution | 0.465 | 0.447 |
| | Distractor | 0.439 | 0.412 |
| | 2x2 Square | 1 | 1 |
| PrediNet | Solution | 0.361 | 0.368 |
| | Distractor | 0.416 | 0.398 |
| | 2x2 Square | 0.551 | 0.551 |
| SimpleFeedForward | Solution | 0.447 | 0.441 |
| | Distractor | 0.456 | 0.447 |
| | 2x2 Square | 0.901 | 0.901 |

Table 6: F1 scores per concept and test sequence length for all 5 baselines in Experiment 2. The test samples for the `2x2 square` concept are the same for both sequence lengths.

# 4 Experiment 3: Complete set of results and additional details

Fig. 12 illustrates some of the train and test samples we created for the *all blue*, *all red*, *vertical alternating red/blue stripes*, and *checkerboard pattern of red/blue* concepts in Experiment 3.
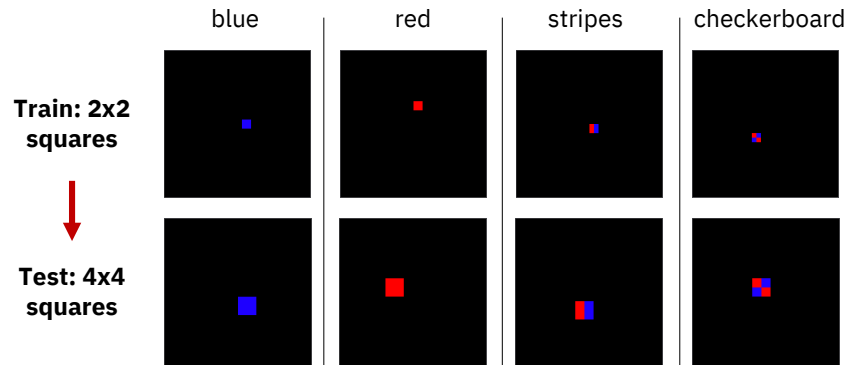


Figure 12: Some examples of the 4 concepts in Experiment 3 (*best viewed in color*).

Table 7 contains the number of samples for the training (2x2 squares) and test (4x4 squares) for each concept.

| Concept | Train (2x2 squares) | Test (4x4 squares) |
|---|---|---|
| all blue | 700 | 350 |
| all red | 700 | 350 |
| stripes | 800 | 350 |
| checkerboard | 800 | 350 |

Table 7: Number of samples per concept for Experiment 3.

Table 8 shows the F1 score per concept and model.

| Model | Concept | F1 score on test set |
|---|---|---|
| ResNet (Pretrained) | Blue | 0.7416 |
| | Red | 0.8295 |
| | Vertical Stripes | 0.3221 |
| | Checkerboard | 0.3807 |
| WReN | Blue | 0.4645 |
| | Red | 0.5293 |
| | Vertical Stripes | 0.10189 |
| | Checkerboard | 0.2415 |
| SimpleConvNet | Blue | 0.97879 |
| | Red | 0.984699 |
| | Vertical Stripes | 0.6614 |
| | Checkerboard | 0.3566 |
| PrediNet | Blue | 0.2388 |
| | Red | 0.3063 |
| | Vertical Stripes | 0.183 |
| | Checkerboard | 0.1833 |
| SimpleFeedForward | Blue | 0.9998 |
| | Red | 1 |
| | Vertical Stripes | 0.5157 |
| | Checkerboard | 0.5074 |

Table 8: F1 scores per concept for all 5 baselines for Experiment 3.

# 5  Experiment 4: Complete set of results and some examples

Fig. 13 showcases the `type 1` and `type 2` classes, as well as the pairs we built with them. Test pairs involve 1 or 2 substitutions.
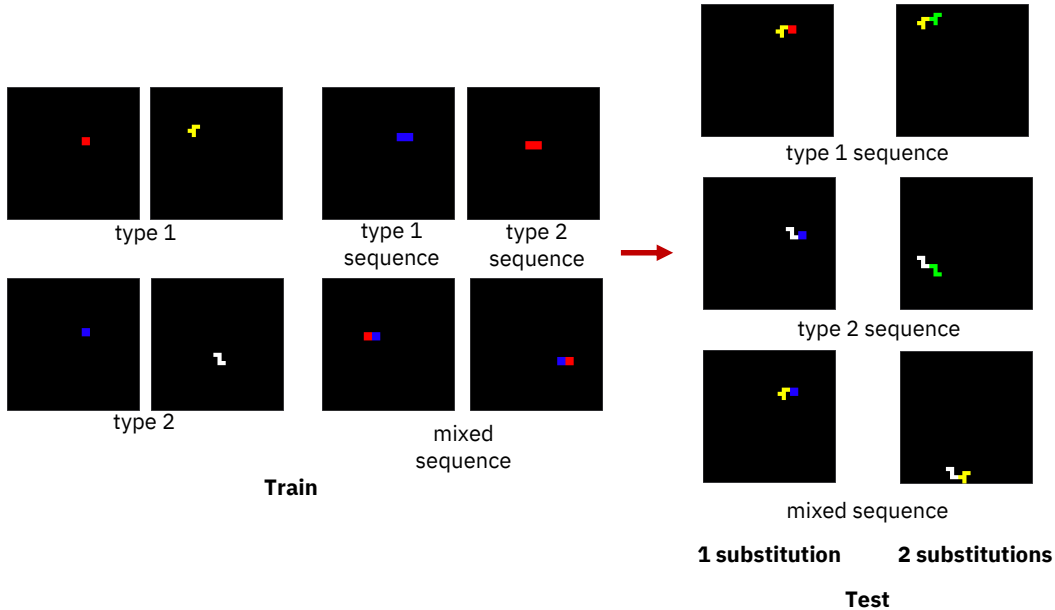


Figure 13: Some examples of the `type 1` and `type 2` classes, as well as the "pure" and mixed sequences built from both types (*best viewed in color*).

The details of the training and test sets are available in Table 9. The `type 1` and `type 2` classes are an union over 2 concepts (i.e. *2x2 red square* and *"F" pentomino* for `type 1`). The number of samples for these 2 classes is thus equidistributed between the 2 concepts.

| Concept | Train | Test | |
|---|---|---|---|
| | | 1 substitution | 2 substitutions |
| type 1 | 400 + 400 | 200 + 200 | |
| type 2 | 400 + 400 | 200 + 200 | |
| type 1 sequences | 800 | 400 | 400 |
| type 2 sequences | 800 | 400 | 400 |
| mixed sequences | 800 | 400 | |

Table 9: Number of samples per concept and number of substitutions for Experiment 4. The `type 1` and `type 2` concepts are independent of the number of substitutions: the same samples are used for both tests.

Table 10 contains the F1 scores per concept, for both test sets, and all baselines.

## 5.1  Curriculum Training

To further study the potential impact of learning the `type 1` and `type 2` concepts concurrently with the higher-order concepts, we tested a curriculum variant. Here, we trained the models until convergence (20 epochs) on `type 1` and `type 2` ($1^{st}$ curriculum stage), then added the `type 1`, `type 2` and `mixed` sequences ($2^{nd}$ curriculum stage). Table 11 contains the F1 scores per concept for each model. The trained models were selected using the validation loss during the $2^{nd}$ curriculum

| Model | Concept | F1 score on test set | |
|---|---|---|---|
| | | 1 substitution | 2 substitutions |
| ResNet (Pretrained) | Type 1 | 0.603 | 0.6008 |
| | Type 2 | 0.5774 | 0.5479 |
| | Type 1 Sequences | 0.0517 | 0 |
| | Type 2 Sequences | 0.0941 | 0.0005 |
| | Mixed Sequences | 0.1159 | 0.001 |
| WReN | Type 1 | 0.582 | 0.5246 |
| | Type 2 | 0.594 | 0.5486 |
| | Type 1 Sequences | 0.2974 | 0.0324 |
| | Type 2 Sequences | 0.2031 | 0.0095 |
| | Mixed Sequences | 0.207 | 0.0457 |
| SimpleConvNet | Type 1 | 0.5806 | 0.5612 |
| | Type 2 | 0.586 | 0.581 |
| | Type 1 Sequences | 0.0908 | 0 |
| | Type 2 Sequences | 0.0603 | 0 |
| | Mixed Sequences | 0.1157 | 0 |
| PrediNet | Type 1 | 0.3656 | 0.3324 |
| | Type 2 | 0.4705 | 0.4235 |
| | Type 1 Sequences | 0.2854 | 0.0759 |
| | Type 2 Sequences | 0.357 | 0.0997 |
| | Mixed Sequences | 0.2023 | 0.081 |
| SimpleFeedForward | Type 1 | 0.3916 | 0.3291 |
| | Type 2 | 0.4083 | 0.3989 |
| | Type 1 Sequences | 0.217 | 0.0039 |
| | Type 2 Sequences | 0.0748 | 0 |
| | Mixed Sequences | 0.1757 | 0.0326 |

Table 10: F1 scores per concept for Experiment 4.

stage. For most models, curriculum training resulted in worse F1 scores compared to non-curriculum training.

| Model | Concept | F1 score on test set | |
| --- | --- | --- | --- |
| | | 1 substitution | 2 substitutions |
| ResNet (Pretrained) | Type 1 | 0.6857 | 0.5919 |
| | Type 2 | 0.7071 | 0.5721 |
| | Type 1 Sequences | 0.2522 | 0.0594 |
| | Type 2 Sequences | 0.3024 | 0.0067 |
| | Mixed Sequences | 0.3785 | 0.244 |
| WReN | Type 1 | 0.4751 | 0.4365 |
| | Type 2 | 0.4327 | 0.4038 |
| | Type 1 Sequences | 0.2935 | 0.0892 |
| | Type 2 Sequences | 0.197 | 0.0224 |
| | Mixed Sequences | 0.2415 | 0.1642 |
| SimpleConvNet | Type 1 | 0.5576 | 0.5559 |
| | Type 2 | 0.5861 | 0.5826 |
| | Type 1 Sequences | 0.001 | 0 |
| | Type 2 Sequences | 0.005 | 0 |
| | Mixed Sequences | 0.0541 | 0 |
| PrediNet | Type 1 | 0.2393 | 0.235 |
| | Type 2 | 0.1333 | 0.1233 |
| | Type 1 Sequences | 0.1834 | 0.1179 |
| | Type 2 Sequences | 0.2958 | 0.232 |
| | Mixed Sequences | 0.138 | 0.1453 |
| SimpleFeedForward | Type 1 | 0.292 | 0.2965 |
| | Type 2 | 0.258 | 0.2476 |
| | Type 1 Sequences | 0.049 | 0.0274 |
| | Type 2 Sequences | 0.0155 | 0.0103 |
| | Mixed Sequences | 0.0253 | 0.0158 |

Table 11: F1 scores per concept for the curriculum variant of Experiment 4.