# Geometric Generalization Based Zero-Shot Learning Dataset Infinite World: Simple Yet Powerful

**Rajesh Chidambaram** [1 2]  **Michael Kampffmeyer** [3 2]  **Willie Neiswanger** [4]  **Xiaodan Liang** [4]
**Thomas Lachmann** [1 5]  **Eric Xing** [4 6]

## Abstract

Raven's Progressive Matrices are one of the widely used tests in evaluating the human test taker's fluid intelligence. Analogously, this paper introduces geometric generalization based zero-shot learning tests to measure the rapid learning ability and the internal consistency of deep generative models. Our empirical research analysis on state-of-the-art generative models discern their ability to generalize concepts across classes. In the process, we introduce *Infinite World*[1], an evaluable, scalable, multi-modal, light-weight dataset and Zero-Shot Intelligence Metric ZSI. The proposed tests condenses human-level spatial and numerical reasoning tasks to its simplistic geometric forms. The dataset is scalable to a theoretical limit of infinity, in numerical features of the generated geometric figures, image size and in quantity. We systematically analyze state-of-the-art model's internal consistency, identify their bottlenecks and propose a pro-active optimization method for few-shot and zero-shot learning.

## 1. Introduction

A remarkable case of zero-shot learning by the human brain can be drawn from one of the greatest minds of the 20th century, Albert Einstein. On the basis of his general theory of relativity, Einstein predicted gravitational waves (Einstein & Rosen, 1937), decades before the actual experiments confirmed its existence (Abbott et al., 2016). Like no other

[1]Center for Cognitive Science, TU Kaiserslautern, Germany [2]Visiting Scholar, Carnegie Mellon University, USA [3]UiT Machine Learning Group, UiT The Arctic University of Norway, Tromsø, Norway [4]Machine Learning Department, Carnegie Mellon University, USA [5]Department of Experimental Psychology, University of Leuven, Leuven, Belgium [6]Petuum Inc, USA. Correspondence to: Rajesh Chidambaram <rajesh1990@live.in>.

[1]The dataset will be open-sourced on GitHub

modern physicist, his discoveries fundamentally altered and expanded our understanding of nature (Fölsing, 1997). The ability to perform description based learning and simulation based planning, are vital aspects of human intelligence. Unlike any other species, human brains can rapidly acquire knowledge through natural language conversations and by reading books. The human brain comprehends textual content by constructing mental models of the text (Woolley, 2011). It is hypothesized that the brain optimizes cost functions in a diverse, region-specific and developmental stage dependent manner (Marblestone et al., 2016).

On the other hand, Machine Learning (ML) methods such as Artificial Neural Network (ANN), powered by deep learning techniques (LeCun et al., 2015; Schmidhuber, 2015) have proven to be powerful function approximators. Their applications span across image classification, audio processing, game playing, machine translation, etc., (Sprechmann et al., 2018). Despite their success, their ability to zero-shot learn via simulation (Ha & Schmidhuber, 2018) and perform human level rapid learning from mere textual descriptions (Chaplot et al., 2017) are experimental. Specifically, when the data availability is scarce, simulation based few-shot learning methods becomes indispensable.

In this paper, we introduce human-level zero-shot tests for text-to-image synthesis models and a Zero-Shot Intelligence Metric ZSI. The standardized tests posed by the dataset *Infinite World*, consists of human-level 2-dimensional geometric generalization tasks.

### 1.1. Generative models for text-to-image synthesis

The advent of Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) in generating realistic images, image restoration, video generation, text-to-image synthesis (Radford et al., 2015; Zhang et al., 2016; Xu et al., 2017) have been impressive. Upon training, these models are able to generate images of faces, birds and room interiors from mere textual inputs. Yet, the ability of such models to generate images on unseen texts that involve human-level reasoning and rapid learning are experimental. In this paper, we systematically analyze the performances of state-of-the-art text-to-image synthesis models on the proposed
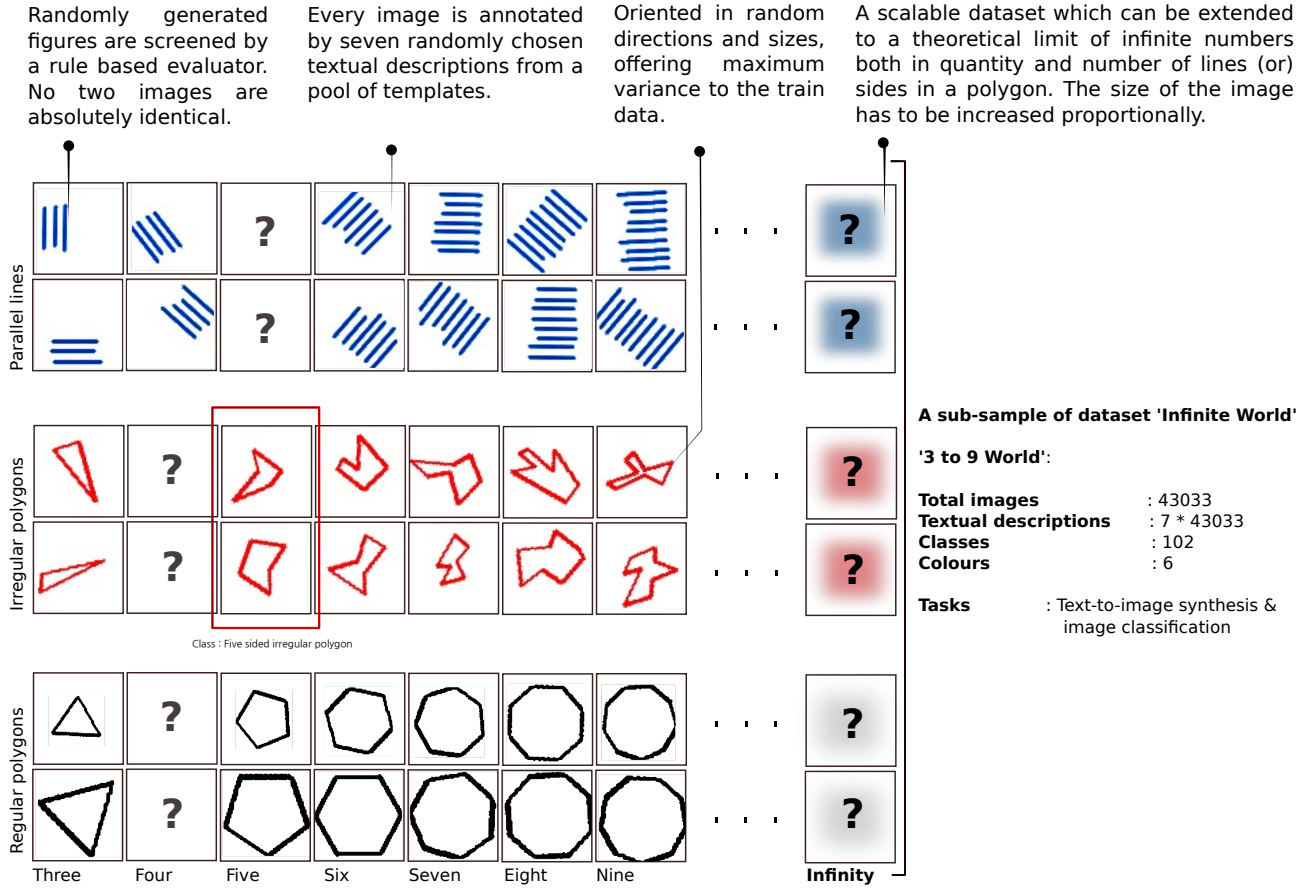
Randomly generated figures are screened by a rule based evaluator. No two images are absolutely identical.

Every image is annotated by seven randomly chosen textual descriptions from a pool of templates.

Oriented in random directions and sizes, offering maximum variance to the train data.

A scalable dataset which can be extended to a theoretical limit of infinite numbers both in quantity and number of lines (or) sides in a polygon. The size of the image has to be increased proportionally.

**A sub-sample of dataset 'Infinite World'**

**'3 to 9 World'**:

| | |
|---|---|
| **Total images** | : 43033 |
| **Textual descriptions** | : 7 * 43033 |
| **Classes** | : 102 |
| **Colours** | : 6 |
| **Tasks** | : Text-to-image synthesis & image classification |

Class : Five sided irregular polygon

**SPATIAL REASONING TASKS:**

Tests ability of the model in generating novel / unseen lines & shapes.

PARALLEL LINES:

TEXT: 'five lines' ?

FULLY CONNECTED SHAPES

TEXT: 'four sided regular polygon' ?

**NUMERICAL REASONING TASKS:**

Tests ability of the model in generating shapes of any sides, given one sample image for the number of lines representation.

TEXT: 'twenty-five fully connected lines' ?

*Figure 1.* Zero-shot learning tests on 2-dimensional geometric generalizations tasks

2d-geometric generalization tasks. Through the Zero-Shot Intelligence metric ZSI ($\psi$) a precise analysis of such text-to-image synthesis models reveal their internal consistency and caveats in image generation.

### 1.2. Zero-Shot Learning

The demand for vast supervised data has made deep neural networks incompatible with several cardinal tasks. Inspired by the rapidly learning nature of human-beings, zero-shot learning aims to imitate such behaviors in machines. Few-shot learning and zero-shot learning are extreme forms of transfer learning (Goodfellow et al., 2016). Zero-shot learning is also called as zero-data learning. Popular few-shot adaptation techniques include bayesian modeling (Tenen-

baum, 1999), fine tuning a pre-trained neural net (Yu et al.), memory augmented networks (Santoro et al., 2016) and meta-learners (Schmidhuber, 1987; Ravi & Larochelle, 2016). In this paper, we shall see the necessity of more advanced optimization techniques for generative models, to perform zero-shot learning on human-level reasoning tasks.

### 1.3. Geometric generalization

A three-sided triangle, a four-sided quadrilateral, and so on to $n$ sides can be generalized to the geometric concept of *polygon*. It is sufficient to teach a human the concept of a polygon for few numbers and yet the person can generalize the concept to unseen numbers. Such a generalization is possible for any $n$ dimensional space ($n > 1$). Whereas,

state-of-the-art machine learning methods demands examples from every new class. In this work, we approach 2d-geometric generalization tasks with ANNs. Utilizing ANNs for such tasks leverages the ability of machine learning algorithms to perform learning with comparatively less human input. Upon developing a domain specific model which can perform zero-shot learning on complicated tasks, the model can be scaled to real world applications by expanding its verbal and visual corpus.

## 2. Dataset: *Infinite World*

In general, existing benchmarks such as the visual genome (Krishna et al., 2017) and the CLEVR dataset (Johnson et al., 2017) proposes to test the reasoning abilities of the machine learning algorithms. Models such as the relational network (Santoro et al., 2017), which successfully solves visual question answering (VQA) problems on the CLEVR datasets, suffers severe limitations (Junkyung Kim, 2018). Alternatively, it is notable that images from physics and geometry were used to demonstrate *why a diagram is sometimes worth ten thousand words* (Larkin & Simon, 1987). On a similar note, we propose geometric generalization based zero-shot learning tests for methods in ML. The proposed tests condenses human-level reasoning tasks to its simplistic geometric forms. This simplicity aids in precisely evaluating the performance of the generative model. The dataset is scalable to a theoretical limit of infinity, in numerical features of the generated geometric figures, image size and in quantity. It is notable that the newly introduced dataset *Infinite World* and Zero-Shot Intelligence metric are not restricted to text-to-image synthesis algorithms, but also for other machine learning tasks such as classification.

*Infinite World* proposes tasks that are comparable to fluid intelligence tests such as Cattell Culture Fair IQ test (Cattell, 1963) and Progressive Matrices (Duncan et al., 1995). While deducing the intelligence of an agent to a single metric is difficult, we formulate a task specific zero-shot intelligence metric $\psi$ to compare the zero-shot performances of different models. Metric $\psi$ is always accompanied by the task in consideration. This offers ANN models a common platform to test their zero-shot performances on tasks including human level spatial and numerical reasoning. The generated images are verified for its accuracy by a rule-based evaluator. For every rejected generation, a new figure is generated for the required textual description of the image. The same rule-based evaluator is later used to evaluate the performance of the text-to-image synthesis model.

### 2.1. *3-9 World*: A subsample of *Infinite World*

The *3-9 World* dataset, a subset of the *Infinite World* dataset generator, considers that only numbers from $3$ $to$ $9$ exists. The sub-sampled dataset consists of $41,000+$ train images

of lines and polygons. The images are of size $64 \times 64$ pixels. Each image consists of seven randomly-chosen textual descriptions from a pool of templates. The dataset consists of 567 distinct train texts and 91 distinct zero-shot test texts. The pool of textual annotations are customizable and scalable. The model is required to generalize the concept of connectedness of shapes and non-connectedness of lines; the zero-shot task is to predict a particular *unseen* shape from the generalized concept. From our preliminary experiments, we observed that the generative models were comparatively better in generating images of same color as in the text. But they developed significant gap in performance while generating disconnected lines and connected shapes. Hence, the dataset was further developed to mainly focus on parallel lines, irregular polygons and regular polygons.

## 3. Zero-Shot Intelligence metric ZSI ($\psi$)

Popular text-to-image synthesis models are tasked for zero-shot generation of realistic images of bedrooms, flowers and human faces. In such tasks, it is difficult to precisely evaluate the ability to reason, zero-shot learn and infer consistency. Existing evaluation metrics such as the inception score (Salimans et al., 2016) and R precision (Xu et al., 2017) do not capture the generative abilities of the model (Barratt & Sharma, 2018). Unlike real world images, since the *Infinite World* condenses complex reasoning tasks to simplistic geometric forms, it is possible to precisely evaluate the model's performance on a zero-shot metric. Subjective nature of evaluation is comparatively minimal in such tasks. Hence, we introduce a new evaluation method for Zero-Shot Intelligence ZSI ($\psi$) to measure the reliability of black-box function approximators.
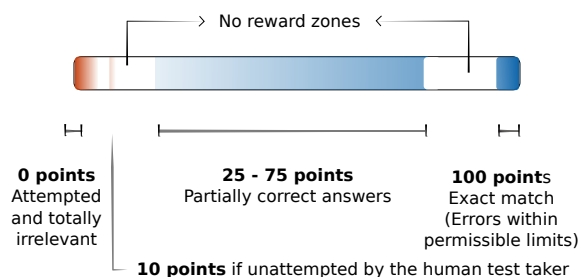


*Figure 2.* ZSI ($\psi$) - scoring method

### 3.1. Computing the $\psi$ score

Since deep neural networks continue to remain mostly as black-box function approximators (Chakraborty et al., 2017), two models that produce the same results with partially correct zero-shot results, cannot be awarded the same score. For the 3-9 World dataset, let us consider that two algorithms on zero-shot testing generate four lines in the place
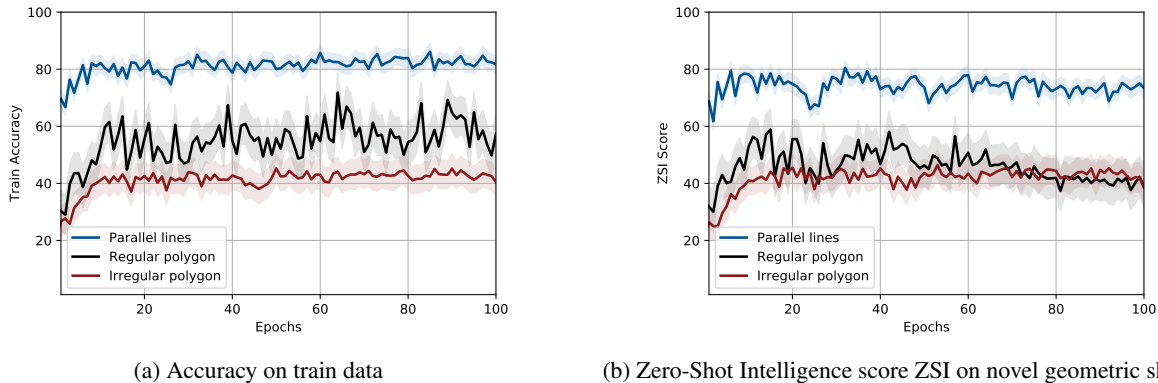
(a) Accuracy on train data

(b) Zero-Shot Intelligence score ZSI on novel geometric shapes

*Figure 3.* Performance of Generative Adversarial text-to-image synthesis on the 3-9 world dataset. The results were averaged over several initializations of internal parameters.

of five lines. Though the answers are partially correct to the same extent, while one algorithm might have generated four lines using regression techniques, the other might have generated four lines through recursive attention methods. While it cannot be predicted which algorithm has reasoned on the task, we awarded a reduced score to both algorithms. Hence, partial scores were awarded in the range $25$ $to$ $75$, while full scores were awarded only when the results are completely correct. Such partial scores were derived in proportion to the closeness of the generated image to the description. For example, for a text description of *'five lines'*, if four lines were generated, a partial score proportional to 4/5 was awarded in the range of $25$ $to$ $75$. Complete correctness of the results were relaxed with permissible error limits, as derived from the experimental results from the top 5 performances of human test takers. $\psi$ score was calculated independently for parallel lines, irregular polygons and regular polygons.

### 3.2. Rule-based evaluator

A rule-based evaluator was used to perform automatized evaluation of every image for its closeness to the given textual description. The evaluator evaluates images during dataset generation, as well as to deduce the $\psi$ score of the given text-to-image synthesis model. Douglas-Peucker algorithm (Douglas & Peucker, 1973) was used to detect the contours of the generated figure. To accommodate minor variations, upon Gaussian blurring, canny edge detectors (Canny, 1987) were used for detecting lines. To identify the color of the figure, the k-means clustering (Hartigan & Wong, 1979) for pixel level RGB values were computed. Since Gaussian blurring is used, it is sufficient for the text-to-image synthesizer to generate a figure of the specified color either in the first dominant or second dominant region of the color cluster.

## 4. Experiments and results

### 4.1. Conditional generative networks: Generative adversarial text-to-image synthesis

Generative Adversarial Text-to-Image Synthesis (Reed et al., 2016) is one of the widely used text augmented conditional generative network for text-to-image synthesis. The performance of this vanilla model on the *3-9 World* dataset precisely illustrates how the model fails to rapidly combine its learned representations for zero-shot generations. Interestingly, the model performs comparatively better in generating disconnected-parallel lines than while connecting such lines to form a meaningful polygon. More specifically, the model struggles in generating a irregular polygon than a regular polygon. This evinces that the convolutional and deconvolutional neural networks of the generative model develop specific filters for regular shapes but fail to map randomly varying irregular polygons. From time to time, the model generates completely connected polygons, illustrating its ability to generate connected shapes. From figure 3, the $\psi$ scores across epochs illustrates that the zero-shot performance is comparatively consistent but poor. Hence, a proactive-rapid optimizer can potentially boost the zero-shot performance of the model.

#### 4.1.1. COSINE DISTANCE BETWEEN THE SENTENCE EMBEDDINGS OF SKIP-THOUGHT VECTORS FOR '3-9 WORLD' DATASET

The similarity between various sentence embeddings were computed. The cosine distance between the skip-thought vectors (Kiros et al., 2015) illustrates that all the skip thought vectors used to generate images were highly distinct and dissimilar. Upon training, the model established strong correlation between similar textual descriptions. For example, the text 'three green colored lines' and 'the image contains three lines that are green in color' generated similar images,
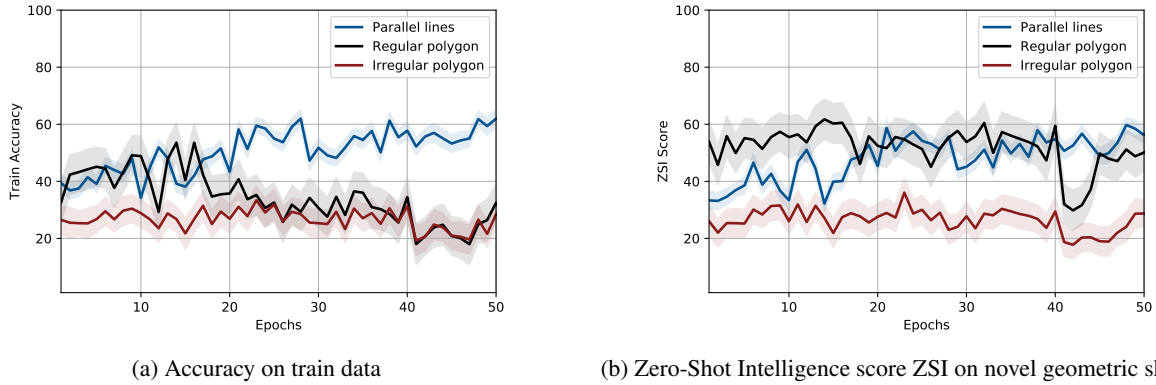
(a) Accuracy on train data

(b) Zero-Shot Intelligence score ZSI on novel geometric shapes

*Figure 4.* Zero-Shot Intelligence performance of AttnGAN on the *3 to 9 world* dataset. The results were averaged over several initializations of internal parameters.
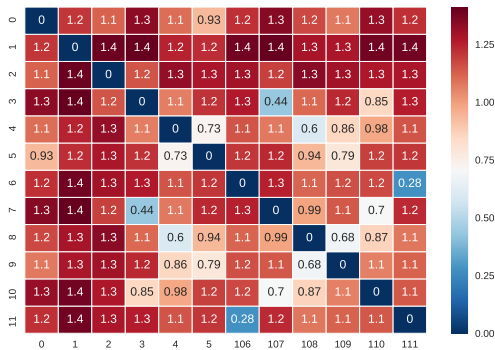


*Figure 5.* Cosine distance between the sentence embeddings of skip-thought vectors for '3-9 World' dataset. Columns 0 *to* 5 represent sentences of the same image while columns 106 *to* 111 represent sentences of the same image from a different class.



*Figure 6.* AttnGAN generated images when trained on CUB-2011 dataset. The images illustrates that the ability of the model to spatially and numerically reason are experimental. Note that the displayed textual inputs are not from completely unseen classes.

as computed by the similarity between their $\psi$ scores.

## 4.2. Recurrently attentive generative networks: AttnGAN

We chose another state-of-the-art text-to-image generation model, Attentional Generative Adversarial Networks (AttnGAN) (Xu et al., 2017). In addition to the multiple generative networks stacked upon one another, its recurrent attention on word embeddings is an interesting feature of this model. As reported in their actual work, upon training on the CUB (Wah et al., 2011) and COCO datasets (Lin et al., 2014), the model out performed all state-of-the-art text-to-image generation methods. The best inception score on corresponding datasets being $4.36 \pm .03$ and $25.89 \pm .47$ respectively. Figure 6 illustrates the zero-shot generations by AttnGAN. Though the images are comparatively more realistic, the generated images do not match the input text by

all means. Hence, an empirical research on the fundamental abilities of the model to spatially and numerically reason are essential.

We trained AttnGAN on the *3-9 World* dataset over resized $256 \times 256$ pixel images. The model was trained on various DAMSM encoders over the $200^{th}$ and $450^{th}$ epochs. Yet, the $\psi$ scores remained same. Notably, despite its success on fine grained image generation, the model was not able to demonstrate an improved $\psi$ score. From the $\psi$ scores in figure 4, it is evident that the recurrent attention method has not improved the generalizing and few-shot learning abilities of the model. Interestingly, the vanilla Generative Adversarial Text-to-Image Synthesis model outperforms AttnGAN in the $\psi$ scores. Such low $\psi$ scores and their tendency to fluctuate across epochs indicates the model's poor generalization ability and inconsistency respectively. Hence, the multiple recurrently attentive method and introduction of noise into

AttnGAN are clearly subjected to further examination.

## 4.3. Analogous experiments on human beings

We observe that the task of zero-shot text-to-image synthesis and the evaluation metric $\psi$ are not exactly applicable to the human test takers. Yet, we conducted analogous experiments on humans (n = 21, mean age = 27.3, s.d = 11.3), only to have an estimate of human-level performance. All test-takers were literate and at-least had a high school degree. The test takers were then asked to draw images of regular and irregular polygons of randomly chosen sides in the range 3 to 30. During the questionnaire, the human test-takers noted on average that they have rarely come across specific polygons that were more than 13 sided. Hence, sketches of polygons that were more than 13 sided can be considered analogous to zero-shot generation. The disadvantage of this approach being, no direct comparison could be drawn to the generative models' zero-shot performances on *3 to 9 world* dataset. The scores were computed manually, partially comparable to the $\psi$ scoring metric.
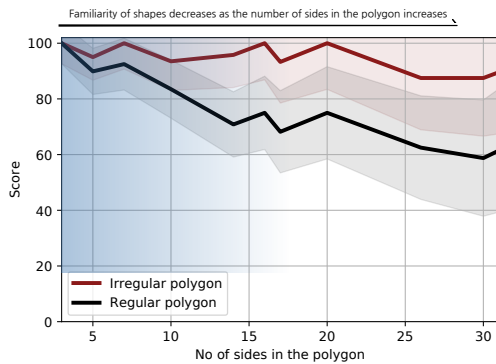


*Figure 7.* Analogous experiments on human beings. The peaks in-between are from even numbered polygons. The drop at around number five was owing to few participants skipping easier ones.

Multiple attempts by a human test-taker were ignored, until the person confirmed his final answer. Inspiration for such an evaluation technique is derived from the multiple attempts that were made during discovery of gravitational waves (Einstein & Rosen, 1937). Despite several editions and updates to his initial predictions, Einstein's ultimate prediction preceded the experimental confirmation. Hence, such multiple attempts can be neglected and the learning can still be considered as zero-shot learning.

The analogous experiment results on humans show that the fundamental way in which humans and machine learning algorithms generate geometric figures differ greatly. While machine learning algorithms perform comparatively better on regular polygons via rote memorization of train images, human beings consistently generate irregular polygons with ease (until they were tired of sketching figures that were

more than 30 sided).

## 5. Discussion

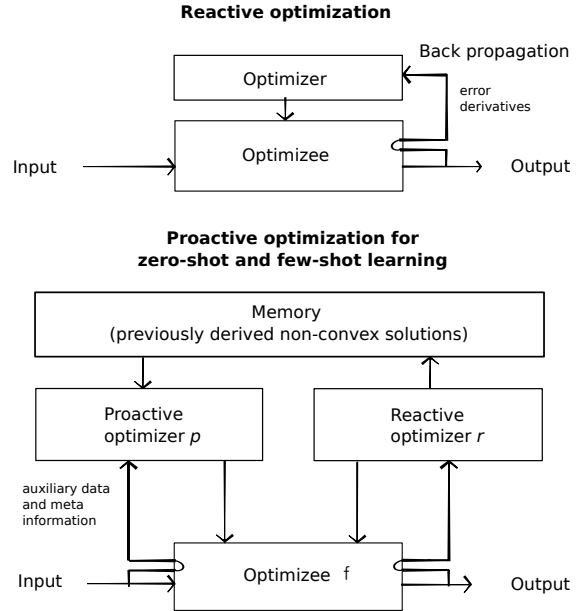### 5.1. Proactive Optimization for Zero-Shot Learning



*Figure 8.* Reactive and proactive optimization

---

**Algorithm 1** Proactive optimization

**Input:** Data $x_n$, explicit or implicit auxiliary information $a_n(x_i, y_i)$, meta-level information $m_i(f, x_i)$ available within the system $f$, reactively optimized parameters $W$, pro-actively optimized parameters $W^*$, loss function $L(y_i, f(x_i, W, W^*))$, proactive optimizer $p$, reactive optimizer $r$, zero-shot input $x_z$, zero-shot ground truth $y_z$:

**repeat**
    Initialize $W$, $W^*$ for $f : X \rightarrow Y$
    **for** $i = 1$ **to** $n$ **do**
        **if** $L(.)_{train\ set} > L(.)_{validation\ set}$ **then**
            $W^*_{t+1} = W^*_t + p(W^*_t, a_i(x_i, y_i), m_i(f, x_i))$
            $W_{t+1} = W_t + r(\nabla_\theta L(W_t))$
        **end if**
    **end for**
**until** $f(x_z, W, W^*) = y_z$

---

Non-convex optimization methods in deep neural networks have increasingly focused on skipping saddle points (Dauphin et al., 2014; Ge et al., 2015; Anandkumar & Ge, 2016; Bottou et al., 2016) and local minima while minimizing a specific loss function $L$. As sparse data necessitates (Duchi et al., 2013), demand for few-shot learning

and generalization of concepts across tasks are essential. Yet, the subject has been rarely focused. Here, we propose pro-active non-convex optimization methods for generalization and zero-shot learning. Given a training set $\{(x_n, y_n), n = 1...N\}$, optimizing the internal parameters $W$ of a system $f$ upon reception of the back-propagated gradients through supervised data $y_n$ is a reactive way of learning $f : X \rightarrow Y$. Proactive methods of optimization includes optimization of internal parameters $W^*$ through explicit or implicit auxiliary data $a_z(x_z, y_z)$ about the unseen class $x_z$, by utilizing the the meta-level information $m_z(f, x_z)$ computed by the system $f$ and uniquely combining previously derived non-convex solutions, that were stored through external memory augmentation (Sprechmann et al., 2018). Recent approaches (Romera-Paredes & Torr, 2015; Munkhdalai & Yu, 2017) to few-shot learning have focused on a fully or partially proactive methods for faster adaptation. The proactive approach to optimization aids in utilizing the auxiliary cues and meta-level information of the task, thereby aiding rapid learning. Both reactive and pro-active optimizers can exist in the same system. Such parallel streams of learning are analogous to the complementary learning systems in the human brain (OReilly et al., 2014).

## 6. Conclusion

In this paper, we proposed a new range of zero-shot learning tasks to evaluate the generative model's internal consistency and generalization abilities. Our empirical research work on state-of-the art text-to-image generation model exhibits a huge gap between human-level rapid learning and the few-shot methods in machine learning. Despite its excellence in fine grained image generation by the recurrently attentive generative adversarial network, the model's internal consistency is lower than the vanilla conditional generative adversarial networks. Hence, the performance of a model on the proposed 2d-geometric generalization tasks, can be used as an effective method to compute the reliability of a model. Here, the reliability measure indicates the ability of the model to consistently perform human-level spatial and numerical reasoning tasks across classes. The dataset *Infinite World* will aid in the development of a new range of optimization methods that can rapidly learn through generalization.

## Acknowledgments

## References

Abbott, Benjamin P, Abbott, Richard, Abbott, TD, Abernathy, MR, Acernese, Fausto, Ackley, Kendall, Adams, Carl, Adams, Thomas, Addesso, Paolo, Adhikari, RX, et al. Observation of gravitational waves from a binary black hole merger. *Physical review letters*, 116 (6):061102, 2016.

Anandkumar, Animashree and Ge, Rong. Efficient approaches for escaping higher order saddle points in nonconvex optimization. In *Conference on Learning Theory*, pp. 81–102, 2016.

Barratt, Shane and Sharma, Rishi. A note on the inception score. *arXiv preprint arXiv:1801.01973*, 2018.

Bottou, Léon, Curtis, Frank E, and Nocedal, Jorge. Optimization methods for large-scale machine learning. *arXiv preprint arXiv:1606.04838*, 2016.

Canny, John. A computational approach to edge detection. In *Readings in Computer Vision*, pp. 184–203. Elsevier, 1987.

Cattell, Raymond B. Theory of fluid and crystallized intelligence: A critical experiment. *Journal of educational psychology*, 54(1):1, 1963.

Chakraborty, Supriyo, Tomsett, Richard, Raghavendra, Ramya, Harborne, Daniel, Alzantot, Moustafa, Cerutti, Federico, Srivastava, Mani, Preece, Alun, Julier, Simon, Rao, Raghuveer M, et al. Interpretability of deep learning models: a survey of results. DAIS, 2017.

Chaplot, Devendra Singh, Sathyendra, Kanthashree Mysore, Pasumarthi, Rama Kumar, Rajagopal, Dheeraj, and Salakhutdinov, Ruslan. Gated-attention architectures for task-oriented language grounding. *arXiv preprint arXiv:1706.07230*, 2017.

Dauphin, Yann N, Pascanu, Razvan, Gulcehre, Caglar, Cho, Kyunghyun, Ganguli, Surya, and Bengio, Yoshua. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in neural information processing systems*, pp. 2933–2941, 2014.

Douglas, David H and Peucker, Thomas K. Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *Cartographica: The International Journal for Geographic Information and Geovisualization*, 10(2):112–122, 1973.

Duchi, John, Jordan, Michael I, and McMahan, Brendan. Estimation, optimization, and parallelism when data is sparse. In *Advances in Neural Information Processing Systems*, pp. 2832–2840, 2013.

Duncan, John, Burgess, Paul, and Emslie, Hazel. Fluid intelligence after frontal lobe lesions. *Neuropsychologia*, 33(3):261–268, 1995.

Einstein, Albert and Rosen, Nathan. On gravitational waves. *Journal of the Franklin Institute*, 223(1):43–54, 1937.

Fölsing, Albrecht. *Albert Einstein: a biography*. Viking, 1997.

Ge, Rong, Huang, Furong, Jin, Chi, and Yuan, Yang. Escaping from saddle pointsonline stochastic gradient for tensor decomposition. In *Conference on Learning Theory*, pp. 797–842, 2015.

Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.

Goodfellow, Ian, Bengio, Yoshua, Courville, Aaron, and Bengio, Yoshua. *Deep learning*, volume 1. MIT press Cambridge, 2016.

Ha, David and Schmidhuber, Jürgen. World models. *arXiv preprint arXiv:1803.10122*, 2018.

Hartigan, John A and Wong, Manchek A. Algorithm as 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1): 100–108, 1979.

Johnson, Justin, Hariharan, Bharath, van der Maaten, Laurens, Fei-Fei, Li, Zitnick, C Lawrence, and Girshick, Ross. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pp. 1988–1997. IEEE, 2017.

Junkyung Kim, Matthew Ricci, Thomas Serre. Not-so-CLEVR: Visual relations strain feedforward neural networks, 2018. URL https://openreview.net/forum?id=HymuJz-A-.

Kahou, Samira Ebrahimi, Atkinson, Adam, Michalski, Vincent, Kadar, Akos, Trischler, Adam, and Bengio, Yoshua. Figureqa: An annotated figure dataset for visual reasoning. *arXiv preprint arXiv:1710.07300*, 2017.

Kiros, Ryan, Zhu, Yukun, Salakhutdinov, Ruslan R, Zemel, Richard, Urtasun, Raquel, Torralba, Antonio, and Fidler, Sanja. Skip-thought vectors. In *Advances in neural information processing systems*, pp. 3294–3302, 2015.

Krishna, Ranjay, Zhu, Yuke, Groth, Oliver, Johnson, Justin, Hata, Kenji, Kravitz, Joshua, Chen, Stephanie, Kalantidis, Yannis, Li, Li-Jia, Shamma, David A, et al. Visual genome: Connecting language and vision using crowd-sourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017.

Larkin, Jill H and Simon, Herbert A. Why a diagram is (sometimes) worth ten thousand words. *Cognitive science*, 11(1):65–100, 1987.

LeCun, Yann, Bengio, Yoshua, and Hinton, Geoffrey. Deep learning. *nature*, 521(7553):436, 2015.

Lin, Tsung-Yi, Maire, Michael, Belongie, Serge, Hays, James, Perona, Pietro, Ramanan, Deva, Dollár, Piotr, and Zitnick, C Lawrence. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.

Marblestone, Adam H, Wayne, Greg, and Kording, Konrad P. Toward an integration of deep learning and neuroscience. *Frontiers in Computational Neuroscience*, 10: 94, 2016.

Munkhdalai, Tsendsuren and Yu, Hong. Meta networks. *arXiv preprint arXiv:1703.00837*, 2017.

OReilly, Randall C, Bhattacharyya, Rajan, Howard, Michael D, and Ketz, Nicholas. Complementary learning systems. *Cognitive Science*, 38(6):1229–1248, 2014.

Radford, Alec, Metz, Luke, and Chintala, Soumith. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.

Ravi, Sachin and Larochelle, Hugo. Optimization as a model for few-shot learning. 2016.

Reed, Scott, Akata, Zeynep, Yan, Xinchen, Logeswaran, Lajanugen, Schiele, Bernt, and Lee, Honglak. Generative adversarial text to image synthesis. *arXiv preprint arXiv:1605.05396*, 2016.

Romera-Paredes, Bernardino and Torr, Philip. An embarrassingly simple approach to zero-shot learning. In *International Conference on Machine Learning*, pp. 2152–2161, 2015.

Salimans, Tim, Goodfellow, Ian, Zaremba, Wojciech, Cheung, Vicki, Radford, Alec, and Chen, Xi. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pp. 2234–2242, 2016.

Santoro, Adam, Bartunov, Sergey, Botvinick, Matthew, Wierstra, Daan, and Lillicrap, Timothy. One-shot learning with memory-augmented neural networks. *arXiv preprint arXiv:1605.06065*, 2016.

Santoro, Adam, Raposo, David, Barrett, David G, Malinowski, Mateusz, Pascanu, Razvan, Battaglia, Peter, and Lillicrap, Tim. A simple neural network module for relational reasoning. In *Advances in neural information processing systems*, pp. 4974–4983, 2017.

Schmidhuber, Jürgen. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117, 2015.

Schmidhuber, Jrgen. *Evolutionary principles in self-referential learning, or on learning how to learn: The meta-meta-... hook*. Diplomarbeit, Technische Universitt Mnchen, Mnchen, 1987.

Sprechmann, Pablo, Jayakumar, Siddhant M, Rae, Jack W, Pritzel, Alexander, Badia, Adrià Puigdomènech, Uria, Benigno, Vinyals, Oriol, Hassabis, Demis, Pascanu, Razvan, and Blundell, Charles. Memory-based parameter adaptation. *arXiv preprint arXiv:1802.10542*, 2018.

Tenenbaum, Joshua Brett. *A Bayesian framework for concept learning*. PhD thesis, Massachusetts Institute of Technology, 1999.

Wah, Catherine, Branson, Steve, Welinder, Peter, Perona, Pietro, and Belongie, Serge. The caltech-ucsd birds-200-2011 dataset. 2011.

Woolley, Gary. Reading comprehension. In *Reading Comprehension*, pp. 15–34. Springer, 2011.

Xu, Tao, Zhang, Pengchuan, Huang, Qiuyuan, Zhang, Han, Gan, Zhe, Huang, Xiaolei, and He, Xiaodong. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. *arXiv preprint arXiv:1711.10485*, 2017.

Yu, Dong, Deng, Li, and Dahl, George. Roles of pre-training and fine-tuning in context-dependent dbn-hmms for real-world speech recognition.

Zhang, Han, Xu, Tao, Li, Hongsheng, Zhang, Shaoting, Huang, Xiaolei, Wang, Xiaogang, and Metaxas, Dimitris. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. *arXiv preprint arXiv:1612.03242*, 2016.

## A. Appendix

### A.1. Exception handling

In the dataset, in addition to the illustrated ways of connecting lines to form a polygon, there can be many other ways in partially or completely forming a polygon. For example, a star polygon is also a regular polygon, but the *Infinite World* dataset does not depict one. Hence, the rule based evaluator

has been assigned to notify whenever it has potentially encountered such an exception. Though no such exceptions were generated in the succeeding experiments, any detected exceptions were assigned to be saved separately for manual evaluation. The selection of such exceptions were based on the unique geometric features of the exception. For example, an interestingly-connected open-ended figure would have more free edges than a irregular or regular polygon.
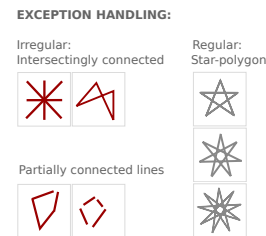


*Figure 9.* Handling of exceptional cases

### A.2. One-shot classification performance of Meta Networks, a partially pro-active optimizer

Since the task of image classification is different from image generation, we performed few-shot classification on state-of-the-art one-shot classifier Meta Networks (Munkhdalai & Yu, 2017). Meta Networks uses a partially proactive and as well as a reactive optimizer for few-shot image classification. Since the proactive optimizer is dependent on the error derivatives from a support-set, the proactive optimizer is not completely independent from the back propagation technique. The model performed one-shot learning at above 70% accuracy (accuracy as defined in the original paper). Note that, only one-shot performance was measured and the ability to numerically reason beyond number 9 was not computed in this experiment.
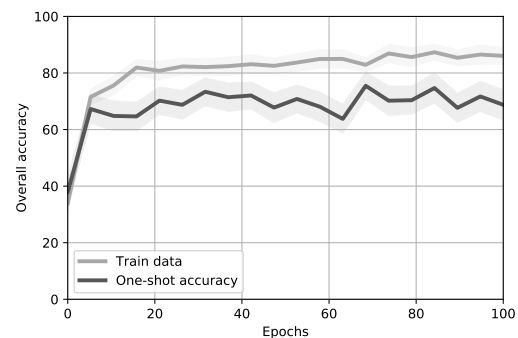


*Figure 10.* One-shot classification performance of Meta Networks, a partially proactive optimization based image classifier, on 3-9 World
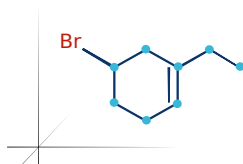
The proposed spatial and numerical reasoning based zero-shot learning tasks are elemental for applications in:



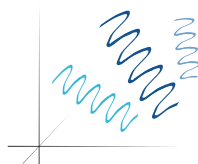Figure 11. Few future applications of description based zero-shot learning. In addition to text-to-image synthesis and classification tasks, the dataset can also be extended to FigureQA (Kahou et al., 2017) type advanced VQA tasks

### A.3. Applications and future work

To perform geometric generalization, it requires both numerical and spatial reasoning abilities. Hence, future applications of such zero-shot learning are enormous. The newly introduced tasks are elemental for applications in simulation based planning for autonomous vehicles, automated simulation of protein conformation, improving reading comprehension through image generation, intrinsic reward prediction for reinforcement learning through simulation of future states, etc., To enable such applications, the dataset can be further extended for more than two spatial dimensions and for optical character association to the nearest geometric features. Further scope of development lies in expanding the verbal and visual corpus to enable such models to acquire knowledge directly from text books.

GENERATED IMAGES ON TRAIN DATA TEXTS AND UNSEEN TEXTS BY 'GENERATIVE ADVERSARIAL TEXT-TO-IMAGE SYNTHESIS' ON '3 TO 9 WORLD' DATASET:



*Figure 12.* Generated images on train texts and unseen texts from Generative Adversarial Text-to-Image Synthesis on 3-9 World dataset. The displayed images were chosen in random.

GENERATED IMAGES ON TRAIN DATA TEXTS AND UNSEEN TEXTS BY 'AttnGAN' ON '3 TO 9 WORLD' DATASET:
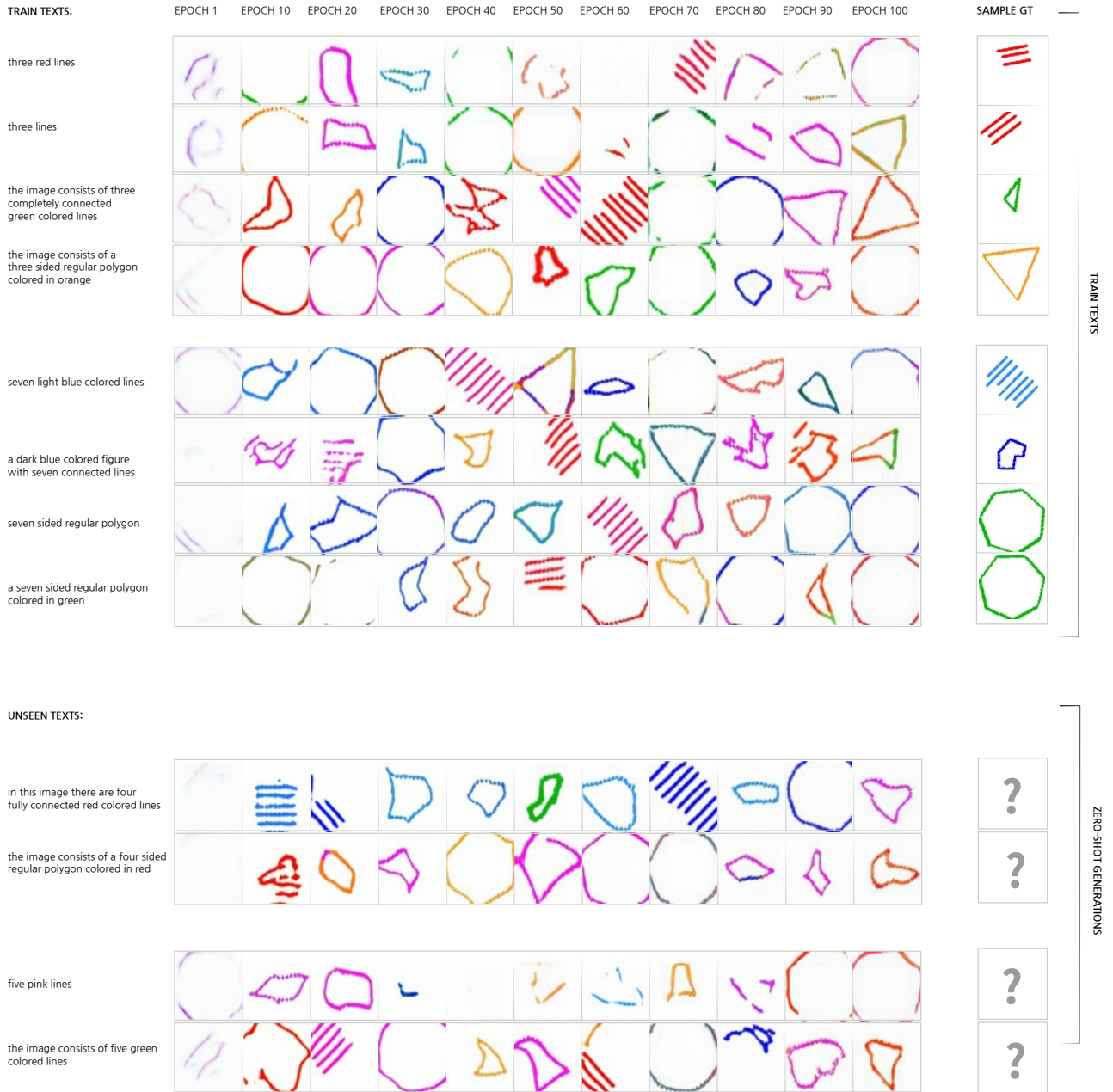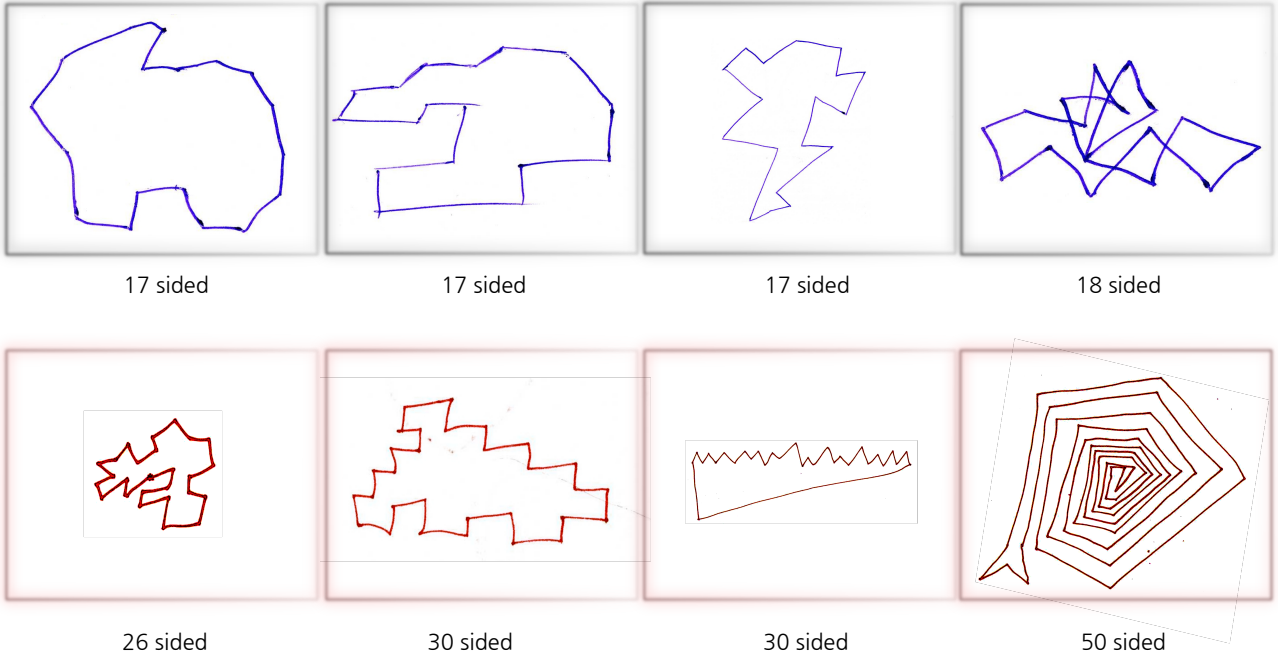


*Figure 13.* Generated images on train texts and unseen texts from AttnGAN on 3-9 World dataset. The displayed images were chosen in random.

IRREGULAR POLYGON



| 17 sided | 17 sided | 17 sided | 18 sided |



| 26 sided | 30 sided | 30 sided | 50 sided |

REGULAR POLYGON

Multiple attempts were considered as single attempt since the
participant realized his / her mistake and redrew the diagram
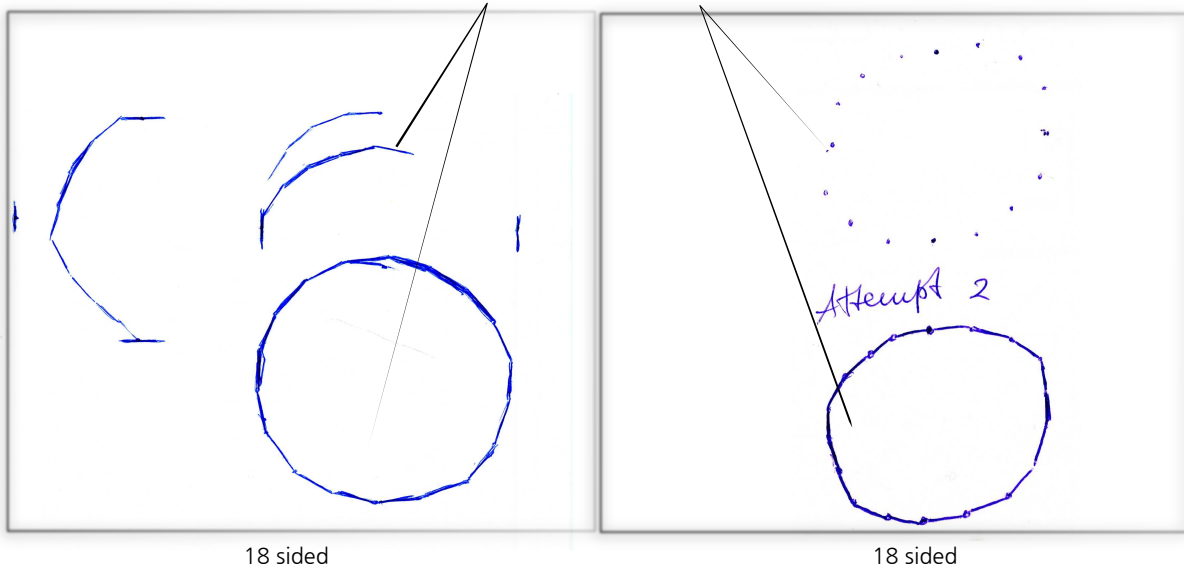


| 18 sided | 18 sided |

*Figure 14.* Analogous 2d-geometric generalization tests for human beings. All displayed results were awarded full scores.