

Tell, Draw, and Repeat: Generating and Modifying Images Based on Continual Linguistic Instruction

Appendix

A. Object detector and localizer network

All of the evaluation metrics for the Generative Neural Visual Artist (GeNeVA) task rely on the object detector and localizer network and hence, it needs to have high detection and localization performance. We report the performance of the trained object detector and localizer network on the test set images of both Collaborative Drawing (CoDraw) and Iterative CLEVR (i-CLEVR) datasets in Table 1.

Dataset	Precision \uparrow	Recall \uparrow	F1-Score \uparrow	NRMSE \downarrow
CoDraw	0.962	0.972	0.964	0.121
i-CLEVR	1.000	1.000	1.000	0.060

Table 1. Mean test set Precision, Recall, and F1-Score for the object detector and localizer network. Normalized Root Mean Squared Error (NRMSE) is the root mean square distance between the localizer’s predicted and ground truth object centroids normalized by the image dimensions. \uparrow : higher is better, \downarrow : lower is better.

B. Relational Similarity metric: rsim

B.1. Additional details

For both CoDraw and i-CLEVR datasets, we determine front-behind and left-right relationships by comparing the coordinates of their centre predicted by the object detector and localizer network. We run the network on both ground truth and generated images to predict the centre coordinates (rather than using perspective coordinates provided by the renderer as these are only available for the ground truth for i-CLEVR).

B.2. Appropriateness for evaluation

The CoDraw and i-CLEVR datasets are constructed such that there is at most one object of each object class per image. Hence, we train the object detector to predict only binary presence of each object class and the localizer regresses only one set of centroid coordinates per class. This design breaks if multiple instances of an object class are generated or if the object detector frequently misclassifies objects. However, qualitatively assessing the generated im-

ages, over-generation is rare and the object detector accuracy is very high (cf. Table 1).

Since all objects in ground truth scenes occur at most once, generations with multiple instances per class are out-of-distribution. The model cannot learn to exploit this flaw, since rsim is not optimized during training. Thus, over-generation is not a failure mode we have observed. Additionally, rsim is position-sensitive: over-generation would not necessarily produce the correct relative positions of objects since the object localizer only localizes one instance per class. For datasets with multiple instances per class, the rsim metric should be modified such that the denominator is the union of ground-truth and predicted detections, which will penalize over-generation.

B.3. Shortcomings

Quantitative measures for attributes like “boy kicking” are currently a missing piece. We share this shortcoming with all text-to-image Generative Adversarial Network (GAN)-based methods and most of the conditional GAN literature. At the moment, conditional GANs are evaluated using Inception Score (IS) and Fréchet Inception Distance (FID), both of which do not account for attributes. An evaluation metric that accounts for attributes will be a valuable contribution for future research.

B.4. Comparison with existing metrics

The Scene Similarity Metric (SSM) used by Kim et al. [1] is well-suited for the setting of predicting object location and attributes. SSM is a weighted score across recall and considers objects that face the wrong direction, incorrect expressions, poses, clip art size, distance between object positions in ground truth and predicted image, and left-right and front-behind relationships. SSM achieves the highest score for exact reconstructions. In our case, we want to not just reward reconstructions but also plausible generations where left-right, front-behind relationships are correct. Our main focus here is to generate complete images instead of predicting object location and attributes. Several attributes, such as boy or girl poses / expressions, or object directions have lower detector accuracy and consequently would reduce metric reliability (cf. Section B.2).

B.5. Qualitative evaluation

We provide generated image examples with scores spread out between the minimum value (0) and maximum value (1) on the rsim metric in Figure 1. This is to provide readers with a more intuitive understanding of how the metric captures which spatial relationships match between the ground truth and the generated image.

C. Generation Examples

We present selected examples generated using our best model (D Subtract) on two datasets. Examples generated for CoDraw are presented in Figure 2 and examples generated for i-CLEVR are presented in Figure 3. We also present random examples from all the models present in the ablation study for a qualitative comparison on the CoDraw dataset. These are shown in Figure 4 (Baseline), Figure 5 (Mismatch), Figure 6 (G prior), Figure 7 (Aux), Figure 8 (D Concat), Figure 9 (D Subtract), and Figure 10 (Non-iterative).

D. Generalization to new background images

GeNeVA-GAN was trained using the empty background image as the initial image. We ran an experiment where we used a different image (intermediate ground truth image from the test set containing objects) as the initial image. We present generated examples from this experiment in Figure 11. The model is able to place the desired object at the correct location with the correct color and shape over the provided image. This shows that the model is capable of generalizing to a background it was not trained on and it can understand the existing objects from just the initial image without any instruction history for placing them.

E. i-CLEVR Dataset Generation

To generate the image for each step in the sequence, an object with random attributes is rendered to the scene using Blender [2]. We ensure that all objects have a unique combination of attributes. Each object can have one of 3 shapes (cube, sphere, cylinder) and one of 8 colors. In contrast to CLEVR, we have a fixed material and size for objects. For the first image in the sequence, the object placement is fixed to the image center. For all the following images, the objects are placed in a random position while maintaining visibility (not completely occluded) and at a minimum distance from other objects.

To generate instructions, we use a simple text template that depends on the instruction number. For example, the second instruction in the sequence will have the following template:

“Add a [object color] [object shape] [relative position:

depth] it on the [relative position: horizontal]”

From the third instruction onward, the object position is described relative to two objects. These two objects are chosen randomly from the existing objects in the scene.

F. Additional implementation details

We use 300-dimensional GloVe¹ word embeddings for representing the words in each instruction q_t . These word embeddings are encoded using a bi-directional-GRU to obtain a 1024-dimensional instruction encoding d_t . All state dimensions for the higher level GRU R are set to 1024. The output of the conditioning augmentation module is also 1024-dimensional.

The code for this project was implemented in PyTorch [3]. For the generator and discriminator optimizers, “betas” was set to (0.0, 0.9) and weight decay was set to 0. The learning rates for the image encoding modules were set to 0.006. Gradient norm was clipped at 50. For each training experiment, we used a batch size of 32 over 2 NVIDIA P100 GPUs.

G. Additional language encoder experiments

We experimented with using skip-thought encoding for sentences instead of training the bi-directional-GRU encoder over GloVe embeddings. For the paper, we chose to use the latter since it performed better.

We also experimented with passing the previous image through the language encoder, but observed that it was easier for the model to generate an accurate image when the previous image features are passed to the Generator directly.

References

- [1] Jin-Hwa Kim, Nikita Kitaev, Xinlei Chen, Marcus Rohrbach, Byoung-Tak Zhang, Yuandong Tian, Dhruv Batra, and Devi Parikh, “CoDraw: Collaborative drawing as a testbed for grounded goal-driven communication,” in *Proceedings of the 57th Conference of the Association for Computational Linguistics*, 2019, pp. 6495–6513. [Online]. Available: <https://www.aclweb.org/anthology/P19-1651>
- [2] Blender Online Community, “Blender - a 3D modelling and rendering package,” 2016. [Online]. Available: <http://www.blender.org>
- [3] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer, “Automatic differentiation in PyTorch,” in *NIPS Autodiff Workshop*, 2017.

¹<http://nlp.stanford.edu/data/glove.840B.300d.zip>

		<p>Objects detected in generated image: cube gray, cube brown, sphere gray, sphere blue</p> <p>Objects detected in ground truth image: cube yellow, sphere yellow, cylinder green, cylinder brown, cylinder cyan</p> <p>Recall: 0.00</p> <p>rsim: 0.00</p> <p>Explanation: None of the correct objects are drawn.</p>
		<p>Objects detected in generated image: cube gray, cylinder gray, cylinder purple, cylinder cyan</p> <p>Objects detected in ground truth image: cube purple, sphere yellow, cylinder gray, cylinder brown, cylinder cyan</p> <p>Recall: 0.4</p> <p>rsim: 0.25</p> <p>Explanation: The cyan and gray cylinders are the only two objects detected in the generated image from the five ground truth objects detected.</p>
		<p>Objects detected in generated image: cube red, cube green, sphere brown, cylinder gray</p> <p>Objects detected in ground truth image: cube red, sphere red, sphere brown, sphere yellow, cylinder brown</p> <p>Recall: 0.4</p> <p>rsim: 0.35</p> <p>Explanation: The red cube and brown sphere are detected common to both images. Most of the relationships of these two and the center are correct.</p>
		<p>Objects detected in generated image: cube gray, cube red, cube yellow, sphere purple</p> <p>Objects detected in ground truth image: cube gray, cube red, cube blue, cube yellow, sphere purple</p> <p>Recall: 0.8</p> <p>rsim: 0.45</p> <p>Explanation: Only the blue cube is not detected in the generated image. Several spatial relationships of the common objects and the center are incorrect.</p>
		<p>Objects detected in generated image: sphere brown, sphere cyan, cylinder blue, cylinder purple, cylinder cyan</p> <p>Objects detected in ground truth image: cube cyan, sphere brown, cylinder blue, cylinder purple, cylinder cyan</p> <p>Recall: 0.8</p> <p>rsim: 0.675</p> <p>Explanation: Cyan cube detected in ground truth image is missing from the generated image. Most spatial relationships of the common objects and center are correct.</p>
		<p>Objects detected in generated image: sphere green, sphere purple, cylinder green, cylinder purple, cylinder cyan</p> <p>Objects detected in ground truth image: sphere green, sphere purple, cylinder green, cylinder purple, cylinder cyan</p> <p>Recall: 1.0</p> <p>rsim: 0.76</p> <p>Explanation: All the objects are detected correctly but some of the spatial relationships are incorrect.</p>
		<p>Objects detected in generated image: cube red, cube blue, cube yellow, sphere brown, cylinder blue</p> <p>Objects detected in ground truth image: cube red, cube blue, cube yellow, sphere brown, cylinder blue</p> <p>Recall: 1.00</p> <p>rsim: 1.00</p> <p>Explanation: All the objects are detected correctly and are in the correct relative positions.</p>

Figure 1. **Column 1:** Generated final image; **Column 2:** Ground truth final image; **Column 3:** List of objects detected in the generated and ground truth image, the recall on object detection, the value of the relational similarity (rsim) metric. The examples have been selected to qualitatively show examples with diverse score values between the minimum (0) and the maximum (1) values of the rsim metric.

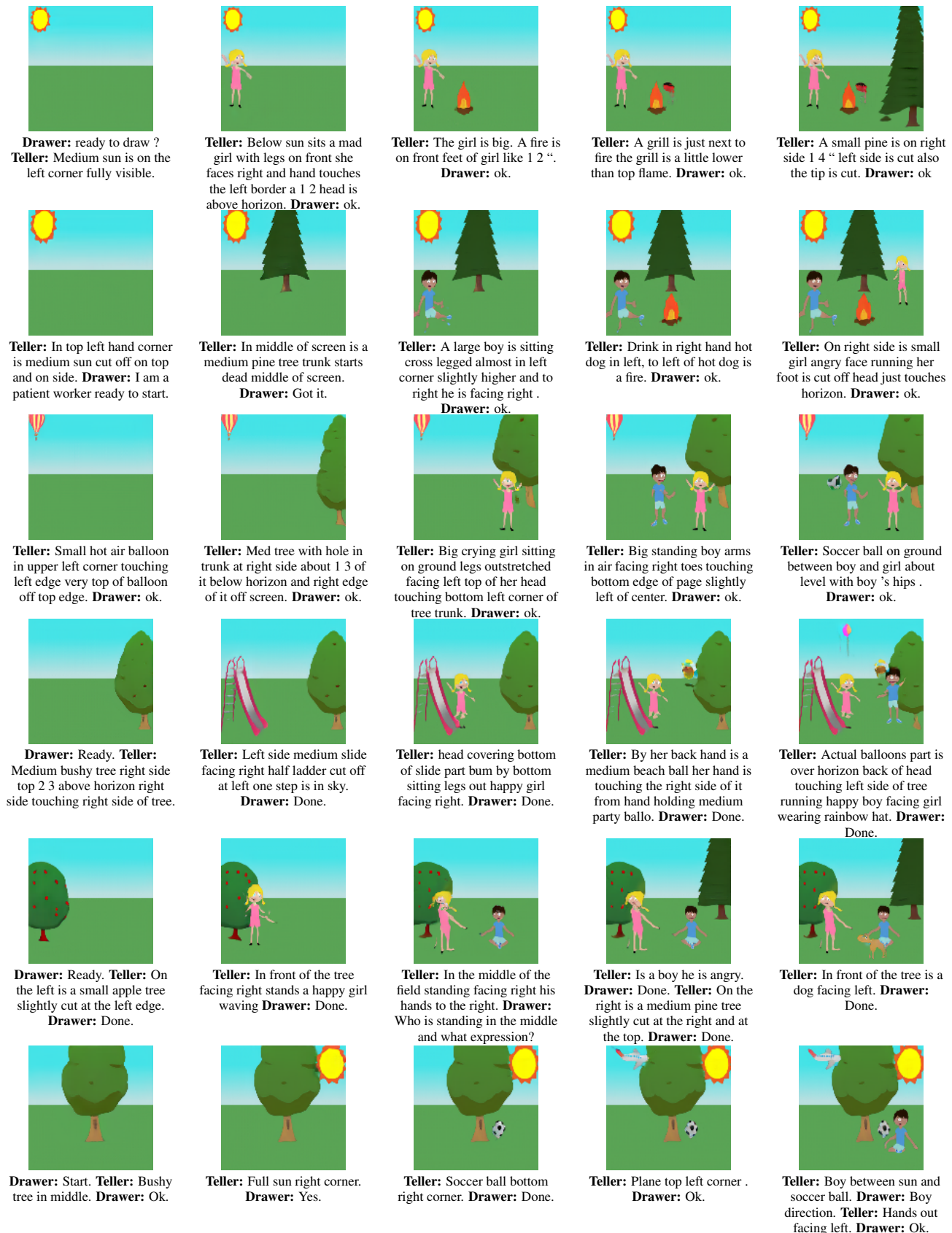


Figure 2. Generation examples from our best model (*D Subtract*) for the CoDraw dataset.

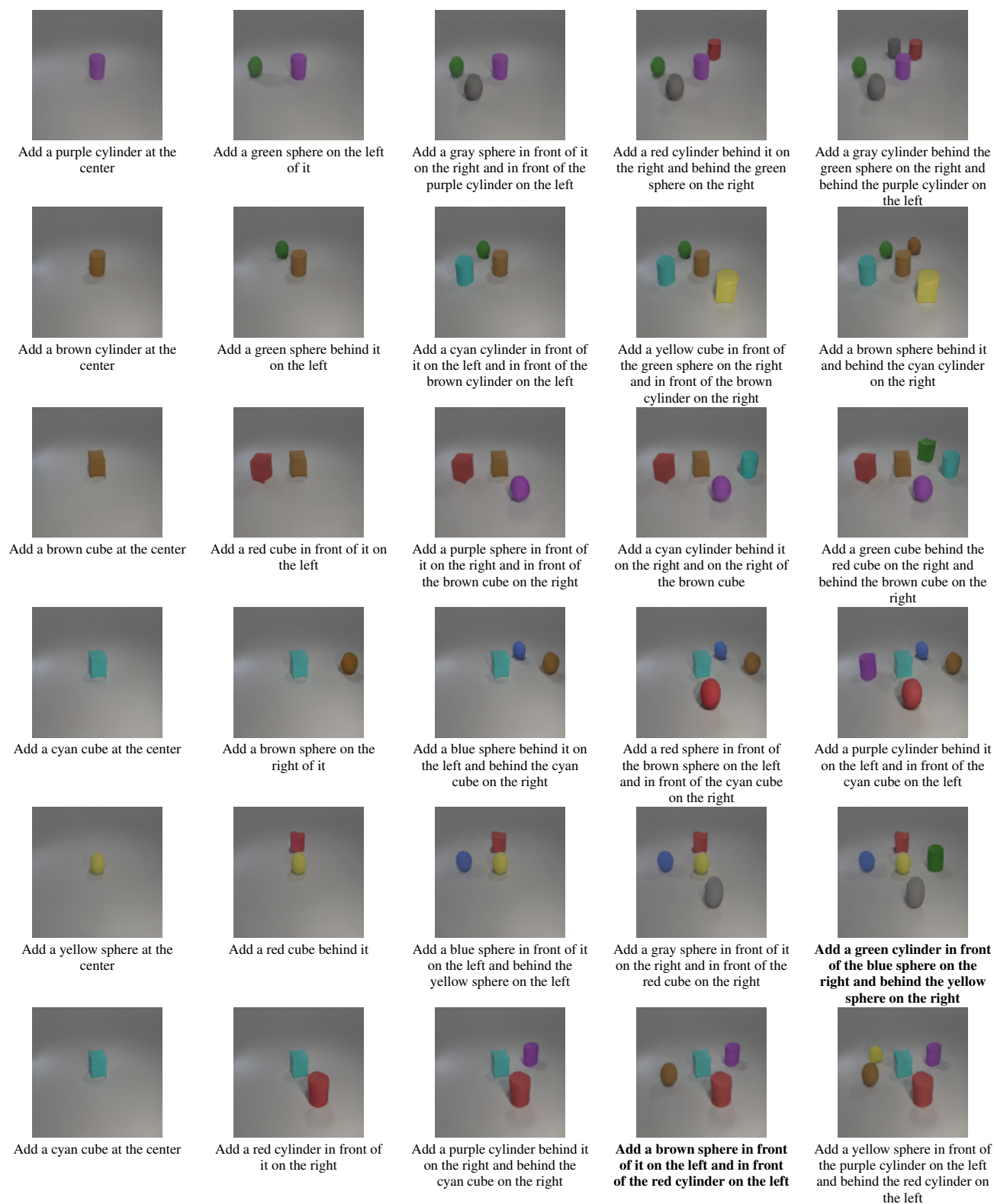


Figure 3. Generation examples from our best model (D Subtract) for the i-CLEVR dataset. Instructions where the model made a mistake are marked in bold.

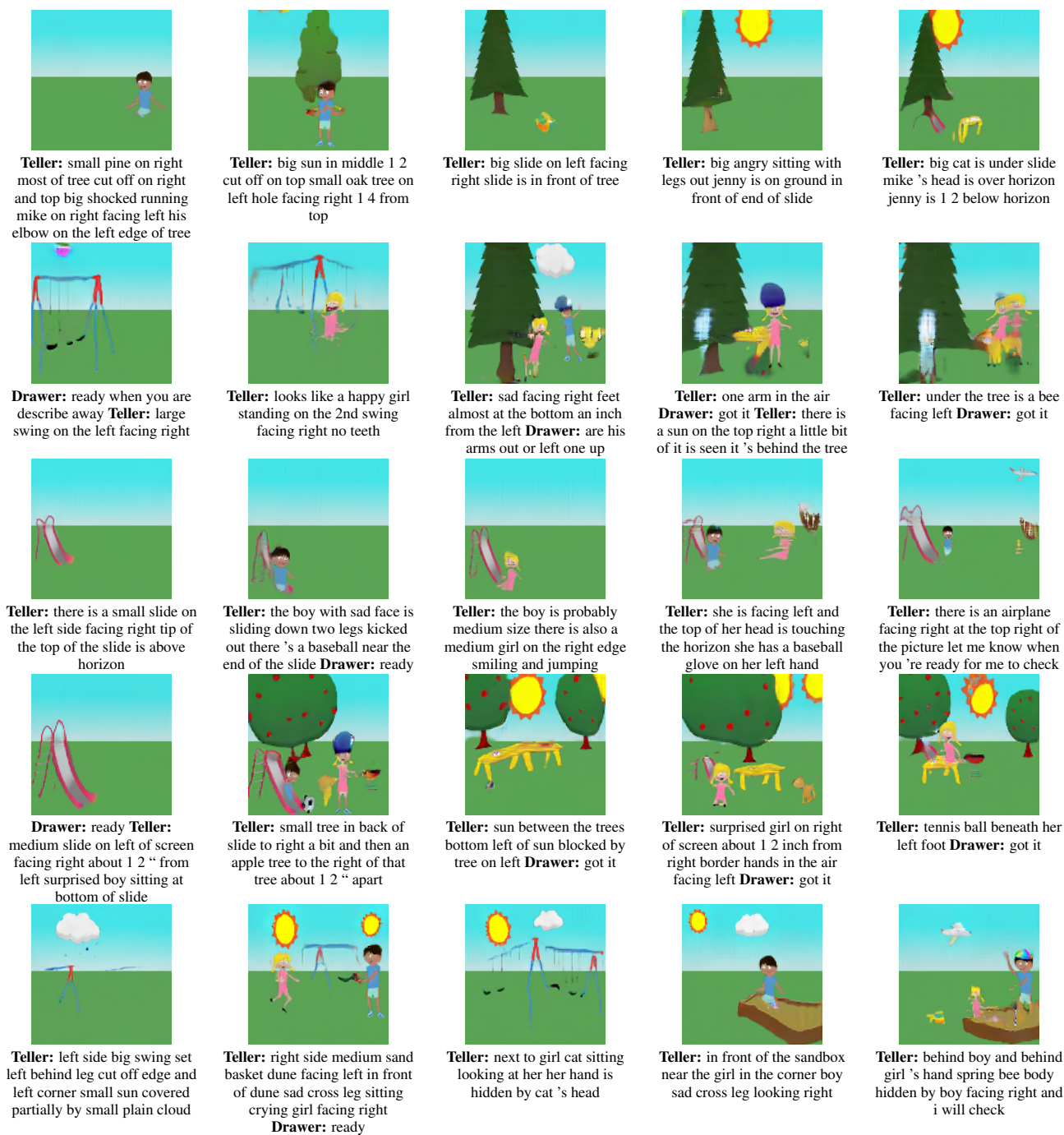
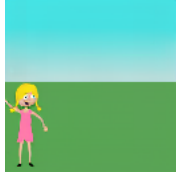


Figure 4. Random selection of examples generated by our Baseline model for the CoDraw dataset.



Drawer: ready **Teller:** 1 girl happy running facing right 0 2 inch from bottom to top and 0 2 inches from left to right with a chef



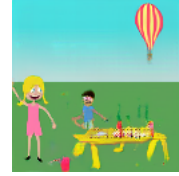
Teller: hat on her a table 1 and a half inches from left to right 1 2 inches from bottom to top with a pizza in the middle



Teller: and and a hot dog facing left to the right of the pizza a fire 0 1 inches from bottom to top 0 4 inches from right to left and above a



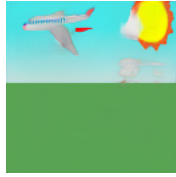
Teller: tent facing left and top of the tent is above the horizon line 0 1 inches and right is cover a little bit and above



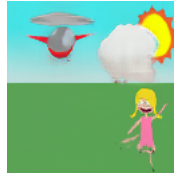
Teller: a air balloon small like 1 2 inches from right to left and 0 2 inches from top to bottom and that 's it



Drawer: what 's in the picture **Teller:** in the top right is a sun covered partially by clouds **Drawer:** is it a large sun



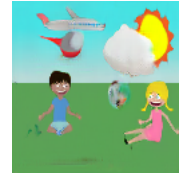
Teller: in the middle left is a helicopter **Drawer:** which way is it facing



Teller: yes a large sun **Teller:** heli is facing to the right **Teller:** tail is to the left **Teller:** on bottom right is a girl in pink with left arm raised



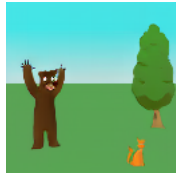
Teller: on his mouth is an o shape **Teller:** to the right of the boy is a dog with a blue collar



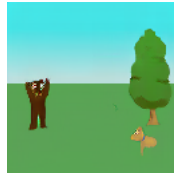
Teller: above the boys left hand in the blue sky is a yellow ball



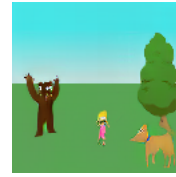
Drawer: may you please tell the first thing to draw **Teller:** there is a small tree on the right side of the scene sort of in the background **Drawer:**



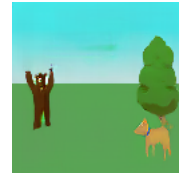
Teller: there is a bear to the left but in the foreground of the tree **Drawer:** what next



Teller: the bear is small **Drawer:** **Teller:** the girl is to the left facing left looking at the bear with her leg out scared facing right **Drawer:**



Teller: there is a small dog below the girl and a angry boy to left facing left with a racket in the left hand **Drawer:** what is next



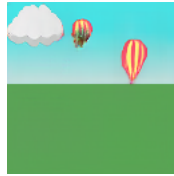
Teller: there is a small helicopter at the top in the sky in between the boy and girl **Drawer:** right above the bear



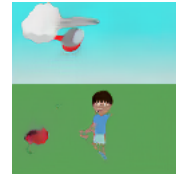
Teller: there is a cloud in the top left corner it is cut off on the top and left sides two puffs on the right and 3 puffs on the bottom



Teller: an inch from the right of the cloud are small balloons the orange balloon is on the right



Teller: airplane 1 4 inch to the right of the balloon the nose of the plane is in line with the yellow balloon it faces left



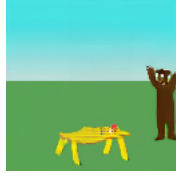
Teller: angry boy sitting below the plane facing left about 1 2 inch grass above and below him



Teller: a girl sits to the left of the boy facing him feet almost touching surprised wearing viking hat top brown part of hat touches horizon



Drawer: ready **Teller:** a small bear close to the right side one finger off picture small sliding on his left side bear left foot touching the sliding



Teller: table medium in the center front



Teller: medium pine tree behind the table one inch



Teller: sad girl standing far left part hand missing sad face medium



Teller: boy sitting on her right side look mad

Figure 5. Random selection of examples generated by our Mismatch model for the CoDraw dataset.

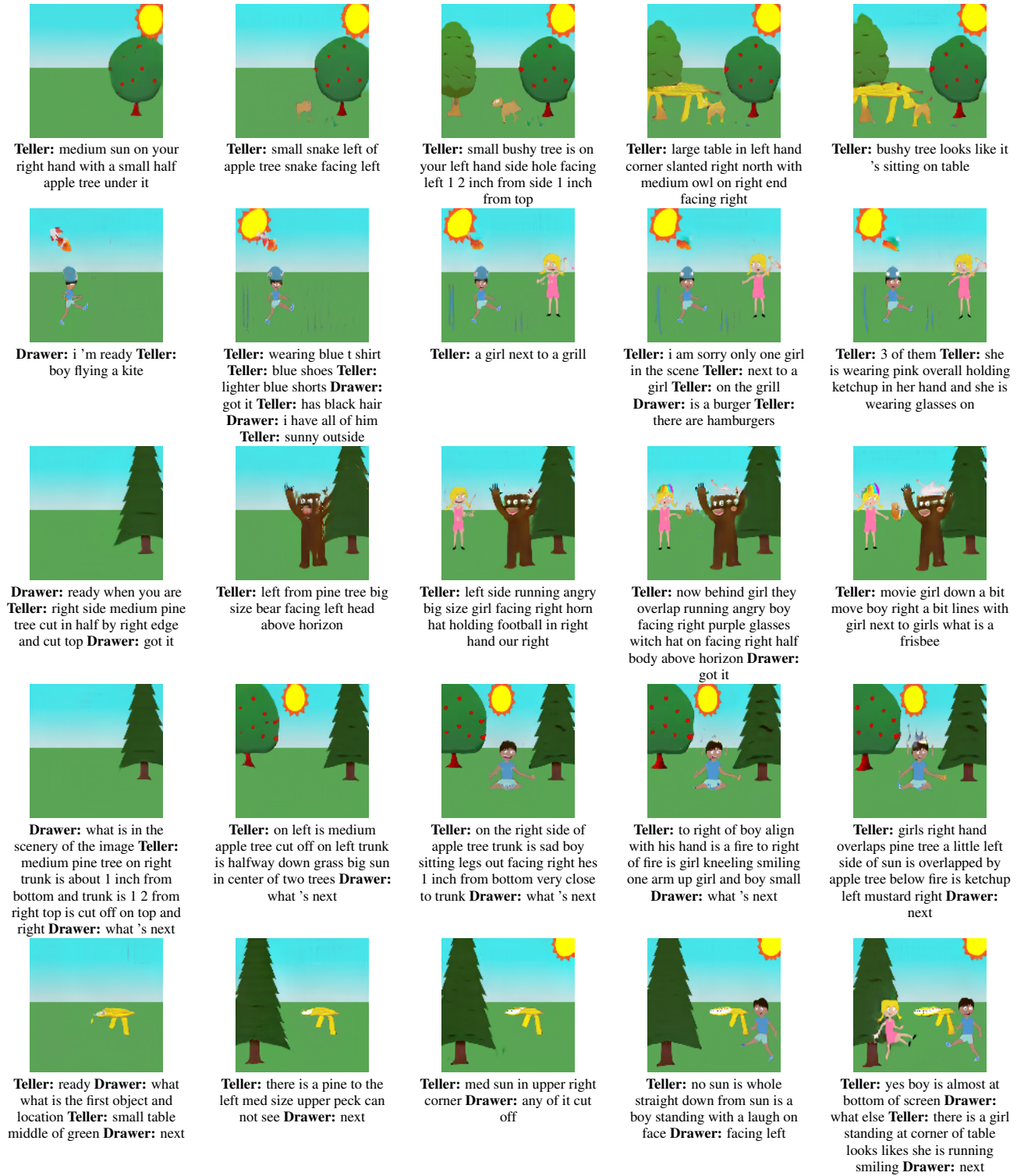


Figure 6. Random selection of examples generated by our G prior model for the CoDraw dataset.



Drawer: go **Teller:** large rain cloud left corner touches side cut off on top drops almost touch grass **Drawer:** next



Teller: large rocket on right tip of cloud flying left with very small girl sad legs out sitting on its upper wing **Drawer:** sitting on rocket the rocket is middle scene



Teller: rocket overcloud large regular cloud on right side cut off on top and side a bit surprised boy legs out facing left under cloud **Drawer:** next



Teller: cat facing boy 1 2 inch to left of his feet **Drawer:** next



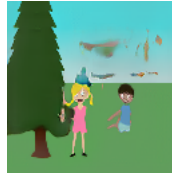
Teller: i will check and send adjustments **Drawer:** yes i do n't have the girl tell me where is the girl **Teller:** she is sitting on the rockets upper wing her back arm is under the window **Drawer:** check



Drawer: where is jenny and mike **Teller:** on left hand side of the screen 2 inches above bottom is a pine tree cut off at top



Teller: straight down from tree is jenny legs crossed facing right right arm in the air



Teller: next to jenny is mike same level standing facing right with arms out mouth open



Teller: on the right hand side inch from the bottom is a duck facing jenny and mike



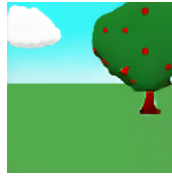
Teller: straight above duck is soccer ball



Teller: large cloud in left corner top and left are off screen



Teller: large apple tree right side top of trunk lines up with horizon right side and top of screen



Teller: in front of trunk little over to the right of trunk by right side is a soccer ball



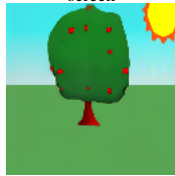
Teller: left side large girl angry face sitting cross legged facing right



Teller: neck on horizon line she 's holding a baseball in her up hand and wearing rainbow hat



Drawer: what do we have **Teller:** med sun right corner **Drawer:** and



Teller: middle of green with trees half in blue half in green is a apple tree med size **Drawer:** and



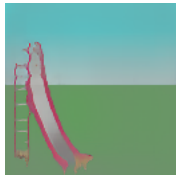
Teller: med boy standing on right of tree **Drawer:** face expression and where are his hands **Teller:** he has a tennis ball in right hand he is smiling showing teeth **Drawer:** and



Teller: right arm sticking out across from him is a girl almost to the left edge left arm down right hand with a ball glove **Drawer:** smiling please give more elaborations



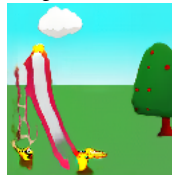
Teller: yes smiling looking right she is standing her middle is on the line of green **Drawer:** and



Teller: large slide to the left facing right with owl sitting on top of platform **Drawer:** top of handles or where we stand



Teller: small cloud in top center a bit to the right owl is on the platform **Drawer:** go



Teller: there is a medium to small apple tree on right half in blue half in green right side of tree cut off a bit **Drawer:** go



Teller: a dog directly under tree facing left **Drawer:** size



Teller: a girl standing in front of slide arms up smiling dog looks small

Figure 7. Random selection of examples generated by our Aux model for the CoDraw dataset.

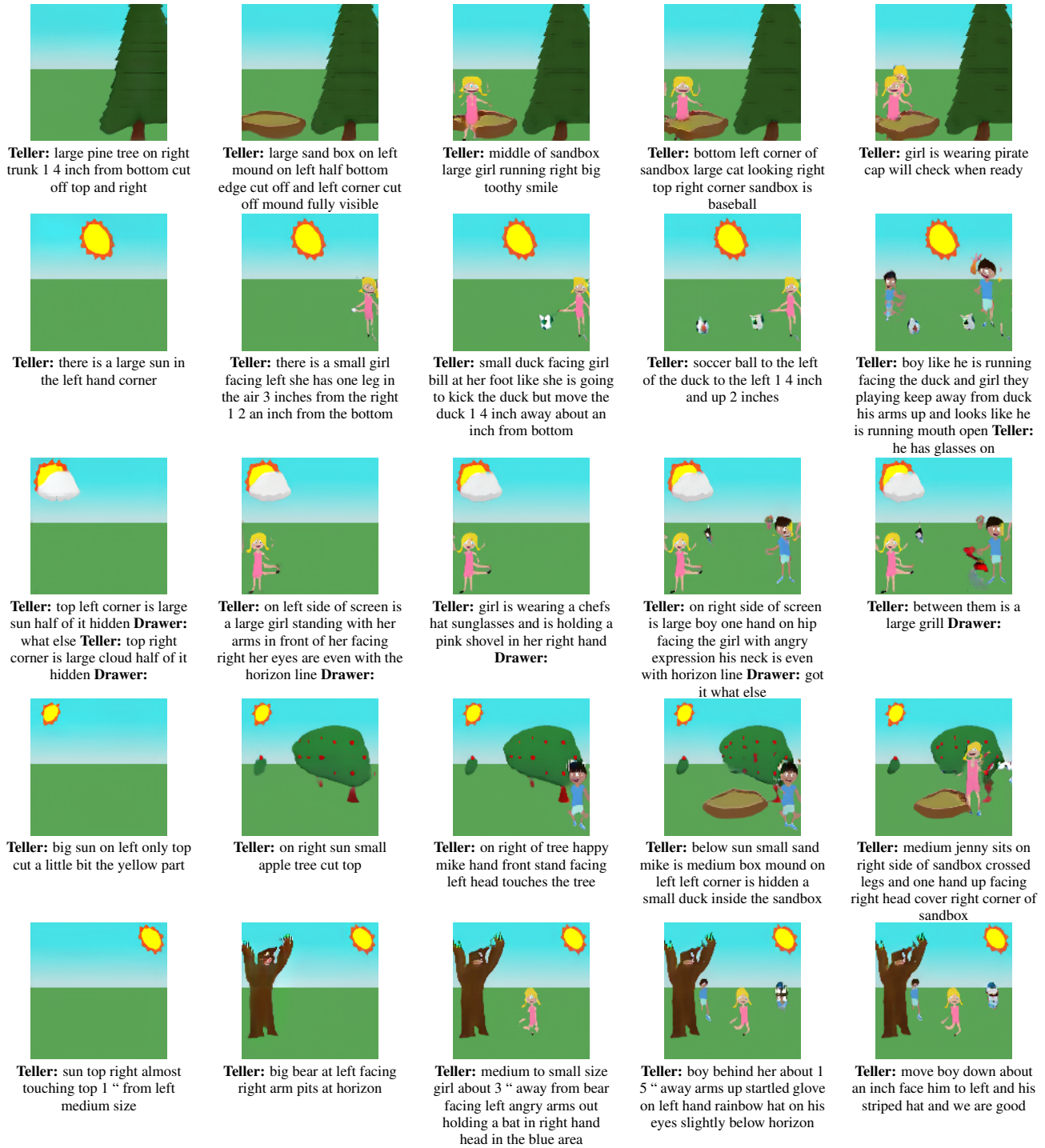
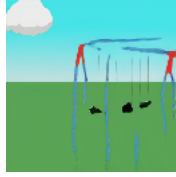


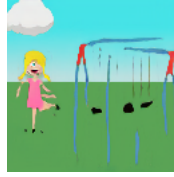
Figure 8. Random selection of examples generated by our D Concat model for the CoDraw dataset.



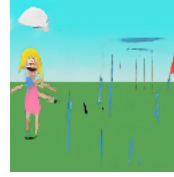
Teller: big cloud top left side
Drawer: got it



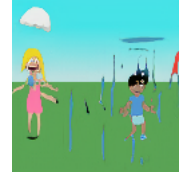
Teller: on the right side is a swing big size
Drawer: any parts cut off



Teller: girl on the left side neg horizon one arm up facing right happy face
Drawer: got it



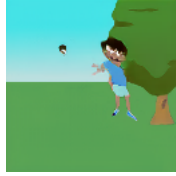
Teller: one part cut from swing just a bit from the right
Drawer: **Teller:** next to the right of the cloud is a basketball a bit over the cloud
Drawer:



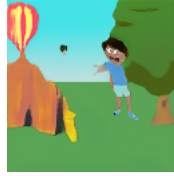
Teller: a boy is on the swing the right sit legs cross surprised face facing left color hat baseball glove
Drawer: got it



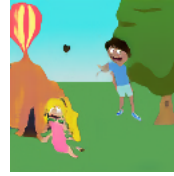
Drawer: ready **Teller:** big oak on right hole facing right hole almost touching horizon
Drawer: if its large it is huge how much is cut off on the right



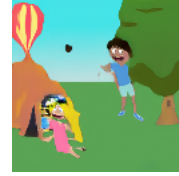
Teller: smiling big hands out mike on right left trunk point touching his back hair above horizon



Teller: small hot balloon on left 1 4 from top 1 in from left big tent on left facing right cut off slightly on back top above horizon



Teller: smiling big hands out jenny is in front of left opening tent



Teller: she has an owl sitting on her left wrist her hand is in the dark opening
Drawer: she is facing him and large owl right



Drawer: go **Teller:** small bushy tree facing left owl on right middle of tree



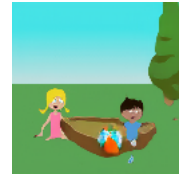
Teller: to the left of tree is a medium sandbox mound on right close to bottom
Drawer: next



Teller: girl sitting in left corner indian style smiling one arm up
Drawer: next



Teller: boy in right corner sitting indian style smiling with arms open both facing right
Drawer: next



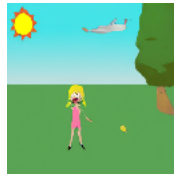
Teller: under boys right hand is cup in sand straw to left to left of cup is medium beach ball



Teller: tree hole facing left cut off from right side a little bit top hiding a bit of the sun



Teller: bumblebee with ear touching the bottom left of tree trunk facing to the right side
Drawer: got it



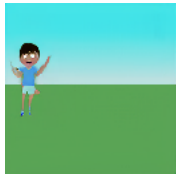
Teller: girl sitting smiling facing right hand behind her one inch from side wearing crown



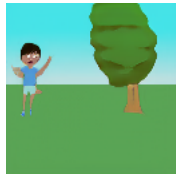
Teller: crown almost touches the horizon
Drawer: got it
Teller: boy faces girl sitting smiling his feet r half inch from hers and raised up a little he wears a beanie with top of it just at horizon



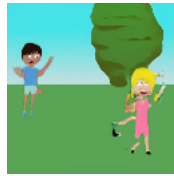
Teller: duck between the two with ducks feet level to boys top foot



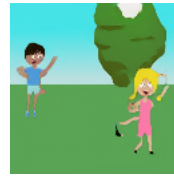
Teller: on the left an inch from the edge is a boy
Drawer: what is he doing
Teller: he is facing right standing one hand up teeth showing and holding a racket with one hand that is in the air



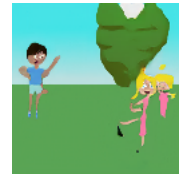
Teller: next to him a medium tree hole facing left head touches the tip of last branch truck aligns with his waist



Teller: on the right 1 inch from edge is a girl sad looking left one hand in the air head aligns with the boy's
Drawer: her left hand cut off



Teller: above her is a small cloud right above her can i check



Teller: no about 2 cm from the edge the hand
Teller: move her 1 more cm from the edge she is holding a yellow small ball in the hand in air

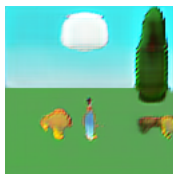
Figure 9. Random selection of examples generated by our D Subtract model for the CoDraw dataset.



Drawer: what is there **Teller:** big thunderbolt on left bolt facing right touching on left edge close to top **Drawer:** and **Teller:** big raindrop cloud to the right of thunderbolt cloud cutting it off on right side a little big shocked jenny with arms up **Drawer:** where is she **Teller:** head right below horizon sad big sitting mike with legs facing right is beside her **Drawer:** and **Teller:** he 's wearing a star hat soccer ball is covering his left foot and shin **Drawer:** and



Teller: ready **Drawer:** and ready **Teller:** upper right corner large sun with right edges a bit cut off and top cut off **Teller:** under sun happy boy standing facing left with right arm up his shoulders just above horizon line **Teller:** he is wearing a pirate hat it touches on of the sun tips on the left side **Teller:** happy girl kicking on left side she is about 1 5 inches in from left side her mouth is at the horizon line **Drawer:** got it **Teller:** just a tiny bit off of girls kicking foot is a beach ball a cloud is over the girl towards the right center **Drawer:** is the cloud on the right or the sun you said sun upper right corner



Teller: boy left side kicking leg facing right his half torso aligns with horizon he is shocked **Drawer:** go **Teller:** finger away from his leg soccer ball its bottom part touches horizon **Teller:** right side medium tree 1 4 cut off right side and trunk half way in grass with slight cut off as well right side hole facing left **Drawer:** go **Teller:** plain cloud top middle top part cut off big size in front of tree dog its legs behind completely cut off and it 's facing left **Teller:** near dog is a big cat its tail cover 's dog 's front leg slightly and facing right and then girl sitting smiling facing right

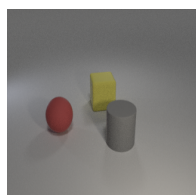


Drawer: ready **Teller:** top left facing left one 1 4 inch from side blade touch top small helicopter facing left **Drawer:** what 's in the left the helicopter **Teller:** nothing it is a 1 4 inch from side flying left **Drawer:** got it it 's tiny right **Teller:** yes **Teller:** below copter is large boy facing right arms out mouth open neck at horizon **Teller:** right of boy his top hand is on first plank is a large picnic table left top corner is highest point pie is there in corner **Drawer:** which side is the pie **Teller:** right of pie is large girl facing left standing with smile no teeth one arm up and one down pie top left corner **Drawer:** where is she to the horizon and she is in front of the table **Teller:** girl in front of table nose at horizon top right corner of table is ketchup **Drawer:** got it **Teller:** 1 2 inch from right side and 1 2 inch from horizon is large grill

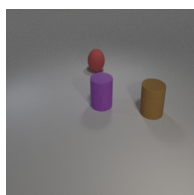


Drawer: ready **Teller:** there is a medium in the center of the sky just below the top edge **Drawer:** medium cloud **Teller:** oh sorry medium sun the medium cloud is down and to the right in the sky **Teller:** there is a small oak tree on the left an inch away from the left edge hole facing right 2 3s of the leaves are above the horizon **Teller:** on the right side the kids are both medium sized and facing left jenny is happy jumping half inch from the right edge

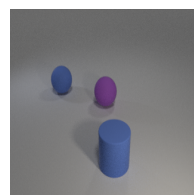
Figure 10. Random selection of examples generated by our Non-iterative model for the CoDraw dataset.



(a) **Left:** Initial Image **Right:** Final Image
Instruction: Add a yellow cylinder behind the gray cylinder on the right and behind the yellow cube on the right



(b) **Left:** Initial Image **Right:** Final Image
Instruction: Add a cyan cube behind the brown cylinder on the left and behind the purple cylinder on the left



(c) **Left:** Initial Image **Right:** Final Image
Instruction: Add a gray sphere behind the blue cylinder on the right and behind the purple sphere on the right

Figure 11. When GeNeVA-GAN is provided with an initial image different from the background image used during training, it still adds the desired object with the right properties at the correct location. The model was not trained in this setting and the success of this experiment demonstrates that it has learnt to preserve the existing canvas, understand the existing objects, and add new objects with the correct relationships to existing objects.