

1 We thank the reviewers for their helpful feedback and positive view of the work. We will aim to incorporate your  
 2 suggestions in future versions of the paper.

3 **Re: dataset release.** We are planning to release the new datasets in the next few weeks.

4 **Reviewer 1**

5 • **Add a speed comparison.** E-GQN usually reaches a given loss significantly (often 10x) faster in wall clock time (in  
 6 addition to having a lower final loss). On the two datasets (rfc, jaco) where final performance is similar, E-GQN is  
 7 somewhat slower due to processing  $\sim 30\%$  fewer samples per second. We will add a speed analysis to the final paper.

8 • **Provide visualization of the attention mechanism.** If accepted, we will add one to the supplemental materials.

9 • **Add an ablation study.** Since the aim of our work is to design & study a new attention mechanism, we hold the  
 10 remainder of the architecture constant, so the GQN numbers show the performance without the attention mechanism.

11 **Reviewer 2**

12 • **Discuss how approaches similar in spirit to EGQN could benefit training of conditional models in general.**  
 13 Good question. Geometric structure could also be exploited in settings with moving cameras. Temporal locality may  
 14 also be exploitable. We will expand on these in the conclusion.

15 • **Have the authors analyzed EGQN’s scene representation?** We agree this would be interesting, but considered it  
 16 out of scope for this paper as our main aim was improving model performance on more challenging datasets.

17 • **Discuss weaknesses or in general cons of the proposed model.** Thanks for pointing this out. The main weaknesses  
 18 of the method are inapplicability when there is little content overlap between context and target frames (which we  
 19 briefly mention in the paper), additional memory requirement relative to GQN, and slower training in terms of samples  
 20 per second (though usually faster training overall as mentioned above). We will expand in the paper.

21 **Reviewer 3**

22 • **Would be good to show other evaluation measures.** We updated the evaluation to include standard deviation,  
 23 RMSE, and ELBO in addition to existing evaluation for all 7 datasets. The new table is included below as **Figure 1**.

24 • **Exemplary images are not given for all datasets.** Thanks for pointing this out. We agree more images would help  
 25 readers better understand the models’ performance. We uploaded additional images from all datasets to the supplemental  
 26 website for the paper, and also included a few of them below as **Figure 2**.

27 • **How was the min Y value in Fig 4 computed?** We compute the minimum ELBO by assuming the KL term to be 0  
 28 and the mean of the output distribution to be the true target image. Per Section 4.2, we use a different output variance  
 29 hyperparameter for our datasets vs. those from the original GQN paper, which causes the different scaling of the ELBO.  
 30 We will add more detail to the description in the paper.

31 • **Would be nice to compare to other methods that build on the GQN, e.g. CGQN.** CGQN addresses the problem  
 32 of multiple simultaneous predictions being inconsistent with one another. Since we focused on improving the quality of  
 33 individual predictions, we considered it out of scope. Combining CGQN and E-GQN could be interesting future work.

Dataset	Mean Absolute Error (pixels)		Root Mean Squared Error (pixels)		ELBO (nats / dim)	
	GQN	E-GQN	GQN	E-GQN	GQN	E-GQN
rrc	7.40 $\pm$ 6.22	<b>3.59 <math>\pm</math> 2.10</b>	14.62 $\pm$ 12.77	<b>6.80 <math>\pm</math> 5.23</b>	0.5637 $\pm$ 0.0013	<b>0.5629 <math>\pm</math> 0.0008</b>
rfc	12.44 $\pm$ 12.89	<b>12.05 <math>\pm</math> 12.79</b>	<b>26.80 <math>\pm</math> 21.35</b>	27.65 $\pm$ 20.72	<b>0.5637 <math>\pm</math> 0.0011</b>	0.5639 $\pm$ 0.0012
jaco	4.30 $\pm$ 1.12	<b>4.00 <math>\pm</math> 0.90</b>	8.58 $\pm$ 2.94	<b>7.43 <math>\pm</math> 2.32</b>	0.5634 $\pm$ 0.0007	<b>0.5631 <math>\pm</math> 0.0005</b>
sm7	3.13 $\pm$ 1.30	<b>2.14 <math>\pm</math> 0.53</b>	9.97 $\pm$ 4.34	<b>5.63 <math>\pm</math> 2.21</b>	0.5637 $\pm$ 0.0009	<b>0.5628 <math>\pm</math> 0.0004</b>
oab	10.99 $\pm$ 5.13	<b>5.47 <math>\pm</math> 2.54</b>	22.11 $\pm$ 8.00	<b>10.39 <math>\pm</math> 4.55</b>	1.2587 $\pm$ 0.0018	<b>1.2569 <math>\pm</math> 0.0011</b>
disco	18.86 $\pm$ 7.16	<b>12.46 <math>\pm</math> 9.27</b>	32.72 $\pm$ 6.32	<b>22.04 <math>\pm</math> 11.08</b>	1.2635 $\pm$ 0.0055	<b>1.2574 <math>\pm</math> 0.0007</b>
rro	10.12 $\pm$ 5.15	<b>6.59 <math>\pm</math> 3.23</b>	19.63 $\pm$ 9.14	<b>12.08 <math>\pm</math> 6.52</b>	1.2573 $\pm$ 0.0011	<b>1.2566 <math>\pm</math> 0.0009</b>

Figure 1: Performance of GQN and E-GQN. Note: ELBO scaling is due to different choices of output variance.

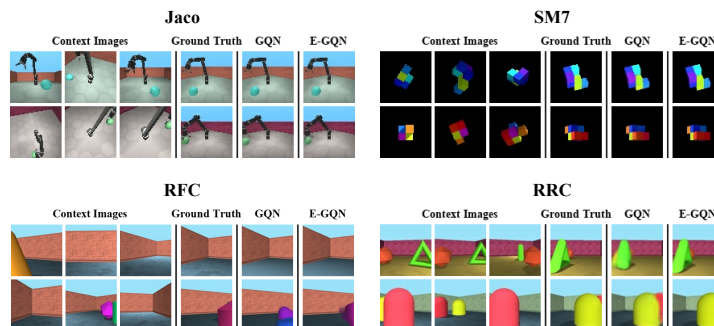


Figure 2: Randomly chosen samples from GQN and E-GQN