

Table 1: Performance of models on COGITAO experiments for the *CompGen* and *EnvGen* studies across different experiment settings (ES). We report the ID (in-domain) and OOD (out-of-domain) results for Grid accuracy (i.e., percentage of perfect match) and the relative drop ( $\Delta = \text{ID} - \text{OOD}$ ) averaged over all the experiments within the experiment setting.

Study	ES	Vanilla-ViT			Grid-ViT			LLaD	
		ID	OOD	$\Delta$	ID	OOD	$\Delta$	ID	OOD
CompGen	C1	28.4	15.7	12.7	<b>71.1</b>	18.7	52.4	54.2	<b>26.2</b>
	C2	34.0	15.6	18.4	<b>69.0</b>	<b>18.6</b>	50.4	49.2	15.2
	C3	15.1	1.0	14.1	75.6	5.1	70.5	<b>78.9</b>	<b>6.4</b>
EnvGen	G1	98.3	76.7	21.6	99.2	<b>90.0</b>	9.2	<b>99.4</b>	89.1
	G2	68.1	0.0	68.1	97.0	<b>73.0</b>	24.0	<b>99.3</b>	64.3
	G3	57.1	17.5	39.6	<b>78.5</b>	<b>27.2</b>	51.3	76.9	<b>27.2</b>
	G4	47.6	21.3	26.3	75.2	13.8	61.4	<b>76.3</b>	<b>35.6</b>
	G5	70.0	0.0	70.0	<b>98.0</b>	0.2	97.8	71.0	<b>14.0</b>