

Table 7: **(Q1)**: Metrics to assess the goodness of *layer 1* and *31* to detect personas (Level 2) in *Granite-7B-Instruct*. Metrics: Silhouette Score (SH), Calinski-Harabasz Score (CH), Euclidean Distance (ED), Davies-Bouldin Score (DB).

Model	Dimension	layer	SH (\uparrow)	CH (\uparrow)	DB (\downarrow)	ED (\uparrow)
<i>Granite-7B-Instruct</i>	AGREE	1	0.420 \pm 0.0312	210.7 \pm 75.01	0.925 \pm 0.225	0.287 \pm 0.019
		31	0.611 \pm 0.0000	969.0 \pm 0.000	0.605 \pm 0.000	17.62 \pm 0.000
	CONSC	1	0.561 \pm 0.0000	578.4 \pm 0.003	0.709 \pm 0.000	0.250 \pm 0.000
		31	0.625 \pm 0.0000	1033.5 \pm 0.002	0.573 \pm 0.000	16.83 \pm 0.000
	OPEN	1	0.491 \pm 0.0000	312.4 \pm 0.000	0.912 \pm 0.000	0.265 \pm 0.000
		31	0.991 \pm 0.0000	37132.7 \pm 0.038	0.004 \pm 0.000	2073.7 \pm 0.000
	EXTRA	1	0.529 \pm 0.0000	397.0 \pm 0.000	0.820 \pm 0.000	0.265 \pm 0.000
		31	0.618 \pm 0.0000	938.7 \pm 0.000	0.595 \pm 0.000	18.58 \pm 0.000
	NEURO	1	0.491 \pm 0.0362	270.3 \pm 156.9	0.707 \pm 0.157	0.279 \pm 0.044
		31	0.637 \pm 0.0000	1128.8 \pm 0.006	0.549 \pm 0.000	18.66 \pm 0.000
	VIRTUE	1	0.533 \pm 0.0001	350.0 \pm 0.002	0.803 \pm 0.000	0.242 \pm 0.000
		31	0.996 \pm 0.0000	164373.4 \pm 0.028	0.000 \pm 0.000	27948.3 \pm 0.000
	RELAT	1	0.491 \pm 0.0000	332.1 \pm 0.000	0.927 \pm 0.000	0.256 \pm 0.000
		31	0.614 \pm 0.000	828.7 \pm 0.000	0.602 \pm 0.000	17.32 \pm 0.000
	DEONT	1	0.462 \pm 0.0023	171.7 \pm 65.54	0.678 \pm 0.241	0.335 \pm 0.048
		31	0.605 \pm 0.0000	871.0 \pm 0.005	0.622 \pm 0.000	17.42 \pm 0.000
	UTILI	1	0.481 \pm 0.0000	286.0 \pm 0.001	0.956 \pm 0.000	0.262 \pm 0.000
		31	0.994 \pm 0.0000	80355.2 \pm 0.057	0.002 \pm 0.000	3690.19 \pm 0.000
	NIHIL	1	0.441 \pm 0.0208	184.3 \pm 65.67	0.730 \pm 0.229	0.300 \pm 0.025
		31	0.602 \pm 0.0000	292.7 \pm 0.000	0.678 \pm 0.000	18.44 \pm 0.000
	CONS	1	0.511 \pm 0.0659	284.2 \pm 160.3	0.683 \pm 0.126	0.276 \pm 0.025
		31	0.993 \pm 0.0000	8423.03 \pm 0.000	0.001 \pm 0.000	21282.9 \pm 0.000
	LIBER	1	0.536 \pm 0.0529	378.8 \pm 236.3	0.624 \pm 0.104	0.275 \pm 0.039
		31	0.993 \pm 0.0000	72452.4 \pm 0.017	0.002 \pm 0.000	3292.78 \pm 0.000
	LGBTQ	1	0.497 \pm 0.0509	290.2 \pm 167.2	0.682 \pm 0.126	0.283 \pm 0.039
		31	0.593 \pm 0.0000	838.5 \pm 0.010	0.643 \pm 0.000	15.95 \pm 0.000
	IMMI	1	0.510 \pm 0.0000	346.1 \pm 0.000	0.871 \pm 0.000	0.257 \pm 0.000
		31	0.617 \pm 0.0000	989.7 \pm 0.000	0.594 \pm 0.000	17.32 \pm 0.000