

Table 6: **(Q1)**: Metrics to assess the goodness of *layer 1* and *31* to detect personas (Level 2) in *Mistral-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>Mistral-7B-Instruct</i>	AGREE	1	0.370 \pm 0.082	176.5 \pm 110.62	0.690 \pm 0.084	0.110 \pm 0.006
		31	0.251 \pm 0.000	164.1 \pm 0.00	1.743 \pm 0.0008	5.403 \pm 0.003
	CONSC	1	0.242 \pm 0.113	96.8 \pm 17.79	1.753 \pm 0.505	0.063 \pm 0.02
		31	0.171 \pm 0.010	114.6 \pm 9.23	2.140 \pm 0.086	4.249 \pm 0.122
	OPEN	1	0.423 \pm 0.020	176.8 \pm 13.84	1.006 \pm 0.596	0.107 \pm 0.023
		31	0.235 \pm 0.000	194.2 \pm 0.01	1.693 \pm 0.000	5.147 \pm 0.000
	EXTRA	1	0.459 \pm 0.041	172.6 \pm 39.78	0.558 \pm 0.096	0.122 \pm 0.020
		31	0.327 \pm 0.000	225.2 \pm 0.000	1.102 \pm 0.000	8.020 \pm 0.000
	NEURO	1	0.392 \pm 0.109	123.9 \pm 7.74	0.877 \pm 0.573	0.107 \pm 0.032
		31	0.203 \pm 0.001	137.4 \pm 3.31	1.873 \pm 0.143	5.182 \pm 0.271
	VIRTUE	1	0.389 \pm 0.098	110.1 \pm 7.74	0.961 \pm 0.666	0.104 \pm 0.032
		31	0.272 \pm 0.000	211.5 \pm 1.18	1.448 \pm 0.056	5.499 \pm 0.079
	RELAT	1	0.468 \pm 0.077	378.1 \pm 226.82	0.794 \pm 0.070	0.119 \pm 0.005
		31	0.226 \pm 0.013	126.5 \pm 13.76	1.370 \pm 0.081	6.602 \pm 0.295
	DEONT	1	0.392 \pm 0.085	201.0 \pm 120.04	0.740 \pm 0.175	0.123 \pm 0.011
		31	0.242 \pm 0.001	155.4 \pm 0.31	1.587 \pm 0.082	5.631 \pm 0.118
	UTILI	1	0.445 \pm 0.064	188.0 \pm 62.33	0.563 \pm 0.087	0.135 \pm 0.012
		31	0.319 \pm 0.000	289.4 \pm 0.000	1.215 \pm 0.000	6.743 \pm 0.000
	NIHIL	1	0.529 \pm 0.000	436.5 \pm 0.000	0.555 \pm 0.000	0.130 \pm 0.000
		31	0.177 \pm 0.011	132.1 \pm 17.07	2.058 \pm 0.139	5.067 \pm 0.278
	CONS	1	0.360 \pm 0.098	98.9 \pm 25.45	1.253 \pm 0.559	0.083 \pm 0.026
		31	0.230 \pm 0.004	152.9 \pm 0.31	1.843 \pm 0.023	4.927 \pm 0.072
	LIBER	1	0.511 \pm 0.026	119.4 \pm 49.65	0.495 \pm 0.100	0.119 \pm 0.002
		31	0.189 \pm 0.026	119.4 \pm 5.45	1.862 \pm 0.326	4.769 \pm 0.798
	LGBTQ	1	0.448 \pm 0.041	136.8 \pm 40.46	0.585 \pm 0.101	0.122 \pm 0.006
		31	0.255 \pm 0.000	215.0 \pm 0.01	1.478 \pm 0.002	6.092 \pm 0.006
	IMMI	1	0.437 \pm 0.034	140.0 \pm 93.28	0.558 \pm 0.046	0.126 \pm 0.009
		31	0.200 \pm 0.006	149.7 \pm 11.71	1.928 \pm 0.080	4.945 \pm 0.155