

Table 5: **(Q1)**: Metrics to assess the goodness of *layer 1* and *31* to detect personas (Level 2) in *Llama3-8B-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>Llama3-8B-Instruct</i>	AGREE	1	0.500 $\pm$ 0.0741	340.6 $\pm$ 163.9	0.731 $\pm$ 0.072	0.403 $\pm$ 0.012
		31	0.792 $\pm$ 0.0000	3264.5 $\pm$ 0.002	0.326 $\pm$ 0.000	27.57 $\pm$ 0.000
	CONSC	1	0.635 $\pm$ 0.0000	718.8 $\pm$ 0.026	0.569 $\pm$ 0.000	0.370 $\pm$ 0.000
		31	0.813 $\pm$ 0.0000	4150.4 $\pm$ 0.003	0.285 $\pm$ 0.000	27.47 $\pm$ 0.000
	OPEN	1	0.602 $\pm$ 0.0000	570.2 $\pm$ 0.0005	0.645 $\pm$ 0.000	0.414 $\pm$ 0.000
		31	0.795 $\pm$ 0.0000	3564.1 $\pm$ 0.0125	0.319 $\pm$ 0.000	27.60 $\pm$ 0.000
	EXTRA	1	0.578 $\pm$ 0.0000	527.5 $\pm$ 0.001	0.705 $\pm$ 0.000	0.382 $\pm$ 0.000
		31	0.788 $\pm$ 0.0000	3176.5 $\pm$ 0.001	0.330 $\pm$ 0.000	27.47 $\pm$ 0.000
	NEURO	1	0.584 $\pm$ 0.0000	615.0 $\pm$ 0.065	0.686 $\pm$ 0.000	0.378 $\pm$ 0.000
		31	0.796 $\pm$ 0.0000	3372.4 $\pm$ 0.001	0.306 $\pm$ 0.000	27.22 $\pm$ 0.000
	VIRTUE	1	0.614 $\pm$ 0.0000	644.7 $\pm$ 0.007	0.639 $\pm$ 0.000	0.402 $\pm$ 0.000
		31	0.797 $\pm$ 0.0000	3471.5 $\pm$ 0.000	0.309 $\pm$ 0.000	27.24 $\pm$ 0.000
	RELAT	1	0.500 $\pm$ 0.0000	405.4 $\pm$ 0.000	0.935 $\pm$ 0.000	0.440 $\pm$ 0.000
		31	0.814 $\pm$ 0.0000	3960.0 $\pm$ 0.002	0.276 $\pm$ 0.000	28.06 $\pm$ 0.000
	DEONT	1	0.523 $\pm$ 0.0230	315.1 $\pm$ 104.9	0.758 $\pm$ 0.000	0.479 $\pm$ 0.054
		31	0.824 $\pm$ 0.0000	4297.7 $\pm$ 0.004	0.254 $\pm$ 0.000	28.17 $\pm$ 0.000
	UTILI	1	0.542 $\pm$ 0.0405	326.0 $\pm$ 160.5	0.645 $\pm$ 0.000	0.454 $\pm$ 0.062
		31	0.798 $\pm$ 0.0000	3587.6 $\pm$ 0.001	0.301 $\pm$ 0.000	27.59 $\pm$ 0.000
	NIHIL	1	0.556 $\pm$ 0.0000	472.8 $\pm$ 0.001	0.810 $\pm$ 0.000	0.421 $\pm$ 0.000
		31	0.808 $\pm$ 0.0000	3909.8 $\pm$ 0.016	0.282 $\pm$ 0.000	27.20 $\pm$ 0.000
	CONS	1	0.685 $\pm$ 0.0001	860.7 $\pm$ 0.024	0.475 $\pm$ 0.000	0.400 $\pm$ 0.000
		31	0.818 $\pm$ 0.0000	4529.9 $\pm$ 0.007	0.266 $\pm$ 0.000	27.17 $\pm$ 0.000
	LIBER	1	0.692 $\pm$ 0.0001	897.2 $\pm$ 0.023	0.469 $\pm$ 0.000	0.387 $\pm$ 0.000
		31	0.803 $\pm$ 0.0000	3809.3 $\pm$ 0.004	0.294 $\pm$ 0.000	26.80 $\pm$ 0.000
	LGBTQ	1	0.634 $\pm$ 0.0000	703.1 $\pm$ 0.014	0.573 $\pm$ 0.000	0.399 $\pm$ 0.000
		31	0.774 $\pm$ 0.0000	2815.3 $\pm$ 0.034	0.352 $\pm$ 0.000	27.62 $\pm$ 0.000
	IMMI	1	0.570 $\pm$ 0.0000	476.5 $\pm$ 0.000	0.730 $\pm$ 0.000	0.424 $\pm$ 0.000
		31	0.806 $\pm$ 0.0000	3756.4 $\pm$ 0.002	0.288 $\pm$ 0.000	27.67 $\pm$ 0.000