

## Revisions made in this submission.

Below we list **all changes introduced since the last review cycle**, grouped by the issues reviewers raised. Each bullet notes the exact additions or rewrites and where they appear in the revised, fully anonymised manuscript.

- **Finer-grained embedding analysis (R-9YM2)** — Added a **mean-centring experiment**. Figure 1 in section 3.1, before plots was cosine-similarity distributions which is moved to appendix and *later* added mean centring plot, and describe the outcome and quantifies how centring removes global bias or ADA/PaLM but not most open-source LLMs.

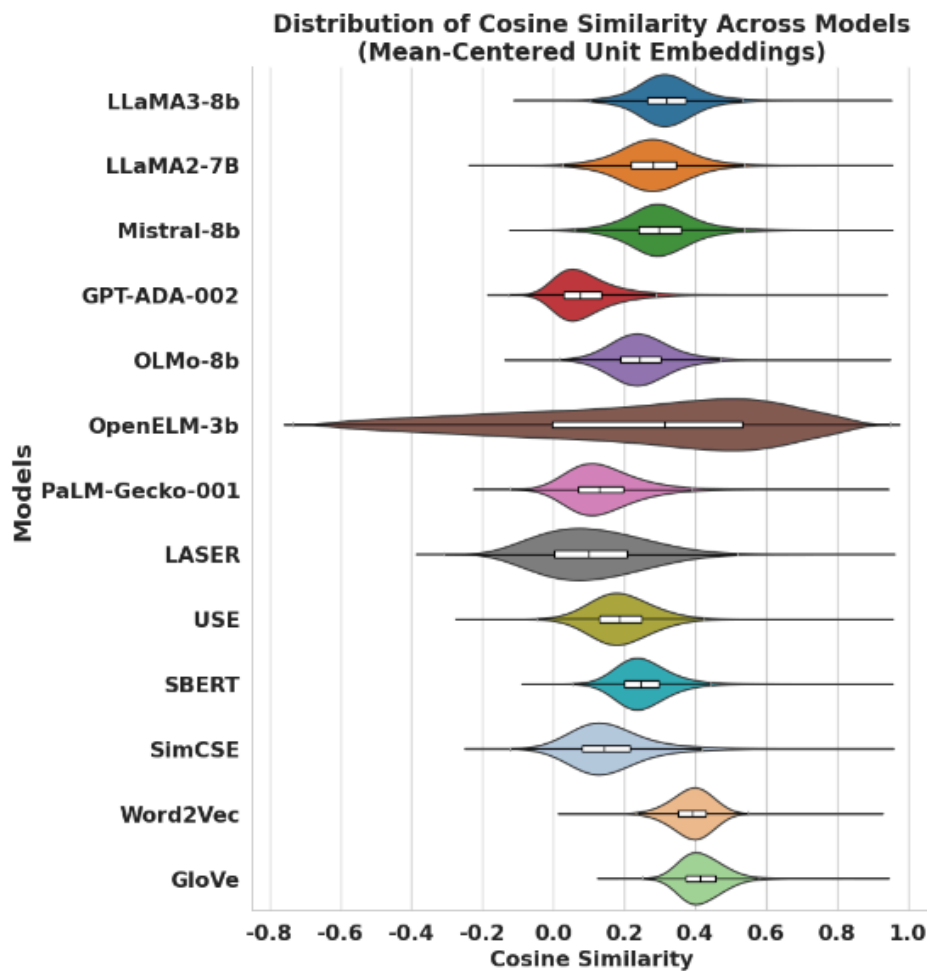


Figure 1: The Mean-Centered Embedding distribution of cosine similarities between all pairs of words.

- **New downstream-task showcase (meta-reviewer)** — introduces **Table 3 and 5** showed the “**Evaluation of LLMs on downstream task.**” The table reports performance on two categories; Reasoning focused and Embedding focused. We discussed the outcome and showed how we use our criteria as a model performance predictor on downstream task.

Model	Sbert	SimCSE	LlaMA2-7b	LlaMA3-8b	Mistral-7b	OLMo-7b	OpenELM-3B
ARC-c	-	-	54.2	<b>79.3</b>	<b>78.6</b>	48.5	35.58
ARC-e	-	-	84.0	<b>92.4</b>	<b>90.8</b>	65.4	59.89
BoolQ	-	-	86.1	<b>87.5</b>	<b>89.3</b>	74.4	67.4
HellaSwag	-	-	78.9	<b>81.8</b>	<b>83.0</b>	76.4	72.44
MMLU	-	-	46.2	<b>66.6</b>	<b>64.0</b>	40.5	26.76
PIQA	-	-	57.8	77.2	<b>80.6</b>	<b>78.4</b>	78.24
SIQA	-	-	77.5	81.6	<b>82.8</b>	78.5	<b>92.7</b>
WinoGrande	-	-	71.7	<b>76.2</b>	<b>77.9</b>	67.9	65.51
Clustering	42.35	29.04	45.24	<b>46.45</b>	<b>54.93</b>	32.0	18.71
Pair classification	82.37	70.33	<b>88.03</b>	87.8	<b>88.59</b>	49.32	56.71
Reranking	<b>58.04</b>	46.47	57.38	<b>59.68</b>	50.15	33.91	37.0
STS	78.9	74.33	<b>83.73</b>	83.58	<b>84.77</b>	27.04	38.31
Summarization	30.81	<b>31.15</b>	28.49	30.94	<b>36.32</b>	20.83	18.71

Table 5: Model Evaluation Results Across Various Tasks. **Blue** is top scorer and **black** is second best.

- A short “Discussions & Final Words” subsection distils practical take-aways for model selection in low-resource settings and outlines plans to release task scripts.
- **Minor corrections** — replace blue finding boxes with simple blue text.