Pre-trained Transformer Models for Standard-to-Standard Alignment Studies

Hye-Jeong Choi^a, Reese Butterfuss^b, Meng Fan^a, and Emily Dickinson^a

^a HumRRO ^b Certiverse

Abstract

The current study evaluated the accuracy of five pre-trained large language models (LLMs) in matching human judgment for standard-to-standard alignment study. Results demonstrated comparable performance across LLMs despite differences in scale and computational demands. Additionally, incorporating domain labels as auxiliary information did 10 not enhance LLMs performance. These 11 findings provide initial evidence for the 12 viability of open-source LLMs to facilitate 13 alignment study and offer insights into the 14 utility of auxiliary information. 15

Introduction

17 Large language models (LLMs) are increasingly in educational psychological 18 used and 19 measurement activities. Their 20 sophistication and ability to represent deep, 60 include examples to provide clarity and context. It 21 contextual semantics make them viable tools to 61 is also not uncommon to have the exact same 22 support subject matter experts (SMEs) in 62 standard 23 reviewing large volumes of text-based context, 63 domains. Understanding how 24 such as educational standards (e.g. Butterfuss & 64 information influences LLMs performance is 25 Doran, 2025; Kim et al., 2023; Kusumawardani & ²⁶ Alfarozi, 2023; Zhou & Ostrow, 2022). However, 27 little guidance exists on the effective use of LLMs 28 in such contexts. Our goal was to compare popular, 67 2 29 pretrained LLMs in a common measurement 30 context (i.e., standard-to-standard alignment) to 33 require extensive review of large bodies of text.

40 to review two sets of standards and determine 41 alignment such that each standard in one set is 42 evaluated against the standards in the second set 43 until any or all standards that capture the same 44 meaning are identified. It is a time-consuming 45 process because it requires evaluation of 46 potentially thousands of possible pairs of content 47 standards. Recently, the potential for NLP and 48 LLMs as a supporting tool in this process has been 49 presented (e.g., Butterfuss & Doran, 2024; Zhou & 50 Ostrow, 2022), but there is a lack of work that 51 provides guidance on which LLMs to choose for 52 such tasks.

This study aimed to address two research 54 questions: (1) how do five popular pre-trained 55 transformer models compare in standards-to-56 standards alignment studies? and (2) does auxiliary 57 information (e.g., domain label) impact LLMs 58 performance? Educational standards, typically evolving 59 presented as brief, abstract statements, often appear under different such auxiliary 65 crucial for developing more effective automated 66 alignment tools.

Methods

68 **2.1** Data

31 provide initial evidence on which LLMs may be 69 The alignment study dataset used for the current 32 particularly useful for measurement tasks that 70 study consisted of individual standards from 33 71 states and aligned each state standard to the Alignment is a critical aspect of validity 72 corresponding the National Assessment of 35 evidence for any assessment (AERA, APA, 73 Educational Progress (NAEP) standard for grades 36 NCME, 2014). Standards-to-standards alignment 74 4, 8, and 12 for science. Each standard was 37 is a process to examine how well two distinct sets 75 classified into one of three domains: life science 38 of content standards target the same content 76 (LS), physical science (PS), and earth & space 39 (Neidorf et al., 2016). In general, it requires SMEs 77 science (ES). The number of potential pairs ranged

79 of standards represented within each state varied. 130 of meaning contribute to cosine similarity values. 80 In the original work, SMEs judged each possible 131 Due to this variability, we extracted embeddings 81 pair of standards as aligned, partially aligned, or 132 for each standard using five different popular 82 not aligned. Thus, we used the SME decision as 133 LLMs: 83 "ground truth" for evaluating the LLMs. More 84 details about the dataset and original alignment 134 85 study can be found in the published report 135 86 (Dickinson et al., 2021).

Pre-trained transformer models

88 We accessed the LLMs via the Hugging Face 89 Transformers library, a popular open-source library 90 that provides a simple and consistent way to use 91 pre-trained models for various NLP tasks. As of 92 2025, the Hub hosts over 50,000 models, many of 93 which are based on Transformer architectures.

The LLMs transform each content standard into 145 95 an embedding, or numeric representation of the 146 96 meaning of the text. Once every standard is 147 97 transformed into an embedding, then the relations 98 among the embeddings can be evaluated using 148 99 cosine similarity. Cosine similarity is a metric used 149 100 to measure how similar two vectors are irrespective 150 101 of their magnitude. It calculates the cosine of the 151 102 angle between two vectors, determining whether 152 103 they point in roughly the same direction. 153 104 Commonly used in text analysis, recommendation 154 105 systems, and information retrieval. While it 155 106 behaves similarly to a correlation in some contexts, 156 107 cosine similarity specifically only measures 108 directional similarity, not linear correlation or 109 magnitude.

In this study we used the cosine similarity between every possible pair of standards that can 112 be made from the two sets. Doing so allows us to gauge which standard pairs are more similar than 163 114 others. Critically, standard pairs that share high 115 semantic overlap (i.e., large cosine similarity 116 values) are more likely to be aligned than standard 165 117 pairs that share little semantic overlap (Butterfuss 166 118 & Doran, 2025).

To calculate cosine similarity, we utilized five 168 120 LLMs which are widely used, including all- 169 distilroberta-v1, all-MiniLM-L6-v2, multi-qa- 170 122 MiniLM-L6-cos-v1, all-mpnet-base-v2, and gtr-t5-123 large. All of these are sentence embedding models 171 2.3 124 that can be used to calculate cosine similarity 172 We employed three threshold setting approaches to between texts. The mathematical formula for 173 pair state and NAEP standards: cosine similarity 126 calculating cosine similarity remains the same 174 value, percentile, and rank order. First, we used 127 across all these models. However, LLMs vary in 175 predetermined cosine similarity values: if the

78 from approximately 18,000 to 60,000. The number 129 emphasize, and thus LLMs differ in which aspects

- all distilroberta v1 (DistilRoBERTa-v1). It is a distilled version of the RoBERTa (Liu et al., 2019) model to cover a wide range of topics and styles. It is a smaller, more efficient model that's designed to be faster and more computationally efficient.
- all MiniLM L6 v2 (MiniLM-L6-v2). MiniLM is designed for efficiency and smaller size. It's useful for text classification, sentiment analysis, or question answering. They are particularly useful for deployment in resourceconstrained environments, such as mobile devices or edge computing platforms (Wang et al., 2020).
- multi qa MiniLM L6 cos v1 (MultiQA-MiniLM-L6). It is a variant of the MiniLM model that is designed for multi-question answering tasks, such as answering multiple questions about a given text passage, identifying relevant passages or sentences that answer multiple questions, and generating answers to multiple questions based on a given text passage.
- all mpnet base v2 (MPNet-Base-v2). A model known for its efficiency and performance on a wide range of NLP tasks, including text classification, sentiment analysis, question answering, and more. It's particularly useful when you need a model that can handle long-range dependencies and contextual relationships in text data.
- gtr t5 large (GTR-T5-Large). It is known as a powerful language model. It can be used text generation and summarization, question and reading comprehension, answering sentiment analysis and opinion mining, and language translation and machine translation.

Three approaches to set a threshold

128 the specific linguistic features their embeddings 176 cosine similarity of two-paired standards was

138

177 greater than the predetermined cosine similarity 178 value, we classified the state-to-NAEP standard 179 pair as aligned. We used three different values as 180 the predetermined value (i.e., 0.4, 0.5, 0.6).

Second, we used a percentile to set the cut score 182 cosine similarity value. As mentioned in the 183 previous section, cosine similarity is a measure of 184 direction but not magnitude. Using percentile can 185 resolve potential scaling issues across LLMs. We used three percentiles (i.e., 70, 80, 90) to obtain the threshold of cosine similarity to pair standards.

Finally, we utilized a rank-order approach to 189 classify aligned standard pairs: if the cosine 190 similarity of standard pairs was within the 191 predetermined top n highest cosine similarity, we 192 classified those standards as aligned. We used top 193 3, 5 and 10. After we classified each pair as either aligned or not aligned based on those criteria, we 227 and F1. Note that we used three different methods 195 evaluated those results with SMEs judgment.

LLMs performance was evaluated using overall 197 accuracy, recall and F1 metrics. Overall accuracy 198 (either hit rate or precision) refers to the probability 199 of capturing the true matches (according to human 200 judgment) within condition. Recall measures the 201 proportion of actual positive instances that were 202 correctly identified by the model. It is a metric used 203 to evaluate the completeness of a classification 204 model's positive predictions. The F1-score is a 205 metric that combines precision and recall into a 206 single value, providing a balance between these 207 two sometimes competing metrics. Precision 208 measures the proportion of correctly identified 209 positive instances among all instances that the 210 model predicted as positive. It's particularly useful 211 when you need a single measurement to evaluate a 212 classification model's performance.

213 3 Results

214 3.1 Comparison of five pre-trained transformer models

216 Table 1 presents descriptive statistics for the 217 cosine-similarity values generated by each LLM 218 for each grade. Overall, the correlations between LLMs were high (higher than .76). In both 235 the true pairs with respect to F1 and accuracy. 220 conditions, the results indicated that the models 221 produced cosine-similarity values that were scaled 222 slightly differently. In particular, means of GTR-T5-Large were higher and standard deviations 239 T5-Large performed differently from other four 224 were smaller than other LLMs.

SMEs with respect to the overall accuracy, recall, 242 four other models.

Table 1 Descriptive statistics of cosine similarity by LLMs for each grade

	MPNet	Distil	Mini	GTR-	Multi				
	MIFINEL	RoBERT	LM	T5	QA				
Grade 4	(N=18,744))							
Mean	.22	.18	.19	.55	.17				
STD	.14	.14	.14	.07	.14				
Grade 8	Grade 8 (N=55,857)								
Mean	.22	.18	.20	.56	.18				
STD	.13	.13	.14	.06	.14				
Grade 12 (N=59,829)									
Mean	.20	.16	.18	.55	.16				
STD	.13	.13	.14	.06	.13				
Note. MF	Note. MPNet = MPNet-Base-v2; DistilRoBERT =								
DistilRoB	DistilRoBERTa-v1; MiniLM = MiniLM-L6-v2; GTR-T5 =								

GTR-T5-Large; MultiQA = MultiQA-MiniLM-L6

Table 2 Overall comparison LLMs results with SMEs rating under different conditions

Model	Stats	Cut Value		Percentile		Rank	
Model	Stats	M	STD	M	STD	M	STD
MPNet-	F1	.24	.21	.29	.13	.35	.13
Base-v2	Recall	.22	.26	.91	.11	.94	.05
base-v2	Accuracy	.98	.01	.86	.10	.89	.07
Distil	F1	.20	.21	.29	.13	.35	.13
RoBERTa	Recall	.16	.21	.90	.11	.94	.06
-v1	Accuracy	.98	.01	.86	.10	.90	.07
MiniLM	F1	.24	.21	.29	.13	.35	.13
-L6-v2	Recall	.20	.24	.90	.12	.94	.06
-L0-V2	Accuracy	.98	.01	.86	.10	.90	.07
CED TO	F1	.19	.19	.29	.14	.35	.14
GTR-T5	Recall	.53	.44	.91	.10	.95	.04
-Large	Accuracy	.79	.30	.86	.10	.89	.07
MultiQA	F1	.21	.19	.27	.12	.34	.13
-MiniLM	Recall	.16	.20	.87	.13	.92	.06
-L6	Accuracy	.98	.00	.85	.10	.90	.07

228 to classify pairs of standards: cosine similarity value, percentile, and rank order. First, notably, the 230 free, open-source models fared nearly as well as the 231 costlier, more computationally intensive model 232 (GTR-T5-Large). Overall, the correlations 233 between LLMs were high (higher than .76). 234 Second, all LLMs performed similarly in capturing 236 However, recall indicates percentile and rank 237 performed much better to identify the true pairs for 238 all five models. When cut score was used, GTR-240 models. That was because cosine similarity from Table 2 summarizes comparison of LLMs to 241 GTR-T5-Large tended different and larger than

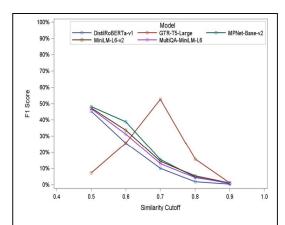


Figure 1 F1 score trend of LLMs across cosine similarity points (Grade 4)

Table 4 Comparison of LLMs at different ranks (Grade 4, N=18,744)

Model	Rank	Accuracy	Recall	F1
MPNet-Base-v2	3	.95	.87	.50
MPNet-Base-v2	5	.90	.95	.34
MPNet-Base-v2	10	.76	.99	.19
DistilRoBERTa-v1	3	.95	.83	.49
DistilRoBERTa-v1	5	.90	.93	.34
DistilRoBERTa-v1	10	.76	.99	.19
MiniLM-L6-v2	3	.95	.85	.49
MiniLM-L6-v2	5	.90	.94	.34
MiniLM-L6-v2	10	.76	1.00	.19
GTR-T5-Large	3	.95	.91	.51
GTR-T5-Large	5	.90	.97	.35
GTR-T5-Large	10	.76	1.00	.19
MultiQA-MiniLM-L6	3	.95	.83	.49
MultiQA-MiniLM-L6	5	.90	.91	.34
MultiQA-MiniLM-L6	10	.76	.99	.19

best when the cosine similarity was set at .50. Table 281 ranging from .75 to .92. ²⁵¹ 4 presents the comparison of LLMs with different ²⁸²

Table 3 Comparison of LLMs at different percentiles (Grade 4, N=18,744)

Model	%ile	Accuracy	Recall	F1
MPNet-Base-v2	70	.73	1.00	.17
MPNet-Base-v2	80	.83	.97	.24
MPNet-Base-v2	90	.92	.83	.37
DistilRoBERTa-v1	70	.73	.99	.17
DistilRoBERTa-v1	80	.83	.94	.24
DistilRoBERTa-v1	90	.92	.83	.37
MiniLM-L6-v2	70	.73	.99	.17
MiniLM-L6-v2	80	.83	.95	.24
MiniLM-L6-v2	90	.92	.84	.37
GTR-T5-Large	70	.73	.99	.17
GTR-T5-Large	80	.83	.97	.24
GTR-T5-Large	90	.92	.89	.40
MultiQA-MiniLM-L6	70	.73	.98	.17
MultiQA-MiniLM-L6	80	.83	.95	.24
MultiQA-MiniLM-L6	90	.92	.79	.35

259 percentile. Again, LLMs performed similarly: the 260 higher percentile, the better in capturing the true 261 pairs with respect to accuracy whereas the lower 262 percentile, the better in terms of recall. Both 263 accuracy and recall indicate LLMs with either 70 264 percentile or rank order 10 well captured the true 265 pairs. In other words, the NAEP standards, with 266 which the state standard is aligned, appear among 267 the top ten or even top five pairs rank-ordered or 70 268 or 80 percentiles by cosine similarity.

Effects of domain information effect on 269 3.2 cosine similarity

271 Table 7 presents cosine similarity distributions 272 when domain labels were added for each grade. 273 Note that the N counts for all grades were slightly Next, we will present the results focusing on 274 larger than the N counts in Table 1. This was 244 Grade 4 as the results with other grades were 275 because the same standard was assigned into 245 similar. Figure 1 depicts F1 across several cosine 276 different domains. The descriptives were similar 246 similarity points from .50 to .90 for all five 277 with ones in Table 1; however, those values were 247 language models for Grade 4. GTR-T5-Large 278 slightly lower. Also, the correlations between performed best when the cosine similarity was set 279 cosine similarity measures for standard pairs with 249 at .70 whereas four languages models performed 280 domain were similar with ones without domain,

Next, we compare how LLMs performed to 252 ranks. As expected, LLMs captured the true pair 283 capture the true pairs compared with cosine 253 with lower ranks (rank=3). However, recall 284 similarity without domain. Again, we present the 254 indicates LLMs captured the true pair with rank=10. 285 results for Grade 4 as the results with other grades 255 In other words, the NAEP standards, with which 286 were similar. Figure 2 depicts F1 scores across 256 the state standard is aligned, appear among the top 287 several cosine similarity points from .50 to .90 for 257 ten pairs rank-ordered by cosine similarity. Table 3 288 all five language models for Grade 4. The results 258 shows the comparison of LLMs with different 289 show a similar pattern with Figure 1; with respect

Table 7 Descriptive statistics of cosine similarity with domain by LLMs for each grade

	MPNet	Distil	Mini	GTR	Multi			
	IVII INCL	RoBERT	LM	-T5	QA			
Grade 4 (N=18,846)								
Mean	.21	.17	.17	.55	.15			
STD	.13	.13	.14	.06	.13			
Grade 8	Grade 8 (N=56,932)							
Mean	.22	.17	.19	.56	.16			
STD	.12	.13	.13	.06	.13			
Grade 12 (N=59,927)								
Mean	.20	.15	.17	.55	.14			
STD	.13	.13	.13	.06	.13			
DistilRoE	Note. MPNet = MPNet-Base-v2; DistilRoBERT = DistilRoBERTa-v1; MiniLM = MiniLM-L6-v2; GTR-T5 = GTR-T5-Large; MultiQA = MultiQA-MiniLM-L6							

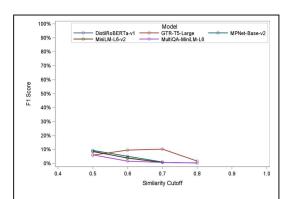


Figure 2 F1 Score by cosine similarity cut for each language model with domain (Grade 4)

to F1, GTR-T5-Large outperformed and performed the best when the cosine similarity was set at .70. Overall, however, all the five LLMs performed slightly worse with domain labels.

Table 6 and 7 show that domain information improved accuracy but reduced recall and F1. Adding domain labels did not enhance overall model performance.

98 4 Summary and discussion

The results of the current study indicated that scaling differences among LLMs in raw cosine similarity values meant that using a raw cosine value threshold may not be feasible, particularly when comparing multiple LLMs. Overall, when percentile or rank order was used, the results suggest that the five LLMs performed comparably with respect to accuracy for standards-to-standards alignment of science content standards. Specifically, the models were generally 90%

Table 6 Comparison of LLMs with domain at different ranks (Grade 4, N=18,876)

Model	Rank	Accuracy	Recall	F1
MPNet-Base-v2	3	.97	.23	.28
MPNet-Base-v2	5	.95	.31	.25
MPNet-Base-v2	10	.88	.44	.17
DistilRoBERTa-v1	3	.97	.21	.26
DistilRoBERTa-v1	5	.95	.28	.24
DistilRoBERTa-v1	10	.88	.43	.17
MiniLM-L6-v2	3	.97	.23	.27
MiniLM-L6-v2	5	.95	.30	.25
MiniLM-L6-v2	10	.88	.44	.17
GTR-T5-Large	3	.97	.26	.30
GTR-T5-Large	5	.95	.33	.25
GTR-T5-Large	10	.87	.47	.17
MultiQA-MiniLM-L6	3	.97	.20	.25
MultiQA-MiniLM-L6	5	.95	.29	.24
MultiQA-MiniLM-L6	10	.88	.45	.17

Table 5 Comparison of LLMs with domain at different percentiles (Grade 4, N=18,876)

Model	%ile	Accuracy	Recall	F1
MPNet-Base-v2	70	.70	.44	.08
MPNet-Base-v2	80	.79	.35	.09
MPNet-Base-v2	90	.88	.23	.10
DistilRoBERTa-v1	70	.70	.44	.08
DistilRoBERTa-v1	80	.79	.33	.08
DistilRoBERTa-v1	90	.88	.23	.10
MiniLM-L6-v2	70	.70	.45	.08
MiniLM-L6-v2	80	.79	.35	.09
MiniLM-L6-v2	90	.88	.23	.10
GTR-T5-Large	70	.70	.46	.08
GTR-T5-Large	80	.79	.38	.09
GTR-T5-Large	90	.89	.25	.11
MultiQA-MiniLM-L6	70	.70	.44	.08
MultiQA-MiniLM-L6	80	.79	.33	.08
MultiQA-MiniLM-L6	90	.88	.21	.09

309 accurate at capturing the "true matches" according 310 to human judgment above the 90 percentile or 311 within the top five highest-cosine pairs. Put another 312 way, for a given state standard, the SME-aligned 313 NAEP standard had appeared either the 90 314 percentile or among the top five cosine similarity 315 pairs.

Using Grade 4 from our real-world alignment study as an example, the current method would reduce the number of pairs that SMEs must compare from nearly 18,000 pairs to around 2,840 pairs. Moreover, the current findings suggest that 321 any of the popular, open-source LLMs we 373 322 compared may yield such benefits. Thus, for 374 323 contexts similar to those in the current work, 375 324 researchers and practitioners may be well-suited to 376 Kim, D. E., Hong, C., & Kim, W. H. (2023, July). 325 choose any of the models we evaluated given their 377 326 comparable performance. Also, the current 378 327 findings highlight a potentially enormous 379 328 efficiency increase by dramatically reducing the 380 329 number of pairs SMEs must consider via 381 Kusumawardani, A., & Alfarozi, M. (2023). Exploring 330 economical LLMs and a relatively simple 382 331 percentile or rank-order approach using cosine- 383

334 did not improve LLMs performance to capture the 386 335 true matches. Subsequent work in this area is 387 336 needed to examine the added benefits of including more contextual information for the LLMs when 389 Wang, W., Wei, F., Dong, L., Bao, H., Yang, N., & 338 extracting embeddings for each standard (i.e., 390 339 content domain descriptions), as well as the 391 340 conditions under which it is useful to include or 392 341 omit accessory information that some content 342 standards include, such as exemplary information 394 Zhou, Z., & Ostrow, K. S. (2022, June). Transformer-343 or explanatory information. The results did not 395 344 show LLMs performance were substantially 345 impacted by grade.

Critically, the current LLM approach does not 347 replace humans in making alignment decisions. 348 Instead, the method provides a simple, economical 349 way to support SMEs in making alignment 350 decisions more efficiently by leveraging an 351 organizational structure based on semantic 352 similarity and constraining the number of viable 353 pairs that must be considered. Overall, the current 354 study represents a judicious, human-centered use 355 of AI in a laborious routine measurement task.

357 References

Educational Research 358 American Association. American Psychological Association, & National 359 Council on Measurement in Education. (2014). 360 Standards for educational and psychological 361 testing. American Educational Research 362 Association.

364 Butterfuss, R., & Doran, E. (2025). Advancing measurement with large language models: 365 Implications for content review. Journal of 366 Educational Measurement, 62(1), 45-63.

Dickinson, E. R., Gribben, M., Schultz, S. R., Spratto, E., & Woods, A. (2021). Comparative analysis of 369 the NAEP science framework and state science 370 standards (Tech. Rep.). National Assessment 371 Board. Retrieved Governing 372

https://www.nagb.gov/content/dam/nagb/en/docum ents/publications/frameworks/science/NAEP-Science-Standards-Review-Final-Report-508.pdf

Efficient Transformer-based Knowledge Tracing for a Personalized Language Education Application. In Proceedings of the Tenth ACM Conference on Learning@ Scale (pp. 336-340).

AI-driven tools for curriculum mapping. International Journal of Educational Technology, 18(2), 110-125.

Unfortunately, adding "domain" to the standard 385 Neidorf, T. S., Binkley, M., Galia, J., & Stephens, M. (2016). Assessing alignment among curriculum, instruction, and assessments: Methodologies and findings. OECD Publishing.

> Zhou, M. (2020). Minilm: Deep self-attention distillation for task-agnostic compression of pretrained transformers. Advances in neural information processing systems, 33, 5776-5788.

Based Automated Content-Standards Alignment: A Pilot Study. In International Conference on Human-Computer Interaction (pp. 525-542). Cham: Springer Nature Switzerland.