

MUSESCORER: Idea Originality Scoring At Scale

Ali Sarosh Bangash*, Krish Veera*, Ishfat Abrar Islam, Raiyan Abdul Baten†

Bellini College of Artificial Intelligence, Cybersecurity, and Computing,
University of South Florida, USA

{alibangash, krishv, ishfatabrar, rbaten}@usf.edu

*Equal contributions †Correspondence: rbaten@usf.edu

Abstract

An objective, face-valid method for scoring idea originality is to measure each idea’s statistical infrequency within a population—an approach long used in creativity research. Yet, computing these frequencies requires manually bucketing idea rephrasings, a process that is subjective, labor-intensive, error-prone, and brittle at scale. We introduce MUSESCORER, a fully automated, psychometrically validated system for frequency-based originality scoring. MUSESCORER integrates a Large Language Model (LLM) with externally orchestrated retrieval: given a new idea, it retrieves semantically similar prior idea-buckets and zero-shot prompts the LLM to judge whether the idea fits an existing bucket or forms a new one. These buckets enable frequency-based originality scoring without human annotation. Across five datasets ($N_{\text{participants}}=1143$, $n_{\text{ideas}}=16,294$), MUSESCORER matches human annotators in idea clustering structure (AMI = 0.59) and participant-level scoring ($r = 0.89$), while demonstrating strong convergent and external validity. The system enables scalable, intent-sensitive, and human-aligned originality assessment for creativity research.

1 Introduction

Assessing creativity at scale remains a central challenge in cognitive science and computational linguistics. Creativity scoring broadly considers two complementary dimensions: the *intrinsic* qualities of ideas (e.g., creative ideas are semantically flexible or diverse) and their *extrinsic* statistical infrequency within a population (i.e., original ideas do not appear very often) (Beketayev and Runco, 2016; Runco and Jaeger, 2012). Recent methods have scaled intrinsic assessments using unsupervised, semi-supervised, and supervised approaches (Organisciak et al., 2023; Beaty and Johnson, 2021; Organisciak and Dumas, 2020). In contrast, frequency-based originality scoring still

depends on manual tabulation of response occurrences (Reiter-Palmon et al., 2019). This requires grouping rephrasings of the same idea into buckets (e.g., ‘hold papers down’ and ‘use as a paperweight’ for a brick), a process that is subjective, fatigue-intensive, and error-prone as annotators track an ever-expanding set of buckets (Acar and Runco, 2014; Baten et al., 2020, 2021, 2022; Buczak et al., 2023). Moreover, the field lacks consensus on what counts as an ‘infrequent’ idea, leaving frequency-based scoring with limited psychometric validation.

We introduce MUSESCORER, a fully automated, psychometrically validated system for frequency-based originality scoring—bringing us closer to comprehensive automated creativity assessment. Bucketing different phrasings of the same idea together is computationally non-trivial: (i) semantic similarity alone is insufficient to distinguish rephrasings from distinct intents, (ii) traditional clustering algorithms struggle with singleton and low-frequency ideas that are critical for infrequency scoring, (iii) real-world idea datasets follow fat-tailed bucket size distributions, defying uniform or Gaussian assumptions, and (iv) bucket count grows as new ideas arrive, rendering ineffective text labeling tools that require label sets apriori. MUSESCORER addresses these challenges via an LLM-as-a-judge framework with externally orchestrated retrieval: for each new idea, it retrieves from its database semantically similar prior buckets as candidates and zero-shot prompts the LLM to decide whether the idea fits an existing bucket or forms a new one. Unlike conventional clustering, this approach replicates the granularity of human bucketing in both structure and resolution.

Our work also contributes to the creativity literature in two ways. *First*, we establish rigorous psychometric validity for frequency-based originality scoring, showing high agreement with human annotations and strong correlations with relevant cognitive traits. In doing so, we elucidate how

‘infrequency’ can be reliably operationalized. *Second*, we release an automated, interpretable scoring pipeline deployable across diverse open-ended ideation tasks, enabling creativity research at scale¹. More broadly, MUSESCORER demonstrates how advanced NLP methods can address long-standing annotation challenges, providing validated tools that adjacent disciplines can adopt with confidence.

2 Related Work

2.1 Computational Assessment of Creativity

Creativity assessment has long relied on divergent thinking tasks like the Alternate Uses Test (AUT), where participants list novel uses for everyday objects (Guilford, 1967). Response sets are then scored for fluency (idea count), flexibility (distinct semantic category count), originality (statistical infrequency relative to a population), novelty (Likert-scale ratings by human judges), and other metrics (Dumas and Dunbar, 2014; Runco and Mraz, 1992).

Several computational methods have been proposed to automate these scores. Unsupervised approaches estimate (i) flexibility by measuring the semantic diversity of an idea set (Snyder et al., 2004; Bossomaier et al., 2009), and (ii) human judges’ novelty ratings by computing an idea’s semantic distance from the task prompt (Beaty and Johnson, 2021; Dumas et al., 2021; Acar and Runco, 2014). Hybrid and supervised methods directly predict novelty ratings using regression and clustering-based pipelines (Organisciak et al., 2023; Stevenson et al., 2020). However, these methods face generalizability issues, with models trained on one task or dataset often performing poorly on another (Buczak et al., 2023). More recently, studies have explored LLMs for zero-shot creativity scoring, but results show little to no correlation with human labels (Chakrabarty et al., 2024).

Unfortunately, computational approaches for scoring ideas by *statistical rarity* remain underexplored. Recent work has addressed related challenges—for example, Lu et al. (2024) contrast AI model outputs against an extrinsic human-generated text corpora using n -gram overlap and Word Mover’s Distance (WMD) to probe the origin of AI creativity. However, using purely lexical (n -gram) or embedding-based matching methods (WMD) for operationalizing social comparison

holds the risk of conflating distinct intents or over-separating true rephrasings, as they privilege surface similarity over conceptual intent (Olson et al., 2021). Our approach addresses these limitations by incorporating a zero-shot LLM in the annotation loop to make the *subjective, intent-sensitive* bucketing judgments, which, in turn, enables frequency-based originality scoring at scale.

2.2 Text Clustering and Annotation

LLMs have recently been explored for zero- and few-shot text clustering and annotation (Xiao et al., 2023b). Deductive clustering methods prompt LLMs to partition a given set of texts, generating categories or groupings directly (Viswanathan et al., 2024; Chew et al., 2023). However, most LLM-based deductive clustering methods assume all clusters are discoverable upfront and perform poorly when the concept space evolves. Inductive annotation methods, on the other hand, present labeled exemplars to classify new instances incrementally (Dai et al., 2023). While current approaches show promise on well-bounded tasks like topic labeling or thematic analysis, it remains unclear how best to navigate fat-tailed distributed datasets, where cluster (i.e., idea bucket) counts grow without bound as data scale increases.

2.3 LLM-as-a-Judge

The LLM-as-a-judge paradigm has emerged as a powerful approach for evaluating, ranking, and filtering outputs across tasks such as summarization, translation, alignment, and reasoning in NLP (Li et al., 2024a; Liang et al., 2023; Zhao et al., 2024). Unlike earlier evaluation approaches (Papineni et al., 2002; Zhang et al., 2019), judge LLMs can assess contextual fit, intent, and subtle distinctions between candidates, using pointwise, pairwise, or listwise formats (Gao et al., 2023; Shen et al., 2024).

Our task combines listwise judgment with decision-making: the LLM determines whether a new idea matches any retrieved exemplar or forms a new semantic bucket, akin to selection-based judgment (Li et al., 2024b; Yao et al., 2023). We adopt a modular retrieval-based framework (Lewis et al., 2020; Izacard and Grave, 2020), where retrieval and codebook management are handled externally (Khandelwal et al., 2020), leaving the stateless LLM to focus on subjective bucketing decisions. This separation improves stability and preserves the interpretability and psychometric auditability critical for creativity research.

¹<https://github.com/cssai-research/MuseScorer>

Dataset	# Participants	# Tasks	# Ideas	# Judges
socialmuse24 (Baten et al., 2024)	109	5	5703	2
beaty18 (Beaty et al., 2018)	171	2	2917	4
silvia17 (Silvia et al., 2017)	141	2	2355	3
beaty12 (Beaty and Silvia, 2012)	133	1	1807	3
mohr16 (Hofelich Mohr et al., 2016)	305 + 284	1 + 1	1930 + 1582	4

Table 1: Dataset summary. Each participant did one task in mohr16. In other datasets, all participants did all tasks.

3 Dataset Acquisition

We use five Alternative Uses Test datasets (Table 1). Each dataset includes one or more tasks where participants generate alternative uses for an everyday object. Responses range from short phrases to full sentences (e.g., for a shoe: “We can use a shoe as a hamster bed” or “As a doorstep”).

3.1 Primary Dataset: socialmuse24

We use the socialmuse24 dataset to establish criterion validity (Baten et al., 2024). Two trained research assistants (H1 and H2) independently bucketed rephrased ideas within each task. The annotators saw the ideas in a random order. They followed the coding rules of Bouchard and Hare (Bouchard Jr and Hare, 1970) and Guilford’s scoring key (Guilford et al., 1978). Each idea thus has two categorical bucket IDs—one per annotator—serving as ground truth for evaluating our method. In the original study, these human-assigned buckets were used to estimate originality without automation; our goal is to replicate this process computationally. The dataset also includes flexibility-based *Creativity Quotient* scores (Snyder et al., 2004; Bossomaier et al., 2009), which we use to assess convergent validity.

3.2 Secondary Datasets

We draw on four publicly available AUT datasets to assess convergent and external validity (Organis- ciak et al., 2023; Beaty and Johnson, 2021). Unlike socialmuse24, these datasets lack human bucketing annotations and thus cannot be used for our primary originality-scoring goal.

The beaty18 dataset (Beaty et al., 2018) includes four judges’ *Creative Quality* ratings on a 1–5 scale, along with measures of: (i) *Creative Metaphor Generation*, where participants produced novel metaphors for two open-ended prompts, each rated 1–5 by four judges (Beaty and Silvia, 2013); (ii) *Big Five Personality*, assessing the participants’ neuroticism, extraversion, openness to experience,

agreeableness, and conscientiousness via standard questionnaires (McCrae et al., 2005); (iii) *Fluid Intelligence*, through sequence-completion tasks with images (Cattell and Cattell, 1960), letters (Ekstrom et al., 1976), and numbers (Thurstone, 1938); and (iv) *Creative Self-concept*, via self-efficacy and self-identity questionnaires (Karwowski, 2014).

The silvia17 dataset (Silvia et al., 2017) provides three judges’ *Creative Quality* ratings, as well as openness-to-experience *Personality* scores (Lee and Ashton, 2004). The beaty12 dataset (Beaty and Silvia, 2012) includes three judges’ *Creative Quality* ratings, plus *Big Five Personality*, *Creative Metaphor Generation*, and *Fluid Intelligence* measures, paralleling beaty18.

Finally, mohr16 (Hofelich Mohr et al., 2016) contains four judges’ ratings of idea *Originality* and *Flexibility*. Here, originality captured the uncommonness, remoteness, and cleverness of responses (1–5 scale) (Silvia et al., 2008), while flexibility was defined as the number of categories in each participant’s responses, averaged across judges.

4 Task Description

4.1 Problem Formulation

Let $\mathcal{P} = \{p_1, p_2, \dots, p_N\}$ denote a corpus of N participants, each completing T ideation tasks. For each task $t \in \{1, \dots, T\}$, participant p_i produces a variable-length set of $n_{i,t}$ free-form textual responses, denoted $\mathcal{I}_{i,t} = \{x_{i,t}^{(1)}, \dots, x_{i,t}^{(n_{i,t})}\}$.

Let $\mathcal{X}_t = \bigcup_{i=1}^N \mathcal{I}_{i,t}$ denote the full idea set for task t . The goal is to induce a task-specific partition $\mathcal{B}_t = \{B_{t,1}, \dots, B_{t,K_t}\}$ over \mathcal{X}_t , where each ‘bucket’ $B_{t,k} \subseteq \mathcal{X}_t$ contains semantically equivalent ideas expressing the same underlying concept.

Let $k(x)$ denote the index of the bucket to which idea $x \in \mathcal{X}_t$ is assigned. We define $m_{t,k}$ as the number of distinct participants contributing at least one idea to bucket $B_{t,k}$. Importantly, the bucketing is performed *within* each task and *across* participants, and no bucket identity is shared across tasks.

4.2 Originality Metrics

We explore 4 frequency-based originality metrics:

(i) **rarity**: Each idea bucket $B_{t,k}$ is scored as $(1 - \frac{m_{t,k}}{N})$, reflecting the bucket’s relative infrequency in the sample (Forthmann et al., 2020, 2017). A participant’s unnormalized rarity score is the sum across their ideas: $R_{i,t}^{\text{rarity}} = \sum_{x \in \mathcal{I}_{i,t}} (1 - \frac{m_{t,k(x)}}{N})$.

(ii) **shapley**: Each bucket $B_{t,k}$ is scored as $\frac{1}{m_{t,k}}$, making a bucket’s marginal value inversely proportional to the number of participants sharing it (Page, 2018). A participant’s unnormalized shapley score is the sum across their ideas: $R_{i,t}^{\text{shapley}} = \sum_{x \in \mathcal{I}_{i,t}} \frac{1}{m_{t,k(x)}}$.

(iii) **uniqueness**: Ideas in singleton buckets ($m_{t,k} = 1$) receive a score of 1, and all others receive 0 (Forthmann et al., 2020; Baten et al., 2021, 2024). A participant’s unnormalized uniqueness score is the count of their unique ideas: $R_{i,t}^{\text{uniqueness}} = \sum_{x \in \mathcal{I}_{i,t}} \mathbb{I}\{m_{t,k(x)} = 1\}$.

(iv) **threshold**: Ideas are scored by a tiered function $S(x)$ based on bucket prevalence (Olson et al., 2021; DeYoung et al., 2008; Forthmann et al., 2020):

$$S(x) = \begin{cases} 3 & \text{if } \frac{m_{t,k(x)}}{N} \leq 0.01, \\ 2 & \text{if } 0.01 < \frac{m_{t,k(x)}}{N} \leq 0.03, \\ 1 & \text{if } 0.03 < \frac{m_{t,k(x)}}{N} \leq 0.10, \\ 0 & \text{otherwise.} \end{cases}$$

A participant’s unnormalized threshold score is the sum of these scores: $R_{i,t}^{\text{thresh}} = \sum_{x \in \mathcal{I}_{i,t}} S(x)$.

To compute a participant’s overall unnormalized score across all tasks, we take $R_i^{\text{metric}} = \sum_{t=1}^T R_{i,t}^{\text{metric}}$. To account for fluency (i.e., the number of ideas, $n_{i,t}$, contributed by participant p_i in task t), we define normalized originality as, $O_{i,t}^{\text{metric}} = \frac{R_{i,t}^{\text{metric}}}{n_{i,t}}$ and $O_i^{\text{metric}} = \sum_{t=1}^T O_{i,t}^{\text{metric}}$.

4.3 Evaluation Strategy

We assess construct validity along two dimensions:

(i) alignment between computational and human idea-to-bucket clustering, and (ii) agreement in participant-level originality scoring.

Bucket-level construct validity. The bucket labels are categorical and arbitrary. Moreover, the bucket sizes follow a fat-tailed distribution with a few highly frequent buckets and many rare ones (see §5.1). Thus, traditional clustering metrics (e.g., Adjusted Rand Index) can be misleading due to being inflated by rare buckets. We therefore adopt

Adjusted Mutual Information (AMI) (Vinh et al., 2010) as our primary metric to evaluate idea-to-bucket clustering alignment between our proposed method and human annotations. This metric adjusts for chance agreement, is robust to label permutation and skewed distributions, and is well-suited for comparing clusterings with different numbers of clusters. For insight development, we also use *Normalized Mutual Information (NMI)* (Vinh et al., 2010), which quantifies mutual dependence between clusterings without chance correction, and *V-measure* (Rosenberg and Hirschberg, 2007), which is the harmonic mean of homogeneity and completeness, reflecting both internal purity and cross-cluster coverage.

Participant-level construct validity. For originality scoring agreement, we use (i) *Zero-order Correlations* (Pearson’s r for linear agreement and Spearman’s ρ for monotonic consistency), (ii) *Intraclass Correlation Coefficient* for consistency across judges (Shrout and Fleiss, 1979), and (iii) *Bland–Altman Plots* to identify systematic, scale-level biases (Bland and Altman, 1986).

Convergent and external validity. *Convergent validity* is assessed by correlating model originality scores with theoretically aligned creativity metrics (e.g., Creativity Quotient and Creative Quality Ratings). *External validity* is evaluated by correlating model scores with established psychological and cognitive variables: personality traits, creative metaphor generation ability, fluid intelligence, and creative self-concept (Beaty and Johnson, 2021).

5 Understanding Human-Annotated Ground Truth Characteristics

5.1 Distributional Properties of Idea Buckets

We assess the structure of idea diversity in socialmuse24 using H1 and H2’s buckets. H1 created more buckets per task (399.6, 95% CI: [354.1, 445.1]) than H2 (230.8 [192.8, 268.8]), indicating finer-grained distinctions.

To examine bucket size (idea frequency) distributions, we fit a discrete power-law model to the bucket sizes for each task and compare it to a lognormal distribution via a likelihood ratio test (Clauset et al., 2009). Both annotators produced fat-tailed distributions, with scaling exponents $\alpha_{H1} = 2.01$ [1.73, 2.28] and $\alpha_{H2} = 1.74$ [1.60, 1.88], consistent with power-law like behavior in linguistic and social systems ($\alpha \approx 2$ to 3) (Newman, 2018). This confirms that a few

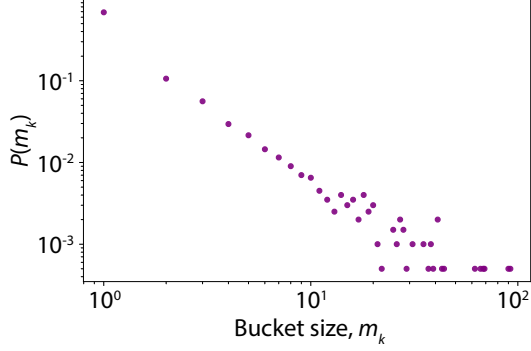


Figure 1: Idea bucket size distribution based on annotator H1’s bucketing. See Figure A1 for H2’s case.

buckets are highly frequent while many are rare (Figure 1). However, the power-law model is not statistically favored over lognormal ($P \geq 0.05$), suggesting that despite being fat-tailed, bucket size distributions are not strictly power-law and may be better described by lognormal or other alternatives.

5.2 Inter Human Annotator Agreement on Idea-level Bucketing

H1 and H2 show a mean AMI of 0.66 [0.64, 0.68], indicating strong alignment beyond what would be expected by random bucketing. NMI elucidates how informative one annotator’s bucketing is about the other’s without adjusting for chance (i.e., NMI is less conservative). As expected, the mean NMI is higher at 0.85 [0.84, 0.87], reflecting strong underlying structure shared across annotators (Table A1).

V-measure also yields a high mean of 0.85 [0.84, 0.87]. Its homogeneity component (0.80) shows that H1’s buckets are reasonably pure with respect to H2, and its high completeness component (0.92) shows that H2’s buckets almost perfectly recover H1’s buckets. This pattern corroborates that H1 split buckets more finely than H2, but both annotators identified similar idea groupings.

Overall, the annotators strongly agreed on their idea bucketing, despite granularity differences.

5.3 Inter Human Annotator Agreement on Participant-level Originality Scoring

We compute participant-level $\{O_i^{\text{metric}}\}$ using H1 and H2’s bucket assignments and assess agreement.

The threshold and shapley metrics show the strongest correlations (threshold: $r = 0.77$ [0.69, 0.84]; shapley: $r = 0.79$ [0.70, 0.85]; both $P < 0.001$). uniqueness and rarity show lower but still good correlations (uniqueness: $r = 0.73$ [0.63, 0.81]; rarity: $r = 0.72$ [0.61, 0.81]; both

$P < 0.001$; see Table A2 for ρ estimates).

The threshold and shapley metrics also show the strongest average consistency across judges: $ICC(3, k) = 0.85$ [0.78, 0.90], $P < 0.001$ for both. uniqueness yields the lowest but good agreement: $ICC(3, k) = 0.8$ [0.71, 0.86], $P < 0.001$ (Table A3). Together, we note strong agreements in originality scoring across the human annotators.

5.4 Insights for MUSESCORER Development

These analyses help establish expectations for machine-based originality scoring. *First*, human-annotated bucket sizes exhibit a fat-tailed structure. Any automated scoring system must account for this characteristic for its bucketing performance to approach the strong AMI baseline of humans.

Second, based on the above evidence, we take the threshold-based normalized scores, $\{O_i^{\text{thresh}}\}$, as our person-level gold standard against which we evaluate machine-based originality scoring. We test for robustness against the other metrics.

6 The MUSESCORER System

6.1 Insights from Early Prototypes

Our first prototype mimicked a human annotator’s workflow by comparing each new idea against an expanding codebook of *all* prior buckets. However, LLM prompts became intractable when K_t exceeded roughly 150. Given the scale-free bucket size distributions, massive corpora can have very large K_t , making exhaustive prompting infeasible. We therefore shifted to a retrieval-based approach, selecting a small candidate set for LLM judgment.

A second prototype let the LLM handle retrieval, decision-making, and codebook updates end-to-end, but this proved brittle—especially with smaller models (e.g., phi4). To improve stability, we offloaded retrieval and codebook management to external modules, leaving the LLM to focus solely on subjective bucketing.

6.2 MUSESCORER System Architecture

Algorithm 1 summarizes MUSESCORER’s workflow. The LLM processes one idea at a time and assigns it to a semantically equivalent bucket or creates a new one. A dynamic codebook is initialized for each ideation task and updated as new ideas arrive. For each idea $x \in \mathcal{X}$, a dictionary of candidate buckets \mathcal{D}_x is constructed via K -NN-based semantic search over the current codebook (Khandelwal et al., 2020). \mathcal{D}_x has a maximum size of K_c .

Algorithm 1 MUESCORER: LLM-Based Incremental Bucketing for a Single Creativity Task

Require: Idea set $\mathcal{X} = \{x_1, x_2, \dots, x_{|\mathcal{X}|}\}$, LLM, candidate dictionary size K_c

Ensure: Partition $\mathcal{B} = \{B_1, \dots, B_K\}$, assignment map $k(x)$

```
1: Initialize empty codebook  $\mathcal{C} \leftarrow \emptyset$ 
2: Initialize bucket index  $K \leftarrow 0$ 
3: for all ideas  $x \in \mathcal{X}$  do
4:   if  $|\mathcal{C}| \leq K_c$  then
5:      $\mathcal{D}_x \leftarrow \mathcal{C}$ 
6:   else
7:     Use  $K$ -NN search to find top- $K_c$  closest entries in  $\mathcal{C}$  to  $x$ 
8:      $\mathcal{D}_x \leftarrow \{(k_j, d_j)\}_{j=1}^{K_c}$ 
9:   end if
10:  Query LLM: “Is  $x$  a rephrasing of any  $d_j \in \mathcal{D}_x$ ? Return  $k_j$  or  $-1$ .” (In CoT prompting, also return a justification sentence)
11:  if LLM returns  $k^* \neq -1$  then
12:    Assign  $k(x) \leftarrow k^*$ 
13:     $B_{k^*} \leftarrow B_{k^*} \cup \{x\}$ 
14:  else
15:     $K \leftarrow K + 1$ 
16:    Create new bucket  $B_K \leftarrow \{x\}$ 
17:    Update codebook  $\mathcal{C} \leftarrow \mathcal{C} \cup \{(K, x)\}$ 
18:    Assign  $k(x) \leftarrow K$ 
19:  end if
20: end for
21: return  $\mathcal{B} = \{B_1, \dots, B_K\}, k(x) \forall x \in \mathcal{X}$ 
```

When the number of existing buckets is smaller than K_c , all of those buckets are taken in \mathcal{D}_x . Each candidate in the dictionary $\mathcal{D}_x = \{(k_j, d_j)\}_{j=1}^{K_c}$ maps bucket IDs to representative descriptions.

We employ two kinds of prompting strategies: (i) In *vanilla* prompting, the LLM determines whether x is a rephrasing of any d_j . If so, it returns the corresponding key k_j ; otherwise, it returns -1 , signaling the creation of a new bucket with x as its description. (ii) In Chain-of-Thought (CoT) prompting, the LLM additionally provides a one-sentence reasoning (Wei et al., 2022). The codebook and bucket assignment are updated accordingly.

It is important to distinguish this design from semantic-similarity-based clustering methods, which can be used to bucket ideas directly (we employ such methods as our computational baselines; see §7.1). Such methods typically attempt to assign ideas to clusters based on embedding dis-

tances, which risks conflating distinct intents or over-separating true rephrasings. In contrast, our use of K -NN-based retrieval only serves to keep the comparison space (and thereby prompt length) tractable; the final decision about bucket membership is made by the LLM. This separation ensures that semantic similarity supports efficiency, while the subjective, intent-sensitive aspect of bucketing remains with the judge LLM.

We fix $K_c = 10$ to allow a manageable prompt length while leaving sufficient margin for retrieval noise, and test robustness against other K_c choices. We experiment with a factorial combination (‘MUESCORER configurations’) of LLM model variants (llama3.3, qwen3, and phi4), prompting strategies (*vanilla* and CoT), and sentence embeddings (Appendix Section A.1).

7 Results and Discussion

7.1 Computational Baselines

We use unsupervised clustering to establish a computational baseline for MUESCORER. We require algorithms that (i) allow clusters of vastly different sizes, including fat-tail distributed ones, and (ii) preserve singleton and rare buckets without dropping them as noise or outliers (§5.4).

These constraints discourage us from using algorithms like DBScan (singleton and rare buckets are likely to be marked as noise) (Ester et al., 1996) and HDBScan (minimum cluster size is 2) (Campello et al., 2013), and our experiments also corroborate their poor performance. K -means clustering is poor at handling imbalanced cluster sizes or shapes, and requires the number of clusters to be close to the number of datapoints to allow many singleton or rare buckets (MacQueen, 1967). Agglomerative hierarchical clustering is a reasonable choice for our constraints (Ward Jr, 1963).

We report results with K -means and agglomerative algorithms. For each algorithm, we automatically search for the optimal number of buckets K_t over the full range of $K_t = 1$ to $|\mathcal{X}_t|$. To facilitate this search, we evaluate structural and semantic criteria using two metrics: (i) *Silhouette Score*, which assesses cluster quality based on geometric compactness and separation, with higher values indicating better-defined clusters (Rousseeuw, 1987); and (ii) *Semantic Score*, which is the geometric mean of coherence (intra-cluster similarity) and exclusivity (inter-cluster distinctiveness), encouraging clusters that are both internally consistent

	Model	AMI	NMI	Pearson's r	Spearman's ρ	ICC(3,1)
MUSESCORER	llama3.3, CoT	0.59 \pm 0.05	0.88 \pm 0.02	0.88 \pm 0.04	0.87 \pm 0.05	0.88 \pm 0.04
	qwen3, CoT	0.56 \pm 0.05	0.87 \pm 0.02	0.79 \pm 0.07	0.78 \pm 0.07	0.77 \pm 0.08
	phi4, CoT	0.54 \pm 0.01	0.83 \pm 0.01	0.78 \pm 0.08	0.76 \pm 0.08	0.72 \pm 0.09
	llama3.3, vanilla	0.59 \pm 0.03	0.86 \pm 0.02	0.83 \pm 0.06	0.79 \pm 0.07	0.81 \pm 0.06
	phi4, vanilla	0.53 \pm 0.02	0.83 \pm 0.01	0.80 \pm 0.07	0.78 \pm 0.08	0.75 \pm 0.08
Baseline	K -means, Silhouette	0.32 \pm 0.09	0.86 \pm 0.02	0.65 \pm 0.11	0.67 \pm 0.11	0.62 \pm 0.12
	K -means, Semantic	0.35 \pm 0.06	0.87 \pm 0.02	0.71 \pm 0.10	0.70 \pm 0.10	0.67 \pm 0.10
	Aggl., Silhouette	0.39 \pm 0.02	0.85 \pm 0.02	0.73 \pm 0.09	0.68 \pm 0.10	0.69 \pm 0.10
	Aggl., Semantic	0.31 \pm 0.05	0.86 \pm 0.02	0.65 \pm 0.11	0.65 \pm 0.11	0.61 \pm 0.12

Table 2: Agreement metrics comparing computational models to H1’s ground truths. Values are means \pm half-width of the 95% C.I. ($N = 109$). See Table A5 for results based on H2’s annotations, which replicate identical takeaways.

and mutually distinct (Mimno et al., 2011).

7.2 Distributional Properties of Computationally-labeled Idea Buckets

We find that K -means and agglomerative algorithms produce an exorbitantly high K_t , with 831 and 797 buckets produced by the K -means algorithm (respectively based on Silhouette and Semantic scores), and 588 and 838 buckets by the agglomerative algorithm. For reference, $|\mathcal{X}_t| \approx 1141$ per task in socialmuse24. These bucket counts are significantly higher than H1 and H2’s annotations ($P < 0.001$; see §5.1). In contrast, the MUSESCORER configurations produce K_t in the range of 255 to 465, overlapping those of the humans. The scaling exponents of K -means and agglomerative are systematically higher than the human baseline ($P < 0.001$), but the MUSESCORER configurations align with humans (Table A4).

7.3 Construct Validity of Idea-level Bucketing

Table 2 and Figure A2 show the AMI and NMI agreements between H1 and machine bucketing. Taking H2 as the reference replicates identical insights (see Table A5). Interestingly, all methods score highly in the less conservative NMI metric and match the H1-H2 agreement, showing reasonable preservation of semantic grouping.

However, when we correct for random chance and penalize mismatch in structure and granularity using the AMI metric, the MUSESCORER configurations sustain human-like performance while the K -means and agglomerative algorithms suffer dramatically and systematically. Specifically, against a human-human AMI of 0.66 [0.64, 0.68], the llama3.3 LLM with CoT prompting achieves the best AMI among the MUSESCORER configurations at 0.59 [0.55, 0.64], while the silhouette-

tuned agglomerative algorithm manages the best AMI among the baseline models at a poor 0.39 [0.36, 0.41]. This is unsurprising, since a drop in AMI implies deviation from the structure and resolution of the human bucketing, which is corroborated by the systematically larger number of buckets K -means and agglomerative algorithms produce. In contrast, the MUSESCORER configurations preserve more of the mutual structures, semantic coherence, and resolution, capturing up to 89% of the fine-grained patterns humans see.

Overall, MUSESCORER shows strong idea-bucketing alignment with the humans, surpassing the performances of clustering-based baselines.

7.4 Construct Validity of Participant-level Originality Scoring

Table 2 and Figure A3 show the participant-level $\{O_i^{\text{thresh}}\}$ score agreements based on H1 and machine bucketing. The results are robust to taking H2 as the reference (Table A5). MUSESCORER with llama3.3 and CoT prompting once again shows the best correlation ($r = 0.89$ [0.83, 0.92], $P < 0.001$). The baselines perform significantly worse, with the silhouette-tuned agglomerative algorithm achieving the best baseline correlation ($r = 0.73$ [0.63, 0.81], $P < 0.001$).

MUSESCORER with llama3.3 and CoT prompting also shows the best $ICC(3,1) = 0.88$ [0.83, 0.92], $P < 0.001$. The clustering baselines reach a maximum of $ICC(3,1) = 0.69$ [0.57, 0.77], $P < 0.001$, with the silhouette-tuned agglomerative model, performing significantly worse than llama3.3 ($P < 0.001$). Based on the above evidence, we pick llama3.3 with CoT prompting as the default configuration for MUSESCORER and use it for the remaining analysis.

We next visualize a Bland-Altman plot to iden-

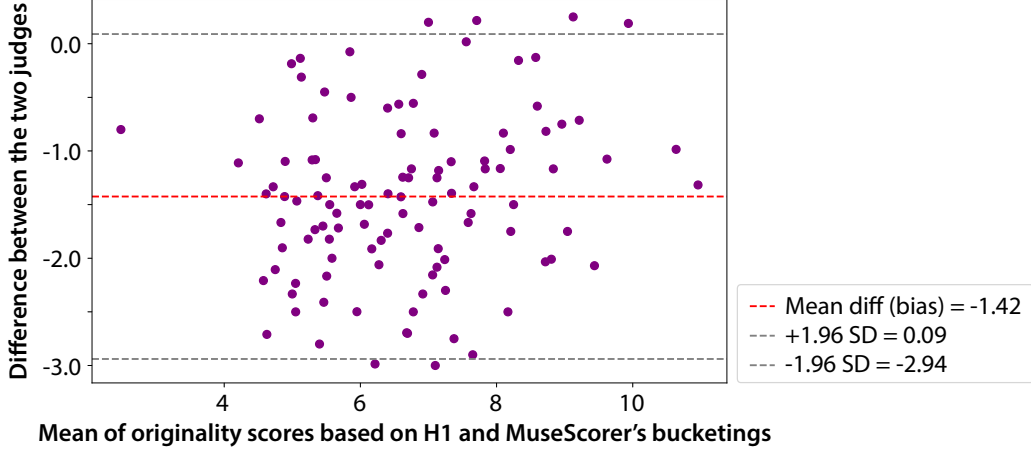


Figure 2: Bland-Altman visualization for bias detection.

tify systematic biases between H1 and MUSESCORER-derived originality scores (Figure 2). 94.5% of the points fall within the limits of agreement (LoA) of ± 1.96 SDs, and so does the mean difference (bias). This shows that MUSESCORER-derived scores stay strongly in line with human scores across the originality spectrum. Although the proportional bias regression slope is slightly positive (0.09), the effect is not statistically significant ($P > 0.05$), suggesting no systematic trend where the machine over- or under-scores ideas as originality increases. This supports the conclusion that MUSESCORER provides stable, human-comparable originality assessments.

Taken together, MUSESCORER shows strong construct validity in originality scoring against human ground truth.

7.5 Convergent and External Validity

We evaluate MUSESCORER for convergent and external validity against established creativity measures. Table 3 summarizes the correlations.

For convergent validity, MUSESCORER’s normalized originality scores $\{O_i^{\text{thresh}}\}$ correlate strongly with Creativity Quotient (CQ) scores in socialmuse24. CQ is a flexibility measure that captures the diversity of semantic categories. However, CQ is unnormalized and confounded by idea fluency. Unsurprisingly, unnormalized $\{R_i^{\text{thresh}}\}$ scores show a stronger correlation with CQ. In rating-based datasets, MUSESCORER’s $\{O_i^{\text{thresh}}\}$ scores correlate highly with human creative quality judgments (beaty18, silvia17, beaty12) and with rating-based originality in mohr16. The latter dataset also contains manually annotated flexibility scores, which do not account for fluency. Unsur-

prisingly, these flexibility scores correlate strongly with unnormalized $\{R_i^{\text{thresh}}\}$. Together, these findings confirm that MUSESCORER captures core constructs of creativity with high fidelity.

For external validity, MUSESCORER’s $\{O_i^{\text{thresh}}\}$ scores correlate systematically with metaphor generation quality, openness to experience, and self-reported creative identity and self-efficacy. We did not observe systematic associations with fluid intelligence or other Big Five traits. Our results largely corroborate previous insights (Beaty and Johnson, 2021), underscoring the system’s broader external validity.

7.6 Robustness

The results depend on LLM, sentence embedding, and prompting strategy choices. We obtain the best MUSESCORER results with a configuration comprising the llama3.3:70b LLM (Meta AI, 2024), e5-large-v2 sentence embedding (Wang et al., 2022), and Chain-of-Thought prompting (Wei et al., 2022) (§A.1). We further probe this configuration’s robustness across $K_c \in \{5, 15\}$, and find results statistically similar to the default $K_c = 10$. To assess ordering effects, we run the configuration with randomly ordered \mathcal{X}_t across 3 seeds. We find the results stable within the bounds reported in Table 2. The main results with the threshold metric are largely reproduced by the other three metrics. But we find that rarity shows proportional bias in the Bland-Altman plot (slope = 0.2, $P < 0.01$), while shapley and uniqueness show no correlation with openness in the silvia17 dataset, losing some external validity. The threshold metric thus emerges as the most robust choice for operationalizing statistical infrequency for originality scoring.

Dataset	Comparison Variable	Correlation
<i>Convergent Validity</i>		
socialmuse24	Creativity Quotient (CQ)	$r = 0.40 [0.23, 0.55], P < 0.001$
socialmuse24	CQ vs. unnormalized R_i^{thresh}	$r = 0.48 [0.32, 0.62], P < 0.001$
beaty18	Creative Quality (mean ratings)	$r = 0.77 [0.71, 0.83], P < 0.001$
silvia17	Creative Quality (mean ratings)	$r = 0.54 [0.41, 0.65], P < 0.001$
beaty12	Creative Quality (mean ratings)	$r = 0.42 [0.27, 0.55], P < 0.001$
mohr16	Rating-based Originality	$r = 0.42 [0.35, 0.49], P < 0.001$
mohr16	Flexibility vs. unnormalized R_i^{thresh}	$r = 0.76 [0.73, 0.80], P < 0.001$
<i>External Validity</i>		
beaty18	Metaphor Generation	$r = 0.17 [0.02, 0.32], P < 0.05$
beaty12	Metaphor Generation	$r = 0.25 [0.08, 0.40], P < 0.01$
beaty18	Openness	$\rho = 0.16 [0.01, 0.30], P < 0.05$
beaty12	Openness	$r = 0.30 [0.14, 0.45], P < 0.001$
silvia17	Openness	$\rho = 0.14 [-0.02, 0.30], P = 0.09$
beaty18	Creative Self-Identity	$r = 0.34 [0.19, 0.48], P < 0.001$
beaty18	Creative Self-Efficacy	$r = 0.29 [0.14, 0.44], P < 0.001$

Table 3: Convergent and external validity of MUSESCORER’s originality scores, $\{O_i^{\text{thresh}}\}$. Values are Pearson (r) or Spearman (ρ) correlations with 95% C.I. and significance levels.

8 Conclusion

This work introduces MUSESCORER, a scalable, zero-shot system for scoring the originality of creative ideas. By combining the LLM-as-a-judge paradigm with externally orchestrated retrieval, our method produces psychometrically aligned, intent-sensitive judgments without requiring task-specific fine-tuning or training data.

Across five distinct AUT datasets, MUSESCORER demonstrates robust and consistent performance despite variation in task structures and idea distributions. Unlike opaque embedding-only approaches, our use of chain-of-thought (CoT) prompting yields interpretable outputs, allowing the system to provide justifications for bucketing decisions transparently.

Our approach is well-suited to support the growing body of research on human and AI creativity, particularly as large-scale, high-throughput studies become increasingly common (Doshi and Hauser, 2024; Chakrabarty et al., 2025; Tanveer et al., 2018; Kelty et al., 2025). By combining reliability, interpretability, and scale, this system expands the practical and methodological toolkit for researchers and opens new avenues for measuring and understanding creative potential in human and AI agents.

Limitations

Several limitations should be considered in future work when extending frequency-based origi-

inality scoring. First, demographic fairness and accessibility remain important concerns. Variations in language use across cultural or educational backgrounds—especially in non-English contexts—may influence bucketing judgments and introduce bias if not carefully monitored.

Second, our validation is confined to AUT-style, text-based divergent thinking tasks. How well the approach generalizes to other creative domains (e.g., design, visual arts) remains an open question.

Third, while externally orchestrated retrieval mitigates some variability, the system remains sensitive to prompt length and phrasing (Liu et al., 2023). Subtle formatting changes can affect judgment quality, suggesting that prompt engineering and robustness testing deserve further study.

Fourth, efficiency may have room for improvement. We process ideas one at a time, which stabilizes performance—particularly for smaller models—but limits throughput. Future work could explore batching or multi-step reasoning to increase efficiency, though at potential cost to stability and computation.

Fifth, we kept the candidate retrieval size small ($K_c = \{5, 10, 15\}$). Larger candidate sets may improve coverage but increase token usage and cost. Similarly, our most effective threshold metric uses a heuristic tiering scheme adopted from prior literature; the robustness of these cutoffs remains to be validated.

Finally, as with all LLM-based systems, hallucination is a risk. In our case, hallucination manifests as misassigning an idea to the wrong bucket. However, MUSESCORER achieves strong alignment with human annotations and passes psychometric validation despite this risk, suggesting the system is reasonably reliable within scope.

Ethical Considerations

We reanalyzed public datasets from prior works (consistent with their intended use) and did not collect any new human data for this research. Given the nature of the research in creative assessment, we do not readily foresee potential harm or risk.

References

- Marah Abidin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, and 8 others. 2024. [Phi-4 technical report](#). *Preprint*, arXiv:2412.08905.
- Selcuk Acar and Mark A Runco. 2014. Assessing associative distance among ideas elicited by tests of divergent thinking. *Creativity Research Journal*, 26(2):229–238.
- Raiyan Abdul Baten, Richard N Aslin, Gourab Ghoshal, and Ehsan Hoque. 2021. Cues to gender and racial identity reduce creativity in diverse social networks. *Scientific Reports*, 11(1):10261.
- Raiyan Abdul Baten, Richard N Aslin, Gourab Ghoshal, and Ehsan Hoque. 2022. Novel idea generation in social networks is optimized by exposure to a “Goldilocks” level of idea-variability. *PNAS Nexus*, 1(5):pgac255.
- Raiyan Abdul Baten, Daryl Bagley, Ashely Tenesaca, Famous Clark, James P Bagrow, Gourab Ghoshal, and Ehsan Hoque. 2020. Creativity in temporal social networks: How divergent thinking is impacted by one’s choice of peers. *Journal of the Royal Society Interface*, 17(171):20200667.
- Raiyan Abdul Baten, Ali Sarosh Bangash, Krish Veera, Gourab Ghoshal, and Ehsan Hoque. 2024. AI can enhance creativity in social networks. *arXiv preprint arXiv:2410.15264*.
- Roger E Beaty and Dan R Johnson. 2021. Automating creativity assessment with SemDis: An open platform for computing semantic distance. *Behavior Research Methods*, 53(2):757–780.
- Roger E Beaty, Yoed N Kenett, Alexander P Christensen, Monica D Rosenberg, Mathias Benedek, Qunlin Chen, Andreas Fink, Jiang Qiu, Thomas R Kwapiil, Michael J Kane, and 1 others. 2018. Robust prediction of individual creative ability from brain functional connectivity. *Proceedings of the National Academy of Sciences*, 115(5):1087–1092.
- Roger E Beaty and Paul J Silvia. 2012. Why do ideas get more creative across time? An executive interpretation of the serial order effect in divergent thinking tasks. *Psychology of Aesthetics, Creativity, and the Arts*, 6(4):309.
- Roger E Beaty and Paul J Silvia. 2013. Metaphorically speaking: Cognitive abilities and the production of figurative language. *Memory & cognition*, 41:255–267.
- Kenes Beketayev and Mark A Runco. 2016. Scoring divergent thinking tests by computer with a semantics-based algorithm. *Europe’s Journal of Psychology*, 12(2):210.
- J. Martin Bland and Douglas G. Altman. 1986. [Statistical methods for assessing agreement between two methods of clinical measurement](#). *The Lancet*, 327(8476):307–310.
- Terry Bossomaier, Mike Harré, Anthony Knittel, and Allan Snyder. 2009. A semantic network approach to the Creativity Quotient (CQ). *Creativity Research Journal*, 21(1):64–71.
- Thomas J Bouchard Jr and Melana Hare. 1970. Size, performance, and potential in brainstorming groups. *Journal of Applied Psychology*, 54(1p1):51.
- Philip Buczak, He Huang, Boris Forthmann, and Philipp Doebler. 2023. The machines take over: A comparison of various supervised learning approaches for automated scoring of divergent thinking tasks. *The Journal of Creative Behavior*, 57(1):17–36.
- Ricardo J. G. B. Campello, Davoud Moulavi, and Jörg Sander. 2013. Density-based clustering based on hierarchical density estimates. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD)*, pages 160–172. Springer.
- Raymond Bernard Cattell and Alberta KS Cattell. 1960. *Measuring intelligence with the culture fair tests*. Institute for Personality and Ability Testing.
- Tuhin Chakrabarty, Philippe Laban, Divyansh Agarwal, Smaranda Muresan, and Chien-Sheng Wu. 2024. Art or artifice? Large language models and the false promise of creativity. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–34.
- Tuhin Chakrabarty, Philippe Laban, and Chien-Sheng Wu. 2025. Can AI writing be salvaged? Mitigating idiosyncrasies and improving human-AI alignment in the writing process through edits. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pages 1–33.

- Robert Chew, John Bollenbacher, Michael Wenger, Jessica Speer, and Annice Kim. 2023. LLM-assisted content analysis: Using large language models to support deductive coding. *arXiv preprint arXiv:2306.14924*.
- Aaron Clauset, Cosma Rohilla Shalizi, and Mark EJ Newman. 2009. Power-law distributions in empirical data. *SIAM Review*, 51(4):661–703.
- Shih-Chieh Dai, Aiping Xiong, and Lun-Wei Ku. 2023. LLM-in-the-loop: Leveraging large language model for thematic analysis. *arXiv preprint arXiv:2310.15100*.
- Colin G DeYoung, Joseph L Flanders, and Jordan B Peterson. 2008. Cognitive abilities involved in insight problem solving: An individual differences model. *Creativity Research Journal*, 20(3):278–290.
- Anil R Doshi and Oliver P Hauser. 2024. Generative AI enhances individual creativity but reduces the collective diversity of novel content. *Science Advances*, 10(28):eadn5290.
- Denis Dumas and Kevin N Dunbar. 2014. Understanding fluency and originality: A latent variable perspective. *Thinking Skills and Creativity*, 14:56–67.
- Denis Dumas, Peter Organisciak, and Michael Doherty. 2021. Measuring divergent thinking originality with human raters and text-mining models: A psychometric comparison of methods. *Psychology of Aesthetics, Creativity, and the Arts*, 15(4):645.
- Ruth B Ekstrom, John W French, Harry H Harman, and D Dermen. 1976. Manual for kit of factor-referenced tests. *Princeton, NJ: Educational Testing Service*, 586:1989–1995.
- Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*, pages 226–231.
- Boris Forthmann, Heinz Holling, Pınar Çelik, Martin Storme, and Todd Lubart. 2017. Typing speed as a confounding variable and the measurement of quality in divergent thinking. *Creativity Research Journal*, 29(3):257–269.
- Boris Forthmann, Sue Hyeon Paek, Denis Dumas, Baptiste Barbot, and Heinz Holling. 2020. Scrutinizing the basis of originality in divergent thinking tests: On the measurement precision of response propensity estimates. *British Journal of Educational Psychology*, 90(3):683–699.
- Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Shiping Yang, and Xiaojun Wan. 2023. Human-like summarization evaluation with ChatGPT. *arXiv preprint arXiv:2304.02554*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Joy Paul Guilford. 1967. *The Nature of Human Intelligence*. McGraw-Hill.
- JP Guilford, PR Christensen, PR Merrifield, and RC Wilson. 1978. *Alternate Uses: Manual of Instructions and Interpretation*. Orange, CA: Sheridan Psychological Services.
- Alicia Hofelich Mohr, Andrew Sell, and Thomas Lindsay. 2016. Thinking inside the box: Visual design of the response box affects creative divergent thinking in an online survey. *Social Science Computer Review*, 34(3):347–359.
- Gautier Izacard and Edouard Grave. 2020. Distilling knowledge from reader to retriever for question answering. *arXiv preprint arXiv:2012.04584*.
- Maciej Karwowski. 2014. Creative mindsets: Measurement, correlates, consequences. *Psychology of Aesthetics, Creativity, and the Arts*, 8(1):62.
- Sean Keltz, Raiyan Abdul Baten, Adiba Mahbub Proma, Ehsan Hoque, Johan Bollen, and Gourab Ghoshal. 2025. [The innovation trade-off: how following superstars shapes academic novelty](#). *Humanities and Social Sciences Communications*, 12(1):926.
- Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through memorization: Nearest neighbor language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Kibeom Lee and Michael C Ashton. 2004. Psychometric properties of the HEXACO personality inventory. *Multivariate Behavioral Research*, 39(2):329–358.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP task. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, and 1 others. 2024a. From generation to judgment: Opportunities and challenges of LLM-as-a-judge. *arXiv preprint arXiv:2411.16594*.
- Dawei Li, Shu Yang, Zhen Tan, Jae Young Baik, Sukwon Yun, Joseph Lee, Aaron Chacko, Bojian Hou, Duy Duong-Tran, Ying Ding, and 1 others. 2024b. DALK: Dynamic co-augmentation of LLMs and KG

- to answer Alzheimer’s disease questions with scientific literature. *arXiv preprint arXiv:2405.04819*.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2023. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. *Lost in the middle: How language models use long contexts*. *Preprint*, arXiv:2307.03172.
- Ximing Lu, Melanie Sclar, Skyler Hallinan, Niloofar Mireshghallah, Jiacheng Liu, Seungju Han, Allyson Ettinger, Liwei Jiang, Khyathi Chandu, Nouha Dziri, and 1 others. 2024. AI as humanity’s Salieri: Quantifying linguistic creativity of language models via systematic attribution of machine text against web text. *arXiv preprint arXiv:2410.04265*.
- J. MacQueen. 1967. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pages 281–297. University of California Press.
- Robert R McCrae, Paul T Costa, Jr, and Thomas A Martin. 2005. The NEO-PI-3: A more readable revised NEO personality inventory. *Journal of Personality Assessment*, 84(3):261–270.
- Meta AI. 2024. LLaMA 3.3-70B-Instruct. <https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct>. Accessed: 2025-05-18.
- David Mimno, Hanna Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. 2011. Optimizing semantic coherence in topic models. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 262–272.
- Mark Newman. 2018. *Networks*. Oxford University Press.
- Jay A Olson, Johnny Nahas, Denis Chmoulevitch, Simon J Cropper, and Margaret E Webb. 2021. Naming unrelated words predicts creativity. *Proceedings of the National Academy of Sciences*, 118(25):e2022340118.
- Peter Organisciak, Selcuk Acar, Denis Dumas, and Kelly Berthiaume. 2023. Beyond semantic distance: Automated scoring of divergent thinking greatly improves with large language models. *Thinking Skills and Creativity*, 49:101356.
- Peter Organisciak and Denis Dumas. 2020. Open creativity scoring. <https://openscoring.du.edu>. [Computer software].
- Scott E. Page. 2018. *The Model Thinker: What You Need to Know to Make Data Work for You*. Basic Books, New York.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318.
- Nils Reimers and Iryna Gurevych. 2019. *Sentence-BERT: Sentence embeddings using siamese BERT-networks*. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Roni Reiter-Palmon, Boris Forthmann, and Baptiste Barbot. 2019. Scoring divergent thinking tests: A review and systematic framework. *Psychology of Aesthetics, Creativity, and the Arts*, 13(2):144.
- Andrew Rosenberg and Julia Hirschberg. 2007. *V-measure: A conditional entropy-based external cluster evaluation measure*. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 410–420, Prague, Czech Republic. Association for Computational Linguistics.
- Peter J Rousseeuw. 1987. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65.
- Mark A Runco and Garrett J Jaeger. 2012. The standard definition of creativity. *Creativity Research Journal*, 24(1):92–96.
- Mark A Runco and Wayne Mraz. 1992. Scoring divergent thinking tests using total ideational output and a creativity index. *Educational and Psychological Measurement*, 52(1):213–221.
- Yanxin Shen, Lun Wang, Chuanqi Shi, Shaoshuai Du, Yiyi Tao, Yixian Shen, and Hang Zhang. 2024. Comparative analysis of listwise reranking with large language models in limited-resource language contexts. *arXiv preprint arXiv:2412.20061*.
- Patrick E Shrout and Joseph L Fleiss. 1979. Intraclass correlations: Uses in assessing rater reliability. *Psychological Bulletin*, 86(2):420.
- Paul J Silvia, Emily C Nusbaum, and Roger E Beaty. 2017. Old or new? Evaluating the old/new scoring method for divergent thinking tasks. *The Journal of Creative Behavior*, 51(3):216–224.
- Paul J Silvia, Beate P Winterstein, John T Willse, Christopher M Barona, Joshua T Cram, Karl I Hess, Jenna L Martinez, and Crystal A Richard. 2008. Assessing creativity with divergent thinking tasks: Exploring the reliability and validity of new subjective scoring methods. *Psychology of Aesthetics, Creativity, and the Arts*, 2(2):68.
- Allan Snyder, John Mitchell, Terry Bossomaier, and Gerry Pallier. 2004. The Creativity Quotient: An objective scoring of ideational fluency. *Creativity Research Journal*, 16(4):415–419.

- C. Stevenson, I. Smal, M. Baas, M. Dahrendorf, R. Grasman, C. Tanis, E. Scheurs, D. Sleiffer, and H. van der Maas. 2020. [Automated AUT scoring using a big data variant of the consensual assessment technique](#). Report Final Technical Report, Modeling Creativity Project, Universiteit van Amsterdam, Amsterdam. Faculty of Social and Behavioural Sciences (FMG), Psychology Research Institute (PsyRes).
- M Iftekhar Tanveer, Samiha Samrose, Raiyan Abdul Baten, and M Ehsan Hoque. 2018. Awe the audience: How the narrative trajectories affect audience perception in public speaking. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–12.
- L. L. Thurstone. 1938. [Primary mental abilities](#). *The Mathematical Gazette*, 22(251):411–412.
- Nguyen Xuan Vinh, Julien Epps, and James Bailey. 2010. [Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance](#). *Journal of Machine Learning Research*, 11(95):2837–2854.
- Vijay Viswanathan, Kiril Gashtevski, Kiril Gashtevski, Carolin Lawrence, Tongshuang Wu, and Graham Neubig. 2024. [Large language models enable few-shot clustering](#). *Transactions of the Association for Computational Linguistics*, 12:321–333.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*.
- Joe H. Ward Jr. 1963. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association*, 58(301):236–244.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023a. [C-pack: Packaged resources to advance general chinese embedding](#). *Preprint*, arXiv:2309.07597.
- Ziang Xiao, Xingdi Yuan, Q Vera Liao, Rania Abdelghani, and Pierre-Yves Oudeyer. 2023b. Supporting qualitative analysis with large language models: Combining codebook with GPT-3 for deductive coding. In *Companion Proceedings of the 28th International Conference on Intelligent User Interfaces*, pages 75–78.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36:11809–11822.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTscore: Evaluating text generation with BERT. *arXiv preprint arXiv:1904.09675*.
- Lirui Zhao, Yue Yang, Kaipeng Zhang, Wenqi Shao, Yuxin Zhang, Yu Qiao, Ping Luo, and Rongrong Ji. 2024. Diffagent: Fast and accurate text-to-image API selection with large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6390–6399.

A Supplementary Materials

A.1 System Component Choices

We experiment with the following system component alternatives:

(i) Large language models: $\mathcal{M} = \{\text{llama3.3:70b-Instruct (Meta AI, 2024; Grattafiori et al., 2024), qwen3:32b (Yang et al., 2025), phi4:14b (Abdin et al., 2024)}\}$. We pick these mid-sized, open-source models for their cost and computation efficiencies.

(ii) Sentence embedding models: $\mathcal{E} = \{\text{all-mpnet-base-v2 (Reimers and Gurevych, 2019), bge-large-en-v1.5 (Xiao et al., 2023a), e5-large-v2 (Wang et al., 2022)}\}$. These models are freely available on Huggingface and have been widely used in recent technological developments.

(iii) Prompting strategies: $\mathcal{P} = \{\text{vanilla, CoT (Wei et al., 2022)}\}$.

In our experiments, we found the combination of llama3.3:70b-Instruct, e5-large-v2, and CoT to give the best performance.

A.2 Experimentation Setup and GPU Usage

We conducted all experiments using (i) an Intel Core i7-based computer with 64GB RAM and an RTX 3070 Ti graphics card, and (ii) three MacBook Pro laptops. All our code and data are available on GitHub. The R&D and final result generation took roughly 100 GPU days.

A.3 LLM Prompts

System Prompt (Vanilla Prompting)

You are an idea bucket annotator for ideas generated for the object {object_name} in Guilford's Alternative Uses Test. You will be given an input_idea to annotate against up to {comparison_k} comparison_ideas, given to you in a dictionary format with key-value pairs of comparison_idea_ID: comparison_idea_description. The keys are integers, and the values are strings. Your goal is to determine if the input_idea is a very obviously rephrased version of one of those comparison_idea_description, or if it is slightly different.
if input_idea is a very obviously rephrased version of a certain comparison_idea_description:
 your_annotation_ID = comparison_idea_ID
key of that comparison_idea_description value
elif input_idea is a slightly different one:
 your_annotation_ID = -1
Your response must be a text string containing exactly: <your_annotation_ID>.
For example: if your_annotation_ID is 6 since the input idea is a very obviously

rephrased version of comparison_idea_ID 6, your response string should be "6". Another example: if your_annotation_ID is -1 because the input idea is not an obvious rephrasing of any comparison_idea_ID, your response string should be "-1".
Absolutely do not provide any extra text.

System Prompt (CoT Prompting)

You are an idea bucket annotator for ideas generated for the object {object_name} in Guilford's Alternative Uses Test. You will be given an input_idea to annotate against up to {comparison_k} comparison_ideas, given to you in a dictionary format with key-value pairs of comparison_idea_ID: comparison_idea_description. The keys are integers, and the values are strings. Your goal is to determine if the input_idea is a very obviously rephrased version of one of those comparison_idea_description, or if it is slightly different.
if input_idea is a very obviously rephrased version of a certain comparison_idea_description:
 your_annotation_ID = comparison_idea_ID
key of that comparison_idea_description value
elif input_idea is a slightly different one:
 your_annotation_ID = -1
You will also provide a reason string containing a single sentence explaining why you gave the input_idea that specific your_annotation_ID.
Your response must be a text string containing exactly:
<your_annotation_ID><SPACE><reason>.
For example: if your_annotation_ID is 6 and the reason is "The input idea is a very obviously rephrased version of comparison_idea_ID 6", your response string should be "6 The input idea is a very obviously rephrased version of comparison_idea_ID 6".
Another example: if your_annotation_ID is -1 and the reason is "The input idea is not an obvious rephrasing of any comparison_idea_ID", your response string should be "-1 The input idea is not an obvious rephrasing of any comparison_idea_ID".
Absolutely do not provide any extra text.

User Prompt Per Idea (Both Conditions)

input_idea: {idea_text}
comparison_ideas: {repr(comparison_ideas)}

A.4 AI Usage

We used Grammarly AI to improve the grammatical accuracy of the manuscript, and ChatGPT to speed up the implementation of standard statistical analysis code.

A.5 Supplementary Tables and Figures

Table A1: Inter-human annotator agreement on idea bucketing in socialmuse24.

Metric	Mean [95% C.I.]
AMI	0.66 [0.64, 0.68]
NMI	0.85 [0.84, 0.88]
V-measure	0.85 [0.84, 0.87]
Homogeneity	0.80 [0.77, 0.82]
Completeness	0.92 [0.89, 0.95]

Table A2: Pearson and Spearman correlations between participant-level normalized O^{metric} scores based on H1 and H2’s bucketing. $N = 109$ in all cases.

Scoring Metric	Correlation Type	Estimate	95% C.I.	P-value
threshold	Pearson’s r	0.77	[0.69, 0.84]	$P < 0.001$
	Spearman’s ρ	0.75	[0.65, 0.82]	$P < 0.001$
shapley	Pearson’s r	0.79	[0.70, 0.85]	$P < 0.001$
	Spearman’s ρ	0.74	[0.64, 0.82]	$P < 0.001$
rarity	Pearson’s r	0.72	[0.61, 0.80]	$P < 0.001$
	Spearman’s ρ	0.64	[0.51, 0.74]	$P < 0.001$
uniqueness	Pearson’s r	0.73	[0.63, 0.81]	$P < 0.001$
	Spearman’s ρ	0.66	[0.54, 0.76]	$P < 0.001$

Table A3: ICC reliability of the participants’ normalized originality O^{metric} scores based on H1 and H2’s bucketing.

Scoring Metric	$ICC(3, k)$	F	$df1$	$df2$	P-value	95% C.I.
threshold	0.85	6.79	108	108	$P < 0.001$	[0.78, 0.90]
shapley	0.85	6.67	108	108	$P < 0.001$	[0.78, 0.90]
rarity	0.83	5.73	108	108	$P < 0.001$	[0.75, 0.88]
uniqueness	0.80	4.97	108	108	$P < 0.001$	[0.71, 0.86]

Table A4: Cluster count K and power-law exponent α for various computational scoring methods.

Model	K [95% C.I.]	α [95% C.I.]
llama3.3, CoT	465.4 [426.8, 504.0]	2.28 [2.14, 2.42]
qwen3, CoT	462.4 [432.7, 492.1]	2.43 [2.20, 2.67]
phi4, CoT	255.0 [207.3, 302.7]	2.39 [1.72, 3.05]
llama3.3, vanilla	367.8 [333.3, 402.3]	2.29 [1.97, 2.61]
phi4, vanilla	275.6 [229.5, 321.7]	2.51 [2.23, 2.78]
K -means, Silhouette	830.6 [729.2, 932.0]	3.12 [2.82, 3.43]
K -means, Semantic	797.4 [757.8, 837.0]	3.12 [2.67, 3.57]
Agglomerative, Silhouette	588.0 [524.9, 651.1]	5.68 [1.26, 10.09]
Agglomerative, Semantic	838.0 [815.9, 860.1]	3.80 [2.63, 4.97]

Table A5: Agreement metrics comparing computational models to H2’s ground truths. Values denote mean \pm half-width of the 95% C.I. ($N = 109$).

	Model	AMI	NMI	Pearson r	Spearman ρ	ICC(3,1)
MUSESCORER	llama3.3, CoT	0.57 ± 0.04	0.84 ± 0.02	0.76 ± 0.08	0.74 ± 0.09	0.74 ± 0.09
	qwen3, CoT	0.54 ± 0.04	0.83 ± 0.02	0.74 ± 0.09	0.73 ± 0.09	0.74 ± 0.09
	phi4, CoT	0.56 ± 0.03	0.79 ± 0.01	0.67 ± 0.10	0.68 ± 0.10	0.67 ± 0.10
	llama3.3, vanilla	0.59 ± 0.03	0.83 ± 0.01	0.76 ± 0.08	0.74 ± 0.09	0.75 ± 0.08
	phi4, vanilla	0.55 ± 0.04	0.80 ± 0.01	0.73 ± 0.09	0.71 ± 0.10	0.73 ± 0.09
Baseline	K -means, Silhouette	0.28 ± 0.07	0.80 ± 0.02	0.59 ± 0.12	0.62 ± 0.12	0.59 ± 0.12
	K -means, Semantic	0.30 ± 0.05	0.80 ± 0.02	0.66 ± 0.11	0.68 ± 0.10	0.66 ± 0.11
	Aggl., Silhouette	0.36 ± 0.03	0.80 ± 0.02	0.65 ± 0.11	0.60 ± 0.12	0.64 ± 0.11
	Aggl., Semantic	0.26 ± 0.05	0.80 ± 0.02	0.60 ± 0.12	0.64 ± 0.11	0.60 ± 0.12

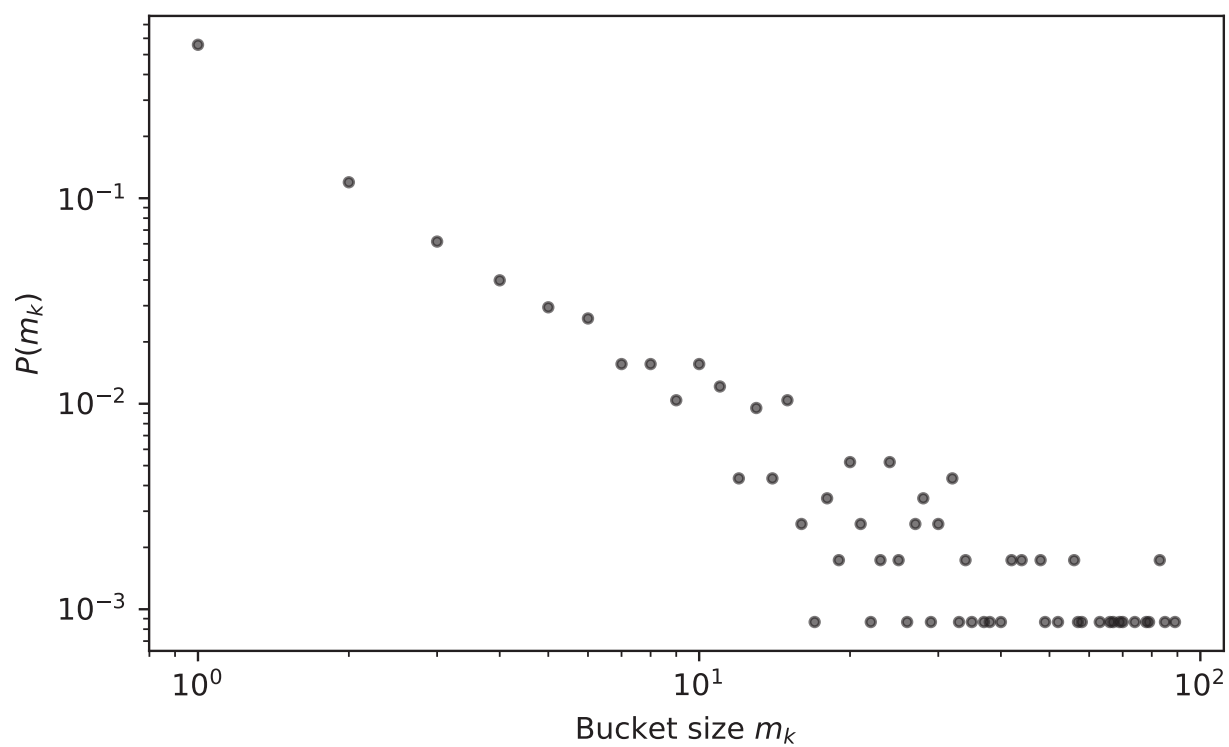


Figure A1: Idea bucket size distribution based on annotator H2's bucketing.

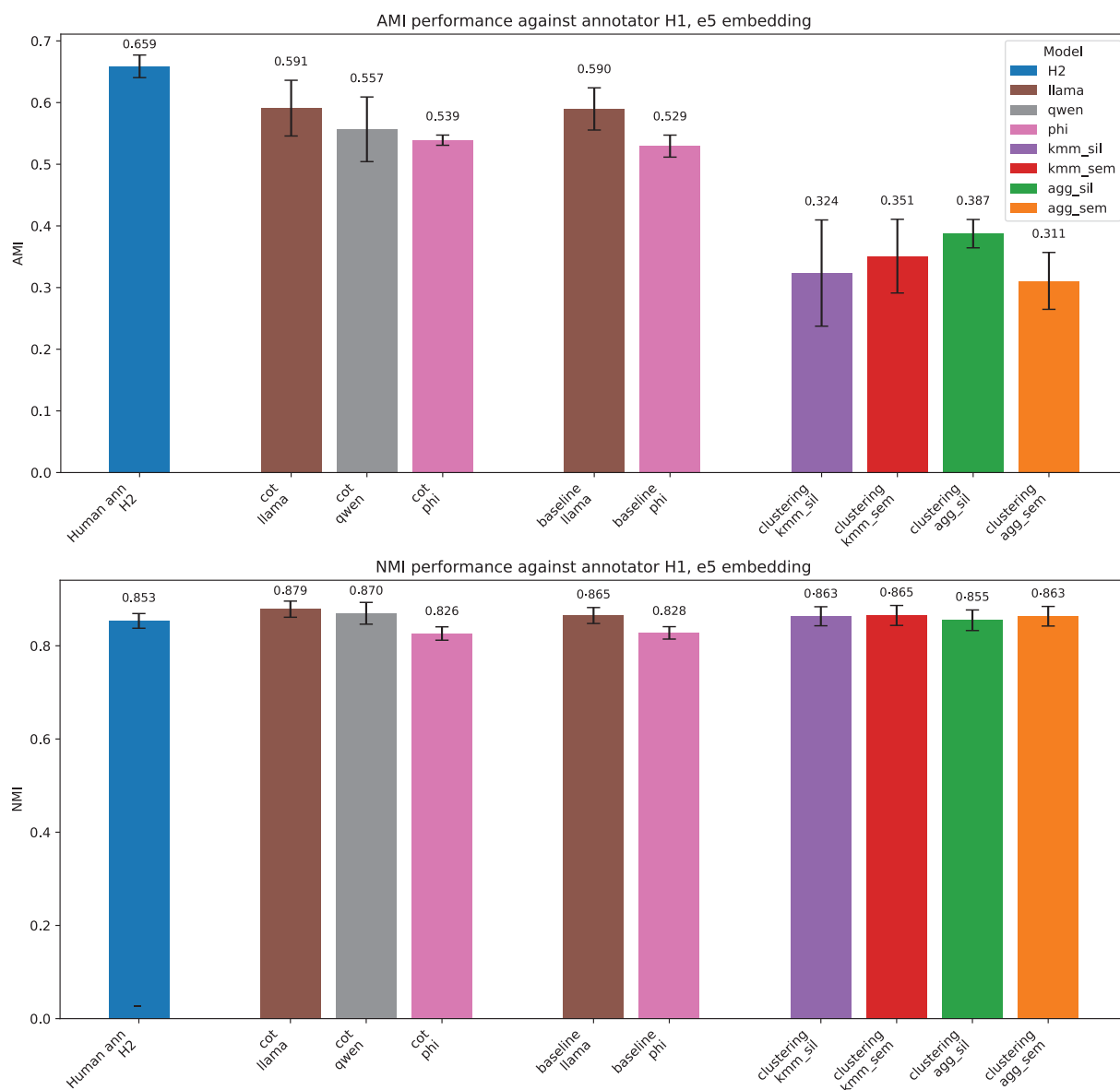


Figure A2: AMI and NMI performance comparison against annotator H1

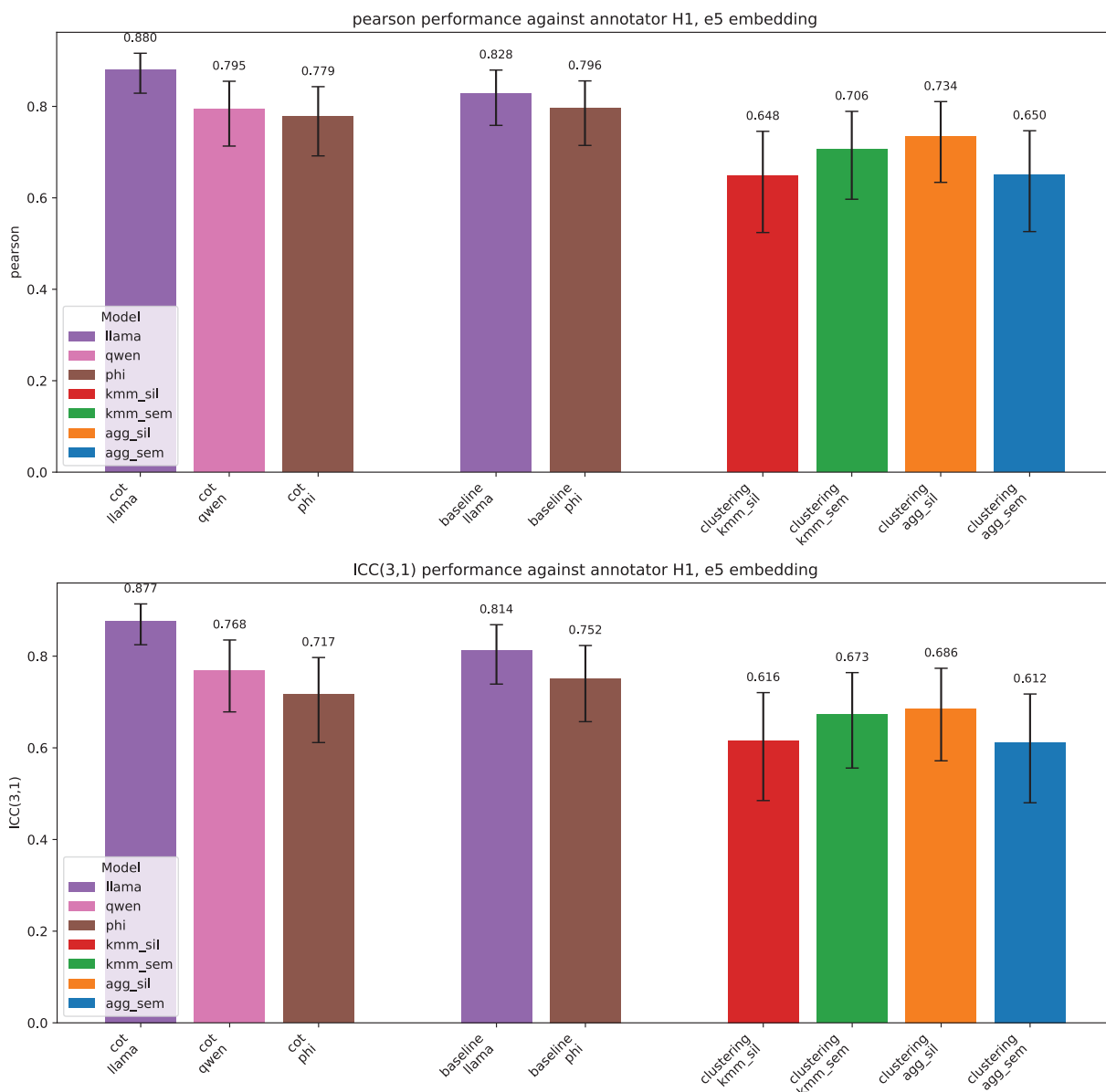


Figure A3: Pearson's r and ICC performance comparison against annotator H1