

# Confident Rankings with Fewer Items: Adaptive LLM Evaluation with Continuous Scores

Anonymous ACL submission

## Abstract

Computerized Adaptive Testing (CAT) has proven effective for efficient LLM evaluation on multiple-choice benchmarks, but modern LLM evaluation increasingly relies on generation tasks where outputs are scored continuously rather than marked correct/incorrect. We present a principled extension of IRT-based adaptive testing to continuous bounded scores (ROUGE, BLEU, LLM-as-a-Judge) by replacing the Bernoulli response distribution with a heteroskedastic normal distribution. Building on this, we introduce an uncertainty aware ranker with adaptive stopping criteria that achieves reliable model ranking while testing as few items and as cheaply as possible. We validate our method on five benchmarks spanning n-gram-based, embedding-based, and LLM-as-judge metrics. Our method uses 2% of the items while improving ranking correlation by 0.12  $\tau$  over random sampling, with 95% accuracy on confident predictions.

## 1 Introduction

Rigorous evaluation of large language models is essential, but current practice faces two challenges. The first is cost: exhaustive evaluation becomes expensive as the number of models, test items, and metrics grows. The second is methodological: score differences are often reported without significance testing, leading many studies to mistake statistical noise with improved performance and produce non-replicable results (Dror et al., 2018). Efficient evaluation methods that maintain statistical validity remain underexplored.

Computerized Adaptive Testing (CAT) dynamically selects informative test items to estimate model capabilities with far fewer evaluations than exhaustive benchmarking, and has emerged as a promising approach for efficient evaluation of large language models (Liu et al., 2024). These methods use adaptive use Item Response Theory (IRT)

as their theoretical foundation (Martínez-Plumed et al., 2016; Maia Polo et al., 2024; Rodriguez et al., 2021), with recent work extending this to fully adaptive item selection (Hofmann et al., 2025).

Existing approaches for LLM evaluation focus exclusively on multiple-choice datasets where responses can only be correct or incorrect. However, many of the tasks that reflect real use cases of LLMs such as summarization, dialogue, instruction following and machine translation must be scored on a continuous scale. These scoring strategies cover traditional metrics such as BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004), and embedding similarity-based scores such as BERTScore (Zhang et al., 2020) and COMET (Rei et al., 2020). LLM-as-judge evaluation (Zheng et al., 2024; Liu et al., 2023), now widely adopted for assessing instruction following and open-ended generation, also produces ordinal ratings that can be normalized and treated as continuous.

We present a continuous extension of IRT-based adaptive testing that accommodates real-valued scores while preserving the mathematical structure of binary CAT. Our key insight is to replace the Bernoulli response distribution with a heteroskedastic normal distribution that maintains the same logistic mean function while introducing variance that mimics the Bernoulli structure. This preserves the natural property that variance is highest when outcomes are most uncertain ( $\mu=0.5$ ) and shrinks at the boundaries where scores are constrained.

Continuous CAT already reduces evaluation cost substantially, but further savings are possible when we consider how evaluation results are typically used. Model evaluation is fundamentally comparative: the primary goal is often to determine whether one model is better than others. Building on continuous CAT, we introduce an adaptive multi-model ranking method that efficiently obtains statistically significant rankings by monitoring uncertainty estimates around model scores. Rather than test-

ing each model to a fixed precision threshold, we stop as soon as models are well-differentiated from each other based on pairwise confidence at a user-specified level. This focuses the testing effort on close competitors while reducing items for clearly separated models. We further observe that when distinguishing two models, testing either one yields information about their relative performance, allowing cost-aware allocation that preferentially tests cheaper models to maximize uncertainty reduction per dollar spent.

We validate our approach on five generation benchmarks that span summarization (GovReport, BioLaySumm2025), named entity recognition (Nemotron PII), question answering (TruthfulQA) and translation (FLORES). For each dataset, we evaluate 5 disjoint sets of 4 hold-out models and multiple metric types (ROUGE, BLEU, BERTScore, COMET, readability indices, LLM-as-judge). Our adaptive ranker achieves 0.73 Kendall’s  $\tau$  correlation with ground-truth rankings while using only 2% of the full evaluation budget, outperforming random sampling by 0.12  $\tau$ . Compared to fixed-length CAT, adaptive stopping provides an additional 32% item reduction and 42% cost savings.

Our contributions are:

1. A principled extension of IRT-based CAT from binary to continuous bounded outcomes via heteroskedastic normal distribution
2. Adaptive multi-model ranking algorithm with pairwise stopping and cost-aware allocation
3. Empirical validation across diverse generation tasks, metrics, and model scales

The continuous CAT framework opens adaptive testing to the full spectrum of modern LLM evaluation, and the adaptive multi-model ranker enables efficient comparisons of generation quality across model candidates. We make our code available on GitHub [link]

## 2 Related Work

Early studies on efficient model evaluation used IRT to analyze benchmark properties such as item difficulty and discrimination (Martínez-Plumed et al., 2016; Lalor et al., 2016; Vania et al., 2021; Rodriguez et al., 2021). Maia Polo et al. (2024) and Kipnis et al. (2025) used IRT to identify representative static subsets, while Hofmann et al. (2025)

developed a fully adaptive framework that dynamically selects items during evaluation. However, all existing approaches assume binary outcomes and cannot accommodate the continuous scores. Lalor et al. (2016) applied IRT to construct evaluation scales for NLP systems, and Prudêncio et al. (2015) and Lalor et al. (2019) showed that IRT models can be fit using response patterns from model ensembles rather than human annotations. Chen et al. (2019) proposed a Beta-distributed IRT model applied to modelling classifier confidence scores.

This follows the tradition in psychometrics, with Samejima (1973) introducing the continuous response model, and Noel and Dauvier (2007) proposing Beta-distributed models that naturally respect  $[0, 1]$  bounds. Both introduce additional complexity to Fisher Information calculations, complicating item selection. Beta models also cannot accommodate exact boundary values, requiring either ad hoc transformations or zero-and-one inflated extensions (Molenaar et al., 2022). Related continuous-response formulations include linear factor-analytic indices for difficulty and information, and Rasch-type models for continuous responses in large-scale learning systems (Ferrando, 2009; Deonovic et al., 2020). These typically assume homoskedastic residuals or prioritize scalable scoring over closed-form CAT information. Mellenbergh (1994) developed generalized linear IRT with simpler closed-form solutions, but assumes constant variance across the score range. Our heteroskedastic normal formulation combines simplicity with appropriate variance structure and maintains direct compatibility with standard CAT algorithms.

## 3 Background

**Item Response Theory (IRT)** (Hambleton et al., 1991) models the probability of a correct response as a function of the test-taker’s latent ability  $\theta$  and item parameters. The simplest model, the 1-Parameter Logistic (1PL) model, defines this probability as:

$$P(X = 1|\theta, a, b) = \frac{1}{1 + \exp(-a(\theta - b))} \quad (1)$$

where  $b$  is item difficulty and  $a$  is a discrimination parameter shared across all items. When ability matches difficulty ( $\theta = b$ ), the probability of success is exactly 0.5. Higher  $a$  yields steeper

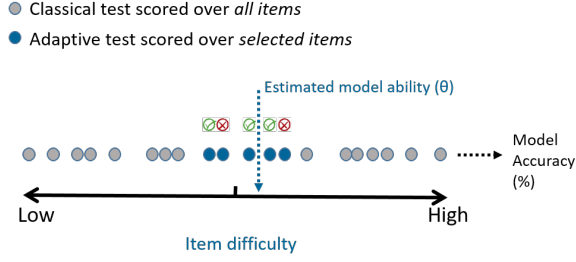


Figure 1: Adaptive testing focuses on items around the model ability, skipping those for which it would most certainly (i.e., uninformatively) get high or low scores.

response curves that better differentiate between ability levels.

**Computerized adaptive testing (CAT)** (Wainer, 2000) dynamically selects items to efficiently estimate a test-taker’s ability. The algorithm initializes with a prior distribution over  $\theta$ , typically  $\mathcal{N}(0, \sigma_0^2)$  where the prior variance is a hyperparameter controlling initial uncertainty. At each iteration, it selects the item that maximizes Fisher Information at the current ability estimate  $i^* = \arg \max_i I(\hat{\theta}|b_i)$ .

After observing the response, the ability estimate and posterior are updated via Bayesian updating. As items accumulate, the standard error decreases:

$$SE(\hat{\theta}) = \frac{1}{\sqrt{\sum_{i=1}^n I_i(\hat{\theta})}} \quad (2)$$

Testing terminates when this standard error falls below a predetermined threshold, yielding a precise ability estimate with fewer items than fixed-length testing.

**Fisher Information** quantifies how much an observation tells us about  $\theta$ . For the 1PL model with binary outcomes, this is derived directly from Bernoulli variance:

$$I(\theta|a, b) = a^2 \cdot P(\theta, b) \cdot (1 - P(\theta, b)) \quad (3)$$

Information is highest when  $P = 0.5$  (maximum uncertainty about the outcome) and approaches zero as  $P$  approaches 0 or 1 (outcome nearly certain). This creates heteroskedastic variance as a function of the predicted probability  $P(\theta, b)$ .

## 4 Continuous CAT

To extend CAT to continuous outcomes, we require a response distribution for scores in  $[0, 1]$  rather than binary  $\{0, 1\}$ . We retain the logistic mean

function from standard IRT, which captures the relationship between ability, difficulty, and expected performance, and replace the Bernoulli distribution with a normal distribution:

$$X|\theta, b, k \sim \mathcal{N}(\mu(\theta, b), \sigma^2(\theta, b)) \quad (4)$$

where the mean follows the logistic function:

$$\mu(\theta, b) = \frac{1}{1 + \exp(-(\theta - b))} \quad (5)$$

We preserve the Bernoulli variance structure  $\text{Var}(X) = P(1 - P)$  by defining:

$$\sigma^2(\theta, b) = k \cdot \mu(\theta, b) \cdot (1 - \mu(\theta, b)) \quad (6)$$

where  $k = 1/a^2$  is a noise parameter capturing measurement precision. We estimate a single  $k$  per dataset-metric combination from calibration data.

The Fisher Information for this model is:

$$I(\theta|b, k) = \frac{(d\mu/d\theta)^2}{\sigma^2} = \frac{\mu(1 - \mu)}{k} \quad (7)$$

Comparing with the 1PL Fisher Information  $I = a^2 \cdot P(1 - P)$ , we see that  $k = 1/a^2$ : high discrimination ( $a > 1$ ) corresponds to low noise ( $k < 1$ ), and vice versa. Both parameterizations capture the same underlying question: how reliably does an observed response reflect true ability? This equivalence indicates that our continuous extension constitutes a natural generalization that unifies the treatment of measurement precision across binary and continuous settings.

A limitation of the normal distribution is its unbounded support: it can technically assign probability to values outside  $[0, 1]$ . However, the heteroskedastic variance structure mitigates this concern: as  $\mu$  approaches the boundaries,  $\sigma^2 = k \cdot \mu(1 - \mu)$  shrinks, concentrating the distribution within the valid range. Alternative distributions such as the Beta are naturally bounded but do not preserve the connection to binary IRT and complicate the Fisher Information derivation. In practice, we find that the normal approximation performs well across the metrics we consider.

The CAT algorithm itself remains unchanged. Item selection still maximizes Fisher Information at the current ability estimate, with the information formula now incorporating  $k$ . Posterior updates follow the same Bayesian logic with normal likelihoods. The standard error formula  $SE(\hat{\theta}) = 1/\sqrt{\sum_i I_i}$  applies as before.

## 5 Parameter Estimation

We estimate both  $b_i$  and  $k$  from historical evaluation data by fitting the model to observed score distributions across a population of models.

**Item Difficulty.** We adopt the 1PL model and estimate item difficulties from observed scores across calibration models. Following standard IRT convention, we center the ability scale such that the mean model ability is zero. Under this parameterization, the maximum likelihood estimate for item difficulty simplifies to:

$$b_i = -\text{logit}(\hat{p}_i) = \log\left(\frac{1 - \hat{p}_i}{\hat{p}_i}\right) \quad (8)$$

where  $\hat{p}_i$  is the average score on item  $i$  across all calibration model-temperature configurations. When scores cluster in a narrow range (common for metrics like ROUGE where most models achieve 0.3–0.5), the resulting difficulties may not span the full ability range. We therefore apply min-max normalization to  $[\epsilon, 1 - \epsilon]$  before the logit transformation.

**Noise Parameter.** We estimate a global  $k$  from the calibration data using method-of-moments. Given item difficulties  $\{b_i\}$  and model ability estimates  $\{\theta_j\}$  (computed as  $\theta_j = \text{logit}(\bar{y}_j)$  where  $\bar{y}_j$  is the model’s average scaled score), we obtain:

$$k = \frac{\sum_{i,j} (y_{ij} - \mu_{ij})^2}{\sum_{i,j} \mu_{ij} (1 - \mu_{ij})} \quad (9)$$

where  $\mu_{ij} = \text{logit}^{-1}(\theta_j - b_i)$  is the predicted mean score for model  $j$  on item  $i$ . This formula directly estimates the variance inflation factor: when  $k = 1$ , observed variance matches the Bernoulli structure exactly;  $k > 1$  indicates noisier measurements than binary responses;  $k < 1$  indicates more precise measurements.

**Filtering Negative Discrimination Items.** Some items may anti-correlate with model ability. On these items, weaker models score higher than stronger models, rendering them effectively uninformative. Because we use a 1PL model with a global discrimination parameter, we filter such items by computing the Pearson correlation between item scores and model abilities across the calibration set, and excluding items with negative correlation from the adaptive item pool.

## 6 Adaptive Ranker

Standard CAT approaches evaluate each model independently to a fixed precision threshold, ignoring the comparative nature of model evaluation. This can be inefficient, since a model that clearly dominates or trails the field may receive the same number of test items as one locked in a close competition with a neighbor. We propose an adaptive multi-model ranking algorithm that directly optimizes for the comparative goal. Given  $M$  models to evaluate on a shared item bank, the algorithm runs independent CAT runs for each model, but monitors pairwise confidence in the ranking and terminates when all adjacent model pairs are statistically well-differentiated at a user-specified confidence level.

**Pairwise confidence.** Given ability estimates  $\hat{\theta}_i$  and  $\hat{\theta}_j$  with standard errors  $SE_i$  and  $SE_j$ , we compute the probability that model  $i$  outperforms model  $j$  under a normal approximation:

$$P(\theta_i > \theta_j) = \Phi\left(\frac{\hat{\theta}_i - \hat{\theta}_j}{\sqrt{SE_i^2 + SE_j^2}}\right) \quad (10)$$

where  $\Phi$  is the standard normal CDF. This quantifies our confidence in the pairwise ordering given current uncertainty in both estimates.

**Adaptive stopping.** Traditional CAT terminates when each model’s standard error falls below a threshold  $\epsilon$ , regardless of how the models compare. We instead terminate when all adjacent pairs in the current ranking satisfy  $P(\theta_i > \theta_j) > 1 - \frac{1-\gamma}{2}$  or  $P(\theta_i > \theta_j) < \frac{1-\gamma}{2}$ , corresponding to a two-sided confidence interval at level  $\gamma$ , or reach a maximum item budget  $n_{\max}$ . This focuses testing effort where it matters: models that are clearly separated require fewer items to establish their relative ordering, while close competitors receive additional testing until their difference reaches statistical significance or the budget is exhausted. If two models remain indistinguishable at the budget limit, they can be reported as tied rather than forced into an arbitrary ordering.

**Cost-aware allocation.** When a pair  $(i, j)$  requires additional testing to reach confidence threshold  $\gamma$ , we must choose which model to evaluate next. Since testing either model reduces uncertainty in the pairwise comparison, we select the

---

**Algorithm 1** Adaptive Multi-Model Ranking

---

**Require:** Models  $\{1, \dots, M\}$ , item bank  $\{(b_i, k_i)\}$ , costs  $\{c_m\}$ , confidence  $\gamma$ , max items  $n_{\max}$

- 1: Initialize  $\hat{\theta}_m \sim \mathcal{N}(0, 1)$ ,  $SE_m = 5$  for all  $m$
- 2: **Warm-up:** Administer  $n_{\text{init}}$  items to each model using MFI selection
- 3: **while** not all adjacent pairs confident or at max items **do**
- 4:   Rank models by  $\hat{\theta}_m$
- 5:   Identify uncertain pairs:  $\mathcal{U} = \{(i, j) : \gamma > P(\theta_i > \theta_j) > 1 - \gamma \text{ and } n_i < n_{\max} \text{ and } n_j < n_{\max}\}$
- 6:   **if**  $\mathcal{U} = \emptyset$  **then**
- 7:     **break**
- 8:   **end if**
- 9:   Collect candidate models:  $\mathcal{C} = \{m : m \in (i, j) \text{ for some } (i, j) \in \mathcal{U}\}$
- 10:   Select model:  $m^* = \arg \max_{m \in \mathcal{C}} \frac{SE_m^2}{(n_m + 1) \cdot c_m}$
- 11:   Select item:  $i^* = \arg \max_i I(\hat{\theta}_{m^*} | b_i, k_i)$
- 12:   Administer item  $i^*$  to model  $m^*$ , observe score  $y$
- 13:   Update  $\hat{\theta}_{m^*}$  and  $SE_{m^*}$
- 14: **end while**
- 15: **return** Ranking with pairwise confidence scores

---

342 model that provides the greatest expected uncer-  
343 tainty reduction per unit cost. We quantify this  
344 as:

345 
$$\text{value}_m = \frac{SE_m^2}{(n_m + 1) \cdot c_m} \quad (11)$$

346 where  $n_m$  is the number of items already admin-  
347 istered to model  $m$  and  $c_m$  is its per-item evaluation  
348 cost. The numerator reflects current uncertainty;  
349 the denominator captures diminishing returns (each  
350 additional item contributes less) weighted by cost.  
351 This criterion naturally allocates more items to  
352 cheaper models when resolving uncertain compar-  
353 isons, yielding additional cost savings beyond item  
354 reduction alone.

355 **Algorithm.** The complete procedure combines  
356 continuous CAT estimation with pairwise stopping  
357 and cost-aware allocation under a total budget con-  
358 straint. Rather than allocating a fixed maximum  
359 number of items per model, we specify a total cost  
360 budget  $B$  that can be flexibly distributed across  
361 models. After a warm-up phase that administers  
362 initial items to all models via MFI selection, the  
363 algorithm iteratively identifies uncertain pairs, se-  
364 lects the most cost-effective model to test, chooses  
365 the maximally informative item for that model, and  
366 updates posteriors. Testing continues until either  
367 all adjacent pairs reach the confidence threshold,  
368 or the total budget  $B$  is exhausted.

## 7 Experiments 369

We evaluate our continuous CAT extension and 370  
adaptive ranker across five generation tasks span- 371  
ning summarization, machine translation, question 372  
answering, and named entity recognition. Our goal 373  
is to approximate ground-truth rankings from full 374  
dataset evaluation using as few items as possible. 375  
We compare our adaptive ranker against random 376  
sampling, and examine whether the confidence es- 377  
timates can reliably tell apart genuine performance 378  
differences from statistical ties. 379

### 7.1 Experimental Protocol 380

We select tasks representing diverse generation set- 381  
tings and metric types (Table 1). This selection 382  
covers n-gram overlap (ROUGE, BLEU), learned 383  
embeddings (BERTScore, COMET), readability 384  
indices (FKGL), LLM-as-judge, and span-level 385  
F1. We exclude BERTScore from GovReport as 386  
all models achieved near-identical scores (range 387  
<0.02), yielding no meaningful ground-truth rank- 388  
ing. Full prompts and scoring details are provided 389  
in Appendix A. 390

We evaluate 21 models from 6 families: Ope- 391  
nAI GPT, Google Gemini, Amazon Nova, Meta 392  
Llama, Mistral, and Qwen (see Appendix A for 393  
full list). For each dataset-metric pair, we collect 394  
exhaustive evaluation scores by running all mod- 395  
els on all items, establishing ground-truth rankings. 396  
Each model is evaluated at four temperature set- 397  
tings ( $T \in \{0.0, 0.4, 0.7, 1.0\}$ ), which yields 84 398  
model-temperature configurations per item. Our 399

Table 1: Evaluation datasets. FKGL: Flesch-Kincaid Grade Level (Kincaid et al., 1975).

Dataset	Size	Split	Task	Metrics
BioLaySumm (Goldsack et al., 2024)	1,376	val	Lay summarization	ROUGE-L BERTScore FKGL
GovReport (Huang et al., 2021)	973	val	Document summarization	ROUGE-L
TruthfulQA (Lin et al., 2022)	817	val	Open-ended QA	LLM-Judge BERTScore (diff.)
FLORES (NLLB Team et al., 2022)	1,012	devtest	Translation (TR→EN)	COMET BLEU
Nemotron-PII (Steier et al., 2025)	2,000	test <sup>1</sup>	NER/PII detection	F1 (span-level)

main experiment uses five-fold cross-validation over models. In each fold, we randomly select 4 models as the hold-out set to be ranked adaptively; all temperature variants of hold-out models are excluded from calibration. The remaining 17 models (68 configurations) provide calibration data for estimating item parameters ( $b_i, k$ ). We fit parameters via maximum likelihood on the calibration set, then simulate adaptive evaluation on the hold-out models at  $T = 0$ . We set the confidence threshold  $\gamma = 0.95$  for pairwise decisions and require a minimum of 10 items per model before stopping.

For each hold-out set, we compare the adaptive ranker to a random sampling baseline which allocates the total budget  $B$  across models by iteratively selecting a random affordable model and a random unseen item to administer. We assess ranking quality using Kendall’s  $\tau$  correlation between predicted and ground-truth rankings computed from point estimates (mean ability  $\hat{\theta}$ ). For efficiency, we report the mean number of items administered per model.

## 7.2 Main Results

Table 2 presents results across all dataset-metric combinations. On 7 of 9 dataset-metric pairs, the adaptive ranker achieves higher Kendall’s  $\tau$  than the random sampling baseline. The largest gains occur on discriminative metrics: +12% on BioLaySumm ROUGE-L (0.957 vs 0.853), +26% on TruthfulQA BERTScore (0.450 vs 0.190), and +22% on FLORES BLEU (0.803 vs 0.580). These results demonstrate that our adaptive ranker provides genuine performance gains over random sampling. The adaptive ranker uses 52–101 items per hold-out set, and achieves 0.73  $\tau$  with the full dataset ranking while using only 2% of the items.

The discrimination parameter  $a$  predicts where adaptive methods provide the most benefit. On high-discrimination metrics ( $a > 3$ ), the adaptive ranker consistently outperforms baseline. On low-discrimination metrics like Nemotron F1 ( $a = 1.42$ ), individual items provide less information about ability, limiting the advantage of intelligent selection.

**Distributional Conformance.** Our continuous IRT extension assumes scores follow a heteroskedastic normal distribution with variance  $\sigma^2 = k \cdot \mu(1 - \mu)$ . We examine how well each metric conforms to this assumption by measuring the  $R^2$  between observed and predicted variance across score bins (see Appendix B for details). Conformance varies substantially: BERTScore and BLEU exhibit moderate fit ( $R^2 = 0.24$ – $0.36$ ), while ROUGE-L and COMET show poor conformance despite high ranking accuracy. We find no positive correlation between conformance and ranking quality ( $r = -0.12$ ); instead, discrimination  $a$  is the stronger predictor ( $r = 0.68$ ). This suggests that violations of the heteroskedastic variance assumption do not degrade ranking performance when items are sufficiently discriminative.

## 7.3 Tie Detection

The  $\tau$  correlation in Table 2 measures ranking accuracy based on point estimates, but point estimates alone do not indicate when a ranking is reliable. Reporting that model A outperforms model B when the difference is due to noise leads to non-reproducible results. Our method uses the IRT model’s posterior to detect statistical ties, where a pair is reported as a tie when confidence falls below 95%

Table 2: Ranking performance across benchmarks. Higher discrimination ( $a$ ) indicates more informative items. % Used is items administered relative to full evaluation.

Dataset	Metric	Size	$a$	Adaptive $\tau \uparrow$	Baseline $\tau \uparrow$	Items $\downarrow$	% Used $\downarrow$
BioLaySumm	ROUGE-L	1376	4.13	<b>0.957</b> $\pm$ 0.13	0.853 $\pm$ 0.20	82 $\pm$ 17	1.5 $\pm$ 0.3%
BioLaySumm	BERTScore	1376	3.40	<b>0.903</b> $\pm$ 0.17	0.743 $\pm$ 0.36	88 $\pm$ 9	1.6 $\pm$ 0.2%
BioLaySumm	FKGL	1376	2.34	<b>0.800</b> $\pm$ 0.25	0.713 $\pm$ 0.40	88 $\pm$ 15	1.6 $\pm$ 0.3%
GovReport	ROUGE-L	973	4.79	0.800 $\pm$ 0.20	<b>0.823</b> $\pm$ 0.22	98 $\pm$ 30	2.5 $\pm$ 0.8%
TruthfulQA	LLM-Judge	817	2.59	<b>0.490</b> $\pm$ 0.38	0.400 $\pm$ 0.49	93 $\pm$ 14	2.9 $\pm$ 0.4%
TruthfulQA	BERTScore	817	2.65	<b>0.450</b> $\pm$ 0.43	0.190 $\pm$ 0.46	92 $\pm$ 7	2.8 $\pm$ 0.2%
FLORES	BLEU	1012	3.12	<b>0.803</b> $\pm$ 0.23	0.580 $\pm$ 0.31	92 $\pm$ 8	2.3 $\pm$ 0.2%
FLORES	COMET	1012	4.07	<b>0.677</b> $\pm$ 0.33	0.503 $\pm$ 0.41	101 $\pm$ 11	2.5 $\pm$ 0.3%
Nemotron	F1	2000	1.36	0.673 $\pm$ 0.19	<b>0.707</b> $\pm$ 0.29	52 $\pm$ 5	0.7 $\pm$ 0.1%
<i>Overall (mean)</i>			3.16	<b>0.73</b>	0.61	88	2.0%

Table 3: Tie detection against ground-truth ties from bootstrap CIs on full evaluation. Adapt/GT Tie% are the fraction of pairs reported as ties for adaptive ranker and ground truth. Tie P/R/F1 are the precision, recall and F1 on tie detection quality of the adaptive ranker compared to the full evaluation. Conf Acc is the accuracy of the adaptive ranker on its non-tie predictions.

Dataset	Metric	Adapt Tie%	GT Tie%	Tie P $\uparrow$	Tie R $\uparrow$	Tie F1 $\uparrow$	Conf Acc $\uparrow$
BioLaySumm	ROUGE-L	14%	3%	0.60	1.00	0.62	1.00
BioLaySumm	BERTScore	42%	13%	0.36	0.99	0.43	1.00
BioLaySumm	FKGL	37%	17%	0.46	0.95	0.54	0.98
GovReport	ROUGE-L	18%	10%	0.39	0.73	0.43	0.94
TruthfulQA	LLM-Judge	76%	47%	0.50	0.89	0.60	0.77
TruthfulQA	BERTScore	100%	73%	0.73	1.00	0.82	–
FLORES	BLEU	58%	30%	0.48	0.99	0.58	1.00
FLORES	COMET	57%	37%	0.52	0.94	0.56	0.89
Nemotron	F1	65%	7%	0.08	1.00	0.14	1.00
<i>Overall (mean)</i>		52%	26%	0.46	0.94	0.52	0.95

471 Table 3 presents tie detection performance of  
472 the adaptive ranker against ground-truth ties de-  
473 rived from bootstrap confidence intervals on full  
474 evaluation. The adaptive ranker is consistently con-  
475 servative: Adapt Tie% exceeds GT Tie% for all  
476 metrics. Tie detection recall is high (0.77–1.00),  
477 which means that our method catches most ground-  
478 truth ties. For pairs where one model is ranked  
479 strictly higher than the other, overall accuracy is  
480 95%. TruthfulQA BERTScore reports 100% ties  
481 because all models score within a 0.7% range, and  
482 the differences are too small to resolve within the  
483 budget constraints.

484 Two metric-dataset combinations show lower  
485 performance: GovReport ROUGE-L (recall 0.73)  
486 and TruthfulQA LLM-Judge (confident accuracy  
487 0.77). In both cases, this can be traced to spe-  
488 cific model comparisons within single holdout sets

489 that have consistently biased estimates across seeds.  
490 For TruthfulQA LLM-Judge, one holdout contains  
491 three models that are tied in the ground-truth rank-  
492 ing, but the adaptive ranker consistently ranks one  
493 model below the others, detecting the ties in only  
494 1–3 of 20 seeds. For GovReport ROUGE-L, a  
495 single ground-truth tie in one holdout is never de-  
496 tected (0/20 seeds). Even with these outliers, the  
497 adaptive ranker achieves 95% overall accuracy on  
498 confidently ranked pairs and 94% recall in detect-  
499 ing ground truth ties, perfectly aligning with the  
500 95% confidence cutoff of the method.

#### 501 7.4 Ablation: Adaptive vs Fixed-Length CAT

502 The adaptive ranker terminates based on pairwise  
503 confidence, allocating more items to close com-  
504 petitors and fewer to well-separated models, and  
505 further manages cost by preferentially choosing

Table 4: Comparison with fixed-length CAT that administers the same number of items to all models.

Dataset	Metric	Adapt $\tau \uparrow$	Fixed $\tau \uparrow$	$\Delta\tau$	Item Savings	Cost Savings
BioLaySumm	BERTScore	0.90	0.82	-0.08	36%	46%
BioLaySumm	FKGL	0.80	0.78	-0.02	31%	38%
BioLaySumm	ROUGE-L	0.96	0.97	+0.02	39%	47%
FLORES	BLEU	0.80	0.78	-0.02	33%	43%
FLORES	COMET	0.68	0.80	+0.12	34%	44%
GovReport	ROUGE-L	0.80	0.85	+0.05	35%	46%
Nemotron	F1	0.68	0.76	+0.08	23%	31%
TruthfulQA	BERTScore	0.45	0.31	-0.14	28%	41%
TruthfulQA	LLM-Judge	0.49	0.52	+0.03	29%	41%
<i>Overall</i>		0.73	0.73	0.00	32%	42%

cheaper models. To isolate the value of this strategy, we compare against fixed-length CAT. For each holdout set, we administer the same number of items  $n$  to all models, where  $n$  is set to the maximum items used by the adaptive ranker for any model in that holdout.

Table 4 shows that the two methods achieve the same average  $\tau$ , but adaptive uses 32% less items and 42% less cost. This shows that the additional items used by the fixed-length CAT are superfluous, and the stopping mechanism leads to significant efficiency gains with no performance degradation.

## 7.5 Generalization to Unseen Model Families

Our main experiments use cross-validation over individual models, but calibration and hold-out models may share architectural similarities within families. To test whether item parameters generalize to genuinely novel architectures, we conduct a family-stratified evaluation: we hold out Mixtral-8x7b-instruct, Qwen3-32b, and Nova-prov1 and calibrating only on OpenAI, Meta/Llama, and Google/Gemini models (14 total).

Table 5 shows results. The adaptive ranker matches or outperforms baseline on 7 of 9 dataset-metric pairs, with substantial gains on TruthfulQA BERTScore (+60%) and GovReport ROUGE-L (+30%). These results suggest that item difficulty parameters transfer across model families: an item that is difficult for Llama models tends to also be difficult for Qwen or Nova models.

## 8 Conclusion

We present a principled extension of IRT-based adaptive testing from binary to continuous bounded scores. Building on this foundation, we intro-

Table 5: Family holdout evaluation. Calibration on OpenAI, Llama, Gemini (14 models); evaluation on hold-out Mistral, Qwen, Nova (3 models).

Dataset	Metric	Base $\tau \uparrow$	Adapt $\tau \uparrow$
BioLaySumm	BERTScore	0.53±0.52	<b>1.00±0.00</b>
BioLaySumm	FKGL	1.00±0.00	1.00±0.00
BioLaySumm	ROUGE-L	<b>1.00±0.00</b>	0.87±0.27
FLORES	BLEU	0.67±0.33	<b>1.00±0.00</b>
FLORES	COMET	0.73±0.33	0.73±0.33
GovReport	ROUGE-L	0.56±0.31	<b>0.85±0.28</b>
Nemotron	F1	<b>0.93±0.20</b>	0.67±0.33
TruthfulQA	BERTScore	0.27±0.55	<b>0.87±0.27</b>
TruthfulQA	LLM-Judge	0.40±0.55	<b>0.60±0.33</b>
<i>Overall</i>		0.68	<b>0.84</b>

duce an adaptive multi-model ranking algorithm with pairwise stopping and cost-aware allocation. Across five generation benchmarks spanning different metrics, our method achieves 0.73 Kendall’s  $\tau$  correlation with ground-truth rankings and 95% accuracy on confident predictions while using only 2% of items. Notably, our method can reliably rank models from entirely unseen families, and remains robust to empirical deviations from the assumed variance structure.

For future work, extending our to model item-specific discrimination parameters can improve item selection. Incorporating variable item costs can enable cost-optimal item selection at a more granular level. Extending to multi-metric evaluation can allow ranking models on several metrics simultaneously in a single adaptive run, further increasing evaluation efficiency. Developing methods for more efficient item parameter estimation such as transfer learning from related benchmarks or few-shot calibration can reduce the cold-start cost for new datasets and enable faster deployment.

## 9 Limitations

Our method requires calibration data from existing model evaluations to estimate item parameters, creating a cold-start problem for new benchmarks. This upfront cost of exhaustive evaluation on calibration models is amortized over subsequent adaptive evaluations, but represents a barrier to immediate deployment on new datasets. Furthermore, item parameters must be re-estimated for each metric separately, which might add overhead when comprehensive multi-metric evaluation is needed.

Our stopping criterion evaluates each adjacent pair independently at significance level  $\alpha$ , which controls the per-comparison error rate but not the family-wise error rate (FWER). With  $M$  models and  $M - 1$  adjacent comparisons, the probability of at least one incorrectly ordered pair under the null hypothesis is approximately  $1 - (1 - \alpha)^{M-1}$ , reaching 18% for five models at  $\alpha = 0.05$ . Corrections such as Holm-Bonferroni could provide FWER control by requiring stricter thresholds for the most confident pairs at the cost of additional test items, potentially negating much of the efficiency gain from adaptive selection.

## 10 Ethical Considerations

Large language model inference carries substantial environmental costs, with energy consumption and carbon emissions scaling with the number of API calls. By reducing evaluation from exhaustive testing to approximately 2% of items, our method directly decreases the environmental footprint of model comparison. As the number of candidate models, benchmarks, and evaluation metrics continues to grow, efficient evaluation methods become increasingly important for sustainable AI development.

Beyond environmental benefits, efficient evaluation has implications for equity in AI research. Exhaustive evaluation across large benchmarks favors well-resourced organizations. By dramatically reducing the number of required items, our approach can lower barriers for smaller laboratories, students and independent researchers to perform meaningful model comparisons. Additionally, our method’s explicit uncertainty quantification and tie detection encourages more honest reporting practices.

One potential concern is benchmark gaming: if certain items are predictably selected as most informative, model developers could optimize specifically for those items rather than general capability.

Periodic re-calibration of data and item banks can mitigate this risk.

## References

- Yu Chen, Telmo Silva Filho, Ricardo B Prudencio, Tom Diethe, and Peter Flach. 2019.  $\beta^3$ -irt: A new item response model and its applications. In *The 22nd international conference on artificial intelligence and statistics*, pages 1013–1021. PMLR.
- Benjamin Deonovic, Maria Bolsinova, Timo Bechger, and Gunter Maris. 2020. [A Rasch model and rating system for continuous responses collected in large-scale learning systems](#). *Frontiers in Psychology*, 11:500039.
- Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. The hitchhiker’s guide to testing statistical significance in natural language processing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1383–1392. Association for Computational Linguistics.
- Pere J Ferrando. 2009. [Difficulty, discrimination, and information indices in the linear factor analysis model for continuous responses](#). *Applied Psychological Measurement*, 33(1):9–24.
- Tomas Goldsack, Carolina Scarton, Matthew Shardlow, and Chenghua Lin. 2024. Overview of the BioLaySumm 2024 shared task on the lay summarization of biomedical research articles. In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, pages 122–131, Bangkok, Thailand. Association for Computational Linguistics.
- Ronald K Hambleton, Hariharan Swaminathan, and H Jane Rogers. 1991. *Fundamentals of Item Response Theory*. Sage Publications.
- Valentin Hofmann, David Heineman, Ian Magnusson, Kyle Lo, Jesse Dodge, Maarten Sap, Pang Wei Koh, Chun Wang, Hannaneh Hajishirzi, and Noah A Smith. 2025. [Fluid language model benchmarking](#). In *COLM*.
- Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. 2021. [Efficient attentions for long document summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1419–1436, Online. Association for Computational Linguistics.
- J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. 1975. Derivation of new readability formulas for Navy enlisted personnel. *Research Branch Report*, 8(75).
- Alex Kipnis, Konstantinos Voudouris, Luca M Schulze Buschoff, and Eric Schulz. 2025. [Metabench: A sparse benchmark of reasoning](#)

665	and knowledge in large language models. In	NLLB Team, Marta R Costa-jussà, James Cross, Onur	721
666	<i>Proceedings of the International Conference on</i>	Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hef-	722
667	<i>Learning Representations</i> .	ernan, Elahe Kalbassi, Janice Lam, Daniel Licht,	723
668	John P Lalor, Hao Wu, and Hong Yu. 2016. Building	Jean Maillard, Anna Sun, Skyler Wang, Guillaume	724
669	an evaluation scale using item response theory. In	Wenzek, Al Youngblood, Bapi Akula, Loic Barrault,	725
670	<i>Proceedings of the Conference on Empirical Meth-</i>	Gabriel Mejia Gonzalez, Prangthip Hansanti, and	726
671	<i>ods in Natural Language Processing</i> , pages 648–657.	20 others. 2022. No language left behind: Scaling	727
672	Association for Computational Linguistics.	human-centered machine translation. <i>arXiv preprint</i>	728
673	John P Lalor, Hao Wu, and Hong Yu. 2019. Learn-	<i>arXiv:2207.04672</i> .	729
674	ing latent parameters without human response pat-	Yvonnick Noel and Bruno Dauvier. 2007. <a href="#">A beta item</a>	730
675	terns: Item response theory with artificial crowds. In	<a href="#">response model for continuous bounded responses</a> .	731
676	<i>Proceedings of the 2019 Conference on Empirical</i>	<i>Applied Psychological Measurement</i> , 31(1):47–73.	732
677	<i>Methods in Natural Language Processing and the</i>	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	733
678	<i>9th International Joint Conference on Natural Lan-</i>	Jing Zhu. 2002. <a href="#">BLEU: a method for automatic eval-</a>	734
679	<i>guage Processing (EMNLP-IJCNLP)</i> , pages 4240–	<a href="#">uation of machine translation</a> . In <i>Proceedings of the</i>	735
680	4250. Association for Computational Linguistics.	<i>40th Annual Meeting of the Association for Compu-</i>	736
681	Chin-Yew Lin. 2004. ROUGE: A package for automatic	<i>tational Linguistics</i> , pages 311–318. Association for	737
682	evaluation of summaries. In <i>Text Summarization</i>	Computational Linguistics.	738
683	<i>Branches Out</i> , pages 74–81. Association for Compu-	RB Prudêncio, José Hernández-Orallo, and Adolfo	739
684	tational Linguistics.	Martinez-Usó. 2015. Analysis of instance hardness	740
685	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	in machine learning using item response theory. In	741
686	<a href="#">TruthfulQA: Measuring how models mimic human</a>	<i>Second international workshop on learning over mul-</i>	742
687	<a href="#">falsehoods</a> . In <i>Proceedings of the 60th Annual Meet-</i>	<i>tiplex contexts in ECML</i> , volume 3.	743
688	<i>ing of the Association for Computational Linguistics</i>	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	744
689	<i>(Volume 1: Long Papers)</i> , pages 3214–3252, Dublin,	Lavie. 2020. COMET: A neural framework for MT	745
690	Ireland. Association for Computational Linguistics.	evaluation. In <i>Proceedings of the 2020 Conference</i>	746
691	Qi Liu, Yan Zhuang, Haoyang Bi, Zhenya Huang,	<i>on Empirical Methods in Natural Language Pro-</i>	747
692	Weizhe Huang, Jiatong Li, Junhao Yu, Zirui Liu,	<i>cessing</i> , pages 2685–2702. Association for Computa-	748
693	Zirui Hu, Yuting Hong, Zachary A Pardos, Haip-	tional Linguistics.	749
694	ing Ma, Mengxiao Zhu, Shijin Wang, and Enhong	Pedro Rodriguez, Joe Barrow, Alexander Miserlis	750
695	Chen. 2024. Survey of computerized adaptive test-	Hoyle, John P Lalor, Robin Jia, and Jordan Boyd-	751
696	ing: A machine learning perspective. <i>arXiv preprint</i>	Graber. 2021. <a href="#">Evaluation examples are not equally</a>	752
697	<i>arXiv:2404.00712</i> .	<a href="#">informative: How should that change NLP leader-</a>	753
698	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang,	<a href="#">boards?</a> In <i>Proceedings of the 59th Annual Meet-</i>	754
699	Ruo Chen Xu, and Chenguang Zhu. 2023. G-eval:	<i>ing of the Association for Computational Linguistics</i>	755
700	NLG evaluation using GPT-4 with better human	<i>and the 11th International Joint Conference on Natu-</i>	756
701	alignment. <i>arXiv preprint arXiv:2303.16634</i> .	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	757
702	Felipe Maia Polo, Lucas Weber, Leshem Choshen,	pages 4486–4503. Association for Computational	758
703	Yuekai Sun, Gongjun Xu, and Mikhail Yurochkin.	Linguistics.	759
704	2024. tinyBenchmarks: evaluating LLMs with fewer	Fumiko Samejima. 1973. <a href="#">Homogeneous case of the con-</a>	760
705	examples. In <i>Proceedings of the 41st International</i>	<a href="#">tinuous response model</a> . <i>Psychometrika</i> , 38(2):203–	761
706	<i>Conference on Machine Learning</i> , volume 235, pages	219.	762
707	34303–34326. PMLR.	Amy Steier, Andre Manoel, Alexa Haushalter, and	763
708	Fernando Martínez-Plumed, Ricardo BC Prudêncio,	Maarten Van Segbroeck. 2025. <a href="#">Nemotron-PII: Syn-</a>	764
709	Adolfo Martínez-Usó, and José Hernández-Orallo.	<a href="#">thesized data for privacy-preserving AI</a> .	765
710	2016. Making sense of item response theory in ma-	Clara Vania, Phu Mon Htut, William Huang, Dhara	766
711	chine learning. In <i>Proceedings of the Twenty-second</i>	Mungra, Richard Yuanzhe Pang, Jason Phang,	767
712	<i>European Conference on Artificial Intelligence</i> , pages	Haokun Liu, Kyunghyun Cho, and Samuel Bowman.	768
713	1140–1148.	2021. Comparing test sets with item response theory.	769
714	Gideon J Mellenbergh. 1994. <a href="#">Generalized linear item</a>	In <i>Proceedings of the 59th Annual Meeting of the As-</i>	770
715	<a href="#">response theory</a> . <i>Psychological Bulletin</i> , 115(2):300–	<i>sociation for Computational Linguistics and the 11th</i>	771
716	307.	<i>International Joint Conference on Natural Language</i>	772
717	Dylan Molenaar, Mariana Curi, and Jorge L Bazán.	<i>Processing (Volume 1: Long Papers)</i> , pages 1141–	773
718	2022. <a href="#">Zero and one inflated item response theory</a>	1158. Association for Computational Linguistics.	774
719	<a href="#">models for bounded continuous data</a> . <i>Journal of Ed-</i>	Howard Wainer. 2000. <i>Computerized Adaptive Testing:</i>	775
720	<i>ucational and Behavioral Statistics</i> , 47(6):693–735.	<i>A Primer</i> , 2nd edition. Lawrence Erlbaum Associ-	776
		ates.	777

778	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-	pair using the devtest split, which contains 1,012	829
779	berger, and Yoav Artzi. 2020. BERTscore: Evalu-	sentence pairs. Each model was prompted with the	830
780	ating text generation with BERT. In <i>International</i>	instruction: “ <i>Translate the following Turkish text</i>	831
781	<i>Conference on Learning Representations</i> .	<i>to English. Output only the translation, nothing</i>	832
782	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	<i>else.”</i> followed by the source Turkish sentence.	833
783	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	This yielded 85,008 total inference responses	834
784	Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang,	(1,012 items $\times$ 21 models $\times$ 4 temperatures).	835
785	Joseph E Gonzalez, and Ion Stoica. 2024. Judging		
786	LLM-as-a-judge with MT-bench and chatbot arena.	<b>GovReport-Summarization.</b> We used the Gov-	836
787	In <i>Advances in Neural Information Processing Sys-</i>	Report dataset (Huang et al., 2021), a long-	837
788	<i>tems</i> , volume 36.	document summarization benchmark containing	838
789	<b>A Experimental Setup</b>	U.S. government reports paired with expert-written	839
790	All datasets were evaluated using the same set of 21	summaries. We evaluated on the validation split,	840
791	models spanning 6 families: OpenAI GPT (GPT-5-	which contains 973 documents. Each model was	841
792	mini, GPT-5-nano, GPT-4.1-mini, GPT-4.1-nano,	prompted with the instruction: “ <i>Summarize the fol-</i>	842
793	GPT-4o-mini), Google Gemini (Gemini-2.5-flash,	<i>lowing government report:”</i> followed by the full	843
794	Gemini-2.5-flash-lite, Gemini-2.0-flash, Gemini-	report text. This yielded 81,732 total inference	844
795	2.0-flash-lite), Amazon Nova (Nova-pro-v1, Nova-	responses (973 items $\times$ 21 models $\times$ 4 tempera-	845
796	lite-v1, Nova-micro-v1), Meta Llama (Llama-4-	tures).	846
797	maverick-17b, Llama-4-scout-17b, Llama-3.3-70b,	<b>Nemotron-PII.</b> We used the Nemotron-PII	847
798	Llama-3.2-3b, Llama-3.1-8b), Mistral (Mistral-	dataset (Steier et al., 2025), a named entity recog-	848
799	7b-instruct-v0.2, Mixtral-8x7b-instruct-v0.1), and	nition benchmark for detecting personally identifi-	849
800	Qwen (Qwen3-32b, Qwen3-coder-30b-a3b). Mod-	able information (PII) and protected health infor-	850
801	els were accessed via their native APIs (OpenAI,	mation (PHI) across 54 entity types. We evaluated	851
802	Google) or Amazon Bedrock (Nova, Llama, Mis-	on a random sample of 2,000 items from the test	852
803	tral, Qwen). Each model was queried at 4 temper-	split. Each model received a system prompt con-	853
804	ature settings: $T \in \{0.0, 0.4, 0.7, 1.0\}$ , result-	taining the complete list of valid entity labels and	854
805	ing in 84 model-temperature configurations per item.	4 few-shot examples demonstrating the expected	855
806	Note that GPT-5 models do not expose a tempera-	output format (EXTRACTED_TEXT   LABEL), fol-	856
807	ture parameter; for these models, the temperature	lowed by the document text. This yielded 168,000	857
808	setting was ignored, effectively treating all 4 con-	total inference responses (2,000 items $\times$ 21 models	858
809	figurations as identical.	$\times$ 4 temperatures). Responses were scored using	859
810	<b>BioLaySumm2025-PLOS.</b> We used the	standard span-level F1, where a predicted entity is	860
811	BioLaySumm2025-PLOS dataset (Goldsack	considered correct only if both the extracted text	861
812	et al., 2024), a biomedical lay summarization	and entity type exactly match a gold annotation.	862
813	benchmark containing scientific articles paired	<b>TruthfulQA.</b> We used the TruthfulQA	863
814	with plain-language summaries. We evaluated	dataset (Lin et al., 2022), a benchmark de-	864
815	on the validation split, which contains 1,376	signed to evaluate model truthfulness on questions	865
816	items. Each model was prompted with the	that commonly elicit false or misleading answers.	866
817	instruction: “ <i>Summarize the following scientific</i>	We evaluated on the validation split of the	867
818	<i>article in plain language for a general audience:”</i>	generation configuration, which contains 817	868
819	followed by the full article text. This yielded	questions. Each model was given a system	869
820	115,584 total inference responses (1,376 items $\times$	message “ <i>You are a helpful assistant.”</i> and	870
821	21 models $\times$ 4 temperatures). After scoring, all	prompted directly with the question. This yielded	871
822	metrics were linearly scaled to $[0, 1]$ based on the	68,628 total inference responses (817 items $\times$ 21	872
823	observed minimum and maximum values across all	models $\times$ 4 temperatures).	873
824	model-temperature configurations for each item.	Responses were scored using two metrics. First,	874
825	<b>FLORES-Turkish-English.</b> We used the	we computed a differential BERTScore using	875
826	FLORES-200 dataset (NLLB Team et al., 2022),	DeBERTa-xlarge-mnli embeddings: for each re-	876
827	a massively multilingual translation benchmark.	sponse, we calculate the maximum BERTScore F1	877
828	We evaluated on the Turkish-to-English language	against all correct answers and the maximum F1	878

879 against all incorrect answers, then compute the dif-  
 880 ferential as  $(\max_{\text{correct}} - \max_{\text{incorrect}} + 1)/2$ , nor-  
 881 malized to  $[0, 1]$ . This rewards responses seman-  
 882 tically closer to correct answers while penalizing  
 883 those closer to known misconceptions. Second, we  
 884 used LLM-as-judge with GPT-4.1-nano, employing  
 885 a 10-point rubric that penalizes answers matching  
 886 known incorrect responses and rewards alignment  
 887 with verified correct answers. Although common  
 888 practice favors strong judge models, preliminary  
 889 experiments showed high correlation between GPT-  
 890 4.1-nano, GPT-4.1-mini, and GPT-4.1 scores when  
 891 ground-truth correct and incorrect answer sets are  
 892 available, justifying the use of the more efficient  
 893 model.

All AI generated text was reviewed and corrected  
 by the authors.

914  
 915

## 894 B Conformance Analysis Details

895 Our continuous IRT extension assumes scores fol-  
 896 low a heteroskedastic normal distribution with vari-  
 897 ance  $\sigma^2 = k \cdot \mu(1 - \mu)$ , where  $\mu$  is the predicted  
 898 mean and  $k$  is the estimated noise parameter. To  
 899 examine how well each metric conforms to this as-  
 900 sumption, we bin items by their predicted mean  $\mu$   
 901 and compute the observed variance within each bin.  
 902 We then measure conformance as the  $R^2$  between  
 903 observed and predicted variance across bins.

904 Table 6 reports conformance alongside discrim-  
 905 ination and ranking quality. Notably, BioLay-  
 906 Summ ROUGE-L achieves the highest  $\tau$  (0.957)  
 907 with strongly negative  $R^2$  ( $-8.2$ ), whereas Bio-  
 908 LaySumm BERTScore has the best conformance  
 909 ( $R^2 = 0.36$ ) but lower  $\tau$  (0.903).

Table 6: Distributional conformance.  $R^2$  measures fit between observed and predicted variance across score bins.

Dataset	Metric	$a$	$R^2$	$\tau$
BioLaySumm	ROUGE-L	4.13	-8.2	0.957
BioLaySumm	BERTScore	3.40	0.36	0.903
BioLaySumm	FKGL	2.34	-666	0.800
GovReport	ROUGE-L	4.79	0.24	0.800
TruthfulQA	LLM-Judge	2.59	-1.2	0.490
TruthfulQA	BERTScore	2.65	-1524	0.450
FLORES	BLEU	3.12	0.34	0.803
FLORES	COMET	4.07	-20.3	0.677
Nemotron	F1	1.36	-20.5	0.673

## 910 C Use of AI Assistants

911 AI Assistants have been used in this paper to para-  
 912 phrase text written by the authors, and to generate  
 913 experimental details from code for Appendix A.