

# Curiosity-Driven LLM-as-a-judge for Personalized Creative Judgment

Vanya Bannihatti Kumar, Divyanshu Goyal, Akhil Eppa, Neel Bhandari

<sup>1</sup>Adobe Inc

## Abstract

Modern large language models (LLMs) excel at objective tasks such as evaluating mathematical reasoning and factual accuracy, yet they falter when faced with the nuanced, subjective nature of assessing creativity. In this work, we propose a novel curiosity-driven LLM-as-a-judge for evaluating creative writing which is personalized to each individual’s creative judgments. We use the Torrance Test of Creative Thinking (TTCW) benchmark introduced in [Chakrabarty et al. 2024], which has stories annotated by expert humans across various subjective dimensions like *Originality*, to test our hypothesis. We show that our method enables models across various sizes, to learn the nuanced creative judgments of different individuals, by showing improvements over baseline supervised finetuning (SFT) method across various evaluation metrics like Pearson correlation, Cohen’s  $\kappa$  and F1 values. Our method is especially useful in subjective evaluations where not all the annotators agree with each other.

## Introduction

Rigorous, standardized evaluation has repeatedly catalyzed progress in machine learning, from ImageNet[Russakovsky et al. 2015] and GLUE[Wang et al. 2019], driving leaps in the fields of computer vision and Natural Language Processing, respectively. The same effect is evident in objective math reasoning, where benchmarks like GSM8K[Cobbe et al. 2021], together with RL-trained reasoning models such as OpenAI’s o1[OpenAI et al. 2024] and DeepSeek-R1[DeepSeek-AI et al. 2025] have obtained strong results on hard contests like AIME and IMO.

While robust evaluation metrics exist for objective tasks such as mathematical reasoning and factual verification, subjective tasks like creativity remain difficult to assess reliably. There are many previous works [Panickssery, Bowman, and Feng 2024a, Wataoka, Takahashi, and Ri 2025] which show that using Large Language Models (LLM) as a judge prefer their own generations making them unreliable. Despite the success of LLMs on objective benchmarks, they still struggle to evaluate creativity in a manner aligned with human judgment. As shown in [Chakrabarty et al. 2024] and Table 12 and Table 2, even state-of-the-art models fall short in consistently evaluating the subjective dimensions of the story as

well as a human expert. This can be attributed to the fact that individual preferences shape creativity and rarely align uniformly across people.

To address this gap, we present an enhanced LLM-as-a-judge that not only learns from a diverse pool of annotations but also adapts its scoring to align with individual annotators or experts. This allows for more faithful and preference-aware evaluation of creativity. We emphasize personalization in our framework because the task of assessing subjective criteria is inherently variable across individuals. To this end, we propose a curiosity-driven LLM-as-a-judge for evaluating creativity in text generation, drawing inspiration from the curiosity-based Reinforcement Learning (RL) framework of [Pathak et al. 2017]. However, unlike the RL setting in [Pathak et al. 2017], we reinterpret curiosity as an *belief-shift signal* for creative evaluation. Specifically, when the model is “surprised” by an expert’s explanation, it signals a mismatch between the LLM’s prior belief and the expert’s preference; conversely, low surprise indicates alignment between the LLM and the expert (see Fig 5). To implement this, we first train an Intrinsic Curiosity Model (ICM) that measures the LLM’s surprise at a given explanation while simultaneously predicting which expert or annotator produced the explanation. The intuition behind predicting the annotator is that the model can learn which annotator caused the belief shift, allowing it to calibrate the curiosity signal for each annotator individually, thereby improving personalization. The resulting *curiosity score* is then fed as an auxiliary, self-supervised signal to improve a supervised fine-tuning (SFT) model (see Fig 1).

In our experiments, we establish a baseline using an SFT model that predicts annotators’ binary judgments from the story and question (see Fig 3a). To evaluate the effect of curiosity, we enhance this baseline with an ICM-derived curiosity score. More concretely we append the curiosity score to story and question in the baseline model. This helps us do a fair comparison on effect of curiosity signal on the final judgment and thereby measure the lift in performance our methodology provides over the baseline.

We conduct extensive experiments across various model sizes to ensure our method scales well with model size. Since the TTCW dataset size is extremely small, we do a 5-fold cross validation in order to ensure that our results are statistically significant. We also test our method in out-of-

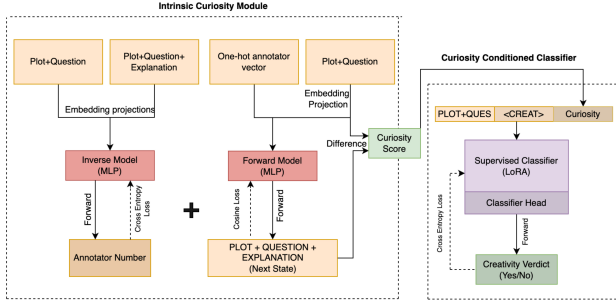


Figure 1: Overview of Architecture during training for Curiosity-Driven LLM-as-a-judge

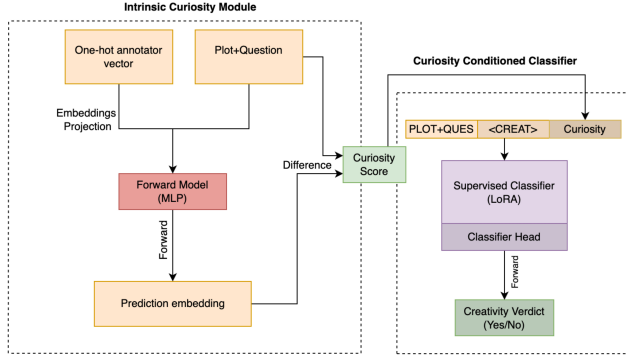


Figure 2: Overview of Architecture during inference for Curiosity Driven LLM-as-a-judge

distribution scenarios to ensure that our method generalizes well. Averaged across model sizes, ICM significantly improves Pearson correlation and F1 scores. More details about the results can be found in Fig 4.

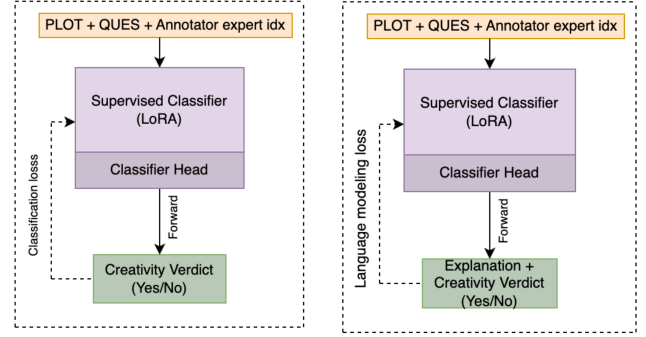
## Methodology

In this section, we describe our curiosity-driven LLM-as-a-judge for evaluating creativity in text generation, which combines belief shift estimation with expert attribution. Our method leverages the TTCW dataset [Chakrabarty et al. 2024], which is based on the Torrance Test of Creative Thinking [Torrance 1966] but adapted for LLMs. We focus on a subset of five creativity dimensions particularly relevant for evaluating the creative judgments of generative language models. We detail the dataset structure, model architecture, loss functions, and the formulation of our curiosity signal.

## Dataset

The TTCW dataset<sup>1</sup> provides expert human-annotated creativity judgments across 14 distinct dimensions. All the distinct dimensions in the TTCW dataset are mentioned in Appendix . For this study, we focus on five dimensions, 3 of which are categorised under Originality and 2 under flexibility: *Originality in Thought*, *Originality in Form*, *Originality in Theme and Content*, *Structural Flexibility*, and *Per-*

<sup>1</sup>Huggingface TTCW dataset



(a) Baseline without using explanations (b) Baseline using explanations

Figure 3: Comparison of baselines with and without using explanations.

*spective and Voice Flexibility*. Our analysis is restricted to these five dimensions, encompassing all dimensions under *Originality* and two representative dimensions from *Flexibility*. We picked these 5 dimensions among the 14 (Table 4) as these are more subjective in nature and hence the most ideal to evaluate our methodology. We defer exploration of the remaining dimensions to future work. Questions associated with each dimension can be found in appendix 6.

## Data Format and Task Setup

Each example in the dataset consists of a story  $S$ , a creativity-focused question  $Q_d$  specific to dimension  $d$ , an expert ID  $z_i$  where  $i \in \{1, 2, 3\}$  for each annotation by an expert, three expert-provided explanations  $\mathcal{E} = \{e_1, e_2, e_3\}$ , and corresponding binary verdicts  $V_i \in \{\text{yes}, \text{no}\}$  for each explanation.

The task is to improve the model’s performance on producing judgments similar to that of a particular expert when the model is presented with the story and the creative question

## Intrinsic Curiosity Model Overview

Our model operates in two stages:

1. **Belief Shift Estimation (Forward Score):** The model measures the impact of an expert explanation on their prediction of creativity.
2. **Expert Attribution (Backward Score):** The model identifies which expert wrote a given explanation.

**Forward Score: Belief Shift via Cosine Loss** We define two states:

- **State A:** Input consisting of the story and question and one-hot vector of the expert ID  $z_i$  represented as  $(S, Q_d, \text{onehot}(z_i))$  where  $i \in \{1, 2, 3\}$  as each story-question pair is annotated by 3 experts.
- **State B:** Input augmented with one expert explanation  $(S, Q_d, e_i)$  where  $i \in \{1, 2, 3\}$ .

Let  $f_{\theta}^{(A)} = f_{\theta}(S, Q_d, \text{onehot}(z_i))$  and  $f_{\theta}^{(B)} = f_{\theta}(S, Q_d, e_i)$ , where  $f_{\theta}$  denote the judge’s scoring function (logit head) with parameters  $\theta$  that maps the input to a scalar judgment logit.

The forward loss is defined as the cosine loss between these two predictions:

$$\mathcal{L}_{\text{forward}} = 1 - \frac{f_{\theta}^{(A)} \cdot f_{\theta}^{(B)}}{\|f_{\theta}^{(A)}\| \|f_{\theta}^{(B)}\|}$$

This loss captures how much the model’s belief about creativity of the story shifts when it incorporates the explanation by the annotator, which we define as the intrinsic curiosity measure.

### Backward Score: Expert Attribution via Cross-Entropy

To help the model to understand the distinct reasoning styles of different experts, we introduce an auxiliary classification task. Given  $(S, Q_d, e_i)$ , the model predicts the identity of the expert  $z_i \in \{1, 2, 3\}$  who authored explanation  $e_i$ :

$$p_{\phi}(z_i | S, Q_d, e_i) = \text{softmax}(g_{\phi}(S, Q_d, e_i))$$

The backward loss is the cross-entropy between the predicted and true expert label:

$$\mathcal{L}_{\text{backward}} = -\log p_{\phi}(z_i | S, Q_d, e_i)$$

**Loss function of Intrinsic curiosity model(ICM)** We define the ICM model’s loss as a weighted combination of the forward and backward components:

$$\mathcal{L}_{\text{curiosity}} = \mathcal{L}_{\text{forward}} + \lambda \cdot \mathcal{L}_{\text{backward}}$$

where  $\lambda$  is a tunable hyperparameter that balances the two objectives. In our experiments we set  $\lambda$  as 1.

**Incorporating the Curiosity Signal to SFT** To evaluate the utility of the learned curiosity signal, we use it as a conditioning input to a supervised fine-tuning (SFT) model trained to predict expert verdicts. For each instance, we append the scalar curiosity score to the original input using a special delimiter token <CREAT>, resulting in the following input format:

$$\text{Input: } Q_d + S + \text{<CREAT>} + \text{Curiosity}_{\text{Score}} \longrightarrow \text{Target: } V_i \quad (1)$$

$$\text{Curiosity}_{\text{score}} = f_{\theta}(S, Q_d, e_i) - f_{\theta}(S, Q_d, \text{onehot}(\text{expert\_idx})) \quad (2)$$

$V_i \in \{\text{yes}, \text{no}\}$  is the binary verdict associated with explanation  $e_i$ . The model uses the  $\text{Curiosity}_{\text{Score}}$  as a signal to predict the verdict of the given annotator. We use cross-entropy loss for training this classifier model

## Inference

During inference(see Fig 2), the story and creativity-focused questions are first passed through the intrinsic curiosity model (ICM) to compute a curiosity score. This score reflects the model’s internal belief shift in response to the input for that particular annotator. The resulting curiosity score is then appended to the original input, using a special delimiter token <CREAT>—and passed to the SFT classifier model. This classifier then predicts the binary creativity verdict (yes or no) for the given story-question pair. .

### Baseline with explanations

For the baseline comparison , we use a standard SFT model that produces the explanation and binary verdict given the input(see fig. 3b). The model input is structured as:

$$\text{Input: } Q_d + S + z_i \longrightarrow \text{Target: } \{V_i, e_i\}$$

At inference time, we provide  $Q_d$ ,  $S$ , and  $z_i$  as input, and the model outputs a JSON structure, from which the predicted verdict is parsed and compared to the ground truth. This baseline is trained using language modeling loss.

### Baseline without explanations

We ensure to compare our method against the baseline SFT in a classification setting rather than a causal language model setting to ensure fairness in comparison(see fig. 3a). Since we set up the baseline SFT in a classification setting, we do not include the explanations as neither part of the input or the output of the classification task. In this classification setting we use the question and the story as part of input and the verdict as part of the output.

$$\text{Input: } Q_d + S + z_i \longrightarrow \text{Target: } \{V_i\}$$

## Evaluation

Evaluating subjective tasks like creativity presents unique challenges, as even human annotators often disagree on what constitutes a "correct" judgment. Rather than attempting to define a universal metric for creativity, our approach embraces this subjectivity by focusing on personalization. We aim to adapt evaluation signals to individual experts by learning from a small number of their labeled examples. This allows us to model subjective preferences more faithfully and use this personalized model to assess creativity in a user-aligned manner. To quantify model performance in capturing individual judgments, we report **Pearson Correlation** [Benesty et al. 2009] and **Cohen’s  $\kappa$**  [Cohen 1960], along with **Precision**, **Recall**, and **F1-score**. These metrics enable us to assess both the predictive accuracy and ranking consistency of our models in aligning with subjective human evaluations.

### Theory: Why Curiosity Beats Using Explanation Text Directly

Let  $e$  denote the expert’s explanation,  $x = Q_d + S$ ,  $s_{\text{base}}(x) = f_{\theta}(S, Q_d, \text{onehot}(z_i))$  the pre-explanation

logit, and  $s_{\text{expl}}(x, e_i) = f_{\theta}(S, Q_d, e_i)$  the post-explanation logit produced by the model when conditioned on  $e$ . The  $\text{Curiosity}_{\text{Score}}$  is defined as the belief shift.

$$\text{Curiosity}_{\text{Score}} = f_{\theta}(S, Q_d, e_i) - f_{\theta}(S, Q_d, \text{onehot}(z_i)),$$

and *discard*  $e$  thereafter. We train a predictor  $\hat{p}_{\theta}(V=1 \mid x, \text{Curiosity}_{\text{Score}}) = \sigma(h_{\theta}(x, \text{Curiosity}_{\text{Score}}))$  where  $V$  is the verdict,  $h$  is the LLM judge model and  $\sigma$  represents softmax. This yields three advantages grounded in standard theory.

**(1) Weight-of-evidence sufficiency.** In logit/Bayesian updates, additional information acts *additively* on log-odds via a log-likelihood ratio (*weight of evidence*) [Agresti 2013]:

$$\log \frac{\Pr(V=1 \mid x, e_i)}{\Pr(V=0 \mid x, e_i)} = \log \frac{\Pr(V=1 \mid x)}{\Pr(V=0 \mid x)} + \underbrace{\log \frac{p(e \mid V=1, x)}{p(e \mid V=0, x)}}_{\text{weight of evidence}} \quad (3)$$

In our methodology,  $\text{Curiosity}_{\text{Score}} = s_{\text{expl}} - s_{\text{base}}$  is an *empirical estimate* of this increment on the log-odds scale, so it preserves the decision-relevant effect of  $e$  while removing lexical/style nuisance. Consequently, conditioning on  $\text{Curiosity}_{\text{Score}}$  approximates the theoretically “right” sufficient update in a logistic decision rule [Agresti 2013].

**(2) Variance reduction via a control-variate effect.** Let  $Z$  be the random quantity we wish to estimate more stably (e.g., per-example loss), and let  $C = \text{Curiosity}_{\text{Score}}$  be the control signal. With Pearson correlation

$$\rho = \text{Corr}(Z, C) = \frac{\text{Cov}(Z, C)}{\sqrt{\text{Var}(Z) \text{Var}(C)}} \in [-1, 1],$$

the classic control-variate construction implies that the optimally adjusted estimator  $Z^* = Z - \alpha^*(C - E[C])$  achieves

$$\text{Var}(Z^*) = \text{Var}(Z) (1 - \rho^2) \quad \text{at} \quad \alpha^* = \frac{\text{Cov}(Z, C)}{\text{Var}(C)}.$$

Thus any nonzero correlation with  $c$  strictly reduces variance [Owen 2013, Ch. 8]. Here,  $Z = \ell_i(\theta)$  (per-example cross-entropy loss) to reduce risk variance. Lower variance improves sample efficiency and stabilizes training.

**(3) Curiosity as a Model of Annotator Behaviour and Generalization** Subjective labels reflect both item difficulty and rater idiosyncrasy. A classic way to formalize this is a random-effects logit [Dawid and Skene 1979, Agresti 2013]:

$$\text{logit } \Pr(V=1 \mid x, z_i) = f(x) + b_{z_i}(x), \quad (4)$$

where  $f(x)$  captures item evidence and  $b_a(x)$  represents the (possibly context-dependent) strictness/leniency of annotator  $a$ . Since the curiosity score is able to model the annotator behaviour without considering the idiosyncrasies of the explanation text, it is able to better generalize to out-of-distribution dimensions for that annotator.

Table 1: ICM method results against the SFT baseline with explanations

Model	Exp.	LoRA $\alpha/R$	Pearson	Cohen’s $\kappa$	F1	Precision	Recall
Qwen0.5B	SFT	256/256	0.170 $\pm$ 0.049	0.155 $\pm$ 0.046	0.382 $\pm$ 0.049	0.452 $\pm$ 0.059	0.334 $\pm$ 0.060
	ICM	32/16	<b>0.524</b> $\pm$ 0.092	<b>0.383</b> $\pm$ 0.076	<b>0.616</b> $\pm$ 0.048	<b>0.494</b> $\pm$ 0.046	<b>0.818</b> $\pm$ 0.067
Qwen1.5B	SFT	256/256	0.170 $\pm$ 0.048	0.155 $\pm$ 0.048	0.402 $\pm$ 0.049	0.432 $\pm$ 0.020	0.383 $\pm$ 0.083
	ICM	32/16	<b>0.587</b> $\pm$ 0.061	<b>0.406</b> $\pm$ 0.065	<b>0.629</b> $\pm$ 0.045	<b>0.506</b> $\pm$ 0.045	<b>0.836</b> $\pm$ 0.056
Qwen3B	SFT	256/256	0.113 $\pm$ 0.083	0.110 $\pm$ 0.081	0.339 $\pm$ 0.051	0.401 $\pm$ 0.067	0.298 $\pm$ 0.060
	ICM	32/16	<b>0.540</b> $\pm$ 0.057	<b>0.356</b> $\pm$ 0.081	<b>0.598</b> $\pm$ 0.054	<b>0.481</b> $\pm$ 0.050	<b>0.794</b> $\pm$ 0.070
Qwen7B	SFT	128/128	0.160 $\pm$ 0.050	0.168 $\pm$ 0.085	0.371 $\pm$ 0.021	0.443 $\pm$ 0.050	0.324 $\pm$ 0.038
	ICM	32/16	<b>0.605</b> $\pm$ 0.083	<b>0.429</b> $\pm$ 0.082	<b>0.643</b> $\pm$ 0.053	<b>0.518</b> $\pm$ 0.051	<b>0.850</b> $\pm$ 0.072

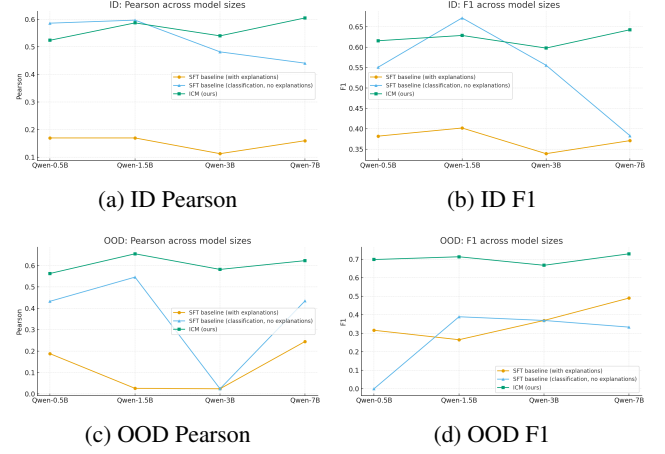


Figure 4: Three-way comparison across model sizes for ICM (ours), SFT baseline (classification, no explanations), and SFT baseline (with explanations). Panels show Pearson and F1 for in-distribution (top) and out-of-distribution (bottom). For exact results of the ID and OOD experiments of *baseline without explanation(classification)*, refer to Table 13 and Table 14

## Experiments

We evaluate our Intrinsic Curiosity Modeling (ICM) approach against a supervised fine-tuning (SFT) baseline (see Section ) across multiple model sizes. For a fair comparison in terms of identical input and outputs, we compare the ICM setup against SFT baseline with explanations. We also compare the ICM setup against FT baseline without explanations in order to ensure the same classification loss is used.

**Dataset** TTCW contains 48 stories annotated on 5 dimensions with three expert judgments per story–dimension pair, yielding 720 examples. We use 5-fold cross-validation with an 80/20 split, giving approximately 576 training and 144 test items per fold. Because individual folds are small, we report means across folds for all metrics (Table 1; see also Section ). Splits are stratified to preserve the positive/negative label ratio.

**Training setup.** The *baseline with explanations* uses a causal language modeling objective and our ICM model uses a classification objective. We align shared hyperparameters—learning rate, LoRA [Hu et al. 2022] rank, and batch size—wherever applicable to ensure compa-

Table 2: Comparison of ICM method against GPT-5 one-shot

Model	Exp.	Pearson	F1	Precision	Recall
Qwen0.5B	ICM	0.524 $\pm$ 0.092	0.616 $\pm$ 0.048	0.494 $\pm$ 0.046	0.818 $\pm$ 0.067
Qwen1.5B	ICM	0.587 $\pm$ 0.061	0.629 $\pm$ 0.045	0.506 $\pm$ 0.045	0.836 $\pm$ 0.056
Qwen3B	ICM	0.540 $\pm$ 0.057	0.598 $\pm$ 0.054	0.481 $\pm$ 0.050	0.794 $\pm$ 0.070
Qwen7B	ICM	0.605 $\pm$ 0.083	0.643 $\pm$ 0.053	0.518 $\pm$ 0.051	0.850 $\pm$ 0.072
GPT-5	ICM	0.2409 $\pm$ 0.1379	0.3467 $\pm$ 0.1592	0.5698 $\pm$ 0.2305	0.2608 $\pm$ 0.1378

rability. The ICM combined loss uses  $\lambda = 1$ . All fine-tuning (ICM and SFT baselines) uses LoRA; full details are in Table 5. For the *baseline without explanations*, which also uses a classification loss, we match all of the ICM hyperparameters.

**Compute and precision.** All runs use a single NVIDIA A100 (80 GB) GPU. Mixed precision with **bfloat16** is enabled when supported. When base models are loaded with 8-bit quantization, matrix multiplies in bitsandbytes execute in FP16 while LoRA heads operate in bfloat16.

**Convergence and reproducibility.** We train to loss convergence in all runs and fix random seeds for data splits and initialization. Hyperparameters and implementation details appear in Table 5.

## Analysis

### Effect of model scale

From Fig 4 we can see that our ICM method improves across model sizes whereas the *baseline classification method with no explanation* degrades with increase in model size for both ID and OOD settings. The reason why the *baseline classification method with no explanation* maybe degrading with scale is because this method primarily overfits on the small dataset with larger model sizes. Although the *baseline with explanation* improves with increase in model size, it remains uniformly low compared to the ICM method.

### Generalization

To understand the generalization ability of the baseline and the ICM models, we use the same setup as earlier but train the model in both methods on 4 dimensions - *Originality in Form*, *Originality in Theme and Content*, *Structural Flexibility*, and *Perspective and Voice Flexibility*, and test these trained models on the held out dimension of *Originality in Thought*. In this way there is absolutely no data leakage since the dimension the model is tested on was never seen during the training. From figure 4, we can see that gains of the ICM method over both the baseline methods are much more in the OOD settings rather than ID settings. This suggests the generalizability of our method because we are essentially allowing the model to understand the user behavior before predicting which is much more generalizable as compared to both baseline SFT methods.

### Comparison with GPT-5

Table 2 has the results of the ICM setup against GPT-5. We can see that even Qwen-0.5B model is able to beat GPT-

Table 3: ICM method results against the SFT baseline with explanations on Out-of-distribution data

Model	Experiment	LoRA $\alpha$ /Rank	Pearson	Cohen's $\kappa$	F1	Precision	Recall
Qwen0.5B	SFT	256/256	0.188	0.147	0.316	0.632	0.211
	ICM	32/16	<b>0.563</b>	<b>0.458</b>	<b>0.698</b>	<b>0.625</b>	<b>0.790</b>
Qwen1.5B	SFT	256/256	0.026	0.023	0.265	0.423	0.193
	ICM	32/16	<b>0.655</b>	<b>0.486</b>	<b>0.713</b>	<b>0.639</b>	<b>0.807</b>
Qwen3B	SFT	256/256	0.024	0.024	0.369	0.413	0.333
	ICM	32/16	<b>0.582</b>	<b>0.403</b>	<b>0.667</b>	<b>0.597</b>	<b>0.754</b>
Qwen7B	SFT	128/128	0.245	0.237	0.490	0.585	0.421
	ICM	32/16	<b>0.623</b>	<b>0.514</b>	<b>0.729</b>	<b>0.653</b>	<b>0.825</b>

5 model across all evaluation metrics except precision. The GPT-5 model was prompted with the same story, question and annotator index along with one shot example (randomly picked from training set) by the same annotator. GPT-5 model was more biased towards the answer "no" and whenever "yes" was predicted, it was almost always wrong. This further proves the effectiveness of our method.

## Conclusion and Future Work

We introduced a curiosity-driven LLM-as-a-judge for evaluating creativity in text generation, addressing the limitations of baseline SFT for inherently subjective tasks. Our approach leverages a two-part curiosity signal, capturing belief shifts via model responses to expert explanations and incorporating expert attribution through a backward prediction task. This signal enhances a SFT setup, leading to stronger alignment with human judgments across multiple creativity dimensions in the TTCW dataset. Experiments show that incorporating curiosity-based modeling consistently improves performance across model scales, surpassing standard SFT baselines in both correlation with human ratings and classification accuracy. Not only does it scale with model size, it also improves the performance in out-of-distribution scenarios, where we test the models on one heldout test dimension by training the models on the other 4 creativity dimension. Future work includes extending the curiosity-driven LLM-as-a-judge to other domains like marketing, evaluating novelty of scientific ideas etc.,. We also plan to use the curiosity signal as a reward signal in RL setup to further improve our current results.

## Literature Review

The evaluation of creativity in language models builds upon decades of work in creativity research, where the Torrance Tests of Creative Thinking (TTCT) assess fluency, flexibility, originality, and elaboration [Torrance 1966], and the Consensual Assessment Technique (CAT) uses aggregated expert judgments, a reliable but labour-intensive process [Patterson et al. 2024]. The authors of [Chakrabarty et al. 2024] adapted TTCT into the Torrance Tests for Creative Writing (TTCW), designing fourteen binary tests and enlisting creative-writing experts to evaluate 48 stories; their study showed that large language models pass these tests three to ten times less often than human writers [Chakrabarty et al. 2024], highlighting a sizable gap in creative competence. Alternative evaluation paradigms, such as the Leap-

of-Thought (LoT) framework for humorous, associative reasoning, argue that step-by-step chain-of-thought prompting can limit creativity and instead encourage non-sequential “leaps” [Zhong et al. 2024]. Efforts to automate creativity scoring (e.g., distributional-semantics proxies for novelty) often align weakly with expert judgments, reinforcing the need for human-aligned signals.

Because creativity judgments are *subjective*, collapsing rater perspectives via majority vote can erase systematic, meaningful disagreement. Following work on multi-annotator modeling, we treat annotators as distributions to be modeled rather than aggregated away [Mostafazadeh Davani, Díaz, and Prabhakaran 2022], rather than use the classical aggregation methods that infer a single latent “truth” [Whitehill et al. 2009, Hovy et al. 2013]. In parallel, recent results caution against naïve *LLM-as-judge* usage: evaluators can recognize and prefer their own generations, introducing self-preference bias [Panickssery, Bowman, and Feng 2024b]. Calibrated autoraters offer a partial mitigation via broad multi-task training and bias auditing [Vu et al. 2024]. These findings motivate rater-aware or human-anchored evaluation signals for creativity.

Intrinsic-motivation signals from reinforcement learning offer a principled lens on novelty seeking. Information-gain and prediction-error formulations—VIME [Houthoofd et al. 2017], ICM [Pathak et al. 2017], and Random Network Distillation [Burda et al. 2018]—are effective for exploration under sparse extrinsic reward. By analogy, curiosity-style signals can inform language evaluation by rewarding “useful novelty” (divergent yet coherent), complementing semantic-distance and rater-based methods. Our work instantiates this by modeling belief shifts when a language model incorporates expert explanations (a prediction-error-like signal) and combining it with expert attribution, yielding a more interpretable and *personalized* measure of creativity.

## References

- Agresti, A. 2013. *Categorical Data Analysis*. John Wiley & Sons, 3rd edition.
- Bai, J.; Bai, S.; Chu, Y.; Cui, Z.; Dang, K.; Deng, X.; Fan, Y.; Ge, W.; Han, Y.; Huang, F.; et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Bellemare, M. G.; Srinivasan, S.; Ostrovski, G.; Schaul, T.; Saxton, D.; and Munos, R. 2016. Unifying Count-Based Exploration and Intrinsic Motivation. *arXiv:1606.01868*.
- Benesty, J.; Chen, J.; Huang, Y.; and Cohen, I. 2009. Pearson correlation coefficient. In *Noise reduction in speech processing*, 1–4. Springer.
- Burda, Y.; Edwards, H.; Storkey, A.; and Klimov, O. 2018. Exploration by Random Network Distillation. *arXiv preprint arXiv:1810.12894*.
- Chakrabarty, T.; Laban, P.; Agarwal, D.; Muresan, S.; and Wu, C.-S. 2024. Art or Artifice? Large Language Models and the False Promise of Creativity. *arXiv:2309.14556*.
- Cobbe, K.; Kosaraju, V.; Bavarian, M.; Chen, M.; Jun, H.; Kaiser, L.; Plappert, M.; Tworek, J.; Hilton, J.; Nakano, R.; Hesse, C.; and Schulman, J. 2021. Training Verifiers to Solve Math Word Problems. *arXiv preprint arXiv:2110.14168*.
- Cohen, J. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1): 37–46.
- Dawid, A. P.; and Skene, A. 1979. Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm. *Journal of The Royal Statistical Society Series C-applied Statistics*, 28: 20–28.
- DeepSeek-AI; Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; Bi, X.; Zhang, X.; Yu, X.; Wu, Y.; Wu, Z. F.; Gou, Z.; Shao, Z.; Li, Z.; Gao, Z.; Liu, A.; Xue, B.; Wang, B.; Wu, B.; Feng, B.; Lu, C.; Zhao, C.; Deng, C.; Zhang, C.; Ruan, C.; Dai, D.; Chen, D.; Ji, D.; Li, E.; Lin, F.; Dai, F.; Luo, F.; Hao, G.; Chen, G.; Li, G.; Zhang, H.; Bao, H.; Xu, H.; Wang, H.; Ding, H.; Xin, H.; Gao, H.; Qu, H.; Li, H.; Guo, J.; Li, J.; Wang, J.; Chen, J.; Yuan, J.; Qiu, J.; Li, J.; Cai, J. L.; Ni, J.; Liang, J.; Chen, J.; Dong, K.; Hu, K.; Gao, K.; Guan, K.; Huang, K.; Yu, K.; Wang, L.; Zhang, L.; Zhao, L.; Wang, L.; Zhang, L.; Xu, L.; Xia, L.; Zhang, M.; Zhang, M.; Tang, M.; Li, M.; Wang, M.; Li, M.; Tian, N.; Huang, P.; Zhang, P.; Wang, Q.; Chen, Q.; Du, Q.; Ge, R.; Zhang, R.; Pan, R.; Wang, R.; Chen, R. J.; Jin, R. L.; Chen, R.; Lu, S.; Zhou, S.; Chen, S.; Ye, S.; Wang, S.; Yu, S.; Zhou, S.; Pan, S.; Li, S. S.; Zhou, S.; Wu, S.; Ye, S.; Yun, T.; Pei, T.; Sun, T.; Wang, T.; Zeng, W.; Zhao, W.; Liu, W.; Liang, W.; Gao, W.; Yu, W.; Zhang, W.; Xiao, W. L.; An, W.; Liu, X.; Wang, X.; Chen, X.; Nie, X.; Cheng, X.; Liu, X.; Xie, X.; Liu, X.; Yang, X.; Li, X.; Su, X.; Lin, X.; Li, X. Q.; Jin, X.; Shen, X.; Chen, X.; Sun, X.; Wang, X.; Song, X.; Zhou, X.; Wang, X.; Shan, X.; Li, Y. K.; Wang, Y. Q.; Wei, Y. X.; Zhang, Y.; Xu, Y.; Li, Y.; Zhao, Y.; Sun, Y.; Wang, Y.; Yu, Y.; Zhang, Y.; Shi, Y.; Xiong, Y.; He, Y.; Piao, Y.; Wang, Y.; Tan, Y.; Ma, Y.; Liu, Y.; Guo, Y.; Ou, Y.; Wang, Y.; Gong, Y.; Zou, Y.; He, Y.; Xiong, Y.; Luo, Y.; You, Y.; Liu, Y.; Zhou, Y.; Zhu, Y. X.; Xu, Y.; Huang, Y.; Li, Y.; Zheng, Y.; Zhu, Y.; Ma, Y.; Tang, Y.; Zha, Y.; Yan, Y.; Ren, Z. Z.; Ren, Z.; Sha, Z.; Fu, Z.; Xu, Z.; Xie, Z.; Zhang, Z.; Hao, Z.; Ma, Z.; Yan, Z.; Wu, Z.; Gu, Z.; Zhu, Z.; Liu, Z.; Li, Z.; Xie, Z.; Song, Z.; Pan, Z.; Huang, Z.; Xu, Z.; Zhang, Z.; and Zhang, Z. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *arXiv:2501.12948*.
- Fisch, A.; Eisenstein, J.; Zayats, V.; Agarwal, A.; Beirami, A.; Nagpal, C.; Shaw, P.; and Berant, J. 2025. Robust Preference Optimization through Reward Model Distillation. *arXiv:2405.19316*.
- Guilford, J. P. 1967. Creativity: Yesterday, Today and Tomorrow. *Journal of Creative Behavior*, 1: 3–14.
- Houthoofd, R.; Chen, X.; Duan, Y.; Schulman, J.; Turck, F. D.; and Abbeel, P. 2017. VIME: Variational Information Maximizing Exploration. *arXiv:1605.09674*.
- Hovy, D.; Berg-Kirkpatrick, T.; Vaswani, A.; and Hovy, E. 2013. Learning Whom to Trust with MACE. In Vanderwende, L.; Daumé III, H.; and Kirchhoff, K., eds., *Proceedings of the 2013 Conference of the North American Chap-*

ter of the Association for Computational Linguistics: Human Language Technologies, 1120–1130. Atlanta, Georgia: Association for Computational Linguistics.

Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; Chen, W.; et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2): 3.

Jiang, D.; Ren, X.; and Lin, B. Y. 2023. LLM-Blender: Ensembling Large Language Models with Pairwise Ranking and Generative Fusion. *arXiv:2306.02561*.

Madotto, A.; Namazifar, M.; Huizinga, J.; Molino, P.; Ecoffet, A.; Zheng, H.; Papangelis, A.; Yu, D.; Khatri, C.; and Tur, G. 2020. Exploration Based Language Learning for Text-Based Games. *arXiv:2001.08868*.

Mostafazadeh Davani, A.; Díaz, M.; and Prabhakaran, V. 2022. Dealing with Disagreements: Looking Beyond the Majority Vote in Subjective Annotations. *Transactions of the Association for Computational Linguistics*, 10: 92–110.

Nath, A.; Jung, C.; Seefried, E.; and Krishnaswamy, N. 2025. Simultaneous Reward Distillation and Preference Learning: Get You a Language Model Who Can Do Both. *arXiv:2410.08458*.

OpenAI; ; Jaech, A.; Kalai, A.; Lerer, A.; Richardson, A.; El-Kishky, A.; Low, A.; Helyar, A.; Madry, A.; Beutel, A.; Carney, A.; Iftimie, A.; Karpenko, A.; Passos, A. T.; Neitz, A.; Prokofiev, A.; Wei, A.; Tam, A.; Bennett, A.; Kumar, A.; Saraiva, A.; Vallone, A.; Duberstein, A.; Kondrich, A.; Mishchenko, A.; Applebaum, A.; Jiang, A.; Nair, A.; Zoph, B.; Ghorbani, B.; Rossen, B.; Sokolowsky, B.; Barak, B.; McGrew, B.; Minaiev, B.; Hao, B.; Baker, B.; Houghton, B.; McKinzie, B.; Eastman, B.; Lugaresi, C.; Bassin, C.; Hudson, C.; Li, C. M.; de Bourcy, C.; Voss, C.; Shen, C.; Zhang, C.; Koch, C.; Orsinger, C.; Hesse, C.; Fischer, C.; Chan, C.; Roberts, D.; Kappler, D.; Levy, D.; Selsam, D.; Dohan, D.; Farhi, D.; Mely, D.; Robinson, D.; Tsipras, D.; Li, D.; Oprica, D.; Freeman, E.; Zhang, E.; Wong, E.; Proehl, E.; Cheung, E.; Mitchell, E.; Wallace, E.; Ritter, E.; Mays, E.; Wang, F.; Such, F. P.; Raso, F.; Leoni, F.; Tsimpourlas, F.; Song, F.; von Lohmann, F.; Sulit, F.; Salmon, G.; Parascandolo, G.; Chabot, G.; Zhao, G.; Brockman, G.; Leclerc, G.; Salman, H.; Bao, H.; Sheng, H.; Andrin, H.; Bagherinezhad, H.; Ren, H.; Lightman, H.; Chung, H. W.; Kivlichan, I.; O’Connell, I.; Osband, I.; Gilaberte, I. C.; Akkaya, I.; Kostrikov, I.; Sutskever, I.; Kofman, I.; Pachocki, J.; Lennon, J.; Wei, J.; Harb, J.; Twore, J.; Feng, J.; Yu, J.; Weng, J.; Tang, J.; Yu, J.; Candela, J. Q.; Palermo, J.; Parish, J.; Heidecke, J.; Hallman, J.; Rizzo, J.; Gordon, J.; Uesato, J.; Ward, J.; Huizinga, J.; Wang, J.; Chen, K.; Xiao, K.; Singhal, K.; Nguyen, K.; Cobbe, K.; Shi, K.; Wood, K.; Rimbach, K.; Gu-Lemberg, K.; Liu, K.; Lu, K.; Stone, K.; Yu, K.; Ahmad, L.; Yang, L.; Liu, L.; Maksin, L.; Ho, L.; Fedus, L.; Weng, L.; Li, L.; McCallum, L.; Held, L.; Kuhn, L.; Kondraciuk, L.; Kaiser, L.; Metz, L.; Boyd, M.; Trebacz, M.; Joglekar, M.; Chen, M.; Tintor, M.; Meyer, M.; Jones, M.; Kaufer, M.; Schwarzer, M.; Shah, M.; Yatbaz, M.; Guan, M. Y.; Xu, M.; Yan, M.; Glaese, M.; Chen, M.; Lampe, M.; Malek, M.; Wang, M.; Fradin, M.; McClay, M.; Pavlov, M.; Wang, M.; Wang, M.; Murati, M.; Bavarian, M.;

Rohaninejad, M.; McAleese, N.; Chowdhury, N.; Chowdhury, N.; Ryder, N.; Tezak, N.; Brown, N.; Nachum, O.; Boiko, O.; Murk, O.; Watkins, O.; Chao, P.; Ashbourne, P.; Izmailov, P.; Zhokhov, P.; Dias, R.; Arora, R.; Lin, R.; Lopes, R. G.; Gaon, R.; Miyara, R.; Leike, R.; Hwang, R.; Garg, R.; Brown, R.; James, R.; Shu, R.; Cheu, R.; Greene, R.; Jain, S.; Altman, S.; Toizer, S.; Toyer, S.; Miserendino, S.; Agarwal, S.; Hernandez, S.; Baker, S.; McKinney, S.; Yan, S.; Zhao, S.; Hu, S.; Santurkar, S.; Chaudhuri, S. R.; Zhang, S.; Fu, S.; Papay, S.; Lin, S.; Balaji, S.; Sanjeev, S.; Sidor, S.; Broda, T.; Clark, A.; Wang, T.; Gordon, T.; Sanders, T.; Patwardhan, T.; Sottiaux, T.; Degry, T.; Dimson, T.; Zheng, T.; Garipov, T.; Stasi, T.; Bansal, T.; Creech, T.; Peterson, T.; Eloundou, T.; Qi, V.; Kosaraju, V.; Monaco, V.; Pong, V.; Fomenko, V.; Zheng, W.; Zhou, W.; McCabe, W.; Zaremba, W.; Dubois, Y.; Lu, Y.; Chen, Y.; Cha, Y.; Bai, Y.; He, Y.; Zhang, Y.; Wang, Y.; Shao, Z.; and Li, Z. 2024. OpenAI o1 System Card. *arXiv:2412.16720*.

Owen, A. B. 2013. Monte Carlo theory, methods and examples.

Pan, S. 2025. Tiny Reward Models. *arXiv:2507.09973*.

Panickssery, A.; Bowman, S.; and Feng, S. 2024a. Llm evaluators recognize and favor their own generations. *Advances in Neural Information Processing Systems*, 37: 68772–68802.

Panickssery, A.; Bowman, S. R.; and Feng, S. 2024b. LLM Evaluators Recognize and Favor Their Own Generations. *arXiv:2404.13076*.

Pathak, D.; Agrawal, P.; Efros, A. A.; and Darrell, T. 2017. Curiosity-driven Exploration by Self-supervised Prediction. *arXiv:1705.05363*.

Patterson, J. D.; Barbot, B.; Lloyd-Cox, J.; and Beaty, R. E. 2024. AuDrA: An automated drawing assessment platform for evaluating creativity. *Behavior Research Methods*, 56(4): 3619–3636.

Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A. C.; and Fei-Fei, L. 2015. ImageNet Large Scale Visual Recognition Challenge. *arXiv:1409.0575*.

Schmidhuber, J. 2010. Formal Theory of Creativity, Fun, and Intrinsic Motivation. *IEEE Trans. on Auton. Ment. Dev.*, 2(3): 230–247.

Sidahmed, H.; Phatale, S.; Hutcheson, A.; Lin, Z.; Chen, Z.; Yu, Z.; Jin, J.; Chaudhary, S.; Komarytsia, R.; Ahlheim, C.; Zhu, Y.; Li, B.; Ganesh, S.; Byrne, B.; Hoffmann, J.; Mansoor, H.; Li, W.; Rastogi, A.; and Dixon, L. 2024. Parameter Efficient Reinforcement Learning from Human Feedback. *arXiv:2403.10704*.

Torrance, E. P. 1966. *Torrance Tests of Creative Thinking: Norms–Technical Manual (Research Edition)*. Princeton, NJ: Personnel Press.

Vu, T.; Krishna, K.; Alzubi, S.; Tar, C.; Faruqui, M.; and Sung, Y.-H. 2024. Foundational Autraters: Taming Large Language Models for Better Automatic Evaluation. *arXiv:2407.10817*.



Wan, Y.; Wu, J.; Abdulhai, M.; Shani, L.; and Jaques, N. 2025. Enhancing Personalized Multi-Turn Dialogue with Curiosity Reward. arXiv:2504.03206.

Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2019. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. arXiv:1804.07461.

Wang, X.; Wei, J.; Schuurmans, D.; Le, Q.; Chi, E.; Narang, S.; Chowdhery, A.; and Zhou, D. 2023. Self-Consistency Improves Chain of Thought Reasoning in Language Models. arXiv:2203.11171.

Wataoka, K.; Takahashi, T.; and Ri, R. 2025. Self-Preference Bias in LLM-as-a-Judge. arXiv:2410.21819.

Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Ichter, B.; Xia, F.; Chi, E.; Le, Q.; and Zhou, D. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv:2201.11903.

Whitehill, J.; Wu, T.-f.; Bergsma, J.; Movellan, J.; and Ruvo, P. 2009. Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise. In Bengio, Y.; Schuurmans, D.; Lafferty, J.; Williams, C.; and Culotta, A., eds., *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc.

Zar, J. H. 2005. Spearman rank correlation. *Encyclopedia of Biostatistics*, 7.

Zhang, Y.; Wang, L.; Fang, M.; Du, Y.; Huang, C.; Wang, J.; Lin, Q.; Pechenizkiy, M.; Zhang, D.; Rajmohan, S.; and Zhang, Q. 2025. Distill Not Only Data but Also Rewards: Can Smaller Language Models Surpass Larger Ones? arXiv:2502.19557.

Zhao, Y.; Zhang, R.; Li, W.; and Li, L. 2025. Assessing and Understanding Creativity in Large Language Models. *Machine Intelligence Research*, 22(3): 417–436.

Zhong, S.; Huang, Z.; Gao, S.; Wen, W.; Lin, L.; Zitnik, M.; and Zhou, P. 2024. Let’s Think Outside the Box: Exploring Leap-of-Thought in Large Language Models with Creative Humor Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 13246–13257.

## Appendix

### Dimensions in dataset

In Table 4, all the dimensions that are part of the TTCW dataset are mentioned.

### More experiment and compute details

#### Limitations

Our study has some limitations that we hope to address in future work. First, the empirical scope is narrow: we evaluate only on TTCW dataset. Our current method is text-only; extending to richer modalities and subjective tasks beyond TTCW remains future work. In addition, the dataset is small (48 stories  $\times$  5 dimensions with three expert judgments per story–dimension, totaling 720 instances). We therefore rely on 5-fold cross-validation and report means and deviation

Table 4: Dimensions of TTCW dataset

Dimension	Facets
Fluency	Understandability & Coherence
	Narrative Pacing
	Scene vs Exposition
	Literary Devices & Language Proficiency
	Narrative Ending
Flexibility	Emotional Flexibility
	Perspective & Voice Flexibility
	Structural Flexibility
Originality	Originality in Form
	Originality in Thought
	Originality in Theme & Content
Elaboration	World Building & Setting
	Character Development
	Rhetorical Complexity

Table 5: Core hyperparameters used in all runs.

max_length	4096
lora_dropout	0.1
target_modules	["q_proj", "k_proj", "v_proj", "o_proj", "gate_proj", "up_proj", "down_proj"]
lr_scheduler	cosine (warmup_ratio = 0.1)
per_device_train_batch_size	4
gradient_accumulation_steps	8
weight_decay	0.01
max_grad_norm	0.5
num_train_epochs	3
seed	42

across 5 folds. Finally, model coverage is limited to one family (Qwen2.5 0.5B–7B), leaving generalization across architectures untested, which we aim to do in future work.



## Question for each dimension

Table 6: Creativity evaluation categories and questions

Category	Question
Originality in Thought	Is the story an original piece of writing without any cliches?
Originality in Form and Structure	Does the story show originality in its form and/or structure?
Originality in Theme and Content	Will an average reader of this story obtain a unique and original idea from reading it?
Perspective and Voice Flexibility	Does the story provide diverse perspectives, and if there are unlikeable characters, are their perspectives presented convincingly and accurately?
Structural Flexibility	Does the story contain turns that are both surprising and appropriate?

## Statistical significance testing

Table 7: Statistical significance test across 5 folds for Qwen-0.5b model

Metric	SFT(with expl) (mean±SD)	ICM (mean±SD)	$\Delta$ (ICM−SFT)	$p$ (paired $t$ )	Statistically significant?
Pearson	0.160 ± 0.055	0.524 ± 0.092	0.364	0.002	Yes
Spearman	0.160 ± 0.055	0.484 ± 0.078	0.324	< 0.001	Yes
F1	0.371 ± 0.054	0.616 ± 0.048	0.245	< 0.001	Yes

Table 10: Statistical significance test across 5 folds for Qwen-7b model.

Metric	SFT(with expl) (mean±SD)	ICM (mean±SD)	$\Delta$ (ICM−SFT)	$p$ (paired $t$ )	Statistically significant?
Pearson	0.170 ± 0.058	0.606 ± 0.084	0.436	0.002	Yes
Spearman	0.170 ± 0.058	0.542 ± 0.089	0.372	< 0.001	Yes
F1	0.381 ± 0.029	0.663 ± 0.058	0.282	< 0.001	Yes

Table 8: Statistical significance test across 5 folds for Qwen-1.5b model

Metric	SFT(with expl) (mean±SD)	ICM (mean±SD)	$\Delta$ (ICM−SFT)	$p$ (paired $t$ )	Statistically significant?
Pearson	0.170 ± 0.058	0.586 ± 0.064	0.416	< 0.001	Yes
Spearman	0.170 ± 0.058	0.522 ± 0.069	0.352	< 0.001	Yes
F1	0.402 ± 0.050	0.629 ± 0.045	0.227	< 0.001	Yes

Table 9: Statistical significance test across 5 folds for Qwen-3b model.

Metric	SFT(with expl) (mean±SD)	ICM (mean±SD)	$\Delta$ (ICM−SFT)	$p$ (paired $t$ )	Statistically significant?
Pearson	0.113 ± 0.092	0.540 ± 0.074	0.427	< 0.001	Yes
Spearman	0.113 ± 0.092	0.494 ± 0.091	0.381	< 0.001	Yes
F1	0.339 ± 0.053	0.618 ± 0.061	0.279	< 0.001	Yes

Table 11: Average passing rate (%) on individual TTCW, based on annotations of 10 creative writing experts across 48 stories; last column reports Fleiss’  $\kappa$  (expert agreement).

Dimension	Test	GPT-3.5	GPT-4	Claude v1.3	New Yorker	Expert $\kappa$
Fluency	Understandability & Coherence	22.2	33.3	55.6	91.7	0.27
	Narrative Pacing	8.3	52.8	61.1	94.4	0.39
	Scene vs Exposition	8.3	50.0	58.3	91.7	0.27
	Literary Devices & Language	5.6	36.1	13.9	88.9	0.37
	Narrative Ending	8.3	19.4	33.3	91.7	0.48
Flexibility	Emotional Flexibility	16.7	19.4	36.1	91.7	0.32
	Perspective & Voice Flexibility	8.3	16.7	19.4	72.2	0.44
	Structural Flexibility	11.1	19.4	30.6	88.9	0.39
Originality	Originality in Form	2.8	8.3	0.0	63.9	0.41
	Originality in Thought	2.8	44.4	19.4	91.7	0.40
	Originality in Theme & Content	0.0	19.4	11.1	75.0	0.66
Elaboration	World Building & Setting	16.7	41.7	58.3	94.4	0.33
	Character Development	8.3	16.7	16.7	61.1	0.31
	Rhetorical Complexity	2.8	11.1	5.6	88.9	0.66
<b>Average</b>		<b>8.7</b>	<b>27.9</b>	<b>30.0</b>	<b>84.7</b>	<b>0.41</b>

Table 12: Correlation between LLM-administered TTCW and expert annotations (Cohen’s  $\kappa$ ) on all 48 stories.

Dimension	Test	GPT-3.5	GPT-4	Claude
Fluency	Understandability & Coherence	-0.01	-0.01	-0.17
	Narrative Pacing	0.05	0.00	-0.22
	Scene vs Exposition	-0.03	-0.08	-0.23
	Literary Devices & Language	0.04	-0.09	-0.11
	Narrative Ending	-0.02	0.02	0.02
Flexibility	Emotional Flexibility	-0.04	0.00	0.09
	Perspective & Voice	0.00	0.26	0.14
	Structural Flexibility	-0.04	0.00	-0.07
Originality	Originality in Form	0.08	0.09	0.03
	Originality in Thought	0.19	0.31	0.15
	Originality in Theme & Content	0.06	-0.01	0.18
Elaboration	World Building & Setting	0.00	0.00	0.09
	Character Development	-0.08	0.02	0.00
	Rhetorical Complexity	0.00	0.00	0.02
<b>Average</b>		<b>0.016</b>	<b>0.035</b>	<b>-0.006</b>

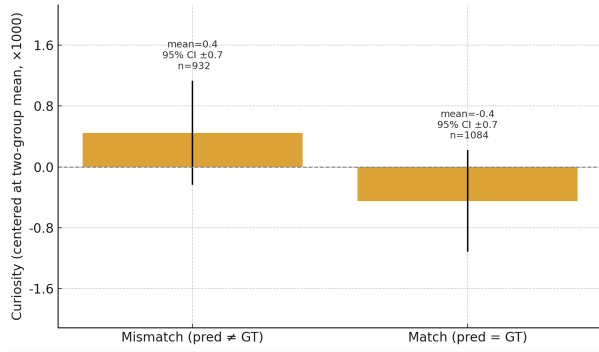


Figure 5: Curiosity scores based on match and mismatch of predictions from Qwen-0.5B base non-finetuned model and the ground truth

### ICM results against SFT baseline without explanations

Table 13: ICM method results against the SFT baseline without explanations (classification). Means $\pm$ SD are shown where SD was available from 5-fold runs.

Model	Experiment type	pearson	precision	recall	f1
Qwen-0.5B (SFT-Classification)	ID	<b>0.586 <math>\pm</math> 0.085</b>	<b>0.769</b>	0.461	0.551 $\pm$ 0.198
Qwen-0.5B (ICM)	ID	0.524 $\pm$ 0.092	0.494	<b>0.818</b>	<b>0.616 <math>\pm</math> 0.048</b>
Qwen-1.5B (SFT-Classification)	ID	<b>0.602 <math>\pm</math> 0.064</b>	<b>0.787</b>	0.602	<b>0.663 <math>\pm</math> 0.070</b>
Qwen-1.5B (ICM)	ID	0.586 $\pm$ 0.064	0.481	<b>0.794</b>	0.629 $\pm$ 0.045
Qwen-3B (SFT-Classification)	ID	0.482 $\pm$ 0.160	<b>0.670</b>	0.573	0.556 $\pm$ 0.094
Qwen-3B (ICM)	ID	<b>0.540 <math>\pm</math> 0.074</b>	0.481	<b>0.794</b>	<b>0.618 <math>\pm</math> 0.061</b>
Qwen-7B (SFT-Classification)	ID	0.441 $\pm$ 0.130	<b>0.535</b>	0.342	0.383 $\pm$ 0.251
Qwen-7B (ICM)	ID	<b>0.606 <math>\pm</math> 0.084</b>	0.518	<b>0.850</b>	<b>0.663 <math>\pm</math> 0.058</b>

**Note.** SDs for *precision* and *recall* were not available in the provided per-fold summaries; once those per-fold values are supplied, I will fill in their  $\pm$  SD as well. Pearson/F1 SDs are computed across 5 folds.

Table 14: ICM method results against the SFT baseline without explanations(classification) on Out-of-distribution data

Model	Experiment type	pearson	precision	recall	f1
Qwen-0.5B(SFT-Classification)	OOD	0.433	0.000	0.000	0.000
Qwen-0.5B(ICM)	OOD	<b>0.563</b>	<b>0.625</b>	<b>0.790</b>	<b>0.698</b>
Qwen-1.5B(SFT-Classification)	OOD	0.604	<b>0.962</b>	0.439	0.602
Qwen-1.5B(ICM)	OOD	<b>0.655</b>	0.639	<b>0.807</b>	<b>0.713</b>
Qwen-3B(SFT-Classification)	OOD	0.546	<b>0.933</b>	0.246	0.389
Qwen-3B(ICM)	OOD	<b>0.582</b>	0.597	<b>0.754</b>	<b>0.667</b>
Qwen-7B(SFT-Classification)	OOD	0.435	0.800	0.211	0.333
Qwen-7B(ICM)	OOD	<b>0.623</b>	<b>0.653</b>	<b>0.825</b>	<b>0.729</b>

### Curiosity scores based on non-finetuned base Qwen-0.5B model’s prediction and ground truth match and mismatch

#### Why is inverse model necessary?

When we ablated for the inverse model in our ICM setup with the given expert annotated data we do not see any dif-

ference in the results with using the inverse model or without using it. But the inverse model becomes necessary when we have a non-expert annotator like GPT-2, since it helps to clearly distinguish such outliers. This shows that our forward model of the ICM is good enough to distinguish between multiple expert annotators but we do need the inverse model for outlier cases. The details of our experiments can be found in Table 15, we used Qwen-0.5B model for this experiment.

Table 15: Inverse model ablations

Method	Annotations	Pearson	Precision	Recall	F1	Cohen's $\kappa$
ICM with Inverse	Without GPT-2	$0.503 \pm 0.014$	$0.552 \pm 0.014$	$0.728 \pm 0.017$	$0.628 \pm 0.015$	$0.347 \pm 0.027$
ICM without Inverse	Without GPT-2	$0.500 \pm 0.027$	$0.551 \pm 0.011$	$0.727 \pm 0.009$	$0.627 \pm 0.010$	$0.346 \pm 0.017$
ICM with Inverse	With GPT-2	$0.151 \pm 0.300$	$0.153 \pm 0.265$	$0.233 \pm 0.403$	$0.185 \pm 0.320$	$0.093 \pm 0.166$
ICM without Inverse	With GPT-2	$0.002 \pm 0.041$	$0.333 \pm 0.577$	$0.001 \pm 0.002$	$0.002 \pm 0.004$	$0.000 \pm 0.004$