

Contrastive Decoding Mitigates Score Range Bias in LLM-as-a-Judge

Yoshinari Fujinuma

Patronus AI

fujinumay@gmail.com

Abstract

Large Language Models (LLMs) are commonly used as evaluators in various applications, but the reliability of the outcomes remains a challenge. One such challenge is using LLMs-as-judges for direct assessment, i.e., assigning scores from a specified range without any references. Using summarization as our primary testbed, we first show that this challenge stems from LLM judge outputs being associated with score range bias, i.e., LLM judge outputs are highly sensitive to pre-defined score ranges. We also show that similar biases exist among models from the same family. We then mitigate this bias through contrastive decoding, achieving up to 11.7% relative improvement in Spearman correlation with human judgments, averaged across score ranges.¹

1 Introduction

Large Language Model (LLM) judges have become an integral component of the evaluation ecosystem (Lin et al., 2022; Chiang and Lee, 2023; Bubeck et al., 2023). In evaluations ranging from direct assessment, where judges evaluate individual outputs by assigning scores (Liu et al., 2023), to pairwise comparisons, where judges compare two outputs and determine which is superior (Zheng et al., 2023; Ye et al., 2024), using LLM as a judge is increasingly deployed to provide automatic, scalable, and cost-effective evaluation across diverse tasks. However, the reliability of such evaluations faces significant challenges, particularly when models assess their own outputs (Zheng et al., 2023) or those from the same model family (Goel et al., 2025). These biases constrain the set of models that can be reliably employed as LLM-as-a-judge for evaluation. But could there be any other biases hidden when using LLMs as judges?

¹Code: https://github.com/akkikiki/contrastive_decoding_llm_judge

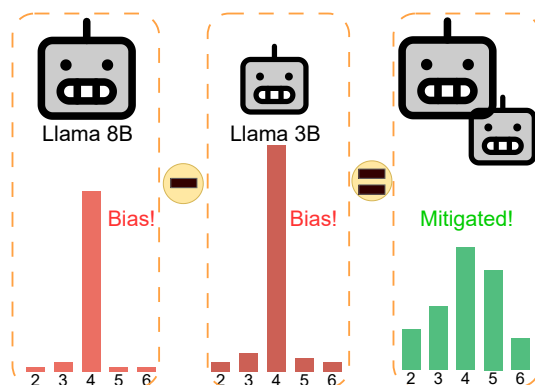


Figure 1: Overview of score range bias in 2-4 range and how contrastive decoding mitigates it through canceling out similar bias across models from the same family.

We reveal another bias in LLM judge outputs: **score range bias**, where LLM judges produce different correlations with human judgments under equivalent but shifted score ranges (e.g., 0-4 vs. 2-6), a phenomenon motivated by prior findings that LLMs struggle in simple arithmetic tasks (Nogueira et al., 2021; Gambardella et al., 2024). Upon identifying such biases, we also explore a mitigation strategy by connecting recent work on contrastive decoding (Li et al., 2023; O’Brien and Lewis, 2023) and family-enhancement bias (Goel et al., 2025), aiming to cancel out similar score range biases encoded across the models from the same family. We use summarization as our primary testbed, as prior work has shown LLM judges to be unreliable on summarization (Ye et al., 2024; Panickssery et al., 2024), and further validate that our findings generalize, a broad benchmark (Kim et al., 2025).

Our summarized contributions are as follows:

- We first show that LLM judges have *score range bias*, a bias observed across different model sizes and families (Llama-3 and Qwen-2.5) when judging on direct assessment.

- We then show that contrastive decoding, motivated by the observation of similar score range biases shared across models from the same family, successfully mitigates these biases.
- We show that these findings generalize beyond summarization to the diverse BigGen-Bench benchmark (§4.2).

2 Related Work

We now review the related work on LLM judges focusing on the judge tasks and their biases.

LLM Judge Tasks LLM judge tasks fall into two categories: *direct* assessment (Jones et al., 2024; Li et al., 2024; Zhu et al., 2025) and *pairwise* assessment (Zheng et al., 2023; Ye et al., 2024). Direct assessment (Liu et al., 2023) involves assigning numerical ratings to a single output example. In pairwise assessment, LLM judges show higher correlation with human preferences than direct assessment (Liu et al., 2024), supporting that the challenge remain in direct assessment, and we therefore focus on experimenting on direct assessment.

LLM Judge Biases One known bias in LLM judges is self-enhancement bias i.e., the tendency to favor their own output (Liu et al., 2023; Zheng et al., 2023; Ye et al., 2024) even in proprietary models like GPT-4 (Wataoka et al., 2024). Extending beyond self-enhancement bias, Goel et al. (2025) reported a family enhancement bias where models favor outputs from the same model family.

Input Sensitivity in LLM Judges A complementary line studies LLM sensitivity to input variations. Wei et al. (2024) identify order/token sensitivity in multiple-choice tasks and mitigate it via dataset-level probability weighting and calibration. We instead study score range bias in direct assessment, where shifting an equivalent 5-point range (e.g., 0-4 vs. 2-6) changes judge-to-human correlation; we further show this bias is shared across same-family models (§4.2) and mitigate it via instance-level contrastive decoding that cancels family-shared bias.

3 Analysis and Mitigation

3.1 Identifying Score Range Bias

An ideal LLM judge should maintain consistent correlation with human judgments across shifted ranges, as they represent the same 5-point scale (e.g., 0-4, 1-5, 2-6, 3-7). Failure to do so indicates **score range bias**: models exhibit different

correlations depending on the score range used, even when evaluating identical content. This bias outputs skewed distributions where models, for example, favor specific scores. We hypothesize that these biases can be mitigated with contrastive decoding, which cancels out shared biases.

3.2 Mitigation by Contrastive Decoding

Contrastive decoding (Li et al., 2023) modifies the model outputs by using two models: a main model and an assistant model. Given the next token probability of a main model p_{main} and an assistant model p_{asst} , the final adjusted score is calculated by subtracting the weighted p_{asst} from p_{main} i.e.,

$$\log p_{\text{main}} - \lambda \log p_{\text{asst}} \quad (1)$$

where $\lambda \in \mathbb{R}$ is the hyperparameter to control the magnitude of assistant model and logit e_i of token i is controlled by temperature $t > 0$ i.e., $p_{\text{asst}} = \frac{e_i/t}{\sum_j e_j/t}$. We depart from Li et al. (2023) by introducing λ because logit magnitudes differ across sizes (e.g., max logit ≈ 25 for Qwen-2.5-3B vs. ≈ 34 for 14B) and this mismatch carries over to log-probabilities; without λ the assistant can under-correct or erase the main model’s signal. We tune λ on a small held-out set (Appendix B).

4 Experiments

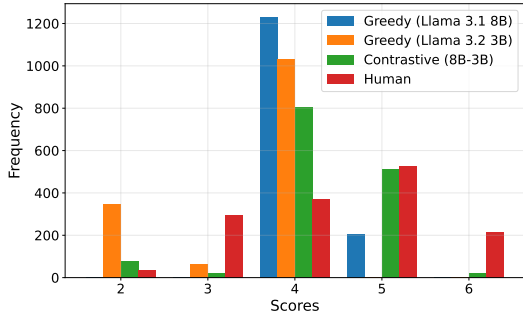
We focus on direct assessment on summarization since prior work reported that LLM judges fall short (Ye et al., 2024) and they are commonly evaluated on summarization (Panickssery et al., 2024).

4.1 Setup

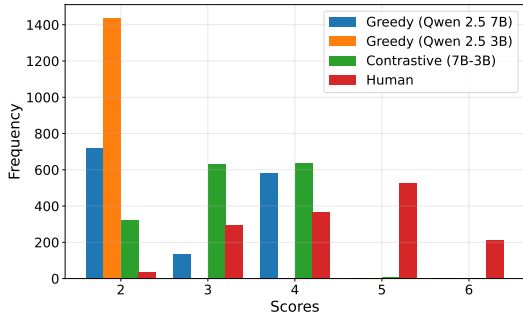
Task and Metrics We focus on the summarization task where LLM judges are commonly used (Liu et al., 2023; Panickssery et al., 2024). The correlations between human annotations are measured using three metrics: Pearson, Spearman, and Kendall correlations.

Score Scale and Ranges We use the 5 points Likert scale (Likert, 1932) on different score ranges (0-4, 1-5, 2-6, 3-7)². If output score parsing fails, we set to the lowest score following Liu et al. (2023) and if the parsed score exceeds the maximum, we clamp to the highest score in the range.

²We stopped at 7 inspired by Likert (1932) showing high correlation between 5 points (1-5) and 7 points (1-7) results.



(a) Llama-3 Family Results



(b) Qwen2.5 Family Results

Figure 2: Coherence score distribution in 2-6 score range with greedy decoding, contrastive decoding, and human annotations. The greedy decoding outputs from both Llama 8B and 3B models are highly skewed towards outputting score of 4. Qwen2.5 3B (see also Figure 3) and 7B models are outputting score of 2 showing similar biases are encoded in these models.

Models We experiment on two model families.³ For Llama-3 family (Grattafiori et al., 2024), we use Llama-3.1-8B-Instruct as the main model, Llama-3.2-3B-Instruct and Llama-3.2-1B-Instruct as the assistant model,⁴ and for Qwen2.5 family (Qwen et al., 2025), we use Qwen-2.5-14B-Instruct and Qwen-2.5-7B-Instruct as the main models, and Qwen-2.5-3B-Instruct as the assistant model. See Appendix A for the prompt used.

Dataset We use SummEval (Fabbri et al., 2021), a summarization benchmark also used by Liu et al. (2023) which contains 100 news articles where

³We leave models like Prometheus (Kim et al., 2024) specifically finetuned on judge tasks as future work since multiple model sizes are not available for contrastive decoding and those models are finetuned towards 1-5 score range.

⁴Although Llama-3.1 and Llama-3.2 are versioned separately, Meta documents them as part of the same Llama 3 herd (Grattafiori et al., 2024) with a shared tokenizer and architecture family. Figure 2 empirically supports it: Llama-3.1-8B and Llama-3.2-3B exhibit the same score range bias, both skewing toward Score 4 in the 2-6 range.

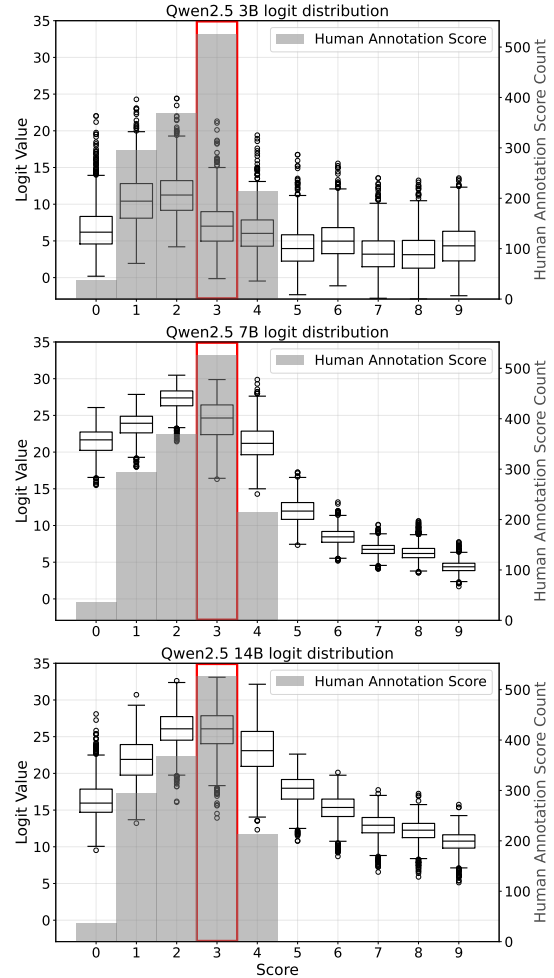


Figure 3: Logit distribution of the first output token in the 0-4 score range for Qwen2.5 3B, 7B, and 14B, with human annotation frequencies overlaid (gray bars). As model size increases from 3B to 14B, Score 3’s logit (red rectangle) grows and moves closer to the human annotations.

each article is associated with 16 summaries with human annotation scores, which sums up to 1600 summaries. 10% of the news articles are used as the held out development set to conduct grid search.

We now first reveal the score range bias of LLM judges by evaluating correlation to human annotations in the same 5 points scale but with different score ranges, and then experiment on mitigating it.

4.2 Reveal and Mitigate Score Range Biases

Similar score range biases exist across models from the same family We first analyze the distribution of the output scores in the 2 to 6 score range (Figure 2). Llama family models (3B and 8B) tend to output score of 4 (Figure 2a) and Qwen 2.5 family models tend to output score of 2 (Figure 2b). By using contrastive decoding, these biases toward

Model	Range	Pear.	Spear.	Kend.
Llama 3.2-1B	0 to 4	.073 _{.036}	.055 _{.035}	.047 _{.016}
Llama 3.2-3B	0 to 4	.056 _{.033}	.089 _{.028}	.073 _{.026}
Llama 3.1-8B	0 to 4	<u>.407_{.029}</u>	<u>.384_{.036}</u>	<u>.327_{.022}</u>
Contra. (8B-1B)	0 to 4	.384 _{.037}	.363 _{.038}	.309 _{.025}
Contra. (8B-3B)	0 to 4	.380 _{.030}	.357 _{.048}	.303 _{.040}
Llama 3.2-1B	1 to 5	.000 _{.000}	.000 _{.000}	.000 _{.000}
Llama 3.2-3B	1 to 5	.151 _{.036}	.201 _{.042}	.163 _{.016}
Llama 3.1-8B	1 to 5	.338 _{.036}	.309 _{.042}	.262 _{.029}
Contra. (8B-1B)	1 to 5	<u>.357_{.021}</u>	<u>.390_{.049}</u>	<u>.327_{.030}</u>
Contra. (8B-3B)	1 to 5	.323 _{.043}	.315 _{.029}	.265 _{.029}
Llama 3.2-1B	2 to 6	.000 _{.000}	.000 _{.000}	.000 _{.000}
Llama 3.2-3B	2 to 6	.008 _{.032}	.031 _{.029}	.026 _{.027}
Llama 3.1-8B	2 to 6	.262 _{.029}	.257 _{.025}	.220 _{.015}
Contra. (8B-1B)	2 to 6	<u>.344_{.038}</u>	<u>.337_{.035}</u>	<u>.288_{.034}</u>
Contra. (8B-3B)	2 to 6	.302 _{.049}	.305 _{.034}	.254 _{.023}
Llama 3.2-1B	3 to 7	-.039 _{.035}	-.051 _{.020}	-.043 _{.028}
Llama 3.2-3B	3 to 7	-.007 _{.035}	-.009 _{.044}	-.008 _{.027}
Llama 3.1-8B	3 to 7	<u>.445_{.027}</u>	<u>.426_{.030}</u>	<u>.352_{.030}</u>
Contra. (8B-1B)	3 to 7	.425 _{.030}	.408 _{.033}	.341 _{.019}
Contra. (8B-3B)	3 to 7	.442 _{.018}	.419 _{.035}	.347 _{.023}
<i>Average across all score ranges</i>				
Llama 3.2-1B		.009	.001	.001
Llama 3.2-3B		.052	.078	.064
Llama 3.1-8B		.363	.344	.290
Contra. (8B-1B)		.378	.375	.316
Contra. (8B-3B)		.362	.349	.292

Table 1: Llama-3 family correlation to humans on summary coherence with 95% confidence interval from bootstrap testing. Max correlation within score range are underlined and max averages are **bolded**. Maximal improvement is observed in the 2-6 score range.

specific ranges are mitigated and making the score outputs closer to human annotations.

Upon analyzing the first output token logit distribution (Figure 3) of the Qwen family models in the 0-4 score range, Qwen-2.5 3B, 7B, 8B, and 14B models encode similar biases where Score 2 is the highest while the most frequent human annotation is Score 3. The bias towards Score 2 gradually decreases as the model size scales from 3B to 14B, but still remains even in the 14B model. Furthermore, the logit range in each model differs e.g., max logit in 3B \approx 25, 7B \approx 30, and 14B \approx 34 (Figure 3), motivating the inclusion of λ in Eq. 1 to align the distributions between these models. This bias on Score 2 in the 3B model helps decrease similar bias encoded in the 7B and 14B models when used as the assistant model.

As a result of score range bias, using Llama 3B or 7B with greedy decoding causes the lowest correlation in 2-6 score range (Table 1). This trend is not limited to the Llama-3 family models and it is also observed in the Qwen-2.5 family models (Table 2). Focusing on greedy decoding, Qwen-2.5 family

Model	Range	Pear.	Spear.	Kend.
Qwen2.5-3B	0 to 4	-.042 _{.038}	-.059 _{.023}	-.051 _{.028}
Qwen2.5-7B	0 to 4	.248 _{.048}	.245 _{.028}	.209 _{.018}
Qwen2.5-14B	0 to 4	.435 _{.041}	.452 _{.019}	.382 _{.021}
Contra. (7B-3B)	0 to 4	.334 _{.033}	.332 _{.048}	.282 _{.023}
Contra. (14B-3B)	0 to 4	<u>.442_{.037}</u>	<u>.458_{.031}</u>	<u>.385_{.025}</u>
Qwen2.5-3B	1 to 5	-.044 _{.042}	-.056 _{.040}	-.048 _{.015}
Qwen2.5-7B	1 to 5	.382 _{.043}	.370 _{.042}	.309 _{.024}
Qwen2.5-14B	1 to 5	.468 _{.040}	.468 _{.026}	.385 _{.025}
Contra. (7B-3B)	1 to 5	.365 _{.042}	.356 _{.056}	.297 _{.019}
Contra. (14B-3B)	1 to 5	.463 _{.039}	<u>.470_{.041}</u>	<u>.387_{.027}</u>
Qwen2.5-3B	2 to 6	-.010 _{.000}	-.013 _{.000}	-.011 _{.000}
Qwen2.5-7B	2 to 6	.362 _{.024}	.355 _{.036}	.293 _{.025}
Qwen2.5-14B	2 to 6	.310 _{.037}	.313 _{.056}	.260 _{.026}
Contra. (7B-3B)	2 to 6	.373 _{.041}	.364 _{.034}	.298 _{.021}
Contra. (14B-3B)	2 to 6	<u>.410_{.028}</u>	<u>.426_{.035}</u>	<u>.359_{.040}</u>
Qwen2.5-3B	3 to 7	.033 _{.016}	.034 _{.000}	.029 _{.012}
Qwen2.5-7B	3 to 7	.341 _{.045}	.349 _{.024}	.289 _{.037}
Qwen2.5-14B	3 to 7	.349 _{.039}	.336 _{.036}	.275 _{.036}
Contra. (7B-3B)	3 to 7	.351 _{.019}	.357 _{.034}	.298 _{.027}
Contra. (14B-3B)	3 to 7	<u>.413_{.027}</u>	<u>.398_{.020}</u>	<u>.327_{.029}</u>
<i>Average across all score ranges</i>				
Qwen2.5-3B		-.016	-.024	-.020
Qwen2.5-7B		.333	.330	.275
Qwen2.5-14B		.391	.392	.326
Contra. (7B-3B)		.356	.352	.294
Contra. (14B-3B)		.432	.438	.365

Table 2: Qwen-2.5 family correlation to humans on summary coherence with 95% confidence interval from bootstrap testing. Max correlation within score range are underlined and max averages are **bolded**. Maximal improvement is observed in the 2-6 score range.

and Llama3.1-3B show clearer trend that the 1-5 score range shows the highest correlation among the experimented score ranges (7B, 1-5: .370, 14B, 1-5: .468), while Llama3.1-8B being the exception that 3 to 7 score range showing highest correlation among all score ranges. These outcomes further raises concern on applying LLM judges beyond the standard 1-5 range.

Contrastive Decoding is a Robust Mitigation Strategy across Different Score Ranges

We use *robust* here to mean reducing variance in judge-to-human correlation across score ranges, rather than uniformly improving at every individual range. Table 1 and 2 show that contrastive decoding exhibits this robustness: correlations remain more consistent across varying score ranges, directly addressing the score range bias observed. While using a single model suffers from decrease in correlation when score ranges are shifted, contrastive decoding maintains more stable correlations with human judgments regardless of the score ranges (Table 1). This robustness is evident in the 2-6 range, where

contrastive decoding on Llama-3 family achieves a Pearson correlation of .310 (compared to .168 for Llama 3.2-3B and .270 for Llama 3.1-8B) and similar improvements in Spearman and Kendall correlations, also seen as 6.7% relative improvement for Llama 8B (.330 \rightarrow .352) and 11.7% for Qwen 14B (.392 \rightarrow .438) on average across all score ranges on SummEval. The stability across different scoring ranges enables search on optimal score ranges beyond the 1-5 range (e.g., 0-4 range showing the best correlation in summary relevance for Qwen family in Appendix D).

Does Assistant Model Choice Matter for Bias Mitigation?

Table 1 shows that the choice of assistant model slightly impacts correlations, with the 1B model marginally outperforming the 3B model. The 1B assistant achieves an average Spearman correlation of .375 compared to .352 for the 3B assistant. Notably, using larger assistant models can degrade performance: our ablation study with Qwen 14B-7B (Appendix H) shows significant degradation in the 1-5 range (Spearman: .282 vs .470 with 3B assistant), showing that larger assistants penalize correct logits from the main model.

Generalization to BigGen-Bench To validate that our findings extend beyond summarization, we replicate our analysis on BigGen-Bench (Kim et al., 2025), a fine-grained LLM evaluation benchmark spanning 77 tasks covering categories such as Reasoning, Dialogue/Generation, Safety, Planning, Tool Usage, Instruction Following, Theory of Mind, and Grounding, each paired with task-specific rubrics and human gold scores. We apply the same judge setup as in our SummEval experiments: the 5-point Likert scale with the four score ranges (0-4, 1-5, 2-6, 3-7), identical prompting, and the benchmark’s human gold scores as the correlation reference. Tables 3 and 4 show that contrastive decoding improves the average Spearman correlation with human judgments from .349 to .352 for Llama 3.1-8B (1B assistant) and from .500 to .517 for Qwen2.5-14B (3B assistant), consistent with the trends observed on SummEval.

5 Conclusion

In this work, we analyze and experiment with LLM-as-a-judge on direct assessment, which reveals two key findings: First, LLM judges exhibit a score range bias across different model families and sizes with a tendency to favor specific scores

Model	Range	Pear.	Spear.	Kend.
Llama 3.1-8B	0 to 4	<u>.384</u> _{.035}	<u>.361</u> _{.034}	<u>.301</u> _{.029}
Contra. (8B-1B)	0 to 4	<u>.395</u> _{.036}	<u>.367</u> _{.034}	<u>.307</u> _{.029}
Llama 3.1-8B	1 to 5	<u>.377</u> _{.035}	<u>.345</u> _{.035}	<u>.289</u> _{.030}
Contra. (8B-1B)	1 to 5	<u>.392</u> _{.036}	<u>.352</u> _{.034}	<u>.300</u> _{.030}
Llama 3.1-8B	2 to 6	<u>.353</u> _{.037}	<u>.336</u> _{.035}	<u>.284</u> _{.030}
Contra. (8B-1B)	2 to 6	<u>.359</u> _{.035}	<u>.349</u> _{.034}	<u>.294</u> _{.029}
Llama 3.1-8B	3 to 7	<u>.373</u> _{.034}	<u>.354</u> _{.034}	<u>.300</u> _{.030}
Contra. (8B-1B)	3 to 7	<u>.357</u> _{.037}	<u>.338</u> _{.036}	<u>.288</u> _{.031}
<i>Average across all score ranges</i>				
Llama 3.1-8B		.372	.349	.294
Contra. (8B-1B)		.376	.352	.297

Table 3: Llama 3.1 family correlation results on BigGen-Bench with 95% bootstrap confidence intervals. Max correlation within score range are underlined and max averages are **bolded**.

Model	Range	Pear.	Spear.	Kend.
Qwen2.5-14B	0 to 4	<u>.545</u> _{.031}	<u>.517</u> _{.030}	<u>.445</u> _{.027}
Contra. (14B-3B)	0 to 4	<u>.571</u> _{.029}	<u>.533</u> _{.030}	<u>.460</u> _{.026}
Qwen2.5-14B	1 to 5	<u>.520</u> _{.033}	<u>.489</u> _{.032}	<u>.427</u> _{.028}
Contra. (14B-3B)	1 to 5	<u>.534</u> _{.033}	<u>.498</u> _{.032}	<u>.435</u> _{.028}
Qwen2.5-14B	2 to 6	<u>.533</u> _{.032}	<u>.501</u> _{.032}	<u>.436</u> _{.028}
Contra. (14B-3B)	2 to 6	<u>.556</u> _{.031}	<u>.530</u> _{.031}	<u>.459</u> _{.027}
Qwen2.5-14B	3 to 7	<u>.528</u> _{.033}	<u>.494</u> _{.033}	<u>.432</u> _{.029}
Contra. (14B-3B)	3 to 7	<u>.539</u> _{.032}	<u>.507</u> _{.030}	<u>.443</u> _{.028}
<i>Average across all score ranges</i>				
Qwen2.5-14B		.532	.500	.435
Contra. (14B-3B)		.550	.517	.449

Table 4: Qwen2.5 family correlation results on BigGen-Bench with 95% bootstrap confidence intervals. Max correlation within score range are underlined, and max averages are **bolded**. Qwen2.5 14B with contrastive decoding (14B-3B) consistently achieves the best performance across all score ranges.

regardless of the quality of the summaries. Second, we show that contrastive decoding effectively mitigates score range bias by leveraging the similar biases present in models from the same family. Our score range bias analysis framework can test arbitrary models and can help unlock the potential to expand beyond the standard 1-5 score range.

Limitations

Inference Time Compute Contrastive learning increases the test time compute due to running forward pass on two models rather than one. On the other hand, using a main model and an assistant model is very common in real world setup to speed up decoding with speculative decoding (Leviathan et al., 2023), and therefore contrastive decoding can be used without an additional forward pass when

speculative decoding is used.

Model Size Our experiments were limited to models with up to 14B parameters due to computational budget constraints.

Language Coverage Our experiments are conducted only on English language, however, we have not exploited linguistic knowledge specific to English.

Acknowledgments

We sincerely thank the anonymous reviewers for their constructive feedback and insightful comments, which helped strengthen the analysis and improve the clarity of this work.

References

- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with gpt-4](#). *Preprint*, arXiv:2303.12712.
- Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large language models be an alternative to human evaluations?](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.
- Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. [SummEval: Re-evaluating summarization evaluation](#). *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Andrew Gambardella, Yusuke Iwasawa, and Yutaka Matsuo. 2024. [Language models do hard arithmetic tasks easily and hardly do easy arithmetic tasks](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 85–91, Bangkok, Thailand. Association for Computational Linguistics.
- Shashwat Goel, Joschka Struber, Ilze Amanda Auzina, Karuna K Chandra, Ponnurangam Kumaraguru, Douwe Kiela, Ameya Prabhu, Matthias Bethge, and Jonas Geiping. 2025. [Great models think alike and this undermines ai oversight](#). *Preprint*, arXiv:2502.04313.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-han Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-badur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-

feng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Couderc, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Barambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damla, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khanelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha

White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangrabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihalescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

Jaylen Jones, Lingbo Mo, Eric Fosler-Lussier, and Huan Sun. 2024. [A multi-aspect framework for counter narrative evaluation using large language models](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 147–168, Mexico City, Mexico. Association for Computational Linguistics.

Seungone Kim, Juyoung Suk, Ji Yong Cho, Shayne Longpre, Chaeun Kim, Dongkeun Yoon, Guijin Son, Yejin Cho, Sheikh Shafayat, Jinheon Baek, Sue Hyun Park, Hyeonbin Hwang, Jinkyung Jo, Hyowon Cho, Haebin Ravi, Hadi Pouransari, Suyeon Kevin Ha, Cheng Chen, Jiwoo Roh, Chirag Raman, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2025. [The BIGGEN BENCH: A principled benchmark for fine-grained evaluation of language models with language models](#). In *Proceedings of the 2025 Conference of the North American Chapter of the Association for Computational Linguistics*.

- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. [Prometheus 2: An open source language model specialized in evaluating other language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4334–4353, Miami, Florida, USA. Association for Computational Linguistics.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org.
- Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, hai zhao, and Pengfei Liu. 2024. [Generative judge for evaluating alignment](#). In *The Twelfth International Conference on Learning Representations*.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023. [Contrastive decoding: Open-ended text generation as optimization](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12286–12312, Toronto, Canada. Association for Computational Linguistics.
- Rensis Likert. 1932. A technique for the measurement of attitudes. *Archives of psychology*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [TruthfulQA: Measuring how models mimic human falsehoods](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Yinhong Liu, Han Zhou, Zhijiang Guo, Ehsan Shareghi, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2024. [Aligning with human judgement: The role of pairwise preference in large language model evaluators](#). In *First Conference on Language Modeling*.
- Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. 2021. Investigating the limitations of transformers with simple arithmetic tasks. *Preprint*, arXiv:2102.13019.
- Sean O’Brien and Mike Lewis. 2023. [Contrastive decoding improves reasoning in large language models](#). *Preprint*, arXiv:2309.09117.
- Arjun Panickssery, Samuel R. Bowman, and Shi Feng. 2024. [LLM evaluators recognize and favor their own generations](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. [Qwen2.5 technical report](#). *Preprint*, arXiv:2412.15115.
- Koki Wataoka, Tsubasa Takahashi, and Ryokan Ri. 2024. [Self-preference bias in LLM-as-a-judge](#). In *Neurips Safe Generative AI Workshop 2024*.
- Sheng-Lun Wei, Cheng-Kuang Wu, Hen-Hsen Huang, and Hsin-Hsi Chen. 2024. [Unveiling selection biases: Exploring order and token sensitivity in large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 5598–5621, Bangkok, Thailand. Association for Computational Linguistics.
- Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, Nitesh V Chawla, and Xiangliang Zhang. 2024. [Justice or prejudice? quantifying biases in llm-as-a-judge](#). *Preprint*, arXiv:2410.02736.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA. Curran Associates Inc.
- Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2025. [JudgeLM: Fine-tuned large language models are scalable judges](#). In *The Thirteenth International Conference on Learning Representations*.

A Judge Prompts

We use the following prompt experimented by Liu et al. (2023).

Score Range {min_range}-{max_range} for Coherence

You will be given one summary written for a news article.

Your task is to rate the summary on one metric.

Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.

Evaluation Criteria:

Coherence ({min_range}-{max_range}) - the collective quality of all sentences. We align this dimension with the DUC quality question of structure and coherence whereby "the summary should be well-structured and well-organized. The summary should not just be a heap of related information, but should build from sentence to a coherent body of information about a topic."

Evaluation Steps:

1. Read the news article carefully and identify the main topic and key points.
2. Read the summary and compare it to the news article. Check if the summary covers the main topic and key points of the news article, and if it presents them in a clear and logical order.
3. Assign a score for coherence on a scale of {min_range} to {max_range}, where {min_range} is the lowest and {max_range} is the highest based on the Evaluation Criteria.

Example:

Source Text:

{{Document}}

Summary:

{{Summary}}

Evaluation Form (scores ONLY):

- Coherence:

What is the coherence of the summary above? Provide only rating and no other text.

Score Range {min_range}-{max_range} for Relevance

You will be given one summary written for a news article.

Your task is to rate the summary on one metric.

Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.

Evaluation Criteria:

Relevance ({min_range}-{max_range}) - selection of important content from the source. The summary should include only important information from the source document. Annotators were instructed to penalize summaries which contained redundancies and excess information.

Evaluation Steps:

1. Read the summary and the source document carefully.
2. Compare the summary to the source document and identify the main points of the article.
3. Assess how well the summary covers the main points of the article, and how much irrelevant or redundant information it contains.
4. Assign a relevance score from {min_range} to {max_range}.

Example:

Source Text:

{{Document}}

Summary:

{{Summary}}

Evaluation Form (scores ONLY):

- Relevance:

What is the relevance of the summary above? Provide only rating and no other text.

Score Range {min_range}-{max_range} for Consistency

You will be given one summary written for a news article.

Your task is to rate the summary on one metric.

Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.

Evaluation Criteria:

Consistency ({min_range}-{max_range}) - the factual alignment between the summary and the summarized source. A factually consistent summary contains only statements that are entailed by the source document. Annotators were also asked to penalize summaries that contained hallucinated facts.

Evaluation Steps:

1. Read the news article carefully and identify the main topic and key points.
2. Read the summary and compare it to the news article. Check if the summary covers the main topic and key points of the news article, and if it presents them in a clear and logical order.
3. Assign a score for consistency based on the Evaluation Criteria.

Example:

Source Text:

{{Document}}

Summary:

{{Summary}}

Evaluation Form (scores ONLY):

- Consistency:

What is the consistency of the summary above? Provide only rating and no other text.

B Hyper-parameters

We conduct grid search over two hyperparameters for contrastive decoding: 1) temperature t and 2) scaling constant λ from the following ranges:

- $\lambda = [0.01, 0.1, 0.5, 1.0]$

- $t = [0.5, 1.0, 2.0, 3.0, 4.0, 5.0]$

The following table shows the hyperparameters setup for each setting:

Main	Asst	Range	λ	t	
Llama 3.1 8B	Llama 3.2 3B	0-4	0.01	1.0	
		1-5	1.0	0.5	
		2-6	1.0	0.5	
	Llama 3.2 1B	3-7	0.01	5.0	
		0-4	0.01	0.5	
		1-5	0.1	5.0	
	Qwen 2.5 7B	Qwen 2.5 3B	2-6	0.1	2.0
			3-7	0.1	2.0
			0-4	0.1	4.0
Qwen 2.5 14B	Qwen 2.5 3B	1-5	0.01	4.0	
		2-6	0.1	1.0	
		3-7	0.1	2.0	

Table 5: Hyperparameter settings for contrastive decoding for each main and assistant model pair from each model family for evaluating summary coherence.

Main	Asst	Range	λ	t	
Llama 3.1 8B	Llama 3.2 3B	0-4	0.01	0.5	
		1-5	0.01	0.5	
		2-6	0.01	0.5	
	Llama 3.2 1B	3-7	0.5	0.5	
		0-4	0.01	0.5	
		1-5	0.1	5.0	
	Qwen 2.5 7B	Qwen 2.5 3B	2-6	0.1	5.0
			3-7	0.01	0.5
			0-4	0.1	5.0
Qwen 2.5 14B	Qwen 2.5 3B	1-5	0.5	1.0	
		2-6	1.0	1.0	
		3-7	0.1	4.0	
		0-4	0.01	3.0	
		1-5	0.1	0.5	
		2-6	0.01	3.0	
		3-7	0.01	3.0	

Table 6: Hyperparameter settings for contrastive decoding for each Llama-3 main and assistant model pair for relevance.

C Score Distribution Across All Ranges

Figures 4, 5, and 6 show the score distribution comparison plots across all four score ranges (0-4, 1-5, 2-6, 3-7) for Llama-3, Qwen2.5 7B, and Qwen2.5 14B family models respectively. These plots show how score range bias manifests differently across ranges and model sizes, and how contrastive decoding consistently mitigates this bias.

Main	Asst	Range	λ	t		
Llama 3.1 8B	Llama 3.2 3B	0-4	0.01	5.0		
		1-5	0.1	2.0		
		2-6	0.1	2.0		
	Llama 3.2 1B	Llama 3.2 3B	3-7	0.1	3.0	
			0-4	0.1	1.0	
			1-5	0.1	0.5	
		Qwen 2.5 7B	Qwen 2.5 3B	2-6	0.1	2.0
				3-7	0.1	1.0
				0-4	0.1	3.0
Qwen 2.5 14B	Qwen 2.5 3B	1-5	0.01	5.0		
		2-6	0.01	2.0		
		3-7	0.01	5.0		
	Qwen 2.5 7B	Qwen 2.5 3B	0-4	0.1	1.0	
			1-5	0.1	5.0	
			2-6	0.1	3.0	
Qwen 2.5 14B	Qwen 2.5 3B	3-7	0.1	3.0		

Table 7: Hyperparameter settings for contrastive decoding for each Llama-3 main and assistant model pair for consistency.

C.1 Invalid Output Analysis

Tables 8 and 9 show the analysis of invalid model outputs across different score ranges for Llama-3 and Qwen2.5 models. Invalid outputs are classified into three categories: (1) *No Pred*: outputs where no numeric pattern could be extracted, and (2) *Below Range*: outputs with valid numbers below the minimum score for the range.

For Llama models, the 8B model achieves near-perfect parsing with 0% invalid outputs across all ranges. The 3B model shows moderate failure rates ranging from 17.4% to 24.0%, primarily due to malformed outputs. Notably, contrastive decoding maintains excellent performance with minimal failures (0.3% in the 2-6 range).

For Qwen models, the analysis reveals that the 3B model exhibits catastrophic failure with approximately 98% invalid outputs across all score ranges, primarily due to malformed outputs that cannot be parsed. In contrast, the 7B and 14B models show significantly better performance, with the 7B model achieving near-perfect parsing (0% invalid) and the 14B model showing minimal failures, primarily in the 3-7 range where 4.6% of outputs fall below the minimum score.

D Relevance and Consistency Results

In Tables 10 and 11, we present the correlation results for summary relevance evaluation across different score ranges for the Llama-3 and Qwen2.5 model families, respectively. In Tables 12 and 13, we present the correlation results for summary consistency evaluation across different score ranges.

Model	Range	Invalid	Failure Types
Llama 3.1 8B	0-4	0.0%	-
	1-5	0.0%	-
	2-6	0.0%	-
	3-7	0.0%	-
Llama 3.2 3B	0-4	23.2%	No Pred: 334
	1-5	17.4%	No Pred: 251
	2-6	24.0%	No Pred: 345
	3-7	18.6%	No Pred: 268
Llama 3.1 8B Contrastive	0-4	0.0%	-
	1-5	0.0%	-
	2-6	0.3%	No Pred: 5
	3-7	0.0%	-

Table 8: Invalid output statistics for Llama-3 models across all score ranges on coherence evaluation. The 8B model shows perfect parsing, while the 3B model exhibits moderate failure rates of 17-24%.

E Model Size and Budget

For all the experiments in this paper, NVIDIA’s A100 GPU was used. The base models used in this paper are licensed under Meta Llama 3 License⁵ for Llama 3 family models and Apache-2.0 license for Qwen 2.5 family models. We followed their intended use case.

F Information About Use of AI Assistants

We have used Claude on this manuscript to enhance the clarity of the paper and fixing grammatical mistakes. We also used it to create the codes to run experiments.

G Potential Risks

As discussed in the limitations section, the experiments are only conducted on English, which may bias the takeaways on English.

H Assistant Model Size Analysis

The paper notes that smaller assistant models sometimes work better than larger ones, but does not explore why. We conducted a follow-up experiment using Qwen2.5-7B as an assistant model (instead of the 3B model) for the Qwen2.5-14B main model on summary coherence evaluation.

Table 14 shows the results for Qwen 14B-7B contrastive decoding across all score ranges. We observe no significant changes in the 2-6 and 3-7 ranges compared to using the 3B assistant. However, there is a notable degradation in the 1-5 range, where the 7B assistant model significantly reduces

⁵<https://www.llama.com/llama3/license/>

Model	Range	Invalid	Failure Types
Qwen2.5 3B	0-4	98.2%	No Pred: 1414
	1-5	98.7%	No Pred: 1421
	2-6	97.6%	No Pred: 1406
	3-7	97.6%	No Pred: 1376, B: 30
Qwen2.5 7B	0-4	0.0%	-
	1-5	0.0%	-
	2-6	0.0%	-
	3-7	0.0%	-
Qwen2.5 14B	0-4	0.0%	-
	1-5	0.0%	-
	2-6	0.1%	B: 1
	3-7	4.6%	B: 66
Qwen2.5 7B Contrastive	0-4	0.0%	-
	1-5	0.0%	-
	2-6	0.0%	-
	3-7	0.0%	-
Qwen2.5 14B Contrastive	0-4	0.0%	-
	1-5	0.0%	-
	2-6	0.1%	B: 1
	3-7	5.6%	B: 81

Table 9: Invalid output statistics for Qwen2.5 models across all score ranges on coherence evaluation. The 3B model shows near-complete failure in generating valid scores, while larger models perform significantly better. B: Below

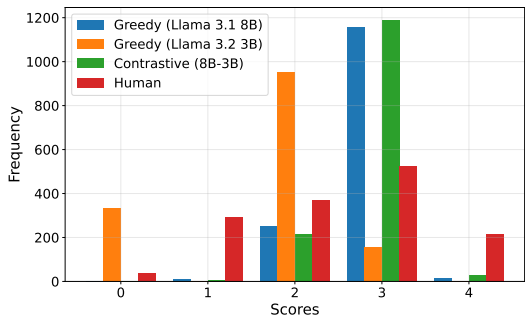
correlation (Spearman: 0.282 vs 0.470 with 3B assistant, see Table 2). This finding supports the hypothesis that using a larger assistant model is more likely to penalize correct logits from the main model, thereby degrading performance.

Model	Range	Pear.	Spear.	Kend.
Llama 3.2-1B	0 to 4	<u>.028</u> _{.041}	<u>.028</u> _{.048}	<u>.024</u> _{.037}
Llama 3.2-3B	0 to 4	<u>.072</u> _{.027}	<u>.111</u> _{.036}	<u>.094</u> _{.025}
Llama 3.1-8B	0 to 4	<u>.429</u> _{.050}	<u>.374</u> _{.027}	<u>.323</u> _{.016}
Contra. (8B-3B)	0 to 4	<u>.407</u> _{.036}	<u>.360</u> _{.038}	<u>.310</u> _{.027}
Llama 3.2-1B	1 to 5	<u>.014</u> _{.054}	<u>.022</u> _{.036}	<u>.019</u> _{.033}
Llama 3.2-3B	1 to 5	<u>.161</u> _{.056}	<u>.179</u> _{.023}	<u>.152</u> _{.023}
Llama 3.1-8B	1 to 5	<u>.391</u> _{.042}	<u>.341</u> _{.025}	<u>.293</u> _{.019}
Contra. (8B-3B)	1 to 5	<u>.375</u> _{.030}	<u>.322</u> _{.035}	<u>.277</u> _{.023}
Llama 3.2-1B	2 to 6	<u>.000</u> _{.000}	<u>.000</u> _{.000}	<u>.000</u> _{.000}
Llama 3.2-3B	2 to 6	<u>.076</u> _{.041}	<u>.107</u> _{.037}	<u>.091</u> _{.041}
Llama 3.1-8B	2 to 6	<u>.378</u> _{.020}	<u>.360</u> _{.033}	<u>.312</u> _{.032}
Contra. (8B-3B)	2 to 6	<u>.393</u> _{.033}	<u>.381</u> _{.030}	<u>.331</u> _{.025}
Llama 3.2-1B	3 to 7	<u>.036</u> _{.018}	<u>.033</u> _{.018}	<u>.029</u> _{.016}
Llama 3.2-3B	3 to 7	<u>.092</u> _{.044}	<u>.095</u> _{.040}	<u>.082</u> _{.037}
Llama 3.1-8B	3 to 7	<u>.438</u> _{.041}	<u>.421</u> _{.021}	<u>.357</u> _{.019}
Contra. (8B-3B)	3 to 7	<u>.471</u> _{.035}	<u>.444</u> _{.023}	<u>.369</u> _{.028}
<i>Average across all score ranges</i>				
Llama 3.2-1B		<u>.026</u>	<u>.028</u>	<u>.024</u>
Llama 3.2-3B		<u>.100</u>	<u>.123</u>	<u>.105</u>
Llama 3.1-8B		<u>.409</u>	<u>.374</u>	<u>.321</u>
Contra. (8B-3B)		.412	.377	.322

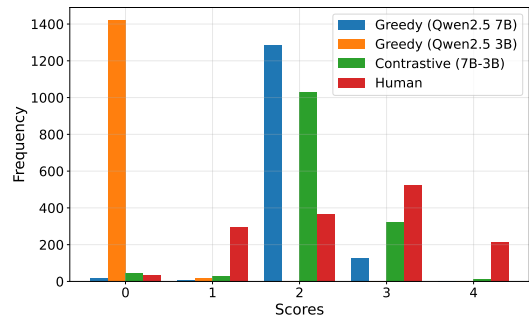
Table 10: Llama-3 family correlation results to human annotations on summary relevance. Max correlation within score range are underlined, max across all ranges are *italicized*, and max averages are **bolded**.

Model	Range	Pear.	Spear.	Kend.
Qwen2.5-3B	0 to 4	<u>.000</u> _{.000}	<u>.000</u> _{.000}	<u>.000</u> _{.000}
Qwen2.5-7B	0 to 4	<u>.323</u> _{.028}	<u>.296</u> _{.041}	<u>.250</u> _{.045}
Qwen2.5-14B	0 to 4	<u>.517</u> _{.033}	<u>.496</u> _{.034}	<u>.416</u> _{.023}
Contra. (7B-3B)	0 to 4	<u>.351</u> _{.031}	<u>.289</u> _{.037}	<u>.236</u> _{.053}
Contra. (14B-3B)	0 to 4	<u>.516</u> _{.031}	<u>.478</u> _{.038}	<u>.406</u> _{.016}
Qwen2.5-3B	1 to 5	<u>.000</u> _{.000}	<u>.000</u> _{.000}	<u>.000</u> _{.000}
Qwen2.5-7B	1 to 5	<u>.345</u> _{.037}	<u>.347</u> _{.043}	<u>.285</u> _{.021}
Qwen2.5-14B	1 to 5	<u>.518</u> _{.036}	<u>.500</u> _{.036}	<u>.427</u> _{.033}
Contra. (7B-3B)	1 to 5	<u>.353</u> _{.045}	<u>.335</u> _{.036}	<u>.283</u> _{.043}
Contra. (14B-3B)	1 to 5	<u>.509</u> _{.024}	<u>.495</u> _{.034}	<u>.422</u> _{.024}
Qwen2.5-3B	2 to 6	<u>.000</u> _{.000}	<u>.000</u> _{.000}	<u>.000</u> _{.000}
Qwen2.5-7B	2 to 6	<u>.415</u> _{.040}	<u>.398</u> _{.024}	<u>.330</u> _{.027}
Qwen2.5-14B	2 to 6	<u>.438</u> _{.034}	<u>.380</u> _{.038}	<u>.319</u> _{.030}
Contra. (7B-3B)	2 to 6	<u>.171</u> _{.020}	<u>.134</u> _{.024}	<u>.112</u> _{.027}
Contra. (14B-3B)	2 to 6	<u>.453</u> _{.036}	<u>.408</u> _{.027}	<u>.344</u> _{.024}
Qwen2.5-3B	3 to 7	<u>.000</u> _{.000}	<u>.000</u> _{.000}	<u>.000</u> _{.000}
Qwen2.5-7B	3 to 7	<u>.468</u> _{.023}	<u>.441</u> _{.035}	<u>.367</u> _{.020}
Qwen2.5-14B	3 to 7	<u>.426</u> _{.034}	<u>.375</u> _{.037}	<u>.313</u> _{.023}
Contra. (7B-3B)	3 to 7	<u>.456</u> _{.035}	<u>.438</u> _{.024}	<u>.365</u> _{.024}
Contra. (14B-3B)	3 to 7	<u>.510</u> _{.026}	<u>.473</u> _{.038}	<u>.394</u> _{.016}
<i>Average across all score ranges</i>				
Qwen2.5-3B		<u>.000</u>	<u>.000</u>	<u>.000</u>
Qwen2.5-7B		<u>.388</u>	<u>.371</u>	<u>.308</u>
Qwen2.5-14B		<u>.475</u>	<u>.438</u>	<u>.369</u>
Contra. (7B-3B)		<u>.333</u>	<u>.299</u>	<u>.249</u>
Contra. (14B-3B)		.497	.464	.392

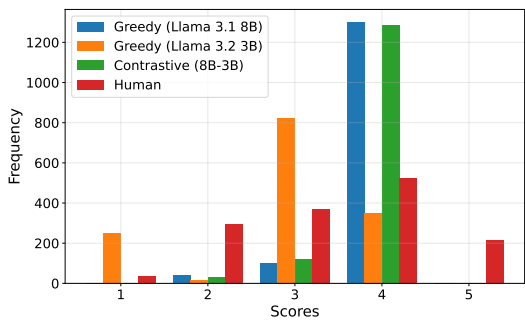
Table 11: Qwen2.5 family correlation results to human annotations on summary relevance. Max correlation within score range are underlined and max averages are **bolded**.



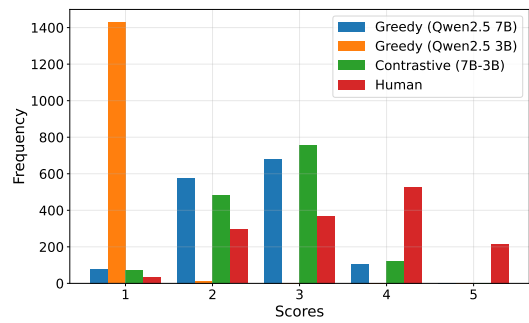
(a) Range 0-4



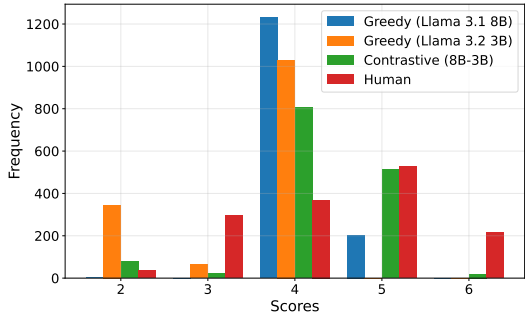
(a) Range 0-4



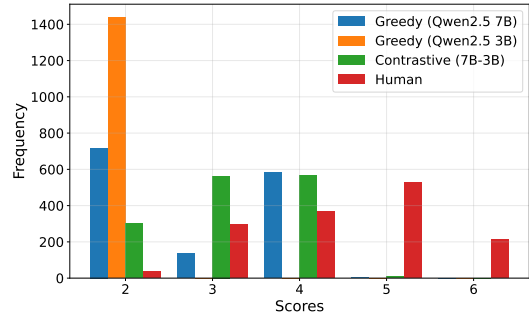
(b) Range 1-5



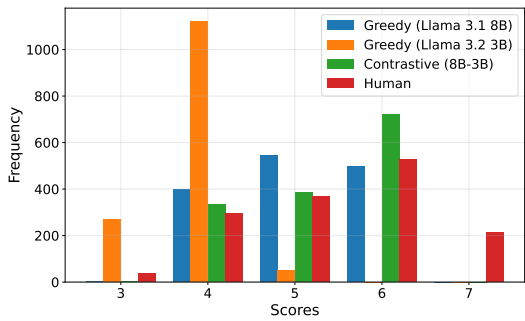
(b) Range 1-5



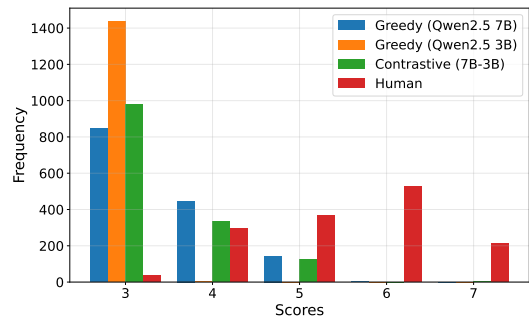
(c) Range 2-6



(c) Range 2-6



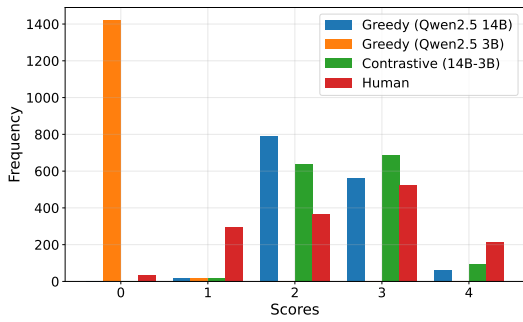
(d) Range 3-7



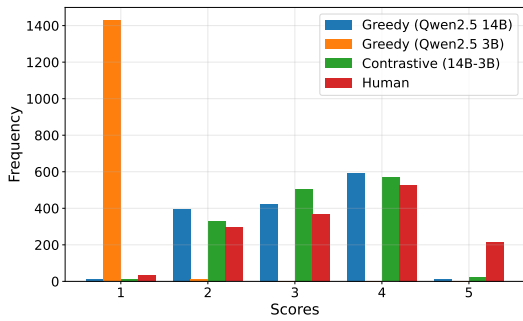
(d) Range 3-7

Figure 4: Predicted score distribution comparison for Llama-3 family across all score ranges on coherence evaluation.

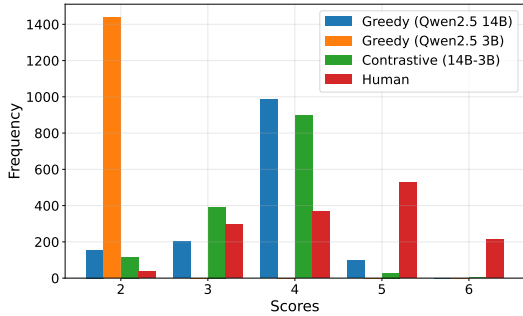
Figure 5: Predicted score distribution comparison for Qwen2.5 family across all score ranges on coherence evaluation.



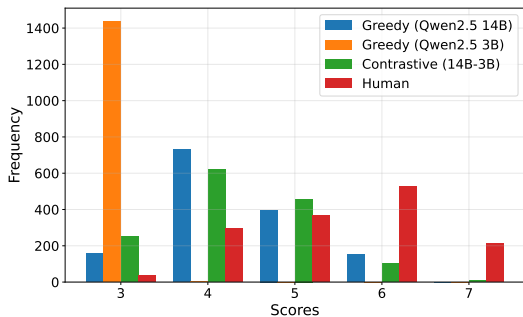
(a) Range 0-4



(b) Range 1-5



(c) Range 2-6



(d) Range 3-7

Figure 6: Score distribution comparison for Qwen2.5 14B family across all score ranges on coherence evaluation.

Model	Range	Pear.	Spear.	Kend.
Llama 3.2-1B	0 to 4	-.022 _{.025}	.003 _{.038}	.003 _{.025}
Llama 3.2-3B	0 to 4	.051 _{.039}	.073 _{.039}	.068 _{.034}
Llama 3.1-8B	0 to 4	<u>.471_{.044}</u>	<u>.411_{.027}</u>	<u>.388_{.030}</u>
Contra. (8B-3B)	0 to 4	<u>.445_{.039}</u>	<u>.381_{.030}</u>	<u>.360_{.027}</u>
Llama 3.2-1B	1 to 5	-.022 _{.007}	-.028 _{.007}	-.027 _{.009}
Llama 3.2-3B	1 to 5	.178 _{.056}	.196 _{.038}	.182 _{.025}
Llama 3.1-8B	1 to 5	<u>.594_{.047}</u>	<u>.487_{.053}</u>	<u>.465_{.032}</u>
Contra. (8B-3B)	1 to 5	<u>.603_{.056}</u>	<u>.518_{.048}</u>	<u>.494_{.035}</u>
Llama 3.2-1B	2 to 6	.000 _{.000}	.000 _{.000}	.000 _{.000}
Llama 3.2-3B	2 to 6	.095 _{.054}	.097 _{.059}	.092 _{.037}
Llama 3.1-8B	2 to 6	<u>.488_{.057}</u>	<u>.431_{.021}</u>	<u>.413_{.045}</u>
Contra. (8B-3B)	2 to 6	<u>.551_{.052}</u>	<u>.469_{.031}</u>	<u>.447_{.041}</u>
Llama 3.2-1B	3 to 7	-.011 _{.040}	-.004 _{.050}	-.004 _{.024}
Llama 3.2-3B	3 to 7	.118 _{.035}	.136 _{.051}	.129 _{.038}
Llama 3.1-8B	3 to 7	<u>.549_{.045}</u>	<u>.484_{.050}</u>	<u>.452_{.050}</u>
Contra. (8B-3B)	3 to 7	<u>.540_{.054}</u>	<u>.496_{.034}</u>	<u>.465_{.028}</u>
<i>Average across all score ranges</i>				
Llama 3.2-1B		-.014	-.007	-.007
Llama 3.2-3B		.111	.126	.118
Llama 3.1-8B		.526	.453	.430
Contra. (8B-3B)		.535	.466	.442

Table 12: Llama-3 family correlation results to human annotations on summary consistency. Max correlation within score range are underlined, max across all ranges are *italicized*, and max averages are **bolded**.

Model	Range	Pear.	Spear.	Kend.
Qwen2.5-3B	0 to 4	-.189 _{.085}	-.192 _{.062}	-.185 _{.051}
Qwen2.5-7B	0 to 4	<u>.396_{.024}</u>	<u>.392_{.021}</u>	<u>.356_{.021}</u>
Qwen2.5-14B	0 to 4	<u>.422_{.023}</u>	<u>.443_{.036}</u>	<u>.410_{.034}</u>
Contra. (7B-3B)	0 to 4	<u>.446_{.025}</u>	<u>.412_{.028}</u>	<u>.369_{.035}</u>
Contra. (14B-3B)	0 to 4	<u>.420_{.038}</u>	<u>.437_{.041}</u>	<u>.404_{.035}</u>
Qwen2.5-3B	1 to 5	-.047 _{.071}	-.063 _{.039}	-.060 _{.055}
Qwen2.5-7B	1 to 5	<u>.542_{.044}</u>	<u>.479_{.028}</u>	<u>.444_{.035}</u>
Qwen2.5-14B	1 to 5	<u>.406_{.051}</u>	<u>.417_{.049}</u>	<u>.394_{.037}</u>
Contra. (7B-3B)	1 to 5	<u>.549_{.044}</u>	<u>.474_{.033}</u>	<u>.441_{.030}</u>
Contra. (14B-3B)	1 to 5	<u>.477_{.051}</u>	<u>.481_{.045}</u>	<u>.459_{.031}</u>
Qwen2.5-3B	2 to 6	-.114 _{.045}	-.095 _{.052}	-.091 _{.062}
Qwen2.5-7B	2 to 6	<u>.489_{.035}</u>	<u>.435_{.032}</u>	<u>.404_{.033}</u>
Qwen2.5-14B	2 to 6	<u>.136_{.034}</u>	<u>.176_{.028}</u>	<u>.162_{.024}</u>
Contra. (7B-3B)	2 to 6	<u>.478_{.021}</u>	<u>.462_{.038}</u>	<u>.432_{.029}</u>
Contra. (14B-3B)	2 to 6	<u>.221_{.029}</u>	<u>.295_{.034}</u>	<u>.269_{.034}</u>
Qwen2.5-3B	3 to 7	-.127 _{.074}	-.126 _{.033}	-.121 _{.049}
Qwen2.5-7B	3 to 7	<u>.391_{.026}</u>	<u>.414_{.033}</u>	<u>.367_{.037}</u>
Qwen2.5-14B	3 to 7	<u>.224_{.046}</u>	<u>.204_{.027}</u>	<u>.188_{.030}</u>
Contra. (7B-3B)	3 to 7	<u>.382_{.028}</u>	<u>.389_{.036}</u>	<u>.345_{.036}</u>
Contra. (14B-3B)	3 to 7	<u>.309_{.033}</u>	<u>.269_{.043}</u>	<u>.246_{.034}</u>
<i>Average across all score ranges</i>				
Qwen2.5-3B		-.119	-.119	-.114
Qwen2.5-7B		.455	.430	.393
Qwen2.5-14B		.297	.310	.289
Contra. (7B-3B)		.464	.434	.397
Contra. (14B-3B)		.357	.370	.345

Table 13: Qwen2.5 family correlation results to human annotations on summary consistency. Max correlation within score range are underlined, max across all ranges are *italicized*, and max averages are **bolded**.

Range	Pearson	Spearman	Kendall
0-4	.419 _{.026}	.444 _{.040}	.370 _{.027}
1-5	.293 _{.027}	.282 _{.023}	.229 _{.035}
2-6	.402 _{.027}	.412 _{.050}	.336 _{.031}
3-7	.418 _{.034}	.402 _{.026}	.330 _{.034}
<i>Average across all score ranges</i>			
	.383	.385	.316

Table 14: Qwen2.5 14B-7B contrastive decoding results on summary coherence with 95% bootstrap confidence intervals. The 1-5 range shows significantly degraded performance compared to using a 3B assistant (Table 2), supporting the hypothesis that larger assistant models can penalize correct logits.