

QoS-QoE Translation with Large Language Model

Yingjie Yu
yyu69@illinois.edu
University of Illinois
Urbana-Champaign
Urbana, Illinois, USA

Mingyuan Wu
mw34@illinois.edu
University of Illinois
Urbana-Champaign
Urbana, Illinois, USA

Ahmadreza Eslaminia
ae15@illinois.edu
University of Illinois
Urbana-Champaign
Urbana, Illinois, USA

Lingzhi Zhao
lz26@illinois.edu
University of Illinois
Urbana-Champaign
Urbana, Illinois, USA

Kaizhuo Yan
kaizhuo2@illinois.edu
University of Illinois
Urbana-Champaign
Urbana, Illinois, USA

Klara Nahrstedt
klara@illinois.edu
University of Illinois
Urbana-Champaign
Urbana, Illinois, USA

A Appendix

A.1 Supervised Fine-tuning Instances

Figure 1 shows abridged chat-style supervised fine-tuning (SFT) instances for the two translation directions. Each instance is constructed from a single source-grounded relationship record. A system message instructs the model to output only the target field in valid JSON format. The user message contains the structured task input, including the task direction, contextual metadata, scenario, parameter mapping and definitions, source-grounded evidence, a partial history log, and a query. To create each training instance, we remove one time point from the original `history_log`. The observed side of that removed record is used as the query, and the corresponding target side is used as the assistant output. In the QoS→QoE setting, the query contains QoS value(s) and the assistant predicts the corresponding QoE value(s). In the QoE→QoS setting, the query contains QoE value(s) and the assistant predicts the corresponding QoS value(s).

A.2 Model Evaluation Details

For continuous-value prediction, we use MAPE and Accuracy@ δ . This section provides the parameter-specific tolerance configuration for Accuracy@ δ . We set the tolerance mode and δ value for each parameter based on its scale, unit, and semantic interpretation in the QoS-QoE translation setting. Because QoS and QoE parameters differ substantially in scale, units, and meaning, we use a parameter-specific δ rather than a single shared threshold. Table 1 reports the canonical parameter names, general definitions, and corresponding tolerance settings used in our evaluation. Notably, these general definitions are intended to provide a unified interpretation for evaluation, while the exact definition of the same parameter may vary slightly across different source papers.

Relative tolerances. Relative tolerances are used for scale-dependent rate-like quantities whose prediction error is more naturally interpreted proportionally than by a fixed absolute difference [11]. For `bitrate`, prior work also evaluates prediction quality using relative bitrate error, reflecting the fact that the same absolute error can correspond to very different levels of deviation across low and high bitrate ranges [10]. We use $\delta = 0.1$ for `bitrate`, which treats predictions within 10% of the ground truth as correct under Accuracy@ δ . This provides a moderate proportional tolerance for

Table 1: Canonical parameter-specific tolerance configuration used for Accuracy@ δ for model evaluation. Here, *rel* and *abs* denote relative and absolute tolerance, respectively.

Parameter	General definition	Mode	δ
<code>bitrate</code>	transmission rate	rel	0.1
<code>qoe_score</code>	generic QoE score	abs	0.3
<code>mos</code>	subjective opinion score	abs	0.3
<code>vmaf</code>	video quality metric	abs	3.0
<code>rebuffering_time</code>	stall duration	abs	0.5
<code>jitter</code>	packet delay variation	abs	5.0
<code>packet_loss_rate</code>	packet loss ratio	abs	0.01
<code>latency</code>	end-to-end delay	abs	10.0
<code>rebuffering_ratio</code>	stall ratio	abs	0.03
<code>bitrate_change</code>	bitrate switching magnitude	rel	0.15
<code>rebuffer_count</code>	stall count	abs	1.0
<code>psnr</code>	peak signal fidelity	abs	1.5
<code>visqol</code>	perceptual quality metric	abs	0.3
<code>ssim</code>	structural similarity	abs	0.02
<code>vifp</code>	visual information fidelity	abs	0.03
<code>vqm</code>	video quality metric	abs	0.05
<code>pesq</code>	perceptual speech quality	abs	0.25
<code>peaq_odg</code>	audio difference grade	abs	0.2
<code>playback_rate</code>	playback speed	abs	0.05
<code>frame_rate</code>	frames per second	abs	2.0
<code>snr</code>	signal-to-noise ratio	abs	2.0
<code>ssim_db</code>	SSIM in decibels	abs	1.5
<code>total_rebuffering_time</code>	cumulative stall duration	abs	1.0
<code>segment_duration</code>	segment duration	abs	0.25
<code>video_delay_p95_ms</code>	p95 video delay	abs	20.0
<code>acr_grade</code>	absolute category rating	abs	0.5

near-correct predictions while preserving separation between materially different bitrate values. We use the same tolerance mode for `bitrate_change`, with $\delta = 0.15$.

Absolute tolerances. Absolute tolerances are used when prediction error is more naturally interpreted on the original scale, making deviations directly meaningful in the native unit or rating scale [3, 11]. In our setup, we use this tolerance mode for bounded

```

117 {
118   "role": "system",
119   "content": "You are a QoS-QoE translation assistant. Given structured input, output only the target field in valid
120   JSON format."
121 },
122 {
123   "role": "user",
124   "content": {
125     "instruction": "Given the task, context, scenario, mapping, evidence, history_log and query, output the QoE
126     value(s) for the queried QoS value(s).",
127     "task": "qos_to_qoe",
128     "context": {
129       "domain": "video_streaming",
130       "protocol": [ "DASH" ],
131       "network_type": [ "wired" ],
132       "device_type": [ "desktop" ],
133       "video_type": [ "2d_vod" ],
134       "user_preference": "low_rebuffer"
135     },
136     "scenario": "A client streams ...",
137     "mapping": {
138       "qos_parameter": [ "initial_loading_delay_level", ... ],
139       "qos_parameter_definition": [ ... ],
140       "qoe_parameter": [ "mos" ],
141       "qoe_parameter_definition": [ ... ],
142     },
143     "evidence": {
144       "data_type": "equation",
145       "relationship": "MOS = 4.23 - 0.0672 L_{ti} - 0.742 L_{fr} - 0.106 L_{tr}",
146       "description": "The equation expresses MOS as a linear function ...",
147       "source": "...",
148     },
149     "history_log": [ ... ],
150     "query": {
151       "time_s": 0,
152       "qos": [ { "metric": "initial_loading_delay_level", "value": "2", ... } ]
153     }
154   },
155 },
156 {
157   "role": "assistant",
158   "content": { "qoe": [ { "metric": "mos", "value": "3.99" } ] }
159 }

```

(a) QoS→QoE SFT instance.

```

175 {
176   "role": "system",
177   "content": "You are a QoS-QoE translation assistant. Given structured input, output only the target field in valid
178   JSON format."
179 },
180 {
181   "role": "user",
182   "content": {
183     "instruction": "Given the task, context, scenario, mapping, evidence, history_log and query, output the QoS
184     value(s) for the queried QoE value(s).",
185     "task": "qoe_to_qos",
186     "context": {
187       "domain": "video_streaming",
188       "protocol": [ "DASH" ],
189       "network_type": [ "wired" ],
190       "device_type": [ "desktop" ],
191       "video_type": [ "2d_vod" ],
192       "user_preference": "low_rebuffer"
193     },
194     "scenario": "A client streams ...",
195     "mapping": {
196       "qos_parameter": [ "initial_loading_delay_level", ... ],
197       "qos_parameter_definition": [ ... ],
198       "qoe_parameter": [ "mos" ],
199       "qoe_parameter_definition": [ ... ],
200     },
201     "evidence": {
202       "data_type": "equation",
203       "relationship": "MOS = 4.23 - 0.0672 L_{ti} - 0.742 L_{fr} - 0.106 L_{tr}",
204       "description": "The equation expresses MOS as a linear function ...",
205       "source": "...",
206     },
207     "history_log": [ ... ],
208     "query": {
209       "time_s": 10.83,
210       "qoe": [ { "metric": "mos", "value": "3.3148" } ]
211     }
212   },
213 },
214 {
215   "role": "assistant",
216   "content": { "qos": [ { "metric": "initial_loading_delay_level", "value": "1" }, ... ] }
217 }

```

(b) QoE→QoS SFT instance.

Figure 1: Abridged chat-style SFT instances for bidirectional QoS-QoE translation. Long text fields and repeated entries are abbreviated with "..." for readability.

quality scores, objective quality metrics, time-based metrics, ratios, counts, and delay-related metrics.

Bounded subjective and perceptual quality scores are defined on fixed rating scales, so their prediction error is more naturally interpreted on the original scale than as proportional error [6, 7]. This category includes qoe_score ($\delta = 0.3$), mos ($\delta = 0.3$), $visqol$ ($\delta = 0.3$), $pesq$ ($\delta = 0.25$), $peaq_odg$ ($\delta = 0.2$), and acr_grade ($\delta = 0.5$). We set these δ values as small absolute tolerances on their scales to allow minor deviations while preserving meaningful quality differences.

Objective fidelity metrics are commonly evaluated on their own metric-specific scales [9, 12]. This category includes $vmaf$ ($\delta = 3.0$), $psnr$ ($\delta = 1.5$), $ssim$ ($\delta = 0.02$), $vifp$ ($\delta = 0.03$), and vqm ($\delta = 0.05$). We set these δ values as small absolute tolerances relative to the range and resolution of each metric.

Signal and native-unit rate metrics are expressed in engineering units such as decibels or frames per second [8, 11]. This category includes snr ($\delta = 2.0$), $ssim_db$ ($\delta = 1.5$), and $frame_rate$ ($\delta = 2.0$). We set these δ values as small absolute tolerances in their native units, allowing modest deviations while preserving materially different signal or temporal characteristics.

Time-based and delay-related metrics are measured in units such as seconds or milliseconds [2, 11]. This category includes $rebuffering_time$ ($\delta = 0.5$), $total_rebuffering_time$ ($\delta = 1.0$), $latency$ ($\delta = 10.0$), $jitter$ ($\delta = 5.0$), $segment_duration$ ($\delta = 0.25$), $video_delay_p95_ms$ ($\delta = 20.0$), and $playback_rate$ ($\delta = 0.05$). We set these δ values as small absolute tolerances in their native units to allow near-correct timing predictions without conflating materially different delay behavior.

Ratio- and count-based parameters are interpreted on their original scales because they represent bounded fractions or discrete events rather than continuously scaled magnitudes [4, 11]. This category includes $packet_loss_rate$ ($\delta = 0.01$), $rebuffering_ratio$

Table 2: Average rating and average confidence of each reviewer in the multi-reviewer evaluation stage.

LLM Reviewers	Avg. Rating	Avg. Confidence
Gemini-2.5-flash-lite	8.48	5.00
Claude-haiku-4-5-20251001	7.24	4.08
Grok-4.20-0309-reasoning	6.20	3.98

($\delta = 0.03$), and $rebuffer_count$ ($\delta = 1.0$). We set these δ values as small absolute tolerances, allowing minor deviations for bounded ratios and a one-event tolerance for counts.

A.3 Multi-LLM reviewer Evaluation Details

To improve data quality, each metadata-enriched record is evaluated by three LLM reviewers: Gemini-2.5-flash-lite [5], Claude-haiku-4-5-20251001 [1], and Grok-4.20-0309-reasoning [13]. Each reviewer returns a rating, a confidence score, and textual feedback. Table 2 summarizes the average rating and average confidence of each reviewer on the first-round LLM review, providing a coarse view of reviewer behavior in our pipeline. The results suggest that the initial extraction quality is already reasonably strong before re-evaluation. In our setup, Gemini-2.5-flash-lite gives the highest average rating and confidence, while Grok-4.20-0309-reasoning is comparatively more conservative. Claude-haiku-4-5-20251001 lies between the other two reviewers in both average rating and confidence.

A.4 Pipeline Prompt Details

The full prompts used in our pipeline are provided in the released codebase.¹ In particular, `relationship_extraction.txt` is used for source-grounded QoS-QoE relationship extraction from parsed

¹<https://github.com/yuyu6969/qos-qoe-translation/tree/main/src/prompts>

papers, `metadata_enrichment.txt` is used for contextual meta-data generation and normalization, and `data_evaluation.txt` is used in the multi-reviewer data evaluation stage. Together, these prompts define the task-specific instructions and structured output requirements used in the dataset construction pipeline.

References

- [1] Anthropic. 2025. Introducing Claude Haiku 4.5. <https://www.anthropic.com/news/claude-haiku-4-5> Official model announcement. Accessed: 2026-04-01.
- [2] J. J. Aranda, M. Cortés, J. Salvachúa, M. Narganes, and I. Martínez-Sarriegui. 2020. *The Quality for Service (QoS) Protocol*. RFC 8802. RFC Editor. <https://doi.org/10.17487/RFC8802>
- [3] Silvia Beddar-Wiesing, Alice Moallem-Oureh, Marie Kempkes, and Josephine M. Thomas. 2025. Absolute Evaluation Measures for Machine Learning: A Survey. <https://doi.org/10.48550/arXiv.2507.03392> arXiv:2507.03392 [cs.LG]
- [4] Florin Dobrian, Vyas Sekar, Asad Awan, Ion Stoica, Dilip Joseph, Aditya Ganjam, Jibin Zhan, and Hui Zhang. 2011. Understanding the Impact of Video Quality on User Engagement. *ACM SIGCOMM Computer Communication Review* 41, 4 (2011), 362–373. <https://doi.org/10.1145/2043164.2018478>
- [5] Google. 2025. Gemini 2.5 Flash-Lite. <https://ai.google.dev/gemini-api/docs/models/gemini-2.5-flash-lite> Official model documentation. Accessed: 2026-04-01.
- [6] International Telecommunication Union. 1996. *Methods for Subjective Determination of Transmission Quality*. Recommendation P.800. ITU-T. <https://www.itu.int/rec/T-REC-P.800-199608-1>
- [7] International Telecommunication Union. 2005. *Application Guide for Objective Quality Measurement Based on Recommendations P.862, P.862.1 and P.862.2*. Recommendation P.862.3. ITU-T. <https://www.itu.int/rec/T-REC-P.862.3-200511-S>
- [8] Alex MacKin, Aaron Zhang, and David R. Bull. 2015. A Study of Subjective Video Quality at Various Frame Rates. In *2015 IEEE International Conference on Image Processing (ICIP)*. 3407–3411. <https://doi.org/10.1109/ICIP.2015.7351436>
- [9] Margaret H. Pinson and Stephen Wolf. 2004. A New Standardized Method for Objectively Measuring Video Quality. *IEEE Transactions on Broadcasting* 50, 3 (2004), 312–322. <https://doi.org/10.1109/TBC.2004.834028>
- [10] Taveesh Sharma, Tarun Mangla, Arpit Gupta, Junchen Jiang, and Nick Feamster. 2023. Estimating WebRTC Video QoE Metrics Without Using Application Headers. In *Proceedings of the 2023 ACM Internet Measurement Conference (IMC '23)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3618257.3624828>
- [11] Juan Terven, Diana-Margarita Cordova-Esparza, Julio-Alejandro Romero-González, Alfonso Ramirez-Pedraza, and E. A. Chávez-Urbiola. 2025. A comprehensive survey of loss functions and metrics in deep learning. *Artificial Intelligence Review* 58, Article 195 (2025). <https://doi.org/10.1007/s10462-025-11198-7>
- [12] Abhinav K. Venkataramanan, Cosmin Stejerean, Ioannis Katsavounidis, and Alan C. Bovik. 2024. One Transform to Compute Them All: Efficient Fusion-Based Full-Reference Video Quality Assessment. *IEEE Transactions on Image Processing* 33 (2024), 509–524. <https://doi.org/10.1109/TIP.2023.3345227>
- [13] xAI. [n. d.]. Grok 4.20 0309 Reasoning. <https://docs.x.ai/developers/models/grok-4.20-0309-reasoning> Official model documentation. Accessed: 2026-04-01.