

Progressive Learning with Human Feedback for Personalized Adaptive Video Streaming

Zhaohui Jiang
Shanghai Jiao Tong University
Shanghai, China
jiangzhaohui@sjtu.edu.cn

Xuening Feng
Shanghai Jiao Tong University
Shanghai, China
cindyfeng2019@sjtu.edu.cn

Tianchi Huang
Tsinghua University
Beijing, China
mythkast@hotmail.com

Ruixiao Zhang
ByteDance Inc.
San Diego, United States
ruixiao.cs.zhang@gmail.com

Paul Weng
Duke Kunshan University
Kunshan, China
paul.weng@dukekunshan.edu.cn

Yifei Zhu*
Shanghai Jiao Tong University
Shanghai, China
yifei.zhu@sjtu.edu.cn

Abstract

Existing quality of experience (QoE)-driven adaptive bitrate (ABR) algorithms either fail to consider personalized QoE or rely on oversimplified QoE models, all resulting in unsatisfactory streaming experiences. Recognizing the wide existence of user feedback schemes in existing streaming applications, we introduce Q+, a framework leveraging progressively gathered personal user opinion scores from multiple interaction sessions for enhanced user-system alignment. Q+ first innovates QoE modeling by incorporating both pairwise ordinal and cardinal preferences constructed from scores. The capturing of both preferences ensures reliable and robust preference representation. Moreover, we design a monotonic neural network as the QoE model to capture the inherent monotonicity property in ABR services, improving model expressivity and generalization ability even with limited human feedback. To align the policy with the progressively updated QoE, we then develop a value-based reinforcement learning (RL) algorithm for bitrate control that integrates reward relabeling and calibrated prioritized experience replay. Extensive experiments reveal that Q+ consistently surpasses state-of-the-art rule-based, control-based, and RL-based baselines within only three sessions, improving QoE by 5.69% to 29.39% across diverse network conditions.

CCS Concepts

• Information systems → Multimedia streaming.

Keywords

Adaptive Video Streaming, Personalized Quality of Experience, Reinforcement Learning from Human Feedback, Online Learning

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
MM '25, Dublin, Ireland.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-2035-2/2025/10
<https://doi.org/10.1145/3746027.3755323>

ACM Reference Format:

Zhaohui Jiang, Xuening Feng, Tianchi Huang, Ruixiao Zhang, Paul Weng, and Yifei Zhu. 2025. Progressive Learning with Human Feedback for Personalized Adaptive Video Streaming. In *Proceedings of Proceedings of the 33rd ACM International Conference on Multimedia (MM '25)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3746027.3755323>

1 Introduction

Video streaming services, constituting about 60% of the Internet traffic [37], have become an indispensable part of modern life. In this context, Dynamic Adaptive Streaming over HTTP (DASH) [42] has emerged as a leading standard, segmenting videos into small chunks on the server side and encoding them at various bitrates. Adaptive bitrate (ABR) algorithms then select appropriate bitrates to ensure high-quality video experiences in fluctuating networks.

While Quality of Experience (QoE)-driven ABR algorithms aim to enhance user satisfaction by maximizing QoE [17, 30, 43, 48], most of them rely on oversimplified *general QoE* general models, which aggregate preferences of a wide population. However, studies [12, 21, 33] reveal significant individual variability in user expectations, highlighting the need for a *personalized QoE*-driven ABR framework that encompasses personalized QoE modeling and corresponding policy learning. Existing personalized approaches [19, 27, 34, 48] rely on manual user configuration and only depend on limited QoE parameters (usually fewer than three). This leads to narrow modeling spaces that fail to model the QoE accurately, resulting in the derived policies being misaligned with actual user experience. Furthermore, current methods for accurate QoE modeling [11, 17, 35, 45] rely on months of dataset preparation, which are too time-consuming and costly to encourage widespread user participation for personalized QoE modeling.

Therefore, to obtain personal user-aligned ABR policies in an accurate and scalable manner, we introduce a progressive personalized QoE-policy refinement process that aligns ABR policies with evolving personalized QoE models iteratively, where in each *training iteration*, the personalized QoE is first updated with new user feedback, and then the policy is fine-tuned to optimize with the revised QoE rewards. Such a progressive learning scheme can be seamlessly integrated into modern streaming platforms (e.g., Netflix [3], YouTube [5], Amazon Prime Video [2]) by leveraging user registration/login and explicit feedback (e.g., star ratings [2], pop-up surveys [5]). In these applications, during each *interaction session*,

which is a distinct period of continuous user engagement with the platform, users may rate the streaming quality for viewed video experiences at their discretion as new feedback.

For personalized QoE modeling, two challenges are: (i) learning reliable QoE from *cross-session feedback* (i.e., feedback from multiple interaction sessions), and (ii) ensuring model expressivity while avoiding misleading predictions with limited feedback. For the first challenge, our analysis reveals that traditional regression-based methods [11, 18, 33, 35] are unreliable for cross-session feedback due to score distortions. In contrast, pairwise ordinal and cardinal comparisons constructed from these scores demonstrate greater consistency, where the two kinds of comparisons capture preferences between video experiences and the strength of these preferences, respectively. We utilize both as learning sources to fully exploit the information inside the human feedback for reliable QoE modeling. For the second challenge, while MLP-based QoE models excel in expressivity [17, 46], they risk generating unreasonable predictions for unseen data, destabilizing policy performance if used as rewards [17]. To address this, we employ monotonic neural networks [36], embedding prior knowledge of monotonic relationships (e.g., rebuffering time vs. perceptual quality) to avoid unreasonable predictions and improve generalization for reward modeling.

For personalized QoE-driven policy learning, the challenge lies in ensuring effective and efficient training. To guarantee that repeating training iterations leads to better policies aligning with personal human expectations, we employ an *online* Reinforcement Learning from Human Feedback (RLHF) framework [9, 14, 16], where feedback is provided on video experiences rendered by the online ABR policy being trained. For sample-efficient fine-tuning under progressive QoE updates, we propose a value-based deep reinforcement learning (DRL) algorithm with two key designs: (i) rewards relabeling for historical trajectories in the replay buffer for reuse, and (ii) calibrated prioritized replay to speed up training.

Our main contributions can be summarized as follows:

- We introduce a user-friendly interaction scheme capturing personal preferences for ABR by gathering user feedback progressively across multiple interaction sessions (Section 3).
- We develop a reliable and efficient QoE modeling method using pairwise ordinal and cardinal preferences to learn from cross-session feedback, enhanced by monotonicity constraints for generalizable results (Section 4).
- We develop an efficient value-based DRL approach that fine-tunes policies along with personalized QoE model updates, progressively aligning with true user preferences (Section 5).
- Integrating our QoE modeling with the personalized QoE-driven ABR policy learning results in a whole framework, named **Q-learning-based Personalized Learning for User-centric ABR System (Q-PLUS, Q+)**, outperforming state-of-the-art rule-based, MPC-based, and DRL-based baselines across diverse network conditions (Section 6).

2 Background and Motivation

Scenarios. We focus on a common scenario where video platforms track user identities, typically through registration or login, enabling a progressive feedback collection to build datasets with *cross-session feedback* gathered from different interaction sessions.

Dataset. Our subsequent analysis employs S_{QoE}-IV [12], the latest published subjective QoE dataset, which is well-organized and comprises 1,350 video experiences, each rated by 31 human evaluators on a 1-100 scale. The 4.5-hour scoring process is split into three *testing sessions* for each evaluator to reduce human fatigue. To mitigate session-specific biases (e.g., mood or cognitive fluctuations) by cross-session score normalization, a *calibration phase* with 10 *overlapping video experiences* is included in each session, with each overlapping video experience presented in two sessions and yielding two scores. The structure of S_{QoE}-IV mirrors the streaming scenarios where *evaluators* as *users*, and *testing sessions* as *interaction sessions*.

Necessity of personalized QoE. While the cardinal aspect of *user heterogeneity*, where scores from different users show significantly varied numerical distributions [12, 17, 33, 46], is widely recognized, our analysis reveals its ordinal aspect as well: users exhibit diverse preference rankings. This is evidenced by Figure 1a, showing low SRCC and PLCC values (mostly <0.2) for different users' scores, indicating weak rank and numerical correlations between different users. Thus, QoEs derived from population-level cardinal scores (i.e., mean opinion score) [11, 18, 33, 35] or ordinal comparisons [17] all fail to align with personal perceptions, underscoring the need for personalized QoE modeling from individual feedback.

Necessity to learn personalized QoE progressively. Subjective score is an ideal form of human feedback for accurate QoE modeling. Still, existing methods [11, 17, 35, 45] rely on datasets that being pre-collected systematically (e.g., S_{QoE}-IV), requiring at least 50 video experiences with more than 25 minutes of time cost per user [33], making them impractical for scalable personalized QoE modeling for large populations. This motivates our *progressive* interaction scheme, where users can provide feedback gradually during many interaction sessions, enabling user-friendly feedback acquisition, therefore scalable for large populations.

Necessity to avoid learning directly from numerical scores. Cross-session subjective scores have significant numerical distortion, as indicated by Figure 1b via various distance measurements on overlapping video experiences's cross-session scores. It reveals that even the same video experience can receive vastly different scores across sessions, suggesting a higher score in one session does not necessarily indicate stronger preferences than lower scores in another. Such distortion is hard to mitigate in real-world scenarios, due to the impracticality of calibration phases. Therefore, traditional methods [11, 17, 18, 35] that model QoE by directly approximating numerical scores are fundamentally flawed with cross-session datasets.

Necessity to learn from pairwise preferences. Despite numerical distortions, there are strong rank and numerical correlations between overlapping video experiences's cross-session scores as shown in Figure 1c, evidenced by high SRCC and PLCC values (mostly >0.6). This motivates the use of cross-session feedback via rank-related preferences for reliable QoE modeling, leveraging preference-based learning techniques that learn from pairwise preferences [9, 13, 25]. Inspired by Huang et al. [17], instead of

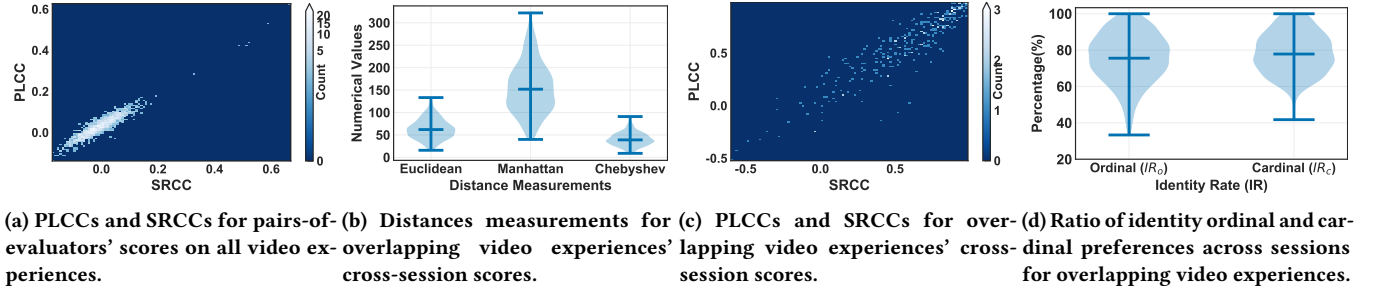


Figure 1: Data analysis of SQoE-IV's subjective scores from all evaluators. Figures 1a and 1c are 2D histograms with colors indicating counts. Figures 1b and 1d are violin plots with bars representing the max, min, and median values, and shadows indicating estimated distribution. Mathematical descriptions can be checked in Appendix A [4].

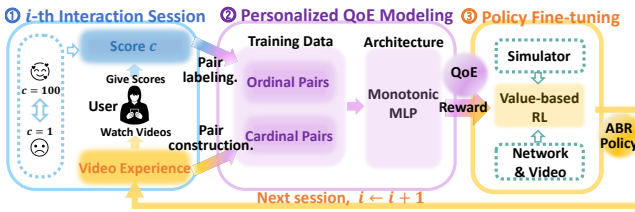


Figure 2: An overview of the proposed framework.

directly letting users provide preferences, a *preference* for a *pair-wise comparison* is determined by numerical responses generated from subjective scores, where an object with a higher response is preferred over one with a lower response. Specifically, an *ordinal comparison* [9] is constructed by a pair of video experiences, with the video experience as the object and its score as the response to determine the *ordinal preference*, and a *cardinal comparison* [13] is constructed by a pair of ordinal comparisons, with the ordinal comparison as the object and its absolute score difference (i.e., the strength of an ordinal preference) as the response to determine the *cardinal preference*. Figure 1d shows that most ordinal and cardinal comparisons receive the same preference labels across different sessions, demonstrating strong consistency and making them appropriate learning sources within the progressive scheme.

3 Progressive Learning Scheme for Adaptive Video Streaming

As shown in Figure 2, Q+ progressively refines the QoE and ABR policy by iterating a three-phase cycle: (i) During each user-system interaction sessions, users are prompted to rate their video experiences, (ii) Personalized QoE modeling based on collected user feedback (see Section 4), and (iii) ABR policy fine-tuning with the updated personalized QoE as a reward model (see Section 5).

Concretely, the system receives new human feedback \mathcal{D}_i^c during the i -th interaction session, consisting of video experiences $\{\tau\}$ generated from the latest ABR policy $\pi_{\theta_{i-1}}$ and corresponding scores $\{c\}$. The user's *training personalized QoE* R_ϕ is then updated to R_{ϕ_i} by training on feedback from all sessions \mathcal{D}^c (Section 4), where $\mathcal{D}^c = \cup_{j \leq i} \mathcal{D}_j^c$. Since the *ground-truth personalized QoE* R^* is inaccessible due to the lack of a systematic subjective test, our

Algorithm 1 Overview: Q+ with Progressive Interaction Scheme

```

1: Initial policy:  $\pi_{\theta_0}$ 
2: Dataset saving user's feedback:  $\mathcal{D}^c \leftarrow \emptyset$ 
3: Total number of interaction sessions:  $N_I$ 
4: for  $i \leftarrow 1$  to  $N_I$  do
5:    $\triangleright$  Repeat training iterations.
6:    $\mathcal{D}_i^c \leftarrow \emptyset$ 
7:   while User  $e$  is using the video streaming service do
8:     Video experiences  $\tau$  generated with  $\pi_{\theta_{i-1}}$ 
9:     if User  $e$  provides a score  $c$  for  $\tau$  then
10:       $\mathcal{D}_i^c \leftarrow \mathcal{D}_i^c \cup \{(\tau, c)\}$ 
11:    end if
12:  end while
13:   $\mathcal{D}^c \leftarrow \mathcal{D}^c \cup \mathcal{D}_i^c$   $\triangleright$  Cross-session feedback.
14:   $\triangleright$  QoE modeling (Algorithm 2 in [4].)
15:  Train personalized QoE  $R_{\phi_i}$  for user  $e$  with  $\mathcal{D}^c$ 
16:   $\triangleright$  QoE-driven policy fine-tuning (Algorithm 3 in [4].)
17:  Obtain  $\pi_{\theta_i}$  by fine-tuning  $\pi_{\theta_{i-1}}$  with  $R_{\phi_i}$ 
18: end for
19: return  $\pi_{N_I}$ 

```

policy learning adaptively adjusts $\pi_{\theta_{i-1}}$ to π_{θ_i} by maximizing R_{ϕ_i} . Repeating such iterations as outlined in Algorithm 1, the policy is expected to align with user preferences gradually, ultimately enhancing user satisfaction.

4 Progressive Personalized QoE Modeling

Formulating a video experience τ with H chunks as $(s_0, a_0, s_1, a_1, \dots, s_H, a_H)$, where s_t refers to the state at timestep t containing information for bitrate selection, and a_t is the target bitrate selected by an ABR policy π given s_t for downloading the t -th chunk. Our QoE modeling aims to learn a QoE model R_ϕ to predict scalar rewards for states such that $R_\phi(\tau) = \frac{1}{H} \sum_{t=0}^H R_\phi(s_t)$ can reflect user preferences for different video experiences. We explain our QoE modeling from the perspective of model training and model architecture design in Section 4.1 and Section 4.2, respectively.

4.1 Preference-Based QoE Training

To maximize the utility of limited human feedback in the progressive learning scheme, we train our QoE model using both ordinal

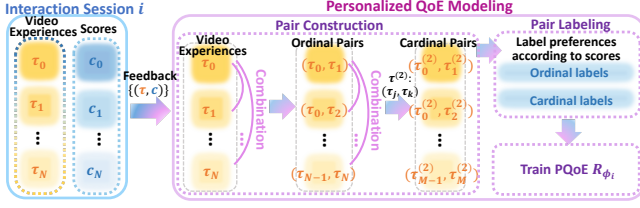


Figure 3: Illustration for the personalized QoE modeling.

and cardinal pairwise comparisons, leveraging the consistency of these preferences across sessions. Our approach is adapted from reward training in preference-based RLHF methods [9, 13], specifically tailored for progressive learning with cross-session datasets.

Pair construction and labeling. An *ordinal comparison* [9] is constructed by a pair of video experiences $\tau^{(2)} = (\tau_j, \tau_k)$, where j and k are used in subscript to distinguish different objects. The corresponding *ordinal preference* $y_o \in \{=, >, <\}$ describes which one of the video experiences is being preferred over another:

$$y_o(\tau_j, \tau_k) = \begin{cases} =, & \text{if } |c_j - c_k| < \delta_o, \\ >, & \text{if } c_j > c_k + \delta_o, \\ <, & \text{if } c_j < c_k - \delta_o, \end{cases} \quad (1)$$

where $>$ (resp. $<$) means τ_j (resp. τ_k) is preferred than τ_k (resp. τ_j), = means equal preference, and δ_o is a positive constant indicating the allowed difference for equal ordinal preference to account for label noise. In addition, a *cardinal comparison* [13] is constructed by a pair of ordinal comparisons $(\tau_j^{(2)}, \tau_k^{(2)})$. The corresponding *cardinal preference* $y_c \in \{=, >, <\}$ describes which ordinal comparison has stronger preference strength based on absolute score differences $|c_j - c_k|$ (denoted as $\Delta_c(\tau^{(2)})$):

$$y_c(\tau_j^{(2)}, \tau_k^{(2)}) = \begin{cases} =, & \text{if } |\Delta_c(\tau_j^{(2)}) - \Delta_c(\tau_k^{(2)})| < \delta_c, \\ >, & \text{if } \Delta_c(\tau_j^{(2)}) > \Delta_c(\tau_k^{(2)}) + \delta_c, \\ <, & \text{if } \Delta_c(\tau_j^{(2)}) < \Delta_c(\tau_k^{(2)}) - \delta_c, \end{cases} \quad (2)$$

where δ_c is a constant indicating the allowed difference for equal cardinal preference. As illustrated in Figure 3, ordinal and cardinal comparisons are formed with video experiences within the same session, and preferences are labeled based on their associated scores.

Training method. Given preference labels, the *target ordinal (resp. cardinal) preference probability* P_o^* (resp. P_c^*) for an ordinal (resp. cardinal) comparison, representing the true probability of one object being preferred over another, is derived as Equation (3) (resp. Equation (4)):

$$P_o^*(\tau_j, \tau_k) = \begin{cases} (0.5, 0.5), & \text{if } \tau_j = \tau_k, \\ (1, 0), & \text{if } \tau_j > \tau_k, \\ (0, 1), & \text{if } \tau_j < \tau_k, \end{cases} \quad (3)$$

$$P_c^*(\tau_j^{(2)}, \tau_k^{(2)}) = \begin{cases} (0.5, 0.5), & \text{if } \tau_j^{(2)} = \tau_k^{(2)}, \\ (1, 0), & \text{if } \tau_j^{(2)} > \tau_k^{(2)}, \\ (0, 1), & \text{if } \tau_j^{(2)} < \tau_k^{(2)}. \end{cases} \quad (4)$$

Then, following Christiano et al. [9], the *modeled ordinal preference probability* is derived using the reward model R based on the

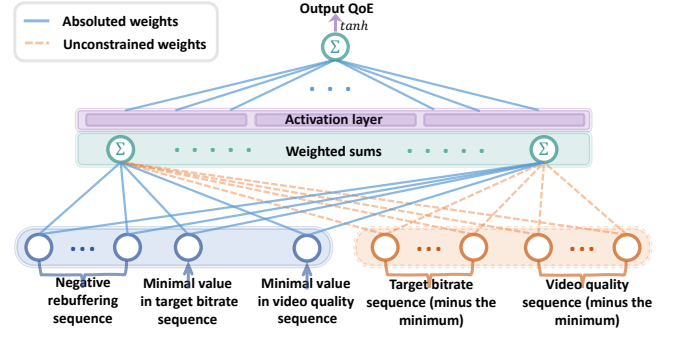


Figure 4: Monotonic MLP as the QoE model, illustrated with one hidden layer for clarity.

Bradley-Terry model [8]:

$$\mathbb{P}(\tau_j > \tau_k; R) = \frac{\exp(R(\tau_j))}{\exp(R(\tau_j)) + \exp(R(\tau_k))}. \quad (5)$$

This defines the *ordinal loss* \mathcal{L}_ϕ^o , which optimizes R_ϕ by minimizing the cross-entropy loss between $\mathbb{P}(\tau_j > \tau_k; R_\phi)$ and $P_o^*(\tau_j, \tau_k)$:

$$\mathcal{L}_\phi^o(\tau^{(2)}) = - \left[P_o^*[0] \log \mathbb{P}(\tau_j > \tau_k; R_\phi) + P_o^*[1] \log \mathbb{P}(\tau_j < \tau_k; R_\phi) \right], \quad (6)$$

where $P_o^*[0]$ and $P_o^*[1]$ refer to the two value of $P_o^*(\tau_j, \tau_k)$.

Following Feng et al. [13], the *modeled cardinal preference probability* is also formulated with the Bradley-Terry model, based on the reward difference $|R(\tau_j) - R(\tau_k)|$ as $\Delta_R(\tau^{(2)})$:

$$\mathbb{P}(\tau_j^{(2)} > \tau_k^{(2)}; R) = \frac{\exp(\Delta_R(\tau_j^{(2)}))}{\exp(\Delta_R(\tau_j^{(2)})) + \exp(\Delta_R(\tau_k^{(2)}))}. \quad (7)$$

This leads to the *cardinal loss* \mathcal{L}_ϕ^c that optimizes R_ϕ by minimizing a cross entropy loss between $\mathbb{P}(\tau_j^{(2)} > \tau_k^{(2)}; R_\phi)$ and $P_c^*(\tau_j^{(2)}, \tau_k^{(2)})$:

$$\mathcal{L}_\phi^c(\tau_j^{(2)}, \tau_k^{(2)}) = - \left[P_c^*[0] \log \mathbb{P}(\tau_j^{(2)} > \tau_k^{(2)}; R_\phi) + P_c^*[1] \log \mathbb{P}(\tau_j^{(2)} < \tau_k^{(2)}; R_\phi) \right]. \quad (8)$$

To effectively combine \mathcal{L}_ϕ^o (Equation (6)) and \mathcal{L}_ϕ^c (Equation (8)) for cross-session feedback, we specifically design our final QoE loss \mathcal{L}_ϕ that enforces concurrent optimization of ordinal and cardinal preferences on sampled cardinal comparisons:

$$\mathcal{L}_\phi = \mathbb{E}_{(\tau_j^{(2)}, \tau_k^{(2)}) \sim \mathcal{D}^c} \left[\mathcal{L}_\phi^o(\tau_j^{(2)}) + \mathcal{L}_\phi^c(\tau_k^{(2)}) + \mathcal{L}_\phi^c(\tau_j^{(2)}, \tau_k^{(2)}) \right]. \quad (9)$$

Pseudo-code for our QoE modeling can be checked in Appendix E.2 [4].

4.2 Monotonic MLP as QoE Architecture

To obtain robust QoE outputs, we introduce *monotonicity constraints* to ensure monotonic relationships between input features and QoE predictions, preventing unreasonable QoE improvements under worsening conditions, thereby aligning with domain knowledge.

Model inputs. The input for our QoE model at timestep t , denoted as s_t , contains recent K chunks' rebuffering time $\mathbf{u}_t = \{u_{t-K+1}, \dots, u_t\}$, target bitrate $\mathbf{q}_t = \{q_{t-K+1}, \dots, q_t\}$, and video quality $\mathbf{v}_t = \{v_{t-K+1}, \dots, v_t\}$ measuring by VMAF [1].

Monotonicity constraints. Then, given two states s'_t (with features $\{\mathbf{u}_{t'}, \mathbf{q}_{t'}, \mathbf{v}_{t'}\}$) and s''_t (with features $\{\mathbf{u}_{t''}, \mathbf{q}_{t''}, \mathbf{v}_{t''}\}$):

- (1) *Monotonicity in rebuffering time:* Given $\mathbf{q}_{t'} = \mathbf{q}_{t''}$, $\mathbf{v}_{t'} = \mathbf{v}_{t''}$, if $\exists j \in \{t - K + 1, \dots, t\}$ such that $\mathbf{u}_{t'}[j] > \mathbf{u}_{t''}[j]$ and $\forall k \neq j, \mathbf{u}_{t'}[k] = \mathbf{u}_{t''}[k]$, then $R_\phi(s'_t) \leq R_\phi(s''_t)$.
- (2) *Monotonicity in target bitrate:* Given $\mathbf{u}_{t'} = \mathbf{u}_{t''}$, $\mathbf{v}_{t'} = \mathbf{v}_{t''}$, if $\min(\mathbf{q}_{t'}) < \min(\mathbf{q}_{t''})$ and $(\mathbf{q}_{t'} - \min(\mathbf{q}_{t'})) = (\mathbf{q}_{t''} - \min(\mathbf{q}_{t''}))$, then $R_\phi(s'_t) \leq R_\phi(s''_t)$.
- (3) *Monotonicity in video quality:* Given $\mathbf{u}_{t'} = \mathbf{u}_{t''}$, $\mathbf{q}_{t'} = \mathbf{q}_{t''}$, if $\min(\mathbf{v}_{t'}) < \min(\mathbf{v}_{t''})$ and $(\mathbf{v}_{t'} - \min(\mathbf{v}_{t'})) = (\mathbf{v}_{t''} - \min(\mathbf{v}_{t''}))$, then $R_\phi(s'_t) \leq R_\phi(s''_t)$.

For rebuffering time, constraint (i) reflects that increasing it for any specific chunk has a non-positive monotonic impact (i.e., either degrade or maintain) on QoE, aligning with empirical evidence that even a single rebuffering incident significantly disrupts user experience [10, 24]. For target bitrate or video quality, increased values of a single chunk do not guarantee a monotonic QoE response, as they can affect playback smoothness and therefore perceptual continuity [7]. Therefore, constraints (ii) and (iii) enforce non-negative monotonicity only for the minimum values in the feature sequence, allowing users to express preferences regarding quality fluctuations.

Monotonicity implementation. Our QoE model, named MonMLP, employs a monotonic neural network [36], integrating the monotonicity constraints seamlessly into standard MLPs by restricting specific first-layer weights and all subsequent-layer weights to be non-negative, as depicted in Figure 4. Concretely, for a state $\{\mathbf{u}_t, \mathbf{q}_t, \mathbf{v}_t\}$, MonMLP assigns absolute non-negative trainable weights to $\{-\mathbf{u}_t, \min(\mathbf{q}_t), \min(\mathbf{v}_t)\}$, ensuring monotonicity. The remaining state information, $\{\mathbf{q}_t - \min(\mathbf{q}_t), \mathbf{v}_t - \min(\mathbf{v}_t)\}$, is also included as inputs but with unconstrained first-layer weights. The QoE predictions are normalized to $(-1, 1)$ using a tanh function. Further implementation details are provided in Appendix E.1 [4].

5 Progressive Policy Learning

With the learning scheme overviewed in Section 3, we develop an RL-based method to effectively adapt $\pi_{\theta_{i-1}}$ (optimized for $R_{\phi_{i-1}}$) into π_{θ_i} in each training iteration, maximizing the updated QoE R_{ϕ_i} by treating QoE as the reward function. Section 5.1 reviews the fundamental algorithm design, serving as the backbone, then Section 5.2 details two novel modifications to the backbone algorithm for efficient progressive fine-tuning across training iterations.

5.1 Value-Based RL as Backbone

Given the frequent policy fine-tuning in each training iteration, sample efficiency (i.e., minimizing environment timesteps to achieve target performance) is crucial for scalability across a large user base. Value-based methods, with their inherent advantage in sample efficiency through trajectory storage and reuse via replay buffers, have become the preferred choice for discrete action tasks [39]. Thus, our ABR policy learning backbone is value-based that trains

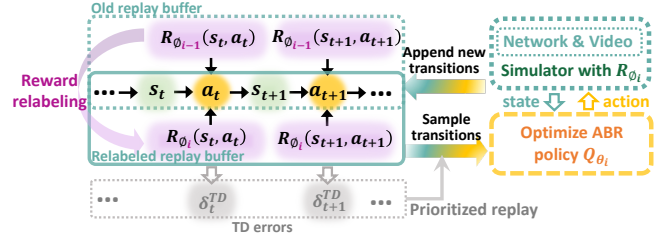


Figure 5: Illustration for the policy fine-tuning with personalized QoE model R_{ϕ_i} learned after the i -th interaction session.

a Q-network Q_θ by minimizing \mathcal{L}_θ based on TD errors δ^{TD} [31]:

$$\delta^{TD}(s_t, a_t, r_t, s_{t+1}) = r_t + \max_{a'} \gamma Q_\theta(s_{t+1}, a') - Q_\theta(s_t, a_t), \quad (10)$$

$$\mathcal{L}_\theta = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{D}^\tau} \left[\left(\delta^{TD}(s_t, a_t, r_t, s_{t+1}) \right)^2 \right]. \quad (11)$$

Here, transitions (s_t, a_t, r_t, s_{t+1}) are sampled from a *replay buffer* \mathcal{D}^τ storing historical trajectories τ , where r_t is the reward for (s_t, a_t) , and the target network Q_θ is updated less frequently than the online network Q_θ . The policy π_θ is derived as $\pi_\theta(s_t) = \arg \max_a Q_\theta(s_t, \cdot)$.

5.2 Progressive Policy Fine-Tuning with Updated Reward

To avoid bad beginning performance, we employ a policy pre-trained with general QoE as the initial policy π_{θ_0} . For the subsequent policy Q_{θ_i} optimizing the updated reward model R_{ϕ_i} , which is fine-tuned from $Q_{\theta_{i-1}}$ without reinitializing θ , we propose two key modifications¹ to adapt the backbone algorithm for progressive fine-tuning, ensuring effectiveness and efficiency (see Figure 5):

Reward relabeling. Unlike standard settings with fixed reward functions, our progressive training scheme dynamically updates the reward function. Directly reusing historical samples labeled with obsolete rewards $R_{\phi_{<i}}$ compromises policy performance due to misleading TD errors (Equation (10)) when optimizing Q_{θ_i} . On the other hand, cleaning \mathcal{D}^τ as an empty buffer wastes historical samples and reduces the sample diversity for future training, necessitating prolonged training to achieve target performance. To address these issues, we relabel historical samples in \mathcal{D}^τ by updating r_t to $R_{\phi_i}(s_t, a_t)$, which can be beneficial to both the policy performance and sample efficiency when fine-tuning the policy.

Calibrated prioritized replay. Building upon our reward relabeling approach, we implement priority relabeling [38], where transition priorities are updated alongside reward relabeling. This allows transitions (s_t, a_t, r_t) with larger TD errors (Equation (10)) under the new reward scheme to be sampled more frequently. Combined with our fine-tuning strategy that retains the old policy parameters θ_{i-1} at the beginning of the i -th training iteration, this calibrated prioritized replay significantly accelerates training.

6 Experiments

We first validate our personalized QoE modeling using SQoE-IV in Section 6.1, showing its superior alignment with human perception.

¹Additional implementation details can be checked in Appendix F [4].

Table 1: Evaluation results for QoE modeling methods, averaged over all evaluators in the SQoE-IV dataset.

General QoE	Methods	Evaluation Metrics (mean \pm std)			
		IR_o	IR_c	SRCC	PLCC
Intuitive	MPC [44]	0.65 \pm 0.04	0.55 \pm 0.04	0.42 \pm 0.11	0.27 \pm 0.10
	BOLA [41]	0.70 \pm 0.05	0.64 \pm 0.05	0.54 \pm 0.13	0.53 \pm 0.10
	Pensieve [30]	0.70 \pm 0.05	0.63 \pm 0.05	0.54 \pm 0.13	0.53 \pm 0.12
	Puffer [43]	0.61 \pm 0.06	0.54 \pm 0.03	0.29 \pm 0.15	0.08 \pm 0.15
Regression	BSQI [11]	0.51 \pm 0.01	0.50 \pm 0.02	0.02 \pm 0.04	0.06 \pm 0.04
	Comyco-Lin [18]	0.63 \pm 0.06	0.56 \pm 0.04	0.35 \pm 0.16	0.09 \pm 0.15
\mathcal{L}_ϕ^o (Equation (6))	Jade-Lin [17]	0.71 \pm 0.05	0.68 \pm 0.06	0.56 \pm 0.13	0.58 \pm 0.11
Personalized QoE (With Ablations)		Evaluation Metrics (mean \pm std)			
	Models	IR_o	IR_c	SRCC	PLCC
Regression	Linear	0.50 \pm 0.11	0.58 \pm 0.06	0.01 \pm 0.31	0.00 \pm 0.32
	MLP	0.50 \pm 0.02	0.50 \pm 0.01	0.00 \pm 0.05	0.01 \pm 0.05
	MonMLP	0.50 \pm 0.02	0.50 \pm 0.01	0.01 \pm 0.05	0.01 \pm 0.06
\mathcal{L}_ϕ^o (Equation (6))	Linear	0.72 \pm 0.05	0.69 \pm 0.06	0.59 \pm 0.13	0.60 \pm 0.11
	MLP	0.84 \pm 0.03	0.81 \pm 0.05	0.83 \pm 0.05	0.79 \pm 0.05
	MonMLP	0.81 \pm 0.05	0.78 \pm 0.04	0.77 \pm 0.09	0.76 \pm 0.05
\mathcal{L}_ϕ (Equation (9))	Linear	0.72 \pm 0.05	0.70 \pm 0.06	0.59 \pm 0.13	0.61 \pm 0.11
	MLP	0.89\pm0.04	0.93\pm0.03	0.88 \pm 0.07	0.88\pm0.03
	MonMLP (Ours)	0.89\pm0.03	0.88 \pm 0.05	0.90\pm0.06	0.85 \pm 0.05

Next, in Section 6.2.1, we evaluate our progressive policy learning framework using both the personalized QoE metric from Section 6.1 and traditional metrics (e.g., rebuffering time, video quality). All experiments use three random seeds per trial, with results averaged for robust statistics. See Appendix B [4] for hyperparameters.

6.1 Evaluations of QoE Modeling in Q+

For each evaluator’s personalized QoE, their subjective scores are randomly split with 80% for training and 20% for evaluation.

Evaluation metrics. For each testing session of an evaluator, we evaluate QoE predictions against subjective scores using the ordinal identity rate IR_o and cardinal identity rate IR_c , measuring the ratio of identical ordinal preferences (Equation (1)) and cardinal preference (Equation (2)), respectively, among all constructed pairwise comparisons. We also calculate the SRCC for rank correlation and the PLCC for numerical correlation. The reported results in the following are averaged over all sessions and evaluators.

Baselines and ablations. As outlined in Table 1, we establish two sets of experiments² to validate our QoE modeling method comprehensively: (i) To compare with other general QoE, we select representative baselines, including those intuitively configured [30, 41, 43, 44], and trained on population-level subjective scores [11, 18] or ordinal comparisons [17]. (ii) Given the lack of prior work specifically focused on learning personalized QoE, we validate our personalized QoE modeling through ablation studies. Specifically, we ablate our training objective \mathcal{L}_ϕ (Equation (9)) by replacing it with two alternatives: traditional regression with MSE minimization [18] and learning from ordinal comparisons [17] with Equation (6). Additionally, we ablate our MonMLP architecture by substituting it with a linear model [17, 18] and a standard MLP [17].

²See implementation details about baselines and ablations in Appendix C [4].

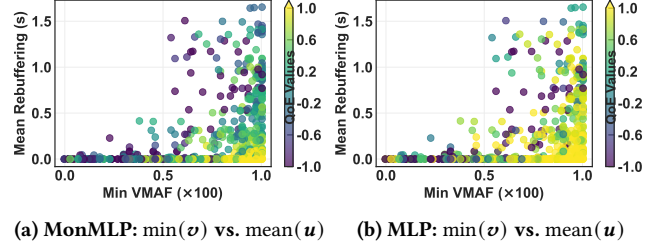


Figure 6: QoE predictions from a MLP or MonMLP trained on an evaluator’s subjective scores with \mathcal{L}_ϕ (Equation (9)).

Performance analysis. As shown in Table 1, our personalized QoE modeling significantly outperforms others, with findings summarized below: (i) Personalized QoE modeling with cross-session feedback requires learning from preferences with \mathcal{L}_ϕ^o or our \mathcal{L}_ϕ , leading to performance exceeding SOTA general QoE. Regression fails due to numerical distortion. (ii) According to ablations on training objective, introducing cardinal loss \mathcal{L}_ϕ^c for QoE training (i.e., \mathcal{L}_ϕ^o vs. \mathcal{L}_ϕ) enhances IR_c while also improving IR_o , SRCC, and PLCC. (iii) According to ablations about model architecture, MonMLP and MLP outperform the Linear model, with MonMLP generating more reasonable predictions that adhere to monotonicity constraints, as demonstrated in Figure 6³. For example, in Figure 6b, MLP predicts unreasonably high values (yellow points) for video experiences with large rebuffering (upper right zone) or poor video quality (lower left zone), while MonMLP produces smooth, monotonic predictions, demonstrating enhanced generalization in Figure 6a, further benefits RL policy training as shown in Section 6.2.4.

6.2 Evaluations of Q+ Policy Learning

6.2.1 Experimental Setups.

Simulator. We conduct policy training experiments in Park [29], a trace-driven virtual player environment for realistic adaptive streaming simulations. Our video dataset includes 83 complete videos from YouTube [18], covering diverse genres such as gaming, cinema, music, sports, and television. Each video is encoded at nine target bitrates ($\{235, 375, 560, 750, 1050, 1750, 2350, 3000, 4300\}$ kbps) and segmented into 4-second chunks, totaling 4,702 segments. Video lengths vary significantly, ranging from 8 to 231 chunks, ensuring diverse streaming scenarios. Our network traces include three trace collections covering varied network conditions: fixed broadband traces (FCC) and dynamic mobile traces (3G/HSDPA) from Pensieve [30], and heavy-tailed traces from Puffer [43]. Videos and network traces are randomly split with 80% for training and 20% for evaluation.

Information in a policy state s_t contains the VMAF of the last selected chunk, current buffer size, estimated bandwidth, delay time for receiving the last chunk, remaining chunks to complete the video, sizes of the next chunk across target bitrates, and the corresponding VMAFs.

To enable large-scale online RL training, the personalized QoEs trained with the whole SQoE-IV dataset using our method, which

³Supplementary figures for additional evaluators can be checked in Appendix H.1 [4].

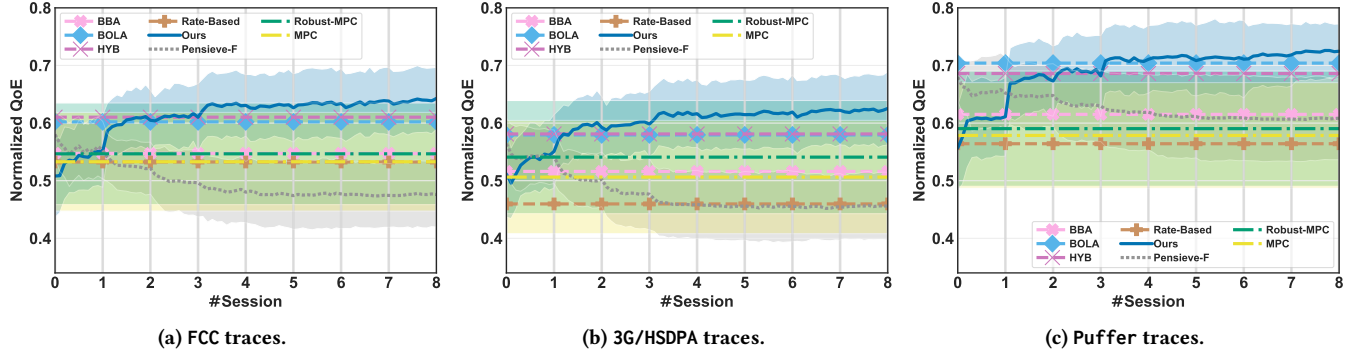


Figure 7: Comparisons Q+ with baselines. Lines and shadows represent average performance and standard deviation, respectively. Vertical lines mark interaction sessions where R_ϕ is updated with $\#F = 5$ new scores, triggering subsequent policy fine-tuning.

outperformed all other approaches as shown in Section 6.1, serve as the ground-truth QoEs for each evaluator e , denoted as R_e^* . Following previous works [23, 26, 28], we synthesize quality assessments by treating $R_e^*(\tau)$ as evaluator e 's numerical rating for segment τ . Such simulation maintains human perceptual alignment while avoiding the prohibitive costs of continuous human annotation.

Performance evaluation. Following previous works [17], the performance of an ABR policy is measured by normalized R_e^* by linear mapped to $[0,1]$, and we show the averaged performance over evaluators as the overall QoE performance. In addition, we also examine Q+ in conventional metrics, e.g., rebuffering time, video quality, and smoothness. Due to the computational cost for multi-user multi-session simulation, 8 evaluators (25% of the SQoE-IV corpus) are randomly selected, and we run 3 seeds for each user.

6.2.2 Compare with Baselines.

Baseline methods. Given the lack of prior ABR methods for the progressive learning scheme, we adapt representative baselines to fit our framework. These adaptations are summarized as follows⁴: (i) For learning-based methods, we extend Pensieve [30], a policy-gradient framework, to our progressive personalized QoE-driven setting, termed Pensieve-F. This extension uses the same experimental configurations as Q+ (see Section 6.2.1) to ensure a fair comparison. (ii) For Model Predictive Control (MPC) methods, we evaluate both the standard MPC and its conservative variant Robust-MPC [44]. Both methods use harmonic-mean bandwidth prediction and finite-horizon QoE optimization. They use the same last QoE model as Q+'s, which is trained with all personal feedback. (iii) For rule-based methods, we consider both buffer-based and rate-based methods. For buffer-based methods, we compare with BBA [20] that determines bitrate via linear interpolation based on buffer levels, and BOLA [41] that uses Lyapunov optimization to select bitrates constrained by buffer occupancy. For rate-based methods, we compare with Rate-Based [22] that selects the maximum bitrate below the predicted throughput estimated via harmonic mean. In addition, we consider HYB [6] that selects the maximum bitrate with chunk size below the estimated available file size based on buffer occupancy and predicted throughput.

⁴Further implementation details can be checked in Appendix D.1 [4].

Performance analysis. For baselines involving personalized QoE-modeling (i.e., Q+, Pensieve-F, MPC, Robust-MPC), we evaluate them under varying *per-session feedback quantities* ($\#F$). Figure 7 compares their QoE performance under a moderate feedback setting ($\#F=5$), where Q+ consistently surpasses all baselines after only 3 sessions with 15 feedbacks. Despite initial lag due to cold-start personalized QoE modeling, Q+ achieves sustained improvements through iterative feedback refinement. Although the initial general QoE-driven policy of Pensieve-F's outperforms Q+'s, Pensieve-F degrades progressively, while Q+'s value-based policy design enhances both the robustness against imperfect reward signals [15, 32, 40, 47] and training efficiency. As for MPC variants, they exhibit higher variance and lower average performance compared to Q+. Their vulnerability to imperfect QoE models highlights the limitations of model-based approaches within the progressive learning scheme. Additionally, Figure 8 compares the breakdown performance of the final policies shown in Figure 7⁵, where Q+ consistently performs within the acceptable upper right regions.

6.2.3 Sensitivity Analysis for Feedback Quantity.

Settings. To study the impact of feedback quantity ($\#F$) on Q+'s adaptation, we conduct experiments with fixed $\#F$ values (3, 5, 7, 10) and a variable scenario where $\#F \sim Normal(\mu = 6, \sigma = 3)$ (samples below 3 are set to 3). The minimum $\#F = 3$ is chosen because \mathcal{L}_ϕ (Equation (9)) requires at least three scores to construct a cardinal comparison. The value $\#F = 10$ represents an ideal, highly engaged user scenario, while 5 and 7 reflect practical feedback levels. $Normal(\mu = 6, \sigma = 3)$ simulates spontaneous user feedback.

Performance analysis. As shown in Figure 9a, Q+ achieves better performance with increased user feedback, surpassing baseline methods. Even with the minimum $\#F = 3$, required to construct one cardinal pair per session, Q+ performs on par with the best-performing baseline, demonstrating its robustness under limited feedback. Furthermore, larger feedback quantities unlock more significant performance improvements, and the distributional $\#F \sim Normal(6, 3)$ setting also shows stable performance, indicating that occasional limited feedback does not hinder overall performance.

⁵Additional figures for 3G/HSDPA and other $\#F$ can be checked in Appendix G.1.2 [4].

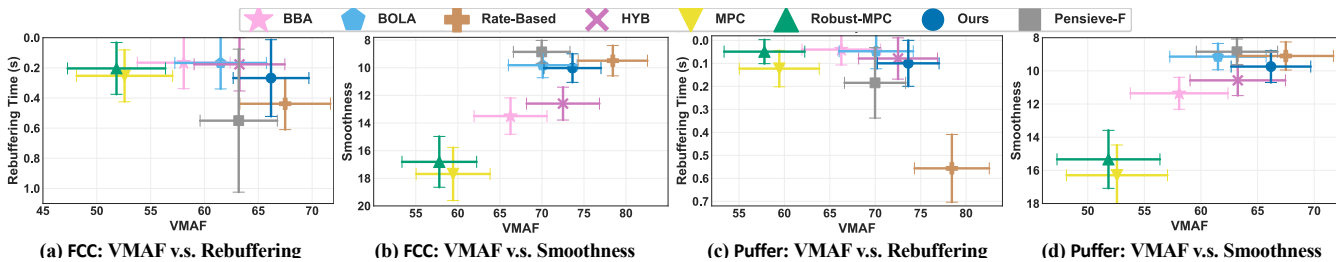


Figure 8: Breakdown performance for the last policies trained with $\#F = 5$. Smaller rebuffering and smoothness values, and larger VMAF values are desired.

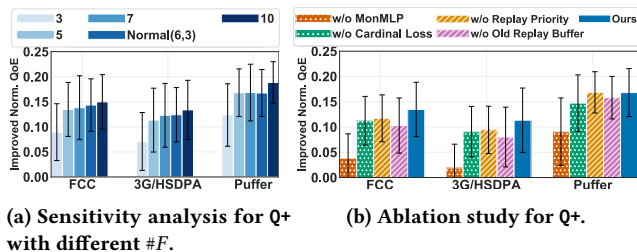


Figure 9: Relative normalized QoE improvement after the final training iteration compared to the initial policy. See complete training curves in Appendix G.2 and G.3 [4].

6.2.4 Ablation Studies for Q+.

Settings. We conduct four ablation experiments to evaluate Q+: two focus on personalized QoE modeling’s impact on progressive policy learning, and two examine our policy training design. For personalized QoE modeling, we compare our method with a variant replacing MonMLP with a standard MLP (w/o MonMLP) and another removing the cardinal loss from \mathcal{L}_ϕ (Equation (9)) and use \mathcal{L}_ϕ^o (Equation (6)) instead (w/o Cardinal Loss). For policy training, we test cleaning the old replay buffer and recollecting new trajectories per training iteration (w/o Old Replay Buffer) and replacing priority replay with uniform replay (w/o Prioritized Replay). See more details for ablations’ configurations in Appendix D.2 [4].

Performance analysis. According to Figure 9b⁶ with $\#F = 5$, omitting MonMLP (w/o MonMLP) causes the most significant decline, emphasizing the necessity of monotonicity constraints for effective personalized QoE modeling in the reward function. Among policy training ablations (w/o Prioritized Replay and w/o Old Replay Buffer), removing the old replay buffer (w/o Old Replay Buffer) has a more detrimental effect, highlighting the importance of retaining historical trajectories.

7 Related Work

Subjective QoE modeling. We review QoE modeling from the aspects of model training and model space. For training, methods

learning from numerical scores via regression [11, 18, 33] or adversarial training [46] fail in cross-session settings due to numerical distortion. Recent work [17] uses pairwise comparisons for ordinal relationships but overlooks cardinal relations, limiting feedback utilization and risking misaligned predictions. In contrast, we learn reliable QoE from both ordinal and cardinal preferences with cross-session datasets. For model space, formula-based models [30, 44] provide interpretability yet misalign with real QoE, non-parametric approaches [11] are limited to numerical regression, and DNNs [17, 46] enhance complexity but struggle to balance expressivity and generalization. Our QoE modeling ensures both expressivity and generalization through monotonic neural networks.

Personalized QoE-driven ABR methods. Quality of Service (QoS)-driven ABR methods are often rule-based [20, 22, 41], focusing on system targets but underperform. Conventional QoE-driven methods optimize subjective perception but require costly retraining for QoE changes [17, 30, 43], while our method avoids this through progressive policy fine-tuning. Unlike policy-gradient DRL methods [17, 30], we use a value-based approach for sample efficiency. While some methods adapt to various QoE models without retraining, they are limited to 2-3 QoE parameters [18, 19, 34, 48] or are noise-sensitive [44]. In contrast, our method supports expressive QoE models that output scalar rewards and does not rely on precise environmental modeling.

8 Conclusion

In this work, we propose Q+, a novel personalized QoE-driven ABR framework that bridges the gap between user preferences and adaptive streaming through iterative human feedback. Our framework introduces a scalable, user-friendly progressive feedback collection scheme, a robust personalized QoE model that unifies ordinal and cardinal preferences via monotonic neural networks, and a calibrated RL algorithm that dynamically aligns bitrate policies with personalized QoE updates. Extensive experiments demonstrate that our framework outperforms SOTA baselines in both QoE accuracy and policy effectiveness, achieving significant improvements in user experience across diverse network scenarios.

Acknowledgments

This work is supported by the National Key R&D Program of China (Grant No. 2024YFC3017100) and the National Natural Science Foundation of China (Grant No. 62302292).

⁶Additional ablation results for $\#F \sim \text{Normal}(6, 3)$ are provided in Appendix G.3 [4], where the findings for $\#F \sim \text{Normal}(6, 3)$ align with those for $\#F = 5$.

References

- [1] [n. d.]. Toward A Practical Perceptual Video Quality Metric | by Netflix Technology Blog | Netflix TechBlog. <https://netflixtechblog.com/toward-a-practical-perceptual-video-quality-metric-653f208b9652-VMAF-blog1>.
- [2] 2025. Amazon Prime Video. <https://www.primevideo.com>
- [3] 2025. Netflix. <https://www.netflix.com/>
- [4] 2025. Q+ Appendix. https://sites.sj.tju.edu.cn/yifei-zhu/wp-content/uploads/sites/19/2025/07/QPLUS_Appendix.pdf
- [5] 2025. YouTube. <https://www.youtube.com/>
- [6] Zahaib Akhtar, Yun Seong Nam, Ramesh Govindan, Sanjay Rao, Jessica Chen, Ethan Katz-Bassett, Bruno Ribeiro, Jibin Zhan, and Hui Zhang. 2018. Oboe: auto-tuning video ABR algorithms to network conditions. In *Proc. ACM SIGCOMM (Budapest, Hu.) (SIGCOMM '18)*. 44–58.
- [7] Avşar Asan, Werner Robitza, Is-haka Mkwawa, Lingfen Sun, Emmanuel Ifeachor, and Alexander Raake. 2017. Impact of video resolution changes on QoE for adaptive video streaming. In *ICME (Hong Kong, China)*. IEEE, 499–504.
- [8] Ralph Allan Bradley and Milton E. Terry. 1952. Rank Analysis of Incomplete Block Designs: I. The Method of Paired Comparisons. *Biometrika* 39, 3/4 (1952), 324–345.
- [9] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *NeurIPS (Long Beach, California, USA) (NIPS'17)*. Curran Associates Inc., Red Hook, NY, USA, 4302–4310.
- [10] 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. In *SIGCOMM (Toronto, Ontario, Canada) (SIGCOMM '11)*. 362–373.
- [11] Zhengfang Duanmu, Wentao Liu, Diqi Chen, Zhuoran Li, Zhou Wang, Yizhou Wang, and Wen Gao. 2023. A Bayesian Quality-of-Experience Model for Adaptive Streaming Videos. *ACM Trans. Multimedia Comput. Commun. Appl.* 18, 3s, Article 141 (feb 2023), 24 pages.
- [12] Zhengfang Duanmu, Wentao Liu, Zhuoran Li, Diqi Chen, Zhou Wang, Yizhou Wang, and Wen Gao. 2020. Assessing the Quality-of-Experience of Adaptive Bitrate Video Streaming. [arXiv:2008.08804 \[eess.IV\]](https://arxiv.org/abs/2008.08804)
- [13] Xuening Feng, Zhaohui JIANG, Timo Kaufmann, Eyke Hüllermeier, Paul Weng, and Yifei Zhu. 2025. Comparing Comparisons: Informative and Easy Human Feedback with Distinguishability Queries. In *ICML*.
- [14] Xuening Feng, Zhaohui Jiang, Timo Kaufmann, Puchen Xu, Eyke Hüllermeier, Paul Weng, and Yifei Zhu. 2025. DUO: Diverse, Uncertain, On-Policy Query Generation and Selection for Reinforcement Learning from Human Feedback. In *AAAI*.
- [15] Leo Gao, John Schulman, and Jacob Hilton. 2023. Scaling Laws for Reward Model Overoptimization. In *ICML (Honolulu, Hawaii, USA)*. 10835–10866.
- [16] Xiao Hu, Jianxiong Li, Xianyuan Zhan, Qing-Shan Jia, and Ya-Qin Zhang. 2023. Query-Policy Misalignment in Preference-Based Reinforcement Learning. [ArXiv abs/2305.17400 \(2023\)](https://arxiv.org/abs/2305.17400).
- [17] Tianchi Huang, Rui-Xiao Zhang, Chenglei Wu, and Lifeng Sun. 2023. Optimizing Adaptive Video Streaming with Human Feedback. In *MM (Ottawa ON, Canada) (MM '23)*. 1707–1718.
- [18] Tianchi Huang, Chao Zhou, Xin Yao, Rui-Xiao Zhang, Chenglei Wu, Bing Yu, and Lifeng Sun. 2020. Quality-Aware Neural Adaptive Video Streaming With Lifelong Imitation Learning. *IEEE Journal on Selected Areas in Communications* 38, 10 (2020), 2324–2342.
- [19] Tianchi Huang, Chao Zhou, Rui-Xiao Zhang, Chenglei Wu, and Lifeng Sun. 2022. Learning Tailored Adaptive Bitrate Algorithms to Heterogeneous Network Conditions: A Domain-Specific Priors and Meta-Reinforcement Learning Approach. *IEEE Journal on Selected Areas in Communications* 40, 8 (2022), 2485–2503.
- [20] Te-Yuan Huang, Ramesh Johari, Nick McKeown, Matthew Trunnell, and Mark Watson. 2014. A buffer-based approach to rate adaptation: evidence from a large video streaming service. In *SIGCOMM (Chicago, Illinois, USA) (SIGCOMM '14)*. 187–198.
- [21] Liangyu Huo, Zulin Wang, Mai Xu, Yong Li, Zhiguo Ding, and Hao Wang. 2020. A Meta-Learning Framework for Learning Multi-User Preferences in QoE Optimization of DASH. *IEEE Transactions on Circuits and Systems for Video Technology* 30, 9 (2020), 3210–3225.
- [22] Junchen Jiang, Vyas Sekar, and Hui Zhang. 2012. Improving fairness, efficiency, and stability in HTTP-based adaptive video streaming with FESTIVE. In *CoNEXT (Nice, France) (CoNEXT '12)*. 97–108.
- [23] Zhaohui JIANG, Xuening Feng, Paul Weng, Yifei Zhu, Yan Song, Tianze Zhou, Yujing Hu, Tangjie Lv, and Changjie Fan. 2025. Reinforcement Learning from Imperfect Corrective Actions and Proxy Rewards. In *ICLR (Singapore)*.
- [24] S. Shunmuga Krishnan and Ramesh K. Sitaraman. 2013. Video Stream Quality Impacts Viewer Behavior: Inferring Causality Using Quasi-Experimental Designs. *IEEE/ACM Transactions on Networking* 21, 6 (2013), 2001–2014.
- [25] Kimin Lee, Laura Smith, Anca Dragan, and Pieter Abbeel. 2021. B-Pref: Benchmarking Preference-Based Reinforcement Learning. In *NeurIPS*. Neural Information Processing Systems Foundation, Curran Associates Inc.
- [26] Kimin Lee, Laura M Smith, and Pieter Abbeel. 2021. PEBBLE: Feedback-Efficient Interactive Reinforcement Learning via Relabeling Experience and Unsupervised Pre-training. In *ICML (Proceedings of Machine Learning Research, Vol. 139)*. PMLR, 6152–6163.
- [27] Yiyun Lu, Yifei Zhu, and Zhi Wang. 2022. Personalized 360-degree video streaming: A meta-learning approach. In *ACM Multimedia*. 3143–3151.
- [28] Jianlan Luo, Perry Dong, Yuexiang Zhai, Yi Ma, and Sergey Levine. 2024. RLIF: Interactive Imitation Learning as Reinforcement Learning. In *ICLR (Vienna, Austria)*.
- [29] Hongzi Mao, Parimarjan Negi, Akshay Narayan, Hanrui Wang, Jiacheng Yang, Haonan Wang, Ryan Marcus, ravichandra addanki, Mehrdad Khani Shirkoobi, Songtao He, Vikram Nathan, Frank Cangialosi, Shailesh Venkatakrishnan, Weihung Weng, Song Han, Tim Kraska, and Dr.Mohammad Alizadeh. 2019. Park: An Open Platform for Learning-Augmented Computer Systems. In *NeurIPS*, Vol. 32. Curran Associates, Inc.
- [30] Hongzi Mao, Ravi Netravali, and Mohammad Alizadeh. 2017. Neural Adaptive Video Streaming with Pensieve. In *SIGCOMM (Los Angeles, CA, USA) (SIGCOMM '17)*. 197–210. doi:10.1145/3098822.3098843
- [31] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *nature* 518, 7540 (2015), 529–533.
- [32] Jane Pan, He He, Samuel R. Bowman, and Shi Feng. 2024. Spontaneous Reward Hacking in Iterative Self-Refinement. [CoRR abs/2407.04549 \(2024\)](https://arxiv.org/abs/2407.04549).
- [33] Leonardo Peroni, Sergey Gorinsky, Farzad Tashtarian, and Christian Timmerer. 2023. Empowerment of Atypical Viewers via Low-Effort Personalized Modeling of Video Streaming Quality. *CoNEXT 1*, CoNEXT3, Article 17 (nov 2023), 27 pages.
- [34] Chunyu Qiao, Jiliang Wang, and Yunhao Liu. 2021. Beyond QoE: Diversity Adaptation in Video Streaming at the Edge. *IEEE/ACM Transactions on Networking* 29, 1 (2021), 289–302.
- [35] Werner Robitza, Marie-Neige Garcia, and Alexander Raake. 2017. A modular HTTP adaptive streaming QoE model – Candidate for ITU-T P.1203 (“P.NATS”). In *2017 Ninth International Conference on Quality of Multimedia Experience (QoMEX)*. Erfurt, Germany, 1–6.
- [36] Davor Runje and Sharath M Shankaranarayanan. 2023. Constrained Monotonic Neural Networks. In *ICML (Honolulu, Hawaii, USA) (Proceedings of Machine Learning Research, Vol. 202)*. PMLR, 29338–29353.
- [37] Sandvine. 2024. The Global Internet Phenomena Report January 2024. <https://www.sandvine.com/blog/sandvines-2024-global-internet-phenomena-report-global-internet-usage-continues-to-grow>
- [38] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. 2016. Prioritized Experience Replay. In *ICLR (San Juan, Puerto Rico)*.
- [39] Max Schwarzer, Johan Samir Obando Ceron, Aaron Courville, Marc G Bellemare, Rishabh Agarwal, and Pablo Samuel Castro. 2023. Bigger, better, faster: Human-level atari with human-level efficiency. In *ICML*. PMLR, 30365–30380.
- [40] Joar Skalse, Nikolaus Howe, Dmitrii Krashennikov, and David Krueger. 2022. Defining and Characterizing Reward Gaming. In *NeurIPS*, Vol. 35. Curran Associates, Inc., 9460–9471.
- [41] Kevin Spiteri, Rahul Uргаonkar, and Ramesh K. Sitaraman. 2020. BOLA: Near-Optimal Bitrate Adaptation for Online Videos. *IEEE/ACM Transactions on Networking* 28, 4 (2020), 1698–1711.
- [42] Thomas Stockhammer. 2011. Dynamic adaptive streaming over HTTP –: standards and design principles. In *Proceedings of the Second Annual ACM Conference on Multimedia Systems (San Jose, CA, USA) (MMSys '11)*. 133–144.
- [43] Francis Y. Yan, Hudson Ayers, Chenzhi Zhu, Sadjad Fouladi, James Hong, Keyi Zhang, Philip Levis, and Keith Winstein. 2020. Learning in situ: a randomized experiment in video streaming. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*. USENIX Association, Santa Clara, CA, 495–511.
- [44] Xiaoqi Yin, Abhishek Jindal, Vyas Sekar, and Bruno Sinopoli. 2015. A Control-Theoretic Approach for Dynamic Adaptive Video Streaming over HTTP. (2015), 325–338. doi:10.1145/2785956.2787486
- [45] Huaizheng Zhang, Linsen Dong, Guanyu Gao, Han Hu, Yonggang Wen, and Kyle Guan. 2020. DeepQoE: A Multimodal Learning Framework for Video Quality of Experience (QoE) Prediction. *IEEE Transactions on Multimedia* 22, 12 (2020), 3210–3223.
- [46] Huanhuan Zhang, Liu zhuo, Haotian Li, Anfu Zhou, Chuanming Wang, and Huadong Ma. 2024. AraLive: Automatic Reward Adaption for Learning-based Live Video Streaming. In *MM*. 11099–11108.
- [47] Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Yuhao Zhou, Limao Xiong, Lu Chen, Zhiheng Xi, Nuo Xu, Wenbin Lai, Minghao Zhu, Haoran Huang, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. Delve into PPO: Implementation Matters for Stable RLHF. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*. Curran Associates Inc.
- [48] Xutong Zuo, Jiayu Yang, Mowei Wang, and Yong Cui. 2022. Adaptive Bitrate with User-level QoE Preference for Video Streaming. In *INFOCOM*. IEEE, London, United Kingdom, 1279–1288.

MM '25, October 27–31, 2025, Dublin, Ireland.

Zhaohui Jiang et al.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

A Mathematical Descriptions for Figure 1

To deliver a clear analysis, we use the following notations to mathematically describe SQoE-IV: τ to denote a video experience, Γ to denote the set of all video experiences, $c_\tau \in \mathbb{R}$ (resp. $C_\Gamma \in \mathbb{R}^{|\Gamma|}$) to denote the score(s) for τ (resp. Γ), $\bar{\tau}$ to specifically denote an overlapping video experience, $\bar{\tau}^{(2)}$ to denote a pair of overlapping video experiences $(\bar{\tau}_j, \bar{\tau}_k)$, $\bar{\Gamma}$ to denote a set of overlapping video experiences from the same session, and the pair(s) of scores provided by an evaluator on $\bar{\tau}$ (resp. $\bar{\Gamma}$) as $c_{\bar{\tau}} \in \mathbb{R}^2$ (resp. $C_{\bar{\Gamma}} \in \mathbb{R}^{2 \times |\bar{\Gamma}|}$). In addition, we use e as the upperscript to c or C to denote scores that come from the evaluator e . In the following, we explain how the scalar results are generated to construct the distributions in subfigures of Figure 1:

Figure 1a: $\forall e_1, e_2 \in \text{SQoE-IV}$

$$(\text{SRCC}(C_{\Gamma}^{e_1}, C_{\Gamma}^{e_2}), \text{PLCC}(C_{\Gamma}^{e_1}, C_{\Gamma}^{e_2})),$$

Figure 1b: $\forall e \in \text{SQoE-IV}$:

$$\text{Euclidean}(C_{\bar{\Gamma}}^e) = \sqrt{\sum_{\bar{\tau} \in \bar{\Gamma}} (c_{\bar{\tau}}^e[0] - c_{\bar{\tau}}^e[1])^2},$$

$$\text{Manhattan}(C_{\bar{\Gamma}}^e) = \sum_{\bar{\tau} \in \bar{\Gamma}} |c_{\bar{\tau}}^e[0] - c_{\bar{\tau}}^e[1]|,$$

$$\text{Chebyshev}(C_{\bar{\Gamma}}^e) = \max_{\bar{\tau} \in \bar{\Gamma}} |c_{\bar{\tau}}^e[0] - c_{\bar{\tau}}^e[1]|.$$

Figure 1c: $\forall e \in \text{SQoE-IV}$:

$$(\text{SRCC}(C_{\bar{\Gamma}}^e[0], C_{\bar{\Gamma}}^e[1]), \text{PLCC}(C_{\bar{\Gamma}}^e[0], C_{\bar{\Gamma}}^e[1])).$$

Figure 1d: $\forall e \in \text{SQoE-IV}$: For a pair of overlapping video experiences $(\bar{\tau}_j, \bar{\tau}_k)$, the evaluator's scores in session i define the *ordinal preference* y_o labels as:

$$y_o(\bar{\tau}_j, \bar{\tau}_k)[i] = \begin{cases} =, & \text{if } |c_{\bar{\tau}_j}[i] - c_{\bar{\tau}_k}[i]| < \delta_o, \\ >, & \text{if } c_{\bar{\tau}_j}[i] > c_{\bar{\tau}_k}[i] + \delta_o, \\ <, & \text{if } c_{\bar{\tau}_j}[i] < c_{\bar{\tau}_k}[i] - \delta_o, \end{cases} \quad (12)$$

where $>$ (resp. $<$) denotes preference for the left (resp. right) video experience, $=$ denotes equal preference, and δ_o is a threshold for equal preference. Additionally, denoting the absolute score difference for two overlapping video experiences $|c_{\bar{\tau}_j}[i] - c_{\bar{\tau}_k}[i]|$ as $\Delta_c(\bar{\tau}_j, \bar{\tau}_k)[i]$, we consider the *cardinal preference* y_c which compares the strength of ordinal preferences (i.e., which one of the two pairs of video experiences is easier to be judged) as:

$$y_c(\bar{\tau}_j^{(2)}, \bar{\tau}_k^{(2)})[i] = \begin{cases} =, & \text{if } |\Delta_c(\bar{\tau}_j^{(2)})[i] - \Delta_c(\bar{\tau}_k^{(2)})[i]| < \delta_c, \\ >, & \text{if } \Delta_c(\bar{\tau}_j^{(2)})[i] > \Delta_c(\bar{\tau}_k^{(2)})[i] + \delta_c, \\ <, & \text{if } \Delta_c(\bar{\tau}_j^{(2)})[i] < \Delta_c(\bar{\tau}_k^{(2)})[i] - \delta_c, \end{cases} \quad (13)$$

where δ_c is the threshold for equal cardinal preference. To measure cross-session consistency of y_o and y_c , we calculate Identity Rate (IR) [11, 15, 20, 31]:

$$\text{IR}_o(\bar{\Gamma}) = \mathbb{E}_{\bar{\tau}_j, \bar{\tau}_k \in \bar{\Gamma}} \left[\mathbb{I}[y_o(\bar{\tau}_j, \bar{\tau}_k)[0] = y_o(\bar{\tau}_j, \bar{\tau}_k)[1]] \right], \quad (14)$$

$$\text{IR}_c(\bar{\Gamma}) = \mathbb{E}_{\bar{\tau}_j^{(2)}, \bar{\tau}_k^{(2)} \in \bar{\Gamma}} \left[\mathbb{I}[y_c(\bar{\tau}_j^{(2)}, \bar{\tau}_k^{(2)})[0] = y_c(\bar{\tau}_j^{(2)}, \bar{\tau}_k^{(2)})[1]] \right], \quad (15)$$

with \mathbb{I} as the indicator function, describing the percentage of identical preferences between two sessions. Figure 1d illustrates the

distribution of ordinal (IR_o) and cardinal (IR_c) identity rates of SQoE-IV, with $\delta_o = \delta_c = 10$ accounting for potential human feedback noise. Additional statistics with different δ_o and δ_c are available in Figure 10, which provides extended statistics compared to Figure 1d on IR_o and IR_c values computed from score vectors on stimuli sets over all evaluators in SQoE-IV. The results indicate that larger thresholds tend to yield higher identity rates as they mitigate noise in human feedback. However, excessively large thresholds for equal-preference checking can reduce the utility of the data for QoE model training. Therefore, in our experiments, we set the thresholds δ_o (resp. δ_c) to 10 (resp. 20), corresponding to 1/10 (resp. 1/5) of the total score range of SQoE-IV (i.e., 100).

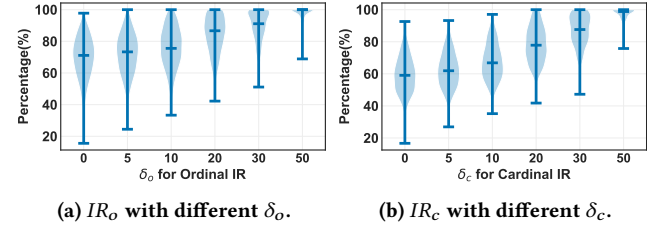


Figure 10: Violin plots show the distribution of IR, using the same configuration as Figure 1d but with adjusted thresholds (δ_o and δ_c) for equal preference in Equations (12) and (13).

B Hyper-Parameters

B.1 QoE Modeling

Table 2 lists the configurations for QoE modeling.

Table 2: Hyperparameters for QoE modeling.

Hyperparameter	Value
Model architecture	
K (number of input chunks)	7
Inner activation	Augmented ReLU, ratio (1:1:1)[41]
Final activation	Tanh
Hidden dimension	256
Hidden layer	2
Model training	
Training epoch	500
Learning rate	1e-3
Batch size	512
Optimizer	Adam [28]
δ_o	10
δ_c	20

B.2 Policy Fine-Tuning

Table 3 lists the configurations for Q+'s policy fine-tuning.

Table 3: Hyperparameters for Q+’s cross-session policy fine-tuning.

Type	Hyperparameter	Value
General	timesteps per session T	1M
	replay size	1M
	replay period timesteps T'	4
	policy updates per replay	4
	training batch size B	32
	target network update interval	2000
	number of atoms for categorical-Q	51
	upper bound for categorical Q-values	10
	lower bound for categorical Q-values	-10
	n-step	1
	α for priority replay	0.5
	β for priority replay	$0.4 \rightarrow 1$
Model Architecture	hidden layers	3
	hidden size	128
Reward	r_{min}	0
	r_{max}	0.05
	discount factor	0.99
Optimizer	type	Adam [28]
	learning rate η	0.0001
	eps	$0.01/B$
	betas	(0.9, 0.999)

C More Details about QoE Baselines and Ablations

C.1 Baseline QoE Models

Previous QoE models were designed for a general user base without incorporating personalization. For the QoE models from BSQI, we directly use their open-sourced model. For other QoE models, the fundamental structure is a linear model, where the QoE state s_t at timestep t is constructed using the video quality measurements $quality_t$ and $quality_{t-1}$ for the t -th and $(t-1)$ -th video chunks, respectively, and the rebuffering time u_t caused by selecting the t -th video chunk. The QoE value for selecting the t -th video chunk is then determined by Equation 16:

$$QoE_{Lin}(s_t) = w_1 \cdot \text{quality}(s_t) + w_2 \cdot u_t \quad (16)$$

$$+ w_3 \cdot |\text{quality}(s_t) - \text{quality}(s_{t-1})| + \quad (17)$$

$$+ w_4 \cdot |\text{quality}(s_t) - \text{quality}(s_{t-1})|_-, \quad (18)$$

where $|\cdot|_+$ (resp. $|\cdot|_-$) measures the positive (resp. negative) smoothness of quality transitions.

As illustrated in Table 4, these baseline QoE models adopt different measurements as the quality function and assign different values to $\mathbf{w} = [w_1, w_2, w_3, w_4]$. The implementation for QoE models of Pensieve, BOLA, and MPC follows Mao et al. [35]’s. Concretely, Pensieve’s pre-defined mapping for q maps target bitrates of [300, 750, 1200, 1850, 2850, 4300] to QoE values of [1, 2, 3, 12, 15, 20], but considering that our bitrate ladder [235, 375, 560, 750, 1050, 1750, 2350, 3000, 4300] differs from Pensieve’s, we assign our minimal and maximal target bitrate to 1 and 20, then perform linear interpolation on the original value map for our bitrate ladder. For BOLA’s, q_{min} represents the minimum value in the bitrate ladder.

For Puffer, Comyco and Jade’s implementation, we follow their original paper.

Table 4: Configurations of baseline QoE models.

QoE Model	Quality Measurement	Weights \mathbf{w}
Pensieve’s [35]	Pre-defined mapping for q	1,-8,-1,-1
BOLA’s [46]	$\log(\frac{q_t}{q_{min}})$	1,-2.66,-1,-1
MPC’s [51]	q_t	1,-4.3,-1,-1
Comyco-Lin [21]	VMAF v_t	0.8469, -28.7959, 0.2979, -1.0610
Jade-Lin [20]	VMAF v_t	0.535, -0.215, -0.13, -1.37
Puffer [50]	SSIM	1,-100,-1,-1

C.2 QoE Ablation Variants

The ablations for our QoE modeling method involve ablating our model architecture MonMLP or training strategy that learns from ordinal and cardinal pairwise comparison with \mathcal{L}_ϕ (Equation (9)). Specifically:

For training strategy. (i) For the regression method, we minimize the Mean Squared Error (MSE) between subjective scores and predicted QoE values. (ii) For the ordinal training method [20], the training procedure is the same as our personalized QoE modeling scheme, except replacing our Equation (9) with Equation (6).

For model architecture. (i) For the linear format, we implement it following previous works as Equation (16), by taking the VMAF value as the quality function in it. (ii) For the MLP format, we use the same number of model layers and hidden sizes, with $\{u_t, q_t, v_t\}$ as the input features and ReLU as the activation.

D More Details about Policy Baselines and Ablations

D.1 Policy Baselines

Pensieve-F. Except for the total number of environmental timesteps, which we need to adjust to ensure a fair comparison with Q+, the hyperparameters for this baseline follow previous methods [20, 21, 35].

Initial policy for Pensieve-F and Q+. For Q+, the initial policy is trained using Rainbow. For Pensieve-F, the initial policy is trained using PPO. Both policies achieve converged performance with Jade’s linear QoE, which is the best-performing general QoE as shown in Table 1. Both policies use the same hyperparameters as their fine-tuned versions, except for the training steps.

MPC and Robust-MPC [51]. Following previous methods [20, 21, 27, 35], the two baselines estimate throughput by calculating the harmonic mean of observed throughput values from the previous 5 chunk downloads. While MPC directly uses the harmonic mean HM_t at timestep t as the predicted future bandwidth PB_t , Robust-MPC is more conservative in the sense that it scales the harmonic mean based on the prediction error PE of the last 5 chunks as $PB_{t+1} = HM_t / (1 + \max(PE_{t-5}, \dots, PE_{t-1}))$, so a larger prediction error leads to a lower PB_t . With PB_t , MPC and Robust-MPC optimize video streaming by predicting future network conditions over a short horizon, for which we keep the same as previous works and

set it to 3. Then it evaluates all possible bitrate combinations within this window, selects the sequence maximizing a predefined QoE metric, and implements only the first step, repeating this process iteratively to balance immediate performance and future outcomes. To ensure comparable evaluation (i.e., with a QoE trained with the same amount of personal feedback), the QoE metric for the two methods to optimize is the personalized QoE model learned from the whole set of feedback collected at the end as Q+.

Rate-Based[25]. Rate-Based select the next chunk with the maximal target bitrate that does not exceed the estimated throughput, with the same throughput estimations as MPC.

BBA [23]. BBA is a buffer-based bitrate adaptation strategy that selects the minimum target bitrate when the buffer is below the RESERVOIR threshold, the maximum target bitrate when the buffer exceeds RESERVOIR + CUSHION, and linearly interpolates the bitrate between the two extremes based on the current buffer size BS relative to the CUSHION range as

$$\lfloor \lceil \mathcal{A} \rceil \cdot (BS - RESERVOIR) / CUSHION \rfloor,$$

casting the result to an integer. Following Kan et al. [27], we set RESERVOIR = 20s and CUSHION = 8s in our experiments for this baseline.

HYB. HYB considers both buffer occupancy B and throughput prediction L, estimating a target file size based on $0.25 \times L \times B$, then selects the largest chunk with file size below the estimation.

BOLA [46]. We utilize the implementation from Kan et al. [27] for this baseline.

D.2 Policy Ablations

In Table 5 we compare the configurations for the two ablations related to testing the effects of personalized QoE modeling designs on Q+, with other configurations related to policy fine-tuning part the same as Q+.

Table 5: Compare the configurations of our ablations for Q+'s online personalized QoE modeling.

Method	Configurations	
	Model Architecture	Training Method
Q+	MonMLP	\mathcal{L}_ϕ (Equation (9))
w/o MonMLP	MLP	\mathcal{L}_ϕ (Equation (9))
w/o Cardinal Loss	MonMLP	\mathcal{L}_ϕ^o (Equation (6))

In Table 6 we compare the configurations for the two ablations related to the policy fine-tuning of Q+, with other configurations related to online personalized QoE modeling part the same as Q+. A more detailed configurations corresponding to Table 6 can be checked in Table 6.

E Implementation Details for our QoE Modeling

E.1 Implementation Details for MonMLP

Letting W , b , and y represent the weights, bias, and the outputs for a fully connected linear layer, $\hat{s}_t \in \mathbb{R}^{K+2}$ and $\tilde{s}_t \in \mathbb{R}^{2 \times K}$ represents the monotonic and unconstrained part of the QoE inputs,

Table 6: Compare the configurations of our ablations for Q+'s policy fine-tuning.

Method	Configurations	
	Old Trajectories	Sample Replay
Q+	Keep & relabel reward	Prioritized (relabel priority)
w/o Old Replay Buffer	Remove	Prioritized (relabel priority)
w/o Prioritized Replay	Keep & relabel reward	Uniform

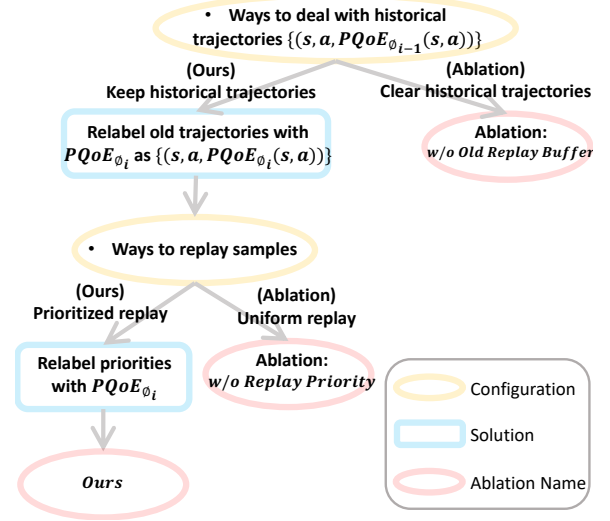


Figure 11: An illustration for the configurations of our ablations for Q+'s policy fine-tuning.

respectively. To guarantee the monotonicity constraints, we set $\hat{s}_t = \{-u_t, \min(q_t), \min(v_t)\}$ and $\tilde{s}_t = \{q_t - \min(q_t), v_t - \min(v_t)\}$. MonMLP computes the first layer's outputs as $y_1 = |\tilde{W}_1^T| \cdot \hat{s}_t + \tilde{W}_1^T \cdot \tilde{s}_t + b_1$. For subsequent layers, it uses $y_j = |\tilde{W}_j^T| \cdot \rho^M(y_{j-1}) + b_j$, where ρ^M denotes the activation layer, \tilde{W}_j are the weights assigned to values requiring monotonicity from the $(j - 1)$ -th layer, and \tilde{W}_j are assigned to other unconstrained values.

The activation layer design follows Runje and Shankaranarayana [41] to preserve model expressivity.

E.2 Pseudo-code for QoE Modeling

See Algorithm 2 for the pseudo-code.

F Implementation Details for Cross-session Policy Fine-Tuning

Suitable resets are required to avoid wrong information stemming from fine-tuning the policy π_{θ_t} from $\pi_{\theta_{t-1}}$ with an updated reward, including: (i) reset the target network by re-initialize $Q_{\bar{\theta}}$ as $Q_{\theta_{t-1}}$, which is the same as the initial policy to fine-tune Q_{θ_t} , (ii) reset the inner parameters of Adam, noisy network, and prioritized replay, with details elaborated in Appendix F.1, Appendix F.2, and Appendix F.3, respectively.

We outline the whole procedure of policy training in Algorithm 3, where our modifications to the original Rainbow are highlighted

Algorithm 2 QoE Modeling (The i -th Session)

```

1: Personalized QoE model with re-initialized  $\phi$ :  $R_\phi$ 
2: User's feedback collected in the  $i$ -th interaction session:  $\mathcal{D}_i^c$ 
3: User's accumulating feedback:  $\mathcal{D}^c (\cup_{j \leq i} \mathcal{D}_j^c)$ 
4: ▶ Construct ordinal pairs
5:  $\mathcal{D}_i^o \leftarrow \{(\tau_1, \tau_2) | \forall \tau_1, \tau_2 \in \mathcal{D}_i^c\}$ 
6: Label  $\mathcal{D}_i^o$  with Equation (1)
7: ▶ Construct cardinal pairs
8:  $\mathcal{D}_i^c \leftarrow \{(\tau_1^{(2)}, \tau_2^{(2)}) | \forall \tau_1^{(2)}, \tau_2^{(2)} \in \mathcal{D}_i^o\}$ 
9: Label  $\mathcal{D}_i^c$  with Equation (2)
10:  $\mathcal{D}^c \leftarrow \cup_{j \leq i} \mathcal{D}_j^c$  ▶ Train with all collected feedback
11: ▶ Training epochs
12: The total number of updating personalized QoE:  $N_R$ 
13: for  $j \leftarrow 1$  to  $N_R$  do
14:   Sample a batch from  $\mathcal{D}^c$ , update  $\phi$  with Equation (9).
15: end for
16: return  $R_{\phi_i}$ 

```

with orange color. Concrete hyperparameters can be checked in Table 3.

Algorithm 3 QoE-Driven Policy Fine-Tuning (The i -th Session)

```

1: Initial policy:  $Q_\theta \leftarrow Q_{\theta_{i-1}}$  ▶ Fine-tune the previous one
2: Personalized QoE updated after the  $i$ -th interaction session:  $R_{\phi_i}$ 
3: Replay buffer labeled with old personalized QoE:  $\mathcal{D}^\tau = \{(s, a, R_{\phi_{i-1}}(s, a))\}$ 
4: ▶ Reset parameters
5: Reset target Q network:  $Q_\theta \leftarrow Q_{\theta_{i-1}}$ 
6: Reset inner parameters for Adam, noisy net, and prioritized replay.
7: ▶ Relabeling with new reward
8: Relabel replay buffer with new personalized QoE:  $\mathcal{D}^\tau \leftarrow \{(s, a, R_{\phi_i}(s, a))\}$ 
9: Relabel prioritise  $P$  (Equation (19)) for samples in replay buffer
10: ▶ Policy training
11: The total number of timesteps for policy training:  $T$ 
12: for  $t \leftarrow 1$  to  $T$  do
13:    $a_t \sim \text{NoisyNet}_{Q_\theta}(s_t)$ 
14:    $\mathcal{D}^\tau \leftarrow \mathcal{D}^\tau \cup (s_t, a_t, R_{\phi_i}(s_t, a_t))$ 
15:   if  $t \% T'$  then ▶ Replay period every  $T'$  timesteps
16:     Sample  $B$  samples from replay buffer with probabilities  $P$  (Equation (19))
17:     Optimize  $\theta$  on samples with  $\mathcal{L}_\theta$  (Equation (11))
18:   end if
19: end for
20: return  $\pi_i$ , derived by  $\arg \max Q_{\theta_i}$ 

```

F.1 Reset Parameters for Optimizer

Following Asadi et al. [6], we reset the internal parameters of Adam, including the update step index and the first- and second-order moments, to zero.

F.2 Reset Parameters for Prioritized Replay

Prioritized replay [43] assigns larger probability P_t to sample transitions with larger TD-errors $\delta_t^{TD} = \delta^{TD}(s_t, a_t, r_t, s_{t+1})$ from \mathcal{D}^τ when optimize Q_θ with \mathcal{L}_θ (Equation (11)):

$$P_t = \frac{|\delta_t^{TD}|^\alpha}{\sum_t |\delta_t^{TD}|^\alpha}, \quad (19)$$

where α is a hyperparameter controls the extent to which $|\delta_t^{TD}|$ are converted into sample priorities. Since the non-uniform sampling used in prioritized replay [43] does not satisfy the requirement of assigning equal weights to experiences in \mathcal{D}^τ as stated in Equation (11), there exists biases introduced by non-uniform sampling that may lead to non-optimal policy convergence. To correct such biases, importance sampling weights w_t are applied to scale the gradient calculated on experience (s_t, a_t, r_t, s_{t+1}) during training:

$$w_t = \left(\frac{1}{|\mathcal{D}^\tau|} \cdot \frac{1}{P_t} \right)^\beta, \quad (20)$$

where $|\mathcal{D}^\tau|$ is the number of experiences in the replay buffer, and β is an annealing parameter that controls the extent of bias corrections. Despite the priority relabeling explained in the main text Section 5.2, we reset the parameter for the importance sampling correction β used in prioritized replay [43]. Concretely, we reset it to the initial small value 0.4 and then increase it to 1 alongside the timesteps of one fine-tuning session.

F.3 Reset Parameters for NoisyNet

NoisyNet [17] injects parametric noise into the weights of neural networks by replacing standard linear layers (i.e., $y = w \cdot x + b$) in Q-Networks with noisy (linear) layers. In these noisy layers, the standard weights w are replaced by noisy weights $w_N \in \mathbb{R}$, parameterized as follows:

$$w_N = \mu_N + \sigma_N \cdot \epsilon_N, \quad (21)$$

where μ_N and σ_N are trainable parameters, and ϵ_N is sampled noise with zero means. By incorporating noisy layers into both the online and target Q-networks with resampled ϵ_N for each forward pass, the parameters of the noisy layers are trained alongside the model parameters. This process updates both the noise level σ_N and the expected value μ_N using gradient descent. Consequently, the system facilitates automatic state-dependent exploration, with less noise for familiar states (where exploration may degrade performance) and more noise for unfamiliar states (where exploration may yield better outcomes).

Therefore, we only reset σ_N , which controls the noise level, at the beginning of each fine-tuning session to ensure appropriate exploration in response to new rewards. The parameters μ_N , which are involved in calculating the expected Q-values, are not reset and remain unchanged.

F.4 Min-Max Scaling for Online QoE

Considering that categorical Q-learning [9], which is one of the design choices in Rainbow, requires pre-determining the range of Q-values to generate the support of categorical distributions, we linearly transform the raw QoE outputs from $[-1, 1]$ to $[r_{min}, r_{max}]$ with $r_{min} = 0$ and $r_{max} = 0.05$ during policy training. This

choice cooperates with the hyperparameters related to categorical Q-learning as shown in Table 3.

G Extra Experimental Results about Policy Learning

In this section, we display more experimental results with different numbers of feedback per session ($\#F$) for different sets of traces.

G.1 Compare with Baselines

G.1.1 Training Polts. In Figure 12, we present the experimental results comparing Q+ with baseline methods under conditions where $\#F = 3$, $\#F = 10$, or $\#F \sim \text{Normal}(\mu = 6, \sigma = 3)$, respectively. The format of these figures mirrors that of Figure 7 presented in the main text. Although Q+ performs similarly with $\#F = 3$ to the best-performing baseline after 8 interaction sessions, in all other cases, Q+ consistently outperforms the baselines within 3 sessions.

G.1.2 Breakdown Performances. Moreover, we show more breakdown performance with different $\#F$ and different traces in Figures 17 to 20, sharing a similar format as Figure 8 in the main text.

G.2 Sensivity Analysis for Q+

Figure 13 illustrates the complete learning curves of our method facing diverse settings for the number of received feedback per session. Figure 13 corresponds to the summarized converged performance shown in Figure 9a. From the learning curves depicted in Figure 13, we observe that a larger $\#F$ enhances both the learning efficiency and the final performance.

G.3 Ablation Studies for Q+

The subfigures in the second row of Figure 14 display the ablation results of Q+ with $\#F \sim \text{Normal}(\mu = 6, \sigma = 3)$. The complete learning curves corresponding to Figure 7 are shown in the first row of Figure 14. The performance analysis for ablations is similar whether $\#F = 5$ or $\#F \sim \text{Normal}(\mu = 6, \sigma = 3)$.

H Extra Experimental Results about QoE Modeling

H.1 Visualization for personalized QoE Outputs

In Figure 6 in the main text we only show two kinds of 2D figures with $\min(\mathbf{q})$ & $\text{mean}(\mathbf{u})$ and $\min(\mathbf{v})$ & $\text{mean}(\mathbf{u})$ as the x & y axes for one of the evaluator (denoted as the 1st evaluator in this section). Here, we show more visualizations with an extra kind of 2D figure with $\min(\mathbf{v})$ & $\min(\mathbf{q})$ as the x & y axes, for 4 evaluators, as shown in Figures 22 to 24 to further support our argument that the MonMLP achieves better generalization ability than MLP.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

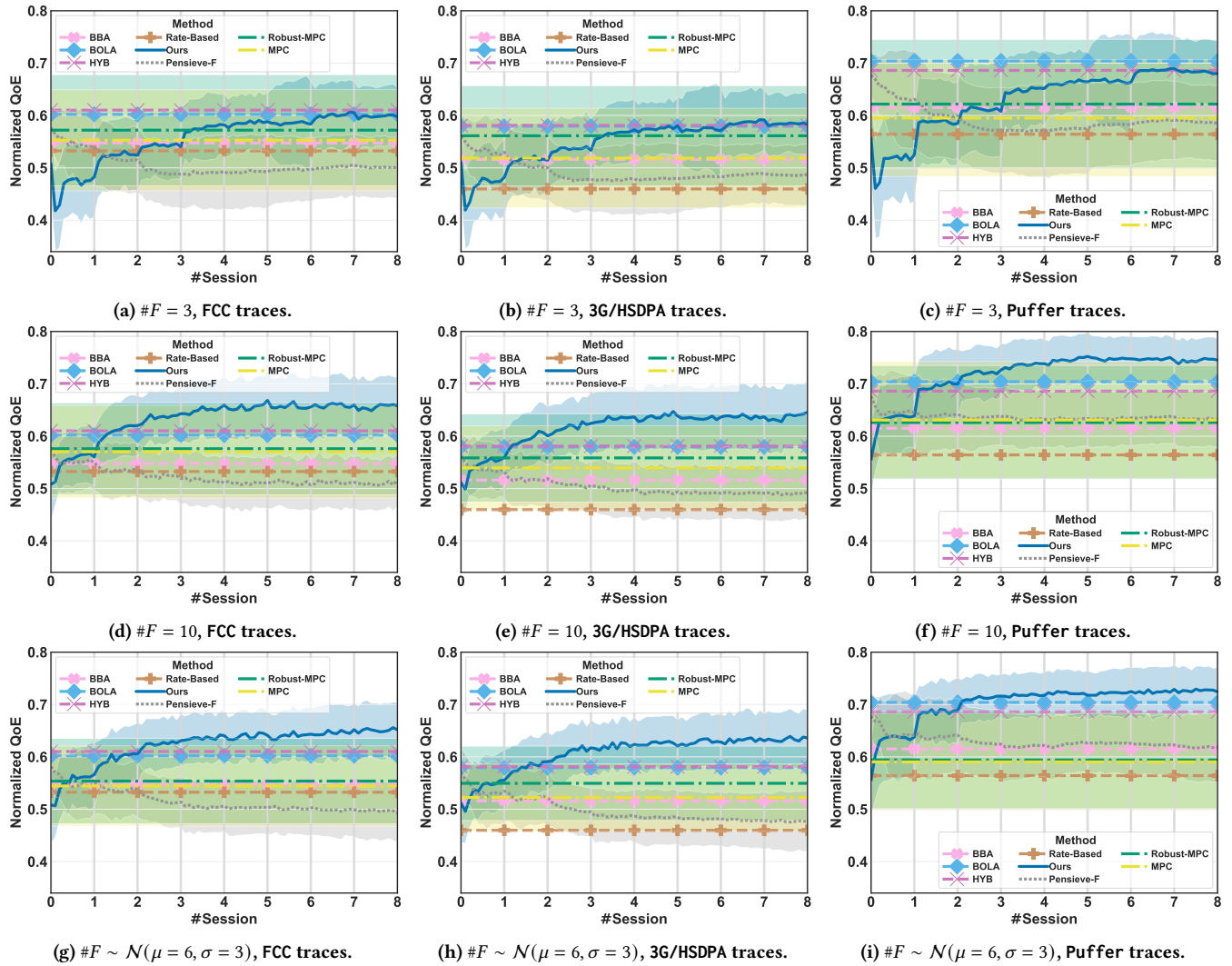


Figure 12: Compare Q+ with other baselines for each set of network traces, where the plots with the $\#F = 5$ setting have been shown in Figure 7. Subfigures in one row correspond to one $\#F$ setting. Vertical grey lines indicate the interaction sessions that R_ϕ will be updated upon receiving a set of $\#F$ new feedback, after which the policy is fine-tuned.

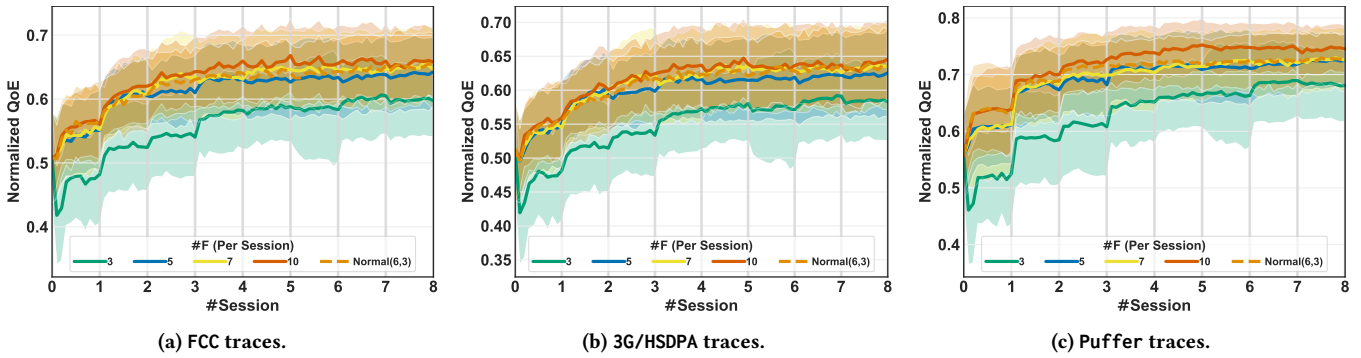


Figure 13: Sensitivity analysis for Q+ with different quantities of feedback per session (#F), with the same annotation meanings as explained in Figure 7. Vertical grey lines indicate the interaction sessions that R_ϕ will be updated upon receiving a set of #F new feedback, after which the policy is fine-tuned. In the legend, $Normal(\mu = 6, \sigma = 3)$ refers to sampling the number of feedbacks from a normal distribution per session, and others refer to constant numbers.

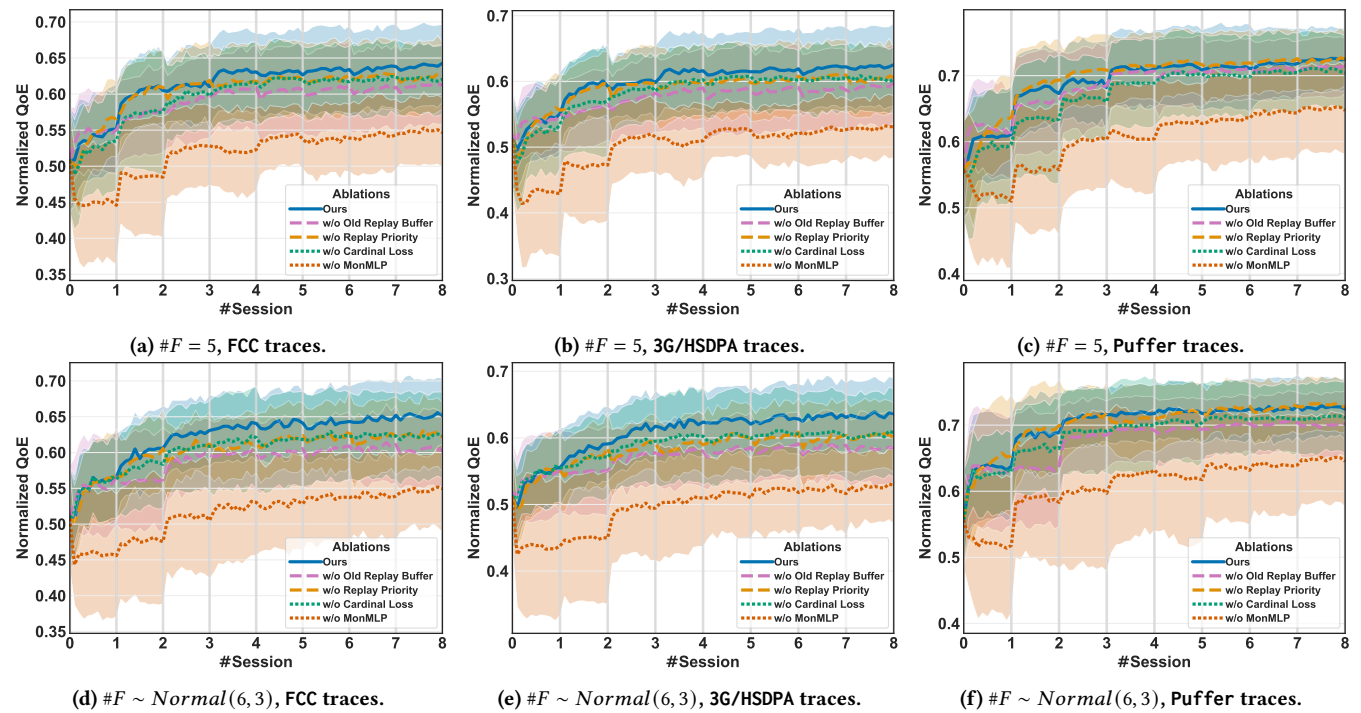


Figure 14: Ablations for Q+ with two representative #F settings #F = 5 (the first row) or #F ~ Normal($\mu = 6, \sigma = 3$) (the second row), with dashed and dotted lines indicating ablations related to policy fine-tuning and personalized QoE modeling, respectively. Vertical grey lines indicate the interaction sessions that R_ϕ will be updated upon receiving a set of #F new feedback, after which the policy is fine-tuned.

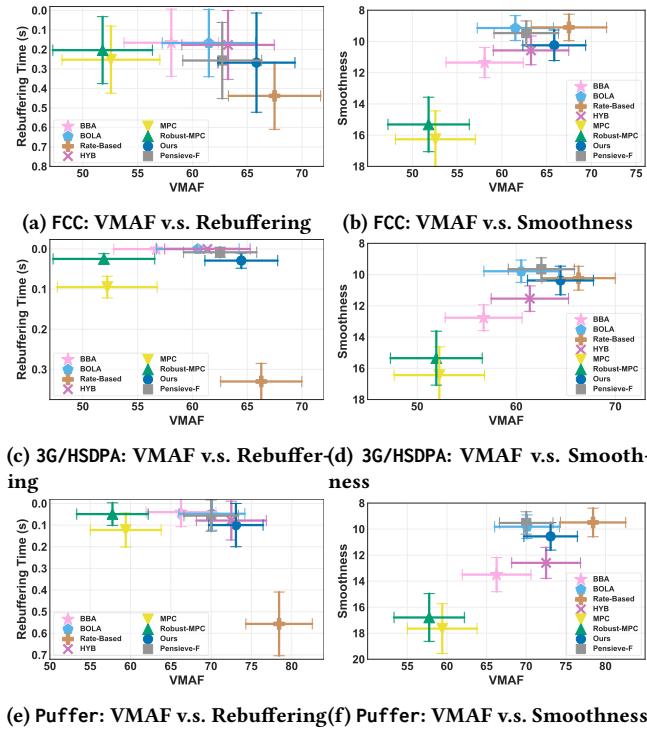


Figure 17: Breakdown performance for the last policies trained with $\#F = 3$.

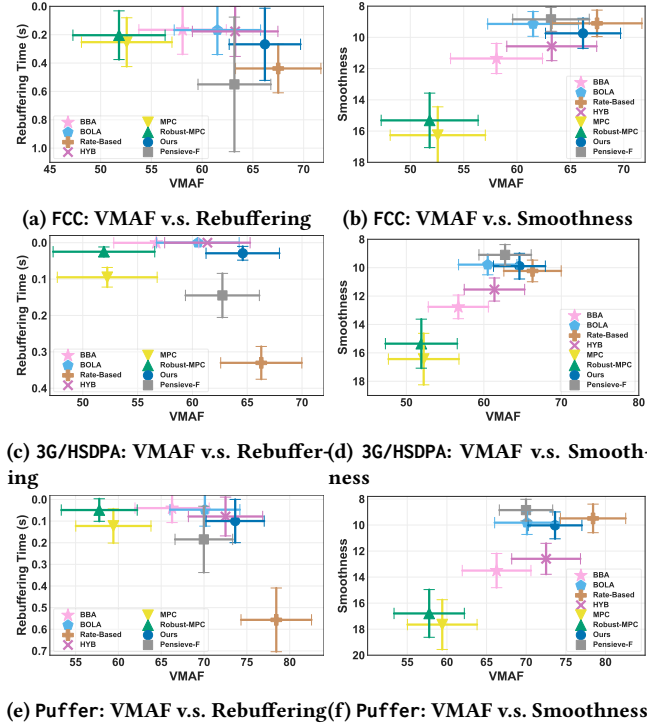


Figure 18: Breakdown performance for the last policies trained with $\#F = 5$.

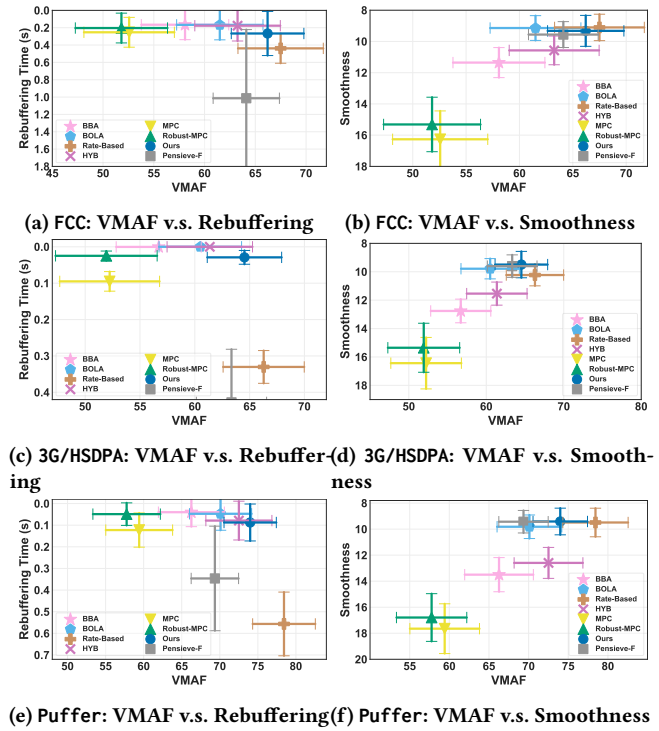


Figure 19: Breakdown performance for the last policies trained with $\#F \sim \text{Normal}(6, 3)$.

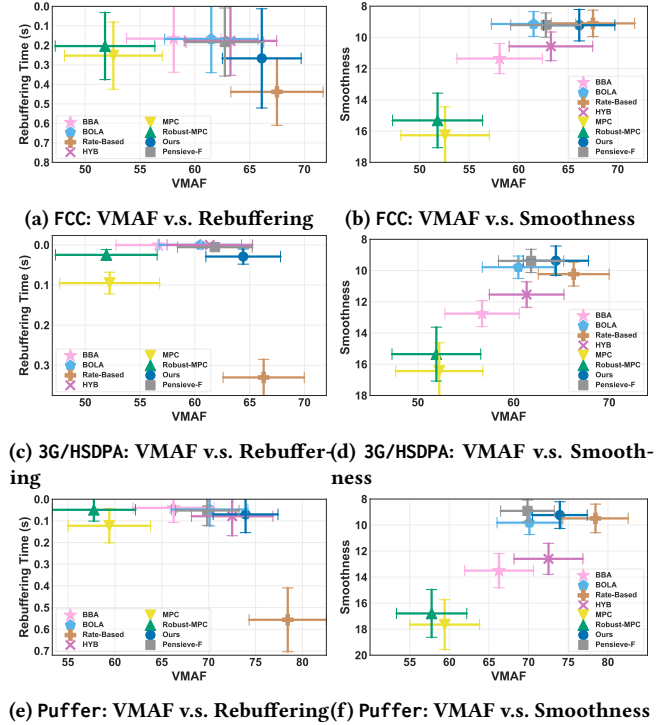


Figure 20: Breakdown performance for the last policies trained with $\#F = 10$.

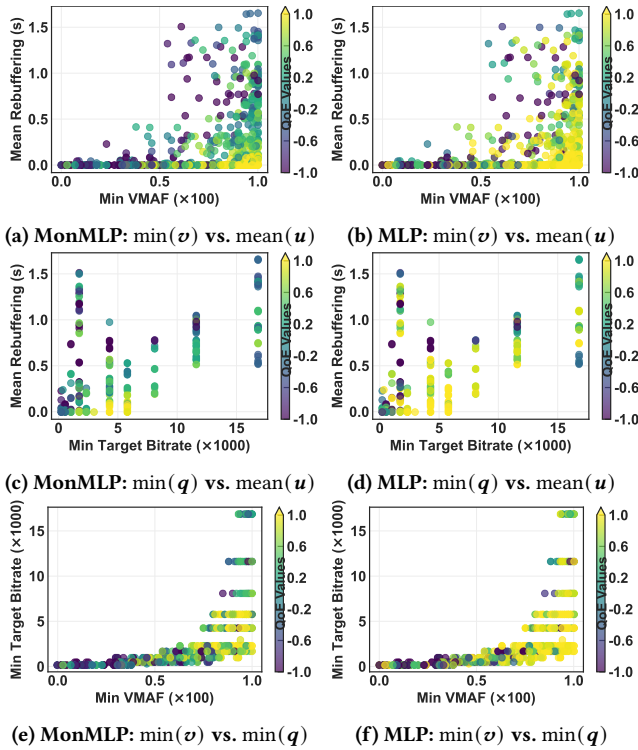


Figure 22: Visualization of the personalized QoE outputs for the first evaluator in SQoE-IV.

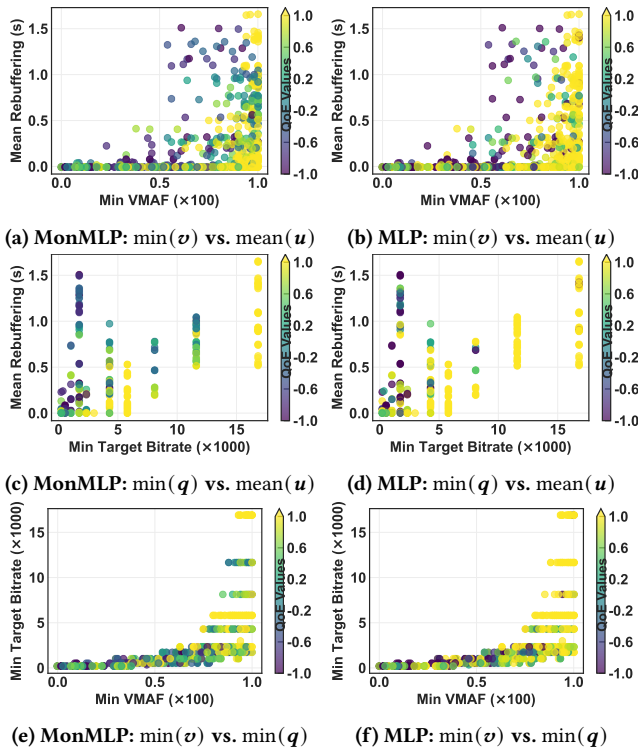


Figure 23: Visualization of the personalized QoE outputs for the second evaluator in SQoE-IV.

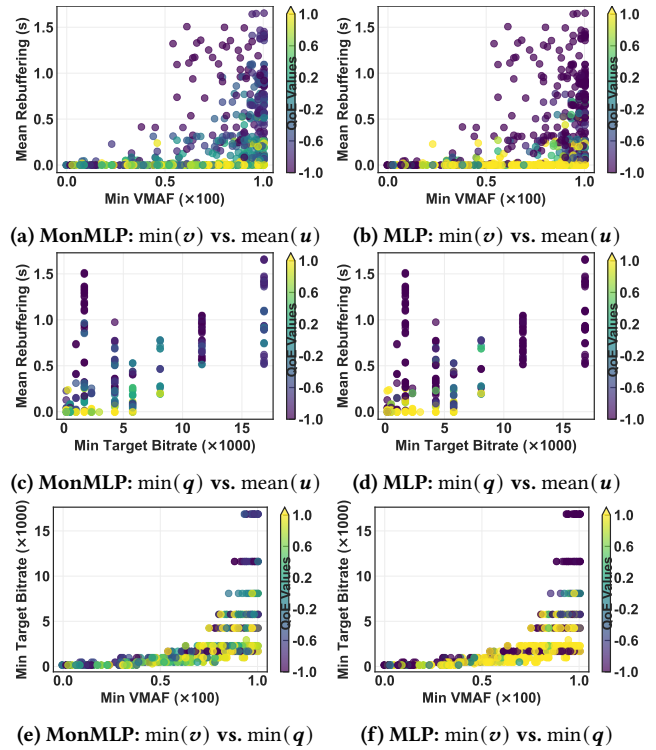


Figure 24: Visualization of the personalized QoE outputs for the third evaluator in SQoE-IV.