

---

# Self-play Reinforcement Learning for Video Transmission (Extended Version)

---

**Tianchi Huang**

Media lab, Tsinghua University.  
htc19@mails.tsinghua.edu.cn

## Abstract

Off-the-shelf learning-based video transmission techniques optimize itself towards the linear-combination of several weighted metrics with mutual restriction rather than deterministic requirements, which might finally generalize a strategy that violates the original demand. To eliminate this concern, *Zwei* aims to utilize the fundamental requirement to update the policy. Considering the given requirement often fails to directly provide any gradients, we propose a novel deep learning-based method which can effectively update the network via a binarized signal, symboling which one is closer to the assigned requirement, between two trajectories sampled from the same environment. To build *Zwei*, we have to develop video transmission simulation environments, design adequate neural network models, and invent training methods for dealing with different requirements on various video transmission scenarios. As expected, trace-driven analysis over three representative tasks shows that *Zwei* optimizes itself according to the assigned requirement faithfully, and rivals or outperforms the best existing scheme.

## 1 Introduction

Thanks to the dynamic growth of video encoding technologies and basic Internet services [9], currently humans are living with the great help of video transmission services. Users often watch interesting videos and live streaming from different video content providers (e.g. YouTube and Kuaishou), or chat with each other via real-time video communication rather than a conventional phone call. In such scenarios, the videos are required to transmit with high bitrate and less rebuffering time or stalling ratio. However, due to the fluctuation and unpredictability of network conditions, blindly achieving high bitrate may heavily increase the probabilities of the rebuffering event. Hence, several rate adaptation methods are proposed to take these factors into consideration. Further, from the video content provider’s perspective, they aim to provide video streaming services with less stalling ratio but lower costs, where it’s also necessary to trade off the stalling ratio against the cost. Thus, the services are required to employ scheduling algorithms for balancing the two. In brief, such aforementioned observations are being left on the horns of a classic dilemma: in the video transmission tasks, both quality of experience (QoE) and quality of service (QoS) are evaluated with contradicted metrics (§2.1). Hence, how to generally bridge the gap between the positive and negative factors?

Unfortunately, as much as the fundamental problem has already been published about *two decades* [22, 32], current approaches, either heuristics or learning-based methods, fall short of achieving this goal. On the one hand, heuristic-based schemes often use existing models [18, 29] or specific domain knowledge [20] as the basic working principle. However, such approaches sometimes require careful tuning and will backfire under the circumstance that violated with presumptions, which results in the failure of achieving acceptable performance under all considered scenarios [26]. On the other hand, learning-based methods [15, 23] leverage deep reinforcement learning (DRL) to train a neural

network (NN) by interacting with the environments without any presumptions, aiming to obtain higher score computed by the reward function, where the function is often defined as a linear-based equation with the combination of weighted sum metrics. Unsurprisingly, the learned policy outperforms heuristics, with the improvements on the overall performance of more than 18% [5]. Nevertheless, in this study, we empirically illustrate that i) an inaccurate reward function may mislead the RL-based algorithm to generalize bad strategies, since ii) the actual requirement can hardly be presented by the linear-based reward function with fixed weights. Moreover, iii) considering the diversity of real-world network environments, we can hardly present an accurate reward function that can perfectly fit *any* network conditions (§2.2). As a result, despite its abilities to gain higher numerical reward score, such training schemes may generalize a strategy that hardly meet the basic rules of the actual requirements.

Taking a look from another perspective, we observe that the aforementioned problem can be naturally written as a deterministic goal or requirement [14]. E.g., in the most cases, the goal of the adaptive bitrate (ABR) streaming algorithm is to achieve lower rebuffering time first, and next, reaching higher bitrate [17, 24]. It is pretty straightforward that it can be easily understood and refined by others. To this end, we attempt to train the NN based on the assigned requirement without reward engineering. Unfortunately, off-the-shelf learning-based algorithms are unable to optimize the policy like this since it lacks the abilities to provide any gradients information to guide the algorithm towards a better performance directly. Hence, we ask if self-play learning [39] can help taming the complexity of video transmission services with the actual requirement and especially, without reward engineering.

Inspired by this opportunity, we envision a self-play reinforcement learning-based framework, known as *Zwei*<sup>1</sup>, which can be viewed as a solution for tackling the video transmission dilemma (§3). The key idea of *Zwei* is to generalize a strategy which can always meet the actual requirement. Specifically, we apply Monte-Carlo (MC) search to making move decisions with autoregressive manner. In the MC search process, many simulated trajectories starting from the starting state  $s_0$  are generated following current policy  $\theta$ , and the expected long-term win rate  $r$  is estimated by averaging the battle results from each of the trajectories. The battle result are determined by a basic question: given two strategies collected from the same environment, which one is closer to the actual demand? Having estimated the long-term win rate, *Zwei* then updates the NN via increasing the probabilities of the winning sample and reducing the possibilities of the failure sample. In this work, we use state-of-the-art policy gradient method proximate policy optimization (PPO) [37] to optimize the NN. Furthermore, inspired by the recent success of state-of-the-art RL methods [27, 37], we further add NN-based baseline into *Zwei* to enhance its overall performance (§3.2).

In the rest of the paper (§4), we attempt to evaluate the potential of *Zwei* using trace-driven analyses of various representative video transmission scenarios. To achieve this, we build several faithful video transmission simulators which can accurately replicate the environment via real-world network dataset. Specifically, we validate *Zwei* on three different tasks (§2.1), including client-to-server, server-to-client, and client-to-client service. Note that each of them has individual requirements and challenges. As expected, evaluation results demonstrate the superiority of *Zwei* against existing state-of-the-art approaches on all tasks. In detail, we show that (i) *Zwei* outperforms existing ABR algorithms on both HD videos and 4K videos, with the improvements on Elo score [11] of 32.24% - 36.38%. (ii) *Zwei* performs well in crowd-sourced live streaming (CLS) scheduling task, reducing the overall costs on 22% and decreasing over 6.5% on the overall stalling ratio in comparison of state-of-the-art learning-based scheduling method LTS [48]. (iii) *Zwei* generalizes well in the real-time communication (RTC) scenario. Comparing the performance of *Zwei* and state-of-the-art heuristics WebRTC [7], we observe that *Zwei* effectively improves 11.34% on average receive rate, but reduces 20.49% on loss ratio and 62.95% on 95-percent round-trip-time (RTT). In particular, we find that *Zwei* perfectly follows the guidance of the given requirement since it demonstrates two distinct stages during the training process.

**Contributions:** This paper makes four key contributions.

- We point out the shortcoming of learning-based schemes in video transmission tasks and present the idea that update networks without reward engineering (§2.2).
- *Zwei* is a novel framework that aims to employ the self-play reinforcement learning method to make the idea practical (§3).

---

<sup>1</sup>*Zwei*: means the number *two*, or *double* in German

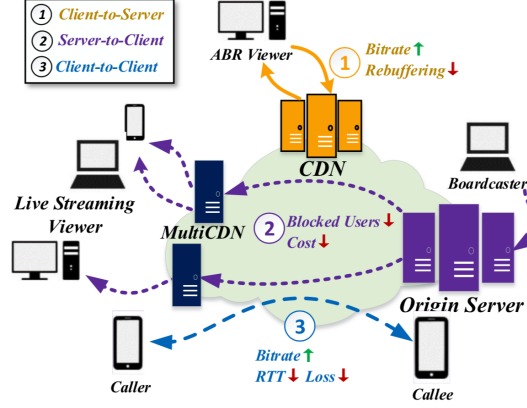


Figure 1: A brief introduction of today’s video transmission service. As shown, the service is mainly composed of client-to-server, server-to-client and client-to-client services (§2.1).

- We implement Zwei into three representative video transmission scenarios, including rate adaptation, transmit scheduling, as well as rate control. For each task, we have to develop new NN representation, determine requirements, and construct a faithful offline simulator (§4).
- We evaluate Zwei under video transmission tasks. Trace-driven experimental results illustrate that Zwei outperforms existing schemes on all considered scenarios (§4).

## 2 Background and Challenges

In this section, we first normally introduce the background of adaptive video streaming scenarios, including client-to-server, server-to-client, and client-to-client service, and then we briefly elaborate the key challenges of each service.

### 2.1 The Diversity of Video Transmission Services

To better understand the limitations that conventional video streaming service suffering from, we plot the general service work-flow on Figure 1. As shown, the service commonly consists of three parts:

1. **Client-to-Server Service.** In this scenario, users often adopts a video player to watch the video on demand. First, video contents are pre-encoded and pre-chunked as several bitrate ladders on the server. Then the video player, placed on the client side, dynamically picks the proper bitrate for the next chunk to varying network conditions. Specifically, the bitrate decisions should achieve high bitrate and low rebuffering on the entire session [6, 17]. We called it *adaptive bitrate streaming* (ABR).
2. **Server-to-Client Service.** Consider, if we were the content provider and currently we had multiple content delivery networks (CDNs) with different costs and performance, how to schedule the users’ requests to the proper CDN, aiming to provide live streaming services withing less stalling ratio and lower cost? In common, we call that *crowd-sourced live streaming* (CLS) [48].
3. **Client-to-Client Service.** Besides, in our daily life, we usually chat with other users instantly via a video call namely *Real-Time video Communication* (RTC). The RTC system consists of a sender and a receiver. The sender deploys a UDP socket channel and sends the encoded video packets to the receiver. The receiver then feedbacks messages to the sender through feedback channel. During the session, the sender adjusts the sending bitrate for next time period, aiming to achieve high bitrate and low round-trip time (RTT) as well as less loss ratio [15].

Besides, we list the requirement and details for each video transmission task in Table 2. Note that we label the underlying metric orderly since each metric has attentive priority. Details of each video transmission scenario are demonstrated in §4.

Table 1: Requirements for each video transmission tasks.

Transmission Type	Is ABR?	Requirement
<b>Client-to-Server</b>	✓	Rebuffering Time <sup>1</sup> ↓, Bitrate <sup>2</sup> ↑
<b>Server-to-Client</b>	×	Stalling Ratio <sup>1</sup> ↓, Cost <sup>2</sup> ↓
<b>Client-to-Client</b>	✓	RTT <sup>1</sup> ↓, Loss <sup>1</sup> ↓, Bitrate <sup>2</sup> ↑

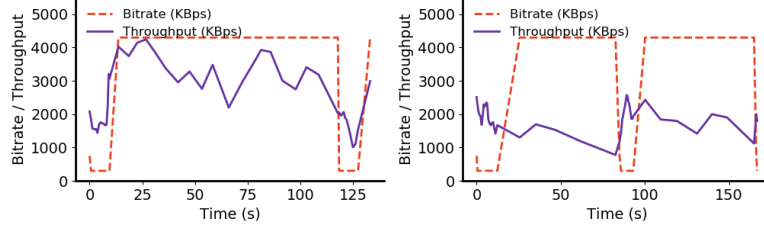


Figure 2: This group of figures shows the drawback of traditional RL-based ABR methods: After trained Pensieve with default parameter settings, we observed that: 1) the policy was simple yet effective; 2) it still maintained a high QoE score. However, it’s clear that the proposed method deviates from the basic rules of ABR.

## 2.2 Challenges

Recent video transmission algorithms mainly consist of two types, i.e., heuristics and learning-based scheme. Heuristic methods utilize an existing fixed model or domain knowledge to construct the algorithm, while they inevitably fail to achieve optimal performance in all considered network conditions due to the inefficient parameter tuning (§5.1,[23]). Learning-based schemes train a NN model towards the higher reward score from scratch, where the reward function is often defined as a linear combination of several weighted metrics [15, 23, 48]. Nevertheless, considering the characteristics of the video transmission tasks above, *we argue that the policy generated by the linear-based reward fails to always perform well on the right track [14]*. In this study, we set up two experiment in the ABR scenario (4.1) to prove this conjecture.

---

**Observation 1.** *Though recent RL-based method outperforms existing schemes, it seems like to explore some tricks for obtaining higher reward score instead of generalized a well-performed strategy.*

---

Considering the ABR process as a Markov Decision Process (MDP), in recent years, many schemes have been proposed to learn ABR algorithms via reinforcement learning (RL) method [8, 23]. Despite the outstanding performances that RL-based ABR algorithms achieve, these schemes suffer from a key limitation: they optimize their neural network via enhancing QoE scores. However, achieving a high QoE score doesn’t necessarily equal to generate an exemplary algorithm. For example, the experimental results of state-of-the-art RL-based scheme Pensieve<sup>2</sup> is illustrated in Figure 2. While Pensieve converges, its policy can be generalized as: *fetching lowest bitrates till the buffer has enough space and time to pick highest bitrates*. That is because though RL methods will successfully train the model under the case which only needs to optimize one metric (e.g., towards higher scores for Atari games [28]), such schemes lack the abilities to tackle the problem which is affected by multiple factors directly. For example, in ABR problems, recent work [23, 40, 47] leverages reward functions (e.g., QoE metrics) with a weighted sum of several underlying metrics to take the place of the multi-factor optimization (Several metrics must be optimized together). Hence, QoE driven RL-based approaches have the abilities to obtain a relative good numerical reward, but they may also provide users with unexpected and unstable viewing experience. This problem imposes critical challenges to RL-based ABR algorithm.

---

<sup>2</sup><https://github.com/hongzimao/pensieve>

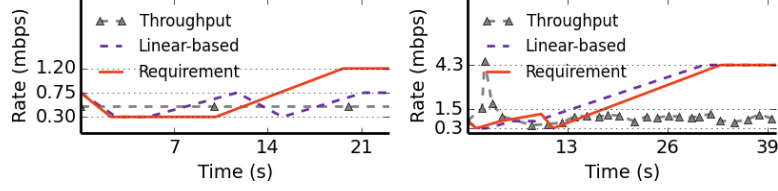


Figure 3: The comparison of linear-based optimal and requirement-based optimal strategy. Results are evaluated on *fixed (0.5mbps)* and *3G* network traces under ABR scenario. We can see the difference between the actual requirement and the optimal trajectory generated by the linear-based reward function.

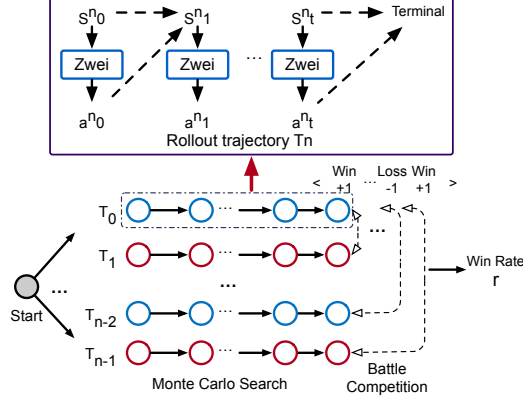


Figure 4: Overview of Self-play Reinforcement Learning Framework (Zwei). The framework is mainly composed of *four* phases: Monte Carlo (MC) searching, battle competition, win rate estimation, and policy optimization.

---

**Observation 2.** *Recent linear-based weighted reward function can hardly map the actual requirement for all the network traces.*

---

Moreover, we attempt to compare the optimal strategy of linear-based with requirement-based under two representative network traces, in which linear-based optimal is the policy which obtains maximum reward, requirement-based optimal stands for the closest strategy in terms of the given requirement. Unsurprisingly, from Figure 3 we observe that linear-based optimal policy performs differently compared with the requirement-based optimal strategy. The reason is that the linear-based optimal policy heavily relies on the given weights, while the weighted parameters are not allowed to be adjusted dynamically according to current network status (e.g., throughput and jitter [5]), which finally leads to the failure of guiding the optimization process on the right track. Generally, we believe that the policy learned by reward engineering might fall into the unexpected conditions.

**Summary.** In general, we observe that no matter how precisely and carefully the parameter of the linear-based reward function tunes, such tuned functions can hardly meet the requirement of any network conditions. E.g., the parameter of stable and unstable network conditions are not the same. To that end, traditional learning-based scheme, which often optimizes the NN via the assigned functions, will eventually fail to provide a reliable performance on any network traces. We therefore, attempt to learn the strategies from the original demands.

### 3 Zwei Design

In this section, we briefly introduce the details of Zwei, including its key principle, the basic training algorithm and the advanced improvements.

### 3.1 Self-play Method

AlphaZero [38] is the most famous self-play algorithm who trains the NN solely based on an agent competes against the others. To tackle the traditional RL methods' problem, we thereby consider following AlphaZero to train the algorithm via self-play RL. The main idea of our proposed framework is to leverage these search operators repeatedly in a policy iteration procedure. However, static games with incomplete information (i.e., video transmission tasks are often defined as partially observable Markov decision process (POMDP) problems [26].) are much dissimilar and more complicated than dynamic games with complete information such as Go or Chess. Thus, we propose a novel self-play reinforcement learning framework, namely Zwei, aiming to tackle a critical challenge: *considering the specific features for video transmission task, how to design a proper general self-play reinforcement learning algorithm?*

### 3.2 Basic Idea

We now start to explaining the fundamental idea for Zwei. Figure 4 presents the main phases in our framework. The pipeline can be summarized as follows:

**Phase1: Monte Carlo (MC) search.** First, we adopt MC search method [44] to sample  $N$  different trajectories  $T_n = \{s_0^n, a_0^n, s_1^n, a_1^n, \dots, a_t^n\}, n \in N$  w.r.t the given policy  $\pi(s)$  under the *same environments* ( $a_t \sim \pi(s_t)$ ) and start point (the gray point in Figure 4). Next, we record and analysis the underlying metric for the entire session. Finally, we store all the sample  $T_n$  into  $D$ . Note that we can select Monte Carlo Tree Search (MCTS) [39], which is widely used in advanced research, to implement the process.

**Phase2: Battle Competition (BC).** To better estimate how the current policy performs, Zwei requires a module to label all trajectories from  $D$ : given two *different* trajectories,  $T_i$  and  $T_j$  which are all collected from the *same environment settings* ( $T_i, T_j \in D$ ), we attempt to identify which trajectory is positive for NN updating, and which trajectory is generated by the worse policy. Thus, we implement a rule called Rule which can determine the *better* trajectory between the given two candidates, in which *better* means which trajectory is closer to the requirement. At the end of the session, the terminal position  $s_t$  is scored w.r.t the rules of the requirement for computing the game outcome  $o$ :  $-1$  for a loss,  $0$  for a draw, and  $+1$  for a win. The equation is listed in Eq. 1.

$$o_i^j = \text{Rule}(T_i, T_j). \quad (1)$$

$$s.t. \quad o_i^j = \{-1, 0, 1\}, T_i, T_j \in D, \quad (2)$$

$$T_i \neq T_j. \quad (3)$$

**Phase3: Win Rate Estimation.** Next, having computed the competition outcome  $o_i$  for any two trajectories, we then attempt to estimate the average win rate  $r_i$  for each trajectory  $T_i$  in  $D$ . The equation is listed in Eq. 4. Notice that the accuracy of the win rate estimation heavily depends on the number  $N$  of trajectories. Since it's impractical to sample infinite number of samples in the real world, we further list the performance comparison of different sample numbers in §4.1.

$$r_i = \mathbb{E}[\text{Rule}(T_i, )] \quad (4)$$

$$= \lim_{N \rightarrow \infty} \frac{\sum_u^N o_i^u}{N}. \quad (5)$$

**Phase4: Policy Optimization.** In this part, given a batch of collected samples and their win rate, our key idea is to update the policy via *elevating* the probabilities of the winning sample from the collected trajectories and *diminishing* the possibilities of the failure sample from the worse trajectories. In other words, the improved policy  $\pi$  at state  $s_t$  is required to pick the action  $a_t$  which produced the best estimated win rate  $r_t$ , i.e.,  $a_t = \text{argmax}_a E[r_t(s_t, a)]$ . Hence, we can employ any policy gradient method (e.g., A3C [27], TRPO [36], ACKTR [46]) to optimize the NN's policy. In this work, We use Proxy Policy Optimization (PPO) [37], a state-of-the-art actor-critic method, as the basic RL algorithm. Briefly, PPO uses clip method to restrict the step size of the policy iteration and update the NN by minimizing the following *clipped surrogate objective*. We list the Zwei's loss function  $\mathcal{L}^{\text{PPO}}(\theta)$  in Eq. (6).

$$\mathcal{L}^{\text{PPO}}(\theta) = \mathbb{E}_t \left[ \min \left( p_t(\theta) r_t, \text{clip}(p_t(\theta), 1 - \epsilon, 1 + \epsilon) r_t \right) \right], \quad (6)$$



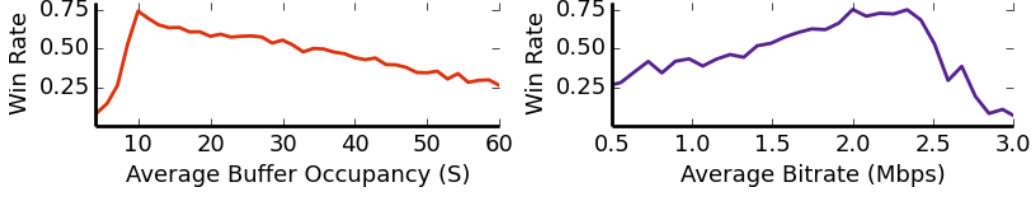


Figure 5: We plot the underlying win rate for different average bitrate and average buffer occupancy.

where  $p_t(\theta)$  denote the probability ratio between the policy  $\pi_\theta$  and the old policy  $\pi_{\theta_{\text{old}}}$ , i.e.,  $\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ ;  $\epsilon$  is the hyper-parameters which controls the clip range. In this paper, we set  $\epsilon = 0.2$  as suggested by the original PPO paper [37].

As listed before, Zwei’s basic training algorithm is derived from a fundamental presumption: the win rate of each state  $s_t$  will be converged to zero, that means, *a draw game*. However, we argue that the distribution of win rate on each state  $s$  is different. Thus, we setup an experiment in the ABR scenario: we select a short video clip encoded with 6 bitrate levels and 5 chunks, and list all the possible trajectory under the same network trace. Then we compute the underlying win rate for each state  $s_t$  via Rule, where  $s_t$  is represented by *Average Bitrate* and *Average Buffer*. This makes sense since recent work have proved that these two metrics are critical feature for ABR task [42]. Experimental results on Figure 5 illustrate that *different states map different win rates*. E.g., the win rate will degrade to 25% if average bitrate is lower than 1 mbps and the average buffer is higher than 40 seconds. At the same time, win rate will achieve the highest score under the conditions that average bitrate is about 2 mbps and the average buffer is in the range of 10 seconds to 20 seconds. What’s more, we also consider the correlation between the average buffer and win rate, as well as the relationship between average bitrate and win rate respectively. Results also indicate the influence on win rate by different values of current states.

To this end, we further add the additional baseline  $V_{\theta_p}(s)$  to avoid the bias which is caused by the high variance in each situation. In details, Zwei with Baselines consists of a policy network and a value network. The gradient of policy network are computed as Eq. 7, where  $\hat{A}_t$  is the advantage function (Eq. 8),  $V_{\theta_p}(s)$  is provided by another value NN. Meanwhile, we update the value network via minimizing the mean square error between  $r_t$  and  $V_{\theta_p}(s)$  (Eq. 9). Moreover, we also add  $H(s_t; \theta)$ , representing the entropy of the policy, to encourage exploration. In this work, we use adaptive entropy decay method to dynamically adjust the entropy weight  $\beta$  from 0.5 to 0.01. For more details, please refer to our repository [13]. To sum up, we summarize Zwei’s loss function  $\mathcal{L}^{\text{Zwei}}$  in Eq. 10.

$$\mathcal{L}^{\text{Policy}} = \mathbb{E}_t \left[ \min \left( p_t(\theta) \hat{A}_t, \text{clip}(p_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right]. \quad (7)$$

$$\hat{A}_t = r_t - V_{\theta_p}(s_t) \quad (8)$$

$$\mathcal{L}^{\text{Value}} = \frac{1}{2} \mathbb{E}_t [V_{\theta_p}(s_t) - r_t]^2. \quad (9)$$

$$\nabla \mathcal{L}^{\text{Zwei}} = \nabla_\theta \mathcal{L}^{\text{Policy}} + \nabla_{\theta_p} \mathcal{L}^{\text{Value}} + \nabla_\theta \beta H(s_t; \theta). \quad (10)$$

### 3.3 Training Methodology Overview

Training procedure of Zwei is summarized in Algorithm 3. As shown, Zwei’s workflow mainly consists of two parts, i.e., exploration and exploitation:

1. *Exploration*. First, we sample  $n$  trajectories roll-outed by the current policy in the same environment. Then we leverage Rule to determine the better strategy via the pairwise comparison. Next, we estimate the win rate w.r.t the given samples.
2. *Exploitation*. We update the gradient by using the loss function  $\mathcal{L}^{\text{Zwei}}$  (Eq. 10).

---

**Algorithm 1** Zwei’s Overall Training Procedure

---

**Require:** Environment Env; Rule Rule.

```
1: Initialize parameters  $\theta$  with random weights;
2: repeat
3:    $D = \{\}$ ;
4:   Env  $\leftarrow$  Sample environment settings, e.g., network trace T, video V, or CDN configuration  $c$ ;
5:   for  $i \in \{1, 2, \dots, N\}$  do ▷ MC search;
6:     Get state  $s_t$  from Env.
7:      $T_i = \{\}$ ;
8:     while not Terminal do
9:       Perform  $a_t$  according to policy  $\pi_\theta(a_t; s_t)$ .
10:       $T_i \leftarrow T_i \cup \{s_t, a_t\}$ .
11:      Receive new state  $s_{t+1}$ 
12:    end while
13:     $D \leftarrow D \cup T_i$ 
14:  end for
15:  for  $\text{pairs}(T_u, T_v) \in D$  do ▷ Battle competition;
16:    Determine win or loss:  $o_u^v = \text{Rule}(T_u, T_v)$ ;
17:  end for
18:  for  $T_u \in D$  do ▷ Win rate estimation;
19:    Estimate win rate:  $r_u = \frac{\sum_i^N o_u^i}{N}$ .
20:    for  $s_t, a_t \in T_u$  do ▷ Policy Optimization;
21:      Update  $\theta$  and  $\theta_p$  with  $\mathcal{L}^{\text{Zwei}}$  (Eq. 10) using collected samples and win rate  $(s_t, a_t, r_u)$ ;
22:    end for
23:  end for
24: until Converged
```

---

### 3.4 Parallel Training

During the training process, we observe that the training progress is inefficient while using a single process. Inspired by the multi-agent training method [27], we modify Zwei’s training in the single agent as training in multi-agents. Multi-agents training consists of two parts, a central agent and a group of forwarding propagation agents. The forward propagation agents only decide with both policy and critic via state inputs and neural network model received by the central agent for each step; then it sends the  $n$ -dim vector containing  $\{s, a, r\}$  to the central agent. The central agent uses the actor-critic algorithm to compute gradient and then updates its neural network model. Finally, the central agent pushes the updated network parameters to each forward propagation agent. Note that this can happen asynchronously among all agents, for instance, there is no locking between agents. By default, Zwei employs 12 forward propagation agents and one central agent.

### 3.5 Implementation

We use TFlearn [43] to implement the NN and leverage TensorFlow [3] to construct Zwei. Zwei consists of two NNs, policy network and value network. Policy network takes a  $n$ -dims vector with *softmax* active function as the output. The value network outputs a value with *tanh* function scaled in  $(-1, 1)$ . Considering the characteristics on video transmission tasks, we construct different NN architectures for each task. In contrast, we use the same set of hyper-parameters for training the NN: sample number  $N = 16$ , learning rate  $\alpha = 10^{-4}$ , and the entropy weight decay  $\zeta = 0.8$ . In addition, we adopt Adam optimizer [21] with default settings. Notice that we will discuss the NN architecture and the effect of different sample number in §4.1.

## 4 Evaluation

In this section, we thoroughly evaluate the performance of the schemes which utilizes Zwei to train the NN. We mainly consider three tasks which have been described before.



Table 2: We list the requirement and details for each video transmission task. Note that we label the underlying metric orderly since each metric has attentive priority. Hence, Zwei will learn the policy according to that sequence.

Task Name	Transmission Type	ABR?	Requirement
ABR (§4.1)	<i>Client-to-Server</i>	✓	Rebuffering Time <sup>1</sup> ↓, Bitrate <sup>2</sup> ↑
CLS (§4.2)	<i>Server-to-Client</i>	×	Stalling Ratio <sup>1</sup> ↓, Cost <sup>2</sup> ↓
RTC (§4.3)	<i>Client-to-Client</i>	✓	RTT <sup>1</sup> ↓, Loss <sup>1</sup> ↓, Bitrate <sup>1</sup> ↑

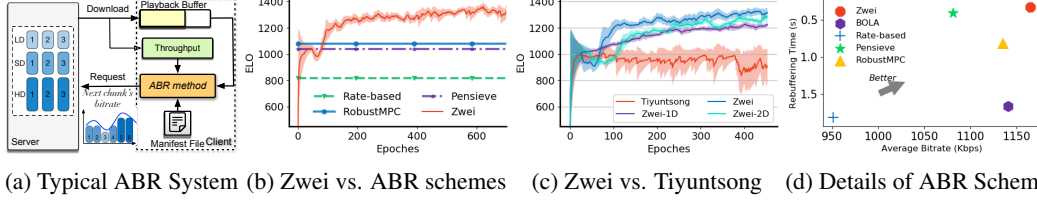


Figure 6: This group of figures show the principle of ABR system, the comparison results of Zwei and other ABR approaches, and the comparison details in terms of average bitrate and rebuffering time.

#### 4.1 ABR Scenarios

**Experimental Setup** In this part, we utilize the standard ABR’s emulation environment, including Mahimahi [30] and Park [25], to evaluate Stick. In details, we train and validate Zwei on various network throughput dataset, such as HSDPA [34], FCC [33] and Oboe [5]. We also use the video dataset provided by Huang et al. [16] during the training process, and adopt *EnvivoDash3*, a video that commonly used in recent work [5, 16, 23], to validate Zwei during the testing process. We divide the dataset into two parts, 50% of the database for training and 50% for testing. Training process lasted approximate 12000 steps, within 2 hours to obtain a reliable result.

**ABR Baselines** In this paper, we select several representational ABR algorithms from various type of fundamental principles:

1. **Rate-based** [19]: the basic baseline for ABR problems. It leverages *harmonic mean* of past five throughput measured as future bandwidth.
2. **BOLA** [40]: the most popular buffer-based ABR scheme in practice. BOLA turns the ABR problem into a utility maximization problem and solve it by using the Lyapunov function. We use BOLA provided by the authors.
3. **RobustMPC** [47]: a state-of-the-art heuristic method which maximizes the objectives by jointly considered the buffer occupancy and throughput predictions. We implement *RobustMPC* by ourselves.
4. **Pensieve** [23]: the state-of-the-art learning-based ABR scheme, which utilizes DRL to select bitrate for next video chunks. We use the pre-trained Pensieve model provided by the authors.
5. **Tiyuntsong** [14]: the first study of multi-objective optimization ABR approach. Tiyuntsong uses actor-critic method to update the NN via the competition with two agents under the same network condition. Note that Tiyuntsong’s learning agent uses DRL to update gradients w.r.t the sample generated by the agent itself, which not only sample inefficient but also reaching the optimal policy.

**Zwei vs. Existing schemes** In this experiment, we compare Zwei with existing ABR schemes under the HSDPA dataset. Results are computed as elo-score [11] and reported in Figure 11(a). We can see that Zwei outperforms recent ABR approaches. In particular, Zwei improves elo-score on 36.38% compared with state-of-the-art learning-based method Pensieve, and increases 31.11% in terms of state-of-the-art heuristic method RobustMPC. Besides, we also illustrate the CDF of the proposed methods on average bitrate and average rebuffering in Figure 11(c)(d). Results show that Zwei can not only achieve highest bitrate but also obtain lowest rebuffering under all network traces. Moreover, we also compare Zwei with previously proposed self-learning method Tiyuntsong on the same experimental settings. During the training process, we record and report the elo-curve on

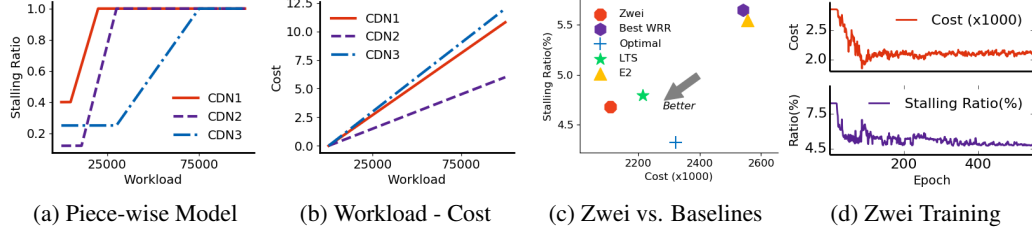


Figure 7: Results of Zwei on the CLS environment. We collect Zwei’s training curve during the training process and observe that there also exists *two stages* on the entire process.

the same validation set in Figure 11(b). As expected, Zwei with any NN architectures outperform Tiyuntsong on average elo-score of 35%.

**Zwei with Different NN Architectures** This experiment evaluates Zwei with all considered NN architectures and compares their performance under the same network setting. Zwei-1D is the standard ABR NN architecture which proposed by Pensieve [23], Zwei-2D stands for the advanced ABR NN architecture that proposed by QARC [15], and Zwei only leverages three fully-connected layers sized {128, 64, 64}. Experimental results are calculated with elo-score and shown in Figure 11(b). Unsurprisingly, when Zwei trains with some complicated NN architecture (Zwei-1D and Zwei-2D), it generalizes poorly and underperforms fully connected NN scheme. This makes sense since ABR is indeed a light-weighted task which can be solved in a practical and uncomplicated manner instead of a NN incorporating some *deep* yet *wide* layers [10].

**Two-Stage Phenomenon** By meticulously observing Zwei’s training curve on Figure 11(b), we find the training process can be split into two stages: 1) Zwei picked a lower bitrate to avoid rebuffering time when the training step is less than 100. Then 2) Zwei realized a new strategy and increased the bitrate with guaranteeing rebuffering events. The reason is Zwei skillfully exploits catastrophic forgetting [12] and abandons far-reaching decisions.

## 4.2 CLS Scheduling

**Testbed Setup.** As suggested by previous work [48], our experiments are conducted on the real-world CLS dataset provided by anonymous <sup>3</sup>, spanning 1 week (6 days for training and 1 day for test). At each time, we select 3 candidates from total 4 different CDN providers, and we fit a separate simulator for them. We train Zwei for about 4 hours to converge in a stable result.

**CDN Configuration** We detail the CDN configuration on Figure 7. Followed by previous work [48], we use a piecewise linear model to characterize this relationship between workload and block ratio (Figure 7a). Note that the following features are extracted by the real CDN dataset. What’s more, we define the CDN pricing model w.r.t various CDN providers in industry, such as Amazon E2, and Tencent CDN, and plot the correlation between the workload and cost on Figure 7b.

**Baselines** Like other scenarios, we also compare Zwei with the following state-of-art scheduling baselines:

1. **Weighted round-robin** [4]. The idea of weighted round-robin (WRR) is: the requests will be redirected to different CDN providers w.r.t a constant ratio. We adopt the algorithm with the *best* parameters.
2. **E2** [20]. Exploitation and Exploration (E2) algorithm utilizes harmonic mean for estimating CDN providers’ performance, and select with the highest upper confidence bound of reward. We use E2 algorithm provided by the authors [20].
3. **LTS** [48]. State-of-the-art CLS algorithm which uses deep reinforcement learning to train the NN towards lower stalling ratio. However, it ignores the trade-off between the cost and the performance. We evaluate LTS by Zhang et. al. [48].

**Zwei vs. State-of-the-art Scheduling Methods** We start to study how well that Zwei achieves under the CLS scenario. As shown in Figure 7c, we find that Zwei stands for the best scheme among

<sup>3</sup>Masked due to the double-blinded policy

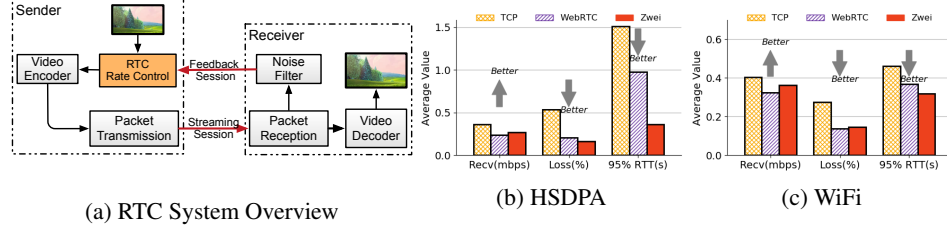


Figure 8: RTC System Overview. We implement the system suggested by prior research [15].

the candidates. Specifically, Zwei reduces the overall costs on 22% compared with state-of-the-art learning-based method LTS, and decreases over 6.5% in terms of the overall stalling ratio. The reason is that LTS takes the weight-summed combination function as the reward, while the function can hardly give a clearer guidance for the optimized algorithm. Moreover, comparing the performance of Zwei with the optimal strategy generated by the reward function, we observe that both optimal policy and Zwei are in the Pareto front. Zwei consumes less pricing costs than the optimal policy since the requirement is to *minimize the cost first*.

**Training Deep Dive** Besides, we also present the training process in Figure 7d. As shown, Zwei converges in less than 600 epochs, which needs about 3 hours. It’s worth noting that Zwei also experienced two stages on the CLS task. The first stage ranges from 0 to 100 epochs, and we can see the goal of Zwei is to minimize the cost without considering the number of stalling ratio. The rest of the stage we find that Zwei attempts to reduce the number of stalling ratio, and at the same time, the cost is converge to a steady state. Such observation also proves that Zwei learns the strategies by following the given requirements.

### 4.3 RTC Task

**Implementation Details** Figure 8 elaborates RTC system overview. We place Zwei as the RTC control module on the sender-side. Followed by recently proposed approaches [15], we have to develop a faithful network emulator which can simulate the network condition according to network saturate traces that collected from different network conditions. The network traces is composed of 3G, 4G/LTE and WiFi network as well as self-collected real-world network environments. Then we train Zwei under the simulator for about 8 hours and test it on the Mahimahi [30].

**RTC Baselines** We compare Zwei with several proposed heuristic methods, including

1. **TCP** [29]. The most popular congestion control algorithm in the Internet. The key insight is the additive-increase, multiplicative-decrease (AIMD) algorithm, which means a linear increasing of the congestion window with an exponential reduction when congestion is detected. We adopt TCP-Reno with *fast recovery*.
2. **WebRTC** [7]. State-of-art RTC scheme on both academic and industry. More precisely speaking, WebRTC leverage delay-based and loss-based module to jointly control the rate adaption problem in the real-time scenario, where its parameters have been carefully tuned for almost one decade. In this work, we adopt WebRTC from its open-sourced repository.

**Zwei vs. RTC Baselines** In this experiment, we compare Zwei with existing proposed heuristic methods (§4.3) on various network conditions such as 3G and Wifi. We collect average receiving bitrate, 95 percent round-trip-time (RTT) as well as average loss ratio every 1 second. Note that these metrics are widely used in RTC evaluation [15]. Results are reported in Figure 8, we demonstrate that Zwei improves 11.34% on average receiving bitrate, reduces 20.49% in terms of loss ratio, and decreases 62.95% on 95-percent RTT compared with WebRTC under HSDPA dataset. Another results on the WiFi network conditions show that Zwei also outperforms WebRTC, with the improvements on sending rate of 12.08% and decreasing in 95% RTT of 13.59%. The comparison of Zwei and TCP illustrates that Zwei can effectively detect the congestion signal and dynamically adjust the sending rate to avoid high loss ratio and RTT.

**Zwei Case Study** To better understand the performance of the selected methods, we **randomly** pick two traces from the validation dataset and compare the sending rate curve of Zwei with TCP and

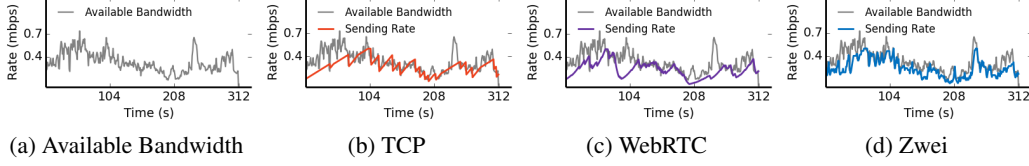


Figure 9: The comparison of Zwei and existing methods on the HSDPA network trace.

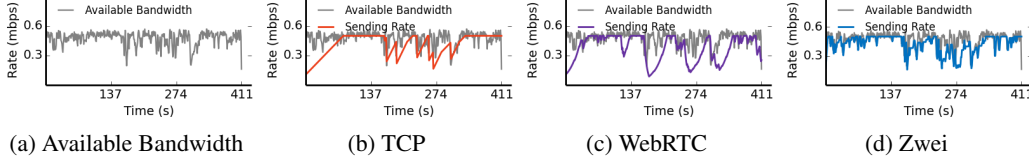


Figure 10: We also compare Zwei with existing methods on the trace collected from the real-world WiFi scenario.

WebRTC. As demonstrated in Figure 9 and Figure 10, Zwei can accurately estimate the bottleneck of the session, as it always selects the high rate with avoiding congestion.

## 5 Related Work

### 5.1 Heuristic Methods

**Adaptive Bitrate Video Streaming** Heuristic-based ABR methods often adopts throughput prediction (E.g., FESTIVE [19]) or buffer occupancy control (E.g., BOLA [40]) to handle the task. However, such approaches suffer from either inaccurate bandwidth estimation or long-term bandwidth fluctuation problems. Then, MPC [47] picks next chunks’ bitrate by jointly considering throughput and buffer occupancy. Nevertheless, MPC is sensitive to its parameters since it relies on well-understanding different network conditions. What’s more, Oboe [5] even proposes an auto-tuning method.

**CSP Selection** Through preliminary measurements, it is widely accepted that the strategies are largely statically configured [4]. In recent years, dynamic scheduling across different CSPs has received more attention. Jiang et. al. [20] uses E2 method to replace traditional model-based methods.

**RTC Rate Adaption** Rate adaption methods are mainly composed of loss-based approach [29], attempting to increase bitrate till packet loss occurs, and delay-based approach [35], aiming to adjust sending rate to control the transmission delay. However, such methods are sensitive to network conditions. Thus, model-based approach, such as WebRTC [7], controls the bitrate based on previous status observed.

### 5.2 Learning-based Schemes

Recent years, several learning-based attempts have been made to tackle video transmission problem. For example, Mao et al. [23] develop an ABR algorithm that uses DRL to select next chunks’ bitrate. LTS [48] is a DRL-based scheduling approach which outperforms previously proposed CLS approaches. Furthermore, QARC [15] optimizes a NN towards higher perceptual video quality and less stalling events. However, such methods fail to achieve actual requirements since they are optimized via a linear-based reward function.

## 6 Conclusion

We propose Zwei, which utilize the comparison results between two samples to enhance itself towards a better *win rate*. Zwei can generalize an outstanding strategy based on the actual requirement, where the requirement is always hard to be described as a linear-based manner. Results show that Zwei outperforms recent state-of-the-art work on average improvements of more than 22% in terms of three fundamental video transmission tasks.

## References

- [1] Dash industry forum | catalyzing the adoption of mpeg-dash, 2019.
- [2] Http live streaming. <https://developer.apple.com/streaming/>, 2019.
- [3] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, 2016.
- [4] V. K. Adhikari, Y. Guo, F. Hao, M. Varvello, V. Hilt, M. Steiner, and Z.-L. Zhang. Unreeling netflix: Understanding and improving multi-cdn movie delivery. In *INFOCOM'12*, 2012.
- [5] Z. Akhtar, P. Adria, M. Mirza, et al. Oboe: auto-tuning video abr algorithms to network conditions. In *SIGCOMM 2018*, 2018.
- [6] A. Bentaleb, B. Taani, A. C. Begen, C. Timmerer, and R. Zimmermann. A survey on bitrate adaptation schemes for streaming media over http. *IEEE Communications Surveys & Tutorials*, 2018.
- [7] G. Carlucci, L. De Cicco, S. Holmer, and S. Mascolo. Analysis and design of the google congestion control for web real-time communication (webrtc). In *MMSys'16*, 2016.
- [8] F. Chiariotti, S. D'Aronco, L. Toni, and P. Frossard. Online learning adaptation strategy for dash clients. In *Proceedings of the 7th International Conference on Multimedia Systems*, 2016.
- [9] Cisco. Cisco visual networking index: Forecast and methodology, 2016-2021, 2017.
- [10] A. Dethise, M. Canini, and S. Kandula. Cracking open the black box: What observations can tell us about reinforcement learning agents. In *Proceedings of the 2019 Workshop on Network Meets AI & ML*, 2019.
- [11] A. Elo. *The rating of chessplayers, past and present*. Arco Pub., 1978.
- [12] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013.
- [13] T. Huang, Zwei. <https://github.com/godka/Zwei>, 2020.
- [14] T. Huang, X. Yao, C. Wu, R.-X. Zhang, and L. Sun. Tiyuntsong: A self-play reinforcement learning approach for abr video streaming. *arXiv preprint arXiv:1811.06166*, 2018.
- [15] T. Huang, R.-X. Zhang, C. Zhou, and L. Sun. Qarc: Video quality aware rate control for real-time video streaming based on deep reinforcement learning. In *Multimedia'18*, 2018.
- [16] T. Huang, C. Zhou, R.-X. Zhang, C. Wu, X. Yao, and L. Sun. Comyco: Quality-aware adaptive video streaming via imitation learning. *arXiv preprint arXiv:1908.02270*, 2019.
- [17] T.-Y. Huang, C. Ekanadham, A. J. Berglund, and Z. Li. Hindsight: Evaluate video bitrate adaptation at scale. In *MMSys'19*, MMSys '19, pages 86–97, New York, NY, USA, 2019. ACM.
- [18] T.-Y. Huang, R. Johari, N. McKeown, M. Trunnell, and M. Watson. A buffer-based approach to rate adaptation: Evidence from a large video streaming service. *SIGCOMM'14*, 44(4), 2015.
- [19] J. Jiang, V. Sekar, and H. Zhang. Improving fairness, efficiency, and stability in http-based adaptive video streaming with festive. *TON*, 22(1), 2014.
- [20] J. Jiang, S. Sun, V. Sekar, and H. Zhang. Pytheas: Enabling data-driven quality of experience optimization using group-based exploration-exploitation. In *NSDI*, volume 1, 2017.
- [21] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [22] C. Lefelhocz, B. Lyles, S. Shenker, and L. Zhang. Congestion control for best-effort service: why we need a new paradigm. *IEEE Network*, 10(1):10–19, 1996.

- [23] H. Mao, P. Adria, M. Mirza, et al. Neural adaptive video streaming with pensieve. In *SIGCOMM'17*, 2017.
- [24] H. Mao, S. Chen, D. Dimmery, S. Singh, D. Blaisdell, Y. Tian, M. Alizadeh, and E. Bakshy. Real-world video adaptation with reinforcement learning. 2019.
- [25] H. Mao, P. Negi, A. Narayan, H. Wang, J. Yang, H. Wang, R. Marcus, R. Addanki, M. Khani, S. He, et al. Park: An open platform for learning augmented computer systems. 2019.
- [26] H. Mao, S. B. Venkatakrisnan, M. Schwarzkopf, and M. Alizadeh. Variance reduction for reinforcement learning in input-driven environments. *ICLR'19*, 2019.
- [27] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML'17*, 2016.
- [28] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [29] J. Mo, R. J. La, V. Anantharam, and J. Walrand. Analysis and comparison of tcp reno and vegas. In *INFOCOM'99*, volume 3, 1999.
- [30] R. Netravali, A. Sivaraman, S. Das, A. Goyal, K. Winstein, J. Mickens, and H. Balakrishnan. Mahimahi: Accurate record-and-replay for {HTTP}. In *ATC'15*, 2015.
- [31] P. G. Pereira, A. Schmidt, and T. Herfet. Cross-layer effects on training neural algorithms for video streaming. In *Proceedings of the 28th ACM SIGMM Workshop on Network and Operating Systems Support for Digital Audio and Video*, 2018.
- [32] R. Rejaie, M. Handley, D. Estrin, et al. Quality adaptation for congestion controlled video playback over the internet. In *SIGCOMM'99*, 1999.
- [33] M. F. B. Report. Raw data measuring broadband america 2016. <https://www.fcc.gov/reports-research/reports/measuring-broadband-america/raw-data-measuring-broadband-america-2016>, 2016. [Online; accessed 19-July-2016].
- [34] H. Riiser, P. Vigmostad, C. Griwodz, and P. Halvorsen. Commute path bandwidth traces from 3g networks: analysis and applications. In *Proceedings of the 4th ACM Multimedia Systems Conference*, 2013.
- [35] D. Rossi et al. Ledbat: The new bittorrent congestion control protocol. In *ICCCN*, 2010.
- [36] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897, 2015.
- [37] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [38] D. Silver, A. Huang, C. J. Maddison, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587), 2016.
- [39] D. Silver, T. Hubert, J. Schrittwieser, et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *CoRR*, abs/1712.01815, 2017.
- [40] K. Spiteri et al. Bola: Near-optimal bitrate adaptation for online videos. In *INFOCOM'16*, 2016.
- [41] K. Spiteri, R. Sitaraman, and D. Sparacio. From theory to practice: improving bitrate adaptation in the dash reference player. In *Proceedings of the 9th MMSys*, 2018.
- [42] Y. Sun and et al. Cs2p: Improving video bitrate selection and adaptation with data-driven throughput prediction. In *SIGCOMM 2016*, 2016.
- [43] Y. Tang. Tf. learn: Tensorflow's high-level module for distributed machine learning. *arXiv preprint arXiv:1612.04251*, 2016.



- [44] G. Tesauro and G. R. Galperin. On-line policy improvement using monte-carlo search. In *Advances in Neural Information Processing Systems*, pages 1068–1074, 1997.
- [45] K. Winstein, A. Sivaraman, and H. Balakrishnan. Stochastic forecasts achieve high throughput and low delay over cellular networks. 2013.
- [46] Y. Wu, E. Mansimov, R. B. Grosse, S. Liao, and J. Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *Advances in neural information processing systems*, 2017.
- [47] X. Yin, A. Jindal, V. Sekar, and B. Sinopoli. A control-theoretic approach for dynamic adaptive video streaming over http. In *SIGCOMM’15*, 2015.
- [48] R.-X. Zhang, T. Huang, M. Ma, H. Pang, X. Yao, C. Wu, and L. Sun. Enhancing the crowd-sourced live streaming: a deep reinforcement learning approach. In *NOSSDAV’19*, 2019.

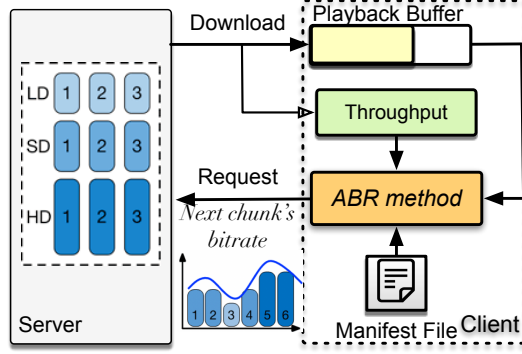


Figure 11: ABR System Overview.

## Appendix

### A Zwei Pseudo-code

#### A.1 Training Methodology.

Training procedure of Zwei is described in Algorithm 3.

---

#### Algorithm 2 Zwei’s Overall Training Procedure

---

**Require:** Environment Env; Discriminator Disc.

- 1: Initialize parameters  $\theta$  with random weights;
  - 2: **repeat**
  - 3:    $B = \{\}; D = \{\};$
  - 4:   Sample network trace  $T$  and video  $v$ ;
  - 5:   **for**  $i \in (1, n)$  **do** ▷ Repeat  $n$  times.;
  - 6:     Roll-out the trajectory  $T_i = \{s_0, a_0, s_1, \dots, a_t\}$  w.r.t policy  $\pi(\cdot; \theta)$  and  $\text{Env}(T, v)$ ;
  - 7:      $D \leftarrow D \cup T_i$
  - 8:   **end for**
  - 9:   **for**  $\text{pairs}(T_u, T_v) \in D$  **do**
  - 10:     Estimate win or loss:  $\hat{r}_u = \text{Disc}(T_u, T_v)$ ;
  - 11:     Store sample  $(S_t, a_t, \hat{r}_u)$  in  $B$ ;
  - 12:   **end for**
  - 13:   Update  $\theta$  with  $L^{\text{Zwei}}$  (Eq. 10) using  $B$ ;
  - 14: **until** Converged
- 

#### A.2 Implementation

We use TFlearn [43] to implement the NN and leverage TensorFlow [3] to construct Zwei. Zwei consists of two NNs, policy network and value network. Policy network takes a n-dims vector with *softmax* active function as the output. The value network outputs a value with *tanh* function scaled in  $(0, 1)$ . Considering the characteristics on video transmission tasks, we design different NN architectures for them. However, we use the same set of hyper-parameters for training the NN: we set sample number  $n = 10$ , learning rate  $\alpha = 10^{-4}$ , and the entropy weight  $\beta = 0.1$ . In addition, we adopt Adam with default settings.

### B Details

In this section, we briefly introduce the fundamental principle of Zwei’s simulation environment.

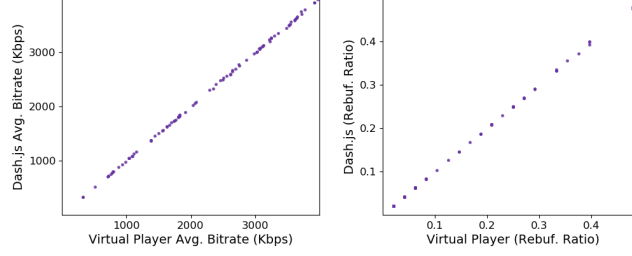


Figure 12: Comparing the strategy between the virtual player and real Dash.js player. Results show the close correlation of strategy between the two players.

## B.1 ABR Tasks

**Background** Due to the rapid development of network services, watching video streaming online has become an upcoming trend. Today, adaptive video streaming, such as HLS (HTTP Live Streaming) [2] and DASH [1], an algorithm that dynamically selects video bitrates via network conditions and client’s buffer occupancy, is the predominant form of video delivery [6]. The traditional video streaming architecture is shown in Figure 11. It consists of a video player client with a constrained buffer length and an HTTP-Server or Content Delivery Network (CDN). The video player client decodes and renders video frames from the playback buffer. Once the streaming service starts, the client fetches the video chunk from the HTTP Server or CDN orderly by an ABR algorithm, and, in the meanwhile, the ABR algorithm, implemented on the client side, determines the next chunk  $N$  and next chunk video quality  $Q_N$  via throughput estimation and current buffer utilization. After finished to play the video, several metrics, such as total bitrate  $b$ , total re-buffering time  $r$  and total bitrate change  $s$  will be summarized as a QoE metric to evaluate the performance. Thus, achieving a high QoE score for video streaming has become a major challenge for ABR algorithms.

**Virtual Player Evaluation** We design a faithful ABR offline virtual player to train Zwei with the given network traces and video descriptions. The player is written in Python3.6 and is closely refers to several state-of-the-art open-sourced ABR simulators including Pensieve, Oboe and Sabre [41]. Besides, We verify the accurateness of the proposed virtual player in Figure 12.

**Zwei’s NN Architecture for ABR Tasks** We now explain the details of the Zwei’s neural network (NN), including its inputs, outputs, network architecture, and implementation.

**> Inputs.** The NN’s inputs is mainly categorized into three parts, network features, video content features and video playback features ( $S_k = \{C_k, M_k, F_k\}$ ). We takes past chunk  $k = 8$  as the input. NN takes past  $t$  chunks’ network status vector  $C_k = \{c_{k-t-1}, \dots, c_k\}$  into NN, where  $c_i$  represents the throughput measured for video chunk  $i$ . Specifically,  $c_i$  is computed by  $c_i = n_{r,i}/d_i$ , in which  $n_{r,i}$  is the downloaded video size of chunk  $i$  with selected bitrates  $r$ , and  $d_i$  means download time for video chunk  $n_{r,i}$ . Besides, we also consider adding video content features into NN’s inputs for improving its abilities on detecting the diversity of video contents. In details, the learning agent leverages  $M_k = \{N_{k+1}\}$  to represent video content features. Here  $N_{k+1}$  is a vector that reflects the video size for each bitrate of the next chunk  $k + 1$ . The last essential feature for describing the ABR’s state is the current video playback status. The status is represented as  $F_k = \{v_{k-1}, B_k, D_k, m_k\}$ , where  $v_{k-1}$  is the perceptual video quality metric for the past video chunk selected,  $B_k, D_k$  are vectors which stand for past  $t$  chunks’ buffer occupancy and download time, and  $m_k$  means the normalized video chunk remaining.

**> Outputs.** We attempt to use discrete action space to describe the output. Note that the output is an  $n$ -dim vector indicating the probability of the bitrate being selected under the current ABR state  $S_k$ . In this work, we set  $n = 6$ , the bitrate ladder is defined as  $\{0.3, 0.75, 1.2, 1.8, 2.8, 4.3\}$  mbps, which is widely used in ABR papers [5, 14, 16, 23, 31].

**> NN Representation.** Zwei’s NN representation is quite simple: it leverages three fully-connected layers, which is sized 128, 64, and 64 respectively, for describing the feature extraction layer. On the one side, the output of the NN’s policy network is a 6-dims vector, which represents the probabilities for each bitrate selected. We utilize *ReLU* as the active function for each feature extraction layer and leverage *softmax* for the last layer. On the other size, NN’s value network outputs a scalar, which uses *tanh* as the active function.

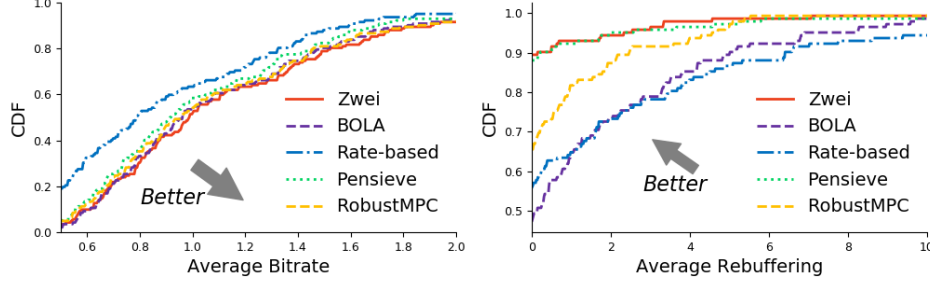


Figure 13: Here is the CDF distribution of average bitrate and rebuffering time for the selected approaches.

**Additional Experiments** We list the experimental details under validation set from Table 3 to Table 5, where the ABR algorithm includes Zwei, BOLA, RobustMPC, Pensieve and Rate-based. The validation set contains several different network scenarios. E.g., *bus-1* represents network trace No.1 which collected from the bus. Besides, we also report the CDF distribution of two underlying metrics in Figure 13. Such results above prove the effectiveness of Zwei on ABR tasks.

## B.2 CLS Tasks.

**CLS NN Representation.** We implement Zwei in CLS as suggested by recent work [48]. More precisely speaking, for each CDN provider  $i$ , Zwei passes past 20 normalized workload and stalling ratio to a single CNN layer with filter=64, size=4, and stride=1. Then several output layers are merged into a hidden layer that uses 64 units to apply the *softmax* activation function. We model the action space as a heuristic way: each CDN provider has 3 choices instead: increase its configuration ratio by 1%, 5%, and 10% based on the previous stalling ratio. And Zwei will reduce the ratio for CDN that obtained worst previous ratio.

## B.3 RTC Scenarios.

**Training with Network Simulator** In this part, we aim to introduce the principle of the RTC network simulator. To train Zwei in the RTC scenario, we first consider training the neural network (NN) model in real-world network conditions, i.e., deploying the model on the edge server. The model will finally converge with the increased number of video session. However, it's hard to converge since online training should explore almost all network status. We then decide to train the model in simulated offline networks. Hence, we are facing a new challenge: *how to design a fast-forward network simulator which can precisely compute the latency with given saturated trace and sending rate?*

We adopt delay gradient instead of one-way delay to train Zwei. Delay gradient is defined as the difference on delay measurement on both side. [7] It is also proposed to measure the latency in un-synchronized clock environment. Delay gradient  $q(t_i)$  is computed as Eq. 11, which is actually used to describe the variety of the queuing delay. If the network is congested, its delay gradient will be greater than zero, vice versa. If the delay gradient equals to zero, we can only infer that network is not congested.

$$q(t_i) = \frac{1}{N} \sum_{i=1}^N [(t_i - t_{i-1}) - (T_i - T_{i-1})] \quad (11)$$

Where  $T_i$  is the time at which the  $i$ -th packet has been sent, and  $t_i$  is the time at which the  $i$ -th packet has been received, and  $N$  represents the packet count in a period.

Our simulator should emulate the process of the packets coming and leaving in different network conditions, and keep track of the timestamps, by which we are able to obtain the corresponding queuing delay gradient. Inspired by recent work [45], we utilize saturated network trace to generate queuing delay data. As shown in Figure 14, assuming the distribution of packets arrival and leave fits closely to the Poisson process [45], we use sending bitrate and bandwidth in saturated network traces as the arriving rate  $\lambda$  and leaving rate  $\mu$ , respectively.

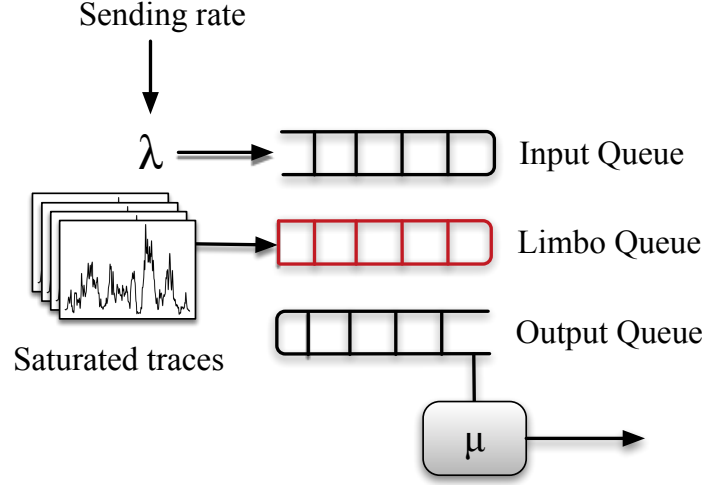


Figure 14: The working principle of the network simulator.

In detail, we regard the channel link between the sender and the receiver as a single device. The device is comprised of three queues, namely *input*, *output* and *limbo*. An algorithm is designed to release packets from each queue corresponding to the network trace dataset in different levels. Each element in the queue is the opportunity of packet-delivery, which means, the time (in milliseconds) at which an MTU-sized packet will be delivered. For each period  $(t_k, t_{k+1}]$ , the offline-network simulator compares the timestamp of the front packet in “input” with that in “limbo”. If timestamp in “input” queue is bigger than or equal to the one in “limbo”, the simulator will push the front packet into the bottom of the “output” queue, and then computes the delay gradient between the two packets. On the contrast, the delivery opportunities will be wasted if the timestamp in “input” queue is smaller than that one in “limbo” queue. The simulator only pops the front packet in “limbo” queue. By this process, the simulator returns mean self-inflicted delay and total bytes that the sum of “output” queue every interval  $t_{k+1} - t_k$ . We implement the network simulator via Python 3.6+.

**NN Architecture Overview** We use past  $k = 10$  sequence, and for each sequence we take past sending rate, past receive rate, delay gradient, round-trip time (RTT), loss ratio as the input. Besides, we output the NN as 11-dims vector  $a$  to control the sending rate, which represents  $\{-0.4, -0.3, -0.2, -0.1, -0.05, 0.0, 0.05, 0.1, 0.2, 0.3, 0.4\}$ . Inspired by AIMD algorithm, we take different logic to increase or decrease the sending rate. Zwei increases the sending rate  $r_t$  as  $r_t = r_t + r_{max} * a_t$ , and reduces the sending rate as  $r_t = r_t * (1 - a_t)$ . We set minimum sending rate  $r_{min} = 0.01$  and  $r_{max} = 0.5$ . The Zwei’s NN representation is equal to previous work [15]. We use the combination of 1D-CNN and Fully connected layers to extract the feature and then merge into a vector. Finally, the NN outputs the vector w.r.t ABR tasks.

**Additional Experiments** To better understand the performance of the selected methods, we show all the traces from the validation dataset and compare the sending rate curve of Zwei with TCP as well as WebRTC. Results are demonstrated from Figure 15 to Figure 53. Recall that we have not tried to pick the result which is good for Zwei.

## C Discussion

### C.1 Quality of Experience (QoE) for ABR

Recall that the goal of ABR algorithm is to select bitrates for the next chunk  $k$  with high bitrate  $r_k$  and less rebuffering time  $t_k$  [24], thus the optimization reward  $q_k$  can be normally written as

$$q_k = r_k - \alpha t_k. \quad (12)$$

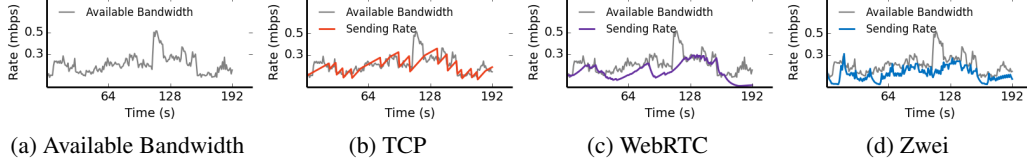


Figure 15: trace-5.

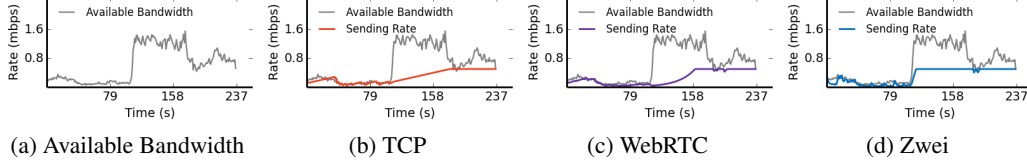


Figure 16: norway-bus-14.

Here  $\alpha$  is the coefficient that adjust the importance of the underlying metrics. In practice,  $\alpha$  is often defined as 3 [47], 4.3 [23], and 6.0 [47]. prior research add additional smoothness metric  $|r_k - r_{k-1}|$  to control the bitrate change of the entire session, while in practice the following metric is neglectable for the ABR algorithm [24]. Hence, in this work we also set the smoothness to zero for better understanding the fundamental performance of Zwei. We believe that Zwei can also solve  $q_k$  with smoothness metric.

## C.2 Discriminator Example

Algorithm 3 is an example of discriminator Disc which used in the ABR scenaio.

---

**Algorithm 3** Discriminator for the ABR task.

---

**Require:** Trajectory  $T_u, T_v$ ;

- 1: Compute undelying metric average bitrate  $r_u, r_v$ , average rebuffering  $b_u, b_v$  from the given trajectories  $T_u, T_v$ . ▷ Requirements: 1) low rebuffering time 2) high bitrate.
  - 2: Initalize Return  $s = \{0, 0\}$ ;
  - 3: **if**  $|b_u - b_v| < 0.01$  or  $|r_u - r_v| < 10$  **then**
  - 4:     Randomly set  $s_0$  or  $s_1$  as 1. ▷ Add noise to improve the robustness.
  - 5: **else if**  $|b_u - b_v| < 0.01$  **then**
  - 6:      $s_i \leftarrow 1, i \leftarrow \operatorname{argmax}_{i \in \{u, v\}} r$ ;
  - 7: **else**
  - 8:      $s_i \leftarrow 1, i \leftarrow \operatorname{argmin}_{i \in \{u, v\}} b$ ;
  - 9: **end if**
  - 10: return  $s$ ;
-



Table 3: Details of ABR Tasks

Trace Name	BOLA		Rate-based		Pensieve		RobustMPC		Zwei	
-	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.
bus 1	2659.6	-	2188.3	-	2480.9	-	2839.4	5.1	<b>2802.1</b>	-
bus 10	1467.0	-	1389.4	4.6	1413.8	-	1427.7	0.1	<b>1490.4</b>	-
bus 11	1327.7	-	1195.7	-	1264.9	-	1336.2	0.6	<b>1363.8</b>	-
bus 12	769.1	4.9	721.3	-	771.3	5.2	769.1	4.0	<b>793.6</b>	-
bus 13	852.1	9.7	<b>529.8</b>	-	623.4	2.6	629.8	5.5	666.0	0.7
bus 14	1675.5	8.3	933.0	-	1586.2	4.6	1610.6	4.6	<b>1638.3</b>	-
bus 15	3919.1	-	3497.9	-	3598.9	-	<b>4141.5</b>	-	3888.3	-
bus 16	3950.0	-	3652.1	-	3857.4	-	<b>4145.7</b>	-	3897.9	-
bus 17	1943.6	-	1844.7	2.2	1777.7	-	2005.3	-	<b>2022.3</b>	-
bus 18	1623.4	-	1559.6	0.4	1596.8	-	1588.3	1.0	<b>1643.6</b>	-
bus 19	1472.3	-	1212.8	-	1461.7	-	1440.4	-	<b>1517.0</b>	-
bus 2	1820.2	-	1656.4	-	1737.2	-	1861.7	-	<b>1919.1</b>	-
bus 20	1726.6	-	1517.0	1.6	1703.2	-	1752.1	-	<b>1795.7</b>	-
bus 21	1669.1	-	1497.9	-	1526.6	-	1759.6	2.1	<b>1730.9</b>	-
bus 22	<b>1743.6</b>	-	1594.7	7.8	1568.1	-	1791.5	3.6	1806.4	3.3
bus 23	1176.6	2.3	985.1	11.7	1076.6	-	1195.7	4.3	<b>1178.7</b>	-
bus 3	1331.9	-	1258.5	-	1283.0	-	1348.9	0.9	<b>1368.1</b>	-
bus 4	2121.3	-	1710.6	-	2103.2	-	2160.6	1.2	<b>2156.4</b>	-
bus 5	950.0	1.9	863.8	14.0	922.3	-	<b>964.9</b>	-	966.0	-
bus 6	1568.1	-	1324.5	0.3	1520.2	-	1606.4	-	<b>1680.9</b>	-
bus 7	<b>2669.1</b>	-	2376.6	-	2460.6	-	2571.3	-	2668.1	-
bus 8	2230.9	-	2008.5	2.9	2093.6	-	<b>2305.3</b>	-	2290.4	-
bus 9	2192.6	-	2014.9	4.0	2079.8	-	2297.9	1.8	<b>2283.0</b>	-
car 1	1440.4	0.8	1413.8	10.1	1467.0	12.7	1447.9	4.0	<b>1434.0</b>	-
car 10	1300.0	-	1127.7	-	1218.1	-	1338.3	2.6	<b>1350.0</b>	-
car 11	920.2	1.7	673.4	1.0	925.5	-	889.4	-	<b>969.1</b>	-
car 12	1291.5	18.2	<b>1283.0</b>	<b>6.4</b>	1291.5	13.7	1308.5	24.9	1325.5	13.5
car 2	1029.8	-	903.2	-	996.8	-	1034.0	-	<b>1060.6</b>	-
car 3	1351.1	1.8	1293.6	-	1334.0	2.0	1309.6	0.8	<b>1327.7</b>	-
car 4	1433.0	1.6	1327.7	3.4	1314.9	-	1423.4	1.7	<b>1476.6</b>	-
car 5	927.7	2.1	836.2	0.1	878.7	-	895.7	-	<b>959.6</b>	-
car 6	1295.7	-	1133.0	-	1263.8	-	1286.2	-	<b>1326.6</b>	-
car 7	710.6	1.3	568.1	11.9	666.0	-	677.7	0.9	<b>697.9</b>	-
car 8	1688.3	-	1448.9	-	1676.6	-	1704.3	-	<b>1748.9</b>	-
car 9	1485.1	-	1366.0	-	1423.4	-	1459.6	-	<b>1526.6</b>	-
ferry 1	1454.3	0.3	1423.4	-	1441.5	-	1455.3	-	<b>1492.6</b>	-
ferry 10	1119.1	0.2	1021.3	0.2	1104.3	0.2	1104.3	0.2	<b>1133.0</b>	<b>0.2</b>
ferry 11	1038.3	1.2	847.9	-	978.7	-	1037.2	0.1	<b>1067.0</b>	-
ferry 12	1525.5	-	1324.5	6.4	1466.0	-	1555.3	-	<b>1586.2</b>	-
ferry 13	1421.3	0.5	1192.6	12.6	<b>1320.2</b>	-	1436.2	5.4	1471.3	1.8
ferry 14	2036.2	-	1850.0	0.3	2074.5	-	2098.9	-	<b>2124.5</b>	-
ferry 15	1551.1	-	1385.1	-	1535.1	-	1486.2	1.1	<b>1573.4</b>	-
ferry 16	2684.0	-	2406.4	-	2480.9	-	2734.0	-	<b>2817.0</b>	-
ferry 17	859.6	0.1	768.1	9.4	818.1	-	826.6	-	<b>850.0</b>	-
ferry 18	724.5	2.8	587.2	4.9	666.0	1.1	667.0	1.0	<b>684.0</b>	-
ferry 19	839.4	4.8	654.3	-	<b>826.6</b>	-	825.5	-	863.8	7.1
ferry 2	1012.8	3.3	951.1	10.6	990.4	-	1001.1	-	<b>1024.5</b>	-

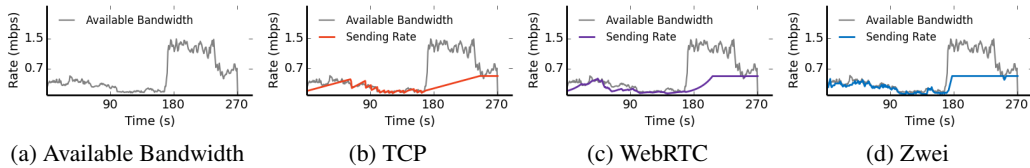


Figure 17: norway-bus-13.

Table 4: Details of ABR Tasks

Trace Name	BOLA		Rate-based		Pensieve		RobustMPC		Zwei	
Trace Name	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.
ferry 20	686.2	6.7	635.1	-	670.2	5.8	672.3	5.0	<b>674.5</b>	-
ferry 3	1211.7	-	1093.6	-	1170.2	-	1231.9	0.7	<b>1252.1</b>	-
ferry 4	562.8	4.2	472.3	1.9	543.6	-	558.5	4.7	<b>563.8</b>	-
ferry 5	1011.7	1.7	901.1	12.0	940.4	-	1018.1	2.2	<b>1027.7</b>	-
ferry 6	2014.9	-	1598.9	-	2039.4	-	2093.6	-	<b>2119.1</b>	-
ferry 7	1367.0	0.5	1231.9	1.4	1309.6	-	1346.8	0.6	<b>1422.3</b>	-
ferry 8	1078.7	3.8	1010.6	-	1076.6	1.7	1078.7	2.5	<b>1105.3</b>	-
ferry 9	964.9	3.8	787.2	5.1	868.1	3.7	<b>898.9</b>	<b>1.9</b>	1004.3	3.3
metro 1	1350.0	-	1141.5	-	1269.1	-	1371.3	-	<b>1389.4</b>	-
metro 10	896.8	-	625.5	3.7	837.2	-	895.7	-	<b>935.1</b>	-
metro 2	1364.9	-	1285.1	2.8	1316.0	-	1381.9	2.5	<b>1409.6</b>	-
metro 3	934.0	-	481.9	-	884.0	-	913.8	-	<b>956.4</b>	-
metro 4	1086.2	1.0	884.0	1.2	968.1	0.6	1030.9	-	<b>1121.3</b>	-
metro 5	1178.7	1.0	1041.5	3.8	1142.6	-	<b>1171.3</b>	-	1229.8	0.5
metro 6	577.7	0.9	386.2	6.6	558.5	-	548.9	-	<b>573.4</b>	-
metro 7	945.7	-	769.1	-	948.9	0.2	961.7	0.1	<b>977.7</b>	-
metro 8	740.4	2.5	443.6	6.4	738.3	-	724.5	-	<b>762.8</b>	-
metro 9	788.3	-	702.1	4.4	735.1	-	788.3	1.5	<b>806.4</b>	-
train 1	1028.7	4.2	863.8	2.3	996.8	-	1008.5	-	<b>1024.5</b>	-
train 10	1838.3	-	1627.7	-	1709.6	-	1894.7	-	<b>1907.4</b>	-
train 11	1256.4	0.4	1163.8	0.2	1171.3	-	1257.4	-	<b>1286.2</b>	-
train 12	1304.3	0.8	911.7	1.3	1114.9	-	1212.8	-	<b>1380.9</b>	-
train 13	1560.6	-	1236.2	-	1334.0	-	1564.9	-	<b>1675.5</b>	-
train 14	1330.9	-	1206.4	-	1259.6	-	1320.2	-	<b>1356.4</b>	-
train 15	1900.0	-	1423.4	-	1650.0	-	<b>1952.1</b>	-	1944.7	-
train 16	1661.7	-	1475.5	3.5	1484.0	-	1575.5	-	<b>1736.2</b>	-
train 17	1620.2	-	1394.7	4.1	1556.4	-	1637.2	-	<b>1706.4</b>	-
train 18	1978.7	-	1725.5	-	1803.2	-	2012.8	-	<b>2091.5</b>	-
train 19	1423.4	-	1230.9	2.5	1302.1	-	1452.1	-	<b>1484.0</b>	-
train 2	635.1	3.5	481.9	11.1	<b>577.7</b>	-	581.9	-	664.9	3.0
train 20	1331.9	-	1052.1	-	1343.6	-	1357.4	-	<b>1373.4</b>	-
train 21	1598.9	-	1421.3	-	1596.8	-	1626.6	-	<b>1652.1</b>	-
train 3	648.9	1.6	529.8	8.6	619.1	-	596.8	-	<b>629.8</b>	-
train 4	1084.0	3.3	868.1	-	1024.5	-	1059.6	0.3	<b>1127.7</b>	-
train 5	1241.5	1.5	946.8	0.4	1155.3	-	1151.1	-	<b>1283.0</b>	-
train 6	969.1	0.9	721.3	-	875.5	-	936.2	0.9	<b>991.5</b>	-
train 7	825.5	8.3	750.0	-	<b>784.0</b>	-	790.4	-	823.4	2.6
train 8	706.4	1.0	510.6	-	666.0	-	693.6	-	<b>730.9</b>	-
train 9	855.3	-	596.8	-	776.6	-	852.1	-	<b>881.9</b>	-
tram 1	635.1	-	472.3	1.6	625.5	-	635.1	-	<b>656.4</b>	-
tram 10	869.1	2.5	673.4	-	850.0	-	<b>878.7</b>	-	871.3	-
tram 11	587.2	4.3	472.3	1.3	558.5	-	548.9	-	<b>580.9</b>	-
tram 12	807.4	-	548.9	-	795.7	-	816.0	-	<b>826.6</b>	-
tram 13	788.3	-	596.8	0.9	739.4	-	787.2	-	<b>831.9</b>	-
tram 14	839.4	2.2	644.7	1.1	808.5	-	797.9	-	<b>844.7</b>	-
tram 15	453.2	6.9	357.4	5.3	414.9	1.8	424.5	-	<b>437.2</b>	-
tram 16	<b>683.0</b>	-	548.9	-	673.4	-	692.6	-	692.6	-
tram 17	<b>654.3</b>	-	424.5	-	625.5	-	635.1	-	663.8	-
tram 18	730.9	2.8	587.2	-	691.5	-	706.4	-	<b>751.1</b>	-

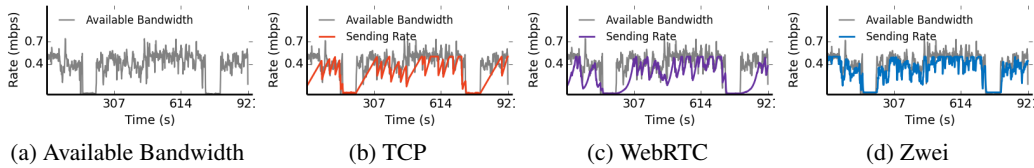
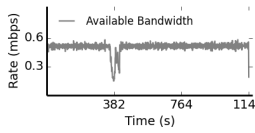


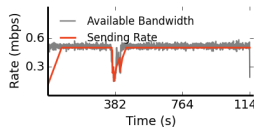
Figure 18: 12318-wifi5.

Table 5: Details of ABR Tasks

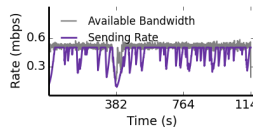
Trace Name	BOLA		Rate-based		Pensieve		RobustMPC		Zwei	
Trace Name	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.	Avg. Bit.	Rebuf.
tram 19	967.0	0.1	788.3	-	907.4	-	973.4	-	<b>985.1</b>	-
tram 2	577.7	3.3	434.0	-	548.9	-	548.9	-	<b>559.6</b>	-
tram 20	930.9	-	730.9	-	885.1	-	928.7	-	<b>988.3</b>	-
tram 21	692.6	0.5	558.5	-	677.7	-	696.8	-	<b>711.7</b>	-
tram 22	510.6	0.3	376.6	-	491.5	-	501.1	-	<b>506.4</b>	-
tram 23	778.7	-	529.8	-	738.3	-	778.7	-	<b>808.5</b>	-
tram 24	788.3	-	616.0	3.4	739.4	-	769.1	0.3	<b>808.5</b>	-
tram 25	907.4	-	711.7	-	813.8	-	888.3	-	<b>950.0</b>	-
tram 26	568.1	0.8	300.0	-	539.4	-	<b>568.1</b>	-	567.0	-
tram 27	<b>692.6</b>	-	491.5	-	673.4	-	692.6	-	698.9	-
tram 28	558.5	0.3	386.2	2.2	510.6	-	<b>539.4</b>	-	540.4	-
tram 29	510.6	3.4	338.3	-	481.9	-	<b>501.1</b>	-	495.7	-
tram 3	577.7	1.3	414.9	-	<b>539.4</b>	-	534.0	1.9	570.2	1.7
tram 30	692.6	-	319.1	-	644.7	-	687.2	-	<b>719.1</b>	-
tram 31	769.1	-	644.7	-	753.2	-	769.1	-	<b>783.0</b>	-
tram 32	<b>788.3</b>	-	558.5	1.7	733.0	-	750.0	-	788.3	-
tram 33	<b>874.5</b>	-	644.7	-	841.5	-	877.7	0.2	880.9	-
tram 34	926.6	-	778.7	1.6	917.0	-	930.9	-	<b>959.6</b>	-
tram 35	548.9	5.0	405.3	1.6	501.1	-	486.2	-	<b>550.0</b>	-
tram 36	850.0	0.3	539.4	-	819.1	-	816.0	-	<b>874.5</b>	-
tram 37	984.0	-	807.4	-	894.7	-	1003.2	-	<b>1026.6</b>	-
tram 38	644.7	9.6	481.9	4.5	<b>591.5</b>	-	625.5	4.9	646.8	2.3
tram 39	855.3	-	721.3	-	814.9	-	826.6	-	<b>868.1</b>	-
tram 4	462.8	5.2	300.0	-	414.9	-	414.9	-	<b>428.7</b>	-
tram 40	797.9	0.6	520.2	0.6	770.2	0.6	772.3	0.8	<b>827.7</b>	<b>0.6</b>
tram 41	769.1	2.8	663.8	-	734.0	-	750.0	0.5	<b>785.1</b>	-
tram 42	1069.1	9.0	692.6	-	961.7	-	1057.4	-	<b>1076.6</b>	-
tram 43	1079.8	-	960.6	0.2	1048.9	-	1084.0	-	<b>1105.3</b>	-
tram 44	869.1	-	759.6	-	801.1	-	874.5	0.7	<b>893.6</b>	-
tram 45	648.9	1.1	472.3	7.0	<b>627.7</b>	<b>0.5</b>	635.1	0.5	627.7	0.5
tram 46	683.0	4.2	414.9	-	606.4	-	662.8	0.1	<b>679.8</b>	-
tram 47	788.3	0.9	539.4	-	705.3	-	764.9	-	<b>795.7</b>	-
tram 48	616.0	0.3	395.7	-	558.5	-	605.3	0.2	<b>619.1</b>	-
tram 49	963.8	1.1	778.7	-	963.8	-	985.1	2.2	<b>1007.4</b>	-
tram 5	654.3	5.2	443.6	2.6	587.2	-	616.0	-	<b>640.4</b>	-
tram 50	863.8	6.8	797.9	-	851.1	-	845.7	-	<b>864.9</b>	-
tram 51	658.5	5.5	587.2	-	<b>638.3</b>	-	647.9	-	674.5	4.6
tram 52	920.2	3.4	788.3	1.7	905.3	-	903.2	-	<b>927.7</b>	-
tram 53	657.4	17.3	443.6	3.8	434.0	0.2	428.7	-	<b>522.3</b>	-
tram 54	1123.4	6.8	862.8	6.6	1059.6	-	<b>1109.6</b>	-	1106.4	-
tram 55	1176.6	0.4	1100.0	-	1141.5	-	1146.8	1.9	<b>1177.7</b>	-
tram 56	971.3	-	797.9	-	944.7	-	957.4	-	<b>992.6</b>	-
tram 6	740.4	-	587.2	-	683.0	-	744.7	-	<b>784.0</b>	-
tram 7	759.6	-	596.8	-	720.2	-	744.7	-	<b>797.9</b>	-
tram 8	501.1	3.5	414.9	-	472.3	-	462.8	0.2	<b>483.0</b>	-
tram 9	606.4	0.2	328.7	-	568.1	-	572.3	-	<b>629.8</b>	-



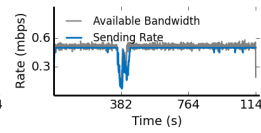
(a) Available Bandwidth



(b) TCP



(c) WebRTC



(d) Zwei

Figure 19: 12314-wifi2.

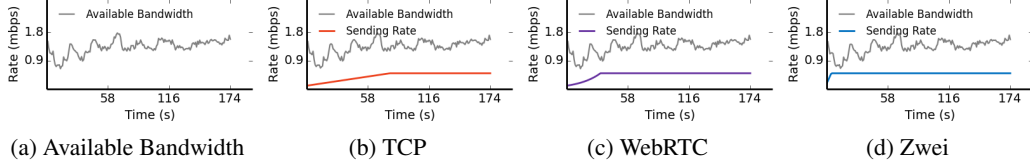


Figure 20: trace-1.

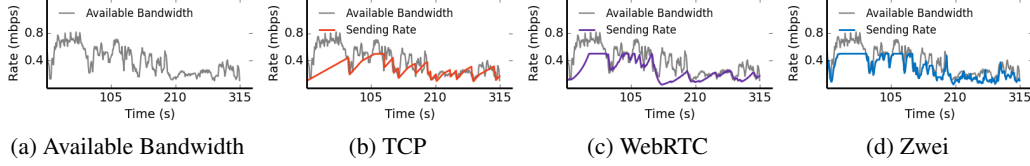


Figure 21: norway-metro-2.

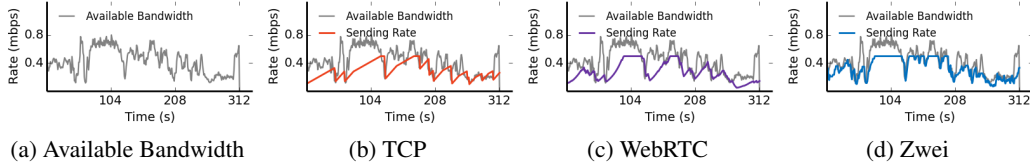


Figure 22: norway-metro-1.

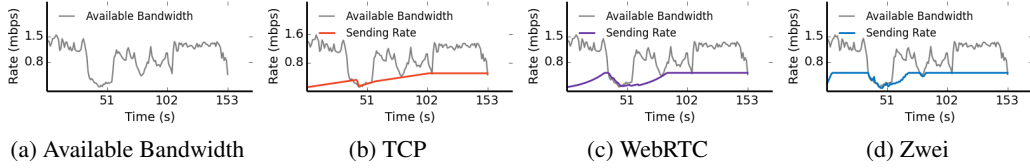


Figure 23: norway-bus-1.

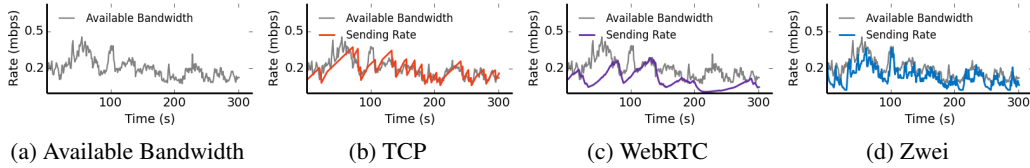


Figure 24: norway-tram-1.

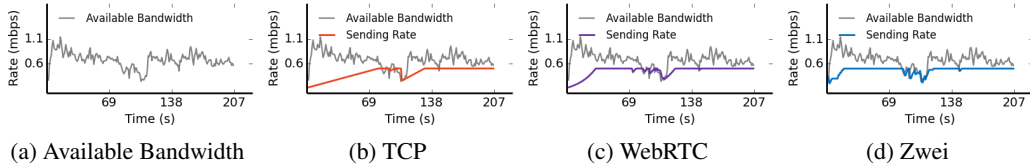


Figure 25: norway-bus-2.

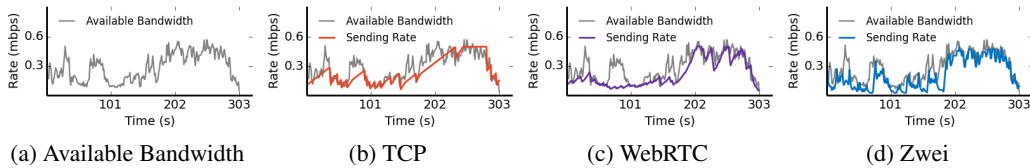


Figure 26: norway-tram-3.

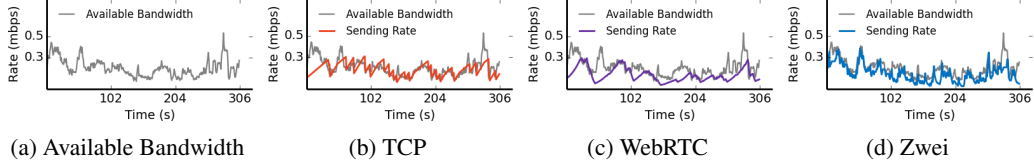


Figure 27: norway-tram-2.

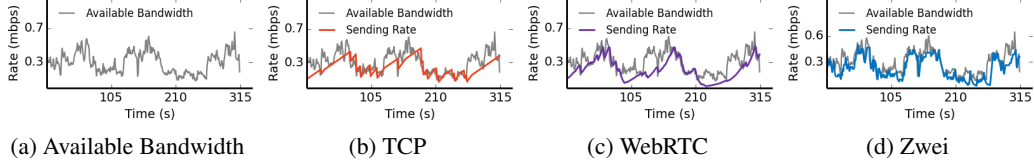


Figure 28: norway-metro-10.

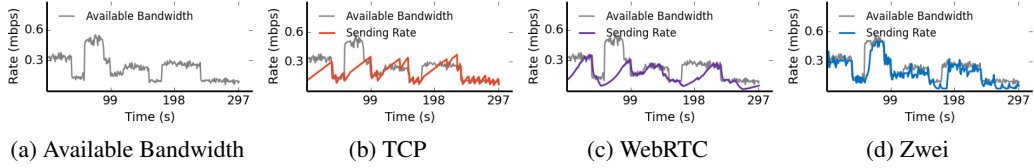


Figure 29: test-1.

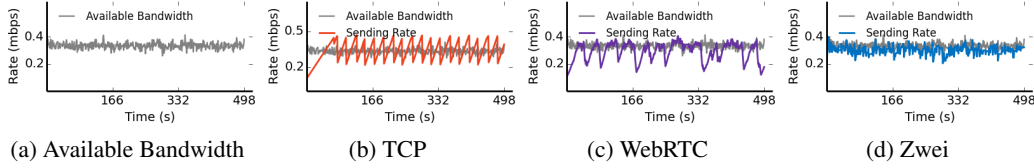


Figure 30: test-0.

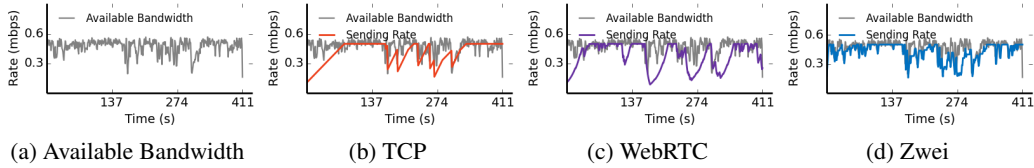


Figure 31: 12320-wifi5.

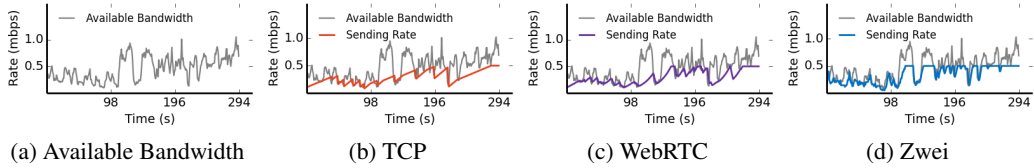


Figure 32: norway-ferry-11.

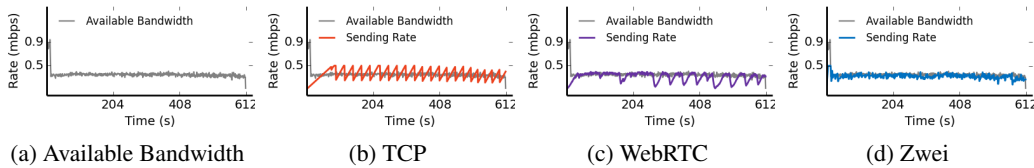


Figure 33: 12315-wifi3.

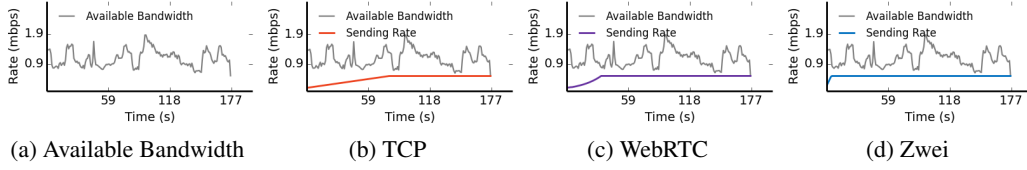


Figure 34: trace-9.

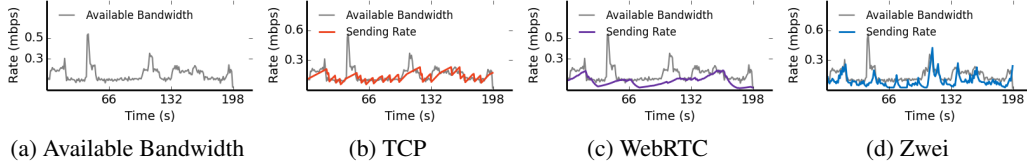


Figure 35: trace-8.

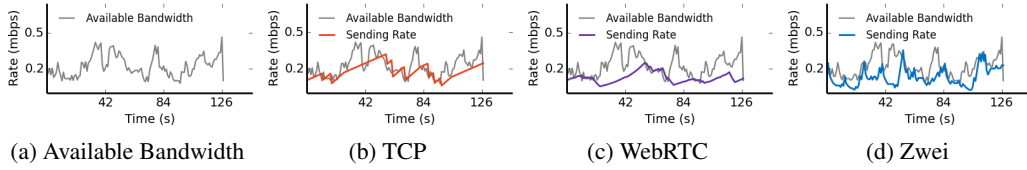


Figure 36: norway-train-2.

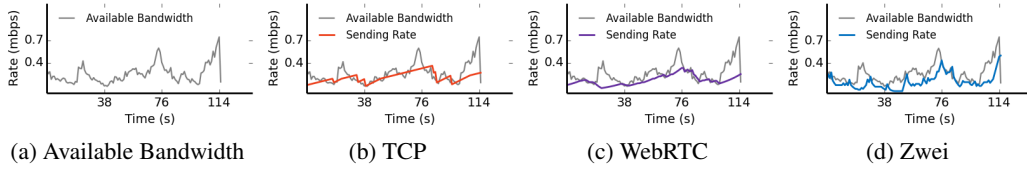


Figure 37: norway-train-3.

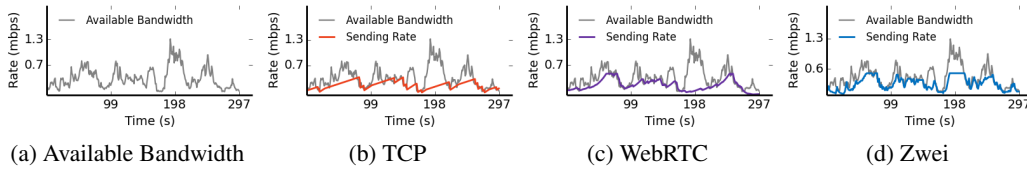


Figure 38: norway-train-1.

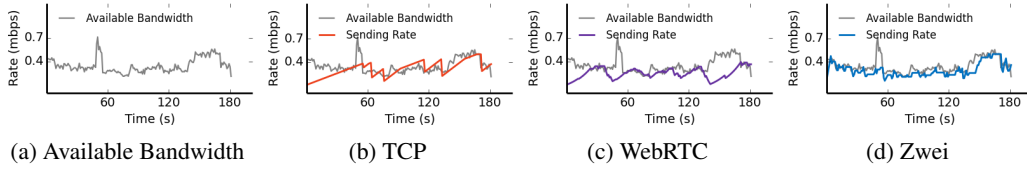


Figure 39: trace-3.

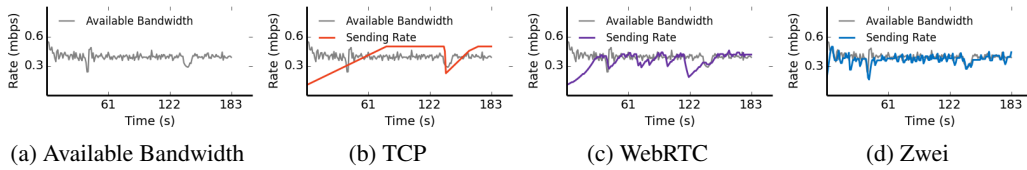


Figure 40: trace-2.



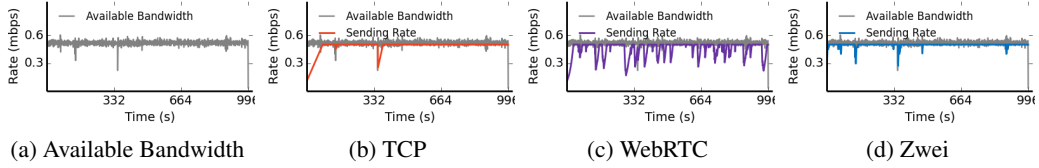


Figure 41: 12338-wifi6.

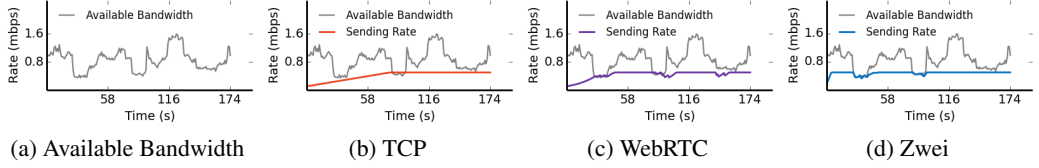


Figure 42: trace-0.

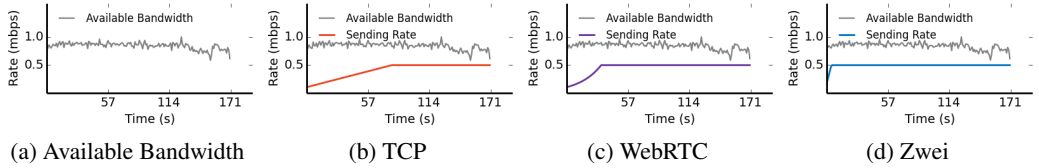


Figure 43: trace-7.

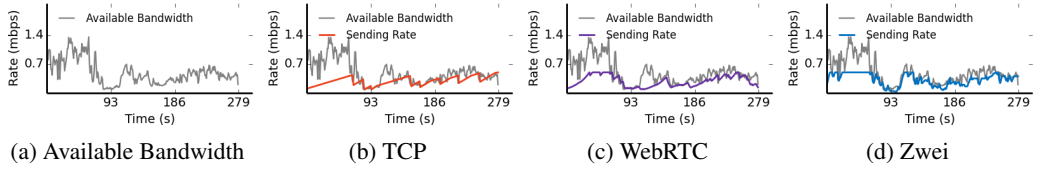


Figure 44: norway-car-1.

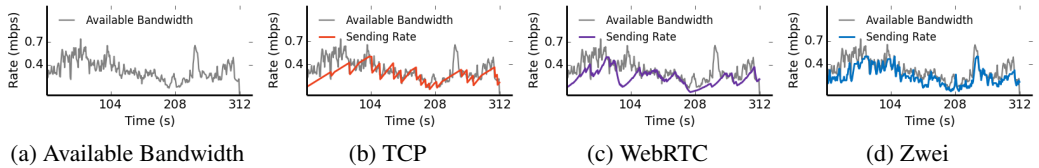


Figure 45: norway-car-2.

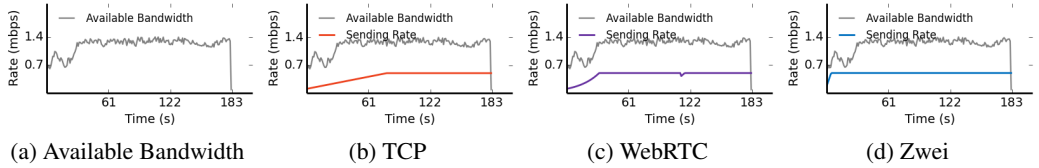


Figure 46: trace-4.

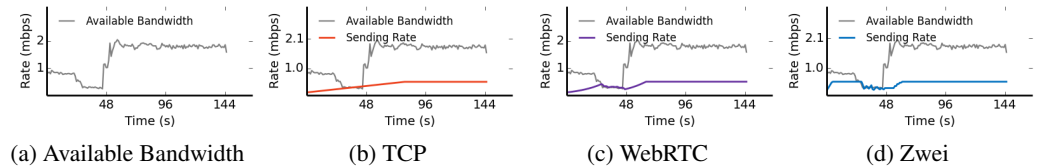


Figure 47: trace-6.

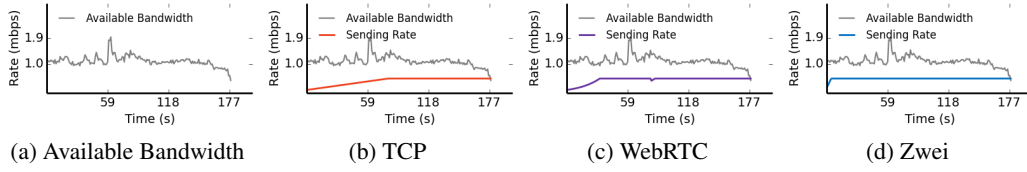


Figure 48: trace-14.

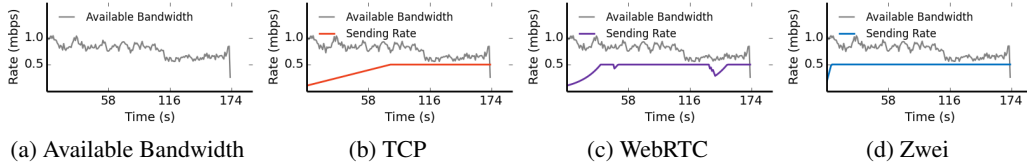


Figure 49: trace-13.

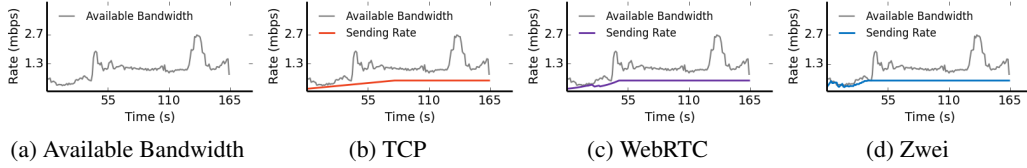


Figure 50: trace-12.

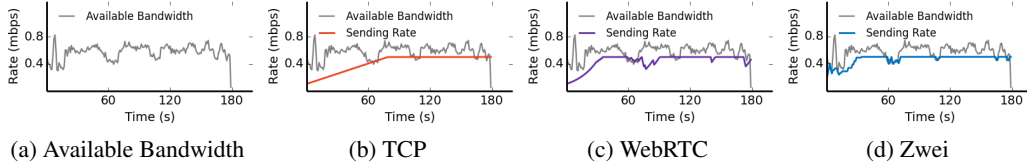


Figure 51: trace-11.

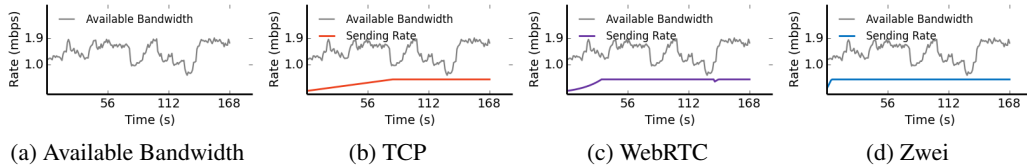


Figure 52: trace-10.

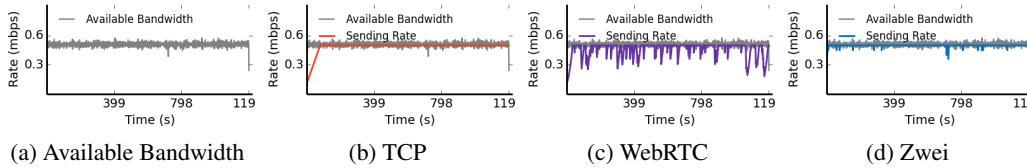


Figure 53: 12314-4G2.