

Requet: Real-Time QoE Metric Detection for Encrypted YouTube Traffic

CRAIG GUTTERMAN, Columbia University

KATHERINE GUO, Nokia Bell Labs

SARTHAK ARORA and TREY GILLILAND, Columbia University

XIAOYANG WANG and LES WU, Nokia Bell Labs

ETHAN KATZ-BASSETT and GIL ZUSSMAN, Columbia University

As video traffic dominates the Internet, it is important for operators to detect video quality of experience (QoE) to ensure adequate support for video traffic. With wide deployment of end-to-end encryption, traditional deep packet inspection-based traffic monitoring approaches are becoming ineffective. This poses a challenge for network operators to monitor user QoE and improve upon their experience. To resolve this issue, we develop and present a system for *REal-time QUality* of experience metric detection for *Encrypted Traffic—Requet*—which is suitable for network middlebox deployment. Requet uses a detection algorithm that we develop to identify video and audio chunks from the IP headers of encrypted traffic. Features extracted from the chunk statistics are used as input to a machine learning algorithm to predict QoE metrics, specifically buffer warning (low buffer, high buffer), video state (buffer increase, buffer decay, steady, stall), and video resolution. We collect a large YouTube dataset consisting of diverse video assets delivered over various WiFi and LTE network conditions to evaluate the performance. We compare Requet with a baseline system based on previous work and show that Requet outperforms the baseline system in accuracy of predicting buffer low warning, video state, and video resolution by 1.12×, 1.53×, and 3.14×, respectively.

CCS Concepts: • **Information systems** → **Multimedia streaming**; • **Networks** → *Network performance analysis*; • **Computing methodologies** → *Classification and regression trees*; • **Networks** → *Network monitoring*;

Additional Key Words and Phrases: Machine learning, HTTP adaptive streaming

ACM Reference format:

Craig Gutterman, Katherine Guo, Sarthak Arora, Trey Gilliland, Xiaoyang Wang, Les Wu, Ethan Katz-Bassett, and Gil Zussman. 2020. Requet: Real-Time QoE Metric Detection for Encrypted YouTube Traffic. *ACM Trans. Multimedia Comput. Commun. Appl.* 16, 2s, Article 71 (July 2020), 28 pages.

<https://doi.org/10.1145/3394498>

A preliminary version of this article was presented at ACM MMSys 2019 [22].

This work was supported in part by NSF grants CNS-1650685, CNS-1910757, and DGE 16-44869.

Authors' addresses: C. Gutterman, E. Katz-Bassett, and G. Zussman, Department of Electrical Engineering, Columbia University, 500 W 120 St, Room 1300, New York, NY 10027; emails: clg2168@columbia.edu, ethan@ee.columbia.edu, gil.zussman@columbia.edu; K. Guo, X. Wang, and L. Wu, Nokia Bell Labs, 600 Mountain Ave bldg 5, New Providence, NJ 07974; emails: {katherine.guo, xiaoyang.wang, les.j.wu}@nokia-bell-labs.com; S. Arora and T. Gilliland, Department of Computer Science, Columbia University, 500 W 120 St, Room 1300, New York, NY 10027; emails: {sa3522, jlg2266}@columbia.edu.

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 ACM 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.

© 2020 Association for Computing Machinery.

1551-6857/2020/07-ART71 \$15.00

<https://doi.org/10.1145/3394498>

1 INTRODUCTION

Video has monopolized Internet traffic in recent years. Specifically, the portion of video over mobile data traffic is expected to increase from 60% in 2016 to 78% by 2021 [2]. Content providers, content delivery networks (CDNs), and network operators are all stakeholders in the Internet video sector. They want to monitor user video quality of experience (QoE) and improve upon it to ensure user engagement. Content providers and CDNs can measure client QoE metrics such as video resolution by using server-side logs [8, 21]. Client-side measurement applications can accurately report QoE metrics such as player events and video quality levels [36, 50].

Traditionally, deep packet inspection (DPI) enabled operators to examine HTTP packet flows and extract video session information to infer QoE metrics [7, 12]. However, to address security and privacy concerns, content providers are increasingly adopting end-to-end encryption. A majority of YouTube traffic has been encrypted since 2016 [4] with a combination of HTTPS (HTTP/TLS/TCP) [9, 18, 40] and QUIC (HTTP/QUIC/UDP) [15, 25]. Similarly, since 2015, Netflix has been deploying HTTPS for video traffic [10]. In general, the share of encrypted traffic was estimated to be greater than 80% in 2019 [5].

Although the trend of end-to-end encryption does not affect client-side or server-side QoE monitoring, it renders traditional DPI-based video QoE monitoring ineffective for operators. Encrypted traffic still allows for viewing packet headers in plain text. This has led to recent efforts to use machine learning (ML) and statistical analysis to derive QoE metrics for operators. These works either provide offline analysis for the entire video session [16, 37] or online analysis using both network- and transport-layer information with separate models for HTTPS and QUIC [35].

Previous research developed methods to derive network-layer features from IP headers by capturing packet behavior in both directions: *uplink* (from the client to the server) and *downlink* (from the server to the client) [26, 35, 37]. However, determining QoE purely based on IP header information is inaccurate. To illustrate, Figure 1 shows a 20-second portion from two example sessions from our YouTube dataset, described Section 4, where each data point represents 100 ms. Both examples exhibit similar patterns in the downlink direction, whereas in the uplink direction, traffic spikes are much higher in Figure 1(b) than in Figure 1(a). However, Figure 1(a) shows a 720p resolution with the buffer decreasing by 15 seconds, whereas Figure 1(b) shows a 144p resolution with the buffer increasing by 20 seconds.

Given this challenge, our objective is to design features from IP header information that utilize patterns in the video streaming algorithm. In general, video clips stored on the server are divided into a number of *segments* or *chunks* at multiple resolutions. The client requests each chunk by individually sending an HTTP GET request to the server. Existing work using chunks either infers QoE for the entire session [31] rather than in real time or lacks insight on chunk detection mechanisms from network- or transport-layer data [16, 28, 42].

To improve on existing approaches that use chunks, we develop *Requet*, a system for *REal-time QUality* of experience metric detection for *Encrypted Traffic* designed for traffic monitoring in middleboxes by operators. Requet is devised for real-time QoE metric identification as chunks are delivered rather than at the end of a video session. Requet is designed to be memory efficient for middleboxes, where the memory requirement is a key consideration. Figure 2 depicts the system diagram for Requet and necessary components to train the QoE models, as well as evaluate its performance. Requet consists of the ChunkDetection algorithm, chunk feature extraction, and ML QoE prediction models. The data acquisition process involves collecting YouTube traffic traces (*trace collection*) and generating ground truth QoE metrics as labels directly from the player (*video labeling*). Packet traces are fed into Requet's ChunkDetection algorithm to determine audio and video chunks. The chunks are then used during the *feature extraction* process to obtain

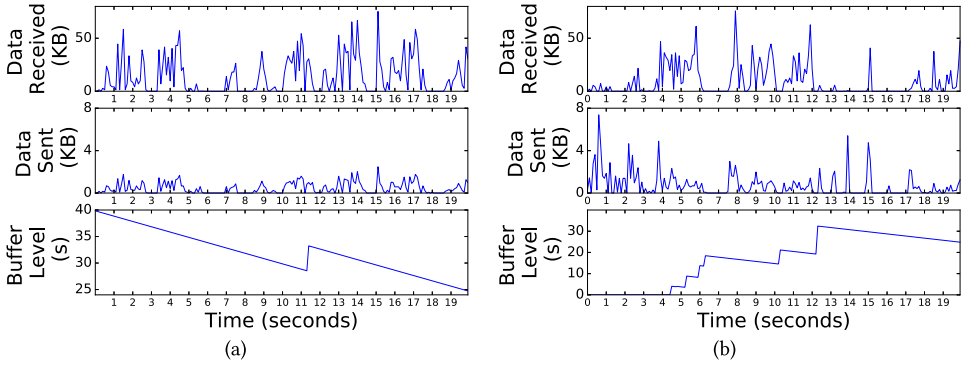


Fig. 1. Amount of data received (KB), amount of data sent (KB), and buffer level (seconds) for two sessions over a 20-second window (100-ms granularity): 720p (a) and 144p (b).

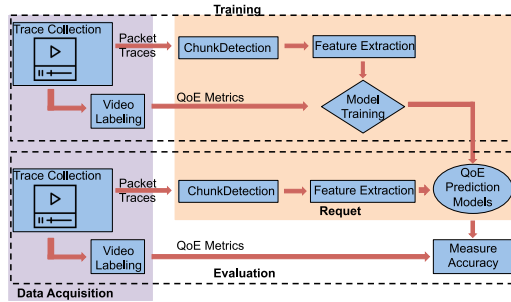


Fig. 2. System diagram. Data acquisition and Requet components: ChunkDetection, feature extraction, and QoE prediction models.

chunk-based features. The chunk-based features from the training data, along with the corresponding QoE metrics, are used to generate QoE prediction models. For evaluation, traffic traces from the testing dataset are fed into the trained QoE models to generate predicted QoE metrics. Accuracy is measured by comparing the predicted QoE metrics and the ground truth labels.

Recent studies have shown that (i) stall events have the largest negative impact on end user engagement, and (ii) higher average video *playback bitrate* improves user engagement [8, 17]. Motivated by these findings, Requet aims to predict the current video resolution and events that lead to QoE impairment ahead of time. This allows operators to proactively provision resources [13, 39]. Requet predicts a low buffer level, which allows operators to provision network resources to avoid stall events [26]. Requet predicts four video states: buffer increase, buffer decay, steady, and stall. Furthermore, Requet predicts current video resolution during a video session in real time. Specifically, Requet predicts video resolution on a more granular scale (144p, 240p, 360p, 480p, 720p, 1080p), whereas previous work predicts only two or three levels of video resolution for the entire video session [16, 31, 35].

We make the following contributions:

- Collect packet traces of 60 diverse YouTube video clips resulting in a mixture of HTTP/TLS/TCP and HTTP/QUIC/UDP traffic. The traces are collected in two distinct settings, with the first set collected from a laptop web browser over WiFi networks from three service providers (two in the United States and one in India) and the second set collected

from the YouTube App on an Android mobile device over LTE cellular networks. This is in contrast to most prior works, which rely on simulation or emulation [26, 35, 46] (Section 4).

- Design Requet components:
 - Develop ChunkDetection, a heuristic algorithm to identify video and audio chunks from IP headers (Section 3).
 - Analyze the correlation between audio and video chunk metrics (e.g., chunk size, duration, and download time) and various QoE metrics, and determine fundamental chunk-based features useful for QoE prediction. Specifically, design features based on our observation that audio chunk arrival rate correlates with the video state (Section 5).
 - Develop ML models to predict QoE metrics in real time: buffer warning, video state, and video resolution (Section 6).
- EvaluateRequet performance:
 - Demonstrate drastically improved prediction accuracy using chunk-based features versus baseline IP layer features commonly used in prior work [26, 35, 37, 47]. For the setting of a web browser over WiFi networks, Requet predicts low buffer warning with 92% accuracy, video state with 84% accuracy, and video resolution with 66% accuracy, representing an improvement of 1.12×, 1.53×, and 3.14×, respectively, over the existing baseline system. Furthermore, Requet delivers a 91% accuracy in predicting low (144p/240p/360p) or high resolution (480p/720p/1080p) in both the web browser over WiFi setting and the YouTube app over the LTE setting (Section 6).
 - Demonstrate that Requet trained in a laboratory environment works on unseen clips with varying lengths from different operators in multiple countries. This evaluation is more diverse than prior work [16, 26, 35, 46] (Section 6).

2 BACKGROUND AND PROBLEM STATEMENT

2.1 Adaptive Bitrate Streaming Operation

A majority of video traffic over the Internet today is delivered using HTTP adaptive streaming (HAS), with its dominating format being dynamic adaptive streaming over HTTP (DASH) or MPEG-DASH [45, 52]. In adaptive bitrate (ABR), a video *asset* or *clip* is encoded in multiple resolutions. Encoding is controlled by multiple parameters, and a given resolution is associated with a fixed setting for quantization, which is then coarsely related to an average bandwidth determined by the source video. A clip with a given resolution is then divided into a number of *segments* or *chunks* of variable length, a few seconds in playback time [33]. Typically, video clips are encoded with variable bitrate (VBR) encoding and are restricted by a maximum bitrate for each resolution. An audio file or the audio track of a clip is usually encoded with constant bitrate (CBR). For example some of the YouTube audio bitrates are 64, 128, and 192 Kbps [49].

At the start of the session, the client retrieves a manifest file that describes the location of chunks within the file containing the clip encoded with a given resolution. There are many ABR variations across and even within video providers [33]. ABR is delivered over HTTP(S), which requires either TCP or any other reliable transport [19]. The ABR algorithm can use concurrent TCP or QUIC/UDP flows to deliver multiple chunks simultaneously. A chunk can either be video or audio alone or a mixture of both.

2.2 Video States and Playback Regions

The client employs a *playout buffer* or *client buffer*, whose maximum value is *buffer capacity*, to temporarily store chunks to absorb network variation. To ensure smooth playback and an adequate *buffer level*, the client requests a video clip chunk by chunk using HTTP GET requests, and

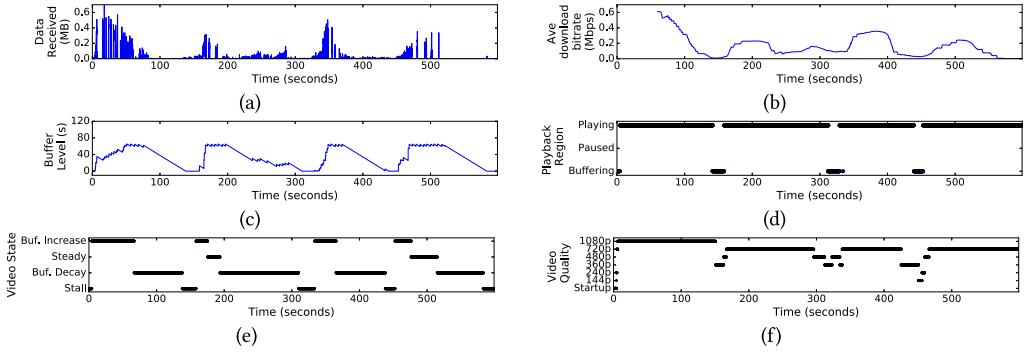


Fig. 3. Behavior of a 10-minute session in 100-ms windows: amount of data received (MB) (a), average download bitrate (Mbps) over the past 60 seconds (b), buffer level (c), playback region (d), video state (e), video resolution (f).

dynamically determines the resolution of the next chunk based on network condition and buffer status.¹

When the buffer level is below a low threshold, the client requests chunks as fast as the network can deliver them to increase the buffer level. We call this portion of ABR operation the *buffer filling* stage. In this stage, the buffer level can increase or decrease. Once the buffer level reaches a high threshold, the client aims to maintain the buffer level in the range between the threshold and buffer capacity. One example of a client strategy is to request chunks as fast as they are consumed by the playback process, which is indicated by the video playback bitrate for the chunk [47]. We call this portion the *steady state* stage. The playback stalls when the buffer is empty before the end of the playback is reached. After all chunks have been downloaded to the client buffer, there is no additional traffic and the buffer level decreases. From the perspective of the buffer level, an ABR session can experience four exclusive video states: *buffer increase*, *buffer decay*, *steady state*, and *stall*.

Orthogonally, from the perspective of YouTube video playback, a session can contain three exclusive regions: *buffering*, *playing*, and *paused*. The buffering region is defined as the period when the client is receiving data in its buffer but video playback has not started or is stopped. The playing region is defined as the period when video playback is advancing regardless of buffer status. The paused region is defined as the period when the end user issues the command to pause video playback before the session ends. In the playing region, video state can be in either buffer increase, decay, or steady state.

Figure 3 shows the behavior of a 10-minute session from our dataset in Section 4 in each 100-ms window with the amount of data received (MB), download throughput (Mbps) for the past 60 seconds, buffer level (seconds), occurrence of three playback regions, occurrence of four video states, and video resolution. At the start of the session and after each of the three stall events, notice that video resolution slowly increases before settling at a maximum level.

2.3 QoE Metrics and Prediction Challenges

This section describes the QoE metrics that we reference and the challenges in predicting these metrics. We focus on metrics that the operator can use to infer user QoE impairments in real time. Specifically, we use three QoE metrics: *buffer warning*, *video state*, and *video quality*. We do not focus on startup delay prediction, as it has been extensively studied in other works [26, 31, 35].

¹The field of ABR client algorithm design is an active research area [24, 34].

The first QoE metric we aim to predict is the current video state. The four options for video state are buffer increase, buffer decay, stall, or steady state. This metric allows for determining when the video level of the user is in the ideal situation of steady state. Video state also recognizes occurrences of buffer decay and stall events, when the operator may want to allocate more resources toward this user given that there are enough resources and the user is not limited by the data plan.

The buffer warning metric is a binary classifier for determining if the buffer level is below a certain threshold $BuffWarning_{\text{thresh}}$ (e.g., under 20 seconds). This enables operators to provision resources in real time to avoid potential stall events before they occur. For example, at a base station or WiFi AP, ABR traffic with buffer warning can be prioritized.

Another metric used is the current video resolution. Video encoders consider both resolution and target video bitrate. Therefore, it is possible to associate a lower bitrate with a higher resolution. One can argue that bitrate is a more accurate indicator of video quality. However, higher resolutions for a given clip often result in higher bitrate values. The YouTube client reports in real-time resolution rather than playback bitrate. Therefore, we use resolution as an indicator of video quality.

ABR allows the client to dynamically change resolution during a session. Frequent changes in resolution during a session tend to discourage user engagement. Real-time resolution prediction enables detection of resolution changes in a session. However, this prediction is challenging, as *download bitrate* to video resolution does not follow a 1-to-1 mapping. In addition, a video chunk received by the client can either replace a previous chunk or be played at any point in the future. Under the assumption that playback typically begins shortly (in the order of seconds) after the user requests a clip, one can associate the average download bitrate with video quality, as higher quality requires higher bitrate for the same clip. However, this is not true in a small time scale necessary for real-time prediction. Network traffic reveals the combined effect of buffer trend (increase or decay) and video playback bitrate that correlates to resolution. During steady state, a video's download bitrate is consistent with playback bitrate. However, when a client is in non-steady state, one cannot easily differentiate between the case in which a higher-resolution portion is retrieved during buffer decay state (Figure 1(a)) and the case in which a lower-resolution portion is retrieved during buffer increase state (Figure 1(b)). Both of these examples exhibit similar traffic patterns; however, the behavior of QoE metrics is dramatically different.

3 CHUNK DETECTION

The fundamental delivery unit of ABR is a chunk [27]. Therefore, identifying chunks instead of relying on individual packet data can capture important player events. Our approach is to explore the fundamental principle of HAS, which is to transmit media in the unit of video and audio chunks. The behavior of chunks over the transmission network is directly associated with the HAS protocol and behavior of the client buffer. This method is able to derive QoE metrics as long as one can (i) detect chunks and (ii) build models associating IP-level traffic information with the client buffer level. Therefore, this method is immune to changes to HAS as long as chunks can be detected from IP-level traffic.

Specifically, the occurrence of a chunk indicates that the client has received a complete segment of video or audio, resulting in an increased buffer level in playback time. An essential component of Requet in Figure 2 is its ChunkDetection algorithm to identify chunks from encrypted traffic traces. Features are extracted from the chunks and used as the input to the ML QoE prediction models. Existing work using chunks either studies per-session QoE metrics [31] instead of predicting QoE metrics in real time or lacks insight into chunk detection mechanisms [16, 28, 42]. In general, there are two approaches to identifying chunks: (i) identify a packet with non-zero payload from

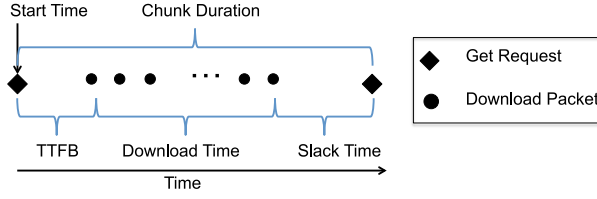


Fig. 4. Definition of chunk metrics (video or audio).

the client to the server as an HTTP request [31] and (ii) use an idle period (e.g., 900 ms is used to separate chunks in a flow of Netflix traffic [30]).

In this section, we first describe metrics capturing chunk behavior. We then develop ChunkDetection, a heuristic algorithm using chunk metrics to identify individual audio and video chunks from IP-level traces. Requet uses ChunkDetection to detect chunks from multiple servers simultaneously regardless of the use of encryption or transport protocol. It relies purely on the source/destination IP address, port, protocol, and payload size from the IP header.

3.1 Chunk Metrics

We define the following metrics for a chunk based on the timestamp of events recorded on the end device (as shown in Figure 4):

- *Start_Time*: The timestamp of sending the HTTP GET request for the chunk.
- *TTFB*: Time To First Byte (TTFB), defined as the time duration between sending an HTTP GET request and the first packet received after the request.
- *Download_Time*: The time duration between the first received packet and the last received packet prior to the next HTTP GET request.
- *Slack_Time*: The time duration between the last received packet and the next HTTP GET request.
- *Chunk_Duration*: The time duration between two consecutive HTTP GET requests. The end of the last chunk in a flow is marked by the end of the flow. Note that a different concept called *segment duration* is defined in standards as playback duration of the segment [6]. For a given chunk, *Chunk_Duration* equals *segment duration* only during steady state.
- *Chunk_Size*: The amount of received data (sum of IP packet payload size) during *Download_Time* from the IP address that is the destination of the HTTP GET request marking the start of the chunk.

Note that for any chunk, the following equation holds: $\text{Chunk_Duration} = \text{sum}(\text{TTFB}, \text{Download_Time}, \text{Slack_Time})$.

3.2 Chunk Detection Algorithm

We explore characteristics of YouTube audio and video chunks. Using the web debugging proxy Fiddler [3], we discovered that audio and video are transmitted in separate chunks, and they do not overlap in time for either HTTPS or QUIC. For both protocols, we notice that at most one video or audio chunk is being downloaded at any given time. Each HTTP GET request is carried in one IP packet with an IP payload size greater than 300 B. Smaller uplink packets are HTTP POST requests regarding YouTube log events, or TCP ACKs.

We propose a heuristic chunk detection algorithm in Algorithm 1 using notations in Table 1. ChunkDetection begins by initializing each IP flow with empty arrays for both audio and video

Table 1. Chunk Notation

Symbol	Semantics
GET_{thresh}	Packet length threshold for request (300 B)
$Down_{thresh}$	Packet length threshold for downlink data (300 B)
$GetTimestamp$	Timestamp of GET request
$GetSize$	Packet length of GET request
$DownStart$	Timestamp of first downlink packet of a chunk
$DownEnd$	Timestamp of last downlink packet of a chunk
$DownSize$	Sum of the payload of downlink packets of a chunk
$GetProtocol$	IP header protocol field
DetectAV	Sorts chunk into audio chunk, video chunk, or no chunk based on $GetSize$, $DownSize$, $GetProtocol$
\vec{Audio}	Audio chunks for an IP flow
\vec{Video}	Video chunks for an IP flow

ALGORITHM 1: Audio Video Chunk Detection Algorithm

```

1: procedure CHUNKDETECTION
2:   Initialize  $\vec{Audio}$  and  $\vec{Video}$  for each IP flow  $I$ 
3:   for each uplink packet  $p$  with IP flow  $I$  do
4:     if uplink( $p$ ) and ( $packetlength(p) > GET_{thresh}$ ) then
5:        $c \leftarrow [GetTimestamp, GetSize, DownStart,$ 
6:          $DownEnd, GetProtocol, I]$ 
7:        $AVflag \leftarrow DetectAV(c)$ 
8:       if  $AVflag == 0$  then
9:         Append  $c$  to  $\vec{Audio}$ 
10:      else if  $AVflag == 1$  then
11:        Append  $c$  to  $\vec{Video}$ 
12:      else
13:        Drop  $c$ 
14:       $GetTimestamp \leftarrow time(p)$ 
15:       $GetSize \leftarrow packetlength(p)$ 
16:       $DownFlag \leftarrow 0$ 
17:      if downlink( $p$ ) and ( $packetlength(p) > Down_{thresh}$ ) then
18:        if  $DownFlag == 0$  then
19:           $DownFlag = 1$ 
20:           $DownStart \leftarrow time(p)$ 
21:           $DownEnd \leftarrow time(p)$ 
22:           $DownSize+ = packetlength(p)$ 

```

chunks. This allows for the ChunkDetection algorithm to collect chunks from more than one server at a time.

ChunkDetection initially recognizes any uplink packet with a payload size greater than 300 B as an HTTP GET request (line 4). This threshold may vary depending on the content provider. For YouTube, we note that GET requests over TCP are roughly 1,300 bytes, whereas GET requests over UDP are roughly 700 bytes. For each new GET request, the $GetTimestamp$ and $GetSize$ are recorded (lines 14–16). After detecting a GET request in an IP flow, the chunk size is calculated by

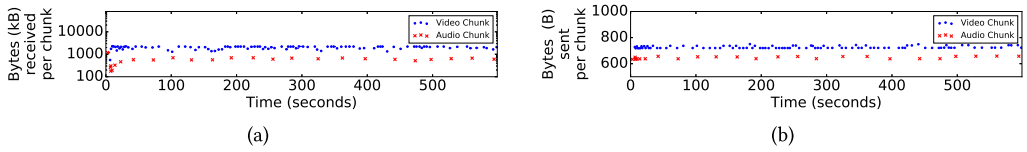


Fig. 5. Individual video/audio chunks in a 10-minute session with highest resolution (V: 1080p, A: 160 Kbps): chunk size (a) and GET request size (b).

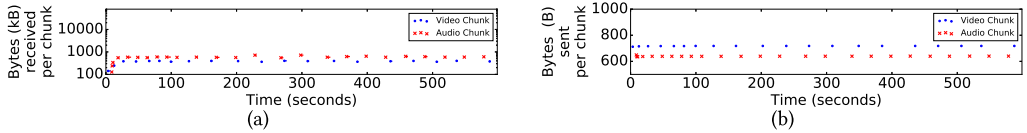


Fig. 6. Individual video/audio chunks in a 10-minute session with lowest resolution (V: 144p, A: 70kbps): chunk size (a) and GET request size (b).

summing up the payload size of all downlink packets in the flow until the next GET is detected (lines 17–22). The last downlink packet in the group between two consecutive GET requests marks the end of a chunk download. The chunk download time then becomes the difference in timestamp between the first and last downlink packet in the group.²

Once the next GET is detected, ChunkDetection records *GetTimestamp*, *GetSize*, download start time *DownStart*, download end time *DownEnd*, the protocol used *GetProtocol*, download size *DownSize*, and the IP flow *I* of the previous chunk (line 5). This allows for the calculation of chunk duration and slack time using the timestamp of the next GET. GET request size and chunk size are used in DetectAV (line 7) to separate data chunks into audio chunks, video chunks, or background traffic (lines 8–11). DetectAV uses the heuristic that the HTTP GET request size for audio chunks is slightly smaller than the request size for video chunks consistently. Figures 5 and 6 plot the HTTP GET request size and subsequent audio/video chunk size in a high (1080p) and a low (144p) resolution session, respectively. It is evident from Figures 5(b) and 6(b) that the HTTP GET request size for audio chunks is slightly smaller than that for video chunks. Through the examination of encrypted YouTube HTTP GET requests for video and audio using Fiddler, we discover that this difference is due to the additional fields used in HTTP GET requests for video content that do not exist for audio content. Furthermore, as can be seen in Figure 5(a), the video chunk size at higher-resolution levels is consistently larger than the audio chunk size. However, as can be seen in Figure 6(a), the video chunk size at lower-resolution levels can be similar to or even smaller than the audio chunk size. We can conservatively set the low threshold for chunk size to be 80 KB for our dataset. Furthermore, if download size is too small (<80KB), DetectAV recognizes that the data is neither an audio or video chunk, and the data is dropped (lines 12 and 13). This allows ChunkDetection to ignore background YouTube traffic.

4 DATA ACQUISITION

As shown in Figure 2, data acquisition provides data for training and evaluation for Requet QoE prediction models, including traffic trace collection and derivation of QoE metrics as ground truth labels associated with traffic traces. We describe additional details about the publicly available dataset in Appendix B. We collect data in two distinct settings: one using YouTube from a browser

²ChunkDetection does not flag TCP retransmission packets and therefore can overestimate chunk size when retransmission happens. ChunkDetection also assumes that chunks do not overlap in time in a given IP flow. If it happens, the individual chunk size can be inaccurate, but the average chunk size over the period with overlapping chunks is still accurate.



Fig. 7. Experimental setup for our trace collection. (a) WiFi experiments conducted in the laboratory on a laptop. (b) Cellular experiments on an Android cellphone.

on a laptop over WiFi networks (Browser-WiFi) and the other using the YouTube app on an Android smartphone over LTE cellular networks (App-LTE). We name the datasets *Browser-WiFi* and *App-LTE*, respectively. The data is collected over two different time periods to illustrate Requet’s performance, since the underlying protocol of YouTube may vary on different devices, over different networks, and over time [33].

4.1 Trace Collection from the Browser over WiFi

For the first set of experiments, we design and implement a testbed (shown in Figure 7(a)) to capture a diverse range of YouTube behavior over WiFi. We watch YouTube video clips using the Google Chrome browser on a MacBook Air laptop. We connect the laptop to the Internet via an access point (AP) using IEEE 802.11n. A shell script simultaneously runs Wireshark’s Command Line Interface, Tshark [1], and a Javascript Node server hosting the YouTube API.

The YouTube IFrame API environment collects information displayed in the “Stats for Nerds” window. From this API, we monitor video playback region (playing, paused, buffering), playback time since the beginning of the clip, amount of video that is loaded, and current video resolution. From these values, we determine the time duration of the portion of the video clip remaining in the buffer. We collect information once every 100 ms and during any change event indicating changes in video playback region or video resolution. This allows us to record any event as it occurs and to keep detailed information about playback progress.

We have two options to collect network-level packet traces in our setup: on the end device or on the WiFi AP. Collecting traces at the AP would limit the test environment to only a laboratory setup. Therefore, we opt to collect traces via Wireshark residing on the end device. This ensures that the YouTube client data is synchronized with Wireshark traces and the data can be collected on public and private WiFi networks. Our traces record packet timestamp, size, and the 5-tuple for the IP header (source IP, destination IP, source port, destination port, protocol). Our dataset contains delivery over HTTPS (9% GET requests) and QUIC (91% GET requests). We do not use any transport-level information. In addition, we record all data associated with a Google IP address. The IP header capture allows us to calculate the total number of packets and bytes sent and received by the client in each IP flow during a given time window.

To generate traces under varying network conditions, we run two categories of experiments: *static* and *movement*. For static cases, we place the end device in a fixed location for the entire session. However, the distance from the AP varies up to 70 feet or multiple rooms away. For movement cases, we walk around the corridor (up to 100 feet) in our office building with the end device while its only network connection is through the same AP.

We select 60 YouTube clips representing a wide variety of content types and clip lengths. Each clip is available in all six common resolutions from YouTube, namely 144, 240, 360, 480, 720, and 1080 p. We categorize them into four groups, where groups *A* and *B* are medium-length clips (8–12 minutes), *C* are short clips (3–5 minutes), and *D* are long clips (25–120 minutes). Table 2 lists the number of unique clips in the group, along with the length of each clip and the session length—that is, the duration for which we record the clip from its start.

Table 2. Clip Distribution in Our Dataset

Group	Clip Length	Session Length	No. of Unique Clips
A	8–12 min	10 min	40
B	8–12 min	10 min	10
C	3–5 min	5 min	5
D	25–120 min	30 min	5

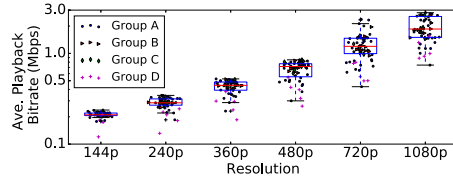


Fig. 8. Average playback bitrate vs. video resolution for clips in our dataset. Clips in all four groups are shown in scatter plots, whereas clips in group A are also shown with box plots.

For group A, we collect 425 sessions in both static (>300) and movement cases (<100) in a laboratory environment in our office building. All remaining experiments are conducted in static cases. For clips in group B, we collect traces in an apartment setting in the United States (set B_1 with 60 sessions) and in India (set B_2 with 45 sessions) reflecting different WiFi environments. We collect traces in sets C and D from the laboratory environment, where each set contains more than 25 sessions. Overall, there are more than 10 sessions for each clip in groups A and B and 6 sessions for each clip in groups C and D.

Clips in both groups A and B range from 8 to 12 minutes in length. In each session, we play a clip and collect a 10-minute trace from the moment the client sends the initial request. We choose this range of length for the client to experience buffer increase, decay, and steady state. Shorter clips with a length close to buffer capacity (e.g., 2 minutes) can sometimes never enter steady state, even when given abundant network bandwidth. In general, when there is sufficient bandwidth to support the clip's requirement, a clip can be delivered in its entirety before the end of the playback happens. On the contrary, when available network bandwidth is not enough to support the clip's requirement, a clip may experience delayed startup and even stall events.

We collect traces over 6 months from January through June 2018, with video resolution selection set to *auto*. This means that the YouTube client is automatically selecting video resolution based on changes in network conditions. For each session, we set an initial resolution to ensure that all resolution levels have enough data points.

Each group includes a diverse set of clips in terms of activity level. It ranges from low-activity types such as lectures to high-activity types such as action sequences. This fact can be seen in the wide range of video bitrates for any given resolution. Figure 8 shows the average playback bitrate for each video resolution for each clip in our dataset. All clips are shown in scatter plots, whereas clips in group A are also shown with box plots.³ One can see that the average video playback bitrate spans overlapping regions. Therefore, this cannot provide a perfect indication of the video resolution even if the entire session is delivered with a fixed resolution.

³For all box plots in the article, the middle line is the median value. The bottom and top lines of the box represent Q1 (25th percentile) and Q3 (75th percentile) of the dataset, respectively. The lower extended line represents $Q1 - 1.5IQR$, where IQR is the inner quartile range ($Q3 - Q1$). The higher extended line represents $Q3 + 1.5IQR$.

In our dataset, we notice that the YouTube buffer capacity varies based on video resolution. For example, it is roughly 60 and 120 seconds for 1080 and 144 p, respectively.

We collect data for each YouTube video session in the Chrome browser as the sole application on the end device. We record all packets between the client and any Google server. The client contacts roughly 15 to 25 different Google servers per session. We examine the download throughput (e.g., see Figure 3(a) and (b)) further by looking at the most commonly accessed server IP addresses for each session sorted by the total bytes received. Our observation is that during a session, the majority of traffic volume comes from a single to a few servers.

4.2 Trace Collection from the YouTube Android App over Cellular

Testing on a mobile device connected to a laptop computer allows us to easily connect to cellular networks, which enables testing outside of a laboratory environment over a cellular network. For the second set of experiments, we design and implement a data acquisition environment (shown in Figure 7(b)) to capture YouTube video playback statistics and encrypted network packet data on an Android device over a cellular network. We use a rooted Motorola Moto G6 smartphone connected to the Internet via Google Fi's cellular networks. A shell script autonomously sets up the testing environment using Android Debugging Bridge (ADB), collects packet data through tcpdump, and collects video playback statistics through the YouTube Android app.

The YouTube app allows for collection of video playback statistics through its "Stats for Nerds" window. This window allows us to easily monitor audio and video resolution, buffer health, and video playback region (playing, paused, and buffering). We copy the information provided by that window every 1 second to a clipboard log using ClipStack, which can then be easily exported from the device.

Because we do not have access to data going through the cellular network, we opt to collect network traffic data on the phone using tcpdump for Android. We conduct tests in multiple cellular conditions, such as in a car driving on the highway, on a Columbia University shuttle bus around upper Manhattan, in a backpack walking up and down the streets of New York City, and during lectures. We collected this set of cellular data over 7 months from June through December 2019. Again, we use the 40 unique medium length clips in group A (8–12 minutes in length). The dataset consists of more than 250 video sessions with resolution ranging from 144 to 1080 p.

5 REQUET ML FEATURE DESIGN

We develop the ML *QoE metric prediction models* for Requet by using packet traces and associated ground truth labels (Section 4). We describe in detail in Appendix A our heuristic algorithm for the video state labeling process to associate each time window with one of the four video states: buffer increase, buffer decay, steady state, and stall. As shown in Figure 2, Requet uses its ChunkDetection component (Section 3) to convert traces into chunks, followed by its feature extraction component to extract associated features.

We develop ML models using random forest (RF) to predict user QoE metrics [23]. We build the RF classifier in Python using the sklearn package. We configure the model to have 200 estimators with the entropy selection criterion and the maximum number of features per tree set to auto. We choose RF for the following reasons. First, ML classification algorithms based on decision trees have shown better results in similar problems [16, 30, 35, 37, 47, 51], with RF showing the best performance among the class [30, 47, 51]. Second, on our dataset, the feedforward neural network and RF result in roughly equal accuracy. Third, RF can be implemented with simple rules for classification in real time, well suited for real-time resource provisioning in middleboxes.

Each session in our dataset consists of (i) IP header trace and (ii) QoE metric ground truth labels generated by our video-labeling process in data acquisition (Section 4). Requet's ChunkDetection

Table 3. Percentage of Chunks in Each State (Set A Browser-WiFi)

Resolution	Video State			
	Stall	Decay	Steady	Increase
Audio	1.2	2.8	40.9	55.1
Video	3.7	5.9	47.6	42.8

Table 4. Percentage of Chunks in Each State (Set A App-LTE)

Resolution	Video State			
	Stall	Decay	Steady	Increase
Audio	2.7	4.0	52.3	41.0
Video	3.6	6.8	49.6	40.0

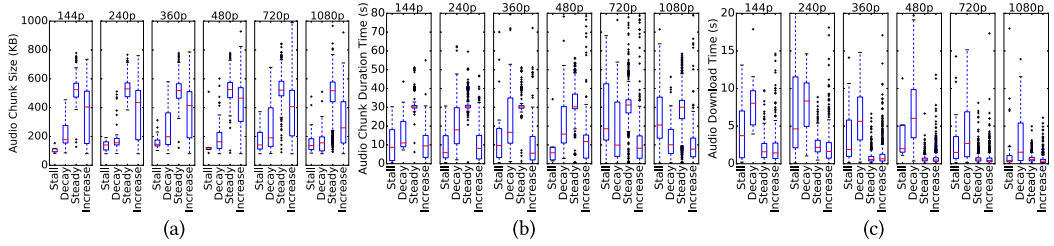


Fig. 9. Chunk metrics for all audio chunks in set A in the Browser-WiFi setting: chunk size (a), chunk duration (b), and download time (c).

(Section 3.2) transforms the IP header trace into a sequence of chunks along with the associated chunk metrics (Section 3.1). The goal of Requet QoE models is to predict QoE metrics using chunk metrics. To train such ML models, it is critical to capture application behavior associated with QoE metrics using chunk-based features. In this section, we analyze chunk behavior in our dataset (Section 5.1), explore how to capture such behavior in chunk-based features (Section 5.2), and explain how to generate baseline features used in prior work that are oblivious to chunk information (Section 5.3).

5.1 Chunk Analysis

We apply the ChunkDetection algorithm (Algorithm 1) of Requet to all sessions from the 40 clips in set A in our dataset in both Browser-WiFi and App-LTE settings.

We examine the correlation between various chunk metrics (audio or video; chunk size; chunk duration; *effective rate*, which we define as chunk size over chunk duration; TTFB; download time; and slack time) to QoE metrics (buffer level, video state, and resolution). In most cases of our dataset, for a given session, audio and video chunks are transmitted from one server. However, in some cases, audio and video traffic comes from different servers. In other cases, the server switches during a session. These findings are consistent with existing YouTube traffic studies [36].

5.1.1 Chunk Analysis in the Browser-WiFi Setting. We list the distribution of audio and video chunks along with video state at the end of chunk download in Table 3. Most of the chunks arrive during steady or buffer increase states. An extremely small fraction (4% audio and 9% video) are associated with stall or buffer decay states. They represent two possible scenarios: (i) bandwidth is limited and there are not enough chunks arriving to increase the buffer level substantially or (ii) the buffer is about to transition into increase state.

Figures 9 and 10 show the box plots for chunk duration, size, and download time for audio and video chunks, respectively. Each plus sign represents an outlier. TTFB reflects the round-trip time from the client to the server and has a median value of 0.05 seconds. This accounts for a tiny portion of chunk duration (median value ≥ 5 seconds). We can safely simplify the relationship between various chunk metrics to (slack time = chunk duration – download time). Notice that the

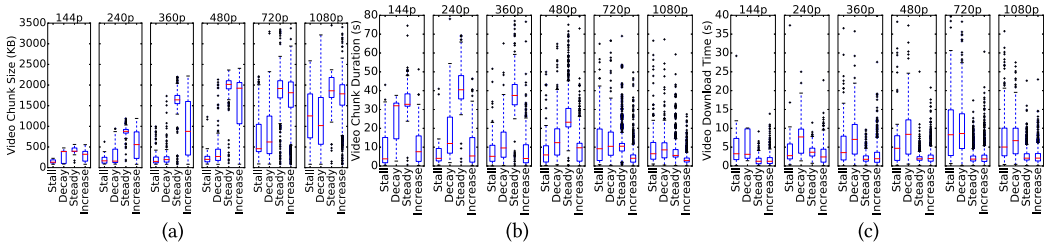


Fig. 10. Chunk metrics for all video chunks in set A in the Browser-WiFi setting: chunk size (a), chunk duration (b), and download time (c).

slack time and effective rate are derivable from the chunk duration, size, and download time. The latter three are the key metrics used in our feature selection for ML models.

Audio is encoded with CBR, however our examination of HTTP requests using Fiddler [3] reveal that in the four video states (steady, buffer increase, decay, and stall), audio chunk size decreases in the same order. This implies that audio chunk playback time also decreases in the same order. This behavior is consistent across all resolution levels (Figure 9(a)) and indicates that the audio chunk size exhibits a strong correlation with video state. Across all resolution levels, Figure 9(b) shows that the median audio chunk duration in steady and buffer increase states is roughly 30 and 10 seconds, respectively, but does not exhibit a clear pattern in stall and buffer decay states. Figure 9(c) shows that the audio chunk download time in steady and buffer increase states are similar in value, both smaller than that of stall state, which is smaller than that of buffer decay state. The longer download time is an indication that the network bandwidth is limited. This is a useful insight that the current bandwidth alone cannot reveal. For example, a specific throughput can be associated to a low resolution with the buffer increasing or a higher resolution with the buffer decreasing. All three audio chunk metrics are clearly correlated with video state.

Figure 10 shows video chunk statistics. There is a large overlap across different resolutions and video states in chunk size (Figure 10(a)) and chunk duration (Figure 10(b)). It reveals that without knowing video state, it would be difficult to determine the video resolution, chunk size, and chunk duration. For example, these statistics are very similar for a 240p chunk in buffer increase state and a 720p chunk in buffer decay state. Using audio chunk statistics to identify video state is critical in separating these two cases.

For video chunks, our examination of HTTP requests using Fiddler also shows that for a clip with a given resolution, the steady state chunk size is larger than that in the remaining three states. Figure 10(a) further shows that the median video chunk size increases as resolution increases from 144 to 480 p and stays roughly the same around 2 MB from 480 to 1080 p. Figure 10(b) shows the median chunk duration in steady state is similar for 144, 240, and 360 p, in the range of 35 to 45 seconds, and decreases from 25 seconds for 480p to 5 seconds for 1080p. To obtain a higher effective rate for higher resolutions, the chunk size levels off; however, to compensate, the chunk duration decreases. Figure 10(c) shows that the median chunk download time exhibits larger values in the stall or buffer decay state, and smaller and similar values in steady or buffer increase states. This is expected, as with limited bandwidth, a session may experience buffer decay or even stall. Both buffer decay and stall periods exhibit larger chunk download times. However, during buffer increase, retrieving smaller chunks faster than steady state results in similar download time as steady state. During steady and buffer increase states, the chunk size and duration combined provide some indication of resolution levels. However, during stall and buffer decay states, no indication can be easily seen from the three metrics.

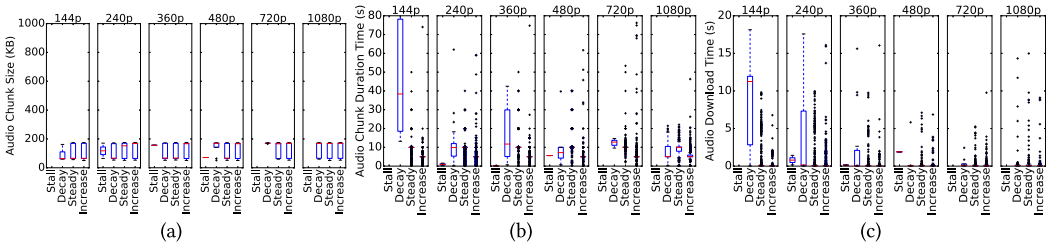


Fig. 11. Chunk metrics for all audio chunks in set A in the App-LTE setting: chunk size (a), chunk duration (b), and download time (c).

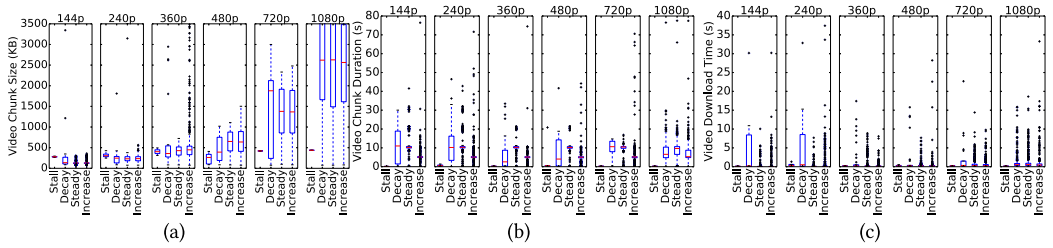


Fig. 12. Chunk metrics for all video chunks in set A in the App-LTE setting: chunk size (a), chunk duration (b), and download time (c).

To summarize, our key observations are as follows. First, without knowing video state, it would be difficult to differentiate between the two cases: (a) a higher-resolution clip in buffer decay and (b) a lower-resolution clip in buffer increase. Second, audio chunk statistics exhibit strong association with the video state. Third, the video chunk size increases and eventually levels off as resolution increases. At the same time, the video chunk duration is higher for lower-resolution levels and decreases as the resolution level increases.

5.1.2 Chunk Analysis in the App-LTE Setting. Similar to the Browser-WiFi setting, most of the chunks in the App-LTE setting arrive during steady or buffer increase states.

For the App-LTE setting, Figures 11 and 12 show the box plots for chunk duration, size, and download time for audio and video chunks, respectively. Each plus sign represents an outlier.

Across all resolution levels, Figure 11(b) shows that the median audio chunk duration is roughly 10 seconds in steady state and 5 seconds in buffer increase state. However, there is no clear pattern in stall or buffer decay states. Audio chunk size is consistent across all resolution levels (Figure 11(a)) and is usually around 70 or 170KB. These patterns are considerably different from the Browser-WiFi setting in Figure 9.

Figure 12 shows video chunk statistics in the App-LTE setting. Again, the pattern is drastically different from the pattern in the Browser-WiFi setting in Figure 10. Figure 12(a) shows that across different states in the same resolution, the chunk size is much more consistent. In addition, there is a clear pattern of increasing chunk size as the resolution increases. The video chunk duration results (Figure 12(b)) show that video chunks arrive roughly every 10 seconds during steady state and roughly every 5 seconds during buffer increase state. This video chunk arrival behavior is similar to that of the audio chunks in the same dataset. Figure 12(c) shows, with a fixed resolution, that the median video chunk download times exhibit larger values in stall or buffer decay states and smaller values in steady or buffer increase states. This is expected, because with a fixed chunk

size, a larger chunk duration is associated with limited bandwidth, which can cause a session to deplete its buffer (enter buffer decay state) or even stall.

5.2 Chunk-Based Features in Requet

Requet identifies chunks using Algorithm 1 executed over all flows during a YouTube session. For each audio or video chunk, it records the following seven chunk metrics: protocol used to send the GET request, start time, TTFB, download time, slack time, chunk duration, and chunk size. However, it does not record the server IP address from which the chunk is delivered to the end device, as it has no relationship with our QoE metrics.

Results from Section 5.1 show that the most important metrics for both audio and video are chunk size, duration, and download time. Chunk arrival is not a uniform process in time, and therefore the number of chunks in a time window vary. This would require a variable number of features. Instead, Requet uses statistics of chunk metrics in different time windows. Specifically, for the 20 windows representing the immediate past 10, 20, ..., 200 seconds, it records the total number of chunks, average chunk size, and download time for each time window, resulting in 60 features each for audio and video, and a total of 120 features.⁴ Regarding video resolution, Requet only makes predictions upon receiving a video chunk. Therefore, beyond the 120 features, it further includes the 7 features associated with the video chunk. By only collecting data on a per-chunk basis, Requet requires a minimal amount of storage of seven fields per chunk in the middlebox. Figures 11(b) and 12(b) show that chunks in the dataset arrive on average once every 5 seconds. The sliding window-based features in Requet make it ideal for middleboxes with a memory requirement of 1,016 bytes for the 127 features (assuming that each feature requires a maximum of 8 bytes).

5.3 Baseline Features

For the baseline system, we remove Requet's ChunkDetection algorithm in Figure 2 and the associated features. We replace Requet and design a baseline system with a set of features that are commonly used in prior work [26, 35, 37, 47]. Specifically, we select features that are used in more than one of these prior works and use time window-based features. We collect basic IP-level features in terms of flow duration, direction, volume (total bytes), burstiness, and transport protocol. For each 100-ms window, we calculate the total number of uplink and downlink packets and bytes, and include a one-hot vector representation of the transport protocols used for each IP address.⁵ The five features for the transport protocol are QUIC, TCP with TLS, TCP without TLS, no packets in that interval, or other. After examining the total downlink bytes of the top 20 flows in a session in our dataset, we decide to include traffic from the top 3 servers in our feature set. The remaining flows have significantly smaller traffic volume and therefore represent background traffic in a session and do not deliver video or audio traffic. By doing so, we effectively eliminate the traffic that is unrelated to our QoE metrics. In addition, we include the total number of uplink/downlink bytes and packets from the top 20 servers for the session.

We calculate the average throughput and the total number of packets in the uplink and downlink direction during a set of time intervals to capture recent traffic behavior. Specifically, we use six intervals immediately proceeding the current prediction window, and they are of length 0.1, 1, 10, 60, 120, and 200 seconds.

⁴We use the past 200 seconds of history, as the YouTube buffer rarely increases beyond 3 minutes.

⁵In natural language processing, a one-hot vector is a $1 \times N$ matrix (vector) used to distinguish each word in a vocabulary from every other word in the vocabulary. The vector consists of 0s in all cells with the exception of a single 1 in a cell used uniquely to identify the word. In our case, each IP address is treated as a word.

Furthermore, during these six windows, we record the percentage of 100-ms slots with any traffic in uplink and downlink separately. These two features are added to determine how bursty the traffic is during the given time window. In addition to the four features for the total network traffic for all servers contacted during the session, the features for each of the top three servers are the following:

- Total bytes in past 100 ms in uplink/downlink
- Total number of packets in past 100 ms in uplink/downlink
- Transport protocol (5 features)
- For each of the windows of length 1, 10, 60, 120, and 200 seconds:
 - average throughput in uplink/downlink
 - total number of packets in uplink/downlink
 - percentage of 100-ms slots without any traffic in uplink/downlink.

To summarize, for each time window, there are up to $4 + 3 \times (4 + 5 + 5 \times 6) = 121$ features for the baseline system.

6 EVALUATION

We evaluate the performance of Requet in both the Browser-WiFi setting and the App-LTE setting. For the Browser-WiFi setting, we compare the accuracy in predicting each QoE metric of Requet versus the baseline system. Both systems predict the current QoE metrics every 5 seconds, except for Requet, which predicts resolution every chunk. Since the collected network traffic transport payload is encrypted, we are unable to evaluate Requet against previous works that use DPI. Data collected as described in Section 4 is used for training, validation, and testing. Out of the four sets of traces in our dataset (Section 4.1), we use group *A*, the largest one to train both systems to predict each QoE metric in real time. We follow the same testing procedure to evaluate the performance of Requet in the App-LTE setting. We then compare the performance differences of Requet in both settings.

We extend the evaluation of Requet in the **Browser-WiFi setting** by testing Requet on smaller groups *B*, *C*, and *D*. Subsequently, we use groups B_1 and B_2 to determine how training in the laboratory environment works on clips with similar length but with different service providers and wireless network conditions. B_1 and B_2 are experiments in residential WiFi settings in the United States and India, respectively. We also use group *A* as the training set for evaluating shorter clips (group *C*) and longer clips (group *D*) in the same laboratory environment as group *A*.

For group *A*, we conduct fourfold cross validation on the 40 clips. Specifically, we divide the 40 clips into four exclusive sets, each with 10 unique clips. In each fold, we first train a model for each QoE metric using RF with features from 30 clips (three of the four sets). We then test the model on the 10 clips from the remaining set. We report each model's average performance over the four folds.

The buffer warning model produces two prediction possibilities. It indicates whether the buffer level is below the threshold $BuffWarning_{thresh}$ or not. The video state model produces four states, and the resolution model produces six resolution levels.

We report *accuracy* of each model as the ratio of the number of correct predictions over the total number of predictions. For each label a model predicts, we further report (i) *precision*, defined as the ratio of true positives to total positives (i.e., the percentage of correct predictions out of all positive predictions of a label), and (ii) *recall*, defined as the ratio of correct predictions to total true occurrences of a label (i.e., the percentage of a label correctly predicted).

6.1 Buffer Warning Prediction

The first metric we examine is buffer warning. We set the threshold for buffer-level warning, $BuffWarning_{thresh}$, to be 20 seconds. This provides ample time to provision enough bandwidth before an actual stall occurs. 二分类

For this metric, each time window in our dataset is labeled with either “no buffer warning” (NBfW) or “buffer warning” (BfW). In group A, significantly more chunks are labeled with NBfW (84%) than BfW (16%). The results in Table 5 show that in the Browser-WiFi setting, both baseline and Requet perform well for this task, with accuracy reaching 85% and 92%, respectively. We see that precision and recall for NBfW are higher than those for BfW in both the baseline and Requet. Given that the current label is BfW, Requet provides significantly higher probability of predicting BfW correctly with recall of 68% over 11% for the baseline. This is because Requet uses chunk features to detect the case when no chunks have recently arrived. However, it is difficult for the baseline system to identify such cases due to the lack of chunk detection. For example, the baseline cannot differentiate packets as being part of a chunk or background traffic.

In the App-LTE setting, Requet shows slightly improved performance compared to Browser-WiFi. Requet achieves a recall of 79.9% for BfW and 99.2% for NBfW. This results in a total accuracy of 97.8%. For the Browser-WiFi dataset, the download time and TTFB of the most recent chunk, the video chunk count, and the average video chunk size of a variety of windows are significant features that are used in the RF model for buffer warning prediction. For the App-LTE dataset, the download time and TTFB of the most recent chunk, the video chunk count, and the audio chunk count of a variety of windows are significant features that are used in the RF model for buffer warning prediction.

6.2 Video State Prediction

The results of video state prediction are shown in Table 6. In the Browser-WiFi setting, Requet achieves overall accuracy of 84%, compared to 55% for the baseline, representing a 53% improvement. Requet also outperforms the baseline in precision and recall for each state.

Stall, buffer decay, buffer increase, and steady state appear in 3.7%, 5.9%, 42.8%, and 47.6% of chunks in group A, respectively (Table 3). The precision and recall for both systems increase in the same order of stall, buffer decay, buffer increase, and steady.

However, baseline achieves below 40% in precision and recall for both stall and buffer decay states. This implies that during these two states, network traffic does not have a significant pattern for baseline to discover. Furthermore, during steady state, there can be gaps of 30 seconds or longer. A long gap also occurs when buffer is in decay state. Baseline features cannot separate buffer decay from steady state.

Examination of the Requet model reveals that audio chunk count for each 20-second window is an important feature to predict the video state. For example, if there are a few audio chunks in the past 20 seconds, it is likely that the buffer is increasing, and if there are no audio chunks in the past 120 seconds, it is likely to be in stall state. This explains the relatively high performance of Requet.

In the App-LTE setting, Requet achieves an overall accuracy of 88.2%. Compared to the Browser-WiFi dataset, Requet in the App-LTE setting achieves an improved performance in predicting the stall state but is worse in predicting buffer decay.

For the App-LTE dataset, the download time and TTFB of the most recent chunk, and the number of video chunks in the time range from 60 to 200 seconds, are significant features that are used in the RF model for state prediction. For the Browser-WiFi dataset, the download time and TTFB of the most recent chunk, the number of video chunks in the time range from 60 to 200 seconds, and the average chunk size are significant features that are used in the RF model for state prediction.

Table 5. Buffer Warning Performance with Data in Group A

Type	Baseline Browser-WiFi		Requet Browser-WiFi		Requet App-LTE	
	Precision	Recall	Precision	Recall	Precision	Recall
BfW	51.0	11.1	79.0	68.7	88.4	79.7
NBfW	86.0	98.1	94.1	96.5	98.5	99.2
Accuracy	84.9		92.0		97.8	

Table 6. Video State Performance with Data in Group A

Type	Baseline Browser-WiFi		Requet Browser-WiFi		Requet App-LTE	
	Precision	Recall	Precision	Recall	Precision	Recall
Stall	31.1	7.6	70.4	51.9	92.2	86.3
Buf. decay	32.0	16.3	78.0	78.7	65.7	25.2
Buf. increase	64.1	57.6	80.2	84.2	88.1	95.8
Steady	57.6	80.2	90.7	92.2	89.6	90.2
Accuracy	55.4		84.2		88.2	

6.3 Video Resolution Prediction

It is extremely challenging for the baseline to predict video resolution even with history of up to 200 seconds. Overall accuracy is only 22%, which is slightly better than randomly picking one out of six choices.

As seen in Figure 8, there is a large overlap of average playback bitrates of video clips of different resolutions due to varying activity levels in the video content. Without any knowledge about the content of the video or the video state, it is extremely difficult, if not impossible, to associate a chunk given its playback bitrate with the resolution with which it is encoded. Furthermore, without knowing video state, there is a large overlap in video chunk size and chunk duration across resolutions, as seen in Figure 10.

By using both audio and video chunks, Requet achieves a 66% accuracy for predicting resolution (six levels) in the Browser-WiFi setting. This result demonstrates that Requet is able to enhance video resolution prediction. By narrowing down the options in resolution to three—small (144p/240p), medium (360p/480p), and large (720p/1080p)—Requet achieves an accuracy of 87%. If the number of options is reduced to two—small (144p/240p/360p) and large (480p/720p/1080p)—the accuracy improves to 91%.

The accuracy of Requet in the App-LTE setting is 80.6%. Requet in the App-LTE setting has improved performance compared to the Browser-WiFi setting in predicting all resolutions except 480p, where it has difficulties differentiating 480p from 360p. This can be caused by the dataset having more data points during 360p, as well as having similar video chunk sizes for these two resolutions.

For both datasets, the most important features are those features related to the most recent chunk and the average video chunk size.

6.4 Performance Comparison of Browser-WiFi vs. App-LTE

The performance of Requet in the App-LTE setting is considerably greater than in the Browser-WiFi setting. As shown in Tables 5, 6, and 7, the accuracy for predicting buffer warning, video state, and video resolution in the Browser-WiFi setting is 92.0%, 84.2%, and 66.9%, respectively, whereas the accuracy for predicting buffer warning, video state, and video resolution in the

Table 7. Video Resolution Performance with Data in Group A

Type	Baseline Browser-WiFi		Requet Browser-WiFi		Requet App-LTE	
	Precision	Recall	Precision	Recall	Precision	Recall
144p	13.0	7.6	80.6	79.9	87.8	86.2
240p	14.6	10.1	68.7	64.3	74.0	81.8
360p	14.1	9.9	49.2	64.4	74.0	79.4
480p	24.7	33.3	64.9	63.8	73.7	57.2
720p	24.5	30.3	60.6	54.5	80.3	83.4
1080p	22.2	20.1	75.0	76.9	91.9	89.4
Accuracy	21.8		66.9		80.6	

App-LTE setting is 97.8%, 88.2%, and 80.6%, respectively. The only exception to this is when predicting the buffer decay state, the accuracy is higher in the Browser-WiFi setting.

There are two potential reasons for the higher accuracy in the App-LTE setting. First, across different states in the same resolution, the chunk size is more consistent in the App-LTE setting. Second, the network conditions are more stable in the App-LTE setting, due to generally good service coverage in our test area. However, in the Browser-WiFi setting, artificially varying network conditions are created from movement experiments during the data collection stage. More stable network conditions naturally lead to less variation in video states once steady state is entered.

6.5 Extended Test over WiFi Networks

Up to this point, we have reported results from our systems trained with part of group A and tested on different clips in group A in both the Browser-WiFi setting and the App-LTE setting. Next, we use group A in the Browser-WiFi setting as the training data for Requet and evaluate with groups B_1 , B_2 , C, and D. We test Requet on 10 clips from groups B_1 and B_2 for residential WiFi settings in the United States and India, respectively, to see how they perform on unseen clips of similar length and unseen WiFi environments. In addition, we use the same laboratory WiFi environment in group A to test Requet on 5 clips of shorter length of 5 minutes in group C and longer length of 25 minutes in group D. Figure 13 reports the average precision and recall of these four tests along with the fourfold cross-validation results from group A.

Depending on the environment and QoE metric, performance of these extended sets of tests either improves or deteriorates compared with results from group A reported earlier in this section. For example, groups B_1 , B_2 , and C have improved precision and recall in predicting stall and buffer decay states. Group D shows lower precision in predicting buffer decay but higher recall for both stall and buffer decay. Improved precision and recall results appear for predicting buffer threshold warning.

Accuracy for video resolution varies from experiment to experiment. Surprisingly, group B_2 has the highest overall accuracy of 70% when training with group A. This is in part due to zero 480p events collected in group B_2 . This resolution level has lower precision than 144, 240, and 1080 p (Table 7) and is extremely difficult for the other test sets to predict as well.

Most precision and recall results for other sets are better than group A, with a few exceptions. This could be due to the fact that group A includes movement experiments, whereas the other groups only contain static ones. A video session naturally exhibits different behavior in different types of environments. In addition, we plan to improve our prediction models by studying how the imbalance in data samples impacts the precision and recall of each model.

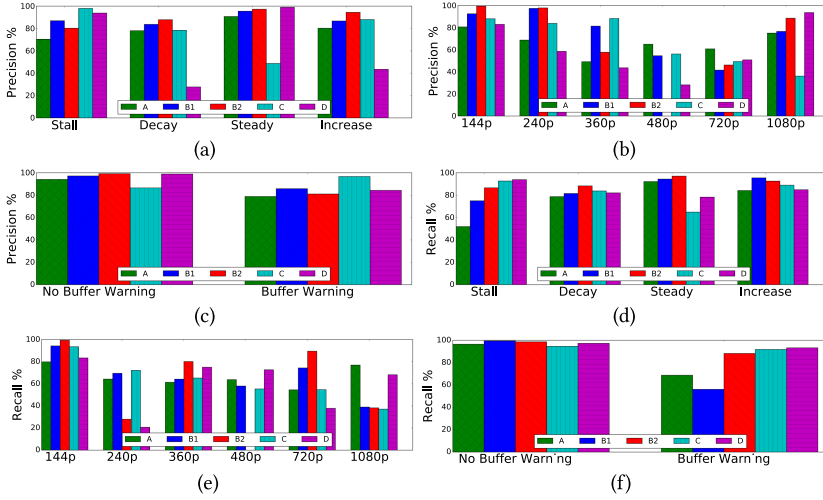


Fig. 13. Accuracy of Requet models trained with group A. (a) Precision of video state. (b) Precision of video resolution. (c) Precision of stall warning. (d) Recall of video state. (e) Recall of video resolution. (f) Recall of stall warning.

7 RELATED WORK

Traditional traffic monitoring systems rely on DPI to understand HTTP request and reply messages. The systems use meta-data to understand ABR and infer video QoE metrics. The MIMIC system estimates the average bitrate, rebuffering ratio, and bitrate switches for a session by examining HTTP logs [32]. Comparatively, BUFFEST builds ML classifiers to estimate the buffer level based either on the content of HTTP requests in the clear or on unencrypted HTTPS requests by a trusted proxy [26]. HighSee identifies HTTP GET requests and builds a linear support vector machine (SVM) [14] model to identify audio, video, and control chunks to separate audio, video, and control flows [20].

For encrypted traffic, proposals fall into two categories. The first category builds session models offline by detecting HTTP requests as in eMIMIC [31], whereas the second category builds ML models to predict QoE metrics either offline or online.

Offline models. The offline approach uses the entire video session traffic to generate features to classify the session into classes. YouQ classifies a session into two to three QoS classes [37]. The system in the work of Dimopoulos et al. [16] builds models to roughly put a session into three categories in terms of stall events (“non-stall,” “0-1 stalls,” or “2-or-more stalls”) or three classes based on average quality level. The system in the work of Orsolic et al. [38] captures IP-level traffic information (which is suitable for both TLS and QUIC traffic) and feeds ML models to predict per-session mean opinion score (MOS) (2 or 3 classes), longest resolution (“sd” vs. “hd”), and stalling occurrences (“yes” vs. “no”). Using simulation, Vasilev et al. [48] build ML models to predict average bitrate, quality variation, and three levels of stall ratio (no, mid, severe) for entire sessions using post-processing. Comparatively, Lin et al. [29] classify a session into two categories (with or without stall events) based on cell-related information collected at the start of a video session.

Focusing on the newly proposed network data analytics function in the 5G architecture, Schwarzmann et al. [43] associate network QoS metrics with the MOS for each video session using ML models. Rather than using actual network traces, the evaluation of the ML models is purely based on simulation.

Online models. The online approach uses traffic from the past time window in the session to generate features to predict QoE metrics specific to that time window. ViCrypt [44] develops ML models to predict stall events both in real time and for the entire video session based on network-level information for separate TCP and UDP flows. However, Wassermann et al. [51] build ML models to purely predict video resolution in real time using network-level information as features. The focus of the study is on feature selection and benchmarking of different ML models. Similarly, Madanapalli et al. [30] simply focus on prediction of the buffer level using features based on network-level information. It only predicts two states—buffering and stable—and discards any transition period in between (i.e., without including these data in training or testing of ML models), whereas Requet predicts video state based on buffer status in much finer granularity (with four exclusive states: buffer increase, buffer decay, steady state, and stall) without discarding any data. The system in the work of Mazhar and Shafiq [35] develops features based on both network- and transport-level information in a 10-second time window to build separate classifiers for HTTPS and QUIC traffic to infer startup delay (below or above a given threshold), stall event occurrence, and video quality level (“low” and “high”). This system uses features based on packet-level information and collects data for time windows of 100 ms. This has a relatively large memory requirement compared to Requet, which only requires network data collected on a per-chunk basis.

The system in the work of Bronzino et al. [11] uses network- and application-level features to infer startup delay and resolution. Similar to Requet, they also identify video chunks.

Flow identification. Identifying video flows from encrypted traffic is orthogonal to the QoE detection problem for given ABR flows. It is an example of the broad encrypted traffic classification problem. The Silhouette system [28] detects video chunks (also named *application data units*) from encrypted traffic in real time for ISP middleboxes using video chunk size, payload length, and download rate threshold values. The real-time system in the work of Reed and Kranch [41] identifies Netflix videos using TCP/IP header information including TCP sequence numbers. This approach relies on a “fingerprint” database built from a large number of video clips hosted by Netflix. The fingerprint is unique for each video title, and therefore it is ineffective in classifying new video titles not seen previously. The system in the work of Tsilimantos et al. [47] classifies an encrypted YouTube flow every 1 second interval into HAS or non-HAS flows in real time. For a HAS flow, it further identifies the buffer states of the video session into filling, steady, depleting, and unclear. The high accuracy to predict the buffer state is partly due to the fact that the entire dataset contains only three clips, with multiple sessions for each clip. This system also uses a feature based on the standard deviation of packet size, which is not feasible for implementation in middleboxes due to the memory requirement.

8 CONCLUSION AND FUTURE WORK

We present Requet, a system for *REal-time QUality* of experience metric detection for *Encrypted Traffic*. We focus on three QoE metrics—buffer warning, video state, and video quality—as they are crucial in allowing network-level resource provisioning in real time. We design a video state labeling algorithm to automatically generate ground truth labels for ABR traffic data.

Requet consists of the ChunkDetection algorithm, chunk feature extraction, and ML QoE prediction models. Our evaluation using YouTube traffic collected over WiFi networks demonstrates Requet using chunk-based features exhibit significantly improved prediction power over the baseline system using IP-layer features.

We demonstrate that the Requet QoE models trained on one set of clips exhibit similar performance in different network environments with a variety of previously unseen clips with various lengths. In addition, by testing with both the browser on WiFi and the YouTube application on LTE settings, we validate that Requet performs well even for different streaming algorithms.

A current limitation of Requet is that it is based on YouTube and needs to be trained separately for each streaming algorithm. Therefore, one direction of our future work includes building a generic model for a wide range of networks and client algorithms for ABR. We plan to evaluate additional services such as Disney+ and Netflix. Another direction of our future work includes using software-defined networking to utilize Requet and investigate the QoE improvements achieved via resource scheduling. We aim to study the joint effect of operator optimization and content provider video optimization mechanisms.

APPENDIXES

A VIDEO STATE LABELING

A goal for predicting video QoE in real time inside the network is to enable real-time resource provisioning to prevent stalls and decreases in video resolution. To enable this prediction, accurate labeling of video state is critical. The four exclusive video states (buffer increase, decay, stall, and steady state) accurately capture the variations in the buffer level. They can be used in combination with the actual buffer level to predict dangerous portions of ABR operation that may lead to QoE degradation. For example, when the buffer level is close to 0, a stall event is likely to happen in the near future. Increasing network capacity for the session may prevent a stall.

As shown in Section 2, playback regions reported by the client ignore buffer-level changes and cannot be used to generate video states. Prior work uses manual examination that is time consuming and can be inaccurate [47]. We opt to automate the process by developing the definition of video states based on buffer-level variation over time followed by our video state labeling algorithm. We define the four video states as follows:

- (1) *Buffer increase*: The buffer level is increasing. It has a slope greater than ϵ per second over time window T_{slope} .
- (2) *Steady state*: The buffer level is relatively flat. The slope of the buffer level is between $-\epsilon$ and $+\epsilon \frac{\text{sec}}{\text{sec}}$ over time window T_{slope} . To be in steady state, the slope needs to be in this range for greater than Thr_{SS} seconds.
- (3) *Buffer decay*: The buffer level is decreasing with a slope less than $-\epsilon \frac{\text{sec}}{\text{sec}}$ over time window T_{slope} .
- (4) *Stall*: The buffer level is less than or equal to δ .

We execute our video state labeling algorithm in Algorithm 2 for each time instance t when buffer information is recorded (every 100 ms) to determine video state for a session according to our definition.

As a chunk arrives at the client, buffer level increases by chunk length in seconds. During playback, the buffer level decreases by 1 second for every second of playback. Looking at short windows or the wrong point of a window would incorrectly determine that the buffer is decreasing. We use a smoothing function to derive a more accurate buffer slope. Specifically, we use a moving median filter over a window around t defined by $[t - T_{\text{smooth}}, t + T_{\text{smooth}}]$. We examine the rate of change of the buffer slope over a window around t defined by $[t - T_{\text{slope}}, t + T_{\text{slope}}]$.

To avoid rapid changes of the stall state, we set δ to 0.08 seconds. This value ensures that small variations in and out of stall state are consistently labeled as being in stall state. If the buffer level is above Buf_{SS} and has a slope between $-\epsilon$ and $\epsilon \frac{\text{sec}}{\text{sec}}$, then we label it as steady state. If these specifications are not met and the slope is negative, we set the state to buffer decay. If the slope is positive, we set the state to buffer increase.

To ensure that video state does not change rapidly due to small fluctuations in the buffer level, we use an additional heuristic of *SmoothState*: steady state has to last longer than Thr_{SS} . This allows

Table 8. Notation Summary

Symbol	Semantics	Defaults
δ	Stall threshold	0.08 sec
ϵ	Buffer slope boundary for steady state	0.15 $\frac{sec}{sec}$
T_{smooth}	Time window for smoothing buffer	15 sec
T_{slope}	Time window to determine buffer slope	5 sec
Buf_{ss}	Minimum buffer level to be in steady state	10 sec
Thr_{ss}	Minimum time window to stay in steady state	15 sec
$MinTime_{ss}$	Time window to look for quick changes out of steady state	10 sec
$MinTime_{stall}$	Time window to look for quick changes out of stall state	10 sec

ALGORITHM 2: Video State Labeling Algorithm

```

1: procedure VIDEOSTATELABELING
2:   Initialize  $\delta, \epsilon, T_{smooth}, T_{slope}$ 
3:   for every  $t$  do
4:     Calculate  $\hat{B}_t \leftarrow median[B_{t-T_{smooth}}, \dots, B_{t+T_{smooth}}]$ 
5:     Calculate  $m_t \leftarrow \frac{\hat{B}_{t+T_{slope}} - \hat{B}_{t-T_{slope}}}{2T_{slope}}$ 
6:     if  $\hat{B}_t \leq \delta$  then
7:        $State_t \leftarrow$  Stall
8:     else if  $-\epsilon \leq m_t \leq \epsilon$  and  $\hat{B}_t > Buf_{ss}$  then
9:        $State_t =$  Steady State
10:    else if  $m_t < 0$  then
11:       $State_t \leftarrow$  Buffer Decay
12:    else
13:       $State_t \leftarrow$  Buffer Increase
14:     $SmoothState(State_t)$ 

```

chunks with playback time longer than this value to arrive at the client. If there are changes out of and then back into stall state that last less than $MinTime_{stall}$, we consider the entire period as stall state. Similarly, if there are changes out of and then back into steady state that last less than $MinTime_{ss}$, we consider the entire period steady state. For clarity, we list all symbols in Table 8, as well as the values that we find to work the best empirically for our dataset.

B DATASET INFO

This appendix provides a description of the dataset acquired in Section 4, used for Requet chunk detection in Section 3, and for evaluation in Section 6.

The dataset can be found in the Github Repository (<https://github.com/Wimnet/RequetDataSet>). The dataset is divided into five *group folders* for data from groups A, B1, B2, C, and D, along with a summary file named *ExperimentInfo.txt* for the entire dataset. Each line in the file describes an experiment using the following four attributes: experiment number, video ID, initial video resolution, and length of experiment in seconds.

A group folder is further divided into two subfolders: one for PCAP files and the other for txt files. Each experiment is described by a PCAP file and a txt file. The PCAP file with a name in

the form of (i) '*baseline_{date}_exp_{num}.pcap*' is for an experiment where the end device is static for the entire duration, whereas a file with name in the form of (ii) '*movement_{date}_exp_{num}.pcap*' is for an experiment where the end device movement occurs during the experiment. The txt file names end with '*merged.txt*'. The txt file contains data collected from YouTube API and summary of PCAP trace for the experiment.

In each '*merged.txt*' file, data is recorded for each 100 ms interval. Each interval is represented as [*Relative Time*, # Packets Sent, # Packets Received, # Bytes Sent, # Bytes Received, [Network Info 1], [Network Info 2], [Network Info 3], [Network Info 4], [Network Info 5], ..., [Network Info 25], [*Playback Info*]].

Relative Time marks the end of the interval. *Relative Time* is defined as the time since the Javascript Node server hosting the YouTube API is started. *Relative Time* for the 0th interval is defined as 0 seconds. It is updated in intervals of 100 ms. TShark is called prior to the Javascript Node server. Therefore, the 0th interval contains Wireshark data up to the start of the Javascript Node server.

Network Info i is represented as [IP_Src, IP_Dst, Protocol, # Packets Sent, # Packets Received, # Bytes Sent, # Bytes Received] for each interval. IP_Src is the IP address of the end device. The top 25 destination IP addresses in terms of total bytes sent and received for the entire session are recorded. For each *i* of the top 25 IP_Dst addresses, the Protocol associated with the higher data volume for the interval (in terms of total number of packets exchanged) is selected, and data volume in terms of packets and bytes for each interval is reported for the IP_Src, IP_Dst, Protocol tuple in [Network Info *i*].

Playback Info is represented as [*Playback Event*, *Epoch Time*, *Start Time*, *Playback Progress*, *Video Length*, *Playback Quality*, *Buffer Health*, *Buffer Progress*, *Buffer Valid*]. From the perspective of video playback, a YouTube session can contain three exclusive regions: *buffering*, *playing*, and *paused*. YouTube IFrame API considers a transition from one playback region into another as an event. It also considers as an event any call to the API to collect data. The API enables the recording of an event and of detailed information about playback progress at the time the event occurs. *Epoch Time* marks the time of the most recent collection of YouTube API data in that interval. *Playback Info* records events occurred during the 100 ms interval.

Each field of *Playback Info* is defined as follows:

- *Playback Event*: This field is a binary array with four indexes for the following states: buffering, paused, playing, and collect data. The collect data event occurs every 100 ms once the video starts playing. For example, an interval with a Playback Event [1,0,0,1] indicates that playback region has transitioned into buffering during the 100 ms interval and a collect data event occurred.
- *Epoch Time*: This field is the UNIX epoch time in milliseconds of the most recent YouTube API event in the 100 ms interval.
- *Start Time*: This field is the UNIX epoch time in milliseconds of the beginning of the experiment.
- *Playback Progress*: This field reports the number of seconds the playback is at epoch time from the start of the video playback.
- *Video Length*: This field reports the length of the entire video asset (in seconds).
- *Playback Quality*: This field is a binary array of size 9 with indices for the following states: unlabeled, tiny (144p), small (240p), medium (360p), large (480p), hd720, hd1080, hd1440, and hd2160. The unlabeled state occurs when the video is starting up, buffering, or paused. For example, a Playback Quality [0, 1, 1, 0, 0, 0, 0, 0, 0] indicates that during the current interval, video playback experienced two quality levels—tiny and small.

- *Buffer Health*: This field is defined as the amount of buffer in seconds ahead of current video playback. It is calculated as

$$\text{Buffer Health} = \text{Buffer Progress} \times \text{Video Length} - \text{Playback Progress.}$$

- *Buffer Progress*: This field reports the fraction of video asset that has been downloaded into the buffer.
- *Buffer Valid*: This field has two possible values: True or -1. True represents when data is being collected from the YouTube IFrame API. The value -1 indicates when data is not being collected from the YouTube IFrame API during the current interval.

REFERENCES

- [1] Wireshark. n.d. About Wireshark. Retrieved May 15, 2020 from <https://www.wireshark.org/about.html>.
- [2] Cisco. n.d. Cisco Annual Internet Report (2018–2023) White Paper. Retrieved May 15, 2020 from <https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/mobile-white-paper-c11-520862.html>.
- [3] Telerik. n.d. Telerik Fiddler, the Free Web Debugging Proxy. Retrieved May 15, 2020 from <https://www.telerik.com/fiddler>.
- [4] Fortune. 2016. How Google Is Making YouTube Safer For Its Users. Retrieved May 15, 2020 from <http://fortune.com/2016/08/02/google-youtube-encryption-https/>.
- [5] Cisco. 2019. Cisco Encrypted Traffic Analytics. Retrieved May 15, 2020 from <https://www.cisco.com/c/dam/en/us/solutions/collateral/enterprise-networks/enterprise-network-security/nb-09-encrytd-traf-anlytcs-wp-cte-en.pdf>.
- [6] 3GPP. 2010. *Transparent End-to-End Packet-Switched Streaming Service (PSS)*. TS 26.234. 3rd Generation Partnership Project. 3GPP.
- [7] Vaneet Aggarwal, Emir Halepovic, Jeffrey Pang, Shobha Venkataraman, and He Yan. 2014. Prometheus: Toward quality-of-experience estimation for mobile apps from passive network measurements. In *Proceedings of ACM Hot-Mobile*.
- [8] Adnan Ahmed, Zubair Shafiq, Harkeerat Bedi, and Amir R. Khakpour. 2017. Suffering from buffering? Detecting QoE impairments in live video streams. In *Proceedings of IEEE ICNP*.
- [9] Johanna Amann, Oliver Gasser, Quirin Scheitle, Lexi Brent, Georg Carle, and Ralph Holz. 2017. Mission Accomplished? HTTPS Security AfterDigiNotar. In *Proceedings of the ACM IMC Conference*.
- [10] Lucian Armasu. 2016. Netflix Adopts Efficient HTTPS Encryption for Its Video Streams. Retrieved May 15, 2020 from <https://www.tomshardware.com/news/netflix-efficient-https-video-streams,32420.html>.
- [11] Francesco Bronzino, Paul Schmitt, Sara Ayoubi, Guilherme Martins, Renata Teixeira, and Nick Feamster. 2020. Inferring streaming video quality from encrypted traffic: Practical models and deployment experience. arXiv:1901.05800.
- [12] Pedro Casas, Michael Seufert, and Raimund Schatz. 2013. YOUQMON: A system for on-line monitoring of YouTube QoE in operational 3G networks. *SIGMETRICS Performance Evaluation Review* 41, 2 (2013), 44–46.
- [13] Giuseppe Cofano, Luca De Cicco, Thomas Zinner, Anh Nguyen-Ngoc, Phuoc Tran-Gia, and Saverio Mascolo. 2016. Design and experimental evaluation of network-assisted strategies for HTTP adaptive streaming. In *Proceedings of ACM MMSys*.
- [14] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine Learning* 20, 3 (1995), 273–297.
- [15] Yong Cui, Tianxiang Li, Cong Liu, Xingwei Wang, and Mirja Kühlewind. 2017. Innovating transport with QUIC: Design approaches and research challenges. *IEEE Internet Computing* 21, 2 (2017), 72–76.
- [16] Giorgos Dimopoulos, Ilias Leontiadis, Pere Barlet-Ros, and Konstantina Papagiannaki. 2016. Measuring video QoE from encrypted traffic. In *Proceedings of ACM IMC*.
- [17] 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 *Proceedings of ACM SIGCOMM*.
- [18] Zakir Durumeric, Zane Ma, Drew Springall, Richard Barnes, Nick Sullivan, Elie Bursztein, Michael Bailey, J. Alex Halderman, and Vern Paxson. 2017. The security impact of HTTPS interception. In *Proceedings of NDSS*.
- [19] Roy T. Fielding and Julian F. Reschke. 2014. *Hypertext Transfer Protocol (HTTP/1.1): Message Syntax and Routing*. RFC 7230. IETF Trust.
- [20] S. Galetto, P. Bottaro, C. Carrara, F. Secco, A. Guidolin, E. Targa, Claudio Narduzzi, and Giada Giorgi. 2017. Detection of video/audio streaming packet flows for non-intrusive QoS/QoE monitoring. In *Proceedings of IEEE MN*.

- [21] Thiago A. Guarnieri, Idilio Drago, Alex Borges Vieira, Ítalo Cunha, and Jussara M. Almeida. 2017. Characterizing QoE in large-scale live streaming. In *Proceedings of IEEE GLOBECOM*.
- [22] Craig Gutterman, Katherine Guo, Sarthak Arora, Xiaoyang Wang, Les Wu, Ethan Katz-Bassett, and Gil Zussman. 2019. Requet: Real-time QoE detection for encrypted YouTube traffic. In *Proceedings of ACM MMSys*.
- [23] Tin Kam Ho. 1995. Random decision forests. In *Proceedings of IEEE ICDAR*.
- [24] 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 *Proceedings of ACM SIGCOMM*.
- [25] Arash Molavi Kakhki, Samuel Jero, David R. Choffnes, Cristina Nita-Rotaru, and Alan Mislove. 2017. Taking a long look at QUIC: An approach for rigorous evaluation of rapidly evolving transport protocols. In *Proceedings of ACM IMC*.
- [26] Vengatanathan Krishnamoorthi, Niklas Carlsson, Emir Halepovic, and Eric Petajan. 2017. BUFFEST: Predicting buffer conditions and real-time requirements of HTTP(S) adaptive streaming clients. In *Proceedings of ACM MMSys*.
- [27] Will Law. 2018. *Ultra-Low-Latency Streaming Using Chunked-Encoded and Chunked-Transferred CMAF*. Technical Report. Akamai.
- [28] Feng Li, Jae Won Chung, and Mark Claypool. 2018. Silhouette: Identifying YouTube video flows from encrypted traffic. In *Proceedings of ACM NOSSDAV*.
- [29] Yu-Ting Lin, Eduardo Mucelli Rezende Oliveira, Sana Ben Jemaa, and Salah-Eddine Elayoubi. 2017. Machine learning for predicting QoE of video streaming in mobile networks. In *Proceedings of IEEE ICC*.
- [30] Sharat Chandra Madanapalli, Hassan Habibi Gharakheili, and Vijay Sivaraman. 2019. Inferring Netflix user experience from broadband network measurement. In *Proceedings of IEEE TMA*.
- [31] Tarun Mangla, Emir Halepovic, Mostafa Ammar, and Ellen Zegura. 2018. eMIMIC: Estimating HTTP-based video QoE metrics from encrypted network traffic. In *Proceedings of IEEE TMA*.
- [32] Tarun Mangla, Emir Halepovic, Mostafa H. Ammar, and Ellen W. Zegura. 2017. MIMIC: Using passive network measurements to estimate HTTP-based adaptive video QoE metrics. In *Proceedings of IEEE TMA*.
- [33] Ahmed Mansy, Mostafa H. Ammar, Jaideep Chandrashekar, and Anmol Sheth. 2014. Characterizing client behavior of commercial mobile video streaming services. In *Proceedings of ACM MoVid*.
- [34] Hongzi Mao, Ravi Netravali, and Mohammad Alizadeh. 2017. Neural adaptive video streaming with Pensieve. In *Proceedings of ACM SIGCOMM*.
- [35] M. Hammad Mazhar and Zubair Shafiq. 2018. Real-time video quality of experience monitoring for HTTPS and QUIC. In *Proceedings of IEEE INFOCOM*.
- [36] Abhijit Mondal, Satadal Sengupta, Bachu Rikith Reddy, M. J. V. Koundinya, Chander Govindarajan, Pradipta De, Niloy Ganguly, and Sandip Chakraborty. 2017. Candid with YouTube: Adaptive streaming behavior and implications on data consumption. In *Proceedings of NOSSDAV*.
- [37] Irena Orsolic, Dario Pevec, Mirko Suznjevic, and Lea Skorin-Kapov. 2016. YouTube QoE estimation based on the analysis of encrypted network traffic using machine learning. In *Proceedings of IEEE Globecom Workshops*.
- [38] Irena Orsolic, Mirko Suznjevic, and Lea Skorin-Kapov. 2018. YouTube QoE estimation from encrypted traffic: Comparison of test methodologies and machine learning based models. In *Proceedings of QoMEX*. IEEE, Los Alamitos, CA, 1–6.
- [39] Stefano Petrangeli, Tingyao Wu, Tim Wauters, Rafael Huysegems, Tom Bostoen, and Filip De Turck. 2017. A machine learning-based framework for preventing video freezes in HTTP adaptive streaming. *Journal of Network and Computer Applications* 94 (2017), 78–92.
- [40] Abbas Razaghpanah, Arian Akhavan Niaki, Narseo Vallina-Rodriguez, Srikanth Sundaresan, Johanna Amann, and Phillipa Gill. 2017. Studying TLS usage in Android apps. In *Proceedings of ACM CoNEXT*.
- [41] Andrew Reed and Michael Kranch. 2017. Identifying HTTPS-protected Netflix videos in real-time. In *Proceedings of CODASPY*.
- [42] Paul Schmitt, Francesco Bronzino, Renata Teixeira, Tiathi Chattopadhyay, and Nick Feamster. 2018. Enhancing transparency: Internet video quality inference from network traffic. In *Proceedings of TPRC46*.
- [43] Susanna Schwarzmann, Clarissa Cassales Marquezan, Marcin Bosk, Huiran Liu, Riccardo Trivisonno, and Thomas Zinner. 2019. Estimating video streaming QoE in the 5G architecture using machine learning. In *Proceedings of the ACM MobiCom Internet-QoE Workshop*.
- [44] Michael Seufert, Pedro Casas, Nikolas Wehner, Li Gang, and Kuang Li. 2019. Features that matter: Feature selection for on-line stalling prediction in encrypted video streaming. In *Proceedings of IEEE INFOCOM Network Intelligence: Machine Learning for Networking Workshop*.
- [45] Thomas Stockhammer. 2011. Dynamic adaptive streaming over HTTP: Standards and design principles. In *Proceedings of ACM MMSys*.
- [46] Dimitrios Tsilimantou, Theodoros Karagkioulos, and Stefan Valentin. 2018. Classifying flows and buffer state for YouTube's HTTP adaptive streaming service in mobile networks. arXiv:1803.00303.

- [47] Dimitrios Tsilimantos, Theodoros Karagkioulos, and Stefan Valentin. 2018. Classifying flows and buffer state for YouTube's HTTP adaptive streaming service in mobile networks. In *Proceedings of ACM MMSys*.
- [48] Vladislav Vasilev, Jérémie Leguay, Stefano Paris, Lorenzo Maggi, and Mérouane Debbah. 2018. Predicting QoE factors with machine learning. In *Proceedings of IEEE ICC*.
- [49] Nick Vogt. 2015. YouTube Audio Quality Bitrate Used for 360p, 480p, 720p, 1080p, 1440p, 2160p. Retrieved May 15, 2020 from <https://www.h3xed.com/web-and-internet/youtube-audio-quality-bitrate-240p-360p-480p-720p-1080p>.
- [50] Florian Wamser, Michael Seufert, Pedro Casas, Ralf Irmer, Phuoc Tran-Gia, and Raimund Schatz. 2015. YoMoApp: A tool for analyzing QoE of YouTube HTTP adaptive streaming in mobile networks. In *Proceedings of EuCNC*.
- [51] Sarah Wassermann, Michael Seufert, Pedro Casas, Li Gang, and Kuang Li. 2019. I see what you see: Real time prediction of video quality from encrypted streaming traffic. In *Proceedings of the ACM MobiCom Internet-QoE Workshop*.
- [52] Nicolas Weil. n.d. The State of MPEG-DASH 2016. Retrieved May 15, 2020 from <http://www.streamingmedia.com/Articles/Articles/Editorial/Featured-Articles/The-State-of-MPEG-DASH-2016-110099.aspx>.

Received December 2019; revised March 2020; accepted April 2020