

IBM Aspera Streaming for Video with FASPStream

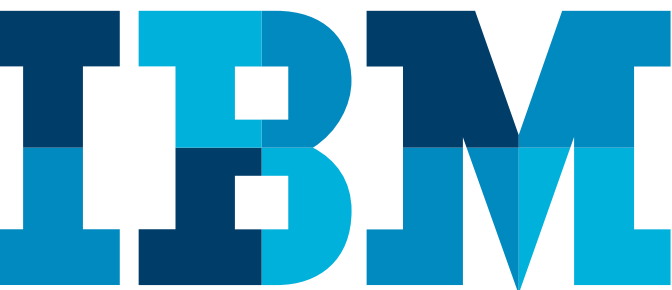
Enabling near-live and live streaming of any format and bit rate over commodity Internet

Contents:

- 1 Overview
 - 2 The challenge
 - 2 The solution
 - 3 The FASPStream model
 - 3 *Enabling start-up delay*
 - 4 *“Glitch” probability and the full mathematical model*
 - 5 *Some expectation values*
 - 5 *Data pipeline*
 - 6 *Skip heatmap*
 - 6 *Benchmarking our improvement over TCP*
 - 9 Conclusion
-

Overview

Live and near-live streaming of broadcast quality video content (10 - 50 Mbps) over IP networks with small start-up delays and glitch-free experiences have traditionally required expensive and specially-provisioned infrastructure. Traditional distribution systems use live satellite feeds from the streaming source or dedicated terrestrial networks with heavy quality of service to ensure low latency and low packet loss rates so as to not degrade the play out quality. In recent years, advances in broadcast quality video compression have also made it possible to distribute lower bit rate versions through web streaming over consumer broadband bandwidths (1-10 Mbps), but even at these bit rates, providers have relied heavily on global content distribution networks (CDNs) to stream the video content from the network edge to the consumer. These CDNs allow the user to take advantage of low round-trip latency and relatively low packet loss rates for a better quality user experience. Note that lower bit rate transcoding does not eliminate the need to ingest a high quality stream to transcoders, which may be remote from the source.



The challenge

At both ends of the spectrum — broadcast quality ingest and remote play out, as well as consumer web streaming — the media enterprise and, indirectly, the consumer pay a heavy premium in infrastructure costs to minimize the network factors that degrade the play out experience, such as network round-trip time and packet loss. This cost is burdensome in any media application or service, but is especially impractical in live and second screen experiences for events that occur only once. One time events, such as sport events, movie premieres, concerts, operas, etc., cannot as easily justify the investment associated in erecting dedicated infrastructure for direct distribution, or amortize the CDN costs over long periods of viewing. Additionally, there exist practical constraints that make it difficult to employ CDNs for direct distribution in many second screen applications; media needs to flow through distant cloud services where scale out properties of the cloud computing platform are necessary for concurrent transcoding of the live stream for several formats. Thus, more often than not, content providers are left to over-provision infrastructure for such live events as a precaution, and pay higher costs.

The need to ingest and distribute live media over low-cost, wide area IP networks, and with the option to go through distant cloud-based transcoding platforms has created a significant technology void. This is not just an opportunity for small incremental advantage solvable through adaptive bit rate “down sampling”, or clever buffering schema — instead, this calls for a fundamental solution.

Traditional TCP-based transport approaches such as adaptive bit rate streaming over HTTP have a significant bottleneck in throughput over commodity Internet WANs. The best case transfer rates for live HTTP streams over a commodity Internet WAN path between South America and Europe (200 milliseconds round-trip time) is 2 Mbps, and with worst case internet packet loss rates ($\geq 2\%$), falls to < 1 Mbps. A 10 Mbps live ingest stream transported with TCP simply isn't possible over these distances. And for higher video bandwidths or more distant delivery, the gap between realizable throughput and required play-out rate widens.

The solution

Aspera FASP transport technology is a patented bulk data transport widely utilized in digital media for achieving highly efficient, high-speed bulk media transfer over IP networks, with efficiency independent of distance and quality (round-trip latency and packet loss). However, the FASP architecture originally had no suitable application interface for transporting live data. In contrast with bulk file-structured data (e.g. VoD), live stream data need be delivered to the play out application in the same order it was passed to the transport tier. This ordered delivery constraint required Aspera to innovate a new byte streamlining capability in its transport platform on top of the FASP datagram delivery protocol. The resulting protocol— “FASPStream” — is a fully reliable bulk data streaming protocol that delivers data and video streams over Internet WANs including minimal buffering or glitches, and with negligible start-up delay. In this paper we describe the underlying protocol design and statistical model that predicts the FASPStream performance and we demonstrate through real world measurements the resulting quality in live video delivery that creates radically new possibilities for live streaming video. Media companies can achieve long distance ingest, remote play out, and even distribution of live video to play out systems running FASPStream without the assistance of CDN edge technology, “game changing” capabilities that could ultimately revolutionize the transport of live video.

Aspera has productized the FASPStream technology as part of the Aspera Streaming for Video solution. It is also available as an embeddable SDK bindings available for C/C++, .NET, and Java for use in third-party applications. The first production use case delivered live video during the World Cup in a pioneering second screen system that captured 14,000 hours of live video from multiple camera angles, ingested this video in real time using Aspera FASP from Brazil to an AWS cloud storage for live transcoding, ultimately yielding approximately 3 million minutes of transcoded content served by broadcasters, and ultimately to consumers. Subsequently for the 2018 World Cup, the Aspera FASPStream technology integrated into Telestream's Vantage and Lightspeed Live products delivered about 4000 hours of HD and UHD live feeds (nearly two petabytes of video) from all 64 World Cup matches in Russia to remote production teams in Los Angeles in near real-time — over unmanaged IP networks.

We begin with a tour of the mathematical model that underlies the protocol, and predicts its efficiency. We then generate live stream data on networks similar to those described above to compare to our expectations—200 ms delay, 2% packet loss— and three representative bit rates (6 Mbps, 10 Mbps, 40 Mbps). These actual transfer data sets confirm both our model, and the transfer protocol as bona fide solutions to the streaming video challenge. We end with some simple visualizations to expose the hard numbers that FASPStream achieves in terms of skips per second of video playback.

The FASPStream Model

Consider a live video stream being “played out” over a wide area IP network. Viewers react to two important measures of “quality” — how long the stream takes to start playing (the start up delay), and whether the stream plays smoothly or “glitches”/“buffers” waiting for the expected next data to arrive (the “glitch” probability). Our formal evaluation of the FASPStream model in this section considers both facets of quality.

Enabling Start-up Delay

Let’s start with quantifying the start-up delay. If all of the packets that comprise the video stream were to always arrive at the play out location on the initial attempt, to determine the delay before play out begins, we would need only ask, what is the One-Way-Transfer time for the network? Simply knowing that time would determine how long the first packet in the stream would take to arrive, and all of the remaining packets would also arrive precisely in time for the player. However, IP networks are imperfect and merely “best effort” by design: packets are lost or delayed in transit and there is no guarantee that any packet will arrive at its destination on the first attempt or at any particular time!

Assuming a typical full duplex intercontinental IP network with a 200 millisecond round-trip time (RTT) and one way propagation delay (OWD) of 100 milliseconds and loss probability of 2% in each direction, for a reliable delivery protocol like HTTP over TCP, the transmission rate has been shown to collapse to an effective rate of less than 1 Mbps, due to the congestion windowing protocol of TCP that reduces the sending rate aggressively in response to packet loss. Thus “live” transmission and play out of data stream rates greater than 1 Mbps using TCP based protocols is impossible at such network distances because start-up delays to compensate are unbearably long.

The Aspera FASP protocol, however, has a reliability algorithm fully decoupled from its congestion avoidance algorithm and independent of round-trip delay. If data is dropped, the FASP receiver requests retransmission of dropped packets by the sender, and does not slow down its the sending rate; packets are retransmitted along with new data at the effective bandwidth rate of the channel. So, if the FASP design principals hold for in-order stream data, we can predict precisely how long we need to wait for a video stream to begin.

For example, we conclude that for a live stream over a WAN with 2% packet loss, on average, only one in every 50 data blocks will be dropped, and for 50 blocks to arrive, we need to wait only the OWD time, plus the time to retransmit a block, 1 round trip time (RTT). While close to accurate already, this doesn’t entirely capture the story. Let’s scale up a bit.

Consider 2500 blocks to be sent. 2% of these blocks are expected to drop, i.e. 50. For each of these 50 blocks, we have two possibilities:

Capability details are highlighted below.

- The message requesting block retransmission sent by the receiver to the sender could be dropped in transit. This may occur with probability 2%
- The retransmitted block might be dropped on its second trip. This may occur with probability 2%.

One can readily compute that we expect 1 of our 50 blocks to be dropped a second time, and 1 of our retransmission requests to be dropped. How would this in turn affect the stream delivery time? Of all blocks, we would expect 2450 to arrive on time, that is, after the OWD time. We expect 48 to arrive after the retransmission request. This is 1 RTT plus 1 OWD. For the two remaining blocks (the double dropped blocks), their arrival delay will suffer another RTT; arrival will occur in 2 RTT + 1 OWD.

Finally, we consider the special case that if the final block’s retransmission request is dropped, the FASP protocol will detect this case and one RTT later, send another retransmission request. Then the block is retransmitted successfully (with 98% likelihood). This entire process takes 1 RTT + 1 RTT + 1 OWD. Thus, we have two blocks arriving after 2 RTT + 1 OWD. To give a sense of this in human time, RTT in our model network is 200ms, and OWD is 100ms.

Thus, in 2500 blocks, all but two arrive within three-tenths of a second, and the last two, in three-fifths of a second. This means that we can ensure a worst-case start up delay of only three-fifths of a second, assuming FASPStream can continue to deliver data at the play out rate and with no data pre-buffered in reserve.

Lest these predictions of block numbers seem sleight-of-hand, consider the actual calculated number of blocks for some sample tests. For 6 Mbps data streams, 547 blocks are sent per second; for 10 Mbps streams, 911 blocks are sent per second, and for 40 Mbps streams, 3641 blocks are sent per second. Compare these numbers with the sample numbers described above, and you'll get an immediate sense on the performance we are expecting. Later we will show some real data to compare, but for now, these are concrete numbers to keep our feet on the ground.

“Glitch” Probability and The Full Mathematical Model

The more interesting question in characterizing the FASPStream performance is how to know the probability that the live stream will “glitch” during playback because the next expected data is not available to the player, and how much additional start up delay is needed to ensure enough data is stored in reserve to prevent this given the network conditions. To formalize this, we set out to compute a probability model for “skipping” a frame. Specifically, we want to know the probability that a skipped block will not be received in time for its expected playback.

First we need consider how the FASP datagram size equates to playback time. Video is delivered in a number of frames per second (the 'framerate') and has a number of bytes per frame. For one second of live video, the number of data bytes the FASPStream will send is framerate multiplied by framesize. In other words, the FASPStream transmission rate in bits per second equals video framerate multiplied by framesize multiplied by 8. We assume also that the FASPStream transport reads the video stream from the video source in “readchunks” of a known size, and delivers these chunks to the receiving player application in order (while individual blocks within the chunks may arrive out of order). It turns out that the size of the chunk has no affect on the formal model.

To determine the probability of a “glitch” in the live transmission, we need to first compute the probability P of waiting a number M RTTs before the next needed chunk arrives and then normalize M for the playback rate of the video.

Note: The absolute baseline for any block to travel the network is 1 OWD. For this reason, we normalize by this OWD and don't include it in all of our computations. When we say “wait M RTTs”, we always mean “wait M Rtt after the first OWD”.

The probability model is an iterated probability of a single block failing. For streaming files, order matters, and thus the delay of playback is the result of a packet being dropped multiple times.

The probability of waiting one RTT for one block is P, i.e. the probability that one block is dropped. The probability of waiting two RTTs for one block is P², or, the probability of that block being dropped, and then being dropped again. This may be iterated for M RTTs.

Hence, the probability of a block not being received in M RTTs is P^{M+1}. The index change is to reflect precisely how many RTTs are needed to receive the block, i.e. the block arrives within M RTTs.

From this we can see that the probability that a block will be received in M Rtos is 1 – P^{M+1}. Now given two blocks, the probability that we receive the first block within M RTTs, AND we receive the second block within M RTTs is the product of their individual probabilities. This is because each block transmission is an independent event. Whence:

$$P(M,2) = (1 - P^{M+1})(1 - P^{M+1}) = (1 - P^{M+1})^2. \quad (1)$$

Now, assume N is the number of blocks in a single readchunk. By this we mean N is the readchunk size divided by the block size. This unit represents the minimal number of blocks to begin playing.

Now, the probability of receiving N blocks in greater than M RTTs is

$$P(M,N) = 1 - (1 - P^{M+1})^N \quad (2)$$

This is a result of the previous fact: that receiving all N packets in M RTTs is the product of their individual probabilities.

Some expectation values

To test our theory we ran a number of experiments to create a dataset consisting of 100 transfers of 30 seconds of video, using two host computers connected over an IP network (10 Gbps Ethernet end-to-end) via a network emulation device configured with 200ms round trip delay and 2% packet loss. We created data streams at 40, 10, and 6 Mbps via netcat to serve as representative video bit rate examples. We assumed that the available bandwidth capacity of the network was greater than the data stream rate and configured the emulated channel bandwidth to 50, 15, and 10 Mbps respectively. Finally, we captured the arrival timestamps of all packets at the FASPStream receiver host and recorded the arrival rate for original and retransmitted packets by block number and by readchunk for our analysis.

For our three cases of 40 Mbps, 10 Mbps, and 6 Mbps we estimate 3641, 911, and 547 blocks per readchunk, where we assume the readchunks are 1 second long, assuming a block size of 1464 bytes. We assume that all FASPStream applications can have a 1 second start-up delay and thus pre-buffer one chunk of one second duration, which translates to 5 RTTs on a 200ms network. We assume a packetloss probability of 2%, and calculate the probability that any 1 chunk will be delayed more than the 1 second of buffer and cause a “glitch” as follows:

$$P(5, 547) = 1 - (1 - .025)^{547} = 0.00000175 \quad (3)$$

$$P(5, 911) = 1 - (1 - .025)^{911} = 0.00000292 \quad (4)$$

$$P(5, 3641) = 1 - (1 - .025)^{3641} = 0.00001165 \quad (5)$$

After running 100 tests of 30 seconds of live data(30 chunks) the probability of a “skipped” chunk is as follows:

$$0.00000175 \times 30 \times 100 = 0.00525 \quad (6)$$

$$0.00000292 \times 30 \times 100 = 0.00875 \quad (7)$$

$$0.00001165 \times 30 \times 100 = 0.03495 \quad (8)$$

This data tells us that with a one second buffer, at 6 and 10 Mbps, we should skip less than 1% of the time, and even at 40 Mbps, we should skip less than 4% of the time.

Video bit rate	Time before expected skip
6 Mbps	6.6 days
10 Mbps	3.96 days
40 Mbps	0.99 days

Our tests of the FASPStream experimentally confirm these bold claims. We see exactly as predicted, zero observable skips, and the majority of frames arrive ahead of time, suggesting the buffer time can be reduced under a second (on the order of 0.5 seconds).

Data pipeline

In our experiment we record the arrival rate of blocks of data including the following five fields:

Timestamp — Time since last block — Full block size — Payload size — Block Number — Original or Retransmission

All times are computed in microseconds. We ran 100 transfers of 30 seconds duration, aggregate the data and calculate the “lateness” of each block (original or retransmission), where the lateness is given as:

$$\text{Lateness} = \text{FinalRexTime} - \text{SenderTime} - \text{OneWayDelay},$$

Where FinalRexTime is the timestamp of the Received Retransmission, and SenderTime is the timestamp when the corresponding original block was sent. Thus, the Lateness is the amount of time passed between the sender originally trying to send, and the final time the block is successfully delivered to the receiver, less the OWD.

The RexCount is computed by counting the number of retransmission requests for that block sent on the receiver side.

We expect that

$$\text{RexCount} \times \text{RTT} \leq \text{Lateness}. \quad (9)$$

From this dataset, we generate our scatterplots (figures 1, 2, and 3). These scatterplots display test number vs. lateness on a per block basis. Some immediate, but important observations:

- There is no statistically significant correlation between test number and lateness.
- The vertical bands in the data, cluster exactly into the expected regions predicted by the FASP protocol.
- Our predicted probabilities are manifest in the dataset.
- We have no inexplicable delays.

Skip heatmap

Rather than visualize the absolute values of lateness on a per block basis, we want to frame this data in terms observable by users. Specifically, we want to look at the lateness on a per-readChunk level. Humans don't see when a block is late – instead they see when a block is so late that the video skips. The readChunk abstraction bridges us from what we see, to what is happening at the protocol level. Heatmaps are a visual technique for associating regions to values, and are an excellent tool for identifying values above a certain threshold. In our case, we will use white—or “hot”—squares to refer to readchunks that are troublesome. Green squares—or, “cool” squares—designate readchunks that arrive close to the minimum possible dictated by network conditions. Our output will appear as a matrix of squares, where rows correspond to individual tests, and columns indicate the index of the readchunk in the video. We will hope to see primarily green; as the squares grade to white, the corresponding readChunks are arriving later. We now need the data in a very different form:

TestNum — ChNum — Lateness

Again, TestNum and ChNum are as before, however Lateness is computed differently. We want to see how long the entire readChunk will take to arrive. We can't simply look at the timestamp of the last block in a readchunk and subtract the timestamp of the first block, because this doesn't respect the possibility that some block

retransmissions will come later. So we look at the latest arrival of any block in the readChunk, and subtract from this time the first arrival time. This tells us how long the block took to arrive. We subtract the readChunk's index times the buffer time to convert this value to a comparison to when it was expected.

In symbols,

$$\text{ChunksLastRexTime} - \text{ChunksFirstTime} - \text{BufferTime} \times (\text{ChNum}) = \text{ChunkLate}. \tag{10}$$

We see excellent performance. In the three test case scenarios we see performance as predicted by the theoretical model. We see in Figure 4 almost exclusively green, and even better, a dark green, close to ideal. Nothing in our testing suggests our model was overoptimistic, and to the contrary, provides us with evidence that we can effectively operate under these assumptions.

Benchmarking our improvement over TCP

The TCP protocol underlying HTTP adaptive bit rate streaming protocols behaves quite differently from FASPStream and is not directly comparable under this model. TCP's congestion avoidance algorithm slows down so aggressively with packet loss that much longer start-up delays are necessary for live streams. For near-live use cases, a huge delay is problematic. To quantify the difference, we show show the kind of startup delay anticipated, for direct comparison to FASPStream.

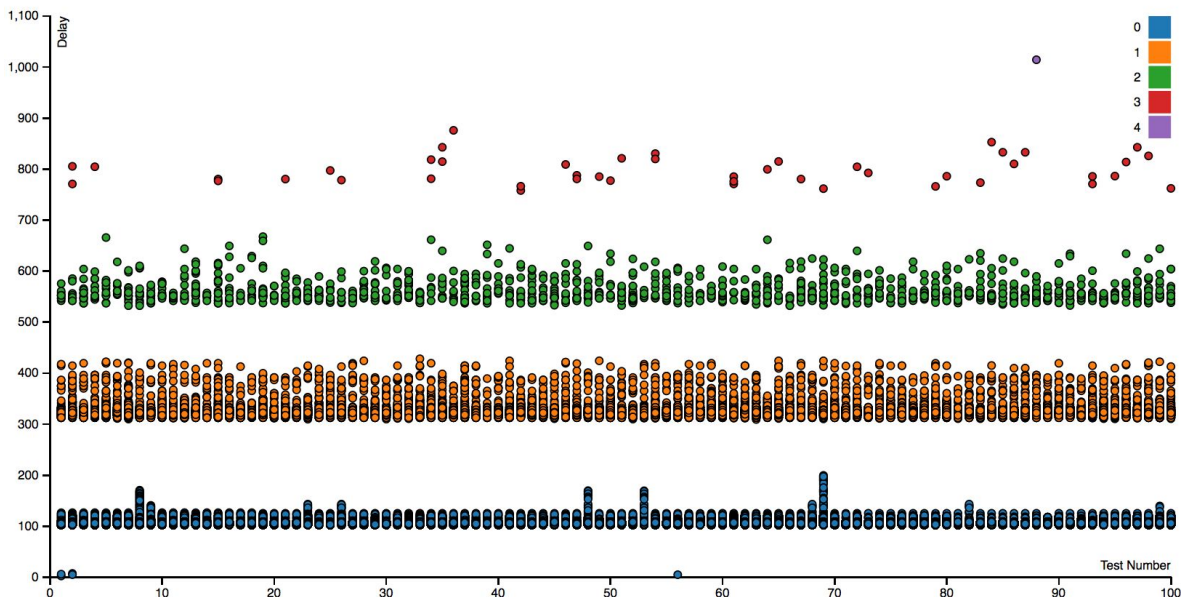


Figure 1: 6 Mbps Lateness Scatterplot

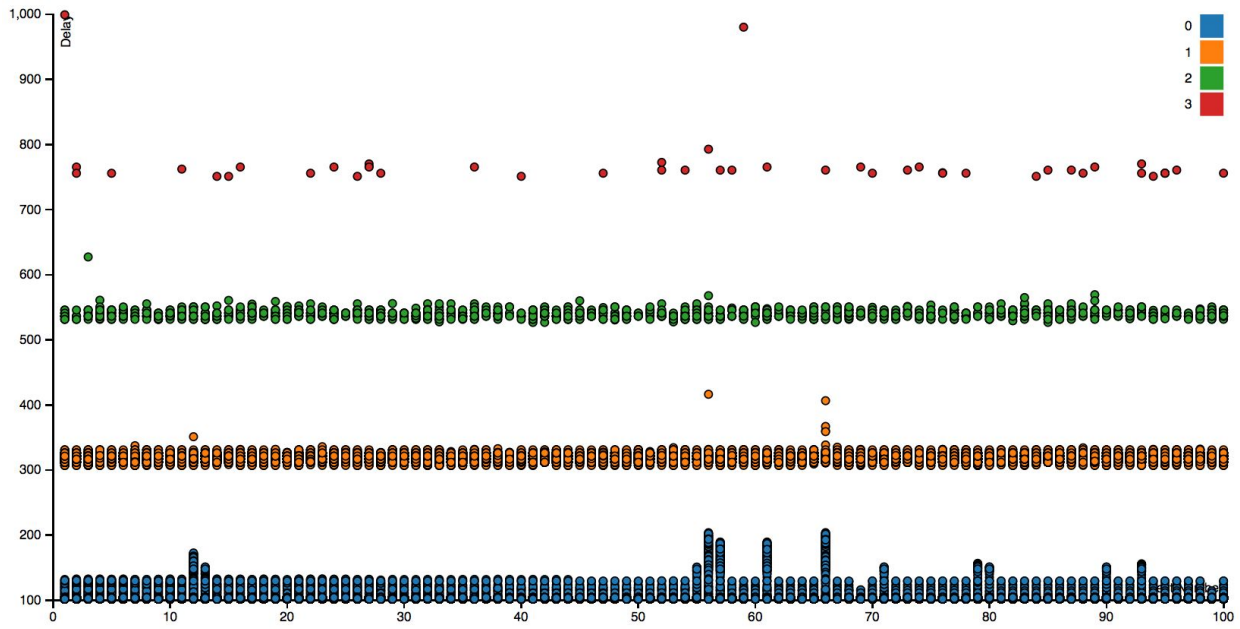


Figure 2: 10 Mbps Lateness Scatterplot

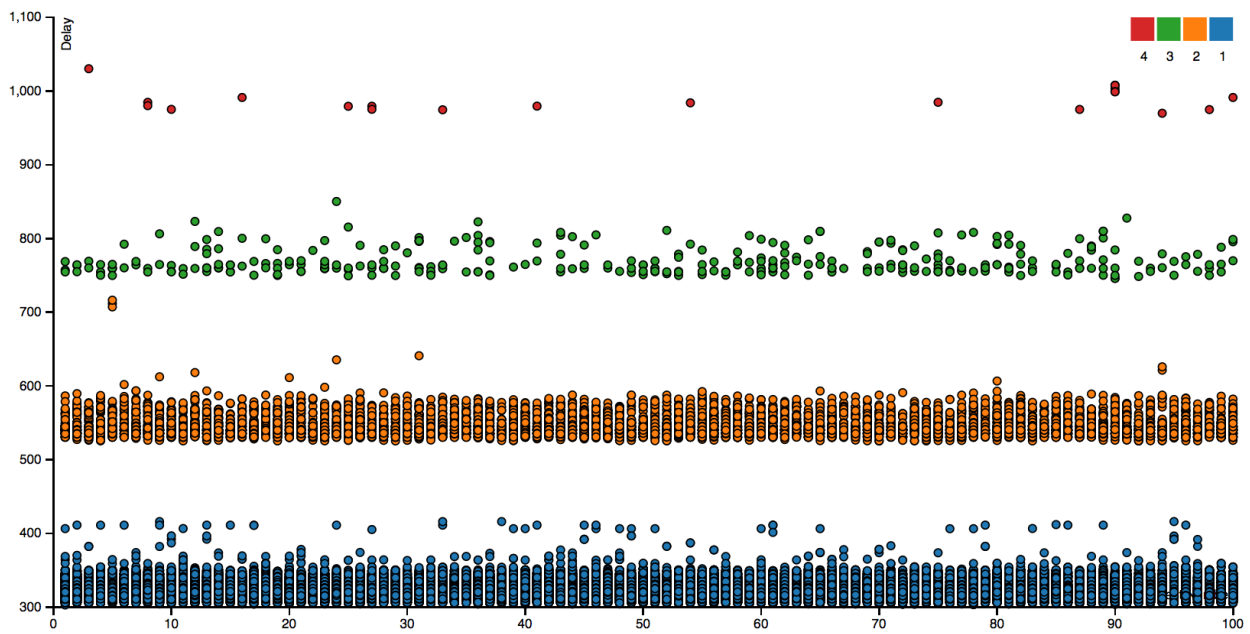


Figure 3: 40 Mbps Lateness Scatterplot

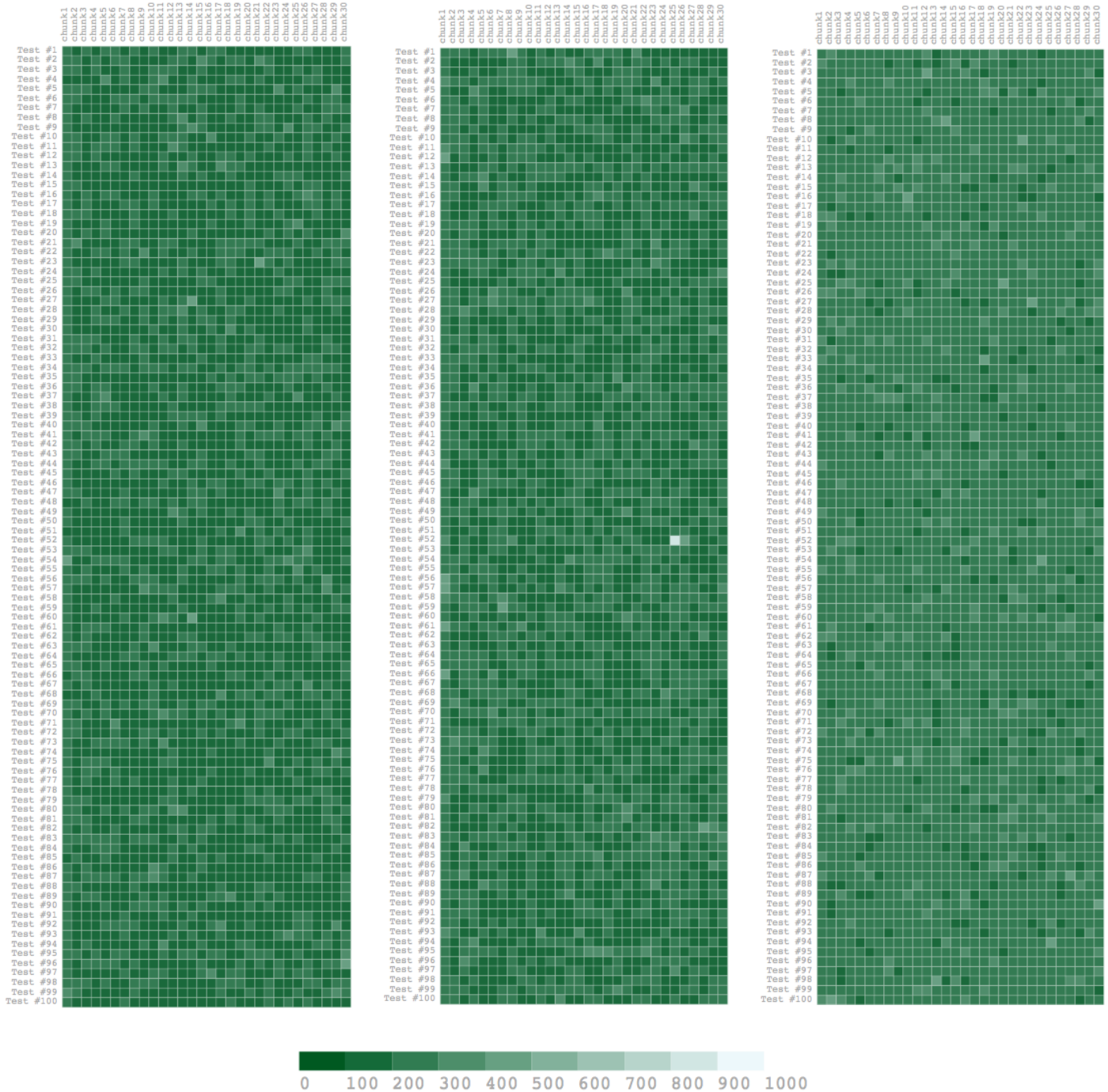


Figure 4: 6, 10, and 40 Mbps Skip Heatmaps (resp.)

IBM Cloud Technical Whitepaper

For the scatterplots (Figures 5, 6, and 7) associated to TCP tests, we plot test number vs. startup delay in milliseconds. The startup delay means the amount of time before the first 1 second bucket, is ready to display. Notice that we can't guarantee continuous playback even from here; this is the minimal time, before we can start playing.

Some immediate take-aways from these scatterplots are:

- The variance in these statistics is extremely high.
- The values themselves are huge, both in comparison to FASP, and for normal use-cases.

As suggested before, delays on the order of 15-120 seconds are totally unacceptable. In any of the cases, we see values well outside the bounds of what we're willing to accept. Furthermore, we don't even see some cases with good performance; all cases are bad.

Conclusion

FASPStream is no more a faster replacement to TCP streaming, than air travel is a faster replacement to American-European driving; you just cannot drive across the Ocean. FASPStream offers minimal playback delays, consistent delivery rates, and high network performance with excellent quality including negligible probability of skipping, opening new game changing possibilities for live and near live streaming at distance. For more information about FASPStream and Aspera Streaming for Video, please visit us at <http://asperasoft.com/software/streaming>.

About Aspera, an IBM Company

Aspera, an IBM company, is the creator of next-generation transport technologies that move the world's data at maximum speed regardless of file size, transfer distance and network conditions. Based on its patented, Emmy® award-winning FASP® protocol, Aspera software fully utilizes existing infrastructures to deliver the fastest, most predictable file-transfer experience. Aspera's core technology delivers unprecedented control over bandwidth, complete security and uncompromising reliability. Organizations across a variety of industries on six continents rely on Aspera software for the business-critical transport of their digital assets.

For more information

For more information on IBM Aspera solutions, please visit <https://www.ibm.com/cloud-computing/products/high-speed-data-transfer/> and follow us on Twitter [@asperasoft](https://twitter.com/asperasoft).



© Copyright IBM Corporation 2018

IBM Corporation
Route 100
Somers, NY 10589

Produced in the United States of America
December 2018

IBM, the IBM logo, ibm.com and Aspera are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both. If these and other IBM trademarked terms are marked on their first occurrence in this information with a trademark symbol (® or ™), these symbols indicate U.S. registered or common law trademarks owned by IBM at the time this information was published. Such trademarks may also be registered or common law trademarks in other countries. A current list of IBM trademarks is available on the Web at “Copyright and trademark information” at:ibm.com/legal/us/en/copytrade.shtml

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

Microsoft, Windows, Windows NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

Other product, company or service names may be trademarks or service marks of others.

This document is current as of the initial date of publication and may be changed by IBM at any time. Not all offerings are available in every country in which IBM operates.

The performance data and client examples cited are presented for illustrative purposes only. Actual performance results may vary depending on the specific configurations and operating conditions. It is the user's responsibility to evaluate and verify the operation of any other products or programs with IBM product and programs. THE INFORMATION IN THIS DOCUMENT IS PROVIDED “AS IS” WITHOUT ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING WITHOUT ANY WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND ANY WARRANTY OR CONDITION OF NON-INFRINGEMENT. IBM products are warranted according to the terms and conditions of the agreements under which they are provided.

¹ The maximum object size supported using the default Aspera software configuration depends on the capabilities of the storage platform. Specific platform limitations are given in Table 1.



Please Recycle