

Investigating Excessive Delays in Mobile Broadband Networks

Natalie Larson¹, Džiugas Baltrūnas², Amund Kvalbein², Amogh Dhamdhere¹, KC Claffy¹, and Ahmed Elmokashfi²

¹CAIDA/UCSD, La Jolla, CA

²Simula Research Laboratory, Oslo, Norway

ABSTRACT

Systematic monitoring of performance characteristics of mobile broadband networks can help users, operators, and regulators understand this increasingly critical infrastructure. We report the results of a measurement study of round trip delays in two mobile networks in Norway using over 200 geographically distributed subscriptions for a period of one month. In general, we find high variation in delay within the same radio access type and RRC state. Across both networks we study, we observe connections with round trip delays of several seconds, often multiple times per hour, that are more frequent when nodes are moving. Correlating these extreme delay events with available metadata, we find they are related to handovers, radio state transitions, and retransmissions at the link and physical layers.

1. INTRODUCTION

The more we rely on cellular mobile broadband (MBB) networks for basic communication, the more important it will become to maintain adequate performance for these networks. Understanding how performance varies across different cellular technologies (e.g., 3G, LTE) will also be important as the industry tries to integrate multiple technologies to facilitate management and convergence as well as optimize performance [14]. One fundamental performance parameter is the packet round trip time (RTT) between a mobile device and the remote end of its communication. Very large or variable RTTs often signify degraded transport and/or application performance. The behavior of this delay parameter is affected by a complex interaction of factors including the radio access protocol (2G, 3G, or LTE), access channel type (dedicated vs. shared) and sub-mode (HSPA, DC-HSPA), channel quality, and traffic level on the radio access network (RAN) and mobile core network (CN).

We report the results of a recent measurement study of end-to-end delay in two operational MBB networks in Norway, in which we quantify and explain observed delay behavior, including extreme delay events. We saw significant variations in delay both across and within connections, but found that RTTs were mainly a function of the type of MBB technology and data channel used. We saw extreme, multi-

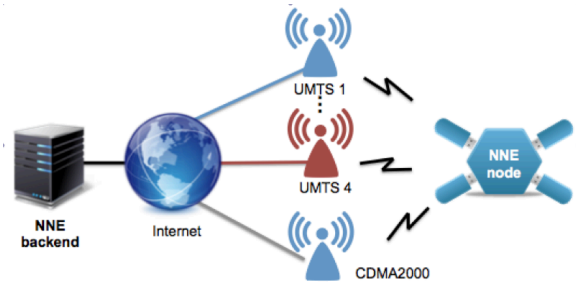


Figure 1: Nornet Edge overview.

second delay episodes, that were more common when our mobile nodes were actually in motion. Correlating these episodes with meta-data captured by our measurement infrastructure indicates that such events are caused by handovers, radio resource state transitions, and heavy retransmissions by physical and link layer protocols.

2. INFRASTRUCTURE AND DATA

2.1 Measurement infrastructure

We used the Nornet Edge (NNE) infrastructure [1] (Fig. 1), a dedicated infrastructure to support measurement of Norwegian MBB networks. NNE consists of several hundred nodes distributed across Norway, and a server-side infrastructure. Measurement nodes are placed indoors in both urban and suburban regions of Norway to reflect the population density of the country. Each node connects with up to four MBB UMTS operator nodes, which are hosted in locations such as schools and government offices, selected to be representative of indoor, stationary users in urban and suburban areas. We also have four nodes on long-distance trains, constantly in motion. Using this diverse, geographically distributed set of nodes allows us to characterize end-to-end delay as a function of radio access and channel type, and to compare performance on stationary and mobile devices.

An NNE node is a custom single-board computer running a standard Linux distribution. Each node is connected with one to four UMTS networks and one CDMA2000 1x EVDO network, using standard subscriptions. For the UMTS networks, connections are through Huawei E353, E3131 3G

USB modems, or E392 LTE USB modems. The former two types of modems are capable of up to HSPA+ (“3.75G”), while the latter is capable of up to LTE (i.e., all standards). In addition to the data service, modems expose a connection’s mode (2G, 3G), signal strength, LAC, and cell id. A daemon on each node manages MBB connections, creating connections on all attached WCDMA modems and retrying connections that break. After successfully establishing a connection, the node creates a PPP tunnel between the user equipment (UE) and the modem, which then establishes a PDP context with the GGSN (P-GW).

NNE nodes perform experimental measurements against backend servers, periodically transferring measurement results to the backend database. The backend servers manage configuration, monitoring, storage, and data post-processing.

2.2 Measurement methodology

To measure end-to-end delay, we send a 20-byte UDP packet to an echo server every second, and record the reply packet. Every request-response pair is logged and later imported into a database. The measurement payload includes a timestamp and an incremental sequence number, which allow duplicate detection and round-trip time calculations. Sending 20 bytes per second does not warrant a dedicated channel (i.e. Cell-DCH), but these connections also transfer periodic maintenance traffic and measurement data, such that our 3G connections spend a considerable fraction of the measurement period on a dedicated channel.

On each node, a measurement script starts automatically on all network interfaces as they become active, and runs as long as the interface is up. When the connection is not available, the measurement script stops sending packets until the connection is active again. We analyze measurements taken throughout November 2014 from the two largest UMTS operators (the other two large operators in Norway do not have their own LTE networks). We do not include results from the CDMA operator, since their CDMA modems provide little connection state and few location attributes.

3. BASIC DELAY STATISTICS

End-to-end delay in MBB is affected by three factors:

1. radio access type (RAT): 2G RTTs are $O(100\text{ms})$, 3G and LTE are closer to $O(10\text{ms})$
2. wireless standard in use, e.g. HSPA, HSPA+, etc.
3. user equipment (UE) radio resource control (RRC) state: (for 3G; LTE only has one state for data transmission). Before transmitting data, the UE attaches to the network and establishes a Packet Data Protocol (PDP) context with the Gateway GPRS Support Node (GGSN). The PDP data structure contains the IP address and other information about the user session. Depending on the traffic pattern, the Radio Network Controller (RNC) then allocates a shared forward access channel

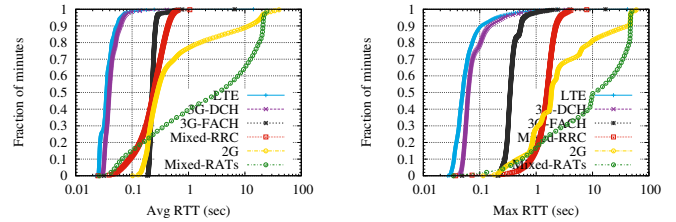


Figure 2: Average and max RTTs for stationary connections on operator A (one-minute bins). Except for bins with mixed-RATs and 2G, RTTs are stable. 3G-DCH and LTE maximum RTTs are double the averages.

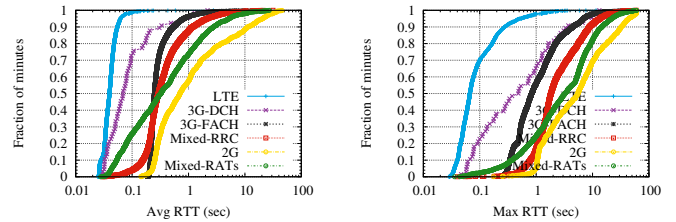


Figure 3: Average and max RTTs for mobile connections on operator A (one-minute bins). 40% of 3G-DCH and 3G-FACH bins exhibit maximum delays ≥ 1 second.

(FACH = $O(100\text{ms})$) or dedicated radio channel (DCH = $O(10\text{ms})$) for the UE.

This section reports typical RTT values, and their variability, for connections using different RATs and radio resource states. For each connection, we divide the measurement duration into 1-minute bins, classify these bins according to RAT and RRC state (legend in Fig. 2), then calculate average and maximum RTTs in each bin. Figure 2 shows significant differences in average and maximum RTTs between different RATs and RRC states. Unsurprisingly, LTE and 3G-DCH exhibit the shortest average and maximum RTTs, with maximum 3G-RCH RTTs slightly higher than maximum LTE RTTs. For 2G connections, average RTTs are less than 500 msec in 70% of bins, but they can be as long as several seconds, consistent with the fact that 2G deployments are often in challenging locations where it is not cost-effective to deploy 3G and LTE.

RRC state transitions affect performance, but bins with inter-RAT changes (handovers) have the highest RTTs, also unsurprisingly since handovers can take seconds, leaving packets buffered until the process finishes. Mean RTTs in bins with RRC state transitions are comparable to FACH (shared channel) bins, suggesting that the presence of FACH data in the bin overshadows delay caused by the state transition.

Maximum RTTs in bins with RRC state transitions, on 2G, and with handovers are more than one second, but are clearly outliers; most packets experience far lower RTTs according to the average RTT distributions. Operator B shows similar results, albeit with small quantitative differences.

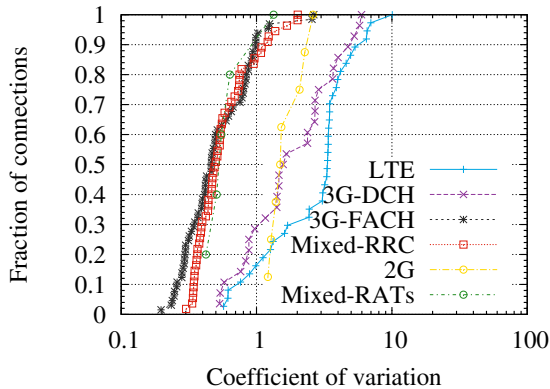


Figure 4: Distribution of maximum RTT coefficient of variation for operator A connections. LTE, 3G-DCH, and 2G connections have highly variable maximum RTTs.

Nodes in motion experienced higher RTTs, likely due to requisite handovers and potential changes in radio access types. Fig. 3 plots the CDFs of average RTTs (left) and maximum RTTs (right) for non-stationary connections on operator A, classified by RAT and radio resource state. We observe three differences between static and mobile connections. First, mobile connections with mixed RATs no longer have the highest RTT profile. This mixed RAT sample for stationary nodes is anomalous, since handovers are not expected; RAT changes must be due to dynamic cell breathing, over-crowded or rural cells, and other network infrastructure problems that are not the focus of our study. Second, 3G-DCH connections exhibit significantly (not just slightly) higher delay than LTE. Third, a large fraction (40%) of 3G-DCH and 3G-FACH bins suffer from maximum delays higher than one second, with approximately one third of LTE bins showing maximum RTTs over 100msec. Perhaps most interesting, however, is the high variability in RTTs we observe for static and moving 3G-DCH and LTE connections, which we investigate next.

Variations in RTT.

To study whether the high variability in RTTs is common across and within connections, we look at each connection separately, calculating the average and standard deviation of maximum RTTs. We include a connection to a given RAT and/or RRC state sample if the connection contributes at least one day of measurements while on that RAT and/or RRC state.

We use the coefficient of variation, which is the outcome of dividing the standard deviation by the mean, to quantify variability. A coefficient of variation less than one indicates low variability, and above one indicates high variability. Fig. 4 shows the CDF of this coefficient for operator A’s connections for different RAT and/or RRC state combinations. Operator B shows similar results; we do not show the graphs due to space limitations. We measure high variability of maximum RTTs for 2G, 3G-DCH and LTE connec-

tions. For instance, the standard deviation is three times the mean for half of operator A’s LTE connections, which is surprisingly unpredictable for stationary measurement nodes! The variability is lower for connections using 3G-FACH or mixed RRC states, consistent with their max-RTT distributions (Fig. 2). Since packets are buffered during RRC state transitions, the maximum RTT in bins with mixed RRC state depends primarily on the length of this transition, which is not highly variable.

FACH (shared access) connections have a much higher baseline, so the magnitude of variability would have to be much higher to yield a large coefficient of variation. For example, if channel quality degradation increases delay by 100msec (i.e. due to link layer retransmissions), delay in a typical FACH connection will increase by $\approx 50\%$ vs. $\approx 200\text{--}300\%$ for 3G-DCH and LTE. Shared (e.g., FACH) connections are more likely to experience congestion delay, but for the two operators we studied we observed much lower packet loss for FACH compared to DCH, which suggests a lack of congestion. One of these operators confirmed privately to us that their FACH channels are over-provisioned. We also found low variations in average RTT, with the exception of 2G, as expected from Fig. 2.

In summary, our measurements show that variations in maximum RTTs while on 3G-DCH and LTE are high even at the same location. However, maximum RTTs for 3G-FACH connections exhibit low variability, which is good for applications that stay mostly on FACH such as M2M. Maximum RTTs can be several seconds especially when nodes are mobile. We also find that average RTTs measured in a window of one-minute duration are stable and predictable, which shows that most observed episodes of increased delays are too brief to affect statistical averages. We take a closer look at these excessive RTTs in the next section.

4. CHARACTERIZING EXCESSIVE DELAYS

To investigate excessive delays, we identify all 3G-DCH one-minute bins with maximum RTTs of at least two seconds. We compare the maximum RTTs to the RTTs immediately before and after. The preceding RTTs are consistent with typical 3G-DCH RTTs, i.e. $O(10\text{ms})$, but the packet with the highest RTT and the following 2-3 packets appear to arrive at the same time. Figure 5 illustrates this RTT behavior, which results in a triangle pattern when plotted against the packet sequence number. We call such episodes “triangle events”. The first four packets do not experience excessive delay, but the following four do. The maximum RTT is 4.3 seconds, and the following RTTs monotonically decrease by one second, eventually dropping to the level before the jump. Note that we send our UDP probes every second, which implies that packets five through nine are buffered and then released simultaneously. We next try to thoroughly characterize triangle events.

Detecting Triangle events

We say that a triangle event has occurred when we see

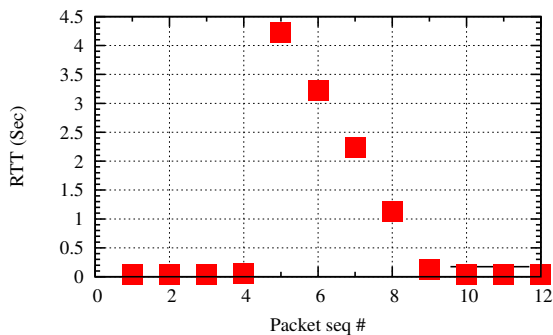


Figure 5: Triangle event illustration.

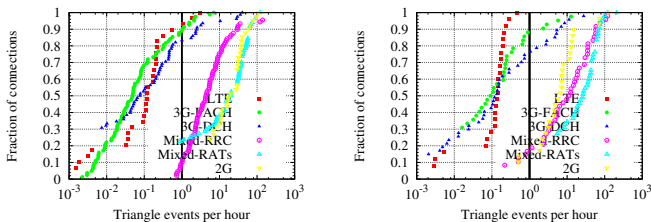


Figure 6: The distribution of the frequency of triangle events for operator A (left) and operator B (right). The solid line marks the fraction of connections that experience more than a single event per hour on average. Approximately one quarter of 3G-DCH connections experience one or more events per hour, on average.

a sequence of at least two packets, the first of which has an RTT of at least two seconds, followed by a sequence of packets with each RTT decreasing 0.8–1.2 seconds, or only increases or decreases slightly (RTTs within 75% of each other) when the RTTs are large (≥ 1 s). We allow for this large-RTT scenario because we occasionally see such RTTs in the middle of otherwise typical triangle event RTT sequences.

In this analysis we discard events that last longer than 60 seconds, involve packet reordering, or involve more than 60 packet losses, i.e. one minute worth of loss. Our goal is to avoid inflating the number of detected triangle events by including periods that bear similarity to triangle events yet exhibit additional evidence of degraded performance.

We detect triangle events that happened for all operator A and operator B connections that contributed at least 10 days of measurements during November 2014. We group events into 5 categories based on RAT and RRC state (when applicable) during the minute in which a triangle event happens: events that occur in bins with RRC state transitions, events in bins with mixed RATs, and events in bins with 2G, 3G-DCH, 3G-FACH and LTE respectively.

Frequency of Triangle Events

Figure 6 plots the CDF of the average number of triangle events per hour for different RAT and RRC state categories for operator A (left) and operator B (right). We calculate av-

Table 1: Triangle event frequency for mobile nodes.

| Category | LTE | 3G-DCH | 3G-FACH |
|-----------------|------|--------|---------|
| Operator A | | | |
| max hourly rate | 1.18 | 32.2 | 13.01 |
| min hourly rate | 0.4 | 4.16 | 5.65 |
| Operator B | | | |
| max hourly rate | 0.72 | 8.51 | 18.51 |
| min hourly rate | 0.25 | 2.3 | 7.39 |

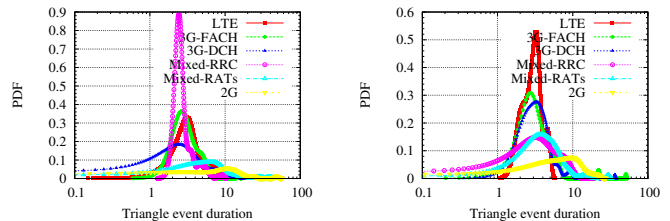


Figure 7: Probability density function of triangle event duration for operators A (left) and B (right).

erages by dividing the total number of events in a category by the total time (normalized to hours) spent in that category. The plot shows that triangle events are frequent during RRC state transitions, RAT changes, and 2G, consistent with the associated packet buffering delays (for the first two) and infrastructure challenges (for 2G). A large fraction of 3G-DCH connections, 30% in operator B and 20% in operator A, experience more than one event per hour. LTE and FACH appear to suffer fewer triangle events.

Table 1 shows the maximum and minimum hourly rates of triangle events for our four mobile nodes for LTE, 3G-FACH, and 3G-DCH. With the exception of LTE, mobile node results are comparable to the worst 5% of stationary connections. Although our mobile sample is small, results for both operators and all RAT/RRC combinations (including those not shown in the table) are consistent: triangle events are more likely when nodes are moving.

Duration of Triangle Events

Figure 7 presents the probability density function of triangle event duration for operators A and B respectively. We observe distinct modes for most connection classes. For both operators, triangle events affecting 3G-FACH, 3G-DCH, and LTE connections generally last between two and four seconds. The two operators have quite different triangle event duration profiles for connections that experience RRC transitions: tightly clustered around 2.5s for operator A, and more evenly dispersed around 3s for operator B. On the other hand, triangle events in bins with RAT changes are longer in operator A than in operator B, with modes of 6.5s and 3.2s, respectively. Triangle duration distributions for mobile nodes are comparable to stationary nodes. These observed differences between operators confirm that our results are not caused by measurement artifacts related to our hardware. Further, they suggest room for improvement in network management, i.e., operators can configure their RRC

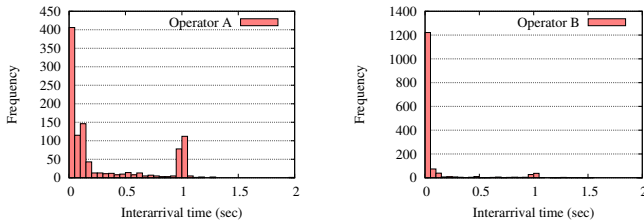


Figure 8: Median packet inter-arrival at server during triangle events for operator A (right) and B (left).

transition and inter-RAT handover times such that they reduce buffering.

Where does the buffering happen?

To test whether triangle events were more common on the uplink vs the downlink we recorded packet arrival times during triangle events at the server over a three-day period. Figure 8 shows histograms of the median packet interarrival times during triangle events for connections on 3G. We find a bimodal distribution: zero and one second between packets that make up triangle events arriving at the server, for both operators, with most packets arriving zero seconds apart. Packets arriving essentially simultaneously at the server signifies that they were buffered somewhere on the uplink and released at once. One second spacing between packet arrivals at the server signifies that the packets were not intercepted on their way to the server, as the client sent them one second apart. Thus, triangle events appear to occur both on the uplink and the downlink, though predominantly on the uplink. User equipment transmits packets on the uplink, whereas a cell tower transmits them on the downlink. Cell towers have greater power and more sophisticated coding schemes, and thus can more easily overcome poor signal quality (E_c/I_o and RSSI), resulting in fewer losses, less buffering, and fewer triangle events on the downlink.

Triangle events and signal quality

While some observed triangle events are a direct consequence of the time required to perform RRC state transitions and RAT switches, triangle events also happen at times during which no transitions occur. To gain insight into what causes such events, we correlate the rate of triangle events seen by each connection with its average Received Signal Strength Indicator (RSSI) and signal-to-noise ratio (E_c/I_o or received energy per code bit / interference level). We find an inverse relationship between signal-to-noise ratio and triangle event frequency. A poor signal-to-noise ratio could induce packet loss and link-layer retransmissions. For example, in the data link layer Automatic Repeat reQuest (ARQ) protocol Selective Repeat, a sender must buffer all packets until they are ACKed. If a series of packets is lost, the sender will maintain a buffer of the un-ACKed packets and release them together, as soon as a packet is ACKed. Such behavior matches the one-second spaced decreasing round trip time

pattern we observed in triangle events.

In summary, triangle events are fairly common when connections are moving, and increase in frequency when the signal is poor. Their duration varies depending on the underlying change. We also spot differences between operators that reveal room for network improvement via tweaks to state transition and handover timers. Further, we observe triangle events on both the uplink and downlink, albeit more frequently on the former.

5. UNDERLYING CAUSES

To dig deeper into the underlying causes of excessive buffering episodes, we conduct a set of controlled experiments that monitor the link-layer and radio activity. We focus on buffering in 3G DCH channels not caused by RAT changes or RRC state transitions, thus, we force our modem to stay on 3G and send large (1000-byte) UDP packets, ensuring that we stay on DCH. Finally, since triangle events tend to occur more frequently when signal strength is poor, we place our equipment in an indoor location with $RSSI \geq -103dBm$ (i.e., the equivalent of one signal bar on most phones). We use QxDM [12] diagnostic software to capture link-level messages and events. To test whether buffering happens on the uplink, downlink, or both, we enable logging of received and transmitted UDP packets on our echo server. We ran the experiment for two days to capture buffering episodes in both networks. Our two main observations are:

1. Buffering episodes on the downlink coincide with cell changes. We observe that triangle events occur both during UE handovers to an already known cell following a procedure called HSDPA soft repoint, and during UE handovers to a newly discovered cell [13]. The latter is a two-step procedure that starts with the access network incorporating the new cell into the list of active cells that the UE can connect with, via an ACTIVE SET UPDATE message, after which an HSDPA soft repoint is performed. In both cases, the UE stops receiving radio link layer protocol (RLC) acknowledgments for outgoing packets, although these packets are received on the server side without buffering. This lack of acknowledgment indicates that the downlink channel is temporarily unavailable. During this temporary downlink unavailability, incoming frames will be buffered if they are sent in Acknowledged Mode, otherwise they will be dropped altogether.
2. Buffering episodes also take place when the UE loses its signaling radio bearer in the midst of the handover. This results in additional negotiation to re-establish the bearer. We find that buffering events are the most severe (longest duration) in these cases.

6. RELATED WORK

Baltrunas et al. [2] use the NNE framework to examine MBB reliability from the perspective of network availability

and packet loss. Our study builds on this work, focusing on extreme packet delays, and includes an analysis of mobile and LTE nodes.

Our work is novel in that it uses a large, geographically distributed infrastructure, dedicated to recording MBB measurements. Two recent MBB studies using dedicated infrastructures include S. Sen et al. [6], who mount laptops on public buses and compare performance across operators, and Z. Koradia et al. [4] who use notebooks to compare the performance of four operators in seven locations in India. Larger studies using crowd-sourced user initiated measurements include Chen et al. [3] who analyze user-perceived latency with QxDM, and Nikraves et al. [5] who use data collected by the apps Mobiperf and Speedometer to study MBB performance over time, comparing differences between operators, access technologies, and regions, ultimately showing that MBB performance has not improved over time.

Our work is also novel in that it is the first to investigate the “triangle delay” pattern. Other studies that have characterized MBB delay in general include J. Huang et al. [7] who observe that LTE has significantly shorter state promotion delays than 3G. Y. Chen et al. [8] find that RTT can vary widely according to geographical location, coverage of the NodeB and distance between the NodeB and RNC. Y. Xu et al. [10] find that large downlink buffers typically deployed in MBB networks can cause high latency when throughput is too low to drain the buffers fast. F. Qian et al. [11] find that RRC state transitions may incur long latencies and that frequent RRC state transitions can cause unacceptably long delays, supporting our findings.

7. DISCUSSION AND FUTURE WORK

We have performed a large scale measurement study of round trip delays in two mobile operators in Norway using both static and moving measurement nodes. Our work highlights the complex interplay between upper layer protocols and radio and link layer operations. We observe high variability in maximum RTTs measured in one minute bins for 3G-DCH and LTE. We also observe that maximum RTTs can reach several seconds, especially when nodes are moving. Investigating these extreme delay episodes, which we term “triangle events,” we find they are caused by excessive buffering on both the uplink and the downlink. Using connection metadata and the QxDM tool, we determine that these episodes occur due to inter-RAT handovers, RRC state transitions, link-layer retransmissions, inter-cell handovers, and reconfiguration of data channels. We believe that eliminating triangle events, particularly during inter-cell handovers, will be crucial in the near future with the increasing deployment of micro and femto cells. Going forward, we plan to investigate the impact of triangle events on TCP performance and user-perceived quality of experience, and ways to mitigate those impacts.

8. REFERENCES

- [1] A. Kvalbein, D. Baltrunas, J. Xiang, K. R. Evensen, A. Elmokashfi, S. Ferlin-Oliveira, The Nornet Edge platform for mobile broadband measurements, *Computer Networks*, 2014.
- [2] D. Baltrunas, A. Elmokashfi, A. Kvalbein, Measuring the Reliability of Mobile Broadband Networks, *IMC* 2014.
- [3] Q. Chen, H. Luo, S. Rosen, Z. Mao, K. Iyer, J. Hui, K. Sontineni, K. Lau, QoE Doctor: Diagnosing Mobile App QoE with Automated UI Control and Cross-layer Analysis, *IMC* 2014.
- [4] Z. Koradia, G. Mannava, A. Raman, G. Aggarwal, V. Ribeiro, A. Seth, S. Ardon, A. Mahanti, S. Triukose, First Impressions on the State of Cellular Data Connectivity in India, *ACM DEV*, 2013.
- [5] A. Nikraves, D. R. Choffnes, E. Katz-Bassett, Z. M. Mao, M. Welsh, Mobile Network Performance from User Devices: A Longitudinal, Multidimensional Analysis, *PAM*, 2014.
- [6] S. Sen, J. Yoon, J. Hare, J. Ormont, S. Banerjee, Can they hear me now?: A case for a client-assisted approach to monitoring wide-area wireless networks, *IMC*, 2011.
- [7] J. Huang, F. Qian, Y. Guo, Y. Zhou, Q. Xu, Z. M. Mao, S. Sen, O. Spatscheck, An in-depth study of LTE: effect of network protocol and application behavior on performance, *CCR*, 2013.
- [8] Y. Chen, N. Duffield, P. Haffner, W. Hsu, G. Jacobson, Y. Jin, S. Sen, S. Venkataraman, Z. Zhang, Understanding the complexity of 3G UMTS network performance, *Networking*, 2013.
- [9] J. Manweiler, S. Agarwal, M. Zhang, R. Choudhury, P. Bahl, Switchboard: a matchmaking system for multiplayer mobile games, *MobiSys*, 2011.
- [10] Y. Xu, Z. Wang, W. K. Leong, B. Leong, An End-to-End Measurement Study of Modern Cellular Data Networks, *PAM* 2014.
- [11] F. Qian, Z. Wang, A. Gerber, Z. M. Mao, S. Sen, O. Spatscheck, Characterizing Radio Resource Allocation for 3G Networks, *IMC*, 2010.
- [12] Qualcomm, QxDM Professional Proven Diagnostic Tool for Evaluating Handset and Network Performance. 2012.
- [13] Chris Johnson, *Radio Access Networks for UMTS: Principles and Practice*, John Wiley & Sons, West Sussex, England 2008.
- [14] 4G Americas’s Recommendations on 5G Requirements and Solutions, October 2014.