

Use Cases for Communicating End-Points in Mobile Network Middleboxes

Humberto J. La Roche
Email: hlaroche [at] cisco.com
Position Paper for IAB MaRNEW Workshop (2015)
September 2015

Abstract

Gi Middleboxes are deployed in mobile networks on the Gi reference point which is where the 3GPP mobile network meets the Internet. They are essential in performing functions that regulate and manage traffic for better QoE but their role has not been devoid of controversy. The ongoing usage of end-to-end encryption significantly impacts the function of currently deployed middleboxes. For example, application detection via identification of signatures in traffic no longer works on encrypted traffic. Media trans-rating and transparent content caching fail as well. Substrate Protocol for User Datagram (SPUD) is a set of tools that enables endpoints (server, device) to safely communicate information about communications with each other and with middleboxes, without compromising privacy of content and under control of an opt-in framework. This position paper discusses the role of cooperative information sharing frameworks, such as SPUD, in facilitating the deployment of middleboxes for their more lofty traffic management and flow regulation purposes while neutralizing some of the more controversial functions. We argue that cooperation in the eco-system of content publishers, mobile network operators, and subscribers will result in a better mobile Internet experience and ultimately result in a radio-friendly transport protocol that may supersede TCP.

Introduction

This position paper presents a viewpoint on addressing the architectural impact of data plane encryption on 3GPP mobile networks. The mechanism is one where end-points (mobile device and origin servers or their proxies) exchange relevant and pertinent information about the path and the ~~application~~ application with middleboxes (routers, base stations, packet core elements, optimization elements, and other in-path elements). The framework discussed is cooperative in nature and based on an opt-in model controlled by the 3GPP policy layer. We don't expound on incentives for cooperation. These can be monetary (e.g. a network user such as a subscriber or content publisher buys a service) or barter-based (a user permits or allows a capability in exchange for a benefit).

At the time of writing, over 50% of mobile network traffic is encrypted on a per session basis using TLS or SSL [1]. The percentage by volume is expected to rise to near 100% with the deployment of HTTP/2 [2] consistent with the IAB recommendations [3] and the increasing ease of obtaining X.509 certificates [4]. The proper functioning of Gi middleboxes is compromised by encryption. Smith [5] and also Moriarty and Morton [6] have documented the impact in Internet Drafts. Most agree encryption is desirable and necessary to protect the privacy of user traffic as well Internet revenue models such as advertising that rely on content being delivered intact. Our position is that middleboxes play a useful role because they are deployed to address the harshness of the 3GPP radio environment. Radio bandwidth is neither plentiful nor cheap, and variability in the form of fades, shadowing, user mobility, handovers, etc. can constrict the channel at a moment's notice. The use of encryption hides information which is used by network operators to provide network management and control functions that benefit overall quality of experience.

Gi Middleboxes are deployed to optimize delivery of packets to subscribers. Some of these trans-rate or compress the content reducing the volume of bytes delivered over the air interface. Others offer

TCP traffic optimization, attempting to repair TCP in radio access (more later). Others still will accelerate the user experience by caching content closer to the subscriber, reducing RTT and WAN congestion. Per session encryption via SSL or TLS will render caching middleboxes inoperable unless a content publisher has explicitly contracted with a CDN operator for PKI services. This type of contract is one example of how content publishers collaborate with middlebox operators to improve the mobile Internet. Other types of middleboxes will offer HTTP enriched headers that can be consumed by advertising customers. The deep packet inspection (DPI) middlebox is frequently deployed to enable path management functions and fair use policies, throttling traffic based on the identified application type. DPI is a controversial technology as it is seen as enabling attacks on privacy.

Depending on where an actor sits in the mobile Internet value chains, some middlebox practices will be controversial. Subscribers may not like fair use policies. A middlebox altering advertisements in any way, or the insertion/replacement of ads, will be objectionable to the online advertiser. A content publisher will not approve of an image compression function that distorts a photograph. Defining a “user” as either a subscriber or a content publisher, it seems reasonable to insist on the basic principle of informed user consent for opt-in or opt-out to a specific capability of the network. The user willingly manifests their disposition to trade something tangible with an entity in the network path for a particular benefit. A good example of this commerce is a subscriber that explicitly opts-in to use a data compression proxy. The benefit would be an improved mobile user experience if the image distortion or video impairments are not excessive and if byte count is reduced. If a user consents, with full disclosure, to a practice, then we will assume the practice is not objectionable. The practice of informed subscriber consent is well developed in the Internet eco-system. Social network operators, commerce sites, search engines, email services, all have privacy policies that must be agreed to by the user.

Here is a summary of how we develop the position. We begin by noting the “ossification” effect that middleboxes are responsible for and then review issues with TCP performance in the radio and attempts to address; our view being that there must be a better suited transport protocol for the radio than TCP. An area of interest for us is coupling the transport layer protocol to the increasingly sophisticated LTE link layer, emphasizing the argument that others have also made[7], that consideration of cross-layer couplings are essential when considering the radio environment. SPUD is an experimental toolkit introduced by Trammel and Hildebrand [8] that can support the mechanisms for cross-layer coupling and we provide a brief overview of it. The use cases for SPUD in mobile networks follow. We close with a summary of the position expressed. Our hope is to provide a reasonably complete set of references. The author would be pleased to have omissions called to his attention.

Ossification, TCP Protocol Evolution, and Radio Access

Gi middleboxes and network middleboxes actually do hinder protocol evolution through a phenomenon known as “ossification” meaning something like “impedance to move except by breakage”. See [9-12] for references on the topic. In the radio environment, the venerable TCP is simply not getting along with the 3GPP radio environment. To address, there must be transport protocol evolution to achieve performance improvements that maximize the efficiency (ability to carry useful bits) of the radio channel. The crux of the issue is that the very middleboxes that are deployed to compensate for poor TCP and radio interaction are preventing new protocols from being introduced (the “ossification” effect).

There are well understood reasons for why TCP and 3GPP mobile radio are at odds with each other [13-15]:

- (1) Reno and Reno-derived TCPs interpret packet loss and delay as the consequence of a buffer overflow event at a location in the path where there is a throughput bottleneck. On the radio,

loss and delay is most frequently the consequence of radio effects such random bit errors (which cause the radio to re-transmit) and delay spikes (caused by retransmissions). Retransmissions can also be at the LTE link layer to correct for packet losses and can also cause delay spikes in the TCP. These delay spikes cause spurious RTO time-outs that are not related to packet loss but force the TCP sender to back-off.

- (2) LTE handovers (X2-based) may result in packet re-ordering (triggers duplicate ACKs) or delay spikes which will also affect the TCP layer.
- (3) Techniques designed to improve battery life will conflict with TCP [16].
- (4) Network asymmetric bandwidth will delay uplink packets causing TCPs round-trip estimate to be larger than it should.

It is worthwhile to briefly review what the industry has been doing about TCP performance and radio. End-points in the Internet implicitly or explicitly expose data about the path to the transport layer. In TCP (Reno and earlier versions) a congestion state is implicitly indicated to the sending end-point through duplicate ACKs or increased round-trip delay between a segment sent and its corresponding ACK. Instead of packet loss (or delay), the network can set Explicit Congestion Notification bits [17], if the sending TCP/IP stack and receiving TCP/IP stack support ECN. Historically, configuring ECN has involved configuring active queue management (AQM), and the difficulty of configuring it has proven to be one of the main obstacles in ECN deployment yet at least one vendor of popular mobile handsets is optimistic [18]. Also with ECN, base stations implementing it would mark “congestion experienced” based on actual radio access state. This base station feature allows TCP to become “radio congestion” aware. Additional developments include the Re-ECN [19] scheme which builds on ECN by adding additional network control in middleboxes in the form of “policers” and “droppers”. Building on ECN, the Congestion Exposure (CoNex) IETF working group [20] is actively pursuing mechanisms for including data about the path in TCP options fields as well as in IPv6. Most recently, and specific to 3GPP radio, a mechanism where an LTE eNB provides “throughput guidance” to origin server endpoints, has been presented at the recent IETF 93 plenary [21] and is accompanied by at least two Internet Drafts [22, 23] with coauthors from Google, Nokia, and Vodafone. Finally, there are new mechanisms to reduce the configuration complexity of active queue management such CoDEL [24] and PIE [25].

Another initiative to improve transport in a radio environment is based on the idea of using Channel Quality Indications (CQI) and Discontinuous Transmission (DTX). The approach, called CQIC, from Lu *et al* [7] derives advantage of the fact that consistent with 3GPP, UEs indicate to the base station the preferred modulation and coding to be used by the sender on the base station. This information can in principle, be offered to the transport protocol stack and when delivered to the sender, used to optimize its performance. The challenge is twofold: (1) that the UEs must expose information about the link layer and currently they don’t and (2) there is no natural protocol framework to convey useful information about the channel to the sender.

On the protocol front, Google’s innovative QUIC [26] is proposed to replace TCP in web applications with a new reliable transport protocol layered over UDP optimized for HTTP/2. QUIC appears at first blush to do little to address the fundamental issues associated with TCPs poor performance in the radio. It does, however, represent an improvement over TCP. For example, QUIC supports F-RTO which helps with spurious retransmission time-outs [27, 28] in radio. QUIC can support a variety of congestion control algorithms and so in this sense, allows for experimentation to determine which particular congestion control is suitable to the prevalent radio channel model. However, we think that the framework afforded by QUIC is too restrictive and does not provide enough flexibility to address the comprehensive set of issues associated to the radio which are mostly addressed with cross-layer couplings. We need a general replacement for TCP that works well with radio.

SPUD

Substrate Protocol for User Datagrams (SPUD) [8] permits cooperative communications between devices on the network path and endpoints so they can share relevant information about the end-to-end conversation. SPUD is carried over UDP and adds new fields allowing the path to declare information that can be consumed by the middlebox or conversely, allowing the middlebox to declare information that can be consumed by the path. The information shared is in binary format and can be any structured data using Concise Binary Object Representation (CBOR) [29]. CBOR can be thought of as a binary form of JSON and has the advantage of consuming few bits and is easy to process in lightweight code

The higher motivation of SPUD is to contribute to transport protocol evolution. To achieve the goal, SPUD inserts itself between API-accessible UDP transports and the new or legacy transport protocol. No kernel modifications in operating systems are needed allowing the SPUD stack to live in user-space. A consequence is that experimentation with real UEs is possible without having to go to the operating system for kernel changes.

It is well worth emphasizing that the benefit of SPUD relies on both middleboxes and end-points implementing sound procedures to both generate information and consume it. Only the interfaces (basically CBOR formatting rules) need to be standardized. Our position is that implementations are an area where innovation will be able to provide differential value among technology suppliers.

Use Cases

The basic idea of use cases is that relevant information related to the communications, either about the path itself or about the application, can be included in SPUD. The objective is to improve the quality of experience for users that opt-in to the cooperative environment. In some cases, it makes sense to store and manage the consent state for the user (a content publisher or an individual subscribers), in the 3GPP policy layer. The 3GPP policy layer is comprised of the PCRF and the SPR [30, 31] and can enforce the opt-in status. For example, the policy layer can use the information to apply a policy personalized to the user (an example is parental controls). More generally, 3GPP subscriber policy is being introduced today to manage subscriber steering in a mobile SDN. The interface “Sp” in the figure provides visibility into the subscriber session record to the PCRF. The PCRF and then use a policy interface such as the Gx or Sd to communicate policy treatment to the user plane. Use of the Sd is optional.

The figure illustrates network scenarios for middlebox communications. A 3GPP PCRF (Policy and Charging Rules Function) and SPR (Subscriber Policy Repository) provide policy services to users of the network. The eNB (base station) supports the radio interface towards the UE and the PDN-GW provides anchoring services by advertising the IP address of the UE to the Internet via the Gi reference point. The entity enabling the information communications (SPUD support) is shown as an orange rectangle. In the diagrams, the middlebox can be a content-server or a proxy for a content server. An example of one such middlebox is a Data Compression Proxy (DCP) such as Google’s “Flywheel” [32]. But, it would need to be modified to support SPUD semantics. Yet another example is a CDN edge node. Another middlebox example can be a virtualized Gi middlebox implementation that include all capabilities of multiple service nodes essentially “in a box” deployed in the cloud [33]. With the advent of NFV and SDN [34], it is possible to build such a cloud-based middle-box. These implementations use the Network Service Header (NSH) [35] to create service function chains internal to the data center. Classification of traffic into service chains may be performed by Flexible Mobile Service Steering (FMSS) [36].

The middlebox is assumed, for convenience, to be IP addressable.

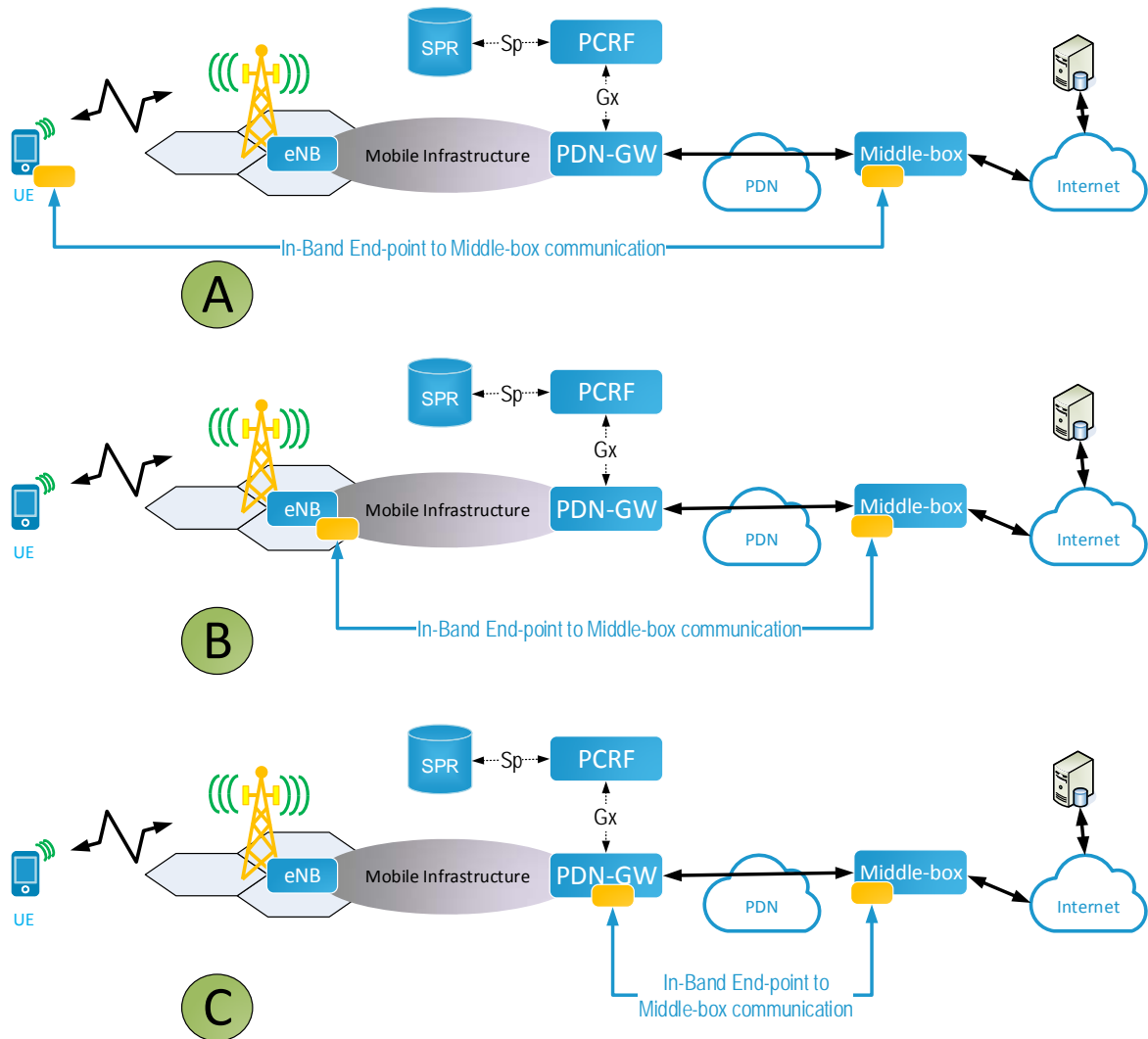


Figure 1: Communicating end-points and middleboxes

Note the scenarios can be deployed concurrently.

- (A) This scenario allows the UE application to exchange information with the middle-box. The UE uses a sockets or sockets-like library which applications can use for networking needs. Some applications will invoke SPUD middlebox communications services whereas others will not. The architecture shown allows the UE to have the opportunity to explicitly declare application-level information which the middlebox can consume and use.
- (B) This scenario allows the eNB, as an in-path element, to declare information about the path such as instantaneous radio interface state, which can be consumed by the middlebox to help it operate optimally. Use cases that use this network architecture require RAN features that are not normally found in existing systems. They can be implemented with the proposed “Mobile Edge Computing” initiative which allows the RAN environment to be highly programmable [31].

- (C) In this scenario, the PDN-GW can communicate with the origin server or the middlebox that proxies it. The importance of this use case is that the PDN-GW has access to the subscriber session record which comprises information about the device type, the access network type, and the mobile identity. In some cases, location information can also be available at the PDN-GW.

We provide a summary of the use cases below. We hope to generate community interest in developing the use cases.

- **RAN-aware TCP Optimization Service.** In this service that can be offered by default or on an opt-in basis, the eNB reports instantaneous radio network state to either a content publisher or a proxy for one (such as a DCP). The proxy offers SPUD-based access *and* conventional TCP towards the network side. The eNB, which knows radio network state and can make near-instantaneous recommendations to a supported sender (again, either a content publisher or proxy for it) on the packet transmission rate that would not congest the RAN. This service is intended to optimize delivery of TCP services in the radio environment, repairing it with explicit information from the RAN that TCP currently lacks. The approach uses SPUD overhead and conceptually, is not that different from ConEx [20] and “Throughput Guidance” [21] approaches which instead of SPUD, use IP or TCP optional fields. Because of the extensibility associated to SPUD, relevant information associated to send window size, sender packet spacing, or power conservation state can be made available to the sender or its proxy. In this sense, the lower layers of LTE can be efficiently cross-coupled with the transport protocol for optimal performance without the inconvenience of mining for fields in IP or TCP or assuming a specific model for the radio channel to determine the congestion control algorithm

RAN information can be derived from the device. It is possible to envision a scheme such as CQIC [7] described earlier providing CQI and DTX information to a sender origin server or middlebox proxy for one. In that case, the information might be embedded as SPUD fields. Interesting to point out: this particular use case involving only exchange about the path between UE and origin server, requires no network support and can be implemented as a differentiated mobile device feature when operating with specific content servers.

- **Advanced Content Rating.** Several operators offer subscriber products that are charged for and rated differently because someone else is paying for the delivery (for example, a content publisher). An example is an all-you-can-consume social networking plan whereby anything used *from* the social network app is charged differently than other traffic. The operator would normally have difficulty distinguishing between a web page accessed from within the app and the same web page launched from outside the app. Since the idea is that traffic from the UE app is rated in a special way, it would mark relevant traffic with the appropriate SPUD header, identifying it as coming from a specific application for which zero rating applies. The 3GPP policy framework would be used to verify enrolment of the user in a specific rating plan, and the content could be directed to the appropriate content servers.
- **Application Traffic Management and Control.** When traffic is encrypted, the transport layer is not able to discern the relative importance of different traffic flows and unhelpfully, not knowing any better, will treat flows as if they have equal priority. To address the issue, a cooperating, properly incented end-point can willingly enrich the SPUD layer with information to tell a network middlebox the network characteristics needed by the application. For example, an email application may would signal it needs best effort, an interactive audio application would signal it needs low delay and low jitter, and a synchronization or backup application might signal it needs less than best effort.. No DPI is

needed and the actual conversation between the end-points can be shrouded by encryption. The 3GPP policy framework would verify the SPUD signaling, and the requested priority allowed or denied.

- **Interworking with middle-box implementation that uses SDN.** Here, we assume that the middlebox is in fact deployed in the cloud and has internal structure built with multiple service nodes deployed as part of multiple service chains. The specifics of the mapping of external flows to encapsulated NSH flows will vary from vendor to vendor but in most cases will include a role for the PCRF and SPR. If the external flow consumed by the virtualized Gi middlebox includes SPUD information, these could be used in: (a) guiding the policy steering decision into a service chain, and (b) used for consumption in any intermediate virtualized service node behaving as a middle-box. Similarly, the virtualized Gi middlebox may publish information in SPUD that can be consumed by an end-point or other path element. The “bump-in-the-wire” middleboxes, which are not IP addressable, can be accommodated in this use case through NSH and SPUD interworking.
- **Bandwidth Management of Adaptive Streaming Video.** Most video delivered today is actually over HTTP (uses HTTP as a substrate protocol) and is adaptive in nature meaning the sender probes for the amount of bandwidth available on the channel, by monitoring RTTs for example, and sends at the highest possible rate, choosing from a collection of segments that are encoded at different bit rates. This is referred to as HTTP Adaptive Streaming (HAS). The difficulty with HAS is that it is greedy. As detailed in [37], greed generates undesired oscillations in streaming bit-rate when clients compete for bandwidth in a bottle-neck link such as the radio environment. The solution is for the eNB to tell the servers, or a proxy for them (such as an optimization middle-box), how much throughput each HAS can consume.
- **Parental Controls** is a very important use case that deserves attention. It is intended to filter inappropriate content for subscribers (minors, for example) who are opted-in to the capability. Challenges with this use case include: (1) need to lock-down the terminal to prevent user from changing options to by-pass the filter and (2) the filter is required to be available in both Wi-Fi and cellular and cellular access. The parental controls filter is a middlebox with an IP address reachable on the Internet. The filtering function would be implemented using criteria derived from a database of named destination IP addresses that are “blacklisted” because they host inappropriate content. The client app in the locked-down device, communicates with the filter middlebox using its IP address. The problem solved by SPUD, is that when traffic is encrypted, the destination URL is not visible to the filter middlebox. The app would pass at least two parameters to the filter as SPUD fields. One would be the destination IP address of the web request (which is different from the IP address of the filter) and two is a parameter akin to the “safe” preference proposed in [38]. The filtering function can then be accomplished. Any app with the ability to provide browsing would need to support the SPUD communications framework to be allowed into the locked-in environment. Other IP apps deemed safe to use such as messaging can be deployed unchanged. If the subscriber consents, a shortened version of the requested URL can be passed to the filter as SPUD overhead potentially allowing for a URL blocking feature (as opposed to a feature that blocks destination IPs). No secret key exchange is necessary.

Summary of the Position

Several schemes that support the exchange of radio congestion information are already in active discussion in the industry as a way of improving TCP performance [17, 21, 39] in the radio. We can anticipate middleboxes will continue to be deployed to adapt content delivery the radio

environment. We are proposing the idea that informed subscriber consent will need to be a continued guiding principle for information exposure. There are cases where middleboxes should have access to data about the channel or the application only when such information is willingly disclosed, typically as part of an exchange in which a user, explicitly opts-in after weighing-in advantages vs. disadvantages. We note the Internet eco-system is replete with cooperative barter as a framework for transactions: individuals allow commerce sites to collect purchase history in exchange for recommendations, opt-into free email in exchange for yielding privacy by allowing indexing of key words to the user account, and search engines provide a valuable capability in exchange for the ability to create targeted advertising. We think it is possible to manage the opt-in using the 3GPP policy framework. The solution for opt-in may extend beyond 3GPP access. Work is in progress to extend 3GPP policy to other access technologies [40]. And what are the consequences of not opting-in? Simply that the user will not see better quality of experience when using the mobile Internet.

To address transport protocol evolution, and specifically for a more radio-friendly approach than TCP can currently provide, we are advocating for an approach that freezes TCP and creates a “substrate” protocol over UDP that allows for fast experimentation associated to addressing a broad set of use cases. In particular, and implicit in the position we have taken, is that cross-layer interactions across the mobile access stack are required to ensure the efficiency of a future reliable transport protocol replacing TCP. We are skeptical that mining TCP and IP for options fields will result in the desired efficiency improvements because of the obstacles placed by the referenced “ossification” effects and posit that, like QUIC, the new transport protocol should be layered on UDP. However, we think that more flexibility than provided for by QUIC is needed to address the challenge attaining a good transport protocol for the radio. In particular, we think that rather than characterizing the radio channel to select among a family of transport protocol congestion control algorithms, it is preferable to near instantaneously, provide “guidance” to the sender. SPUD appears to be a reasonable way of implementing the needed cross-layer couplings.

Neither middleboxes nor exposure of path and application information are inherently evil. The mobile network itself is a cascade of middleboxes exchanging information about subscriber mobility. There is long-standing consensus in the industry that data exchange around path congestion improves user experience. What we propose is to consider SPUD as a way of realizing these use cases. We do believe middleboxes and the privacy supported by encrypted traffic can co-exist and further, with “substrate” frameworks such as SPUD, can facilitate an accelerated approach, via experimentation, to a better transport protocol evolution with better alignment to the radio environment.

References

- [1] D. Naylor, A. Finamore, I. Leontiadis, Y. Grunenberger, M. Mellia, M. Munafò, *et al.*, "The Cost of the “S” in HTTPS," presented at the ACM CoNEXT 2014, Sydney, Australia.
- [2] Internet Engineering Task Force (IETF). (2015). *RFC7540, Hypertext Transfer Protocol Version 2 (HTTP/2)*. Available: <https://tools.ietf.org/html/rfc7540>
- [3] Internet Architecture Board (IAB). *IAB Statement on Internet Confidentiality*. Available: <https://www.iab.org/2014/11/14/iab-statement-on-internet-confidentiality/>
- [4] Let's Encrypt. (2015). *Let's Encrypt is a new Certificate Authority: It's free, automated, and open*. Available: <https://letsencrypt.org/>
- [5] K. Smith, "Network management of encrypted traffic," *Internet Draft, Work in progress*, 2015.
- [6] K. Moriarty and A. Morton, "Effect of Ubiquitous Encryption," *Internet Draft, Work in progress*, 2015.

- [7] F. Lu, H. Du, A. Jain, G. M. Voelker, A. C. Snoeren, and A. Terzis, "CQIC: Revisiting Cross-Layer Congestion Control for Cellular Networks," presented at the Proceedings of the 16th International Workshop on Mobile Computing Systems and Applications, Santa Fe, New Mexico, USA, 2015.
- [8] J. Hildebrand and B. Trammell, "Substrate Protocol for User Datagrams (SPUD) Prototype," *Internet Draft, Work in progress*, 2015.
- [9] B. Trammell and J. Hildebrand, "Evolving Transport in the Internet," *Internet Computing, IEEE*, vol. 18, pp. 60-64, 2014.
- [10] B. Hesmans, F. Duchene, C. Paasch, G. Detal, and O. Bonaventure, "Are TCP extensions middlebox-proof?," presented at the Proceedings of the 2013 workshop on Hot topics in middleboxes and network function virtualization, Santa Barbara, California, USA, 2013.
- [11] M. Honda, Y. Nishida, C. Raiciu, A. Greenhalgh, M. Handley, and H. Tokuda, "Is it still possible to extend TCP?," presented at the Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference, Berlin, Germany, 2011.
- [12] Internet Architecture Board (IAB). (2015). *IAB Workshop on Stack Evolution in a Middlebox Internet (SEMI)*. Available: <https://www.iab.org/activities/workshops/semi/>
- [13] H. Balakrishnan, S. Seshan, E. Amir, and R. H. Katz, "Improving TCP/IP Performance over Wireless Networks," in *1st ACM Int'l Conf. on Mobile Computing and Networking (Mobicom)*, 1995.
- [14] H. Balakrishnan, V. N. Padmanabhan, S. Seshan, and R. H. Katz, "A comparison of mechanisms for improving TCP performance over wireless links," *Networking, IEEE/ACM Transactions on*, vol. 5, pp. 756-769, 1997.
- [15] H. Balakrishnan and V. N. Padmanabhan, "How network asymmetry affects TCP," *Comm. Mag.*, vol. 39, pp. 60-67, 2001.
- [16] J. Erman, V. Gopalakrishnan, R. Jana, and K. K. Ramakrishnan, "Towards a SPDY'ier mobile web?," presented at the Proceedings of the ninth ACM conference on Emerging networking experiments and technologies, Santa Barbara, California, USA, 2013.
- [17] K. Ramakrishnan, S. Floyd, and D. Black. (2001). *RFC 3168, The Addition of Explicit Congestion Notification (ECN) to IP*.
- [18] P. Lakhera and S. Cheshire, "Your App and Next Generation Networks," in *Apple Worldwide Developers Conference*, San Francisco, 2015. Available: <https://developer.apple.com/videos/wwdc/2015/?id=719>
- [19] B. Briscoe, A. Jacquet, C. D. Cairano-Gilfedder, A. Salvatori, A. Soppera, and M. Koyabe, "Policing congestion response in an internetwork using re-feedback," *SIGCOMM Comput. Commun. Rev.*, vol. 35, pp. 277-288, 2005.
- [20] Internet Engineering Task Force (IETF). *Congestion Exposure (conex)*. Available: <http://datatracker.ietf.org/wg/conex/charter/>
- [21] P. Szilágyi and A. Terzis, "Mobile Content Delivery Optimization based on Throughput Guidance," in *IETF 93*, Prague, 2015. Available: <https://www.ietf.org/proceedings/93/slides/slides-93-icrg-3.pdf>
- [22] A. Jain, A. Terzis, N. Sprecher, S. Arunachalam, K. Smith, and G. Klas, "Requirements and reference architecture for Mobile Throughput Guidance Exposure," *Internet Draft, Work in progress*, 2015 2015.
- [23] A. Jain, A. Terzis, H. Flinck, N. Sprecher, S. Arunachalam, and K. Smith, "Mobile Throughput Guidance Inband Signaling Protocol," *Internet Draft, Work in progress*, 2015.
- [24] K. Nichols and V. Jacobson, "Controlling Queue Delay," *Queue*, vol. 10, pp. 20-34, 2012.

- [25] R. Pan, P. Natarajan, C. Piglione, M. Prabhu, V. Subramanian, F. Baker, *et al.*, "PIE: A Lightweight Control Scheme To Address the Bufferbloat Problem," *Internet Draft, Work in progress*, 2013.
- [26] R. Hamilton, J. Iyengar, I. Swett, and A. Wilk, "QUIC: A UDP-Based Secure and Reliable Transport for HTTP/2," *Internet Draft, Work in progress*, 2015.
- [27] P. Sarolahti, M. Kojo, K. Yamamoto, and M. Hata, "RFC 5682, Forward RTO-Recovery (F-RTO): An Algorithm for Detecting Spurious Retransmission Timeouts with TCP," *Internet Engineering Task Force (IETF)*, 2009.
- [28] P. Sarolahti, M. Kojo, and K. Raatikainen, "F-RTO: an enhanced recovery algorithm for TCP retransmission timeouts," *SIGCOMM Comput. Commun. Rev.*, vol. 33, pp. 51-63, 2003.
- [29] C. Bormann and P. Hoffman. (2013). *RFC 7049, Concise Binary Object Representation (CBOR)*. Available: <http://www.rfc-editor.org/info/rfc7049>
- [30] 3rd Generation Partnership Project (3GPP). *TS 23.203, Policy and charging control architecture*. Available: <http://www.3gpp.org/DynaReport/23203.htm>
- [31] European Telecommunications Standards Institute. (2014). *Mobile-Edge Computing* [White Paper]. Available: <http://www.etsi.org/technologies-clusters/technologies/mobile-edge-computing>
- [32] V. Agababov, M. Buettner, V. Chudnovsky, M. Cogan, B. Greenstein, S. McDaniel, *et al.*, "Flywheel: Google's Data Compression Proxy for the Mobile Web," presented at the Proceedings of the 12th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2015.
- [33] W. Haeffner, J. Napper, M. Stiernerling, D. Lopez, and J. Uttaro, "Service Function Chaining Use Cases in Mobile Networks," *Internet Draft, Work in progress*.
- [34] European Telecommunications Standards Institute. Network Functions Virtualization White Papers. Available: https://portal.etsi.org/nfv/nfv_white_paper.pdf, https://portal.etsi.org/NFV/NFV_White_Paper2.pdf, https://portal.etsi.org/Portals/0/TBpages/NFV/Docs/NFV_White_Paper3.pdf
- [35] P. Quinn, J. Guichard, S. Kumar, M. Smith, W. Henderickx, T. Nadeau, *et al.*, "Network Service Header," *Internet Draft, Work in progress*, 2015.
- [36] 3rd Generation Partnership Project (3GPP). (2014). *TR 22.808, Study on Flexible Mobile Service Steering (FMSS)*. Available: <http://www.3gpp.org/DynaReport/22808.htm>
- [37] L. Zhi, Z. Xiaoqing, J. Gahm, P. Rong, H. Hao, A. C. Begen, *et al.*, "Probe and Adapt: Rate Adaptation for HTTP Video Streaming At Scale," *IEEE Journal on Selected Areas in Communications*, vol. 32, pp. 719-733, 2014.
- [38] M. Nottingham, "The "safe" HTTP Preference," *Internet Draft, Work in progress*, 2015.
- [39] B. Briscoe, R. Woundy, and A. Cooper, "RFC 6789: Congestion Exposure (ConEx) Concepts and Use Cases," 2012.
- [40] Broadband Forum. *TR-300, Policy Convergence for Next Generation Fixed and 3GPP Wireless Networks*. Available: <https://www.broadband-forum.org/technical/download/TR-300.pdf>