



CAP on Mobility Control for 4G LTE Networks

Yuanjie Li¹ Zengwen Yuan¹ Chunyi Peng² Songwu Lu¹

¹University of California, Los Angeles ²The Ohio State University
{yuanjie.li,zyuan,slu}@cs.ucla.edu, chunyi@cse.ohio-state.edu

Abstract

The CAP theorem [1] exposes the fundamental tradeoffs among three key properties of strong consistency, availability and partition tolerance in distributed networked systems. In this position paper, we take the CAP perspective on 4G mobility control. We view the control-plane management for mobility support as a distributed signaling system. We show that the impossibility result of the CAP theorem also holds for mobility control: It is impossible for *any* mobility control to guarantee sequential consistency, high service availability, and partition tolerance simultaneously. Unfortunately, the current 4G system adopts its mobility scheme with the notion of sequential consistency. Our empirical study further confirms that, the incurred data unavailability (i.e., data service suspension) time is comparable to that induced by wireless connectivity setup. We argue that the desirable mobility control for the upcoming 5G networks should take a paradigm shift. We discuss our early effort on re-examining the consistency notion for higher availability and fault tolerance.

1 Introduction

Mobility support is critical to offer seamless data services to billions of smartphone users. To date, the 4G¹ cellular network is the only large-scale infrastructure which offers wide-area mobility management in practice. From the system perspective, an ideal mobility solution should offer all three properties of high availability, correctness, and fault tolerance. In the LTE context, it should offer always-available data access, and enforce correct packet forwarding under user-specific policies (e.g. radio access control, billing, QoS), and handle failures with elegant degradation.

In recent years, there have been extensive efforts on offering highly available wireless connectivity [2, 3] and more flexible infrastructure upgrade (e.g. NFV). In this work, we look into the problem on how to improve service availability during user mobility in the LTE network. Towards this end, our first observation is that, the current mobility management is complex and time con-

suming. On one hand, it boasts as the only large-scale networked system, in par with the wired Internet, that offers wide-area mobility support. On the other hand, this essential utility is achieved via rounds of control-plane operations, which span on multiple network nodes (the base station, gateway, mobility controllers, etc.) and the end user device with a variety of cellular protocols involved. Upon mobility, multiple micro-level control procedures are triggered, such as the radio connectivity setup, routing path update, data billing and QoS policies configuration, to name a few. Fundamentally, the control-plane mobility support is a distributed signaling system.

In this position paper, we offer a fresh view on 4G mobility control from the CAP [1] perspective. Given the distributed signaling system, we show that a variant of the CAP theorem holds for mobility control. It is impossible to ensure all three properties of sequential consistency, availability and partition tolerance in a generic setting. Sequential consistency states that, in the distributed system, there exists a total order on operations (update/migration) for the device's control states, such that every operation looks as if it were performed by a single network entity. Unfortunately, we observe that the 4G mobility scheme adopt the sequential consistency notion. As a result, it cannot achieve high availability for data access during mobility. Our empirical study confirms that, the sequential consistency model incurs non-negligible data suspension comparable to that required by wireless connectivity setup. It contributes to 29.8%–35.6% of the total data suspension on average, and up to 69.6% in the worst case. This leads to 193.8 ms to 6.4 s extra data suspension in mobility.

We further argue that a paradigm shift on mobility design is needed for the upcoming 5G network in order to better balance the inherent tradeoffs among availability, consistency, and partition tolerance. The extensive experiences from the Internet service show that, high availability is typically favored over strict notion of consistency upon failures. Our preliminary study indicates that, the sequential consistency notion adopted by 4G mobility seems *more than necessary* for mobile data services. It is thus possible to design highly available, partition-tolerant mobile data service with relaxed, yet user-acceptable consistency requirement.

The rest of the paper is organized as follows. §2 reviews the control-plane procedures for 4G mobility support. §3 elaborates on the three desirable properties of availability, consistency and partition tolerance, proves the impossibility result of the CAP Theorem for mobility, and explains its implications for 4G mobility. §4 users experimental traces to assess the negative impact. §5 describes our early thinking on relaxing the consistency model to ensure high availability. §6 discusses the related effort, and §7 concludes this position work.

¹We use 4G and LTE interchangeably for a slight abuse of notions.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.
HotWireless'16, October 03-07, 2016, New York City, NY, USA
DOI: <http://dx.doi.org/10.1145/2973750.2980120>.

2 Background on 4G Mobility

The 4G mobility control subsystem is fundamentally a distributed signaling system. It spans multiple infrastructure nodes, including base stations, gateways, mobility controllers² and user-profile database servers³. Upon user mobility, several micro-level control procedures are triggered, e.g., the radio connectivity setup, routing path update, data billing, and quality-of-service policies configuration, to name a few. To enforce correct data forwarding during mobility, on the control plane, both the device and the involved network nodes should maintain *proper control states* for device location, radio access control list, billing status, QoS profile, *etc.*. Whenever user mobility is observed, these control-plane states should be migrated to new serving network nodes, and updated based on the latest device location and possibly new policies.

Figure 1 illustrates the detailed process (adopted from [4, 5]). Initially, the device stays within location domain 1⁴ and has available mobile data service (by retaining data sessions from the device to the gateway). Upon moving to a new domain, the device and the network perform a series of control procedures and resume the data service thereafter. Data service is suspended in between. Specifically, the phone first establishes the radio connectivity (P1), and reports its arrival to the mobility controller in the new location domain (P2). The mobility controller then transfers the data session states from the old domain's controller (P3), sequentially updates routing path with the base station and gateways and configures billing and quality-of-service policies (P4). Next it updates the user profile database with the mobile device's latest location information (P5). Last, the controller notifies the success of location update (P6), and resumes the device's mobile data service.

3 CAP Perspective on 4G Mobility

3.1 Desirable Properties

Appropriate mobility control subsystem should offer *always-on* data access to mobile users even upon failures. It thus should possess all three properties from the standpoint of distributed systems:

- **Availability:** Data service should be always available to mobile devices. As long as the route is physically available, every packet from/to a mobile device should be eventually delivered. Note that this definition does not stipulate on how fast the packet is delivered.
- **Consistency:** The control states at all involved nodes should remain consistent at all times, given that each involved node may take different operations on the states simultaneously or asynchronously. Intuitively, from the global network perspective, there exists a total order on operations (update/migration) for the session states, such that every operation looks as if it were performed by a single network node.
- **Partition (Fault) tolerance:** Upon network partition on the control plane, the network should still provide data services to the device as long as the data route is still physically connected. Note that the network nodes may not communicate among one another on the control plane upon partitions. In this work, we focus on fail-stop control-plane failures during mobility. These failures occur for various reasons, such as broken control channels, controller crash, incompatibilities between heterogeneous network nodes, *etc.* The

²e.g., Mobility management entity (MME) in 4G LTE.

³e.g., Home subscriber server (HSS) in 4G LTE.

⁴For scalability, 4G-LTE splits its the infrastructure into multiple location domains, each of which is managed by a mobility controller and includes multiple base stations and gateways.

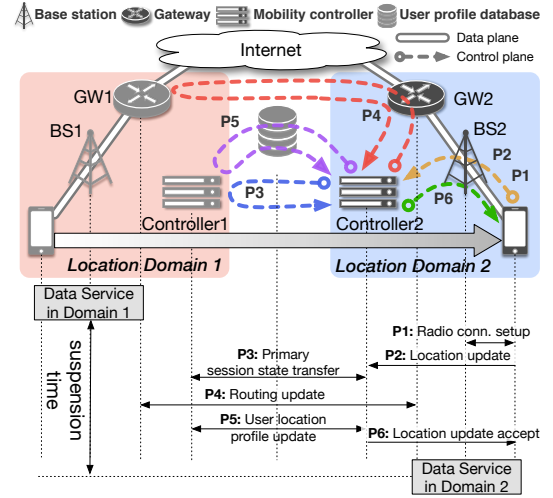


Figure 1: Sequence of procedures in 4G mobility control.

data-plane partitions are not considered since we assume always-connected data path.

3.2 CAP Impossibility Result for Mobility

In this section, we prove that, *any* mobility control (not limited to 4G) cannot guarantee all three properties of sequential consistency, availability and partition tolerance simultaneously.

Our result follows the arguments from the CAP theorem [1, 6] but adapts to the mobility case. We consider the scenarios where the physical data path from the device to the Internet is *always connected* (i.e., connected radio link between device and base station, and connected link between base station and gateway). The device-experienced data service availability is thus determined by the mobility control procedures.

To start with, we model the mobile infrastructure as a graph $G(V, E)$, where each node $v \in V$ is a network element (radio base station, gateway, mobility controller or user profile database), and each undirected edge $(u, v) \in E$ implies u and v can communicate. Each base station connects to at least a gateway, which further connects to the external Internet. Each base station and gateway connect to a mobility controller. Each mobility controller manages a set of base stations and gateways (called a location domain). All network nodes belong to the control plane, while base stations and gateways also belong to data plane (i.e., responsible for user data forwarding).

During the mobility, the end device interacts with the network to retain its data service. Each user device is associated with a set of base station, gateway and mobility controller. A *data session* is maintained from the device to the external Internet (through the associated network nodes). Upon mobility, the data session should be migrated from the old serving network nodes to new ones via a series of control procedures.

The main impossibility result is described as follows.

THEOREM 1 (CAP FOR MOBILITY CONTROL). *Given any mobility control, it is impossible to always guarantee all three properties of sequential consistency, availability, and partition tolerance concurrently.* \square

Proof: We prove the above result by contradictions in a concrete scenario. Consider the setting of Figure 1. A mobile device that was served by a set of network nodes (BS1, GW1 and Controller1) roams to a new location domain (BS2, GW2 and Controller2). However, the new domain is partitioned from the old one, so that

the session state cannot be migrated from those network nodes in the old domain. Note that the old nodes may still update their session states (e.g. QoS). In presence of the partition, no matter how long the new network nodes wait for, it is impossible for them to differentiate the following two cases: (1) The old network nodes have updated the one session state to x_1 ; (2) The old network nodes have updated the state to x_2 ($\neq x_1$). Therefore, they cannot determine whether to forward the user's data packet under the constraint of session state x_1 or x_2 . Two possible outcomes will follow afterwards. The data-plane nodes (base station and gateway) either eventually forward data packet under a non-guaranteed session state (thus violating sequential consistency), or refuse to forward any data (thus violating availability). \square

3.3 On 4G Mobility

Now we elaborate on results on 4G mobility control scheme. Given the mobility subsystem described in §2, we can readily arrive at the following conclusion: 4G mobility control design follows the *sequential consistency* model, because it mandates the sequential execution of micro-level control procedures in mobility. Therefore, it adopts the strong consistency notion among all control state update and migration during mobility.

Indeed, the choice of sequential mobility management has valid rationale. As a distributed system, the 4G control functions must be performed correctly before offering/resuming the mobile data service. Otherwise, the network may not correctly forward data to the right location, under/over-bill the user's data usage, and/or even become vulnerable to security attacks. By mandating sequential consistency, correct mobile data service is indeed enforced. Specifically, two techniques are used by LTE to enforce the above consistency model (see Figure 1). First, when fetching the user states, only the copy from the old mobility controller is migrated to the new one (P3). This eliminates conflicts among multiple copies. Second, the update operations between different nodes are *sequential* by design (P4–P6). However, such practice is more than necessary. We next analyze how the sequential consistency negatively affects 4G's mobile data service availability.

4 Empirical Validation

4.1 The Cost of Sequential Consistency

While sequential consistency does ensure correct control-plane state updates during mobility, it also comes with hidden costs. Notably, it reduces data service availability according to the CAP theorem. Specifically, in 4G LTE it incurs two main aspects of inefficiency, both of which prolongs the mobile data access suspension.

Issue 1: Prolonged network unavailability. To enforce strong consistency, the data service has to be suspended until all network nodes reach consensus on user states. However, not all nodes are involved in data forwarding. For example, in Figure 1, the synchronization on user locations between two mobility controllers and the user-profile database server requires 3 round trips. But none of them is on the data-forwarding path. Even though the data-forwarding path has been updated and available earlier, the data service is still suspended. Moreover, since location updates are sequential, any delay during the step-by-step procedures would block both uplink/downlink data services.

Issue 2: Lower degree of fault tolerance. The data service suspension at the device would be further lengthened when control-plane failures occur. Migrating the *single-copy* of the primary device's states suffers from the single point of failure. Once the state migration fails, the location update procedure would be blocked,

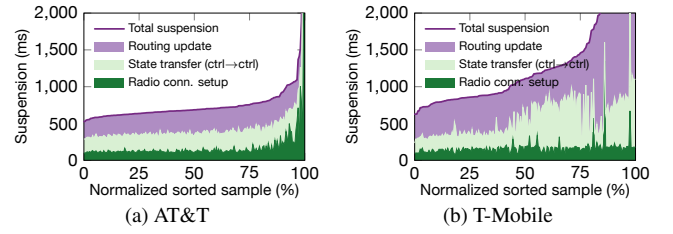


Figure 2: The breakdown of the 4G data suspension time in mobility.

even though replicas of states are still available at other network nodes and the device. In practice, recovering from such control-plane failures can take tens of seconds, during which the device loses data access. In response to failures, the LTE device would either be rejected by network for data service, or experience timeout and re-attach to the network.

4.2 Empirical Assessment

We next use the operational 4G LTE traces to validate and quantify the service unavailability (suspension) caused by sequential consistency. We performed a preliminary user study under in a 8-month period (07/31/15 – 04/22/16). It involves 6 phones (seven Android/iOS phone models, including Samsung Galaxy S4/S5, Huawei Nexus 6P, LG Optimus 2, LG Tribute and Apple iPhone 6 Plus) and five operators (China Mobile, AT&T, T-Mobile, Sprint and Google Project Fi [7]). For each mobile device, we collected 4G RRC/NAS signaling messages using MobileInsight [8]. We performed the breakdown analysis for the traces from commercial 4G network, and quantify the elapsed time spent on each control procedure (Figure 1). Since some control procedures are not visible to end device, we correlate them with the device-side traces based on cellular standards [4, 5, 9].

4G LTE suspension time in mobility. Figure 2 shows the breakdown of samples (sorted by total suspension latency) for 4G mobility control in T-Mobile and AT&T (no failures). Other operators have similar results. It verifies that, the sequential control incurs long mobile data suspension. The suspension due to session state migration fluctuates from 193.8 ms to 6.4 s. The average percentile ranges between 29.8% to 35.6% among operators, and the maximum percentile is 69.6%. Note that these suspensions are typically longer than the radio connectivity setup, which implies that *sequential consistency for mobility control may become the main bottleneck of the mobile data service availability*.

4G LTE control-plane failures in mobility. We further quantify the impact of control-plane failures during location updates. Note that such failures block the user-state consensus, and they lengthen user-experienced data suspensions. Using the collected traces, we compute the failure probability and quantify its impact on data suspension. Given each location update, we identify failures from both the explicit failure messages (e.g. location update reject) and the timeout events (e.g. timeout for location update request). They serve as indicators for node failures (e.g. mobility controller shutdown) and link failures (e.g. user state transfer failures). We exclude failures due to other reasons (e.g. roaming to an unauthorized network). For each location update with failures, we divide the total unavailable time into two parts: the elapsed time before failure message/timeout, and the remaining part. The longer the first part is, the greater impact the failure has on availability.

Figure 3a plots the data unavailable time in the presence of failures, and Figure 3b presents the failure probability. Two observations can be made. First, control-plane failures during location

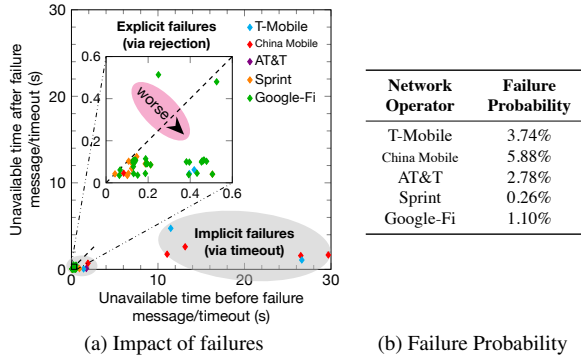


Figure 3: Control-plane failure features in LTE mobility.

updates do occur, and the maximum probability is 5.88% (China Mobile) among the five carriers. Second, whenever a failure is observed, much longer data suspension can be expected than the normal case without failures. The effect is particularly noticeable when atomic consistency is enforced. Upon failures, we also observe following two types of actions taken by the network.

◦ *Explicit failure notification*: The network may explicitly reject the device with proper causes (e.g. “implicit detach”, “sync failure”). In this case, the device re-initiates the registration process (*re-attach*) to the network. It does not require user-state transfer from the previous location domain. Our study shows that, this accounts for up to 500 ms extra latency in data unavailability.

◦ *Implicit failure notification*: The network may not respond to the device. This causes much longer service unavailability. The device has to wait until timeout (15 s by default [5, 10]) before re-initiating its registration. In our experiment, we observed up to 30 s latency in data service unavailability.

5 Consistency-Availability Tradeoff in Mobility

The CAP theorem unveils the fundamental tradeoffs among consistency and availability in the presence of partition failures. In 4G mobile networks, the service availability of data access is typically favored by the end user. Therefore, we believe that relaxing the notion of consistency while improving high availability is a worthwhile direction to pursue. Fortunately, we further reason that, for LTE mobility, sequential consistency on user states is not always necessary to ensure correctness. Instead, the consistency level can be relaxed according to the given user state, while still retaining correctness. This may lead to a novel control context-aware consistency (CCAC) model for mobility. While we have yet completed our effort, the early results is promising. In the following, we offer some preliminary results on a few critical control states, while study on other control states is ongoing.

User location for data forwarding. To correctly deliver packets to roaming users, only a *subset* of network nodes should have updated view on user locations.

• *Uplink data*: To deliver packets originated from the device, only the serving base station (that the device currently connects to) and forwarding gateways need to learn the latest user location. More importantly, the device-originated packets can implicitly notify those data-plane nodes about the new user location. The data-plane packets thus simultaneously update the control-plane location state. There is no need to suspend data delivery while the location state is being updated.

• *Downlink data*: To forward packets destined to the device, the old gateways that relay packets to the device should know the latest user location. Otherwise, packets continue to be forwarded to nodes in the previous location domain before they learn the user location.

Radio access control. In 4G networks, base stations and location domains do not authorize radio access to certain users (e.g. roaming or private Femtocell users). This is done through the per-user radio access control list (RACL), which is kept at both the device (in SIM card) and the network (in base station, user profile data base and mobility controller). To grant radio access, the base station and MME have to wait until the user-specific RACL is available when roaming to a new location domain, thus blocking data service and hurting availability.

We observe that LTE’s radio access control is *group based*. It turns out that, each LTE RACL mandates a group of users with same access rights (e.g. “users that can access a private Femto-cell”, “users that can access a specific location domain”, “all AT&T roaming users”), but a device may belong to multiple groups (e.g. it can access multiple roaming networks). The groups a device belongs to remain largely unchanged, determined by the subscription profile and data plan of the device. Rather than knowing the complete RACLs, the base station and mobility controller only need to verify whether a device belongs to a given group. This does not necessarily require user-specific state on RACL, and can be performed before reaching consistency.

Note that our proposal performs no worse than the current LTE practice⁵. The device still keeps the per-group RACL. The update of RACL takes no longer than the current LTE network. To reduce the prolonged data unavailability, verification of RACL can be possibly done with the radio connectivity setup. This way, no extra round-trip latency is incurred.

Data billing. In the 4G network, each device is billed based on its data usage volume and per-flow charging rules (via the traffic flow template [11]). Each packet is thus counted and mapped to the proper flow for (possibly differentiated) charging. In the current LTE design, charging only starts after the forwarding gateway retrieves the user-specific charging rules.

Given the LTE practice, we can also relax the consistency model for the user state on billing. Note that, packets can be counted without user-specific rules, and the charging rules are determined by the data plan, thus having limited choices. As long as the user-specific charging rules are installed, data accounting can be *correct* by eventually mapping the data usage to the proper charging policy. Similar to radio access control, when the adopted flow templates are of small number, the forwarding gateway can learn the billing group the device belongs to. This way, data billing can be done *in parallel with* synchronizing user states.

QoS. The LTE network offers multiple classes of QoS. Proper QoS is enforced by two metrics. The *QoS class* metric specifies the pre-defined QoS profile, in terms of packet delay/loss/priority. Guaranteed data rate can also be offered to a given user. For each device, its user-specific QoS profile has both metrics. In the current LTE design, the forwarding gateway has to install all user-specific QoS profiles before starting data service.

The class-based QoS metric can also be inferred from the associated group. This can be done before reaching complete consistency on the user state, similar to radio access control. If QoS class is updated in mobility⁶, the network may unnecessarily reserve more

⁵Most operators choose to per-location domain and per-roaming network access control only, to reduce overhead.

⁶Updating the QoS profile for a user is uncommon in practice. We

resource than the device deserves without consistency. For user-specific guaranteed/maximum data rates, we argue that temporary rate mismatch is acceptable as long as they are *no larger than* the user-specific ones.

Given the above relaxed consistency model, we believe that we can reduce the data unavailability/suspension time and handle control-plane failures. Specifically, we are exploring several techniques along this direction. For example, to reduce data unavailability during mobility, we attempt to accelerate the data-path update and activation upon location updates when the device crosses location domains. Given the relaxed consistency model for the location state, we may speed up data forwarding with other control-plane state updates. Second, we plan to ensure consensus among control states at involved nodes. To tolerate control-plane failures, it exploits user-state replicas that are readily available yet underutilized in the current LTE network.

6 Related Work

Improving the service availability in mobile network has been actively studied in recent years. Major research focuses on optimizing radio resource control [2], improving wireless coverage [3], and proposing fast handoff in mobility [12]. This work complements to these efforts by introducing an alternative dimension. Our study is inspired by the seminal work on the CAP theorem [1, 6, 13], but extends the static graph-based model with the introduction of device mobility. We note that, [6] also discusses applying the CAP theorem to wireless networks in the future but focuses on the radio-link failures. In contrast, we mainly examine the core network side.

7 Conclusion

The Internet access is going mobile. With the explosive growth of smartphones and other wearable devices, it becomes increasingly important to offer non-disruptive, always-on network service to the mobile users. Service availability thus turns a premier objective for mobile data access. Unfortunately, there is no free lunch. The CAP theorem states that, in a generic distributed system, it is impossible to always ensure all three properties of availability, consistency, and partition tolerance concurrently. In this work, we show that the current 4G network also fits into this category. The current 4G mobility control adopts the sequential consistency model, and it is impossible to guarantee sequential consistency, availability and partition-tolerance simultaneously. As a result, the current 4G mobility solution also incurs long suspension for data access upon roaming events. Our preliminary empirical study also confirms this finding.

While the impossibility result seems to produce the negative results, it also sheds lights on what to do next. The bottom line is that, like all distributed systems, the LTE network needs to carefully balance the tradeoffs among consistency, availability, and fault tolerance. While the mobile networking community has made extensive efforts on improving the first/last-hop wireless access, issues related to the LTE core network have been overlooked to certain extent. As LTE evolves to the upcoming 5G technology, we believe renovations on both fronts are needed. As we explore more flexible consistency notions to mobility, we may significantly improve the availability and fault tolerance while slightly relaxing the consistency requirement. To this end, we argue that the upcoming

have always observed identical QoS profiles for devices in the same carrier.

5G design should make the paradigm shift away from the legacy sequential consistency. The work reported in this position paper, despite with many open issues to be addressed, represents our early effort towards this direction. Our study hopefully may stimulate more community interests on this important topic for the mobile Internet.

8 References

- [1] S. Gilbert and N. Lynch. Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services. *ACM SIGACT News*, 33(2):51–59, 2002.
- [2] J. Huang, F. Qian, Z. M. Mao, S. Sen, and O. Spatscheck. Radiophet: Intelligent radio resource deallocation for cellular networks. In *Passive and Active Measurement*, pages 1–11. Springer, 2014.
- [3] K. Sundaresan, M. Y. Arslan, S. Singh, S. Rangarajan, and S. V. Krishnamurthy. Fluidnet: a flexible cloud-based radio access network for small cells. *IEEE/ACM Transactions on Networking*, 24(2):915–928, 2016.
- [4] 3GPP. TS23.401: General Packet Radio Service enhancements for Evolved Universal Terrestrial Radio Access Network (E-UTRAN) access, Dec. 2015.
- [5] 3GPP. TS24.301: Non-Access-Stratum for EPS, Jun. 2013.
- [6] S. Gilbert and N. A. Lynch. Perspectives on the cap theorem. *Computer*, 45:30–36, 2012.
- [7] Google. Project Fi. <https://fi.google.com/>.
- [8] Mobileinsight. http://metro.cs.ucla.edu/mobile_insight.
- [9] 3GPP. TS36.331: Radio Resource Control, Mar. 2015.
- [10] 3GPP. TS24.008: Core network protocols; Stage 3, Jun. 2014.
- [11] 3GPP. TS23.125: Overall High Level Functionality and Architecture Impacts of Flow Based Charging, Mar 2006.
- [12] I. Ramani and S. Savage. Syncscan: practical fast handoff for 802.11 infrastructure networks. In *INFOCOM 2005. 24th Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings IEEE*, volume 1, pages 675–684. IEEE, 2005.
- [13] A. Panda, C. Scott, A. Ghodsi, T. Koponen, and S. Shenker. Cap for networks. In *Proceedings of the second ACM SIGCOMM workshop on Hot topics in software defined networking (HotSDN)*, pages 91–96. ACM, 2013.