

# 5GCoreLite: Scalable and Resource Efficient Next Generation Cellular Packet Core (Extended Abstract - NSDI 2019)

Vasudevan Nagendra  
Stony Brook University  
vnagendra@cs.stonybrook.edu

Arani Bhattacharya  
Stony Brook University  
arbhattachar@cs.stonybrook.edu

Anshul Gandhi  
Stony Brook University  
anshul@cs.stonybrook.edu

Samir R Das  
Stony Brook University  
samir@cs.stonybrook.edu

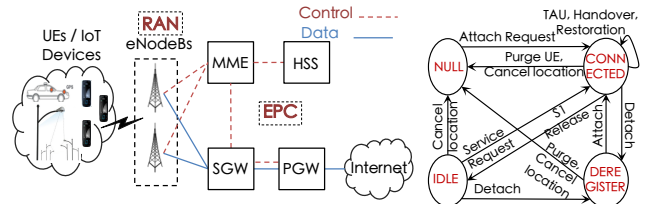
## 1 Problem statement

One of the grand challenges in the design of future cellular core network (LTE/5G) is its *resource-efficient scaling* with the projected growth of signaling or control traffic. Much of this growth is expected to come from the tremendous rise in IoT devices ( $\approx 12$  billion by 2022 [1]). Compared to traditional smartphones, IoT devices generate at least twice the volume of control messages, growing 50% faster than data traffic [2]. This represents a significant overhead as control messages do not directly contribute to the service provider’s revenue. Moreover, the traffic characteristics and performance requirements of cellular-based IoT devices have much greater diversity than traditional user equipments (UEs) like smartphones or laptops [3]. Efficiently managing resources in the presence of this diverse traffic is challenging.

An immediate concern now is the scalability and efficient resource utilization in the cellular core network (also called *Evolved Packet Core* or EPC in connection with LTE networks as shown in Figure 1). Designing an efficient and scalable EPC for 5G requires addressing at least the following key challenges:

*Elasticity:* IoT applications create bursty traffic [3, 4], necessitating dynamic capacity provisioning. Insufficient capacity at any of the EPC core elements (i.e., Mobility Management Entity (MME), Serving Gateway (SGW), Packet Gateway (PGW), and Home Subscriber Server (HSS)) may lead to connection failures and rejections, triggering retry messages and further increasing the load on the EPC. Worse, UEs and all entities inside the EPC maintain stateful contextual information (*static bindings*), making it difficult to migrate connections to other EPC components in case of scale-out or scale-in. Not surprisingly, the current practice is to simply over-provision the EPC, resulting in an expensive and wasteful design [5].

*Customizability:* IoT devices can have very different control and data traffic characteristics and performance requirements [3, 6]. For example, IoT devices in smart cars require stringent Service Level Objectives or SLOs to react to changing traffic conditions, while smart home IoT devices may



(a) LTE architecture with key components. (b) MME state machine.  
Figure 1: Overview of LTE architecture and MME states.

simply require IP connectivity. Unfortunately, today’s cellular networks make use of monolithic EPC devices which are rigid and do not offer any functional or performance flexibility.

*Scalability:* A key bottleneck for large-scale networks is the centralized load balancing mechanism that must immediately assign incoming connections to MME and other EPC components such as HSS, SGW, and PGW nodes. Given the heterogeneity in the entire ecosystem, traditional approaches such as round-robin or least number of connections is no longer effective. While recent approaches based on consistent hashing (CH) distribute connections uniformly [7], they are unable to quickly scale resources in response to bursty IoT traffic, making them vulnerable to “hot spots”, where a few MME hosts are overloaded. Meeting user-specified SLOs while being scalable and resource efficient thus requires a careful reconsideration of load balancing decisions in the network.

## 2 Proposed Solution Overview

To address the above challenges, we propose 5GCoreLite, an agile EPC architecture that exploits recent advances in NFV. We propose to build 5GCoreLite that is *stateless* and *functionally decomposed* design. The statelessness is achieved by externalizing each UE-specific state in shared memory inside each of the EPC nodes, thus decoupling the EPC from the UE contextual information. This stateless design enables dynamic provisioning of EPC nodes responsive to traffic changes, without incurring the overhead of state migration.

To address customizability, we *decompose* the EPC functionality into a set of microservices (or NFs) based on the specific control and data plane procedures they handle, such as attach, service, handover, migrate, packet inspection, billing,

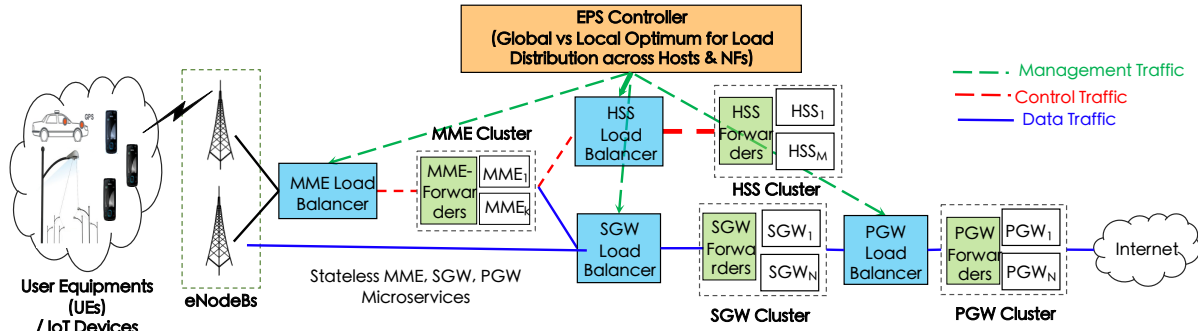
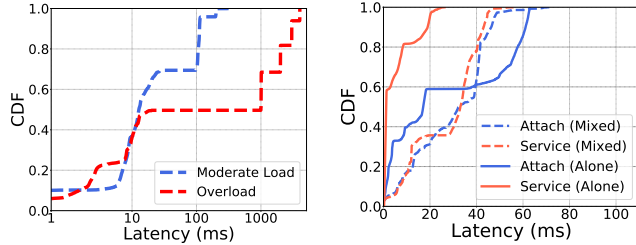


Figure 2: 5GCoreLite: Incrementally deployable, Resource Efficient Scalable Next Generation Cellular Core.



(a) Delays in data transfer instantiation during MME/SGW overloads. (b) Interference among LTE control procedures at line rate traffic.

Figure 3: Experiments demonstrating the limitations of existing control and data plane design with DPDK-based industrial-grade prototype.

etc. This control procedure and function-specific decomposition, facilitated by our microservices design, allows us to cater to specific functional and SLO requirements of individual UEs in a resource-efficient manner. This is in contrast to existing protocol-based decomposition approaches [8] that allow flexibility but fail to provide fine-grained (UE-specific) SLO control.

To address elasticity and scalability, we introduce a multi-level, SLO-aware load balancer and forwarder architecture that optimizes the resource utilization within and across each of the EPC elements by calculating the local and global optimizations for handling each of the connections from UEs. Unlike existing approaches that aim to balance connections across EPC nodes [7,9], we purposely unbalance load to meet SLO requirements and facilitate dynamic scaling. We evaluate the benefits of our SLO-aware load balancer in the context of stateless and functionally decomposed EPC nodes, and contrast it with traditional stateful models.

### 3 Research Goals

From the above discussed 5GCoreLite platform, we plan to achieve following research goals:

- **Traffic Aware resource scheduling:** The lack of awareness about traffic characteristics and IoT device type allows the resources to be assigned redundantly to IoT devices. For example, a cellular-enabled street lamp might not need mobility support from the EPC core. Hence, by bringing

traffic awareness the resources in the EPC core could be efficiently and dynamically assigned.

- **Multi-level Adaptive SLO-aware Load balancing:** Communication within EPC can be complex. It is challenging to provide resource efficient load balancing across all the EPC nodes (i.e., MME, HSS, SGW and PGW nodes) that adheres to the stringent SLO requirements. Hence, we propose to build lightweight optimizations that takes into consideration the local optimizations at each of the EPC to take optimum decision globally for distributing the connections to each of the network functions. To improve the load balancing with in each host and across multiple hosts, we propose to implement controller infrastructure that actively monitors the SLO violations that occurs with in the requirements of each IoT device’s traffic and effectively instantiate and migrate the connection to different stateless SFC.
- **Microservice prioritization:** To improve the SLO requirements of each IoT device, the latency can be further optimized by effectively setting the NF or microservice *priorities* when sharing CPU resources for handing specific control procedure with in any decomposed microservices. Dynamically enforcing priorities is key challenge that need to be addressed for efficiently handling SLO requirements.

### References

- [1] Gartner Reveals Top Predictions for IT Organizations and Users in 2017 and Beyond. 2017. <http://www.gartner.com/newsroom/id/3482117>.
- [2] Nokia Siemens Networks: Signaling is growing 50% faster than data traffic. October 2017. <https://goo.gl/oTbTmM>.
- [3] Muhammad Zubair Shafiq, Lusheng Ji, Alex X. Liu, Jeffrey Pang, and Jia Wang. A First Look at Cellular Machine-to-machine Traffic: Large Scale Measurement and Characterization. In *Proceedings of the 12th ACM SIGMETRICS/PERFORMANCE Joint International Conference on Measurement and Modeling of Computer Systems*, SIGMETRICS '12, pages 65–76, New York, NY, USA, 2012. ACM.
- [4] M. Laner, P. Svoboda, N. Nikaein, and M. Rupp. Traffic Models for Machine Type Communications. In *ISWCS 2013; The Tenth International Symposium on Wireless Communication Systems*, pages 1–5, Aug 2013.
- [5] Study on provision of low-cost machine-type communications (mtc) user equipments (ues) based on lte, 3gpp spec: 36.888, 2017. <https://portal.3gpp.org/desktopmodules/Specifications/SpecificationDetails.aspx?specificationId=2578>.

- [6] Vasudevan Nagendra, Himanshu Sharma, Ayon Chakraborty, and Samir R. Das. LTE-Xtend: Scalable Support of M2M Devices in Cellular Packet Core. In *Proceedings of the 5th Workshop on All Things Cellular: Operations, Applications and Challenges*, MobiCom Workshop, ATC '16, pages 43–48, New York, NY, USA, 2016. ACM.
- [7] Arijit Banerjee, Rajesh Mahindra, Karthik Sundaresan, Sneha Kasera, Kobus Van der Merwe, and Sampath Rangarajan. Scaling the LTE Control-plane for Future Mobile Access. In *Proceedings of the 11th ACM Conference on Emerging Networking Experiments and Technologies*, CoNEXT '15, pages 19:1–19:13, New York, NY, USA, 2015. ACM.
- [8] Diomidis S Michalopoulos, Mark Doll, Vincenzo Sciancalepore, Dario Bega, Peter Schneider, and Peter Rost. Network slicing via function decomposition and flexible network design. 2017.
- [9] Daniel E. Eisenbud, Cheng Yi, Carlo Contavalli, Cody Smith, Roman Kononov, Eric Mann-Hielscher, Ardas Cilingiroglu, Bin Cheyney, Wentao Shang, and Jinnah Dylan Hosein. Maglev: A Fast and Reliable Software Network Load Balancer. In *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*, pages 523–535, Santa Clara, CA, 2016.