

## Accelerating Network Function Virtualization and Service Function Chain Processing for Emerging 5G Services and Edge Computing

Zhi-Li Zhang ([zhzhang@cs.umn.edu](mailto:zhzhang@cs.umn.edu))  
Department of Computer Science & Engineering  
University of Minnesota

Network Function Virtualization (NFV), coupled with Software Defined Networking (SDN), promises to revolutionize networking by allowing network operators to dynamically modify and manage networks. Operators can create, update, remove or scale out/in network functions (NFs) on demand, construct a sequence of NFs to form a so-called service function chain (SFC) and steer traffic through it to meet various policy and service requirements. In the emerging 5G technologies – besides innovations in radio technologies such as 5G new radio (NR), NFV will be a key enabling technology underpinning the envisioned 5G “Cloud RANs” (radio access networks), MECs (mobile edge clouds) and packet core networks for support of network slicing and diverse services ranging from enhanced mobile broadband (eMBB) to massive machine type communications (mMTC) and ultra-reliable low latency communications (URLLC). For example, upon a request for a service (e.g., from a mobile user or a machine, say, an autonomous vehicle or an industrial controller), a SFC will be dynamically constructed using a series of virtualized network functions (vNFs) such as firewalls, mobility managers, network address translators, traffic shapers and so forth that are deployed on demand at appropriate locations within a (dynamic) network slice to meet the desired service requirements

Despite all the hype, realizing many touted advantages of NFV is a daunting challenge in practice, especially when applied to emerging 5G networks, where high scalability, availability and performance will be critical. Given that vNFs are implemented in software and run on commodity servers that are shared compute resources, delivering performance that can match conventional hardware “middleboxes” is a nontrivial task. Moreover, traffic traversing virtualized SFCs may suffer from reduced throughput and increased latency. The flexibility afforded by the combination of SDN and NFV will also likely result in increasingly longer SFCs as networks become ever more highly automated – making this challenge ever more relevant. Besides leveraging specialized hardware and smart NIC capabilities such as DPDK that are increasingly adopted by modern servers – and further exploiting parallelism – to speed up SFC packet processing within a single multi-core server, scale-out is a major (software) technique afforded by NFV for circumventing the performance challenge: by distributing NFs across multiple servers dynamically, it is possible to significantly increase the overall system throughput and reduce SFC processing latency. Unfortunately, as most of NFs of interest is stateful, this poses many challenges in automatically and elastically scaling of NFV across multiple servers while ensuring the correctness of SFC processing.

**Past Research:** We have investigated how to improve the latency of SFC by leveraging parallel packet processing among NFs, in contrast to the conventional serial processing. We have developed *ParaBox*, a novel hybrid packet processing architecture, that dynamically distributes packets to NFs in parallel when possible and merge their outputs to guarantee the correctness of network and service policies. We implement a prototype as a proof-of-concept to demonstrate the feasibility of *ParaBox*. Our preliminary experiment results show that it can not only significantly reduce the service chaining latency but also improve the throughput.

**Ongoing Research:** Building upon the basic idea of *Parabox*, we are currently developing a

novel distributed parallelization framework, dubbed HydraNF, for accelerating NFV service function chain processing at scale. *This research project is funded and supported in part by an NSF/EU ICE-Project, titled “Accelerating NFV Service Function Chain Processing at Scale.”* HydraNF is designed specifically to simultaneously tackle the performance and auto-scaling challenges in real-world large scale deployment of NFV by taking full advantage of a cluster of multi-core servers for dynamic and elastic scale-out. Unlike existing research that mostly focus on enhancing NFV performance for a single server, we recognize that emerging 5G cloud RANs, virtualized EPC networks and edge cloud computing will likely operate in cluster environments with multi-core servers. Leveraging the software nature of vNFs, HydraNF carefully analyzes the configurations, operational rules and state variables of NFs to identify both opportunities and constraints for parallel and distributed SFC packet processing, and decomposes a “monolithic” SFC into a (fine-grained) SFC processing graph (SPG). Based on vNF performance profiles, server capacities and NF placement requirements, HydraNF automatically scales out SFC processing through distribution across multiple servers, and parallelizes the NFV packet processing pipelines within each server by utilizing multiple cores: this is done by exploiting parallelism at both the network function level and traffic level. HydraNF can also significantly enhance SFC availability via appropriately provisioning backup vNFs using its built-in mirror and merging capabilities. In the context of developing HydraNF, we have also starting investigating how the multi-core server architecture affects the performance of SFC execution models by conducting extensive experiments on a multi-core server cluster testbed.

**Planned Efforts:** As part of our planned efforts, we will continue the development of key components. These include NF behavior profiling, performance analysis, SFC decomposition algorithms, NFV execution and runtime systems, and auto scaling mechanisms. In particular, we are studying the key performance bottlenecks and novel ways to scale up/out NFV/SFC performance to meet increasing line speeds. In particular, as we move from 10/40 Gbps to 100/400 Gbps line rates, it will be increasingly challenging to build an NF execution framework that can deliver high performance at the maximum line speed using commodity servers, while providing scalability and flexibility afforded by software. Our current ongoing research reveals that existing NFV frameworks and platforms will unlikely be able to keep up with 100 Gbps or beyond line speeds – this is because the average per-packet processing time will be within 10 nanoseconds, the same speed that current L1/L2 cache dedicated to individual cores. Therefore optimizing the operations of each NF to minimize L1/L2 cache misses is crucial. However, the stateful nature of most NFs, especially, when they are chained in a service function, creates many challenges. This calls for new modular NFV architectures and frameworks that are fine-grained, granular, far flexible with deep visibility to NF operations and NF state. Toward these goals, we plan to develop fine-grained, declarative programming models for NFs and their compositions, design "architecture-aware", highly optimized yet flexible compiler, runtime systems and execution environments,

**International Collaboration:** In conjunction with our NSF/EU ICE-T project on accelerating network function virtualization and service function chains, we are collaborating Dr. Arturo Azcorra, Professor at Universidad Carlos III de Madrid (UC3M) and Director of IMDEA Networks Institute in Madrid, Spain, and his research team at UC3M/IMDEA Networks. Our NSF/EU ICE-T collaborative research project will be carried out in conjunction with two recently funded 5G-EVE and 5G-VNNI projects under the EU Horizon 2020 ICT Programme.