

# Experimental Demonstration of Scalable and Low Latency Crowd Management Enabled by 5G and AI in an Accelerated Edge Cloud

J. C. Borromeo<sup>(1)</sup>, K. Kondepu<sup>(2)</sup>, S. Fichera<sup>(1)</sup>, P. Castoldi<sup>(1)</sup>, L. Valcarengi<sup>(1)</sup>

<sup>(1)</sup>*Scuola Superiore Sant'Anna, Pisa, Italy, [luca.valcarengi@santannapisa.it](mailto:luca.valcarengi@santannapisa.it);*

<sup>(2)</sup>*Indian Institute of Technology Dharwad, Dharwad, India*

**Abstract:** This demo shows how crowd management can benefit from 5G connectivity and accelerated edge cloud AI-based computation to achieve low latency and scalability with respect to a mobile device-based physical distancing computation. © 2021 The Author(s)

**OCIS codes:** (060.0060) Fiber optics and optical communications; (060.4250) Networks;

## 1. Overview

The widespread deployment of edge micro data centers (i.e., edge cloud) [1] together with the support of 5G communication, enables new applications based on data elaboration exploiting Artificial Intelligence/Machine Learning (AI/ML) [2], which can be impractical in end user devices. Indeed, although user equipment (UE) computational capabilities are continuously increasing, such heavy elaborations would compromise their battery duration. For example, video surveillance, search and rescue operation in a hostile environment, traffic management in smart cities, and emergency management by means of Unmanned Aerial Vehicles (UAV) represents realistic scenarios that can benefit from elaboration offloaded to the edge cloud to spare UAV energy [3].

Recently, the COVID-19 pandemic has fostered the development of applications related to crowd management. Social or also called physical distancing is one of the main countermeasures to fight the diffusion of SARS-CoV-2 virus. Applications performing physical distancing verification based on AI/ML are therefore emerging [4]. Often, local police might resort to portable battery powered devices, such as portable cameras or smartphones, to run such applications. Although the computational power of such devices is continuously increasing, physical distancing applications, if locally run in the portable devices could rapidly exhaust battery power and slow them down. This demo shows that offloading the computationally intensive and energy demanding physical distancing application, based on AI/ML, to an accelerated edge cloud, enabled by 5G connectivity, can provide shorter response time and lower energy consumption than the ones achievable by running the application in portable devices.

The proposed demonstration is depicted in Fig. 1. A smartphone acquires the images of an area (e.g., part of the conference premises) where people are present. The raw images are sent to an edge cloud through a 5G connection. The edge cloud is equipped with CPUs, GPUs and FPGAs (i.e., it is accelerated). The edge cloud hosts 5G virtualized Radio Access Network (vRAN) and Next Generation Core (vNGC) network functions. Moreover, it hosts an application for physical distancing computation based on AI/ML. The images conveyed by the 5G fronthaul and backhaul network are elaborated (and also anonymized and not stored for privacy issues) by the physical distancing application and a reply is sent back to the UE, where either a green (i.e., physical distancing fulfilled) or red (i.e., physical distancing not fulfilled) symbol is visualized. Visitors to the demo will experience the capability of the accelerated edge cloud and 5G network to provide a real time physical distancing verification of the area they are filming by offloading the image elaboration to the accelerated edge cloud.

## 2. Innovation

The demo combines very recent innovations in the communication and service fields. First, the demo exploits huge potentials of accelerated edge cloud. Indeed, accelerated edge cloud can be used not only to run computationally intensive functions at the application level but also computationally intensive 5G functions such as the 5G mobile network stack physical functions. The demo also includes the implementation of a disaggregated Next Generation eNB (gNB) where functions are split into Radio Unit (RU), Distributed Unit (DU), and Central Unit (CU) as currently proposed in the standards [5]. The gNB split Option 2 is considered where the RU includes RF functions only, the DU, which is connected to the RU by an optical fronthaul performs all the functions from the Low-PHY to the Radio Link Control (RLC), and the CU performs all the functions above RLC. Moreover, some of the RAN functions are virtualized and accelerated. In particular, the DU Low-PHY functions are offloaded onto a Field Programmable Gate Array (FPGA) present in the accelerated edge micro data center. Finally, the physical distancing application, based on a slightly modified version of an open source one based on Yolo, exploits a Graphical Processing Unit (GPU) available

in the accelerated edge data center. In summary, the demo shows the benefits of exploiting an accelerated edge cloud for supporting low latency applications.

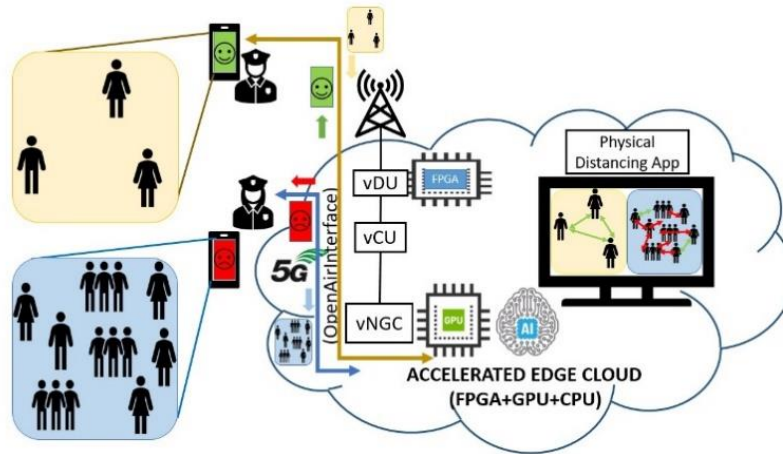


Fig. 1. Demo setup and architecture

### 3. OFC relevance

The demo presents a multifold relevance to OFC 2021. In general, the demo shows how communications networks can contribute to fighting COVID-19 pandemic.

Moreover, the demo shows that, by exploiting the low latency and broadband connectivity provided by 5G and the computational power offered by accelerated edge clouds, computationally intensive applications can be made available also to end user devices that are not computationally powerful, and battery powered without comprising the application response time.

Furthermore, the demo contributes to shedding some light on the type of vertical applications that can fully benefit by elaboration offloaded onto an accelerated edge cloud.

The demo also shows how the integration of high speed optical fronthaul networks with 5G mobile function virtualization and acceleration can benefit end user applications.

Finally, the demo paves the way to other possible applications supported by an accelerated edge cloud. In the considered demo a 1:1 offloading of the end device applications to the edge cloud is considered. However, the proposed approach can be even more beneficial when images shall be elaborated in a synchronous manner for crowd management. Indeed, in this case, a distributed elaboration by the end devices could be difficult, if not impossible, to achieve. Consider, for example, a situation where different security cameras are monitoring different parts of a smart city. If an accident happens in an area (e.g., a shooting, a fire, etc.), citizens shall be directed how to free the area itself and the nearby areas in the fastest possible ways. In this scenario, a crowd management app deployed in the edge cloud receives the images from the different cameras, elaborates them in real-time by counting the number of people in each area and it correlates the images coming from neighboring areas. Then the app sends notification to the citizen smartphones or to street smart displays or smart traffic signs or smart kiosks to direct them away from the accident area as fast as possible.

### 4. Implemented software components and Demo description

A single location demo setup is shown in Fig. 2. One or more smartphones are connected to the 5G network running OpenAirInterface (OAI) [6] as mobile software stack. OAI Core Network (CN) is utilized for implementing the NGC functions. Docker Container based virtualized version of DU (vDU), CU (vCU), and NGC (vNGC) components are deployed in the accelerated edge cloud using Dell PowerEdge R740 servers. Ettus X310 Universal Software Radio Peripherals (USRPs) is utilized as gNB radio front-end performing Radio Frequency (RF) functions. In the accelerated edge cloud the vDU Low-PHY functions are offloaded and accelerated onto an FPGA using the OpenCL Framework [7]. The physical distancing app is deployed as well in the accelerated edge cloud, beyond the NGC, and runs in a GPU. An NVIDIA Tesla T4 data center GPU featuring a 320 NVIDIA Turing tensor cores is utilized. NVIDIA T4 supports all AI framework with 50x higher energy efficiency compared to the CPUs due to the low profile PCIe [8].

