

# Performance and mobility in the mobile cloud

William Tärneberg

Dept. of Electrical and Information Technology

Lund University

Ole Römers väg 3, 223 63 Lund, Sweden

Email: william.tarneberg@eit.lth.se

Jakub Krzywdą

Dept. of Computing Science

Umeå University

SE-901 87 Umeå, Sweden

Email: jakub@cs.umu.se

**Abstract**—In a mobile cloud topology the cloud resources are geographically dispersed throughout the mobile network. Services are actively located with close proximity to the user equipment. Geographically migrating a service from data centre to data centre with its user equipment imposes a load on the affected data centres. Consequently, user equipment mobility provides a fundamental problem to the mobile cloud paradigm. This paper determines the fundamental service performance issues in system of mobile users with dispersed data centres, in relation to the placement of the mobile cloud host nodes and explores the user equipment and provider utility of subscribing to a mobile cloud node at a certain network depth.

**Keywords**—Cloud, Mobility, Mobile infrastructure, User experience consistency, Omnipresent Cloud, Infinite cloud, Edge cloud, Latency, Throughput, Virtualization, Geo-distributed resources, VM migration

## I. INTRODUCTION

Mobile services and user equipment<sup>1</sup> functions are at an increasing rate being virtualized and augmented to the cloud. Rich Mobile Applications [20] will soon, more often than not, be seamlessly executed, partially or fully in the cloud. Alongside applications, fundamental user equipment resources, such as storage and CPU, are being augmented to the cloud. In this resource paradigm, the border between what is being executed locally and remotely is blurred as developers are given more powerful tools to tap into remote ubiquitous generic virtual resources. Additionally, the advent of the internet of things will contribute with vast number of new types of wireless devices, actuators, and sensors querying and connecting to remote virtual and augmented resources.

As we begin to rely more on remote ubiquitous resources we also grow more dependant on the quality of the intermediate WAN network and by the geographical separation of the user equipment and the data centre [10]. Latency sensitive applications and cognitive augmentation services, such as process controls, latency sensitive storage, real time video game rendering, and augmented reality video analysis will quickly falter if subject to a communication delay.

Virtual resources are accessed through increasingly congested mobile access networks. More devices are crowding the mobile networks and applications are generating and receiving more data, this congestion contributes to communication latency [15]. Additionally, the geographic discrepancy between the user equipment and the data centre introduces a propagation delay, bounded by the speed of light.

The mobile cloud paradigm, put forward by [3], [9], [17], [20], [26], attempts to remedy the aforementioned congestion and latency performance inhibitors by locating cloud resources at the edge of, and adjacent to, the mobile access network. In the ad-hoc scenario, resources are shared amongst user equipments where each connected user equipment surrenders its available resources to its peers. In its centralized form, data centre resources are proposedly located at the edge of the network, adjacent or integrated into an radio base station, catering for the user equipments located within its cell coverage. Alternatively, or complimentary, data centres are proposedly integrated with resources in the common administrative nodes of the proposed virtualized radio access networks. The scale and the degree of dispersion can be optimized for each application, given the applications resource tiers and its users mobility behaviour.

Round trip time, is arguably proportional to the geographic distance between the user equipment and the data centre. Services hosted in the mobile cloud are migrated with the user equipment, through the network, to minimize this incurred latency and congestion on the adjacent WAN. In practice, services, or rather the VMs that host the services, are migrated to a data centre that, is available, provides the lowest service latency, and incurs least global network congestion. Doing so might minimize the experienced delay for the user equipment, but will incur a migration overhead in the hosting data centre and in the network over which the VM is migrated or duplicated. Conceivably, various provisioning schemes and cost functions can be deployed to minimize both the delay experienced by the user and the added resource strain on the data centre and the intermediate network.

User equipment mobility is a key differentiator between traditional cloud computing with distant data centres and the mobile cloud, and is a fundamental dynamic property of a mobile cloud. In order to be able to optimize the mobile cloud topology, it is essential to understand how user equipment mobility affects the perceived service performance and what load it imposes on the network.

The topology paradigms of tomorrows all-IP (Internet Protocol) mobile networks [8], [14] are hot topics of research, but one can assume that they will be influenced by the notion of virtualized resources [6], [11]. Large portions of radio base stations can proposedly be virtualized and centralized to a common data centre with a locally-bounded service domain, shared by several radio base stations, leaving the radio base stations, in principal, with just the radio interface [21]. The degree of centralization is conceivably geographically

---

<sup>1</sup>Any user client device accessing a service, such as a mobile phone

bounded by propagation delay and signal attenuation, and is resource hampered by the aggregated traffic that passes through the dedicated data centre. There is to our knowledge, very little research exploring future mobile Telecom infrastructure topologies with the mobile cloud in mind. There is on the other hand, extensive research directed at exploring relevant economic and IT models of how to integrate existing Telecom services to the cloud and how to apply Telecom-grade SLAs to existing cloud services [1], [8], [22]. These services are frequently proposed to reside in the network and managed by the Telecom operators.

The concept of geo-distributed cloud resources has received some research attention over the past few years, but has had a clear research focus on storage and shared data. The authors of [5] present a method to geographically migrate shared data resources globally, not only to minimize the distance between the user equipment and the data centre, and thus service latency, but also to globally load-balance the hosting data centres on which the observed service is distributively hosted. Their results reveal a significant reduction in service latency, inter-data centre communication, and contributed WAN congestion. Their proposed control process runs over long time periods and operate on a global scale with relatively geographically static users. Although sharing some fundamental dynamics, albeit at different scales, in contrast, the mobile cloud paradigm, user equipment movement is more rapid and hand-overs between radio base stations is likely to occur during a service session. Additionally, mobile cloud virtualized resources are assumed to be universal and do not just include storage, they vary in size and capabilities, are deployed by the Telecom operators, and are based on local needs and demand.

The field of mobile cloud has much in common with field of geo-distributed cloud resources, but is dominated by the notions of augmenting user equipments through virtualizing their resources [4] and reducing service response times through geo-cascaded data caching [3], [24]. As a result, much of the research is concerned with coping with specific dynamics, and do thus not address the generic case of small geo-distributed data centres, serving a local mobile subscriber populous. There are, to our knowledge, significant research gaps in how cloud services perform when hyper-dispersed and rapidly migrated. Additionally, there is little research on how the mobile cloud can be accommodated in and optimized for future network topologies.

In this paper we investigate the fundamental effects of user equipment mobility on the mobile cloud by observing data centre utilization, the proportion of data centres resources spent on migrating services, and how service performance is affected by migration. In addition, we propose a simulation model built around the fundamental dynamics that contribute to package latency and VM utilization, and is designed to examine the fundamental and generic resource problems in a mobile cloud of mobile user equipments. The models include a generic mobile network, populated with user equipments subscribing to a number of services, served by a number of locally geo-distributed data centres. The simulation model is subjected to multiple scenarios in constellations of varying number of users, services, and data centre clustering.

The simulated scenarios reveal ...

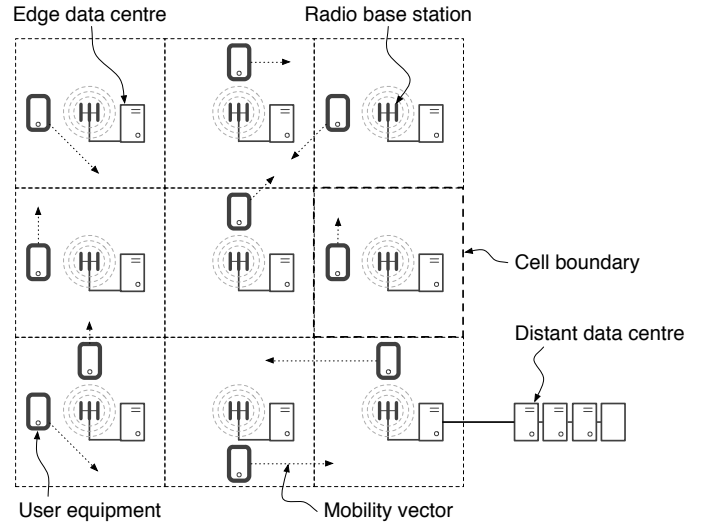


Fig. 1: System model

In this paper, Section II details which aspects and abstractions of the mobile cloud topology that are included in our experiments. Furthermore, the simulation model is specified in Section III. Section IV details the specifics of the simulation experiments. Lastly, Sections V and VI present the results and the consultations drawn from the experiments.

## II. DESIRED MODEL

The desired model shall provide a setting for which we can explore fundamental resource and performance properties of the mobile cloud system paradigm with mobile user equipments. The mobile user equipments, radio access network, and service application will subject the data centres to a load characteristic for generic mobile phone traffic and the type of services that plausible might be deployed to the mobile cloud. Additionally, the model shall capture the systems fundamental parametrizable properties.

As the topology of any future mobile cloud or proposed forthcoming mobile networks is yet to be determined, in this paper we propose a generic Telecom infrastructure model that disregards generational specific properties such as those found in the physical layer and radio resource load-balancing disciplines. These properties are not system variables at the abstraction level the mobile cloud needs to be modelled in this paper. Nevertheless, conceivably, and in order to confine the geographic domain of the model, the model adheres to current general LTE cell planing practices [25], see Figure 1.

In order to be able to explore the fundamental effects of mobility on the performance of an mobile cloud service in the generic case, the model does not adhere to any socio-demographic patterns or urban topologies. In the absence of any geographic bias, the mobile network base stations are uniformly distributed across its 2-dimensional domain.

Similarly, in order to represent the variety of possible services, the service model shall generate traffic that is characteristic for an active, generic, user equipment. Additionally, the generated traffic shall be provided by a stochastic process that is independent of location.

The mobility model, the service model, and the uniformly distributed mobile network will provide the modelled data centres with a characteristic workload. It is worth reiterating that the traffic load and the characteristics of the service are more relevant to our investigation than specific topological and network properties.

The data centre model will host multiple service in VMs that will process the arriving requests corresponding to its service commitment. Additionally, when a VM is migrated between data centres it incurs a load on both data centres. Furthermore, the resources within a data centre are shared amongst the hosted VMs.

### III. SIMULATION MODEL

#### A. Service

Most mobile applications use HTTP as a means to communicate with remote services, often through a web interface [12], [19]. We will model our service application to that of a stateless web service catering to the subscribers in the local network. The HTTP traffic model in [18] provides an open loop traffic model with a long tailed session size distribution, representative of the diversity of mobile the requests.

Each session is separated in time with a poisson process of  $\lambda_{ses}$ . Each sessions produces  $N_{req}$ , sampled from an inverse Gaussian distribution, where each request is separated in time by log-normal distributed delay  $D_{req}$ .

Each service adheres to the same properties, and are only distinguished by the VM in which they are running. Additionally, the properties of the service model are independent of user equipment state and mobility mode.

#### B. Mobility

As user equipments traverse the mobile cloud network, the service(s) they subscribe to will migrate to accommodate the changing distribution of user equipments in the network. The 2-dimensional, multi modal, mobility model detailed in [7] provides us with a uniform distribution of users, with a realistic rate of mobility.

The model defines the fundamental timing and mobility properties of user equipment movement, such as the speed, acceleration, and direction the user equipment is moving in, as well as for how long and when in time to turn next.

#### C. Mobile access network

Forthcoming cell planing practices aim to increase area energy efficiency by favouring smaller cells in urban areas [13], [27]. The model will therefore employ a small homogeneous mobile network composed of  $N_{rbs}$  equidistantly distributed radio base stations. The domain which the network serves is populated by a homogeneous group of user equipments, with a uniform service subscription distribution. A user equipment is handed over between base stations at the point where they cross the cell boundary distinguishing two independent radio base stations defined by the width of the rectangular cells  $d_{rbs}$ . The mobile access network model does not take into account the physical layer, channel provisioning, and cell load balancing. Additionally, the radio access network

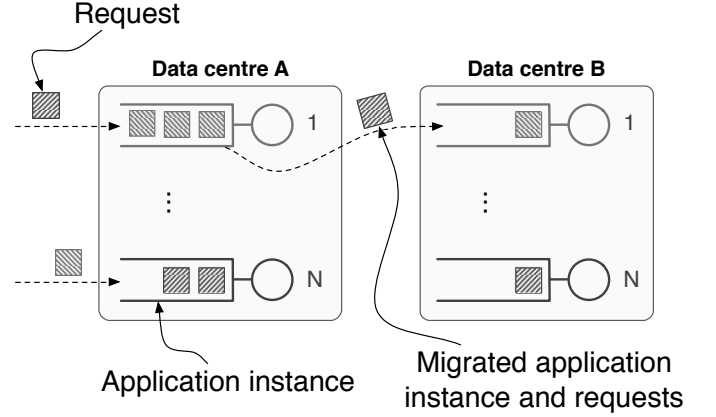


Fig. 2: Data centre model

functions as a mechanism to associate user equipments with data centres, propagation and system processing delays are thus not modelled.

The network is populated by  $N_{userequipment}$  each subscribing to one of the  $N_{ser}$  available services. For the sake of model simplicity ambient users and traffic have not been modelled.

#### D. Core network

The delay induced by the core network is modelled with a Weibull delay  $T_{net}$  in multiples of the number of network nodes between the source and the destination, in accordance with [23]. The distance between radio base stations is equal to the cell dimension  $d_{rbs}$ . Associated radio base stations are equidistant to their common data centre, and are for the sake of simplicity assumed to be separated by one router.

#### E. Data centre

We model a data centre using two parameters, the number of CPUs and their speed (in FLOPS or MIPS). We assume that all CPUs in one data centre are identical, we do not represent that they are located in separate physical machines (servers) and we do not consider memory, storage and intra data centre network.

1) *Hosting applications in a data centre:* Applications run on the resources located in the data centres. The time that is necessary for serving an application request depends on the CPU speed of a hosting machine. The capacity of a data centre is determined by the number of available CPUs. New instances of application are not accepted when the capacity of a data centre is reached.

2) *VM initialization and termination:* When a decision of deploying a new service in a data centre is taken, a new VM will be started there to host that service. Due to the startup time, the newly admitted VM will not be able to start processing requests for a period of  $T_{vm\_init}$ . Nevertheless, the new VM will start using resources of the data centre from the time of admission. Because of that, the service time,  $T_{i,vm}$ , for each of the VMs hosted in that data centre will be recalculated to reflect the temporary load scenario. Similarly, after finishing serving the last request, the VM will still be using the resources of the data centre for time  $T_{vm\_release}$ .

3) *VM activation scheme*: Definitions: service, session, client subscribed to service

**Private VM** — IaaS like, a VM is exclusively used by one user for offloading computations from his user equipment. User provides the executable program that is loaded to an edge data centre from an user equipment or a remote data centre.

**per-session** — a VM is initialized upon receiving the first request from a client and terminated just after finishing processing the last request of the session.

**client-within-cell** — a VM is initialized when a client that is subscribed to a service enters a cell and is kept alive as long as he stays within the cell.

**Shared VM** — SaaS like, a VM hosts a service that can be concurrently accessed by many users.

**any-client-running** — a VM is initialized upon receiving the first request from a client and terminated when there are no more requests to serve (waiting queue is empty and all sessions are finished).

**any-client-within-cell** — a VM is initialized when the first client that is subscribed to a service enters a cell and is kept alive as long as any subscribed client stays within the cell.

4) *VM migration*: Migration has an influence on the service performance as well as on the resource availability on both source and destination data centre. On the source side, apart from ordinary resource requirements due to serving requests, VM consumes resources for sending its image to the destination data centre. In a case of postcopy live migration, VM at the source side still uses some resources even after the workload is redirected to the new location. That is because the VM at the destination side pulls remaining memory pages from the source data centre.

To model the overhead of migration on the service performance additional delay in the response time should be introduced. Primarily, for the time of transferring the image of VM,  $T_{vm\_transfer}$ , the time of serving a request should be lengthen by  $D_{vm\_transfer}$ , because of using a part of resources for I/O operations. Moreover, during the time of switching the execution from the source to the destination,  $T_{vm\_downtime}$ , VM is not able to serve any requests. Additionally, in the case of the postcopy migration, delay  $D_{memory\_pull}$  occurs for some number,  $N_{memory\_pull}$ , of the first requests after redirecting the workload to the new data centre, due to the remote memory calls.

When VMs compete for resources, running additional VM on the destination side introduces an overhead by increasing the response time of other collocated VMs.

The statefulness of the service dictates the amount information that needs to be migrated to the succeeding data centre. The  $H$  parameter in Equation 3, is the state entropy of the communication between a user equipment and a data centre which thus dictates the scale of the incurred overhead when migrating that users state the receiving data centre.

## F. Fundamental system dynamics

The fundamental dynamics of the system can be expressed with if parameters as expressed in Equations 2 & 3. Where  $LU$  and  $LS$  are vectors containing a user equipments all successive radio base station and data centre associations, respectively, throughout the simulation time.

$$D_{M,LU(0)} + D_{N,LU(0) \rightarrow LS(0)} + \sum_{i \in LS[1,N-1]} (D_{Mig} + T_{Q,i} + D_{N,i \rightarrow i+1s}) + T_{Q,LS(N)} + T_{S,LS(N)} + D_{S,LS(N) \rightarrow LU(N)} + D_{M,LU(N)} \quad (1)$$

$$\sum_{i=1}^{N_{init}} T_{init} + \sum_{i=1}^{N_{term}} T_{term} + \sum_{i=1}^{N_{pktser}} T_{ser} + \sum_{i=1}^{N_{usrmig}} \left( \sum_{j=1}^{N_{pktmig}} T_{mig} + S \cdot T_{vmmig} \right) + T_{idle} \quad (2)$$

## IV. EXPERIMENTS

The mobile cloudmodel was implemented in Java employing simjava [2] as the event driven framework.

### A. User equipment, mobility and network

To reveal the effect of varying load on the data centres, the number of user equipment service subscribers  $N_{userequipment}$ , placement of the data centres, and the number for services were varied independantly throughout as many simulation scenarios.

In proportion to the range of movement and to be evenly divisible by  $N_{dc}^2$ , the simulation domain spans 16 radio base stations  $N_{rbs}$ . The cell dimension  $d_{rbs}$  adhere to a proposed maximum cell radius of 750 m, as detailed in [27]. Although rectangular, its area equates to the area of a circular cell with a radius of 750 m, which results in a cell width and hight  $d_{rbs}$  of 1300 m. The population of user equipments subscribing to a service  $N_{ue}$  ranged from 10 to 1600 user equipments, doubling with each increment. The number of data centres  $N_{dc}$  was varied between 16, 4 and 1, with each data centre serving 1, 4 or 16 radio base stations respectively. Additionally, for each run of  $N_{ue}$  the number of services  $N_{ser}$  was varied between 1 and 5.

The simulation time is set to 8 hours,  $T_{sim} = 28800$  seconds, leaving sufficient time for each user equipment to on average visit half of the radio base stations.

### B. Data centre and VM parameters

Times for resource allocation and release and log normal distributed with means  $T_{vm\_init}$  and  $T_{vm\_release}$  respectively, and are modelled after Amazon EC2 measurements for m1.small instance (1 vCPU, 1.7 GB of RAM and 160 GB disk) [16].

The service time for each data centre is set heterogeneously, proportional to the mean request generation rate over the number of radio base stations and the number of VMs running,

| Parameter         | Value                         |
|-------------------|-------------------------------|
| $N_{ue}$          | 10–500 (step: 10)             |
| $N_{rbs}$         | 16                            |
| $N_{dc}$          | {1, 4, 16}                    |
| $N_{ser}$         | 1–5                           |
| $T_{ser}$         |                               |
| $T_{sim}$         | 28800 seconds                 |
| $d_{rbs}$         |                               |
| $T_{net}$         |                               |
| $N_{vm,limit}$    |                               |
| $T_{vm,init}$     | avg 82 s, min 69 s, max 126 s |
| $T_{vm,release}$  | avg 21 s, min 18 s, max 23 s  |
| $T_{vm,transfer}$ |                               |
| $D_{vm,transfer}$ |                               |
| $T_{vm,downtime}$ |                               |
| $D_{memory,pull}$ |                               |
| $N_{memory,pull}$ |                               |

TABLE I: Simulation parameter values

| Component | Distribution | Parameters                |
|-----------|--------------|---------------------------|
| $S_f$     | Pareto       | $K=133000 \alpha=1.1$     |
| $S_r$     | Pareto       | $K=1000$                  |
| $D_r$     | Weibull      | $\alpha=1.46 \beta=0.382$ |
| $D_s$     | Pareto       | $K=1 \alpha=1.5$          |

TABLE II: Service model components

see Equations 4 & 5. Where  $\lambda_{sys}$  is the aggregate arrival rate to the system, and  $K$  is a scaling factor set to 1. Each VM is provisioned to handle 20 concurrent active subscribers ( $N_{userequipment} = 20$ ). The mean service time for each VM is thus  $1/\mu$ .

$$\mu = K \cdot \lambda_{ue} \cdot N_{ue} \cdot \frac{N_{services}}{N_{dc}} \quad (3)$$

$$\lambda_{ue} = \frac{\bar{N}_{requests}}{\bar{N}_{requests} \cdot \bar{T}_{request} + \bar{T}_{session}} \quad (4)$$

As resources are assumed to be ubiquitous we wish to observe the effects of an overloaded data centre. For modelling simplicity, data centre capacity is constrained by a limit of how many VMs  $N_{vm,limit}$  it simultaneously can host. Our experiment regards two fundamental provisioning scenarios. When the VM limit  $N_{vm,limit}$  has been reached, either any new VMs are denied or the additionally needed data centre capacity is shared equally amongst the  $N_{i,vm}$  VMs by increasing the service time  $T_{i,ser}$  with a factor  $K_{over}$ , see Equation 6.

$$K_{over} = \frac{\max(N_{i,vm}, N_{vm,limit})}{N_{vm,limit}} \quad (5)$$

Simulation model parameters used in the experiments can be found in Table I, likewise the service parameters are declared in Table II.

### C. Measurements

To ensure statistical accuracy, each simulation scenario was independently replicated 10 times.

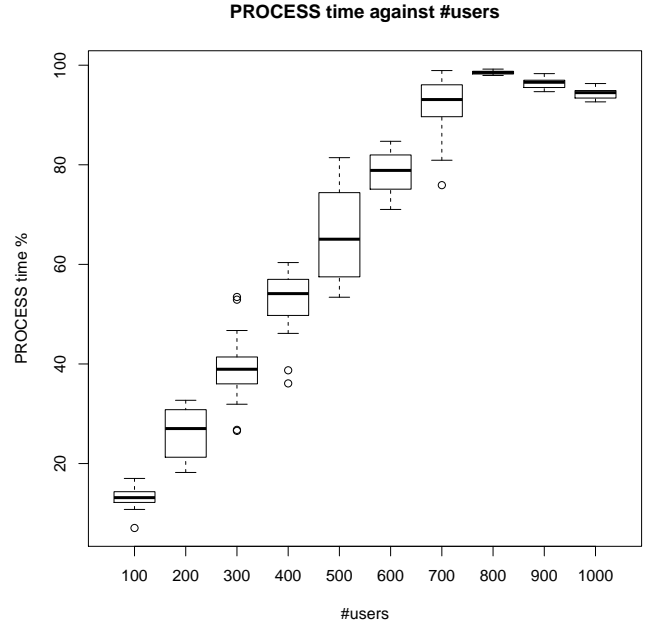


Fig. 3: Amount of time spent processing vs. the number of user equipments in the entire network

For each above mentioned simulation scenario, every request is recorded with where it was processed, if it was terminated, the time spent queuing, processing, and propagating. Similarly, for each VM the proportion of time spent in each state is recorded.

## V. RESULTS

In this section we present the results of our simulation scenarios. The simulation reveals ...

### A. Data centre utilisation

Data centre utilization markedly correlates with the number of potential service subscribers in network  $N_{userequipment}$ . The effect is illustrated Figure 3, which shows a strong growth of VM utilization when approaching maximum stable load at  $N_{userequipment} = 800$ , measured in the percentage of time it spends processing requests. Nevertheless, the process utilization starts to decay once we pass the number of stable subscribers, as more resources are now need to migrate differed users. Conversely, Figure 4 shows how, as a result of increased parallel session residency, the proportion of time spent idle decreases dramatically.

Similarly, compounded by an increased likelihood of congestion in any given VM, the amount of requests needing to be migrated increases near exponentially with a growing number of subscribers  $N_{userequipment}$ . Figure 5 illustrates the growth in the amount of time spent on migrating sessions.

#### 1) Data centre dispersion:

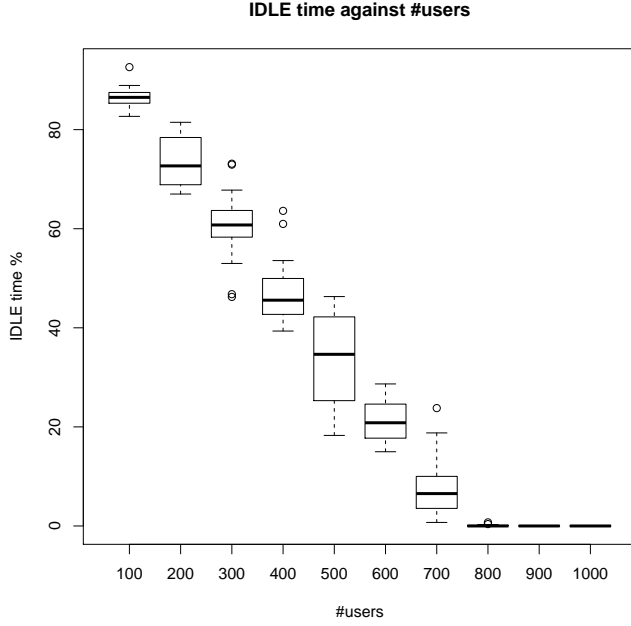


Fig. 4: Amount of time spent in idle state vs. the number of user equipments in the entire network

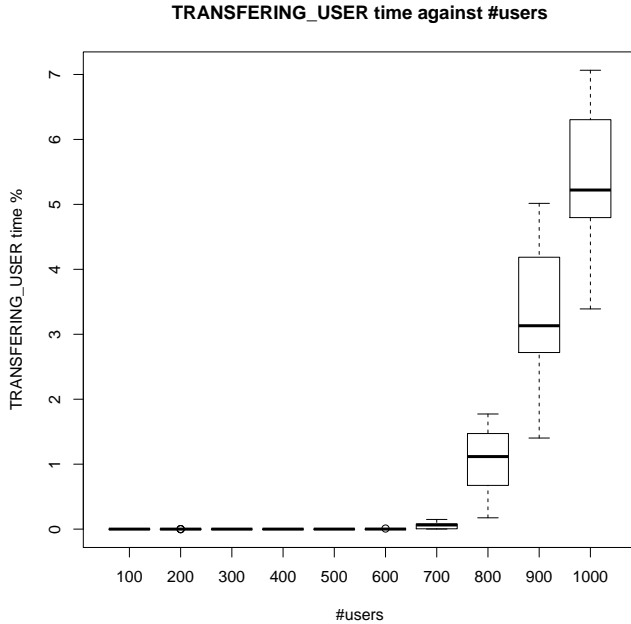


Fig. 5: Amount of time spent transferring requests vs. the number of user equipments in the entire network

#### B. Constrained data centre resources

#### C. Service performance

#### D. Properties of migration

##### 1) Session completion grade per visited VM:



Fig. 6: Number of times a session is migrated vs. the number of user equipments in the entire network

#### E. Inter-Data centre communication

### VI. CONCLUSIONS

### VII. FUTURE RESEARCH

#### A. VM migration schemes

- **Precopy** — the whole memory of VM is copied preemptively before switching the execution to the new data centre.
- **Postcopy** — the memory of VM is copied after switching the execution to the new data centre as is needed to serve the incoming requests.

#### B. VM placement and data centre provisioning schemes

To understand the effects of various migration schemes on the data centre and service performance, the following migration schemes were deployed:

- A VM for service  $S_j$ , if active, resides in the data centre with the largest number of subscribers. If this criteria were to change the the hosting VM will migrate to the resulting data centre.
- Each data centre that hosts a user equipment that subscribes to  $S_j$  hosts an instance of a service  $S_j$  VM. If users disperse, the VM for service  $S_j$  will duplicate to the receiving data centre.

Differentiation between short term jobs (online processing) and long term (offline processing after upload, getting statistics, big data)

Use different service models or change parameters of distributions in current one. Model specific applications (youtube like, facebook like, etc.) and compare them, or show implications of different abstract and extreme configurations.

#### C. Multi-tiered service placement schemes in the mobile cloud

### REFERENCES

- [1] The telecom cloud opportunity. Whitepaper, Ericsson, 2012.
- [2] simjava, 02 2014. Available online at <http://www.icsa.inf.ed.ac.uk/research/groups/hase/simjava/>.

- [3] Ericsson AB. Ericsson and akamai establish exclusive strategic alliance to create mobile cloud acceleration solutions. Press Release, February 2011. <http://www.ericsson.com/news/1488456>.
- [4] S. Abolfazli, Z. Sanaei, E. Ahmed, A. Gani, and R. Buyya. Cloud-based augmentation for mobile devices: Motivation, taxonomies, and open challenges. *Communications Surveys Tutorials, IEEE*, 16(1):337–368, First 2014.
- [5] Sharad Agarwal, John Dunagan, Navendu Jain, Stefan Saroiu, Alec Wolman, and Harbinder Bhogan. Volley: Automated data placement for geo-distributed cloud services. In *NSDI*, pages 17–32, 2010.
- [6] Fabio Baroncelli, Barbara Martini, and Piero Castoldi. Network virtualization for cloud computing. *annals of telecommunications-Annales des télécommunications*, 65(11-12):713–721, 2010.
- [7] Christian Bettstetter. Smooth is better than sharp: A random mobility model for simulation of wireless networks. In *Proceedings of the 4th ACM International Workshop on Modeling, Analysis and Simulation of Wireless and Mobile Systems, MSWIM '01*, pages 19–27, New York, NY, USA, 2001. ACM.
- [8] G. Caryl, T. Rings, J. Gallop, S. Schulz, J. Grabowski, I. Stokes-Rees, and T. Kovacicova. Grid/cloud computing interoperability, standardization and the next generation network (ngn). In *Intelligence in Next Generation Networks, 2009. ICIN 2009. 13th International Conference on*, pages 1–6, 2009.
- [9] Abhishek Chandra, Jon Weissman, and Benjamin Heintz. Decentralized edge clouds. *Internet Computing, IEEE*, 17(5):70–73, 2013.
- [10] Baek-Young Choi, Sue Moon, Zhi-Li Zhang, Konstantina Papagiannaki, and Christophe Diot. Analysis of point-to-point packet delay in an operational network. *Computer networks*, 51(13):3812–3827, 2007.
- [11] NM Mosharaf Kabir Chowdhury and Raouf Boutaba. Network virtualization: state of the art and research challenges. *Communications Magazine, IEEE*, 47(7):20–26, 2009.
- [12] Hossein Falaki, Dimitrios Lymberopoulos, Ratul Mahajan, Srikanth Kandula, and Deborah Estrin. A first look at traffic on smartphones. In *Proceedings of the 10th ACM SIGCOMM conference on Internet measurement*, pages 281–287. ACM, 2010.
- [13] AJ. Fehske, F. Richter, and G.P. Fettweis. Energy efficiency improvements through micro sites in cellular mobile radio networks. In *GLOBECOM Workshops, 2009 IEEE*, pages 1–5, Nov 2009.
- [14] M.A.F. Gutierrez and N. Ventura. Mobile cloud computing based on service oriented architecture: Embracing network as a service for 3rd party application service providers. In *Kaleidoscope 2011: The Fully Networked Human? - Innovations for Future Networks and Services (K-2011)*, *Proceedings of ITU*, pages 1–7, 2011.
- [15] Ningning Hu, Li Li, Zhuoqing Morley Mao, Peter Steenkiste, and Jia Wang. A measurement study of internet bottlenecks. In *INFOCOM 2005. 24th Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings IEEE*, volume 3, pages 1689–1700. IEEE, 2005.
- [16] A Iosup, S. Ostermann, M.N. Yigitbasi, R. Prodan, T. Fahringer, and D. H J Epema. Performance analysis of cloud computing services for many-tasks scientific computing. *Parallel and Distributed Systems, IEEE Transactions on*, 22(6):931–945, June 2011.
- [17] Niko Kiukkonen, Jan Blom, Olivier Dousse, Daniel Gatica-Perez, and Juha Laurila. Towards rich mobile phone datasets: Lausanne data collection campaign. *Proc. ICPS, Berlin*, 2010.
- [18] Zhen Liu, Nicolas Niclausse, and César Jalpa-Villanueva. Traffic model and performance evaluation of web servers. *Performance Evaluation*, 46(2):77–100, 2001.
- [19] Gregor Maier, Fabian Schneider, and Anja Feldmann. A first look at mobile hand-held device traffic. In *Passive and Active Measurement*, pages 161–170. Springer, 2010.
- [20] Verdi March, Yan Gu, Erwin Leonardi, George Goh, Markus Kirchberg, and Bu Sung Lee. ucloud: Towards a new paradigm of rich mobile applications. *Procedia Computer Science*, 5(0):618 – 624, 2011. The 2nd International Conference on Ambient Systems, Networks and Technologies (ANT-2011) / The 8th International Conference on Mobile Web Information Systems (MobiWIS 2011).
- [21] Jordan Melzer. Cloud radio access networks.
- [22] S. Pal and T. Pal. Tsas; customized telecom app hosting on cloud. In *Internet Multimedia Systems Architecture and Application (IMSAA), 2011 IEEE 5th International Conference on*, pages 1–6, 2011.
- [23] Konstantina Papagiannaki, Sue Moon, Chuck Fraleigh, Patrick Thiran, and Christophe Diot. Measurement and analysis of single-hop delay on an ip backbone network. *Selected Areas in Communications, IEEE Journal on*, 21(6):908–921, 2003.
- [24] L. Ramaswamy, Ling Liu, and A. Iyengar. Cache clouds: Cooperative caching of dynamic documents in edge networks. In *Distributed Computing Systems, 2005. ICDCS 2005. Proceedings. 25th IEEE International Conference on*, pages 229–238, June 2005.
- [25] J Salo, M Nur-Alam, and K Chang. Practical introduction to lte radio planning. *A white paper on basics of radio planning for 3GPP LTE in interference limited and coverage limited scenarios, European Communications Engineering (ECE) Ltd, Espoo, Finland*, 2010.
- [26] Mahadev Satyanarayanan, Paramvir Bahl, Ramón Caceres, and Nigel Davies. The case for vm-based cloudlets in mobile computing. *Pervasive Computing, IEEE*, 8(4):14–23, 2009.
- [27] Suhail Najm Shahab, Tiong Sieh Kiong, and Ayad Atiyah Abdulkafi. A framework for energy efficiency evaluation of lte network in urban, suburban and rural areas. *Australian Journal of Basic and Applied Sciences*, 7(7):404–413, 2013.