

x-cloud challenges

Jakub Krzywda
Dept. of Computing Science
Umeå University
SE-901 87 Umeå, Sweden
Email: jakub@cs.umu.se

William Tärneberg
Dept. of Electrical and Information Technology
Lund University
Ole Römers väg 3, 223 63 Lund, Sweden
Email: william.tarneberg@eit.lth.se

Abstract—The abstract goes here.

I. INTRODUCTION

Services accessed from mobile devices is increasingly provided hosted in the cloud. Migrating services to the cloud compliments and virtualizes a mobile devices resources. Nevertheless, the delay as a result of executing code or storing resources fully or partially in the cloud introduces and unwanted delay that can disrupt the desired seamless user experience.

Cloud services are traditionally hosted in aggregated data centres scattered throughout the globe. Theirs ability to cost effectively host services is in general fundamentally derived by their size and energy efficient.

The delay introduced when externalising a service to a data centre is a product of the geographic distance to the data centre, the congestion on the intermediate core network, the mobile access network, and the performance of the data centre.

As the number of cloud service increase and the number of devices rely on cloud services increase so will the congestion in the access network and thus also delay. Moreover, as more traditional hardware or device local service are virtualized to the cloud, the demand to low latency response will increase. Conceivably, storage can be fully virtualized, critical control processes can be migrated to the cloud. In the advent of the emergence the internet of things, data from an vast number of sensors, actuators, and peripheral interaction points will flood the internet with traffic, ranging as vastly in size and QoS needs.

In order to be able to reduce latency and be able to formulate relevant SLAs the service hosting nodes will need to reduce their geographic proximity to the consumption device or process.

We refer to proposed paradigm of migrating cloud service hosting and execution closer to the consumption device as the x-cloud . The distribution of cloud data centre hardware and the virtualization of resources can proposedly coexist with future virtualized mobile networks.

II. THE CASE FOR THE X-CLOUD

The motivations for the adoption of the x-cloud are multi-dimensional. ...

The geographic distance separating the device and the host introduces a propagation delay. The intra-continental propagation delay is in order of 20-30 ms, which conceivably

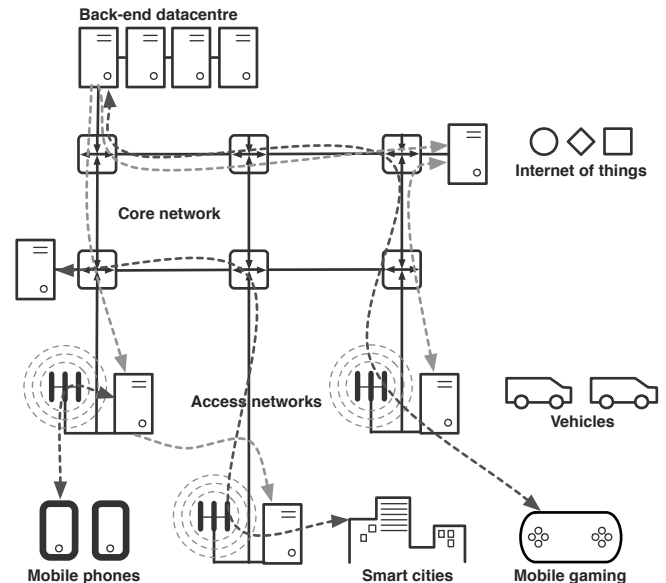


Fig. 1: x-cloud

constitutes a conscious geographic placement of the data centre to minimize the propagation delay, on the scale of hundreds of kilometres. ...

The number of devices accessing a service in a geographic area contributes to congestion on the access links. More and more devices will access the internet and so also cloud services the wireless medium. As a result, the access links will be suffer from an increasing level of congestion. ...

The forthcoming mobile networks are protectedly distributed with virtualized centralized resources. The scale of which is proportional to the acceptable distance between the antenna and the aggregated centralized resources. ...

All of the above factors provide just as many data centre placement possibilities. It is the purpose of our our research to ...

A. The bandwidth case for x-cloud

- Internet structure, latency, and bandwidth: [?]

B. The latency case for x-cloud

The intermediate latency between a client and a data centre is a product of propagation, modulation, and network routing

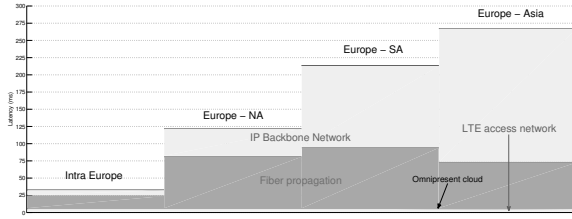


Fig. 2: IP Internet latency in western Europe [?] over LTE [?]

and traffic shaping. Propagation is a clear physical obstacle to reducing latency, and there is very little evidence to suggest that information will propagate faster than $\frac{2}{3}$ of the speed of light, at scale, in the near future. Furthermore, the delay in the backbone network is incurred to the most part by routing. A full point to point network where the propagation speed is the only limit, is not economically viable and would dissolve the fabric of the Internet. As such, we can always expect a certain amount of network contributed latency and jitter. At best, an LTE mobile access network adds about 5 ms of latency [?]. Radio access network latency can be expected to diminish over the next few generations of mobile networks.

Moving the cloud data centres closer to the IP backbone networks eliminates some of the additive latency on one side of the connection. Doing so, not only eliminate the propagation delay, but will over time, add more complexity to peripheries of the backbone as more servers nodes make their home there.

The x-cloud remedies this latency challenge in a more sustainable way. By moving compute resources to the mobile networks, IP backbone network propagation and routing delays are eliminated without disrupting the Internet topology. The resulting distributed infrastructure is capable of delivering content and services at latencies less than 10 ms.

The x-cloud will thus enable latency-sensitive services to be migrated to the cloud, such as, gaming, financial trading, process control, and most real-time human-machine interaction process.

C. The infrastructure case for x-cloud

Distributed virtualized mobile networks will rely on centralized compute nodes for higher level link management. One node will proposedly host multiple base stations, to which they connect over a network link, much like the Ericsson Radio Dot System [?], but at a larger scale. The size of these virtualization resource nodes is proportional to the maximum distance they can reside from the radio nodes, given the induced propagation delay. Supposedly these virtualization resource nodes will be placed in the vicinity of the core IP network. The virtualization resource nodes can be seen as to define geographic areas whose boundaries are defined by the reach of the mobile network which it serves. Depending on the level of desired provision and load balancing flexibility, these geographic domains will overlap to varying degrees.

The virtualization resource nodes are conceivably constructed of generic x86 or ARM servers, hosting VMs or

containers within which the virtualized mobile network infrastructure is executed. Given the placement of the virtualization resource node, any free or designually excess capacity can be used for other services.

The topology is designed to optimize the use of radio resources, the geographic domains which the virtualization resource nodes constitute do not necessarily overlap or map the demographic area which x-cloud services operate.

D. Services in the x-cloud

Placing capacity in the capillaries of the network does not necessarily have to be a means to reduce latency and congestion, but also a vessel within which localized services can be hosted. Localized services can exist purely in the x-cloud, serve only one user and geographically migrate with the user, be fully or partially hosted in the x-cloud. Different tiers can be proposedly hosted at distributed, depending on the user's mobility behaviour. Conceivably, such services could range from compute offloading, game rendering, emergency services, sensor and traffic monitoring, low-latency process control, such as surgical robotics.

III. DESIARD MODEL

The below models are intended to cater for all research topics.

- no sociogeographic model, however the model will acknowledge the primes of service sociogeographic domain, and will thus be geographically bounded appropriately.
- ...

IV. PROPOSED RESEARCH TOPICS

A. Placement of edge data centres Ericsson

1) *Proposed research:* There are no clear directions as to what degree the coming mobile networks will be virtualized. The degree of virtualization will determine the distribution of compute resources in the network, bounded by properties such as propagation delay, and cell resource provisioning.

The geographic domain in which users move, the location and size of the x-cloud hosting entity defines the bounds in which the service can perform optimally.

We propose a service performance study into the placement of the x-cloud resources. The study will contrast the service delay performance with the placement of the x-cloud resources ranging from the radio base station to the adjacent core network, where there the x-cloud node caters for multiple base stations.

The service delay will be determined by the additive latency in the mobile network, the level of congestion in the x-cloud node, and the resource shift instigated by user mobility.

We propose to contract the latency experienced by the user with the additional load imposed on the x-cloud infrastructure, as a form of utility.

2) *Related research:*

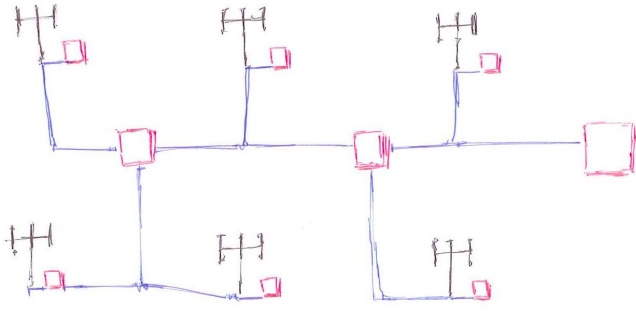


Fig. 3: x-cloud placement

B. VM placement and migration decisions Umeå

1) *Proposed research:* In the x-cloud a service will either exist only locally or distributed, but purposely serves a geographically and demographically bounded populous. The units on which the service is hosted are of limited capacity and cannot be universally virtualized as a traditional data center, with relatively unlimited capacity over time. Give the load on each of the nodes and their geographic relevance to the demographic populous they server, services will conceivably need to be mobile, migrating, dispersing, and contracting to minimize such properties as cost, load, and traffic, also taking into account the cost of the migration itself.

We propose research into an appropriate centralized or distributed, load balancing cost function, taking into account :

- Incurred migration load on host and receiver
- Incurred migration induced network congestion
- Network congestion
- Delay RTT to client to aggregate client base. In other words, minimize delay to aggregate delay/latency to all its served clients. *Perhaps separate research topic preceding this one*
- Client mobility

A property such as energy consumption is perhaps irrelevant to a topology of small data centres as domain of movement is fairly limited and bounded by the demographic service, the rate of which a service migrates is to some extent bounded by the mobility of its users. Nevertheless, energy can become a relevant parameter if the service is allowed to migrate between the x-cloud and a traditional data centre or if the energy profiles of the local x-cloud hosts, is heterogeneous. As the serviced domain is bounds each service by its sociogeographic profile and the fact that latency gains are fairly small accounting for thermal emissions and reuse would be counter productive, and should be dealt with optimally by each node independently. When optimizing for latency, inherent thermal efficiency will conceivably seldom correlate with the service sociogeographic domain.

2) *Related research:* Existing research in this area is mainly focuses on load balancing and provisioning between larger distributed data centres with static users.

C. How many devices can be handle by the current infrastructure (latency/bandwidth limits or other requirements) Umeå

1) *Proposed research:* The number of devices accessing the cloud services is bound to grow significantly with the advent of the internet of things. Most of these devices will communicate over the mobile network. Although the added devices will not necessarily contribute as much data traffic as a user-interaction device, they will significantly add to the congestion to the radio access networks and the adjacent core network. In such a case the, the congestion might result in increased latency in the mobile access network and the adjacent core network.

2) *Related research:*

D. Influence of huge number of (mobile) devices and internet of things Umeå

1) *Proposed research:*

2) *Related research:*

E. Size of edge data centres (number of CPU, memory, etc.) Ericsson

1) *Proposed research:*

2) *Related research:*

F. What type of traffic should be directed to edge data centres? Lund

1) *Proposed research:*

2) *Related research:*

G. Topology of mobile network (antennas detached from BTS) CRAN Lund

1) *Proposed research:*

2) *Related research:*

H. How mobility of user affects network and cloud computing? Lund

1) *Proposed research:*

2) *Related research:*

V. SIMULATION

A. Constituent models

1) *Data Centre:* Cloud model : [?] Web serve model : [?]

2) *Radio access network:* Because the mobile access network is a service access qualifier, the mechanisms of the network is relatively irrelevant to the primary research topics. The network can appropriately be modelled with a series of delays.

It is a constituent objective of our research to determine relevant X-Cloud/NGN ¹ symbiotic topologies. These topologies will be feed into the simulation model and will conceivably

¹Next Generation Network

encompass, resource placement and dimension, cell sizes, and radio resource provisioning [?], [?], [?].

Our basic research topics will require a homogeneous, equidistant, and equirange cell topology. Although it is reasonable to assume that future networks hosting an X-Cloud will be distributed.

3) *Base station*: The base station can conceivably be modelled with a queue and a delay proportional to its propagation distance to its associated "C-RAN" node.

4) *Core network*: The essential property of the core network is bandwidth and delay. Both of which can be modelled with queues.

- Latency
 - "Point-to-Point" core network delay model [?]
 - "One-hop" core network router queue delay model [?]

5) *Mobility*: A smooth random walk, unobstructed, bounded, edge-aware mobility model will provide a uniformly distributed dispersion of users across the simulation domain [?]. The model is two-dimensional and provides pedestrian, bicycle, and auto mobile mobility modes. The model is uniform and does thus not take into account any socio-demographic variations, and local clusters. Nevertheless, exploring specific demographic and urban settings is beyond the scope of our basic research topics. Furthermore, in the absence of a socio-demographic and urban scenarios, an aggregate mobility mode will be deployed.

6) *Service*: There is a multitude of appropriate service models.

- Light weight 1-tier web service model from 1998 [?]
- Modern light weight 1-tier web service model [?]
- YouTube workload generator [?]
- 3-tiered open-loop web service model [?]
- Web browsing behaviour : [?]
- Cloud service usage patterns : [?]

B. Simulation framework

Below are the candidate simulation tools and frameworks proposed during the third Cloud Control Workshop. [?]

1) *SimJava*:

2) *SimPy*: Python and SimPy has the ability to run powerful statistical analyses with R [?], interact with a MATLAB workspace [?], and bind NS-3 modules [?]. Nevertheless, not able to confirm whether or not you can call uncompiled MATLAB SimEvent modules.

3) *CloudSim*: [?]

CloudSim adaptations:

- NetworkCloudSim [?]
- CloudAnalyst [?]

4) *GreenCloud*: [?]

5) *iCanCloud*: [?]

6) *MDCSim*: [?]

7) *SimGrid*: [?]

• [?]

• [?]

8) *CoolSim*: [?]

9) *ns-3*: [?]

10) *Matlab+SimEvent*: (and TrueTime)

VI. PAPERS

A. Comparison of existing simulators from the perspective of x-cloud

A survey of existing simulators with comparison of their capabilities (and limitations) to simulate x-cloud .

Simulators of:

- Data centers,
- BTS,
- Network (BTS — DC),
- Mobile network,
- Mobile devices,
- Users (mobility).

What is different in operation/simulation of x-cloud ?

B. Limitations of current infrastructure & the setup/structure of x-cloud

1) *Limitations of current infrastructure*: Simulate the current infrastructure (mobile network + remote/big data centers) and show the limits of it.

- What will happen when the number of mobile devices increases by order(s) of magnitude? The influence on a network connection between a base station and a big/remote data center.
- How many mobile devices can be handled by the current infrastructure (depending on a latency limits)?

2) *The setup/structure of x-cloud* : The x-cloud consists of antennas, small (edge) data centers and big (remote) data centers. Small data centers are located close to the antennas and can host both virtualized base station software and VMs with applications. Big data centers are located far away from users. Small data centers have smaller amount of resources than big ones (maybe also performance is lower) and running applications there is more expensive. However, latency is much lower than in a case of big (remote) data centers.

Questions about the setup/structure of x-cloud :

- how many antennas should be associated with one small (local) data center? (probably this will be limited by the latency between an antenna and a small data center)
- how big should small (local) data centers be (#CPUs etc)?

C. x-cloud model

D. Throughput and bandwidth limitations in the x-cloud

E. Virtual Machine placement and migration in x-cloud

Regarding placement of Virtual Machines (VMs) in the edge data centers:

- Should a VM that serves all users (even these outside of the range of the directly connected antennas) be placed in an edge data center or should it be rather an additional instance that serves users that are in the close proximity (duplicating a VM in a big data center)?
- When a VM should be placed/duplicated in a small data center?
- While users are moving from one antenna to another when VM should be migrated from one edge data center to another one?

F. Other thoughts

- Different workload patterns?
- How to perform monitoring? (System is very distributed)
- Maybe new metrics to monitor (eg. distance from antenna, velocity, direction, etc.)
- Changes in the architecture of mobile applications

REFERENCES

- [1] Cloudsim. Available online at <http://www.cloudbus.org/cloudsim/>.
- [2] Coolsim. Available online at <http://www.coolsimsoftware.com/>.
- [3] Greencloud. Available online at <https://greencloud.gforge.uni.lu/>.
- [4] icancloud. Available online at <http://www.arcos.inf.uc3m.es/~icancloud/Home.html>.
- [5] Ns-3. Available online at <http://www.nsnam.org/>.
- [6] Ns-3 python interaction. Available online at <http://www.nsnam.org/docs/manual/html/python.html>.
- [7] pymatlab. Available online at <https://pypi.python.org/pypi/pymatlab>.
- [8] R project. Available online at <http://www.r-project.org/>.
- [9] Simgrid. Available online at <http://simgrid.gforge.inria.fr/>.
- [10] Monthly network summary, 05 2014.
- [11] Paul Barford and Mark Crovella. Generating representative web workloads for network and server performance evaluation. *ACM SIGMETRICS Performance Evaluation Review*, 26(1):151–160, 1998.
- [12] 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.
- [13] T Blajić, D Nogulić, and M Družijanić. Latency improvements in 3g long term evolution. *Mipro CTI, svibanj*, 2006.
- [14] Laurent Bobelin, Arnaud Legrand, David A González Márquez, Pierre Navarro, Martin Quinson, Frédéric Suter, and Christophe Thiéry. Scalable multi-purpose network representation for large scale distributed system simulation. In *Proceedings of the 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (ccgrid 2012)*, pages 220–227. IEEE Computer Society, 2012.
- [15] Jianhua Cao, M. Andersson, C. Nyberg, and M. Kihl. Web server performance modeling using an m/g/1/k*ps queue. In *Telecommunications, 2003. ICT 2003. 10th International Conference on*, volume 2, pages 1501–1506 vol.2, Feb 2003.
- [16] H. Casanova, A. Legrand, and M. Quinson. Simgrid: A generic framework for large-scale distributed experiments. In *Computer Modeling and Simulation, 2008. UKSIM 2008. Tenth International Conference on*, pages 126–131, April 2008.
- [17] 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.
- [18] S.K. Garg and R. Buyya. Networkcloudsim: Modelling parallel applications in cloud simulations. In *Utility and Cloud Computing (UCC), 2011 Fourth IEEE International Conference on*, pages 105–113, Dec 2011.
- [19] Phillipa Gill, Martin Arlitt, Zongpeng Li, and Anirban Mahanti. Youtube traffic characterization: A view from the edge. In *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement, IMC '07*, pages 15–28, New York, NY, USA, 2007. ACM.
- [20] H. Khazaei, J. Misić, and V.B. Misić. Performance analysis of cloud computing centers using m/g/m+m+r queueing systems. *Parallel and Distributed Systems, IEEE Transactions on*, 23(5):936–943, May 2012.
- [21] Raymond Kwan, Rob Arnott, Robert Paterson, Riccardo Trivisonno, and Mitsuhiro Kubota. On mobility load balancing for lte systems. In *Vehicular Technology Conference Fall (VTC 2010-Fall), 2010 IEEE 72nd*, pages 1–5. IEEE, 2010.
- [22] Jeongeun Julie Lee and Maruti Gupta. A new traffic model for current user web browsing behavior. *Intel Corporation*, 2007.
- [23] Seung-Hwan Lim, B. Sharma, Gunwoo Nam, Eun Kyoung Kim, and C.R. Das. Mdcsim: A multi-tier data center simulation, platform. In *Cluster Computing and Workshops, 2009. CLUSTER '09. IEEE International Conference on*, pages 1–9, Aug 2009.
- [24] Chao Liu, Ryen W White, and Susan Dumais. Understanding web browsing behaviors through weibull analysis of dwell time. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*, pages 379–386. ACM, 2010.
- [25] Xue Liu, J. Heo, and Lui Sha. Modeling 3-tiered web applications. In *Modeling, Analysis, and Simulation of Computer and Telecommunication Systems, 2005. 13th IEEE International Symposium on*, pages 307–310, Sept 2005.
- [26] 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.
- [27] Andras Racz, Andras Temesvary, and Norbert Reider. Handover performance in 3gpp long term evolution (lte) systems. In *Mobile and Wireless Communications Summit, 2007. 16th IST*, pages 1–5. IEEE, 2007.
- [28] Venugopalan Ramasubramanian, Dahlia Malkhi, Fabian Kuhn, Mahesh Balakrishnan, Archit Gupta, and Aditya Akella. On the treeness of internet latency and bandwidth. *SIGMETRICS Perform. Eval. Rev.*, 37(1):61–72, June 2009.
- [29] 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.
- [30] Bhathiya Wickremasinghe, Rodrigo N Calheiros, and Rajkumar Buyya. Cloudanalyst: A cloudsim-based visual modeller for analysing cloud computing environments and applications. In *Advanced Information Networking and Applications (AINA), 2010 24th IEEE International Conference on*, pages 446–452. IEEE, 2010.
- [31] Gansen Zhao, Jiale Liu, Yong Tang, Wei Sun, Feng Zhang, Xiaoping Ye, and Na Tang. Cloud computing: A statistics aspect of users. In *Cloud Computing*, pages 347–358. Springer, 2009.
- [32] Wei Zhao, Yong Peng, Feng Xie, and Zhonghua Dai. Modeling and simulation of cloud computing: A review. In *Cloud Computing Congress (APCloudCC), 2012 IEEE Asia Pacific*, pages 20–24. IEEE, 2012.