



UNIVERSITY OF  
NOVI SAD

FACULTY OF  
TECHNICAL SCIENCES



# **Dynamic formation of the distributed micro clouds**

- Ph. D. Thesis -

Supervisor:

Goran Sladić, PhD, full  
professor

Candidate:

Miloš Simić

Novi Sad, 2021.

Miloš Simić: Dynamic formation of the distributed micro clouds

SRPSKI: Dinamičko formiranje distribuiranog mikro okruženja u računarstvu u oblaku

SUPERVISOR: Goran Sladić, PhD, full professor

LOCATION:  
Novi Sad, Serbia

DATE:  
Month 2021

**UNIVERZITET U NOVOM SADU**  
**OBRAZAC – 5a**  
**FAKULTET TEHNIČKIH NAUKA**

**KLJUČNA DOKUMENTACIJSKA INFORMACIJA<sup>1</sup>**

Vrsta rada:	Doktorska disertacija
Ime i prezime autora:	Miloš Simić
Mentor (titula, ime, prezime, zvanje, institucija)	dr Goran, Sladić, redovni profesor, Fakultet tehničkih nauka
Naslov rada:	Dinamičko formiranje distribuiranog mikro okruženja u računarstvu u oblaku
Jezik publikacije (pismo):	Engleski jezik, latinica
Fizički opis rada:	Stranica: 256 Poglavlja: 7 Referenci: 195 Tabela: 10 Slika: 30 Grafikona: 0 Priloga:0
Naučna oblast:	Elektrotehničko i računarsko inženjerstvo
Uža naučna oblast (naučna disciplina):	Distribuirani sistemi
Ključne reči / predmetna odrednica:	distribuirani sistemi, računarstvo u oblaku, višestruko računarstvo u oblaku, mikroservisi, softver kao servis, ivično računarstvo, mikro računarstvo u oblaku, veliki podaci, infrastruktura kao kod

---

<sup>1</sup> Autor doktorske disertacije potpisao je i priložio sledeće Obrasce:

5b – Izjava o autorstvu;

5v – Izjava o istovetnosti štampane i elektronske verzije i o ličnim podacima;

5g – Izjava o korišćenju.

Ove Izjave se čuvaju na fakultetu u štampanom i elektronskom obliku i ne koriče se sa tezom.

Rezime na jeziku rada:	U sklopu disertacije izvršeno je istraživanje u oblasti distribuiranih sistema. Predstavili smo dinamičku organizaciju geo-distribuiranih čvorova u mikro centre za obradu podataka koji formiraju mikro okruženja računarstva u oblaku kako bi pokrili bilo koje proizvoljno područje i proširili kapacitet, dostupnost i pouzdanost. Koristili smo organizaciju računarstva u oblaku kao inspiraciju, sa adaptacijama za drugačije okruženje sa jasnom podelom nadležnosti, i modelom aplikacija koji može da iskoristi novoformirani sistem. Jasna podela nadležnosti, model aplikacija i dinamička organizacijom čvorova, čine da se predstavljeni model ponude kao i bilo koji drugi uslužni servis. Takođe dajemo formalne modele za sve protokole koji se koriste za stvaranje takvog sistema.
Datum prihvatanja teme od strane nadležnog veća:	
Datum odbrane: (Popunjava odgovarajuća služba)	
Članovi komisije: (titula, ime, prezime, zvanje, institucija)	<p><b>Predsednik:</b> dr Milan Vidaković, redovni profesor, Fakultet tehničkih nauka Novi Sad</p> <p><b>Član:</b> dr Silvia Gilezan, redovni profesor, Fakultet tehničkih nauka Novi Sad</p> <p><b>Član:</b> dr Miroslav Zarić, vanredni profesor, Fakultet tehničkih nauka Novi Sad</p> <p><b>Član:</b> dr Aleksandar Milosavljević, vanredni profesor, Elektronski fakultet Niš</p> <p><b>Mentor:</b> dr Goran Sladić, redovni profesor, Fakultet tehničkih nauka Novi Sad</p>
Napomena:	

**UNIVERSITY OF NOVI SAD  
FACULTY OF TECHNICAL SCIENCES**

**KEY WORD DOCUMENTATION<sup>2</sup>**

Document type:	Doctoral dissertation
Author:	Miloš Simić
Supervisor (title, first name, last name, position, institution)	Goran Sladić, Ph.D., Full Professor, Faculty of Technical Sciences
Thesis title:	Dynamic formation of the distributed micro clouds
Language of text (script):	English language, latin script
Physical description:	Pages: 256 Chapters: 7 References: 195 Tables: 10 Illustrations: 30 Graphs: 0 Appendices: 0
Scientific field:	Electrical and computer engineering
Scientific subfield (scientific discipline):	Distributed systems
Subject, Key words:	distributed systems, cloud computing, multi cloud, microservices, software as a service, edge computing, micro clouds, big data, infrastructure as code

---

<sup>2</sup> The author of doctoral dissertation has signed the following Statements:

56 – Statement on the authority,

5B – Statement that the printed and e-version of doctoral dissertation are identical and about personal data,

5r – Statement on copyright licenses.

The paper and e-versions of Statements are held at the faculty and are not included into the printed thesis.

Abstract in English language:	<p>This thesis presents research in the field of distributed systems. We present the dynamic organization of geo-distributed edge nodes into micro data-centers forming micro clouds to cover any arbitrary area and expand capacity, availability, and reliability. A cloud organization is used as an influence with adaptations for a different environment with a clear separation of concerns, and native applications model that can leverage the newly formed system. With the separation of concerns setup, edge-native applications model, and a unified node organization, we are moving towards the idea of edge computing as a service, like any other utility in cloud computing. We also give formal models for all protocols used for the creation of such a system.</p>
Accepted on Scientific Board on:	
Defended: (Filled by the faculty service)	
Thesis Defend Board: (title, first name, last name, position, institution)	<p><b>President:</b> Milan Vidaković, PhD, Full Professor, Faculty of Technical Sciences, Novi Sad</p> <p><b>Member:</b> Silvia Gilezan, PhD, Full Professor, Faculty of Technical Sciences, Novi Sad</p> <p><b>Member:</b> Miroslav Zarić, PhD, Associate Professor, Faculty of Technical Sciences, Novi Sad</p> <p><b>Member:</b> Aleksandar Milosavljević, PhD, Associate Professor, Faculty of Electronic Engineering, Niš</p> <p><b>Supervisor:</b> Goran Sladić, PhD, Full Professor, Faculty of Technical Sciences, Novi Sad</p>
Note:	

# Acknowledgements

First of all, I would like to express my sincere gratitude to my mentor professor Goran Sladić, for his unselfish help during the process of writing this thesis and related papers.

My research would have been impossible without the aid and support of colleagues Ivan Prokić and Jovana Dedeić from Chair of Mathematics. Together we managed to finish our research and publish our first paper to the journal, and we were invited to present our research to eminent colleagues from the **Imperial College in London**. Before working with them, I had a different view of what research is, and now my view is changed for the better. I hope that we will continue to collaborate in the future, publishing more great papers, and working on interesting and challenging projects and ideas. I would also like to thank all committee members for their helpful suggestions, making this thesis better.

I am profoundly grateful to Dejan Mijić, who first introduced me to formal models, and their importance in the research area of the distributed systems. Before our talks, I had no idea how to actually test, and prove my work. I would blindly try all scenarios that come to my mind, and ultimately prove probably nothing.

I would like to thank all my friends, relatives and other people who showed me how (not) to act, and how to be a better person. To my university professors, teaching assistants and colleagues from Chair of Informatics, and Chair of Applied Computer Science who pushed me to reach my full potential, to evolve, and ultimately finish my PhD. I

would like to express my gratitude to Biljana Dudaš for all the good deeds she has done over the years, Nataša Petrović and Vanja Bogićević for doing proofreading of this thesis.

Special thanks for my *MetaDoktor* teammates Vladimir Ivančević, Marko Knežević from Faculty of Technial Sciences and Danica Mandić from Faculty of Medicine, our collaboration sparked my interest for doing serious research. I have learned so much from you all, and enjoyed the time we spent working on *Dr Warehouse* project were we ultimately won **best paper award**.

In the end, I would like to dedicate this work to my grandparents, answering their most interesting question — I have finished school finally (Unless I decide to pursue post-doc)! The journey was hard, and not so pleasant most of the time, I had much more those who congratulate, than those who support. But, "*if I have seen further it is by standing on the shoulders of Giants.*"<sup>1</sup>.

---

<sup>1</sup>Sir Isaac Newton letter to Robert Hooke, 1675



# Abstract

Cloud computing is facing some serious latency issues due to huge volumes of data that need to be transferred from the place where data is generated to the cloud. For some types of applications, this is not acceptable.

One of the possible solutions to this problem is the idea to bring cloud services closer to the edge of the network, where data originates. This idea is called edge computing, and it is advertised that it dramatically reduces the network latency as a bridge that links the users and the clouds, and as such, it makes the foundation for future interconnected applications.

Edge computing is a relatively new area of research and still faces many challenges like geo-organization and a clear separation of concerns, but also remote configuration, well defined native applications model, and limited node capacity. Because of these issues, edge computing is hard to be offered as a service for future real-time user-centric applications.

This thesis presents the dynamic organization of geo-distributed edge nodes into micro data-centers and forming micro-clouds to cover any arbitrary area and expand capacity, availability, and reliability. We use a cloud organization as an influence with adaptations for a different environment with a clear separation of concerns, and native applications model that can leverage the newly formed system.

We argue that the presented model can be integrated into existing solutions or used as a base for the development of future systems.

Furthermore, we give a clear separation of concerns for the proposed model. With the separation of concerns setup, edge-native applications model, and a unified node organization, we are moving towards the idea of edge computing as a service, like any other utility in cloud computing.

The first chapter of this thesis, gives motivation and problem are that this thesis is trying to resolve. It also presents research questions, hypotheses and goals based on these questions.

The second chapter gives an introduction to the area of distributed systems, narrowing it down only the parts that are important for further understanding of the other chapters and the rest of the thesis in general.

The third chapter shows related work from different areas that are connected or that influenced this thesis. This chapter also shows what the current state of the art in industry and academia is, and describes the position of this thesis compared to the related research as well.

The fourth chapter proposes a model that is influenced by cloud computing architectural organizations but adapted for a different environment. We present how we can separate the geographic area into micro data-centers that are zonally organized to serve the local population, and form them dynamically. This chapter also gives formal models for all protocols used for the creation of such a system with separation of concerns, applications models, and presents limitations of this thesis.

The fifth presents an implemented framework that is based on the model described in chapter three. We describe the architecture, and in detail every operation a framework can do, with all existing limitations.

The sixth chapter presents the usability of the proposed model, with possible applications that could be implemented based on the model. We also present one example of COVID-19 area traffic control in the city of Milan, Italy.

The seventh and the last chapter concludes this thesis and presents future work that should be done.

**Key words:** distributed systems, cloud computing, multi cloud, microservices, software as a service, edge computing, micro clouds, big data, infrastructure as code.



# Rezime

Razvojem i različitom primenom softverskih i hardverskih sistema menja se i način na koji ljudi komuniciraju, uče i realizuju svoje aktivnosti. Kao posledica sve veće primene ovih sistema u različitim oblastima obim obrade i količina podataka znatno su se povećali [1]. Ovo povećanje za posledicu je imalo sve veću upotrebu distribuiranih sistema da bi se ti poslovi mogli obaviti uspešno.

Proširena stvarnost (AR), igre preko mreže, automatsko prepoznavanje lica, autonomna vozila ili internet stvari (IoT) proizvode izuzetno velike količine podataka što značajno utiče na opterećenje kao i na kašnjenje [1]. Ovakvi zahtevi su izvan onoga što se centralizovanim računarskim modelima, poput računarstva u oblaku, može ponuditi [1]. Čak i mali problemi mogu dovesti do velikog zastoja u komunikaciji aplikacija i uslugama od kojih ljudi zavise.

Primer koji se nedavno desio još jedan je u nizu otkaza na Amazon Web Services (AWS) platformi [2]. Ovim je platforma bila nedostupna korisnicima i aplikacijama, a kao rezultat nedostupnosti platforme velika količina aplikacija i servisa, koji se izvršavaju preko interneta, postaje potpuno nedostupna korisnicima i na kraju neupotrebljiva.

Da bi se razumelo kako je problem nastao, potrebno je objasniti šta je zapravo računarstvo u oblaku, predstaviti organizaciju i definisati arhitekturu ovog računarskog modela.

Računarstvo u oblaku definisano je kao skup računarskih resursa koji se mogu ponuditi korisnicima kroz takozvani uslužni softver [3]. Hardver i softver u velikim centrima za obradu podataka pružaju

usluge svojim korisnicima preko interneta [4]. Resursi poput CPU-a, GPU-a, skladišta podataka i mreže mogu se koristiti za nekakvu obradu podataka, ili osloboditi i to na zahtev i po potrebi korisnika [5].

Prednost računarstva u oblaku su razni servisi koji su ponuđeni korisnicima kao usluge ili uslužni softver [3]. Tradicionalnim modelom računarstva u oblaku pružaju se ogromni procesni i skladišni resursi i to po potrebi, na zahtev korisnika, kako bi se podržale različite potrebe aplikacija. Ovo svojstvo odnosi se na sposobnost računarstva u oblaku da dozvoli korisnicima alokaciju dodatnih resursa ili oslobađanje postojećih, kako bi se podudarali sa radnim opterećenjima aplikacija [6].

Jedna od problema u ovakvim sistemima nastaje kada je potrebno da se (veliki) podaci prebace sa svog izvorišta u oblak. Ovim dolazi do velike latencije ili kašnjenja u sistemu [7]. Na primer, Boeing 787 generiše pola terabajta podataka po jednom letu, dok autonomni automobil generiše dva petabajta podataka tokom samo jedne vožnje [8].

Međutim, propusni opseg nije dovoljno veliki da bi podržao takve zahteve [8]. Prenos podataka nije jedini problem sa kojim se računarstvo u oblaku susreće. Aplikacije kao što su autonomni automobili, bespilotne letelice ili balansiranje opterećenja u električnim mrežama, zahtevaju obradu podataka u realnom vremenu da bi se ispravno donosile odluke i reagovalo na razne promene [8].

Centralizovana arhitektura računarstva u oblaku, sa ogromnim kapacitetima centara za obradu podataka, stvara efikasno upravljanje resursima. Ovom strategijom dolazi se do smanjenja administrativnih troškova celokupnog sistema [9]. Međutim, kada takav sistem dođe do svojih granica, centralizacija uzrokuje mnoge druge probleme [2, 10]. Uprkos svim prednostima ovog modela, servisi i usluge vremenom se suočavaju sa ozbiljnom degradacijom kvaliteta odziva i performansi usled velike propusnosti i kašnjenja [11]. To može dovesti do nesagledivih posledica po poslovanje, ali potencijalno i teže efekte poput uticaja na ljudske živote.

Razne organizacije koriste usluge računarstva u oblaku i oslanjaju se na njega kako bi izbegle izuzetno velike infrastrukturne investicije [12]

poput pravljenja i održavanja sopstvenih centara za obradu podataka. Oni koriste resurse koje su obezbedili drugi pružaoci usluga [13] i plaćaju shodno tome koliko vremenski koriste usluge — *a pay as you go* model.

Cilj ove teze je predstavljanje i upotreba formalnih modela na osnovu kojih se mogu opisati, i formalno verifikovati protokoli, kao i implementirati radni okvir za distribuirani sistem, koristeći geografski rasprostranjena okruženja nalik računarstvu u oblaku. Opisani sistem mogu koristiti ne samo obični korisnici, već ga i pružaoci usluga računarstva u oblaku mogu integrisati u svoju platformu i svoje servise kako bi se minimalizovao zastoj kritičnih sistemskih segmenata. Čitav sistem se može posmatrati kao skup distribuiranih mikro oblaka ili sloj obrade podataka, koji u oblak šalje samo važne podatke, smajujući troškove korisnicima, ali i obezbeđujući veću dostupnost usluga računarstva u oblaku.

Distribuirani softverski sistemi su prilično složeni za modelovanje i implementaciju. Jedan od problema u ovakvim sistemima često nastaje zbog problema u komunikaciji čvorova preko mreže koja nije sigurna i pouzdana. Poruke mogu da kasne, mogu da stignu u različitom redosledu ili da ne stignu. Takođe, čvorovi u sistemu mogu prestati sa radom potpuno nasumično stvarajući dodatne komplikacije. James Gosling i Peter Deutsch, kreirali su listu problema za mrežne aplikacije poznate kao *8 zabluda distribuiranih sistema* [14]:

- (1) **Mreža je pouzdana.** Uvek će se nešto neplanirano desiti sa mrežom koja je prilično nepouzdana - prekid napajanja, prekid kabla, problemi u okruženju, itd.
- (2) **Kašnjenje ne postoji.** Lokalno kašnjenje nije problem, ali se situacija vrlo brzo pogoršava kada se komunikacija oslanja na internet, i slučajeve gde se koristi izuzetno kompleksna mrežna komunikacija računarstva u oblaku.

- (3) **Propusnost je beskonačna.** Iako se širina propusnog opsega stalno povećava i sve je bolja i bolja, srazmerno tome raste i količina podataka koja se prebacuje na obradu ili skladištenje.
- (4) **Mreža je sigurna.** Trendovi internet bezbednosti pokazuju izuzetno veliki rast napada, a ovo još više postaje problem u računarstvu u oblaku javnog tipa.
- (5) **Topologija se ne menja.** Mrežna topologija obično je izvan kontrole korisnika, a topologija mreže stalno se menja usled brojnih razloga - dodati ili uklonjeni novi uređaji, serveri, prekidi u komunikaciji itd.
- (6) **Postoji samo jedan administrator.** Danas postoje brojni administrativi za veb servere, baze podataka, keš memoriju i slično, ali, takođe, kompanije sarađuju sa drugim kompanijama ili pružaocima usluga računarstva u oblaku.
- (7) **Troškovi transporta ne postoje.** Ova tvrdnja nikako nije tačna iz prostog razloga što moramo serijalizovati informacije i podatke koje šaljemo što dodatno troši resurse i povećava ukupno kašnjenje. Ovde nije problem samo u kašnjenju, već u tome što svaka serijalizacija informacija zahteva dodatno vreme i dodatne resurse.
- (8) **Mreža je homogena.** Danas je homogena mreža izuzetak, a ne pravilo. Postoje različiti serveri, sistemi, klijenti koji komuniciraju. Implikacija ovoga je da, pre ili kasnije, mora se pretpostaviti da je potrebna interoperabilnost između ovih sistema. Mogu se koristiti i zaštićeni protokoli koji nisu javno dostupni, koji mogu dodatno trošiti vreme. Ovi protokoli mogu ostati bez podrške, pa ih treba izbegavati.

Iako razvoj distribuiranih sistema traje već nekoliko decenija, problemi koji se javljaju prilikom njihovog razvoja su i dalje identični.



Prilikom razvoja distribuiranih sistema, programeri i dizajneri često zaboravljaju na definisane probleme, što neretko dovodi do izuzetno velikih poteškoća. Način da se to u ranim fazama otkrije jeste korišćenje formalnih matematičkih metoda za opisivanje i modelovanje ovih sistema. Ove metode sačinjavaju razne tehnike koje služe za specifikaciju i verifikaciju kompleksnih sistema i koje su zasnovane na matematičkim i logičkim principima.

Kao odgovor na probleme koji mogu nastati kao posledica kašnjenja, nedostupnosti usluga u računarstva u oblaku, malicioznih napada, ali i usled kvara nekog od resursa na mreži [2], nastala je nova paradigma tzv. ivično računarstvo (eng. edge computing - EC) [13].

EC je model u kome se procesne i skladišne mogućnosti računarstva u oblaku prebacuju u blizini izvora podataka [13]. Kao posledicu toga, računarstvo u oblaku prošireno je novim mogućnostima. Smanjuje se kašnjenje što onda dovodi do novih mogućnosti za aplikacije buduće generacije [15].

Tokom prethodnih godina, pojavili su se razni modeli koji spuštaju obradu i skladištenje podataka bliže izvoru, poput fog računarstva [16], cloudlet-a [12] i mobilnih ivičnih računara (MEC) [17]. U ovoj disertaciji, svi pomenuti modeli se smatraju ivičnim računarstvom, a njihovi cvorovi ivičnim čvorovima.

Svi pomenuti modeli koriste koncept prenosa skladišnih i procesnih mogućnosti iz oblaka bliže izvorima podataka, [18] dok su zahtevnije obrade i dalje zadržane u oblaku iz vrlo prostog razloga — dostupnost znatno veće količine resursa [15]. EC modeli uvode male servere koji se arhitekturno nalaze između izvora podataka i oblaka. Tipično je za ove servere da imaju manje mogućnosti u poređenju sa serverima u oblaku [19].

Prednost malih servera je u tome što se oni mogu nalaziti na različitim lokacijama, na primer u baznim stanicama [17], gradskim centralama, restoranima, bolnicama, školama, kompanijama, ili mogu biti rasprostranjeni po geografskim regionima, a sve to kako bi se izbeglo kašnjenje i povećala propusnost [12].

Oni mogu poslužiti kao zaštitni sloj [20] ili kao nivo obrade pre nego što podaci budu poslani u oblak. Sa druge strane, korisnici dobijaju jedinstvenu mogućnost dinamičke i selektivne kontrole informacija koje bivaju poslate u oblak.

Još jednu prednost ovih servera predstavili su Aroca i saradnici [21]. Naime, njihovi rezultati su pokazali da mali serveri zadržavaju dobre performanse prilikom pokretanja zahtevnih aplikacija i klsterskog okruženja. Malo slabije performanse pokazali su u slučaju trenutno dostupnih skladišta podataka, ali to može biti podsticaj da se polje istraživanja skladišta podataka dopuni novim modelima, optimizovanim za male servere.

Jedan model teško će odgovarati potrebama svih aplikacija u budućnosti, tako da računarstvo u oblaku ne bi trebalo da bude naša konačna granica i jedina opcija. Razni modeli, nastali na bazi malih servera, pokazuju mogućnost da se obrada podataka može obaviti bliže izvoru, dok teški proračuni mogu ostati u oblaku zbog veće dostupnosti resursa. U oblak treba slati samo informacije koje su ključne za druge usluge ili aplikacije [22], a ne sve kako predlaže standardni model oblaka.

Ideja malih servera sa različitim računskim, skladišnim i mrežnim resursima pokreće zanimljive istraživačke ideje i, kao takva, motivacija je za ovu tezu. Korišćenje resursa, koji su organizovani lokalno kao mikro oblaci, oblaci zajednice ili ivični oblaci, [23] predlažu Riden i saradnici.

Usled problema koji mogu nastati u doglednoj budućnosti korišćenjem sve većeg broja računarskih sistema koji su povezani na internet, kao i zbog ograničenja računarstva u oblaku u trenutnoj izvedbi, akademska zajednica kao i industrija počele su da istražuju i razvijaju održiva rešenja. Neka istraživanja više su usredsređena na prilagođavanje postojećih rešenja zahtevima EC-a, dok druga eksperimentišu sa novim idejama i rešenjima.

U svom radu [24] Greenberg i saradnici ističu da se mikro centri za obradu podataka ( $\mu$ DCs) koriste prvenstveno kao čvorovi u mrežama

za distribuciju sadržaja i u drugim “loše distribuiranim” aplikacijama.  $\mu$ DC su zanimljiv model u području brzih inovacija i razvoja. Autori uvode koncept  $\mu$ DC-a kao centra za obradu podataka koji se nalazi u blizini velike populacije, smanjujući pritom fiksne troškove tradicionalnih centara za obradu podataka. Samim tim, minimalna veličina  $\mu$ DC-a definisana je potrebama lokalnih korisnika [24, 25], pružajući agilnost kao ključnu karakteristiku. Ovde agilnost podrazumeva sposobnost dinamičkog rasta i smanjenja potrebe za resursima kao i upotrebe resursa sa optimalne lokacije [24].

Satyanarayanan i saradnici u svom radu [20] pokazuju da  $\mu$ DC-ovi mogu poslužiti kao zaštitni sloj. Simić i saradnici u svom radu [22] opisuju takav sistem kao nivo obrade podataka na njihovom izvoru, dok korisnici dobijaju jedinstvenu mogućnost dinamičkog i selektivnog upravljanja informacijama koje se šalju u oblak.

Zonska organizacija malih servera, koju su predstavili Guo i saradnici [26] u primeni kod pametnih vozila daje zanimljivu perspektivu o EC-u. Autori su pokazali kako modeli koji dele oblast na zone omogućavaju kontinuitet dinamičkih usluga i smanjuju primopredaju veze. Takođe, pokazali su kako da se pokrivenost malim serverima prenese na veću zonu, čime se proširuju računarska snaga i kapacitet skladištenja podataka.

EC potiče iz peer-to-peer sistema, kako su to pretpostavili Lopez i saradnici, ali ga proširuju u novim pravcima i pružaju mogućnost integracije sa računarstvom u oblaku [10].

U svom radu, Kurniawan i saradnici [27] pokazali su vrlo lošu skalabilnost u centralizovanim modelima mreža za isporuku sadržaja(CDN) u oblaku. Autori su predložili decentralizovano rešenje koristeći nano centre za obradu podataka koje čine mrežni uređaji u kući [27]. Ovi centri za obradu podataka opremljeni su, takođe, sa nešto skladišnog prostora. Pokazana je moguća upotreba nano centara za obradu podataka čak i za neke velike primene sa jednom izuzetno bitnom prednošću - mnogo manjom potrošnjom energije.

$\mu$ DC-ovi sa zonskom organizacijom servera dobra su polazna osnova

za izgradnju EC-a (koja može biti ponuđena kao servis korisnicima), ali i mikro računarstva u oblaku jer se može relativno jednostavno proširiti računarska snaga i skladišni kapacitet koji opslužuju lokalne korisnike. Međutim, da bi se to postiglo, potreban nam je dostupniji i elastičniji sistem sa manje kašnjenja.

Dizajn računarstva u oblaku je takav da svaki deo doprinosi otpornijem i skalabilnom sistemu. Regioni ili centri za obradu podataka izolovani su i nezavisni jedni od drugih, a takođe sadrže resurse koji su potrebni aplikacijama za nesmetan rad. Regioni su sačinjeni od nekoliko dostupnih zona [28] i ako neka od zona postane nedostupna, druge zone preuzimaju opsluživanje korisničkih zahteve i time ceo sistem, kao celina, nastavlja neometan rad. Po uzoru na ovakvu arhitekturu, EC bi mogli koristiti vrlo sličnu strategiju formirajući model mikro računarstva u oblaku ( $\mu C$ ), gde se mali serveri ili čvorovi grupišu u *klaster*e, a više klastera u veću logičku celinu nazvanu *region*, povećavajući dostupnost i pouzdanost sistema i njegovih aplikacija.

U tom kontekstu, pod  $\mu C$ -om smatraju se geografski rasprostranjeni distribuirani sistemi koji organizaciono liče na računarstvo u oblaku, ali se nalaze u neposrednoj blizini korisnika i opslužuju njihove zahteve. Ta razlika pruža drugačiju organizaciju nego u standardnom modelu računarstva u oblaku.

Koncept regiona u računarstvu oblaka je fizički element [28], dok se u  $\mu C$ -u pojam *region* može koristiti za opisivanje skupova *klaster*a čvorova preko proizvoljne geografske oblasti. Regioni se sastoje od najmanje jednog klastera, ali mogu se sastojati i od više njih tako da se postigne otporniji, skalabilniji i dostupniji sistem.

Da bi se osiguralo manje kašnjenje u sistemu, u normalnim okolnostima treba izbegavati veliku udaljenost između klastera. U tradicionalnom modelu računarstva u oblaku proširenje regiona zahteva fizičko povezivanje novih modula sa ostatkom infrastrukture [29], što može izazvati nedostupnost tog regiona neko vreme.

U  $\mu C$ -u regioni mogu prihvatiti nove ili osloboditi postojeće klaster

*dinamički* bez direktnog povezivanja novih modula. Više regiona čine novi logički sloj — *topologiju*.

Topologija se sastoji od najmanje jednog regiona, a može se prostirati i na više regiona. Prilikom dizajniranja topologije, posebno ako regionu treba da dele informacije ili da na neki način sarađuju, poželjno je izbegavati veliku udaljenost između regiona ako je to moguće, da bi se smanjilo kašnjenje u sistemu.

Primenom *klastera*, *regiona* i *topologije* moguće je pokriti bilo koju geografsku oblast sa sposobnošću da se smanje ili prošire postojeći *klasteri*, *regioni* pa čak i *topologije*.

Organizacija *klastera*, *regiona* i *topologija* u  $\mu$ C-u isključivo je stvar potrebe korisnika, i kao takva slična je modeliranju u sistemima velikih podataka [30, 31]. Na primer, klasteri mogu obuhvatiti čitav jedan grad ili manji koji obuhvataju uređaji u jednom domaćinstvu i sve između ova dva ekstrema. Grad bi mogao predstavljati jedan region, sa delovima grada koji su organizovani u klastere. Topologija grada može se formirati tako što se grad podeli na više regiona koji sadrže više klastera. Topologija države može se formirati tako što se kroz prirodne i/ili administrativne regije definišu regioni, iako je i bilo koja druga podela moguća.

Čvorovi unutar svakog klastera treba da izvršavaju neki od protokola za održavanje definicije klastera odnosno pripadnosti čvorova klasteru. Neki od *Gossip* protokola poput *SWIM*-a [32] mogu se koristiti u saradnji sa mehanizmima replikacije podataka [33, 34, 35] čineći ceo sistem otpornijim na potencijalne greške. Treba prihvatiti činjenicu da će čvorovi u takvom sistemu iz raznih razloga biti nedostupni. To se ne može izbeći, ali se sistem može projektovati tako da servisi ipak budu ipak budu dostupni koristeći pritom neki od kopija servisa.

U modelu koji opisuje razne resurse kao usluge [36] EC i mikro računarstvo u oblaku nalaze se između CaaS-a i PaaS-a, u zavisnosti od potreba korisnika.

Prethodno definisan model  $\mu$ C mogao bi se ponuditi kao usluga korisnicima kao i bilo koji drugi resurs računarstva u oblaku. U slučaju kada

je potrebno više resursa na jednoj strani, moguće ih je reorganizovati da se iskoriste gde su oni stvarno potrebni. Sa druge strane, kompanije koje pružaju usluge računarstva u oblaku mogu integrisati predloženi model u svoj postojeći sistem, skrivajući nepotrebnu složenost iza nekog komunikacionog interfejsa ili predloženog modela aplikacije ili usluge.

Da bi se i EC modeli mogli iskoristiti na ovaj način potrebna je jasna dinamička organizacija geografski raspoređenih čvorova, dobro definisan model aplikacija i jasno razdvajanje nadležnosti u sistemu. Kao takvi, bili bi jako složeni da se ponude kao usluga korisnicima.

EC sistemi obično postoje nezavisno jedni od drugih, rasuti bez međusobne povezanosti i saradnje. Nude ih pružaoci usluga koji koriste uglavnom ograničavaju na usluge unutar sopstvenog ekosistema, često bez mogućnosti izbora servisa van njihovog kataloga usluga.

Grupisani čvorovi treba da budu organizovani lokalno, čineći sistem kompletnim, a aplikacije dostupnijim i pouzdanijim, proširujući resurse izvan pojedinačnog čvora ili male grupe čvorova. Takav sistem treba da održava dobre performanse za izgradnju servera i klastera [21].

Da bi se postiglo takvo ponašanje, neophodno je imati dinamičko upravljanje resursima i upravljanje uređajima. Potrebno je uvek imati dostupne informacije o resursima, konfiguraciji i zauzetosti čvorova [37, 17] i klastera u celini. Tradicionalni centri za obradu podataka predstavljaju dobro organizovan i povezan sistem. Sa druge strane,  $\mu$ DC-ovi se sastoje od različitih uređaja koji to nisu [38]. Ovaj problem dovodi nas do problema kojim se bavi ova teza.

Da bi opisali fizičke usluge, Jin [39] i saradnici predlažu tri osnovna koncepta i preciziraju njihove odnose. Ovi koncepti su: **(1)** uređaji, **(2)** resursi i **(3)** servisi.

Podela nadležnosti bitan je deo svakog sistema, posebno ako se stvara platforma koja se nudi korisnicima kao usluga. Model podele nadležnosti, koji ova teza predlaže, zasnovan je na ovim konceptima, prilagođen drugačijem slučaju korišćenja i podeljen u tri sloja : **(1)** uređaje, **(2)** resurse, i **(3)** servise.

Prvi sloj čine različiti uređaji odnosno kreatori podataka i korisnici

usluga odnosno servisa. Drugi sloj predstavlja resurse, koji imaju prostorne karakteristike i ukazuju na mogućnosti za obradu odnosno skladištenje podataka čvorova na kojima se izvršavaju [39]. Programeri u bilo kom trenutku moraju znati iskorišćenost resursa kao i stanje i dostupnost aplikacija.

Resursi predstavljaju EC čvorove i, da bi čvor bio deo sistema, mora zadovoljiti četiri jednostavna pravila:

- (1) Mora biti sposoban da pokrene operativni sistem sa sistemom datoteka;
- (2) Mora biti u mogućnosti da pokrene neki od dostupnih alata za izolaciju aplikacija, na primer *container* ili *unikernel*;
- (3) Mora imati dostupne resurse za korišćenje (npr. CPU, GPU, disk itd.);
- (4) Mora imati internet vezu.

Treći sloj cine servisi koji pružaju resurse aplikacijama putem definisanog interfejsa i čine ih dostupnim preko interneta [39]. Oni odmah odgovaraju na klijentske zahteve, ako je to moguće, ili mogu da skladište pronađenu informaciju za neke buduće korisničke upite [40, 41]. Servisi koji se izvršavaju u oblaku treba da budu u stanju da prihvate unapred obrađene podatke i odgovorni su za obradu i skladištenje podataka čiji kapacitet prevazilazi mogućnosti EC čvorova. Ovi servisi takođe treba da budu zamenska opcija u slučaju da prethodno definisani sistem bude nedostupan iz bilo kog razloga.

Razdvajanjem nadležnosti modela i objedinjenjem organizacije čvorova, teži se ka pristupu EC-a kao usluge i mogućnosti dinamičkog formiranja distribuiranih mikro oblaka koji bi mogli da obrade podatke na samom njihovom izvoru.

Međutim, infrastruktura za takav sistem neće biti potpuno efikasna sve dok proces konfigurisanja i korišćenja ne bude značajno pojednostavljen [40]. Ručno podešavanje čvorova predstavlja naporan i

dugotrajan proces, naročito kada se uzme u obzir geografska rasprostranjenost dostupnih čvorova.

Model koji se predlaže u ovoj tezi rešava gorepomenuti problem pomoću dinamičkog podešavanja i formiranja klastera, regiona i topologija i oslanja se na četiri protokola:

- (1) **Provera stanja čvora** - protokol obaveštava sistem o stanju svakog čvora;
- (2) **Formiranje klastera** - protokol formira nove klastere, regione i topologije;
- (3) **Provera idempotencije** - protokol proverava da li klaster, region ili topologija postoje, i da li je potrebno pokrenuti protokol za formiranje;
- (4) **Pregled detalja** - protokol prikazuje trenutno stanje sistema korisniku kroz razne nivoe detalja.

Za formalni opis server ili čvorovi (u disertaciji ovi pojmovi se koriste naizmenično) u sistemu, korišćena je teorija skupova. Prethodno definisane protokole moguće je formalno modelirati upotrebom [42] proširenje *multiparty asynchronous session types* (MPST) [43] - klasa tipova ponašanja formirana za opisivanje distribuiranih protokola oslanjajući se na asinhronu komunikaciju.

Ova matematička teorija se može iskoristiti i za verifikaciju, da li definisani protokoli zadovoljavaju MPST sigurnost (nema dostupnog stanja greške) i napredak (akcija se na kraju izvršava, pod pretpostavkom poštenja).

Proces modelovanja odvija se u dva koraka:

- (1) **Prvi korak** u modeliranju komunikacija sistema pomoću MPST teorije predstavlja definisanje *globalnog tipa*, što predstavlja globalni opis celokupnog protokola sa neutralne tačke posmatranja.



- (2) **Drugi korak** u modeliranju komunikacija sistema pomoću MPST teorije je pružanje sintaksičke projekcije protokola na svakog učesnika u komunikaciji iskazane kao *lokalni tip*, koji se zatim koristi za proveru tipa i implementacije krajnje tačke.

Na osnovu prethodno opisanih ideja i mogućnosti, definiše se problem koji ova teza obrađuje kroz sledeća tri istraživačka pitanja:

- (1) Kako se može definisati formalan proširiv model sa jasnom pode-  
lom nadležnosti, po ugledu na organizaciju računarstva u oblaku,  
koji bi bio prilagođen drugačijem okruženju izvršavanja sa jasnim  
aplikativnim modelom koji će moći da iskoristi novu, prilagođenu  
arhitekturu?
- (2) Kako je moguće ovako organizovane čvorove predstaviti kao us-  
lugu korisnicima za razvoj budućih aplikacija, po poznatom mod-  
elu naplate po utrošku (eng. pay as you go)?
- (3) Kako se može definisati formalno ispravan model koji će biti lak  
za proširenje?

Ako su prethodna istraživačka pitanja potvrдна, onda proširenje  
nalik na oblak proširuje resurse van granica pojedinačnog čvora što  
čitav sistem, kao i same aplikacije koje bi se izvršavale u njemu, čini  
dostupnijim i pouzdanijim.

Na osnovu prethodno definisanih istraživačkih pitanja i motivacija,  
definisane su sledeće hipoteze:

- (1) **Hipoteza:** *Moguće je organizovati čvorove na standardni način,  
zasnovan na arhitekturi računarstva u oblaku i prilagođen dru-  
gačije rasprostranjenom geografskom okruženju pružajući koris-  
nicima mogućnost da na najbolji mogući način organizuju čvorove  
i klastere po raznim geografskim oblastima kako bi opsluživali samo  
lokalne korisnike u neposrednoj blizini — model distribuiranog  
mikro okruženja računarstva u oblaku.*

- (2) **Hipoteza:** *Model distribuiranog mikro okruženja računarstva u oblaku se može iskoristiti za upravljanje infrastrukturnim resursima u blizini korisnika. Ako je potrebno više resursa na jednoj strani, može se uzeti određena količina resursa i reorganizovati u skladu sa potrebama sistema, ili se mogu reorganizovati na bilo koji drugi željeni način.*
- (3) **Hipoteza:** *Moguće je predstaviti jasnu podelu nadležnosti za budući sistem, koji bi se pružio korisnicima kao usluga, i uspostaviti dobro organizovan sistem u kojem svaki deo ima intuitivnu i jasnu ulogu.*
- (4) **Hipoteza:** *Moguće je predstaviti objedinjeni model, koji podržava heterogene čvorove, sa jasnim setom tehničkih zahteva koje budući čvorovi moraju ispuniti ako žele da postanu deo sistema.*
- (5) **Hipoteza:** *Moguće je predstaviti jasan aplikativni model, intuitivan korisnicima, kako bi se mogao iskoristiti puni potencijal predstavljenog rešenja.*

Iz prethodno definisanih hipoteza izvode se primarni ciljevi ove teze pri čemu očekivani rezultati uključuju sledeće:

- (1) *Definisanje formalnog proširivog modela sa jasnom podelom nadležnosti, po ugledu na organizaciju računarstva u oblaku, koji bi bio prilagođen okruženju gde se obrada i skladištenje vrše bliže korisnicima sa jasnim aplikativnim modelom koji će moći da iskoristi novu prilagođenu arhitekturu. Ovaj cilj odnosi se na prvo istraživačko pitanje, a definisano je kroz poglavlje 4.*
- (2) *Definisani model je dostupniji i elastičan sa manje kašnjenja u poređenju sa pojedinačnim malim serverima i može se koristiti kao bilo koja druga usluga u oblaku. Ovaj cilj odnosi se na prvo istraživačko pitanje, a definisano je kroz poglavlja 4 i 6.*

- (3) *Implementacija i verifikacija prototipa kojim bi se potvrdila praktična primenljivost navedenog modela i identifikovale sve njegove prednosti i eventualni nedostaci. Ovo se odnosi na treće istraživačko pitanje i tema je poglavlja 4 i 5.*

Ovim istraživanjem predstavlja se jedno moguće rešenje za organizaciju geografski rasprostranjenog  $\mu C$ -a sa EC čvorovima, zasnovano na organizaciji računarstva u oblaku, prilagođeno obradi i skladištenju podataka u blizini korisnika, a organizovano kroz tri koncepta: *topologiju, region i klaster*.

Primenom ovih koncepata moguće je opisivanje bilo kog geografskog područja u cilju dobijanja dostupnijeg i pouzdanijeg sistema. Organizacija i reorganizacija ovih elemenata vrši se dinamički, opisom željenog stanja sistema, a veličina regiona, klastera i topologija određuje se prema potrebama klijenata koje se opslužuje.

Ishod istraživanja je proširenje nalik na računarstvo u oblaku koje proširuje resurse van granica pojedinačnog čvora što čitav sistem kao i same aplikacije koje bi se izvršavale u okviru njega čini dostupnijim i pouzdanijim. Primarna primena jeste u organizaciji geografski distribuiranih računarskih resursa na efikasan način tako da opslužuje korisnike u neposrednoj blizini.

Prototipsko rešenje bazirano na prethodnom modelu se razvija kao alat otvorenog koda tako da je njegova primena moguća u različitim sistemima od strane drugih istraživačkih ili razvojnih timova. Takođe, može se koristiti kao samostalno rešenje tamo gde se kasnije mogu dodati potrebni podsistemi, ali takođe pruža mogućnost integracije u postojeća rešenja.

Predstavljeno rešenje se može koristiti kao pomoć računarstvu u oblaku prilikom obrade velike količine podataka, kao sloj koji bi vršio preliminarnu obradu na samom izvoru podataka. Skladištenje i obrada se vrši samo za podatke od značaja, što je jako bitno za tipove aplikacija koje bi se izvršavale u realnom vremenu zato što sistem lokalno može da reaguje znatno brže nego udaljeni sistem računarstva u oblaku.

Predstavljeno je preslikavanje računarstva u oblaku na EC, i prikazan je formalni model sistema sa jasnom podelom nadležnosti i modelom aplikacije za budući EC koji bi bio ponuđen korisnicima kao usluga. Dati su primeri domena gde bi sistem mogao da se koristi zajedno sa primerima aplikacija od kojih bi korisnici imali benefit. Teza je organizovana u pet poglavlja.

U **poglavlju 1** dat je opis motivacije sa jasno definisanim istraživačkim pitanjima i hipotezama na koje želimo da odgovorimo ovom tezom.

U **poglavlju 2** predstavljen je kratak uvod u temu distribuiranih sistema, sa fokusom na područja koja su važna za razumevanje ove teze i svih njenih delova.

Pokazano je šta su distribuirani sistemi ili kako neki sistem možemo opisati ili posmatrati kao distribuirani sistem. Predstavili smo probleme koje ovi sistemi stvaraju i zašto ih je tako teško implementirati, koristiti i održavati.

Takođe, predstavljeno je nekoliko primera distribuiranih računarskih aplikacija koje možemo primeniti za efikasno iskorišćavanje velikog broja čvorova u distribuiranom sistemu. Dalje, pokazano je šta je skalabilnost i zašto je ona važna za distribuirane sisteme sa nekoliko primera organizacionih mogućnosti, poput peer-to-peer i master-slave sistema, kao i protokola za opis grupa ili zajednica čvorova koji sarađuju, a koji su važni u distribuiranom okruženju iz različitih razloga. Dati su primeri raznih varijanti računarstva u oblaku koje možemo iskoristiti za svoje potrebe.

Zatim, pokazano je nekoliko tehnika virtuelizacije koje se mogu koristiti i raspoređivanje kako aplikacija, tako i infrastrukture. Prikazane su razne tehnike bitne za raspoređivanje aplikacija i infrastrukture u okruženju računarstva u oblaku, ali i razliku između distribuiranih sistema i nekoliko modela koji se često smatraju distribuiranim poput paralelnog i decentralizovanog računarstva.

Poglavlje **poglavlju 3** daje pregled relevantne literature iz oblasti disertacije.

Prikazali smo različite platforme, gde autori menjaju ili prilagođavaju postojeća rešenja (kao što su Kubernetes ili OpenStack) da rade u oblastima poput ivičnog računarstva i mobilnog računarstva. Dalje smo predstavili implementacije nekoliko platformi koje koriste čvorove, a koje su korisnici ponudili na dobrovoljnoj bazi, da bi se izvršila nekakva obrada ili skladištenje podataka na njima, kao na primer drop computing i Nebule između ostalih.

Pokazali smo kako čvorovi mogu biti organizovani po geografskim područjima na zone, ali i kako mikro centri za obradu podataka mogu da pomognu računarstvu u oblaku da prihvata zahteve lokalnih klijenata koji koriste resurse u neposrednoj blizini. Dalje smo opisali različite tehnike kako se zadaci sa mobilnih uređaja mogu prebaciti na ivične čvorove, ali takođe i različite modele primene koji bi mogli iskoristiti ove tehnike.

Na kraju ovog poglavlja dato je poređejne rezultata i doprinosa istraživanja ove teze u odnosu na slična istraživanja.

**Poglavlje 4** čini srž ove teze. U ovom poglavlju razdvojili smo sve najvažnije aspekte koje treba da zadovoljimo kako bi se omogućilo rešavanja problema kao što su kašnjenje i obrada podataka, posebno u doba mobilnih uređaja i IoT-a.

Predloženi model zasnovan je na  $\mu$ DC-ima koji su zonski organizovani i koji će opsluživati lokalne korisnike ili korisnike u blizini. Predstavili smo model koji se zasniva na računarstvu u oblaku, ali je prilagođen za drugačiji scenario i slične slučajeve korišćenja.

Pokazali smo kako možemo dinamički formirati nove klastere, regione i topologije i kako ih možemo koristiti zajedno sa mobilnim uređajima i aplikacijama poput internet stvari (IoT). Ovaj novoformirani sistem oslanja se na jasan model podele nadležnosti, usvojen iz postojećih istraživanja i prilagođen za novu troslojnu arhitekturu. Formirani model može da služi kao sloj za obradu podataka na samom izvoru ili skoro na samom izvoru, kao sloj za zaštitu privatnosti korisnika i kontrolu sadržaja koji se šalje u oblak. Predstavljeni sistem izuzetno je prilagodljiv i podložan proširivanju prema različitim

dimenzijama, odnosno potrebama i zahtevima.

Predstavljeni model može da obhvata veće ili manje geografske regione. Predstavili smo kako programeri mogu iskoristiti novu infrastrukturu i koji sve modeli aplikacija mogu postojati, ali i kako administratori mogu rasporediti razvijene servise na novoformiranu infrastrukturu koristeći opisni ili deskriptivni model, umesto eksplicitnog slanja komandi i koraka sistemu.

Pred kraj ovog poglavlja, prikazano je kako se isti može koristiti kao sastavni deo postojećih sistema (kao skladište informacija o čvorovima) ili se može koristiti kao novi model u kom možemo razviti nove podsisteme. Predstavili smo protokole za stvaranje takvog sistema i modelirali ih koristeći formalne matematičke metode ili, konkretno, teoriju asinhronih tipova sesija.

U **poglavlju 5** pokazali smo implementirani zasnovan na prethodno predstavljenom modelu. Ovde smo takođe detaljno opisali operacije koje se mogu obaviti u prototipu, ali i kako se implementirani model uklapa i gde mu je mesto u prethodno opisanom modelu podele nadležnosti.

Dalje smo izneli rezultate naših eksperimenata u kontrolisanom okruženju kao i ograničenja implementiranog radnog okvira u trenutnoj fazi razvoja. Takođe, opisali smo moguće primene ovog sistema, ali i to gde bi ovaj model mogao da se koristi kada bi ušao u upotrebu.

**Poglablje 6** prikazuje upotrebljivost predloženog modela, sa mogućim scenarijima primene. Predstavljen je i jedan konkretan primer kontrole saobraćaja i usmeravanje ambulantnih vozila i medicinskog osoblja do najbliže bolnice u gradu Milanu, Italija, zahvaćenim COVID-19 virusom, i kako bi predloženi sistem mogao da se iskoristi da pomogne medicinskom osoblju pri lečenju pacijenata, ali i istraživačima u boljem razumevanju samog virusa.

**Poglavlje 7** predstavlja poslednje poglavlje ove teze. U ovom poglavlju sumirali smo doprinose ove teze, dali smo ograničenja predloženog modela i, ujedno, ove teze, ali i svega onoga čega moramo biti svesni **ako** ako takva tehnologija bude korišćena u realnim situacijama.

Na samom kraju poglavlja dali smo pregled šta se može uraditi u pogledu budućih pravaca istraživanja.

**Ključ reči:** distribuirani sistemi, računarstvo u oblaku, višestruko računarstvo u oblaku, mikroservisi, softver kao servis, ivično računarstvo, mikro računarstvo u oblaku, veliki podaci, infrastruktura kao kod.





# Table of Contents

<b>Abstract</b>	<b>i</b>
<b>Rezime</b>	<b>v</b>
<b>List of Figures</b>	<b>xxix</b>
<b>List of Tables</b>	<b>xxxi</b>
<b>Listings</b>	<b>xxxiii</b>
<b>List of Equations</b>	<b>xxxv</b>
<b>List of Abbreviations</b>	<b>xxxvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem area . . . . .	2
1.2 Motivation and Problem Statement . . . . .	3
1.3 Research Hypotheses, and Goals . . . . .	6
1.4 Structure of the thesis . . . . .	7
<b>2 Distributed computing</b>	<b>9</b>
2.1 Distributed systems . . . . .	10
2.1.1 Scalability . . . . .	12
2.1.2 Cloud computing . . . . .	17
2.1.2.1 Multi-cloud and sky computing . . . . .	20
2.1.3 Membership protocol . . . . .	21

2.1.4	Mobile computing . . . . .	24
2.2	Distributed programming . . . . .	25
2.2.1	Big Data . . . . .	25
2.2.2	Microservices . . . . .	27
2.2.2.1	Distributed Queries . . . . .	33
2.2.3	Observability . . . . .	35
2.3	Distribution Models . . . . .	37
2.3.1	Peer-to-peer . . . . .	37
2.3.2	Master-slave . . . . .	39
2.3.3	Replication . . . . .	40
2.4	Similar computing models . . . . .	41
2.4.1	Parallel computing . . . . .	41
2.4.2	Decentralized systems . . . . .	42
2.5	Transactions . . . . .	43
2.5.1	Distributed transactions . . . . .	43
2.5.2	Sagas . . . . .	44
2.6	Garbage collection . . . . .	45
2.7	Virtualization techniques . . . . .	46
2.8	Deployment . . . . .	48
2.9	Infrastructure as software . . . . .	53
2.9.1	Infrastructure as code . . . . .	54
2.10	Development roles . . . . .	55
2.11	Concurrency vs parallelism . . . . .	57
2.11.1	Actor model . . . . .	57
<b>3</b>	<b>Research review</b>	<b>59</b>
3.1	Nodes organization . . . . .	59
3.2	Platform models . . . . .	61
3.3	Task offloading . . . . .	65
3.4	Application models . . . . .	66
3.5	Infrastructure management . . . . .	68
3.6	Thesis position . . . . .	70
<b>4</b>	<b>Micro clouds</b>	<b>71</b>

4.1	Configurable Model Structure . . . . .	72
4.2	Separation of concerns . . . . .	78
4.3	Applications Model . . . . .	80
4.3.1	Execution models . . . . .	82
4.3.2	Packaging options . . . . .	84
4.4	As a service model . . . . .	85
4.5	Immutable infrastructure . . . . .	86
4.5.1	Deployment in micro clouds . . . . .	88
4.6	Formal model . . . . .	90
4.6.1	Multiparty asynchronous session types . . . . .	91
4.6.2	Health-check protocol . . . . .	93
4.6.3	Cluster formation protocol . . . . .	99
4.6.4	Idempotency check protocol . . . . .	107
4.6.5	List detail protocol . . . . .	114
4.7	Long-lived transactions in micro clouds . . . . .	118
4.7.1	Garbage collection in micro clouds . . . . .	119
4.8	System observability . . . . .	120
4.9	Access pattern . . . . .	121
4.10	Auto scaling micro clouds . . . . .	122
4.11	User data flow in micro clouds . . . . .	123
4.12	Extendability . . . . .	124
4.13	Repercussion . . . . .	125
<b>5</b>	<b>Proof of concept</b>	<b>127</b>
5.1	Platform implementation . . . . .	127
5.1.1	Technologies . . . . .	131
5.1.2	Node daemon . . . . .	133
5.1.3	Separation of concerns details . . . . .	136
5.1.4	Long-lived transactions details . . . . .	136
5.1.5	Garbage collection details . . . . .	137
5.1.6	Limitations . . . . .	137
5.2	Operations . . . . .	138
5.2.1	Query . . . . .	138
5.2.2	Mutate . . . . .	140

5.2.3	Queueing . . . . .	143
5.2.4	List . . . . .	145
5.2.5	Logs . . . . .	146
5.3	Experimental Results . . . . .	146
5.3.1	Experiment . . . . .	147
5.4	The existing solutions enhancement . . . . .	148
<b>6</b>	<b>Model usability</b>	<b>151</b>
6.1	Applications . . . . .	151
6.2	Area traffic control example . . . . .	152
<b>7</b>	<b>Conclusion</b>	<b>155</b>
7.1	Contributions of the thesis . . . . .	155
7.2	Limitations . . . . .	158
7.2.1	Discussion . . . . .	159
7.3	Future work . . . . .	160
	<b>Bibliography</b>	<b>165</b>

# List of Figures

2.1	Difference between scaling vertically (left) and horizontally (right) . . . . .	13
2.2	Difference between cloud options and on-premises solution. . . . .	15
2.3	Difference between cloud options and on-premises solution. . . . .	18
2.4	Direct and indirect ping in SWIM protocol. . . . .	23
2.5	V's of Big Data. . . . .	26
2.6	CQRS pattern diagram. . . . .	34
2.7	API Composition pattern diagram. . . . .	35
2.8	<i>RequestX</i> path through the processes in a distributed system. . . . .	36
2.9	P2P network and client-server network. . . . .	38
2.10	Handling requests master-slave and peer-to-peer . . . . .	40
2.11	Architectural difference between DC and parallel computing. . . . .	41
2.12	Saga transactions separated in sub-transactions. . . . .	44
2.13	Architectural differences between VMs, containers and unikernels. . . . .	48
2.14	Difference between mutable and immutable deployment models. . . . .	52
2.15	Reconciliation loop diagram. . . . .	55
4.1	Three tier architecture, with the response time and resource availability . . . . .	77
4.2	ECC as a service architecture with separation of concerns. . . . .	80

4.3	$\mu$ C applications model tides traditional clouds and edge computing applications . . . . .	81
4.4	Rolling update flow chart . . . . .	89
4.5	Recreation strategy update flow chart . . . . .	89
4.6	Canary update 1:4 ration update chart . . . . .	90
4.7	Low level health-check protocol diagram. . . . .	94
4.8	Low level cluster formation communication protocol diagram. . . . .	100
4.9	Low level view of idempotency check communication. . . . .	109
4.10	Low level view of list operation communication. . . . .	115
4.11	State diagram for cluster formation message . . . . .	119
4.12	New worker addition to the system, with reconcile loop. . . . .	125
5.1	Proof of concept implemented system. . . . .	130
5.2	Low level communication protocol diagram of query operation. . . . .	139
6.1	Conceptual architecture model for COVID-19 area traffic control example . . . . .	154

# List of Tables

2.1	Differences between horizontal and vertical scaling. . .	13
2.2	Downtime for different classes of nines. . . . .	16
2.3	Comparison of public, private and hybrid cloud capabilities. . . . .	19
2.4	Common examples of SaaS, PaaS, and IaaS. . . . .	20
2.5	Differences between the monolith and microservices architecture. . . . .	28
2.6	Idempotent and non-idempotent operations. . . . .	31
2.7	Idempotent and non-idempotent operations. . . . .	32
2.8	Differences between DevOps and SREs. . . . .	56
2.9	Differences between actor model and CSP. . . . .	58
4.1	Similar concepts between cloud computing and ECC. .	73





# Listings

5.1	Actor system hierarchy. . . . .	134
5.2	Daemon configuration file . . . . .	135
5.3	Structure of stored key-value element. . . . .	140
5.4	Example of mutate file using YAML. . . . .	142
5.5	Structure of stored key-value element. . . . .	144



# List of Equations

2.1	Availability percentage formula . . . . .	15
2.2	Availability class formula. . . . .	16
2.2	Commutative formula. . . . .	17
2.2	Associative formula. . . . .	17
2.2	Idempotent formula. . . . .	17
2.3	Idempotency law formula . . . . .	32
4.1	Cluster size limits . . . . .	75
4.2	Position od micro clouds. . . . .	76
4.3	Function composition. . . . .	83
4.4	Global type construction . . . . .	92
4.5	Local type representation syntax. . . . .	93
4.6	Extending free servers set. . . . .	96
4.7	Extending label set. . . . .	97
4.7	Health-check global protocol. . . . .	98
4.7	Health-check global protocol projection. . . . .	99
4.8	Query selector formation. . . . .	102
4.11	Server selector. . . . .	103
4.12	Extending task queue set. . . . .	103
4.12	Cluster formation global protocol. . . . .	105
4.12	Cluster formation global protocol projection. . . . .	107
4.13	Idempotency structure. . . . .	108
4.13	Idempotency check global protocol. . . . .	112
4.13	Idempotency check global protocol projection. . . . .	114
4.13	List detail global protocol. . . . .	117
4.13	List detail global protocol projection. . . . .	118



# List of Abbreviations

<b>API</b>	Application programming interface
<b>AWS</b>	Amazon Web Services
<b>c12s</b>	Constellations
<b>CC</b>	Cloud computing
<b>CDN</b>	Content delivery networks
<b>CQRS</b>	Command Query Responsibility Segregation
<b>CRDTs</b>	Conflict-free replicated datatypes
<b>CSP</b>	Communicating Sequential Processes
<b>CaaS</b>	Container as a service
<b>DBaaS</b>	Databases as a service
<b>DCs</b>	Data centers
<b>DC</b>	Distributed computing
<b>DHT</b>	Distributed Hash Table
<b>DNS</b>	Domain Name System
<b>DSL</b>	Domain specific language

<b>DS</b>	Distributed systems
<b>Dp</b>	Distributed programmig
<b>ECC</b>	Edge-centric computing
<b>EC</b>	Edge computing
<b>ES</b>	Edge servers
<b>GC</b>	Garbage collection
<b>IaC</b>	Infrastructure as code
<b>IaS</b>	infrastructure as software
<b>IaaS</b>	Infrastructure as a service
<b>IoT</b>	Internet of Things
<b>k8s</b>	Kubernetes
<b>LLTs</b>	Long-lived transactions
<b>LLT</b>	long-lived transaction
<b>MANETs</b>	Mobile Ad-Hoc Networks
<b>MA</b>	Evolutionary Memetic Algorithm
<b>MCC</b>	Mobile cloud computing
<b>MEC</b>	Mobile edge computing
<b>MPST</b>	Multiparty asynchronous session types
<b>NaCl</b>	Native Client
<b>NoSQL</b>	Not Only SQL

<b>OS</b>	Operating system
<b>P2P</b>	Peer-to-peer
<b>PaaS</b>	Platform as a service
<b>QoE</b>	Quality of experience
<b>QoS</b>	Quality of service
<b>RAID</b>	Redundant Array of Inexpensive Disks
<b>REST</b>	Representational state transfer
<b>RPC</b>	Remote procedure call
<b>SDN</b>	Software-defined networks
<b>SEC</b>	Strong Eventual Consistency
<b>SPOF</b>	Single Point of Failure
<b>SRE</b>	Site Reliability Engineering
<b>SaaS</b>	Software as a service
<b>SoC</b>	Separation of concerns
<b>TTL</b>	Time to live
<b>VANETs</b>	Vehicular Ad-Hoc Networks
<b>VM</b>	Virtual machine
<b>XaaS</b>	Everything as a Service
<b>μCs</b>	Micro clouds
<b>μDCs</b>	Micro data-centers





# Chapter 1

## Introduction

Various software systems have changed the way people communicate, share information, learn, and run businesses. The interconnected computing devices have numerous positive applications in everyday life. In the past decade or so, the volumes of data that is collected, stored, and computed have grown dramatically [1].

New age applications that might include augmented reality, massive online gaming, face recognition, autonomous drones and vehicles, or the Internet of Things (IoT) produce enormous amounts of different kinds of data.

Such workloads require that latency is below a few tens of milliseconds [1], or even less. These requirements fall just right outside of what a standard centralized model like cloud computing (CC), for example, could offer [1]. Even the smallest problems can contribute to largely unplanned downtime of applications and services people and other services may depend on. A most recent example is yet another outage that happened to the Amazon Web Services (AWS), and as a result, a large amount of internet becomes unavailable [2].

This thesis aims to provide formal models based on which we can model and implement distributed systems (DS) for organizing cloud-like geo-distributed environments for users or CC providers. We can look at the whole system as a micro clouds ( $\mu$ Cs) or pre-processing

layer. The responsibility of such a system is sending only necessary and data to the cloud. This strategy reduces cost for users and ensures the availability of CC services. Such a system could lead to lowering the downtime of critical services. Here we can give a formal definition of geo-distributed environments like:

**Definition 1.0.1.** *Geo-distribution represent organizing processing and storage resources in proximity to some large populations, where  $\mu$ Cs are formed dynamically, serving users requests **locally first**.*

The section starts by describing the general problem area that our work addresses in Section 1.1. Section 1.2, specifies the exact problem that our work addresses, and Section 1.3 describes our hypothesis and research goals. Section 1.4 presents the structure of the thesis.

## 1.1 Problem area

To lower the administration cost, cloud providers create an effective economy of scale [9] by housing data-centers (DCs) with huge capacities. However, such a model does not come without the cost. When such a system grows to its physical limits, a centralized model brings more harm than good [2, 10].

Despite all the benefits that centralization can provide, it is inevitable for the CC services to suffer from serious problems [11]. Over time, and due to the high bandwidth and latency, these services face degradation that we cannot overlook. This serious services degradation can have an enormous consequence on the human business and potentially lives as well [44].

To avoid large investments [12], like creating and maintaining their own DCs, organizations use cloud services created by others [13]. They consume resources and pay for their usage time. This model is known as — *pay as you go model*.

CC requires data transfer to the DCs from data sources. This operation is problematic because it creates a high latency in the system [7].

We can observe some examples like data collection from planes and autonomous cars. Boeing 787s per single flight generates half a terabyte of data, while a self-driving car generates two petabytes of data per single drive. On the other hand, bandwidth is not large enough to support such requirements [8].

Data transfer is not the only problem CC is facing. Some applications require real-time processing for proper decision-making [8]. For example, self-driving cars, delivery drones, or power balancing in electric grids. Such applications might face serious issues if a cloud service becomes unavailable due to whatever reason [2].

Over the years, research led to new computing areas and models in which computing and storage utilities are in proximity to data sources [13]. To overcome cloud latency issues centralized CC model is enhanced with some new ideas [15].

## 1.2 Motivation and Problem Statement

In their work [24] Greenberg et al. point out that micro data-centers ( $\mu$ DCs) are used primarily as nodes in content distribution networks and other “embarrassingly distributed” applications.

One size rarely suits all needs, so the CC should not be our final computing shift. Various models presented in 2.1.4 show a promising possibility of how computing could be done closer to the data sources, to lower the latency for its clients by contacting the cloud only when needed, while the heavy computation could remain in the cloud, because of more available resources. Send to the cloud only information that is crucial for other services or applications [22]. Not ingest everything as the standard cloud model proposes.

A zonally-based organization of servers combined with  $\mu$ DCs shows a great possibility for building  $\mu$ Cs and EC as a service. To achieve such a behavior, a few more abstractions and layers are needed, to make the whole system more available, resilient, and with less latency. By their nature, EC originates from P2P systems [10] as suggested by

López et al., but expands it into new directions and blends it with the CC.

This is very interesting, because there is much research and knowledge available for P2P systems that could be used for inspiration. One extension of the P2P system leads to geo-distributed deployments, and going from node to node is a time-consuming process. Satyanarayanan et al. stated that infrastructure deployment will not happen until the whole process is relatively easy [40].

Assuming that we have a well-defined system, and so that infrastructure can be easily deployed and operated with, we have a system that could be offered like any other resource in the CC – *as a service*. Such a system or service could be offered to various types of users, from researchers to developers to create new types of applications. A well-defined system that is easy to operate with, will be able to move resources from one place to another with no problem. Some cloud providers, though, might choose to integrate the system into their existing CC platform to reduce the load or avoid bottlenecks and single points of failure [45].

The idea of small-scale servers introduced by edge computing (EC), with heterogeneous, compute, storage, and network resources, raise interesting research ideas and it is the main motivation for this thesis. Taking advantage of resources organized locally as  $\mu$ Cs, community clouds, or edge clouds [23, 46] suggested by Ryden et al., to help power-hungry servers reduce traffic [47]; contacting the cloud only when needed [22]; sending to the cloud only information that is crucial for other services or applications; not ingesting everything as the standard CC model proposes.

To achieve such behavior, dynamic resource management, and device management is essential. At any given time, available resources, configuration, and utilization of the system must be known [37, 17]. Traditional DCs is a well organized and connected system, built upon years of experience. On the other hand, these  $\mu$ DCs consist of various devices, including ones presented in 2.1.4 that are not [38]. This idea

brings us to the problem this thesis addresses.

Currently, existing EC and  $\mu$ DCs models lack dynamic, well defined geo-organization structure, native applications model, and a clear separation of concerns model. As such they cannot be offered as a service to the users, and it is hard to form  $\mu$ Cs on them. Usually, these systems exist independently from one another, scattered without any communication between them. Providers who build and maintain these systems, usually lock users in their ecosystem, frequently integrate tightly with their cloud services, giving users small or even no other options. Nearby EC nodes could be organized locally, making the whole system more available and reliable, but at the same time extending resources beyond the single node or group of nodes, maintaining good performance to build servers and clusters [21].

This cloud extension strengthens our understanding of not just DS but also CC as a system and field of research. With EC native applications model, separation of concerns well defined, and a unified node organization strategy, we are moving towards the idea of EC as a service and  $\mu$ Cs [46].

Based on this, the problem this thesis is trying to solve is defined by the three research questions:

- (1) *Can geo-distributed EC nodes be organized in a similar way to the cloud, but adopted for the different environment, with clear separation of concerns and familiar applications model for users forming  $\mu$ C a model?*
- (2) *Can these organized nodes ( $\mu$ Cs) be offered as a service based on the cloud pay as you go model, to the developers and researchers so that they can develop new human-centered applications?*
- (3) *Can a model be made in such a way that is formally correct, easy to extend, understand and reason about?*

This  $\mu$ C model makes both system and applications more available and reliable, while resources are extended further beyond the single node or

even  $\mu$ DCs. In their work [20] Satyanarayanan et al. shows that  $\mu$ DCs can serve as firewalls, while users get a unique ability to dynamically and selectively control their information that will be uploaded to the cloud. Simić et al., in [22] uses a similar idea as a pre-processing tier for the cloud. Years after its inception, EC is no longer just an idea [20] but a must-have tool for novel applications to come.

### 1.3 Research Hypotheses, and Goals

Based on research questions presented on the page 5, and the motivation presented in section 1.2, the hypothesis around which this thesis is based, is derived. They can be summarized as follows:

- (1) **Hypothesis:** *It is possible to organize EC nodes in a standard way based on cloud architecture, adapted for EC geo-distributed environment –  $\mu$ Cs. Giving users the unique ability to organize nodes, descriptively, in the best possible way to serve the population nearby.*
- (2) **Hypothesis:** *It is possible to offer a newly formed  $\mu$ C model to researchers and developers as a service based on the cloud pay as you go model to create new human-centered applications, but sustain the ability to rearrange resources as needed.*
- (3) **Hypothesis:** *It is possible to form a clear separation of concerns for the future  $\mu$ C model (EC as a service model) and establish a well-organized system with an intuitive role for every part.*
- (4) **Hypothesis:** *It is possible to present a unified model that will be able to support heterogeneous small-scale servers (EC nodes). This unified model will rely on a set of technical requirements that nodes must fulfill to join the system.*

- (5) **Hypothesis:** *It is possible to present a clear and familiar application model so that users can use the full potential of newly created infrastructure but familiarly and intuitively.*

From the previously defined hypotheses, the primary goals of this thesis can be derived, where the expected results include:

- (1) *The construction of a model with a clear separation of concerns for the model influenced by cloud organization, with adaptations for a different environment, and with a model for EC applications utilizing these adaptations. This addresses the first research question and is the topic of Chapter 4.*
- (2) *The constructed model is more available, resilient with less latency, and as such it can be offered to the general public as a service like any other service in the cloud. This addresses the second research question and is the topic of Chapters 4 and 6.*
- (3) *The constructed model is described formally well, using solid mathematical theory, but also easy to extend both formally and technically, easy to understand and reason about. This addresses the third research question and is the topic of Chapters 4 and 5.*

### 1.4 Structure of the thesis

Throughout this introductory Chapter 1, the motivation for this work is defined, with problems that this thesis addresses. The rest of the thesis is outlined here.

Chapter 2 presents the necessary background details, informations, and areas for understanding and supporting the rest of the thesis.

Chapter 3 presents the literature review, where different aspects of existing systems and methods important for the thesis are examined. Existing organizational abilities for nodes in both industry and academia frameworks are analyzed, as well as solutions to address the

first research question. Platform models from industry and academia tools and frameworks are examined further to address the second research question. And last but not least, current strategies to offload tasks from the cloud are examined. All three parts address the third research question.

Chapter 4 details the model, how it is related to other research and where it connects to other existing models and solutions. The solution, as well as protocols required for such a system to be implemented formally are further described.

Chapter 5 presents implementation details of a framework implemented to test defined hypotheses, and answer research questions. The framework is implemented based on the protocols formally described in chapter 4. This chapter also shows results after conducting experiments.

Chapter 6 presents usability of the proposed model, and possible applications that could benefit from such a system. There are also examples of how existing infrastructure could be used, as well as familiar application model for developers.

Chapter 7 is the final chapter, and it concludes the work of this thesis. It summarizes the contributions of this thesis, presenting the limitations of the proposed model, and outlines the opportunities and directions for further research and development in this area.



# Chapter 2

## Distributed computing

An overview of the distributed computing (DC) topics that are of significant importance for the rest of the thesis is going to be given in this section, since the thesis is heavily based on these topics.

Sections [2.1](#) and [2.2](#) describe the theoretical background behind the problem, where we examine distributed systems (DS) and distributed programming (DP), focusing on design details, communication patterns, and organizational structure. Section [2.4](#) describes similar models that might be a source of confusion, and how they are different than DS or DC, and how some concepts can fit in the bigger picture. Section [2.5](#) describes different transaction models used for different applications. Section [2.6](#) describes basics of garbage collection techniques and why it is important. Section [2.7](#) describes different virtualization methods that are used in CC for systems and/or applications. Section [2.8](#) describes different architecture and application model and how deployment can be done in large DS. Section [2.9](#) describes infrastructure as software model that allows abstracting infrastructure to software level. Section [2.10](#) describes different development roles in the modern complex software environment, with focus on technical roles, while Section [2.11](#) describes the difference between concurrency and parallelism and introduces an actor system, that will be used later on in the thesis.

## 2.1 Distributed systems

There are various definitions of DS, but we can think of DS as a system where multiple entities can communicate to one another in some way, but at the same time, they can perform some operations. In [48, 49] Tanenbaum et al. give two interesting assumptions about DS:

- (1) “A computing element, which we will generally refer to as a node, can be either a hardware device or a software process”;
- (2) “A second element is that users (be they people or applications) believe they are dealing with a single system. This means that one way or another the autonomous nodes need to collaborate”.

These two assumptions are useful and powerful when talking about DS. As such, in this thesis, we will adopt and use them rigorously.

Three significant characteristics of distributed systems are [49]:

- (1) **Concurrency of components**, refers to the ability of the DS that multiple activities are executed at the same time. These activities take place on multiple nodes that are part of a DS;
- (2) **Independent failure of components**, this property refers to a nasty feature of DS that nodes fail independently. They can fail at the same time as well, but they usually fail independently for numerous reasons;
- (3) **Lack of a global clock**, this is a consequence of dealing with independent nodes. Each node has its notion of time, and as such we cannot assume that there is something like a global clock.

In [48] authors give formal definition “distributed system is a collection of autonomous computing elements that appears to its users as a single coherent system”.

When talking about DS, we usually think about computing systems that are connected via network or over the internet. But DS is not

exclusive to the domain of computer science. They existed before computers started to enrich almost every aspect of human life. DS have been used in various different domains such as: **telecommunication networks, aircraft control systems, industrial control systems** etc. DS are used anywhere where the number of users is growing rapidly so that a single entity cannot respond to the demands in (near) real-time.

Distributed systems (in computer science) consist of various algorithms, techniques, and trade-offs to create an illusion that a set of nodes act as one. Algorithms and techniques used in the DS may include the following: **(1)** replication, **(2)** consensus, **(3)** communication, **(4)** storage, **(5)** processing, **(6)** membership, **(7)** scheduling etc.

DS are hard to implement because of their asynchronis and faulty nature. James Gosling and Peter Deutsch both fellows at Sun Microsystems at the time created a list of problems for network applications known as *8 fallacies of Distributed Systems* [14]:

- (1) The network is reliable;** there will always be something that goes wrong with the network — power failure, a cut cable, environmental disasters, etc.;
- (2) Latency is zero;** locally latency is not an issue, but it deteriorates very quickly when you move to the internet and CC scenarios;
- (3) Bandwidth is infinite;** even though bandwidth is constantly getting better and better, the amount of data we try to push through it rise as well;
- (4) The network is secure;** Internet attack trends are showing growth, and this becomes a problem even more in public CC;
- (5) Topology doesn't change;** network topology is usually out of user control, and network topology changes constantly for numerous reasons — added or removed new devices, servers, breaks, outages, etc.;

- (6) **There is one administrator**; nowadays there are numerous administrators for web servers, databases, cache and so on, and companies collaborate with other companies or CC providers;
- (7) **Transport cost is zero**; we have to serialize information and send data over the wire, which takes resources and adds to the total latency. The problem here is not just latency, but that information serialization takes time and resources;
- (8) **The network is homogeneous**; today, a homogeneous network is the exception, rather than a rule. We have different servers, systems, clients that interact. The implication of this is that we have to assume interoperability between these systems sooner or later, which we must be aware of. We might also have some proprietary protocols that might also take time to send on and they may stay without support, so we should avoid them.

These fallacies were introduced over a decade ago, and more than four decades since we started building DS, but the characteristics and underlying problems remain pretty the same. Because of these problems, DS are hard to implement correctly and they are hard to maintain and test properly.

### 2.1.1 Scalability

Scalability is the property of a system to handle a growing amount of work by adding resources to the system [50]. When talking about computer systems, scalability can be represented in two ways:

- **Scaling vertically** means upgrading the hardware that computer systems are running on — adding more CPU, RAM, storage. Vertical scaling can increase performance to what the latest hardware can offer, and here we are limited by the laws of physics and Moor's law [51]. A typical example that requires this type

of scaling is a relation database server. These capabilities are insufficient for moderate to big workloads;

- **Scaling horizontally** means that we scale our system by adding more and more computers, rather than upgrading the hardware of a single one. With this approach, we are (almost) limitless on how much we can scale. Whenever performance degrades we can simply add more computers (or nodes). These nodes are not required to be some high-end machines.

Figure 2.1 show difference between scaling vertically and horizontally.

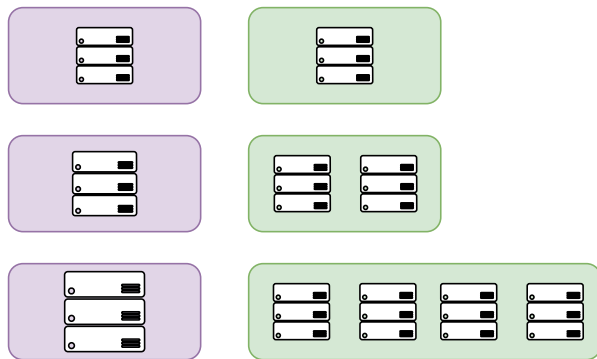


Figure 2.1: Difference between scaling vertically (left) and horizontally (right)

Table 2.1 summarizes differences between horizontal and vertical scaling.

Feature	Scaling vertically	Scaling horizontally
Scaling	Limited	Unlimited
Management	Easy	Complex
Investments	Expensive	Affordable

Table 2.1: Differences between horizontal and vertical scaling.

Scaling horizontally is a preferable way for scaling DS. Not because we can scale easier, or because it is significantly cheaper than vertical scaling (after a certain threshold) [50], but because this approach comes with few more benefits that are especially important when talking about large-scale DS. Adding more nodes gives us two important properties:

- **Fault tolerance** means that applications running on multiple places at the same time are not bound to the fail of a node, cluster, or even DCs. As long as there is a copy of the application running somewhere, the user will get a response back. As a consequence of running multiple copies of a service and on multiple places, we have that service is more **available**, than running on a single node no matter how high-end that node is. Eventually, all nodes are going to break, and if we have multiple copies of the same service we have a more resilient and more available system to serve user requests;
- **Low latency** refers to the idea that the world is limited by the speed of light. If a node running application is too far away, the user will wait too long for the response to get back. If the same application is running in multiple places, the user request will hit the node that is closest to the user.

Nodes are usually organized into clusters of machines. Buyya et al. describes a cluster as a processing system, which consists of a collection of interconnected stand-alone computers cooperatively working together as a single, integrated computing resource. [52].

But despite all the obvious benefits, for a DS to work properly, we need the writing software in such a way that is able to run on multiple nodes, as well as that it **accepts the failure and deals with it**. This turns out to be not an easy task. For example, users need to be aware when using DS which of them is related to distributed data storage systems. Storage implementations that rely on vertical scaling to ensure scalability and fault tolerance, have one nasty feature.

This nasty feature is represented in theorem called **CAP theorem** presented by Eric Brewer [53], and proven after inspection by Gilbert et al. [54]. The CAP theorem states that it is impossible for a distributed data store to simultaneously provide more than two out of three guarantees shown in Figure 2.2.

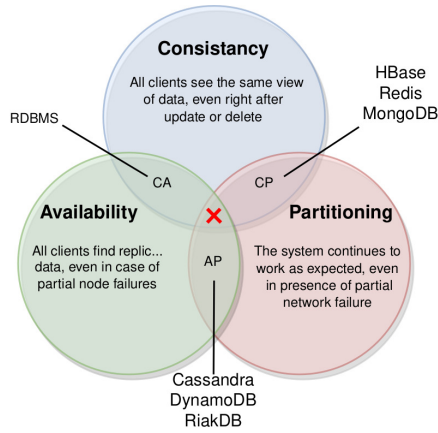


Figure 2.2: Difference between cloud options and on-premises solution.

- (1) **Consistency**, which means that all clients will see the same data at the same time, no matter which node they are connected to. Clients may not be connected to the same node since data could be replicated on many nodes in different locations;
- (2) **Availability**, which means that any client issued request will get a response back, even if one or more nodes are down. DS will not interpret this situation as an exception or error. Availability is represented in percentage, and it describes how much downtime is allowed per year. This can be calculated using formula:

$$Availability = \frac{uptime}{(uptime + downtime)} \quad (2.1)$$

The industry is using measuring availability in “class of nines”. Availability class is the number of leading nines in the availability figure for a system or module [55]. This metric relates to the amount of time (per year) that service is up and running. Table 2.2 show different classes of nine and their availability and unavailability in minutes per year (**min/year**) for some examples [55].

Type	Availability	Unavailability
<b>Unmanaged</b>	90%	50,000
<b>Managed</b>	99%	5,000
<b>Well-managed</b>	99.9%	500
<b>Well-managed</b>	99.9%	500
<b>Fault-tolerant</b>	99.99%	50
<b>High-availability</b>	99.999%	5
<b>Very-high-availability</b>	99.9999%	0.5

Table 2.2: Downtime for different classes of nines.

We can calculate availability class if we have system availability  $A$ , the system’s availability class is defined as [55]:

$$e^{\log_{10} \frac{1}{(1-A)}} \quad (2.2)$$

It is important to notice that even a 99% available system gives almost four days of downtime in a year, which is unacceptable for services like Facebook, Google, AWS, etc. And when the service is down, companies are losing customers;

- (3) **Partition tolerance**, which means that the cluster must continue to work despite any number of communication breakdowns between nodes in the system. It is important to state that in a distributed system, partitions cannot be avoided.



Years after CAP theorem inception, Shapiro et al. prove that we can alleviate CAP theorem problems, but only in some cases, and offers **Strong Eventual Consistency (SEC) model** [56]. They prove that if we can represent our data structure to be:

- **Commutative**  $a * b = b * a$
- **Associative**  $(a * b) * c = a * (b * c)$
- **Idempotent**  $(a * a) = a$

where  $*$  is a binary operation, for example: *max*, *union*, or we can rely on SEC properties,

### 2.1.2 Cloud computing

Vogels et al. describe CC as an “aggregation of computing resources as a utility, and software as a service” [3]. Big DCs provide hardware and software services for their users over the internet [4]. Cloud providers offer various resources like CPU, GPU, storage, and network as utilities that can be used and released on-demand [5].

The key strength of the CC is reflected in the offered services [3]. To support the various application needs, the traditional CC model provides enormous computing and storage resources elastically. This property refers to the cloud ability to allow services to allocate additional resources or release unused ones to match the application workloads on-demand [6].

Services usually fall in one of three main categories:

- **Infrastructure as a service (IaaS)** allows businesses to purchase resources on-demand and as-needed instead of buying and managing hardware themselves;
- **Platform as a service (PaaS)** delivers a framework for developers to create, maintain and manage their applications. All resources are managed by the enterprise or a third-party vendor;

- **Software as a service (SaaS)** delivers applications over the internet to its users. These applications are managed by a third-party vendor.

Figure 2.3 shows the difference in control and management of resources between different cloud options and on-premises solutions.

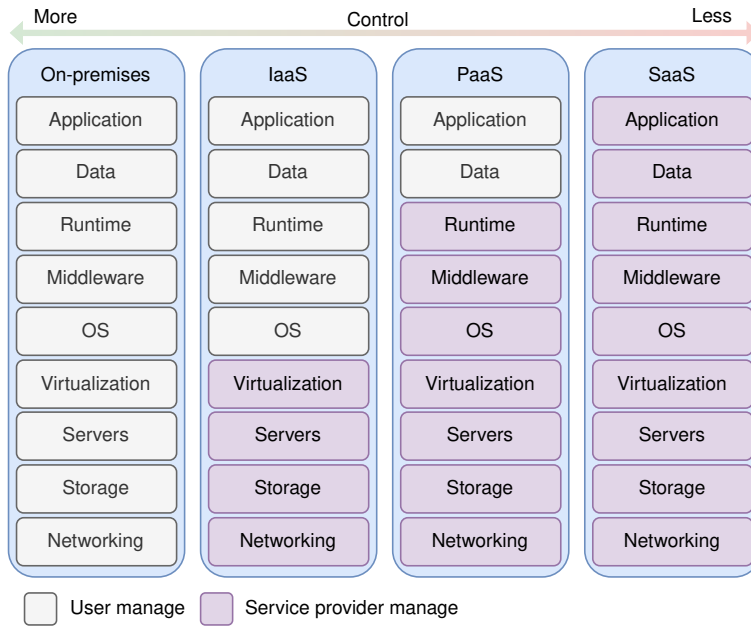


Figure 2.3: Difference between cloud options and on-premises solution.

The user can choose a single solution, or combine more of them if such a thing is required depending on preferences and needs.

By the ownership, CC can be categorized into three categories:

- **Public cloud** is a type where CC is delivered over the internet and shared across many many organizations and users. In this type of CC, architecture is built and maintained by others. Users and organizations pay for what they use. Examples include AWS EC2, Google App Engine, Microsoft Azure, etc.;

- **Private cloud** is a type where CC is dedicated only to a single organization. In this type of CC, architecture is built by an organization that may offer their solution or services to the users or other organizations. These services are in the domain of what the organization does, and that organization is in charge of maintenance. Examples include VMWare, XEN, KVM, etc.;
- **Hybrid cloud** is such an environment that uses both public and private clouds. Examples include IBM, HP, VMWare vCloud, etc.

Table 2.3 shows comparison of public, private and hybrid cloud capabilities.

Capabilities	Public cloud	Private cloud	Hybrid cloud
<b>Data control</b>	IT enterprise	Service Provider	Both
<b>Cost</b>	Low	High	Moderate
<b>Data security</b>	Low	High	Moderate
<b>Service levels</b>	IT specific	Provider specific	Aggregate
<b>Scalability</b>	Very high	Limited	Very high
<b>Reliability</b>	Moderate	Very high	Medium/High
<b>Performance</b>	Low/Medium	Good	Good

Table 2.3: Comparison of public, private and hybrid cloud capabilities.

In the rest of the thesis, if not stated differently when CC term is used it denotes public cloud.

CC has been the dominating tool in the past decade in various applications [13]. It is changing, evolving, and offering new types of services. Resources such as container as a service (CaaS), database as a service (DBaaS) [57] are newly introduced. The CC model gives us a few benefits. Centralization relies on the economy of scale to lower the cost of administration of big DCs. Organizations using cloud services avoid huge investments, like creating and maintaining their own DCs.

They consume resources usually created by others [13] and pay for usage time — *pay as you go model*.

Centralization gives us a few really hard problems to solve. As already stated in section 1.1 data is required to be moved to the cloud from data sources, which introduces a high latency in the system [7].

There are a few notable attempts to help data ingestion into the cloud. Remote Direct Memory Access (RDMA) protocol makes it possible to read data directly from the memory of one computer and write that data directly to the memory of another. This is done by using *specialized hardware* interface cards, switches, and software as well, and operations like read, write, send, receive, etc. do not go through the CPU. With these characteristics, RDMA has low latencies and overhead, and as such reaches better throughputs [58]. This new hardware may not be cheap, and not every CC provider uses them for every use-case. This may not be enough, especially with the ever-growing amount of IoT devices and services.

Over the years there are more service options available, forming **everything as a service (XaaS)** model [36]. This model proposes that any hardware or software resource can be offered as a service to the users over the internet.

Table 2.4 shows common examples of SaaS, PaaS, and IaaS applications.

Platform	Common Examples
IaaS	AWS, Microsoft Azure, Google Compute Engine
PaaS	AWS Elastic Beanstalk, Azure, App Engine
SaaS	Gmail, Dropbox, Salesforce, GoToMeeting

Table 2.4: Common examples of SaaS, PaaS, and IaaS.

### 2.1.2.1 Multi-cloud and sky computing

In recent years there has been one extension of CC from a series standpoint called **multi-cloud** [59, 60] or sky computing [61] (terms

are going to be used interchangeably).

It is such an environment where an enterprise uses more than one cloud platform, with at least two or more public cloud providers that each delivers a specific application or service.

A multi-cloud can be comprised of any model presented on the page 19. This model relies on the possibility that if one cloud provider fails for whatever reason, the next one will be able to serve user requests.

This strategy allows creation of a single heterogeneous architecture, allowing distribution of cloud assets and workloads across multiple providers (active-active), or deploy a single workload on one provider, with a backup on another (active-passive) [62].

CC gives a user an illusion that he is using a single machine, while the background implementation is fairly complicated and consists of various elements that are composed of countless machines. CC is a typical example of a horizontally scalable system presented in 2.1.1.

### 2.1.3 Membership protocol

At the beginning of this section DS were introduced, and two interesting assumptions by Tanenbaum et al. were presented [48, 49]. If one more look is taken at (2) assumption, we will see that users of the DS whether they are users or applications perceive DS as a single unit. Inside this single unit, nodes need to collaborate, so that they are able to do various kinds of tasks.

The most basic of all these tasks is that nodes need to know which group they belong to, and who are their peers in the group they will collaborate with. This might sound like a trivial idea, but when we include 8 fallacies of the DS 2.1 into the equation, things start to be not so trivial after all. In the setup where nodes are connected over the local network or internet, and they need to communicate, things will go wrong for various reasons.

To resolve the problem who their group peers are, a membership protocol comes to help. These protocols need to ensure that each process of one group updates its local list of **non-faulty** members of

the group, and when a new process joins or leaves the group, the local list for every process needs to be updated. This is the most basic idea behind membership protocols.

Processes in the group of nodes in a group will ping each other in different ways, and using different strategies to figure out which nodes are dead and which are alive. There are a few existing algorithms that do this job, and they are (usually) based on the way epidemics spread or how gossip is spread in a human population. Because of this feature, these algorithms are usually called *Gossip* style protocols.

Every membership protocol has some properties that will ensure efficiency and scalability:

- (1) **Completeness**, this property must ensure that every failure in the system is detected;
- (2) **Accuracy**, in an ideal world, there should be no mistakes when detecting failures, but In a real-life scenario, we need to reduce false positives as much as we can;
- (3) **Failure detection speed**, all failures needs to be treated as fast as possible, in order to remove the node from the group and reschedule the tasks from the dead node to alive ones;
- (4) **Scale**, with this property we must ensure that the network load that is generated should be distributed equally between all processes in the group.

The easiest idea to implement this protocol would be **heartbeating** technique where process  $P_i$  will send a heartbeat message to all his peers in the group or **multicast**. After some time if process  $P_j$  did not receive a heartbeat message from  $P_i$ , it will mark him as failed. This idea is easy to understand, and implement but the downsides are that its process is not that **scalable**, especially for large groups, and this will introduce huge network traffic.

To resolve this problem, Das et al. [32] introduced **Scalable Weakly-consistent Infection-style Process Group Membership** protocol (**SWIM** for short). This protocol divides the membership problem into two parts:

- (1) **Failure detection**, this component works so that one node will select a random node in the group, and it will send it *ping* message, expecting *ack* message in return — **direct ping**. If such message is not received, it will pick  $n$  nodes to probe through a *ping-req* message — **indirect ping**. If this fails, the node will be marked as *suspected*, and it will be marked as *dead* after some timeout. If the node gets alive, it will ping some other node and it will get back into the group. Figure 2.4 show message passing in **direct** (left), and **indirect** (right) ping in SWIM protocol;
- (2) **Information dissemination**, with previous strategy, information can be disseminated by **piggybacking** the data on multiple messages (*ping*, *ping-req* and *ack*), and avoid using the multi-cast solution.

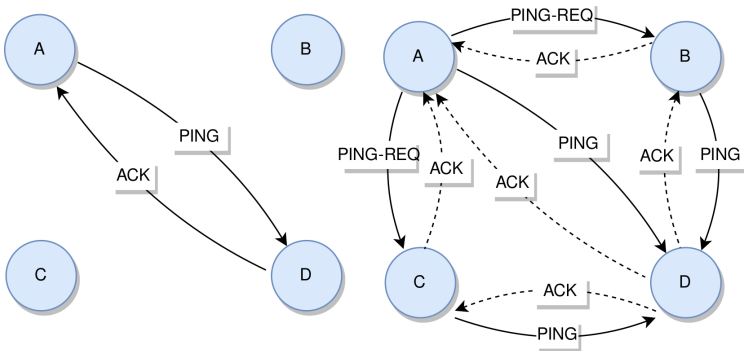


Figure 2.4: Direct and indirect ping in SWIM protocol.

Over the years, researchers found ways to improve the protocol, for example Dadgar et al. presented *Lifeguard protocol* [63] for more

accurate failure detection, and there are other implementations to fine-tune the SWIM, but the base idea is still there. Today SWIM or SWIM-like protocols are standard membership protocols whenever some node clustering is done.

### 2.1.4 Mobile computing

The first idea that introduced task offloading from the cloud [64, 65] was Mobile cloud computing (MCC). Mobile devices run small client software and interact with the cloud over the internet, while heavy computation remains in the cloud.

The cloud is usually far away from end devices because DCs are built on specific locations in the world to target as many users nearby as possible. This sparse deployment will most likely lead to high latency, and bad quality of experience (QoE) [65] for most users. Latency-sensitive applications especially will have a hard time. As a model, MCC is not much different from the standard CC model. The good thing is that the cloud has been relaxed a little bit, and a small number of tasks has been moved from the cloud. But this model opens the door for the next-generation models.

The development led to new computing areas like EC, osmotic computing [66, 67], sky computing [61], etc. EC is a next-generation model where computing and storage resources are in proximity to data sources [13]. This idea might overcome cloud latency issues and known MCC problems. The main strength of the EC lays in the CC enhancements with new processing ideas, for the next-generation use-cases [15].

EC has brought a few different models over the years. Models like fog [16], cloudlets [12], and mobile edge computing (MEC) [17] emerged. This thesis will refer to all these models as edge computing, and nodes used by them as edge nodes. Different EC models rely on the concept of data and computation offloading from the cloud closer to the ground [18]. Only heavy computation remains in the cloud because of more available resources [15], compared to edge nodes.



EC models introduced small-scale servers that operate between data sources and the cloud. These small-scale servers have much fewer capabilities compared to the cloud servers [19]. To avoid latency and huge bandwidth [12], EC nodes can be dispersed in various locations, for example, base stations [17], restaurants, or over arbitrary geographic regions.

## 2.2 Distributed programming

DP can be defined as the use of a DS to solve one large problem, by breaking it down into several smaller parts, where each part is computed in the individual node of the DS and coordination is done by passing messages to one another [49]. Computer programs that use this strategy and run on top of the DS are called **distributed programs** [68, 69].

Similar to CC in Section 2.1.2, to a normal user, DC systems appear as a single system similar to one the user uses every day on his/her personal computer. DC shares the same fallacies to DS presented in 2.1.

### 2.2.1 Big Data

Term big data means that the data is unable to be handled, processed, or loaded into a single machine [70]. That means that traditional data mining methods or data analytics tools developed for centralized processing may not be able to be applied directly to big data [71].

New tools and methods that have been developed rely on DS and one specific feature **data locality**. Data locality can be described as a process of moving the computation closer to the data, instead of moving large data to computation [72]. This simple idea minimizes network congestion and increases the overall throughput of the system.

Two examples of how huge generated data could be have already been given in 1.1, and when other IoT sensors and devices are included these numbers will just keep getting bigger and bigger [73].

Contrary to relational databases that mostly deal with structured data, Big Data is dealing with various kinds of data [70, 71, 72]:

- **Structured** data is a kind of data that have some fixed structure and format. A typical example of this is data stored inside a table of some database. Organizations usually have no huge problem extracting some kind of value out of the data;
- **Unstructured** data is a kind of data where there is not any kind of structure at all. These data sources are heterogeneous and may contain a combination of simple text files, images, videos, etc. This type of data is usually in raw format, and organizations have a hard time deriving the value out;
- **Semi-structured** data is the kind of data that can contain both previously mentioned types of data. An example of this type of data is XML files.

Big Data have recognizable feature called **V's** of Big Data [74], shown in Figure 2.5. Name is derived from initial letters of the big Data features.

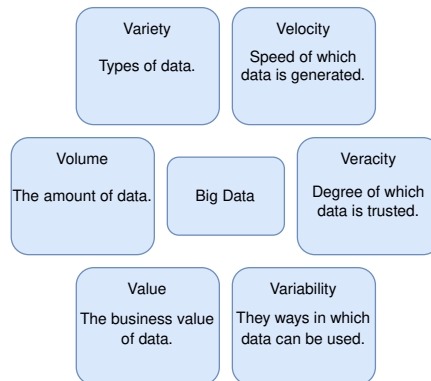


Figure 2.5: V's of Big Data.

Processing in big data systems can be represented as [75, 76]:

- **Batch processing** represents a data processing technique that is done on a huge quantity of the stored data. This type of processing is usually slow and requires time;
- **Stream processing** represents a data processing technique that is done as data get into the system. This type of processing is usually done on a smaller quantity of the data **at the time** [77], and it is faster;
- **Lambda architectures** represents a processing technique where *stream processing* and handling of massive data volumes in a *batch* are combined in a uniform manner, reducing costs in the process [76].

Big data systems, are not processing and value extracting systems. Big data systems can be separated into several categories: (1) data storage, (2) data ingestion, (3) data processing, and analytics. All these systems aid to properly analyze ever-growing requirements [78].

Despite a promise that big data offers to derive value out of the collected data, this task is not easy to do and requires a properly set up system filtering and removing data that contains no value. To aid this idea, data could be filtered and a little bit preprocessed on close to the source [22], and as such sent to data lakes [79].

### 2.2.2 Microservices

There is no single comprehensive definition of what a microservice is. Different people and organizations use different definitions to describe it. A working definition is offered in [80] as “a microservice is a cohesive, independent process interacting via messages”. Despite the lack of a comprehensive definition, all agree on a few features that come with microservices:

- (1) They are small computer programs that are independently deployable and developed;
- (2) They could be developed using different languages, principles, and using different databases;
- (3) They communicate over the network to achieve some goal;
- (4) They are organized around business capabilities [81];
- (5) They are implemented and maintained by a small team.

The industry is migrating much of their applications to the cloud because CC offers to scale their computing resources per their usage [82]. Microservices are small loosely coupled services that follow UNIX philosophy “do one thing and do it well” [83], and they communicate over well defined API [80].

This architecture pattern is well aligned to the CC paradigm [82], contrary to previous models like monolith whose modules cannot be executed independently [80, 84], and are not well aligned with the CC paradigm [84]. Table 2.5 summarizes differences between the monolith and microservices architecture.

Feature	Monolith	Microservices
<b>Structure</b>	Single unit	Independent services
<b>Management</b>	Usually easier	Add DS complexity
<b>Scale/Update</b>	Entire app	Per service
<b>Error</b>	Usually crush entire app	App continue to work

Table 2.5: Differences between the monolith and microservices architecture.

Since its inception, microservices architecture has gone through adaptations, and modern-day microservices are extended with two new models, each with its unique abilities and problems:

- **Cloud-native applications** are specially designed applications for CC. They are distributed, elastic, and horizontally scalable systems by their nature, and composed of (micro)services that isolate state in a minimum of stateful components [85]. These type of applications are self-contained, could be deployed independently, and they are composed of loosely coupled microservices that are packaged in lightweight containers. They have Improved resource utilization, and they are centered around APIs;
- **Serverless applications** is a computing model, where the developers need to worry only about the logic for processing client requests [86]. Logic is represented as an event handler that only runs when a client request is received, and billing is done only when these functions are being executed [86]. **Cold start** is one of the features of serverless computing, and we can define it as user requests need to wait until a new container instance is up and running before it can do any processing at all. Most providers have 1–3 second cold starts, and this is important for certain types of applications where latency is a concern. Cold start is only happening when there are no *warm* containers available for the request, meaning there is no single instance to server request. Other features include: **(1)** simplified services development, **(2)** faster time to market, **(3)** and lower costs;
- **Service Mesh** is designed to standardize the runtime operations of applications [87]. As a part of the microservices ecosystem, this dedicated communication layer can provide several benefits, such as: **(1)** observability, **(2)** providing secure connections, or **(3)** automating retries and backoff for failed requests. With these features, developers only focus on the implementation of business logic, while operators gain out-of-the-box traffic policies, observability, and insights from the services. Advocates

of the microservice movement, nowadays recommend using service mesh architecture when running microservices in production environments.

Previous models are not explicitly different, they all can be viewed as cloud-native applications. The enumeration is given for the sake of pointing out their different models and aspects of working.

Microservices communicate over a network to fulfill some goal using message passing techniques and technology-agnostic protocols such as HTTP. They can be implemented as:

- Representational state transfer (REST) services [88], is an architectural style with a set of constraints that users can create web services and interoperability between computer systems on the internet. It is based on HTTP routes to define resources and use HTTP verbs to represent operations over these resources. It relies on textual based communications, and payload could be represented using *JSON*, *XML*, *HTML*, etc.;
- Remote procedure calls (RPC) represent an architectural way to design services that can call subroutines that are located in different places, usually on another machine. The client calls these operations like they are located locally in his address space;
- Event-driven services are services where communication between services is done using events. Events are sent on some channel and other read messages that are received on another channel. These channels could be implemented either like message queues or message topics. Services connect to message queue or subscribe to the specific topic, and when messages arrive, they can act according to the message type.

They are well aligned with text-based protocols like HTTP/1 using *JSON* for example, or binary protocols such as HTTP/2 using *protobuf* and *gRPC* for example, and even new faster version like

HTTP/3 over new *QUIC* protocol, designed by Google. HTTP/3 is the latest version of the conventional and trusted HTTP protocol. It is very similar to HTTP/2, but it also offers a few important new features.

Table 2.6 shows important difference between versions of HTTP protocol.

Feature	HTTP/1	HTTP/2	HTTP/3
<b>Transport</b>	text	binary	binary
<b>Parallelism</b>	No	Yes	Yes
<b>Protocol</b>	TCP	TCP	QUIC
<b>Space</b>	OS level	OS level	User level
<b>Server push</b>	No	Yes	Yes
<b>Compression</b>	Data	Data/Headers	Data/Headers

Table 2.6: Idempotent and non-idempotent operations.

To ensure a wider range of devices that can communicate with the rest of the systems, developers usually have a gateway into the system that is REST service, and other services could be implemented differently.

It is important to point out, that all flavors of microservices applications rely on continuous delivery and deployment [89]. This is enabled by lightweight containers, instead of virtual machines [90], and orchestration tools such Kubernetes [91]. These concepts will be described in more detail in Section 2.7.

Microservices architecture is a good starting point especially for being built as a service applications model, and applications that should serve a huge amount of requests and users, especially with the benefits of CC to pay for usage, and the ability to scale parts of the system independently. They are not necessarily easy to implement properly, and there is more and more critique to the architecture model [92]. Microservices are rely upon and use parts of the DS, and as such, they inherit almost all problems DS has.

One particular thing that users need to be aware of is **idempotency**. In microservices applications, developers are dealing with inconsistencies in the distributed state, and their operations should be implemented as idempotent. An operation is idempotent if it will produce the same results when executed over and over again. It is a strategy that means that operations with side effects like creation or deletion can be called any number of times while guaranteeing that side effects only occur once. Idempotency is a term that comes from mathematics, and can be represented by simple idempotency law for operation  $*$  like [93]:

$$\forall x, x * x = x \quad (2.3)$$

Not all Create, Read, Update, Delete (CRUD) operations are idempotent by default. Developers need to make effort to make all of them idempotent, to prevent bad outcomes and inconsistent states.

Table 2.7 shows list of idempotent and non-idempotent for standard CRUD operations:

Operation	Idempotent	Non-idempotent
Create		x
Read	x	
Update	x	
Delete	x	

Table 2.7: Idempotent and non-idempotent operations.

*Create* operation is **not** idempotent by default, but to make it idempotent there are multiple strategies how to do so. The most common way is to create **idempotency key** that will be sent in the request, and based on that request server can decide if this operation is already invoked or not. If a server has already “seen” specified idempotency key than the operation is already done and we can return just the response that the operation is done but no operation will be done over the state of



the service or application. If the server sees the idempotency key for the first time, that is the signal that this request is a new one, and it should be done.

Idempotency key could be stored in any kind of storage, it is not uncommon that these keys are stored in cache storage with some time to live (TTL) policy that will automatically remove the key after a specified time.

Another option that is commonly used is hashing user specified actions. It is useful to know which part of the action set is already done and which is not. This strategy is used in scenarios where we must preserve the order of actions.

### 2.2.2.1 Distributed Queries

Applications built using microservices architecture propose different strategies for data storage. One common and recommended technique is *database per service* pattern [94].

The result of this technique is that overall state of the system will be distributed across multiple data stores, accessible only from their own microservices. This creates two important topics to think about: **(1)** transactions that span over multiple services can be implemented using *Sagas* for example (see page 44), and **(2)** the complex queries which require data from multiple databases.

The complex queries will involve data available in multiple databases, and a client can access all these microservices and aggregate data, however, this is not the recommended solution because the client does not have full understanding how the system manages the data.

In microservices architecture, there are two standard ways to solve this problem [94]:

1. **Command Query Responsibility Segregation (CQRS)** separates responsibility for *modifying* data (command), and *reading* the data (query), making the logic clearer and easier to optimize

different parts of the system [94, 95]. This increases the complexity of entire system, but it supports multiple denormalized views that are scalable and performant [94]. When the data in one service is modified, the service emits the event, which will change the service responsible for complex query. The service that will serve the queries will keep the *read-only* replica of the data. Figure 2.6 shows an example diagram for CQRS pattern;

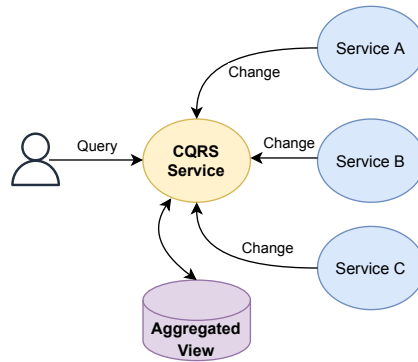


Figure 2.6: CQRS pattern diagram.

2. **API composition** is a simple way to query data in a microservice architecture [94, 96], alternative and more lightweight solution than CQRS. The main difference between CQRS is that this patterns does not have its own data storage. When a request comes in, it accesses every single microservice containing data, combines the results, and then returns the combined result to the client. Making it easier to implement, because database does not need to be refreshed every time when some change occurs in the system. On the other hand, it may yield a slower response, depending on how many services we need to contact for the information, their availability, and time to merge and prepare data in memory. Figure 2.7 shows an example diagram for API Composition pattern.

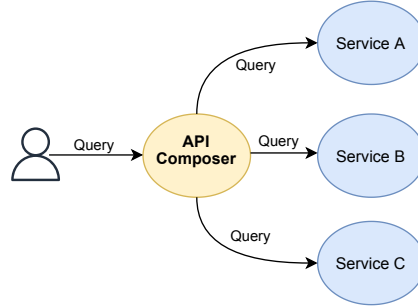


Figure 2.7: API Composition pattern diagram.

The best chance to succeed when implementing a microservices architecture is to simply follow existing patterns and use existing solutions with proven quality.

### 2.2.3 Observability

Observability is an integral part of any real-world computing system, and especially in the DS where observing what is happening in a network of processes is difficult [97], due to its nature.

In case of errors, fails or misbehavior of the system, some insight can be gained into what causes failure or which set of parameters in which circumstances. The three pillars of the observability are: **(1)** logs aggregation, **(2)** distributed tracing, and **(3)** alerting.

The logging operation, gives developers ability to store various arbitrary informations, that will provide more details for those who are investigating the failure. One thing we must be aware of is not to store any sensitive pieces of information in the log because this can cause a bunch of problems. Another thing we must be aware of is that we do not log too much and too often to slow down the business logic and execution of the function.

In monolithic applications logging and monitoring is a little bit easier to implement, because we have the whole application state in one place. When we come to the field of DS and microservices, our state is

scattered across multiple elements or services. In DS, the monitoring involves interactions among concurrently executing processes [98].

The solution to this problem is to use a centralized logging service – **logs aggregation**, that collects logs from each service [99]. This is beneficial because users can search and analyze the logs as a whole state of the system. To do this properly the log must be stored very reliably [100]. Users can then configure the log server for some alerts that are triggered when certain messages appear in the logs. The log of DS usually does not contain enough information to regenerate the timeline of execution, and this is one reason that logs in DS are so hard to interpret [99].

To resolve this problem of DS execution timeline, Google developed a new technique called **tracing** [101]. The trace represents a single execution timeline or execution of one request. Trace will create a tree, and the tree is used to establish order. Every node in the tree represents a unit of work and it is called span. A tree unites all the elements needed to carry out an originating request. In every span or unit of work, we can attach more details about that particular execution element. Figure 2.8 shows the simple example of *RequestX* path through the processes in a distributed system, where each service **call** could be RPC, HTTP or some other request.

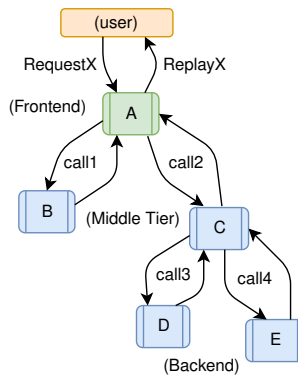


Figure 2.8: *RequestX* path through the processes in a distributed system.

And last, but not the least option is **alerting**. Alerting of a monitored system can be represented as a set of rules that performs actions based on changes in specific metric. Alerting enables a system to notify users when something important happens or (probably) is going to happen.

DS logging, tracing and alerting represents the important role of any system, and as such, it should not be neglected especially in the DS environment. Every user request should be traced and logged from an infrastructure perspective, but we should allow users to store logs from their applications.

## 2.3 Distribution Models

The role of distribution models is to determine the responsibility for the request, or to answer the fundamental question “who is in charge” for a specific request. There are two ways to answer this question: **(1)** all nodes in the system, or **(2)** single node in the system.

### 2.3.1 Peer-to-peer

Peer-to-peer (P2P) communication is a networking architecture model that partitions tasks or workloads between peers [102]. All peers are created equally in the system, and there is no such thing as a node that is more important than others.

Every Peer has a portion of system resources, such as processing power, disk storage, or network bandwidth, directly available to other network participants, without the need for central coordination by servers or stable hosts [102]. P2P nodes are connected and share resources without going through a separate server computer that is responsible for routing.

Figure 2.9 shows difference in network topology between P2P networks (*left*) and client-server architecture (*right*).

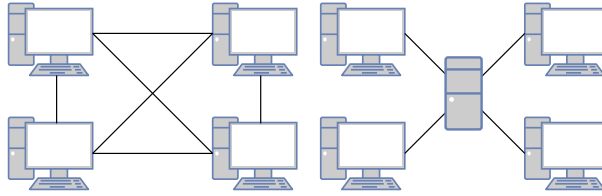


Figure 2.9: P2P network and client-server network.

Peers are creating a sense of virtual community. This community of peers can resolve greater tasks, beyond those that individual peers can do. Yet, these tasks are beneficial to all the peers in the system [103]. When a request comes to such a network, a node that accepts the request is usually called **coordinator**, because it then is trying to find the right peer to send a request to.

Based on how the nodes are linked to each other within the overlay network, and how resources are indexed and located, we can classify networks as [104]:

- **Unstructured** do not have a particular structure by design, but they are formed by nodes that randomly form connections [105]. Their strength and weakness at the same time is the lack of structure. These networks are robust when peers join and leave the network. But when doing a query, they must find more possible peers that have the same piece of data. A typical example of this group is a Gossip-based protocol like [32];
- **Structured** peers are organized into a specific topology, and the protocol ensures that any node can efficiently search the network for a resource. The famous type of structured P2P network is a Distributed Hash Table (DHT). These networks maintain lists of neighbors to do a more efficient lookup, and as such, they are not so robust when nodes join or leave the network. DHT is commonly used in resource lookup systems [106], and as efficient resource lookup management and scheduling of applications, or as

an integral part of distributed storage systems and NoSQL[107] databases;

- **Hybrid** combine the previous two models in various ways.

P2P networks are a great tool in many arsenals, but because of their unique ability to act as a server and as a client at the same time, we must be careful and pay more attention to security because they are more vulnerable to exploits [108].

### 2.3.2 Master-slave

In the master-slave architecture, there is one node that is in charge – **master**. This node accepts requests, and we usually do not communicate to the rest of the nodes or **slaves**. The master node is usually better and more expensive or even specialized hardware such as redundant array of inexpensive disks (RAID) to lower the crash probability. The cluster can also be configured with a **standby** master, and this node is continually updated from the master node.

But no matter how specialized hardware master runs on, it is prone to fail for various reasons, so it is a **single point of failure (SPOF)**. If crash happens, then standby master could continue to the server as a master, or new **leader election** protocol [109] is initiated to pick a new master node.

The master node is responsible for processing any updates to that data. If the master fails, then the slaves can still handle **read** requests. Failure of the standby master node to take over from the master node is a real problem if we want to achieve a high-availability system.

Figure 2.10 shows difference between master-slave (*left*) and peer-to-peer (*right*) request handling.

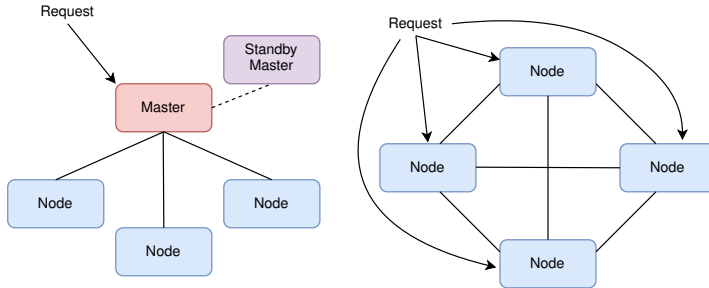


Figure 2.10: Handling requests master-slave and peer-to-peer

Using the right distribution model usually depends on the business requirements. High availability requires a P2P network because of no SPOF. If we could manage data using batch jobs that run in off-hours, then the simpler master-slave model might be the solution.

### 2.3.3 Replication

Replication is a technique used to ensure the availability of redundant resources closer to the place where it is accessed, with a goal to improve reliability, fault-tolerance, or accessibility [48, 49].

On the other hand, replication leads to better performance by balancing the load between components [48]. This is of importance, especially in geo-distributed systems, because it can lower down the communication latency. Users can access a copy of data nearby [49].

Replication requires information about **replication factor**. The replication factor tells us the number of nodes to replicate the data. Replication could be implemented inside a single cluster or multiple regions, increasing the overall availability of the system further.

Storing data in several places requires that system deals with the **consistency** of stored data. The users can access different replicas at different times, leading the data to an inconsistent state. Consistency is tunable parameter, and user can choose between **strong** and **eventual** consistency.



## 2.4 Similar computing models

In this section, we are going to shortly describe models that are similar to the DS, and as such, they may be the source of confusion.

### 2.4.1 Parallel computing

DC and parallel computing seem like models that are the same, and that may share some features like simultaneously executing a set of computations in parallel. Broadly speaking, this is not far from the truth [68].

Differences between the two can be presented as follows: in parallel computing, all processor units have access to the shared memory and have some way of the faster inter-process communication, while in DS and DC all processors have their memory on their machine and communicate over the network to other nodes which are significantly slower.

These models are similar, but they are not identical, and the kinds of problems they are designed to work on are different. Figure 2.11 visually summarizes the architectural differences between DC (*up*) and parallel computing (*down*).

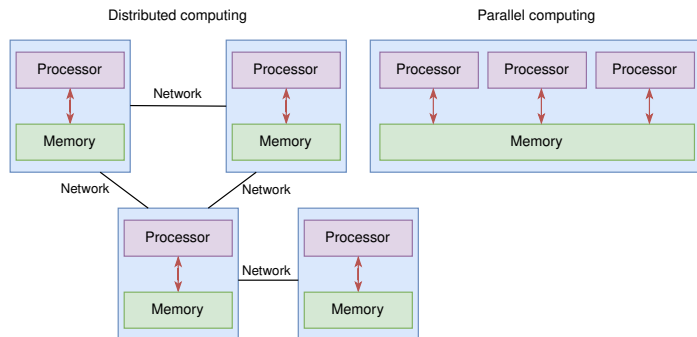


Figure 2.11: Architectural difference between DC and parallel computing.

Parallel computing has often used the strategy with problems that due to their nature or constraints must be done on multi-core machines simultaneously [110]. It is often that some big problems are divided into smaller ones, which can then be solved at the same time.

Several tasks require parallel computing like simulations, computer graphics rendering, or different scenarios in scientific computing.

## 2.4.2 Decentralized systems

Decentralized systems are similar to DS. But if we take a closer look, these systems **should not** be owned by a single entity. CC, for example, is a perfect example of DS, but it is not decentralized by its nature. It is a centralized system by the owner like AWS, Google, Microsoft, or some other private company because all computation needs to be moved to big DCs [7].

By modern standards, when we are talk about decentralized systems, we usually think of blockchain or blockchain-like technology [111], since here we have distributed nodes, that are scattered and there is no single entity that owns all these nodes. But even if this technology is run in the cloud, it is loses the decentralized feature. This is the caveat we need to be aware of. These systems are facing different issues because any participant in the system might be malicious and they need to handle this case.

Nonetheless, CC can and should be decentralized in the sense that some computation can happen outside of cloud big DCs, closer to the sources of data. These computations could be owned by someone else, and big cloud companies could give their solution to this as well to relax centralization and problems that CC will have especially with ever-growing IoT and mobile devices.

## 2.5 Transactions

Transactions are keeping data consistent even in the presence of highly concurrent data accesses and despite all sorts of failures [112]. Transactions are trying to resolve this problem in a generic way, in such a way that is invisible to the application logic.

The main goal of transactions is to maintain system integrity in the consistent state, by ensuring that all operations on the system are either all completed successfully or all canceled successfully. Transactions are typically used in systems that needs to preserve some state (e.g. database or some filesystems).

In their book Gray et al. make a good parallel with the contract law, saying that transactions give us the ability to *clean up the situation*, if something does not work right. [113].

Transactions guarantee following four properties: **(1)** atomicity, **(2)** consistency, **(3)** isolation, and **(4)** durability, also known as *ACID* properties [114].

### 2.5.1 Distributed transactions

Because of the nature of DS, more network hosts are involved which significantly complicates the transaction mechanism. Distributed transactions are required to have all four *ACID* properties. This might not be so easy to achieve amongst other due to *CAP* theorem (see page 15).

In their book, Morgan et al. claim that here exists no distributed commit protocol that can guarantee independent process recovery in the presence of multiple failures (e.g., network partitionings) [112].

Distributed transactions include few protocols such as two-phase commit (2PC), three-phase commit (3PC), Paxos, and various other approaches to quorum giving programmers facade of global serializability [115].

Avoiding distributed transactions allows a much simpler, more robust and efficient solution.

## 2.5.2 Sagas

In 1987, Molina et al. presented *Sagas* [116] and their work in the area of *long-lived transactions (LLTs)*, types of transactions that hold resources for long periods, and as such delay shorter and more common transactions.

These transactions are increasingly relevant and important in the current technological landscape with the distribution of components and microservices architectures.

The saga transaction is composed of sub-transactions, executed in an atomic way so that either all or none of the sub-transactions take effect. However, no isolation is necessarily guaranteed between the sub-transactions of different sagas.

One transaction  $T$  is composed of multiple sub-transactions  $T_1, \dots, T_n$ , and every sub-transaction  $T_i$  has an associated *compensating transaction*  $C_i$  or how the effects of the transaction can be rolled back. The saga transaction executes sub-transactions sequentially.

Figure 2.12 shows an example of one transaction separated into multiple sub-transactions where *green* arrows represent success path, while *red* arrows represent rollbacks, structure of the saga transaction and at least the previous and next element in the chain.

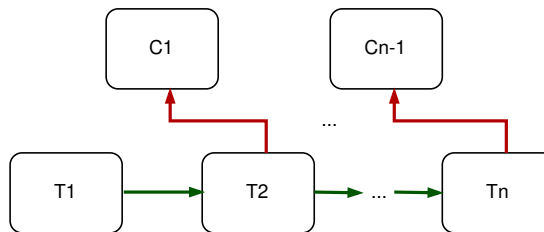


Figure 2.12: Saga transactions separated in sub-transactions.

Sagas can be implemented using two patterns: **(1)** orchestration where one centralized element is responsible for all the coordinating, or **(2)** choreography, where each service trigger local transactions in other services.

This isolation can be added in the application layer [114] i.e. *semantic lock*. This strategy will change state so that some data are in process and should be treated differently.

For example, an order can be in one of the following states: **(1)** *PENDING*, **(2)** *COMPLETED*, or **(3)** *CANCELED* state. Other transactions will not make use the data if status is not set to *COMPLETED*.

## 2.6 Garbage collection

Most modern programming languages nowadays allow programmers to allocate and free memory. This process can be in total responsibility of the programmer, or language can handle this process for the programmer. When there is some automatic process involved to release unused resources, that process is called **garbage collection** or GC.

In [117], Jones, et al. describe *garbage collection* as the automatic management of dynamically allocated storage. However, the term *garbage collection* is not exclusive to programming languages. It is widely used to refer to all forms of automatic management of dynamically allocated resources.

Even with the rapid growth of memory sizes, and lowering the overall cost of memory is not inexhaustible. Like any other limited resources, it requires careful consideration and recycling. Even tools like Kubernetes, have some form of garbage collection that is used to remove unused items and their references in the system.

Over the years, different algorithms are used to do reference counting and tracing methods, to discover unused resources, and to free them. One of the most common algorithms for garbage collection is *mark-sweep* algorithm [118] developed in 1960. The variety on the topic

exists today, but the essence of the algorithm remains widely used today.

Traditional automated garbage collection is usually seen as slow, and disruptive to executing programs. Modern implementations of the garbage collection substantially reduced the overhead of the system [117].

## 2.7 Virtualization techniques

Virtualization as a technique started long ago in time-sharing systems, to provide isolation between multiple users sharing a single system like a mainframe computer [119].

In [120] Sharma et al. virtualization is described as technologies that provide a layer of abstraction of the physical computing resources between computer hardware systems and the software systems running on them.

Modern virtualization differentiates several different tools. Some of them are used as an integral part of the infrastructure for some flavors like IaaS, while others are used in different CC flavors as well as microservices packaging and distribution format, or are new and still are looking for their place. These options are:

- **Virtual machines (VM)** are the oldest technology of the three. They are described as a self-contained operating environment consisting of guest operating system and associated applications, but independent of the host operating system [120]. VMs enable us to pack isolation and better utilization of hardware in big DCs. They are widely used in IaaS environment [121, 122] as a base where users can install their own operating system (OS) and require software tools and applications;
- **Containers** provide the almost same functionality to VMs, but there are several subtle differences that make them a go-to tool

in modern development. Instead of the guest OS running on top of host OS, containers use tools that are in a Linux kernel like *cgroups* that limit process resource usage so that single process can not starve other processes and use all the resources for itself, and *namespaces* to provide isolation and partitions kernel resources so that single process sees node resources like it only exists there. Containers reduce time and footprint from development to testing to production, and they utilize even more hardware resources compared to VMs and show better performance compared to the VMs [123, 90]. Containers provide an easier way to pack services and deploy and they are especially used in microservices architecture and service orchestration tools like Kubernetes [91]. Google has stated several times in their online talks that they have used container technology for all their services, they even run VMs inside containers for their cloud platform. Even though they exist for a while, containers get popularized when companies like Docker and CoreOS developed user-friendly APIs;

- **Unikernels** is the newest addition to the virtualization space. Unikernels are defined as small, fast, secure virtual machines that lack operating systems [124]. Unikernels are comprised of source code, along with only the required system calls and drivers. Because of their specific design, they have a single process and they contain and execute what it absolutely needs to nothing more and nothing less [125]. They are advertised as new technology that will save resources and that they are *green* [126], meaning they save both power and money. When put to the test and compared to containers they give interesting results [125, 127]. Unikernels are still a new technology and they are not widely adopted yet. But they give promising features for the future, especially **if** properly ported to ARM architectures, and various development languages. Unikernels will probably be used as a user application and function virtualization tool, because of

their specific architecture, especially for serverless applications presented in 2.2.2.

With every virtualization technique, the ultimate goal is to pack as many applications on existing hardware as possible, so that there are no resources that are left not used – we are trying to achieve high resource utilization.

Figure 2.13 represents architectural differences between VMs, containers, and unikernels.

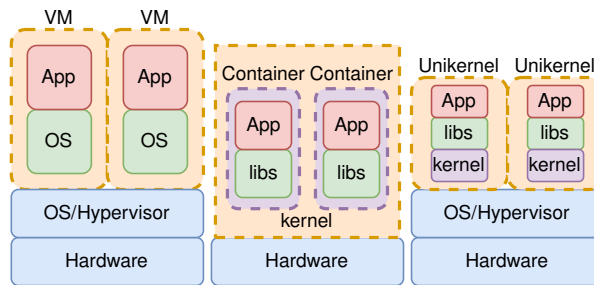


Figure 2.13: Architectural differences between VMs, containers and unikernels.

## 2.8 Deployment

Over the years different approaches evolved how to deploy infrastructure and applications. The difference just gets more amplified, when CC and microservices get into the picture, where frequent deployment is very common.

Deployments in such complex environment can be separated by how they handle changes on existing infrastructure or applications on:

- **A mutable model** is a model where we have in place changes which mean that the parts of the existing infrastructure or applications get updated or changed to do an update. In place change can produce some problems, and has:



- (1) More risk because in-place change may not finish, which puts our infrastructure or the application in a possible bad state. This is especially a problem if we have a lot of services and multiple copies of the same service. The possibility that our system is not on is a lot higher;
  - (2) High complexity, this is a direct implication of the previous feature. Since our change might not get fully done, it cannot be guaranteed that our infrastructure or application is transitioned from one version to another – change is not **discrete**, but **continues** since we might end up in some state in between where we are now and where we want to be.
- **An immutable model** is a model where no in-place changes on existing infrastructure or application are done whatsoever. In this model, the previous version is replaced completely with a new version that is updated or changed compared to the previous version. The previous version gets discarded in favor of the new version. Compared to the previous model, immutable deployment:
    - (1) Has less risk, since the existing infrastructure or the application, is not changed, but a new one is started and the previous one is shut down. This is important especially in DS where everything can fail at any time;
    - (2) Has less complexity of the mutable deployment model. This is a direct implication of the previous feature since the previous version is shut down and fully replaced with the new one. This is a **discrete** version change and atomic deployment with deferring deployments with fast rollback and recovery processes;
    - (3) Requires more resources [128], since both versions must be present on the node for this process to be done. The second

problem is the data that the application has generated should not be lost. The problem is solved by externalizing the data. We should not rely on local storage but store that data elsewhere, especially when the parts of the system are volatile and changed often. The key advantage of this approach is avoiding downtime experienced by the end-user when new features are released.

Immutability is a simple concept to understand and simplifies a lot especially in DS [128]. Write down some data, and ensure that it never changes. It can never be modified, updated, or deleted [129]. When this is combined with the promise that downtime can be avoided especially in complex DS, it is clear why the immutable model is gaining more and more popularity (especially with the arrival of containers).

Immutable infrastructure deployment offers several benefits on how to deploy changes. Even in production, it is easier to switch to a whole new version. These strategies include:

- **Blue-Green deployment**, this strategy requires two separate environments: (1) *Blue* current running version, and (2) *Green* is the new version that needs to be deployed. When there is satisfaction that the green version is working properly, the traffic can be gradually rerouted from the old environment to the new one, for example by modifying the Domain Name System (DNS). This strategy offers near-zero downtime;
- **A canary update** is a strategy where a small subset of requests is directed to the new version — the canary, and the rest of them are directed to an old version. If the change is satisfactory, the number of requests can be increased, and it should be monitored how the service is working with increasing load, if there are errors, etc.;
- **Rolling update** strategy updates large environments, a few nodes at the time. The setup is similar to blue-green deployment,

but here there is a single environment. With this strategy, the new version gradually replaces the old one. If for whatever reason the new version is not working properly on the larger number of nodes, rolling back to the previous version can always be done.

With mutable infrastructure, these strategies would be hard to implement, and maybe it is not possible at all. Besides infrastructure deployment, there is another side that we must be considered, and that is how to describe these deployments. Two different strategies can be considered here:

- **Imperatively**, with this option users have to write code or specific instructions step by step what the specific tool needs to do so that the application or infrastructure is properly set up. In this approach, we have a *smart* user who describes *dumb* machine what is needed to be done and in what order to achieve the desired state;
- **Declaratively**, with this option a user has to describe the end state or what is his desired state, and the tool needs to figure out the way how to do this. Here we have a *smart system* that will find a way how to achieve the **desired state**, and we have a user who *does not care* in what order actions need to be done — that is what the system needs to do. Users do not need to worry about timing, this simplifies the whole process and the code always represents the latest state. With this type of deployment, we can offer users, two different models:
  - (1) Using existing platform independent formats that users are familiar with, like *JSON*, *YAML*, *XML*, etc.;
  - (2) Using a domain-specific language (DSL) that users need to learn, but we might be able to optimize description.

So far, the first option is preferred by many companies, because users are already familiar with these formats, and does not require developers

time to develop new DSL. If the platform becomes too complicated then it makes sense to develop DSL for this purpose.

Figure 2.14 summarizes the difference between mutable and immutable deployment models.

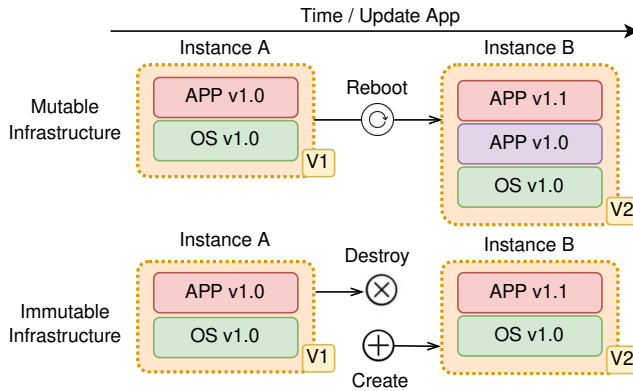


Figure 2.14: Difference between mutable and immutable deployment models.

With the introduction of *LinuxKit*, Linux subsystems are based around very secure containers. With *LinuxKit*, every part of the Linux subsystem runs inside the container, so a Linux subsystem can be assembled with services that are needed. As a result, systems created with *LinuxKit* have a smaller attack surface [130] than general-purpose systems.

This is important not only from the security point of view but also for infrastructure deployment as well. Specific OSs can be created, based on the containers for a different purpose. They can be updated, changed, and adapted for every machine or purpose needed.

Deployment is based on changing parts of the OS, and services that run inside the containers. As a result, everything can be removed or replaced. It's highly portable and can work on desktops, servers, IoT, mainframes, bare metal,  $\mu$ DCs at the edge, and virtualized systems.

## 2.9 Infrastructure as software

The infrastructure needs to be constantly deployed and maintained, so it would be beneficial to view the infrastructure as software (IaS) [131]. The benefit of this approach lies in the already available tools, principles, and techniques (e.g. reuse, testing, modeling, and evaluation) that can equally be used for the infrastructure definitions [132, 131].

In his work [132], Osterweil et al. argue that software and software processes share some similar characteristics: **(1)** both are executed, **(2)** both have requirements that need to be understood, **(3)** both benefit from being modeled, and **(4)** both can be guided by measurement.

Fitzgerald et al. [131] argue that process programming, defined by Osterweil [132] et al., is applicable in the infrastructure domain. The authors claim that it is useful to extend the “software process as software” paradigm to include infrastructure as software (IaS). The move towards multi-tier systems that involve cloud and cyber-physical systems will only stimulate the connection between software and reliable infrastructure systems.

Looking at the infrastructure configuration as the software allows application developers to venture into the “infrastructure programming” [131], allowing platforms and infrastructure to be managed in a similar way as the software is.

In a cloud environment, it is an essential technique to properly implement continuous deployment, giving us tools to automate the configuration and provisioning process of infrastructure using cloud instances [133]. Wittig et al. describe it as a process of managing and provisioning computer data centers through machine-readable definition files, rather than physical hardware configuration or interactive configuration tools [134].

### 2.9.1 Infrastructure as code

In modern cloud and microservices era, a new strategy to manage and deploy complicated infrastructure elements – Infrastructure as code (IaC). In his book [134] Wittig et al. describe it as a process of managing and provisioning computer data centers through machine-readable definition files, rather than physical hardware configuration or interactive configuration tools.

IaC has grown in popularity in recent years because it applies the same kind of version control and repeatability to the orchestration of the infrastructure as developers use for applications source code [135]. The second benefit is that configuration is decoupled from the system, it can more readily be deployed on a similar system elsewhere.

IaC is a continuation of earlier works of Osterweil et al. [132], and Fitzgerald et al. [131], trying to move infrastructure to level of software, keeping existing tools, best practices and techniques. It relies on the **reconciler pattern** [136], widely used in scheduling systems like Kubernetes.

This pattern enables tracking of resources using two simple states:

- (1) **Expected state** represents the desired state provided by the user. This is the state that must be (if possible) achieved at any give point in time;
- (2) **Current state** refers to the actual system state. This state must be refreshed periodically, so that system know should it be corrected or not.

The reconciler pattern runs a **reconciliation loop** that ensures that the *current state* remains the same as the *expected state*.

Figure 2.15 presents reconciliation loop diagram.

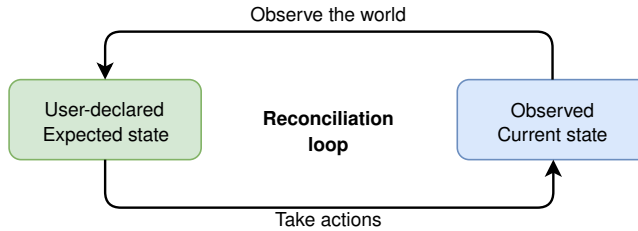


Figure 2.15: Reconciliation loop diagram.

Every node must provide its current state regularly, and when some state is divergent from the desired state, the system must act to ensure that this situation is corrected automatically. The node can send its state over existing channels e.g., **health-check** pings to minimize load and data transferred over the network, or a dedicated channel just for this purpose may exist.

Extension of such system can be done just by writing new *controller* that will react on changes for newly added kind of workload.

## 2.10 Development roles

In our modern internet-connected world of applications, we have several distinctive development roles. Each of them plays an important role so that modern software runs smoothly, and with less downtime. These development roles could be separated into a few categories (focus is on the technical roles):

- **A developer** is usually a person responsible for developing software for the user. A developer is responsible for maintaining existing, and/or developing new features. There could be different sub-roles, dealing with specific parts of the complex software systems;
- **DevOps** role represents a multidisciplinary organizational effort to automate application deployments through continuous delivery

of new software updates [137]. It combines software development with technology operations [138] to shorten the development life cycle;

- **Site Reliability Engineer (SRE)** role is responsible for availability, latency, performance, efficiency, change management, monitoring, emergency response, and capacity planning [139]. It is a software engineering role and needs to have an understanding of the fundamentals of computing [140], applied to the infrastructure and operations problems.

DevOps and SREs seem to be similar roles, they are both trying to bridge the gap between development and operations. As such they have a very large conceptual overlap in how they operate [141], but also have some differences.

Table 2.8 sums the differences between the DevOps and SREs.

Feature	DevOps	SREs
<b>Task</b>	Scaling, uptime, robustness	Development pipeline
<b>Essence</b>	Practices and metrics	Mindset and culture
<b>Team structure</b>	Wide range of roles: QA, developers, SREs etc.	SREs with operations and development skills
<b>Focus</b>	Development and delivery continuity	System availability and reliability
<b>Goal</b>	Bridge the gap between development and operation	

Table 2.8: Differences between DevOps and SREs.

In modern complex software development, SREs should keep things running smoothly, while DevOps principles should be used to improve processes.



So it is not either/or, but it seems that a combination of approaches may provide the best results. However, in some smaller organizations, these roles can overlap.

### 2.11 Concurrency vs parallelism

People usually confuse these two concepts. Even though look similar, they are a different way of doing things. In his speech, Rob Pike [142] gives a great explanation and examples on this topic. He also gives great definitions of these concepts like:

- **Concurrency** is composition of independently executing things. Concurrency is about dealing with a lot of things at once;
- **Parallelism** is the simultaneous execution of multiple things. Parallelism is about doing a lot of things at once.

These things are important, especially when building applications and systems that should achieve very high throughput. They must be built with a good structure and a good concurrence model. These features enable possible parallelism, but with communication [142]. These ideas are based on Tony Hoare work of Communicating Sequential Processes (CSP) [143].

#### 2.11.1 Actor model

An actor model, the main idea is based around **actors** which are small concurrent code, that communicate independently by sending messages, removing the need for lock-based synchronization [144]. This model proposes a similar idea to Tony Hoare's in his work with CSP [143], and actors are often confused with CSP.

Table 2.9 gives differences between an actor model and CSP. Actors do not share a memory, and they are isolated by a nature. An actor can create another actor/s and even watch on them in case they

Feature	CSP	Actor model
<b>Fault tolerance</b>	Distributed queue	Supervisors hierarchy
<b>Process identity</b>	Anonymus	Concrete
<b>Composition</b>	Not applicable	Applicable
<b>Communication</b>	Queue	Direct
<b>Message passing</b>	Sync	Async

Table 2.9: Differences between actor model and CSP.

stop unexpectedly. When an actor finished its job, and it is not needed anymore, it disappears. These actors can create complicated networks that are easy to understand, model, and reason about, and everything is based on a simple message passing mechanism.

Every actor has a designated message box. When a message arrives, the actor will test the message type and do the job according to the message type he received. In this way, we are not dependent on lock-based synchronization that can be hard to understand, and it can cause serious problems.

The actor model is fault-tolerant by design. It supports crash to happen because there is a “self-heal” mechanism that will monitor actor(s), and when the crash happens it will try to apply some strategy, in most cases just restart actor, but other strategies could be applied. This philosophy is useful because it is hard to think about every single failure option.

# Chapter 3

## Research review

In this chapter, we present the the results of the research reviews addressing issues and limits of CC discussed earlier are presented – both academia, and the industry researching and developing viable solutions to help cloud in the future.

Few directions are feasible: **(1)** focusing on adapting existing solutions to fit EC model, **(2)** experiment and develop new ideas and solutions to maybe fit more the nature of EC, and **(3)** try to combine both ideas to cover more research area.

Existing nodes organizational abilities are reviewed in Section 3.1. Section 3.1 reviews platform models from both academia and industry. Section 3.3 reviews cloud offloading techniques, whereas section 3.4 reviews some application models. Section 3.5 reviews some models used in cloud infrastructure management. Section 3.6 concludes this chapter, and gives the position of this thesis, compared to research previously reviewed.

### 3.1 Nodes organization

Guo et al. [26], gives a promising model based on a zone-based organization of edge nodes in the smart vehicles application. The authors

show how zone-based nodes organization enables continuity of dynamic services and reduces the connection handovers. They prove it is possible to enlarge the coverage of edge servers to a bigger zone, but at the same time, the computing power and storage capacity of edge servers could be expanded. EC could be based on the geo-distributed workloads, and this zone-based organization could benefit the EC model in various ways.

In their research [145], Baktir et al. explored the capabilities of software-defined networks (SDN) from a programming standpoint. Their findings show that SDN can be used to simplify the management of the network in a cloud-like environment. Networking in such a complex environment is not an easy task to achieve. The authors show how SDN hides the complexity of the heterogeneous environment from the end-users. As such, SDNs represent a good candidate for networking tasks in complex, cloud-like environments. In [146] Sherwin et al. investigate the programmatic control over the resources and functions of a network to make the network more dynamically configurable.

Sayed et al. in their work [44] show that EC systems will perform actions before connecting to the cloud and how they are easier to integrate with other wireless networks like mobile ad-hoc networks (MANETs), vehicular ad-hoc networks (VANETs), intelligent transport systems (ITSs) and the IoT to mitigate network-related and computational problems.

Content delivery networks (CDN) in centralized delivery models like CC have bad scalability, as Kurniawan et al. [27] argue in their research. To overcome these centralized problems and bad scalability, the authors proposed a different solution, a decentralized solution. To achieve such tasks, authors were using a network of gateways equipped with some storage as well, for internet services at home [27] forming even smaller DCs – nano DCs (*n*DCs). Authors present a possible usage for these *n*DCs in some large scale applications with much less energy consumption than traditional DCs.

In their paper [147], Ciobanu et al. introduce an interesting idea

called *drop computing*. The authors show the possibility for ad-hoc EC platform composition using a decentralized model over multilayered social crowd networks. This idea gives us the ability for collaborative computing that can be formed dynamically. The authors present an idea that we can form a computing group ad-hoc by employing the nearby devices in the mobile crowd, that is fully capable of quick and efficient access to resources, instead of sending requests to the cloud. Crowd nodes could also be an interesting idea for backup nodes, in case more computing power or storage is needed and there no more available resources to use. Forming platforms from crowd resources and ad-hoc, raises a few concerns: (1) crowd nodes availability, and (2) offered resources.

Greenberg et al. [24] introduce the idea of  $\mu$ DCs as DCs that operate in proximity to a big population compared to  $n$ DCs that serve a lot smaller population, for example, a single household.  $\mu$ DCs are an interesting model and area of rapid innovation and development, and because they are close to some population, they are minimizing the costs and the latency for end-users [24]. Their minimum size is defined by the needs of the local population [24, 25], as such, they are reducing traditional DCs fixed costs. The main feature that  $\mu$ DCs are relying on is agility. The authors describe agility as  $\mu$ DCs ability to dynamically grow and shrink resources to satisfy the resource demands and usage from the most optimal location [24]. In [148] Shao et al. present a possible  $\mu$ DCs structure serving only the local population, in the smart city use-case.

## 3.2 Platform models

Kubernetes (k8s) [91] is a system originally developed by Google influenced by their orchestrator platform called Borg [149]. Various other companies joined in developing this system, and now it is de facto standard in the cloud environment for microservices and cloud-native applications. By design, k8s is not intended to operate in

a geo-distributed environment, because it operates on a single cluster [91, 149, 150]. As such it might not be best suited for EC and geo-distributed  $\mu$ Cs. Nonetheless, it is a super valuable tool in the CC because it enables health checking, restarting, and orchestration at scale. Another potential problem with k8s is its relatively complicated deployment concept that might be too complicated for EC workloads. It is developed to connect microservices across the CC environment. It could be used as it is, a cloud-native orchestrator to run the master process for EC and  $\mu$ Cs and also cloud-native applications that will accept streams coming from  $\mu$ C applications. K8s relies upon a lot to work properly. Existing technologies and ideas worth exploring, such as loosely coupling elements with labels and selectors, should not be ignored.

In their research, Rossi et al. [150] present a solution based on k8s. Authors adapted k8s for workloads that are geo-distributed and they had used reinforcement learning (RL) techniques, to learn a suitable scaling policy from past experience. There are potential downsides to this approach. The first might be that machine learning implementation could be potentially slow due to the required model training, but also we need to somehow describe what a good decision, and what a bad decision is, in order to enable some algorithm to learn this. This might be a problematic thing, especially in urgent situations. The second potential problem is that, even though k8s is a promising solution and de-facto standard in the CC environment, as previously stated, it might not be the best proposal for EC and geo-distributed  $\mu$ C environment. Despite all these potential problems, researchers show great work in adapting k8s architecture to work for geo-distributed workloads.

Ryden et al. [23] present a platform for distributed computing, that is more oriented towards user-based applications. Compared to other similar systems, their goal was not to develop a solution that will do resource management policy, on the contrary, their focus is more oriented to give flexibility to the users for application development. Users can develop their applications using exclusively the Javascript

programming language, with some embedded native code for a more efficient solution. The authors rely on a bunch of volunteer nodes to run all the tasks and applications, similar to work presented in [147]. The main difference between these two solutions is that Ryden et al. make a split on which nodes are storage nodes, and which nodes are used for calculation and processing tasks. Their application environment is protected from malicious code, using sandboxing techniques. This presents an interesting work to show how users can develop applications and run them in an EC environment.

In [151] Lèbre et al. show an interesting solution based on extending and adopting the OpenStack system. OpenStack is a free and open standard cloud computing IaaS platform for CC use cases in both public and private clouds. The authors tried to manage both cloud and edge resources using a NoSQL database. Massively distributed multi-site IaaS, using OpenStack is a challenging task [151] to implement, because the communication between nodes of different sites can be subject to important network latencies [151]. On the other hand, if it could be done properly we would gain one major advantage that users of the IaaS solution can continue using the same familiar infrastructure for both cloud and edge/fog use-cases.

Based on the literature survey, the Ning et al. presents current open issues of EC platforms [15]. In their work, authors focus on different aspects of EC systems, and they outline the importance of EC and CC tight collaboration. The CC needs to be unloaded and EC nodes could provide data pre-processing. On the other hand, EC needs massive storage and a strong computing capacity of CC. The authors illustrate the usage of edge computing platforms to build specific applications and conclude that with CC and EC integration, both sides will benefit.

In [130] the de Guzmán et al. present solution based on k8s that use k8s Deployment Manifests to reuse successful principles from k8s by creating a virtual machine for each Pod using Linuxkit. Their solution is based on the immutable infrastructure pattern, and instead of containers, they use the virtual machines as the unit of deployment.

Authors prove that the attack surface of their system is reduced since LinuxKit only installs the minimum OS dependencies to run containers. It represents interesting usage of LinuxKit to deploy OS dependencies and immutable infrastructure patterns, but VMs might be a bit problem for small devices, and ARM nodes as well as the complex flow of the k8s application model. Nonetheless, it is an interesting extension of the k8s framework and proves that LinuxKit can be used for immutable infrastructures with custom OS.

In [152] Sami et al. show an interesting platform for dynamic services distribution over Fog nodes using volunteer nodes. Their platform is tuned for container placement with relevance and efficiency on volunteering fog devices, near users with maximum time availability and shortest distance. They do this *on the fly* with improved QoS.

Besides academy efforts, the industry as well introduced a few interesting platforms and frameworks for EC. For example, Amazon introduced their framework Greengrass [153] that can run on various hardware to do some processing. Amazon turns to the option that their framework is deeply connected to the rest of the AWS cloud ecosystem. KubeEdge [154] is a lightweight extension of the k8s framework, to operate in an EC environment. The same as regular k8s, all workloads are done in the domain of a single cluster which might not be the best solution for geo-distributed  $\mu$ Cs. Both Greengrass and KubeEdge are frameworks that are mainly used for user-based applications. On the other hand, there is General Electric with its Predix [155] platform. Predix is a scalable platform used for industrial IoT applications.

Osmotic computing is a relatively new infrastructure paradigm aiming to decompose applications into services dynamically tailored for the smart environments exploiting resources in edge and cloud infrastructures [66]. The applications developed for this computing paradigm are deployed in the cloud and edge systems equalizing the microservices concentrations on both sides. This thesis can serve as a base for osmotic computing, allowing dynamic and efficient management of  $\mu$ C infrastructure to avoid application breakdown and degradation



of QoS.

### 3.3 Task offloading

As already mentioned in 2.1.4, EC nodes rely on the concept of data and computation offloading from the cloud closer to the ground [18], while heavy computation remains in the cloud because of resource availability [15].

Offloading is an effective strategy when using cloud services. Relying only on cloud services may be prone to introducing long latency, which some applications cannot tolerate. On the other hand, mobile devices and sensors do not have sufficient battery energy for task offloading [156]. The computation performance may be compromised due to insufficient battery energy for task offloading, so these devices might send their data to nearby EC nodes.

In literature, there are few platforms proposing task offloading [157, 18, 19, 65, 38, 156] to the nearby edge layer. These offloading techniques are based on different parameters, options, and techniques to put tasks to different sets of nodes in such a way that it won't drain mobile devices and sensors battery. After the computation is done, this edge layer sends pre-processed data to the cloud for further analysis, storage, etc.

When using task offloading techniques, it is very important to have good QoS that users can rely on. In [152] authors used Evolutionary Memetic Algorithm (MA) to solve their multi-objective container placement optimization problem to achieve better QoS.

One of the key challenges in the area of computation offloading is in the mismatch between how devices demand and access computing resources and how cloud providers offer them [157]. For example, it takes around 27 seconds to start single VM instance on the AWS, but the leasing time for the single VM instance is one hour. On the other hand, execution delay is stable in the CC, while in the EC is not due to the computational and transmission delay [158].

## 3.4 Application models

Sayed et al. in their work [44] describe that EC follows a decentralized architecture model and that data processing is at the edge of the network, thus it enables nodes to make autonomous decisions. So the applications written for EC can perform actions locally before connecting to the cloud at all. This will have some benefits like reducing network overhead issues as well as the security and privacy issues.

As already mentioned on page 62, Ryden et al. [23] present an interesting user-oriented platform for distributed computing called Nebula. In this section, their research is going to be dissected, but from a different angle. Nebula allows users to develop their applications using Javascript exclusively, due to the usage of Google Chrome Web browser-based Native Client (NaCl) sandbox [159] that can run Javascript code only. Restriction on a single language might be a problem for some users and use-cases even though Javascript is a popular language at the moment. On the other hand, virtual machines tend to be too resource-demanding packing stuff that might not be needed, so a solution using containers or unikernels might provide better resource utilization and pack more services per node than virtual machines.

In [40] Satyanarayanan et al. represent an interesting view on cloudlets as a “data center in a box”. They give an example that cloudlets should support a wide range of users, with minimal constraints on their software. They put emphasis on transient VM technology. The emphasis on transient VMs is because cloudlet infrastructure is restored to its pristine software state after each use, without manual intervention. At the time when they conducted their research, containers might not have been working solution or it might have been hard to use them. By modern standards, containers may even fit better, and pack more user software on the same hardware. This may be the case for the unikernels, once they reach a wider adoption rate and stable products.

Various k8s variants like [154, 150], give users the possibility to run

different applications like web servers and databases even on smaller devices creating green DC [21].

Satyanarayanan et al. [20] propose the concept of edge-native applications that will separate space into 3 layers or tiers. Tier **(1)** represents various mobile, IoT devices autonomous vehicles, etc, and these devices produce a lot of data. Tier **(2)** represents applications running in cloudlets or other EC models, that will be able to do some pre-processing, or data filtering before it goes further. Finally, tier **(3)** represents classic cloud applications that will accept pre-processed and filtered data from the previous tier, do more processing, react on some values, or store for future use. This idea represents an interesting concept and gives wide space for users and application development.

In [160] Beck et al. argue that applications should use message bus, streams, or topics because most mobile or edge applications are expected to be event-driven. The message bus system is an interesting proposition because the virtualized applications can subscribe to message streams, i.e., topics, and act only when data arrive. Applications might not be alive the whole time. And if for some reason mobile edge applications cannot reach a close EC server, it can always send data to the cloud. So cloud applications should be changed so slightly, just to cover this edge case.

In [45], Jararweh et al. show how integration between EC with CC principles will create more complex services and applications at the edge of the network opening new possibilities for applications to reduce the load on the centralized cloud model but also avoid bottlenecks and single points of failure.

In [161], authors propose solution that relies on the modularity of the microservices: **(1)** one application instance is at the edge, making the system robust to network partitions (local requests can still be satisfied), and **(2)** collaboration can be programmed with service mesh.

## 3.5 Infrastructure management

In the era of distributed systems, cloud computing, and microservices the open-source community and different companies provided various tools for the purpose of abstracting infrastructure at the level of software [133, 131]. These tools can be separated into two subgroups, based on how users send instructions to the systems on (1) declarative, and (2) imperative [133].

Newly created tools like Terraform, Pulumi, or CFEngine are representative of the declarative movement, using platform-independent language to specify configuration, policies, security, and much more. Users specify their intentions usually through some existing formats like *JSON* or *YAML*, while others use their domain-specific language to achieve the same goal.

In both cases, users do not specify explicit commands that the system needs to execute. Instead, they declare what they want to achieve, and leave it for the system to determine the optimal way of achieving it [133]. The main difference between an existing format, and a domain-specific language is how verbose specification is. Existing formats may be a little bit more verbose, but users are most likely already familiar with them. Developing domain-specific language takes more time and requires users to learn them from scratch, but the final outcome is usually a more concise, expressive and optimized set of instructions.

These tools are turning out to be very important in a multi-cloud environment [61]. The users need to specify artifacts, independently from the cloud provider, and let the system deal with the cloud provider specifics.

On the other hand, already existing and well-known tools like Chef, Ansible, and Puppet usually rely on some specific language, and the user needs to code the instructions that must be done to achieve the same or similar job [133].

This is more prone to error, since users may introduce a bug in the system that might be hard to debug and find. At the same time, these tools have existed for a long time, and there are existing best practices and a lot of available examples for users to utilize.

Every major cloud provider offers a proprietary solution, deeply integrated into their ecosystem. AWS offers CloudFormation, a configuration tool that allows users to code their infrastructure to automate deployments.

Microsoft Azure offers a Resource Manager that allows users to define the infrastructure and dependencies in templates, and organize dependent resources into groups that can be deployed or deleted in a single action, access control, and more.

Google Cloud offers a Deployment Manager that offers many features to automate cloud infrastructure stack. Users can create templates using *YAML* or *Python* programming language. They can preview changes before deployment, in a console user interface, and much more.

In terms of adoptability in  $\mu$ DCs on the edge, already existing tools might be used to set up infrastructure in both places, in the cloud and at the edge, if they allow extensibility. Declarative tools would be harder to adopt since they are built with different goals in mind. While imperative tools might be easier to utilize, they may introduce unnecessary complications to the system.

On the other hand, a newly created IaC system can communicate with existing tools via some standard interfaces to set up the cloud infrastructure, while a new, custom made engine could handle organizing geo-distributed nodes into  $\mu$ DCs [162].

With this strategy, we are getting the best of both worlds. The new user application model would be based on the existing cloud application model, offering users a fast and elegant way to develop new human-centered applications.

IaC has grown in popularity in recent years because it applies the same kind of version control and repeatability to the orchestration of the infrastructure as developers use for applications source code [135].

The second benefit is that configuration is decoupled from the system, it can more readily be deployed on a similar system elsewhere.

## 3.6 Thesis position

In the previous sections different aspects of EC and integration with the CC, influential research, interesting concepts, and implementations were described.

The focus of this thesis is to make the connection between CC and EC stronger, represented as a system that can descriptively and dynamically organize geo-distributed nodes that already exist over an arbitrary vast area into one coherent system —  $\mu$ Cs, that could be offered to users as a service. This approach is not fully addressed in the other solutions.

This thesis is influenced by the organization of CC and their big DCs, but adapted for a different environment such as EC and  $\mu$ Cs.

Adaptations that are required for such tasks must be followed by a clear Separation of concerns (SoC) model and intuitive applications model so that users can fully use new-formed infrastructure properly.

All these adaptations made to traditional CC model will make it possible to push the whole solution more towards EC as a service and  $\mu$ C model that can help CC with latency issues with new-age applications.

# Chapter 4

## Micro clouds

In this section, the core idea of this thesis will be presented. We will introduce the model of distributed  $\mu$ Cs and show how such a system could be dynamically formed and formally described. The  $\mu$ Cs is based on the traditional CC organizational model but adapted for different environments, allowing the creation of distributed clouds.

Throughout this section, we are going to rely on the research presented in Chapters 1 and 3, and make a connection with the existing CC model.

In Section 4.1 we present a high overview of the system – an architecture that is influenced by the standard CC model but adapted for a different environment. Section 4.2 introduces the separation of concerns model for previously defined system architecture. Section 4.3 presents a possible application model and how users can utilize newly created architecture fully. Application models are based on existing development models, so that transition is fairly easy. Possibilities how this system could be offered as a service, and which options can be offered to the potential users for developing their applications is discussed in Section 4.4. Section 4.5 presents the desired option for infrastructure and application deployment, and discusses why immutability is important in a  $\mu$ C geo-distributed environment. Section 4.6 discusses the importance of formal models in computer science, and DS in particular.

This Section also presents formal models for all protocols used for the proper system formation and operation. Section 4.7 shows transactions that appear in such complex system, and garbage collection done on resources that are not used anymore. Section 4.8 shows system observability mechanisms important for geo-distributed environments like  $\mu$ Cs. Section 4.9 shows access patterns that could be used in  $\mu$ Cs and geo-distributed environments. Section 4.10 presents the possibility how micro clouds could be scaled automatically, while Section 4.11 presents flow of user data in micro clouds model. This chapter is concluded with Section 4.13 which presents repercussion of the proposed model, and how it can be used as a stand-alone model where other features could be implemented on top of that or used as service for other, existing, systems.

## 4.1 Configurable Model Structure

In their work [20] Satyanarayanan et al. propose a new architecture pattern and separation into the three tiers, where every tier or layer has a distinct role. This idea is a very powerful one because applications can now be split into parts and optimized for the specific role in the global picture. On the other hand, too many moving parts mean more problems, and the whole system is potentially more prone to errors and failures.

If the previous model is taken and combined with  $\mu$ DCs and a zonally based server organization, a good starting point for building  $\mu$ C infrastructure, and EC as a service is achieved. This extension is an interesting move because the computing power and storage capacity that serves the nearby population can then be extended. This base model is just a starting point that is promising, but to make a fully functional model, it needs to be made more available and resilient with less latency. To achieve such behavior, these concepts have to be extended and adapted for different usage scenarios, but a geo-distributed idea from our vision should not be lost.



To extend the system in a new direction, we can look for some inspiration in existing systems that are proven and working. This is especially important in DS, so we want to lower down the complexity and avoid known problems by sticking to existing models and algorithms. When the CC design is observed, it can be seen that every single part in that system contributes to a more resilient and scalable system. On the high level, the CC architecture is separated into few building blocks that make the whole system lot easier to understand, maintain and operate.

The first building block of CC architecture is a cluster, and a cluster can be defined as a set of nodes or servers that operate as a single unit to achieve some goal. This is where resources are, and this is where user applications are running on. The next building block that consists of multiple clusters is called Regions (or DCs). Regions are isolated and independent from each other, and they contain resource application needs. These resources come in form of clusters. Regions are usually composed of a few availability zones [28]. These zones are the defense against the fail. If for whatever reason, one zone fails or goes offline, there are still more of them to serve user requests, there is a better availability, scalability, and resilience.

If we now observe  $\mu$ Cs as geo-distributed systems with dispersed users [163], we can use a similar model with some adaptations. Rely on a similar proven strategy, do not build the entire model, but adapt the existing one for the different use-case.

Table 4.1 presents similar concepts between CC and edge-centric computing (ECC) concepts. The accent is here put on the difference between the physical logical concepts in the two models.

Edge centric computing	Cloud computing
Topology (logical)	Cloud provider (logical)
Region (logical)	Region (physical)
Cluster (physical)	Zone (physical)

Table 4.1: Similar concepts between cloud computing and ECC.

Since we are talking about geo-distributed systems, our scenario is a little bit different than the one used in the standard CC model. We can still combine multiple EC nodes into clusters, that is what  $\mu$ DCs already propose. If we want to go a little bit further, we can define a cluster as a:

**Definition 4.1.1.** *A  $\mu$ C cluster is a group of nodes that are virtually and/or geographically separated, that works together to provide the same service to clients.*

Multiple node clusters could be combined into the next bigger logical concept of *region*. A region will increase the availability and reliability of both the system and applications.

But the region in CC and ECC is not fully the same thing. In the standard CC model, the region is a physical thing or element of the existing infrastructure [28], while in the ECC a region could be viewed differently, not as a physical element but rather as a logical element. A formal definition of a *region* in ECC can now be given as:

**Definition 4.1.2.** *In a geo-distributed environment like ECC and  $\mu$ Cs, a concept of region is used to describe a set of clusters (that could be) scattered over an arbitrary geographic region. Regions are fully capable to accept or release clusters in the same way that clusters can accept or release nodes.*

In  $\mu$ DCs, a cluster is as big as the population nearby that is using it [24]. If this is combined with the previous definition, then the minimum size of an ECC region can formally be defined as:

**Definition 4.1.3.** *Geo-distributed regions are composed of at least one cluster but can be composed of much more to achieve a more resilient, scalable, and available system.*

With the previous definition, we have to be careful not to introduce huge latency in the system. To lower the region latency, the vast distances between clusters should be strongly avoided in normal circumstances.

In the CC model, new nodes have to be brought and connected physically to the rest of the infrastructure [29], while in the ECC, just by changing the definition to which the specific cluster belongs, we can extend the region.

Multiple regions should be able to form a second logical layer – *topology*. In the geo-distributed systems such as ECC or  $\mu$ Cs, topology can formally be described as:

**Definition 4.1.4.** *In geo-distributed systems, topology represents the highest logical concept that is composed of a minimum of one region and could span over multiple regions. Topology is fully capable to accept new or release the existing regions.*

With these simple abstractions, any geographical region can be easily covered with the ability to shrink or expand clusters, regions, and topologies. Size and formation of *clusters*, *regions*, and *topologies* should be a matter of need, agreement, and usages of nearby population similar to the size of  $\mu$ DCs [24], and modeling in Big Data systems [30, 31]. This separation and organization gives one interesting feature. A more natural administrative division of some region can be followed, and resources can be organized by population usages.

The organization of these concepts should be fully optional. So for example, we could fit clusters in an interval of  $n$ DCs [27] and  $\mu$ DCs [24] or as wide as the whole city or as small as all devices in a single household and everything in between. Formally, the size of some cluster  $\mu$ C cluster ( $\mu$ Cc) can be represented like:

$$\mu Cc \in [ nDCs, \mu DCs ] \quad (4.1)$$

If we go a little bit further, we can represent the city as one region, where parts of the city are organized into clusters. A city topology can even be formed by splitting the city into multiple regions containing multiple clusters, and ultimately a country topology can be formed by

splitting the country into regions, with cities being clusters.

Nodes that belong to the same cluster should run some form of membership protocol presented in 2.1.3. Gossip style protocols, like SWIM [32] (see page 23), are standard in the cloud environment. The same strategy could be applied here, used in conjunction with replication mechanisms [33, 34, 35] making the whole system more resilient. It could be upgraded to better reflect current network conditions [164], or to enable **self-healing** properties in the system [165].

Replication could be used not only in nodes inside the cluster, but data can also be replicated in clusters inside the region and even in regions inside topology. This property should exist, but it should be controlled by users depending on how resilient and available the system he/she wants and needs. In [166] Simić et al. take a look from a theoretical point of view on CRDTs usage, to achieve SEC in EC. The authors conclude that CRDTs could be a natural fit to EC as long as we are aware of the potential pitfalls of CRDTs.

Single *topology* reflects one CC provider, so multiple topologies are forming  $\mu$ Cs that can help CC with huge latency issues, pre-processing in huge volumes of data, and relax and decentralize strict centralized CC model.

These  $\mu$ Cs have much fewer resources compared to standard clouds, but they are much closer to the user meaning they have a much faster response. In the case of storage, if data is not present at the time of user request, they can pull data from the cloud and cache it for later. Formally, the position of  $\mu$ Cs in between CC and EC like:

$$\mu Cs \in ( \textit{Edge computing}, \textit{Cloud computing} ) \quad (4.2)$$

Exclusiveness in the previous formula means that the position of  $\mu$ Cs is moveable in between EC and CC depending do we want our model to be more CC-like or more EC-like, but it should not become one of

them. It could be an integral part of both of these models (even at the same time) [59, 61].

To achieve such elasticity, we must abstract the infrastructure to the level of software, creating "infrastructure programming" [131], allowing  $\mu$ Cs infrastructure to be managed similarly as the software is.

Three-tier architecture with numerous clients at the bottom,  $\mu$ Cs in the middle, and cloud on the top seem to resemble cache level architecture in CPU [167].

Figure 4.1. shows the three-tier architecture, with the response time and resource availability.

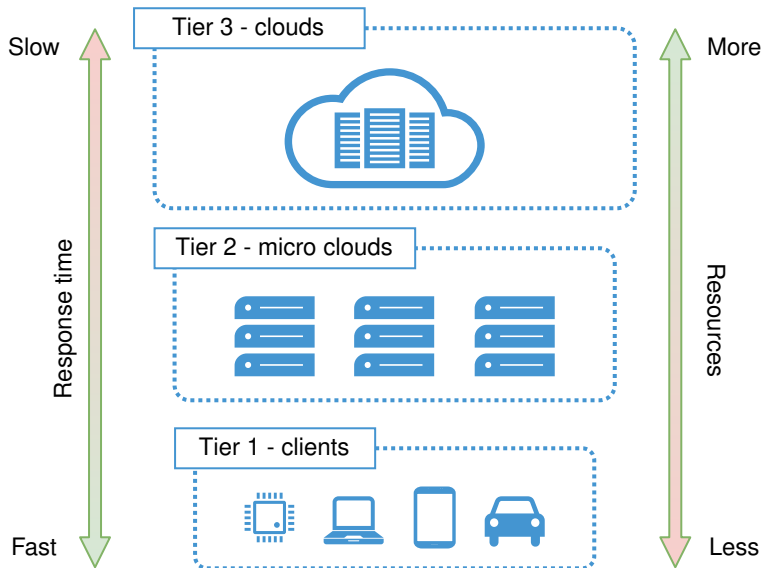


Figure 4.1: Three tier architecture, with the response time and resource availability

On lower levels, response time is the fastest, since data is processed closer to the source. At the same time, there is a limited storage capacity and processing power. As we go on the upper tiers, there is more and more storage capacity and processing power, but the response

time is higher and higher, especially when distance and huge volumes of data that need to be moved to the cloud are considered.

In everything as a service model [36], ECC as a service fits in between CaaS and PaaS, depending on the user needs.

## 4.2 Separation of concerns

In his work, Jin et al. [39] introduces three core concepts to fully describe physical services, and specifies their relationships. Concepts that authors propose are: **(1)** Devices, **(2)** Resources, and **(3)** Services. This separation is interesting because we can rely on it in a geo-distributed environment, to describe SoC for  $\mu$ Cs.

One of the most important part of every system, is the SoC model. This is especially important if a platform to be offered as a service is going to be created. With some adaptations, our SoC model can be based on concepts proposed by Jin et al. [39] in three layers, as depicted in Figure 4.2.

The layer on the very bottom of the three-tier architecture (see Figure 4.2) consists of various devices. These devices are important because they represent main *data creators* but at the same time, they are main *services consumers*.

The layer in the middle (see Figure 4.2) represents resources or EC nodes. These resources have a very important spatial feature, and as such, they indicate the range of their hosting devices [39]. This means that the developers at any given time must know the topology of the system, the resource spread, and utilization across clusters. Besides that main information, users must know the state and health of every application. If some EC node desires to be a part of the system, a node must obey four simple rules:

- (1) A node must run an operating system with a usable file system;
- (2) A node must be able to run some isolation engine for applications, for example, containers or unikernels;

- (3) A node must have available resources for utilization so that applications can be run or data stored;
- (4) A node must have internet connection.

One thing that is important to notice, is that all nodes in clusters, regions, and topologies are **equal** and there are **no special nodes**. Every node that joins the system must be able to store and process information.

Last but certainly not least layer represents services. These services expose resources through some interface and make them available over the internet [39]. These services respond to the client requests immediately, if possible, or cache information [40, 41] for future use. Unlike the previous two layers, services have one specific feature and span over two tiers of the system (see Figure 4.2):

- (i) Services that exist in the  $\mu\text{C}$ , and they are responsible for filtering and data pre-processing before sending it to the cloud. Or cache information that was not available on a previous user request for some future use;
- (ii) Services in the standard cloud that should be able to accept pre-processed data, and they are responsible for computation and storage that is beyond the capabilities of ECC nodes. Services in the cloud should be able to take direct requests from the clients in a case when something catastrophic happens to the  $\mu\text{C}$  that is close to the user, and it is not able anymore to accept user requests.

This kind of services separation creates new application model that we present in detail in the section 4.3.

Figure 4.2. shows the proposed SoC for every layer of the ECC as a service model.

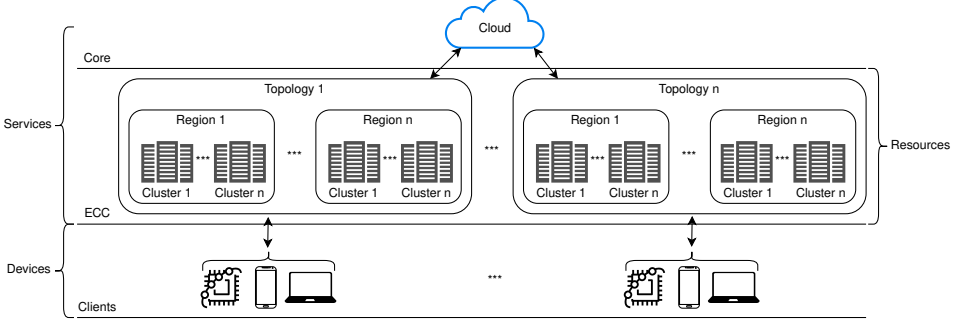


Figure 4.2: ECC as a service architecture with separation of concerns.

### 4.3 Applications Model

Modern-day applications that should run in the cloud are advised to follow the cloud-native model [168]. With this approach, we get applications that are easier to scale, more available, and less error-prone when compared to traditional web applications [168].

Satyanarayanan et al. present the edge-native applications model [20] that should use the full potential of EC infrastructure, and keep good features of their cloud counterparts.

When we introduced SoC for the  $\mu$ C model (see page 79), we described services that have one specific feature, and that is that these services span over two tiers of the system. This information is crucial for application development because we want to use the full potential of formed infrastructure, and keep good features of their cloud counterparts.

To satisfy the previously defined SoC model (see page 79), existing applications can be split into the **front** and **back** processing services. The front processing service is an edge-native application running inside some previously formed cluster to minimize latency, while the back service runs in the traditional cloud as a cloud-native application to leverage greater resources.



The front processing service will handle user requests coming to the nearby cluster, and communicate with the back processing services when needed to synchronize the pieces of information, cache data, or pass filtered or pre-processed data. Separation like that gives developers better flexibility and large design space when creating and optimizing their workloads. At the same time, it ties together the traditional clouds and edge computing into one coherent system. Figure 4.3 shows  $\mu$ C applications model that ties edge and traditional cloud applications.

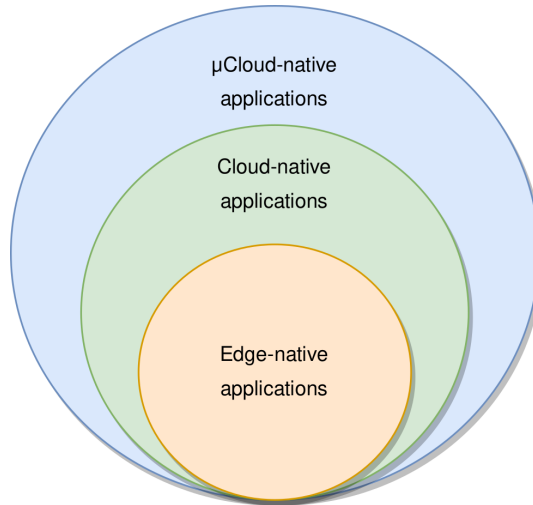


Figure 4.3:  $\mu$ C applications model ties traditional clouds and edge computing applications

With this model, we venture even deeper into understanding and applying the concept of **data locality** (see page 25). Since we have front and back processing services, we are committed to doing processing data closer to their source, instead of moving data to the cloud. This reduces latency and saves some users money on storing and processing unnecessary data.

### 4.3.1 Execution models

Frontend services model should be packed in some standard way [2.7](#) and deployed in the wild, as an event-driven application, with the subscription policy to message streams, or infinite sequences [\[169\]](#) using topics [\[160\]](#). The processing strategy is in the developer's hands, depending on the nature of the use-case.

If a service is subscribed to some stream using a topic, for example, it is natural that events appear in their stream. There are two strategies for building a large-scale time-ordered event system:

- (1) **Fan out on write**, with this strategy, every service has some form of inbox, and when an event appears for some topic, that data is copied to every service that is subscribed to that topic. With these options, reads are fast, but writes are not. The more subscribers to the topic, the longer it takes to persist all updates;
- (2) **Fan in on read**, with this strategy, every topic has a sort of outbox where it can store data. When services read their streams, the system will read the recent data from the outboxes. Writes are fast, and used storage is minimal, but reads are difficult because to do this properly on a request-response deadline is not a trivial task to do.

To implement those ideas without complicated synchronization, CRDTs could be used [2.1.1](#). Companies like Soundcloud and Bet365 are already using them for the same or similar tasks.

Some examples of applications may include:

- (1) **Events** that notify users if some value is above or below some defined threshold;
- (2) **Stream** or processing data as it comes to the system;
- (3) **Micro-batch** processing delivers data more slowly than real-time data processing but faster than typical batch processing. It is

performed in smaller batches, than traditional batch processing. This type of processing is for applications that require high-throughput processing over large data streams [170];

- (4) **Batch** does processing in predefined times over some bigger collection of data;
- (5) **daemons** or processing that are that do some tasks or executions in the background without explicit user intervention, usually their execution could be time defined but it is not mandatory;
- (6) **Services** or applications that would operate over standard request-response model or some variation of that protocol. For example, the NATS messaging system has a request-reply form implemented over topics;
- (7) **Other**, something that falls outside these models, or it is a composition of multiple operations at once. This can be represented as a function composition in mathematics:

$$(f \circ g)(x) = f(g(x)) \tag{4.3}$$

This type should get events from some topic as they arrive, and a user can define his strategy what needs to be done and how it needs to be processed. Users can contact other existing services or combine them (eq. 4.3), and the user is responsible for optimization;

These types of applications could be implemented in many ways like those discussed in section 2.2.2, or some adapted variant of those models, and should have a clear communication interface offered to users or applications and other services.

Users should be given the option to develop their applications in various available languages, not forcing them to use a strict one. Users must also be advised about outcomes of their choice.

For example, some languages might be slower or use more resources than others due to virtual machine execution, or some other tooling that is required to be started as well. The second important thing would be that users can do proper logging and tracing of their services. If something fails, the user must be able to go over the logs and traces to find problems.

### 4.3.2 Packaging options

Because of their nature,  $\mu$ Cs could be most likely composed of ARM devices. These devices in many cases are not able to run full VMs because of their hardware restrictions. In recent years there have been advances in VMs technology and its ability to run VMs on ARM devices. In [171] Ding et al. show such possibility to run VMs on ARM devices.

But even if VMs are fully compatible with ARM devices, we still inherit VMs large footprint, already discussed in section 2.7 if we try to use them *asis*.

On the contrary, containers, and unikernels give us **more or less** the same functionality, but use fewer resources, which means that more services can be run in containers and even more in unikernels. Until unikernels are fully ready to be used, though, they will fall in nice to have a category, and we should stick to containers. But even with containers we need to be aware of their limitations and pitfalls and know that there is no comprehensive solution and there is not just one solution for all scenarios.

At the moment, containers are a more mature solution than unikernels and require fewer resources than VMs. On top of that, there are numerous tools already existing using containers that could be utilized. Knowing all this, at the first stages of  $\mu$ Cs, containers should be an option to go.

In [172] Simić et al. show benefits of using containers in large scale edge computing systems from the theoretical point of view, by looking into architecture difference. Authors focus on differences between VMs

and containers in cases where services need to run on ARM devices with limited resources.

In terms of boot times, throughput, and memory consumption, unikernels are famous for providing excellent performance [173]. In the future and when unikernels are more matured and tested, they could be used for particular use-cases and applications, especially like events or serverless implementations. The containers will probably not be fully replaced, but they can co-exist with unikernels in some cases where more control over running a single function is needed.

Like any other system, users can create variants of the systems and different flavors optimized for certain solutions. In that cases, they may favor one solution over the other one. In general, containers and unikernels should be the preferred way to package, run, and distribute user applications in a  $\mu$ Cs environment.

### 4.4 As a service model

Users should be able to develop their applications using familiar models, depending on their needs. Depending on how much control a user requires over the process of service scheduling, resource selection, etc. three models can be distinguished:

- (1) **Micro PaaS or  $\mu$ PaaS**, where the platform does all the management and offers a simple interface for developers to deploy their applications. This model is similar to PaaS, and the only difference is that this runs in  $\mu$ C and should synchronize with CC;
- (2) **Micro CaaS or  $\mu$ CaaS**, if users require more control over resource requirements, deployment, and orchestration decisions. This model is similar to CaaS (or even IaaS), and the only difference is that this runs in  $\mu$ C and should synchronize with CC;

- (3) **Micro SaaS or  $\mu$ SaaS**, with this option, data is not synchronized with the CC, but all processing and storage is done in the  $\mu$ C. As such, this option should be used vigilantly.

It is important to note, that both variants, **mCaaS** and **mPaaS**, should **not** stand alone, at least not for now. Using **mSaaS** is not being advised for now, since that would require that the whole application be running **only** in the  $\mu$ Cs. In the future, when EC nodes gain more power and storage, this might change. **mSaaS** as an option should more investigated in the future.

Both options, **mCaaS** and **mPaaS**, could be included, and a part of any cloud model (see Section 2.1.2), multi-cloud or offered separately.

## 4.5 Immutable infrastructure

This thesis proposes a DS model that is built with a three-tier architecture that should operate in the geo-distributed environment.

We believe that **immutable deployment** model would be a good fit for such environment. It will simplify the deployment process, since we want to rely on *atomic operations* and do not want to leave the misconfigured system at any level [162]. If something like that happens, there will be a problem that would be hard to properly track, address and resolve.

Geo-distributed  $\mu$ Cs model that is described in this chapter should operate in two levels of deployment that are built one on top of the other:

- (1) **Infrastructure deployment**, update and change should be atomic and immutable. Users should do changes declaratively – mutations of the system, by telling the system what the new state should be, and let the system figure out the best way how to do user-specified changes. In this category, we can account for any change that is doing on the cluster, region, or topology

that user(s) operate on. For example create the new cluster, region, topology, or doing configurations of the setup system. The only change that could be done by imperative strategy updates on the nodes themselves. But even for this strategy, it would be beneficial if we could use declarative way **if possible**. It is important to notice that **mutation** does not mean in place change, but just the name of the operation. This deployment strategy is reserved for operations people or (eg. DevOps or SREs), but if the company or team is small any developer could do this. Developers should not be dealing with this part of the deployment;

- (2) **Services deployment**, makes sense only if the previous action is taken. We must have infrastructure already set up, to put any sort of services (applications) into the system. Like the previous model, this should be done declaratively as well, and all changes should be done immutably without an in-place change. The user should specify his new state or “view of the word” declaratively and let the system do all the changes he wants. All user services should be packed as described in section 4.3.2 because this simplifies the way services are put to the nodes. When done properly, this allows operations people to do faster changes with almost zero downtime deployments with all strategies already discussed in section 2.8. This part of the deployment should be done by developers since they did implementation and testing. They know how many resources they need for their service, what type of service they had developed. This deployment could be done in collaboration between operations and developers if a company is big or time is properly separated.

It is important to notice, that both deployments should be closely followed, for possible errors and problems so that users can act accordingly. These deployment messages, logs, and traces [101] should be

stored in a centralized log system, for convenient lookup, alerting, and reporting.

Separation like this simplifies deployment and usage for both application development spectrums:

- (i) **Operations** who should be dealing with infrastructure deployment, tooling set up, applications deployment, monitoring, and in the general health of applications and infrastructure;
- (ii) **Developers** who should be dealing with the development of the services, their interactions, and cloud to  $\mu$ C and vice-versa synchronization.

Only with tight collaboration with those two development roles, such a complex system like one presented in this chapter can be alive, well, and serving user requests without collapse.

It is important to note that every action in the system should be logged and traced properly. Since we are dealing with multi-tier architecture, a chance that something will fail is increased. Logs and traces should be available to operations people who are responsible for the infrastructure maintenance.

At the same time, actions should be traced by the company that offers these  $\mu$ Cs. For security reason, any of these logs and traces should not be visible to others **but** responsible individuals, and the level of details and personal pieces of information should be different.

### 4.5.1 Deployment in micro clouds

As described in section 2.8, there are a few options when it comes to setup and deployment of infrastructure and/or applications. Let us explore interesting options how both (1) infrastructure, and (2) services could be deployed in  $\mu$ Cs.

*Rolling update* strategy minimize the downtime, at the cost of update speed.



Figure 4.4 shows flow chart for *rolling update* in  $\mu$ Cs for both services and infrastructure.

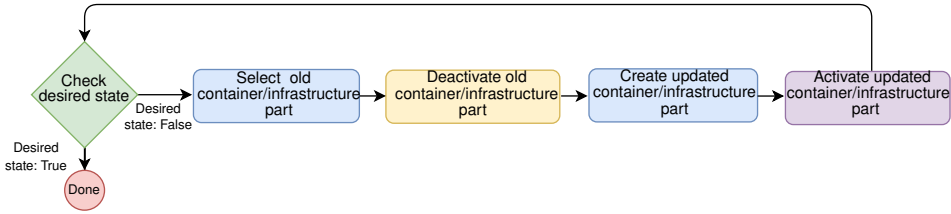


Figure 4.4: Rolling update flow chart

On the other side of the spectrum, the users should be able to choose the *recreation* strategy as well. This strategy will causes higher downtime, but it will do updates quickly.

Figure 4.5 shows flow chart for *recreation strategy* in  $\mu$ Cs for both services and infrastructure.

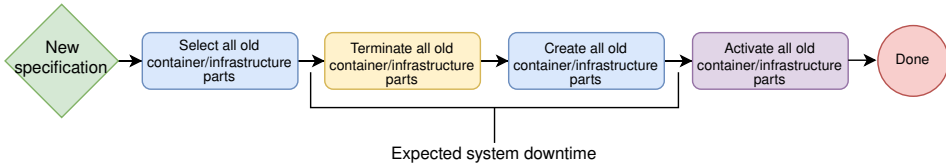


Figure 4.5: Recreation strategy update flow chart

The third strategy is the *canary update* strategy. This strategy offers a partial update process allowing to test new environment version , without a commitment to a full rollout.

This strategy comes handy when we want to test how your services perform in real-life scenarios with production data and users, allowing the revert to previous version fast.

In this strategy, the deployment creates a few new artifacts while keeping most artifacts on the previous version. The ration between new and previous version should be provided by the user. After the

test, we can go either direction: **(1)** full rollout, **(2)** revert to previous version. Both directions are predicated by the results of the test.

Figure 4.6 shows chart for the *canary update* strategy deployment.

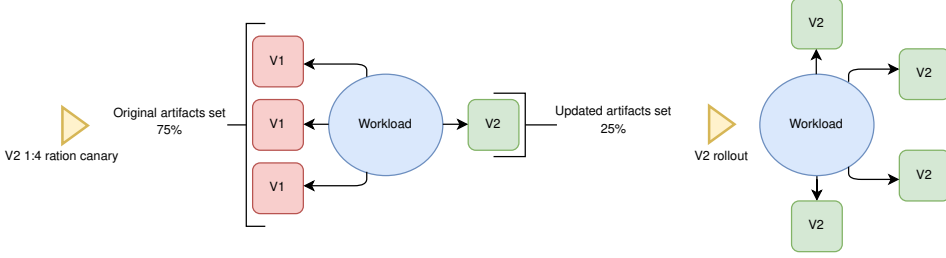


Figure 4.6: Canary update 1:4 ration update chart

## 4.6 Formal model

Ensuring the reliability and correctness of any DS is a very difficult task, and should be mathematically based. Formal methods are techniques that allow us to create specifications and verification of complex (software and hardware) systems based on mathematics and formal logic. There are several options for how to formally describe DS: TLA+ [174], *pi* calculus [175], combinational topology [176], asynchronous session types/multiparty session types (MST) [43], etc. This theis rely on *explicit connection actions* [42], an extension of MST.

Unfortunately, because of their nature DS cannot always be formally described by any of the existing techniques [43, 174, 175, 176]. There are a lot of variables that could influence this. A formally described and correct model can save hours, days, and even months of hard debugging, testing to reveal all bugs and problems in the system that may only happen in some specific circumstances that are hard to initiate. Therefore, it is recommended and beneficial to try to formally describe system that is developing.

Infrastructure deployment will not happen overnight, and it might take years. It might not be started at all until the whole process is trivial

[40], and this is a complicated task [45]. Because of those properties, the key problem that needs to be resolved is how to simplify ECC or  $\mu$ C management. The naive approach would require going to every node and do it manually. This process is super tedious and time-consuming, especially if we consider a geo-distributed environment.

In such a complex environment, formal models are of great help if we can model and prove that protocols that the system relies on are correct. The system we propose tackles this issue using remote configuration and it relies on four formally modeled protocols:

- (1) **Health-check protocol** informs the system about state of every node (see Section 4.6.2);
- (2) **Cluster formation protocol** forms new clusters dinamically (see Section 4.6.3);
- (3) **Idempotency check protocol** prevents system from creating existing infrastructure (see Section 4.6.4);
- (4) **List detail protocol** shows the current state of the system to the user (see Section 4.6.5);

These three protocols are based on the geo-distributed Infrastructure deployment.

### 4.6.1 Multiparty asynchronous session types

Communication protocols that operate from node to the system can be modeled using [42], an extension of *MPST* [43] – a class of behavioral types specifically designed for describing distributed protocols that rely on asynchronous communications.

The type specifications give us one additional benefit. They are not only useful to formally describe protocols, but also reliable for a modeling-based approach developed in [42] to validate our protocols are they satisfy multiparty session types safety (there is no reachable

error state) and progress (an action is eventually executed, assuming fairness).

The modeling process is done in two steps.

- (1) **The first step** in modeling the communications of a system using MPST theory is to provide a *global type*, that is a high-level description of the overall protocol from the neutral point of view. Following [42], the syntax of global types are constructed by:

$$G ::= \{p \uparrow q_i : \ell_i(T_i).G_i\}_{i \in I} \mid \mu t.G \mid t \mid \text{end} \quad (4.4)$$

where  $\uparrow \in \{\rightarrow, \Rightarrow\}$  and  $I \neq \emptyset$ . In the above,  $\{p \uparrow q_i : \ell_i(T_i).G_i\}_{i \in I}$  denotes that *participant*  $p$  can send (resp. connects) to one of the participants  $q_i$ , for  $\uparrow = \rightarrow$  (resp.  $\uparrow = \Rightarrow$ ), a *message*  $\ell_i$  with the *payload* of *sort*  $T_i$ , and then the protocol continues as prescribed with  $G_i$ .  $\mu t.G_1$  is a recursive type, and  $t$  is a recursive variable, while **end** denotes a terminated protocol. We assume all participants are (implicitly) disconnected at the end of each session (cf. [42]).

The advance of using approach of [42], when compared to standard MPST (e.g., [43]), is in a relaxed form of choice (a participant can choose between sending to different participants), and,  $\Rightarrow$ , that explicitly connects two participants, hence (possibly) dynamically introducing participants in the session. Both of these features will be significant for modeling our protocols (we will return to this point);

- (2) **The second step** in modeling protocols by MPST is providing a syntactic projection of the protocol onto each participant as a local type, that is then used to type check the endpoint implementations. We use the definition of projection operator given in [42, Figure 2.9]. In essence, the projection of global type  $G$  onto

participant  $p$  can result in  $S_p = q!\ell(T) \dots$  (resp.  $S_p = q!!\ell(T) \dots$ ) when  $G = p \rightarrow q:\ell(T) \dots$  (resp.  $G = p \twoheadrightarrow q:\ell(T) \dots$ ), and, dually,  $S_p = q?\ell(T) \dots$  (resp.  $S_p = q??\ell(T) \dots$ ) when  $G = q \rightarrow p:\ell(T) \dots$  (resp.  $G = q \twoheadrightarrow p:\ell(T) \dots$ ), while the projection operator “skips” the prefix of a global type if participant  $p$  is not mentioned neither as sender nor as receiver. Furthermore, a local type must be represented by the following syntax:

$$S ::= +\{q_i\alpha\ell_i(T_i).S_i\}_{i \in I} \mid \mu t.S \mid t \mid \text{end} \quad (4.5)$$

where  $\alpha \in \{!, !!\}$  or  $\alpha \in \{?, ??\}$  (in which case  $q_i = q_j$  must hold for all  $i, j \in I$ , to ensure consistent external choice subjects, cf. [42, Page 6.]), and  $I \neq \emptyset$ . Interested reader can find details in [42].

For simplicity reasons, we will consider that all participants are communicating in a single private session. All sent messages, but not yet received, are buffered in a single queue that preserves their order. The order is preserved only for pairs of messages having the same sender and receiver, while other pairs of messages can be swapped since these are asynchronously independent.

### 4.6.2 Health-check protocol

Nodes that exist in the clustered environment usually have a channel where they can send metrics and other data in a form of a health-check mechanism. This channel could be used this channel to reach nodes to send some actions to them, for example, a cluster formation message.

In figure 4.7. a low-level health-check protocol between a single node and the rest of the system can be seen. This process involves the following participants: Node, Nodes, State, and Log.

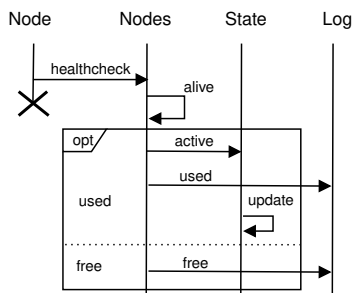


Figure 4.7: Low level health-check protocol diagram.

These participants included in Figure 4.7 follow the health-check protocol that is informally described below:

- (1) **Node** sends a health-check signal to the Nodes service;
- (2) **Nodes** accept health-check signals from every node, update node metrics and if node is used in some cluster, inform that cluster about the node state;
- (3) **State** contains information about nodes in the clusters, regions and topologies;
- (4) **Log** contains records of operations. Users can query this service.

Nodes service will be informed about node existence on his health-check ping. However, the system state will not be updated, changed, or even informed about this ping if the node is not used in some cluster.

In the following algorithm 1 steps required by the system to determine if the node is free or used, and how node information is stored, are described.

---

**Algorithm 1:** Health-check data received

---

```

input : event, config
1 if isNodeFree(event.id) then
2   if exists(event.id) then
3     renewLease(event.id, config.leaseTime);
4     updateData(event.id, event.data);
5   else
6     leaseNewNode(event.id, config.leaseTime, event.data);
7     saveMetrics(event.id, event.metrics);
8   end
9 else if isNodeReserved(event.id) then
10  updateData(event.id, event.data);
11 else
12  renewLease(event.id, config.leaseTime);
13  updateData(event.id, event.data);
14  saveMetrics(event.id, event.metrics);
15  sendNodeACK(event.id);
16 end

```

---

To formally describe servers or nodes (terms are used interchangeably) properties in the system, set theory could be used. At the beginning, the system will have an empty server set  $S$  denoted with  $S = \emptyset$ . To determine node state, we could use existing node properties. One approach may be that we use a node-id structure, for example.

Whatever approach is used, free nodes can be formally described as follows:

**Definition 4.6.1.** *Nodes are free if and only if (henceforth iff) they do not belong to any cluster.*

If **node-id** is taken, for example, when the new health-check message from the particular node is received, **node-id** structure will determine

the node state. When the node is free, it will home some random or user-defined id, while when it is used in some cluster, the node-id structure will reflect this.

If, for example, there are  $n$  free nodes in the wild, this can be denoted with  $s_i$ , where  $i \in \{1, \dots, n\}$ . If all of them send the health-check ping to the system, we need to determine their state. If a node  $s_i$  is free, we should add it to the server set, and thus we have:

$$S_{new} = S_{old} \cup \bigcup_{i=1}^n \{s_i\} \quad (4.6)$$

It is important to notice, that the order in which messages arrive is not important. The only thing that is important is that every message eventually comes into the system. Node roles in the system can now formally be defined like:

**Definition 4.6.2.** *All nodes in the system are equal, no matter if they are a part of some cluster or if they are not.*

The previous definition gives us a strong background in the further formal definition of the system because the only thing that should be cared about is that node is alive and well, and that it is ready to accept some jobs.

We described in algorithm 1 how the system stores the node data, but also how to determine if the node is free or used using the node-id structure. Every node  $s_i$  in the system that is part of the server set  $S$  could be described as a tuple  $s_i = (L, R, A, I)$ , where:

- $L$  is a set of ordered key-value pairs:  $L = \{(k_1, v_1), \dots, (k_m, v_m)\}$  where  $k_i \neq k_j$ , for each  $i, j \in \{1, \dots, m\}$  such that  $i \neq j$ .  $L$  represents node labels or server-specific features. Labels were based on Kubernetes [150] labels concept, which is used as an elegant binding mechanism for its components.



- $R$  is a set of tuples  $R = \{(f_1, u_1, t_1), \dots, (f_m, u_m, t_m)\}$  representing node resources, where  $f_i, u_i, t_i$ , for  $i \in \{1, \dots, m\}$  are as follows:
  - $f_i$  is the free resource value,
  - $u_i$  is the used resource value, and
  - $t_i$  is the total resource value.
- $A$  is a set of tuples  $A = \{(l_1, r_1, c_1, i_1), \dots, (l_m, r_m, c_m, i_m)\}$ , representing running applications, where  $l_j, r_j, c_j, i_j$ , for  $j \in \{1, \dots, m\}$ , are as follows:
  - $l_j$  represents labels, the same way we used for node labels,
  - $r_j$  is the resource set application requires,
  - $c_j$  is the configuration set application requires, and
  - $i_j$  is the general information like name, port, developer.
- $I$  represents a set of general node information like: name, location, IP address, id, cluster id, region id, topology id, etc.

If we want to assign  $m$  (fresh) labels to the  $i_{th}$  server, we start with empty labels set  $s_i[L] = \emptyset$ , then we add labels to server. Therefore, we have

$$s_i[L]_{new} = s_i[L]_{old} \cup \bigcup_{j=1}^m \{(k_j, v_j)\} \quad (4.7)$$

We can now formally define the number of labels per single server  $s_i$  in the system like:

**Definition 4.6.3.** *Every node from the server set  $S$  **must have** non-empty set of labels. The number of labels for every server  $s_i$  in the server set  $S$  may vary.*

Since labels are an important part of the system (more in future sections), they should be picked carefully and agreed on upfront. It should also be possible to change them if such a thing is required.

Labels should stick out some distinctive features of the node, that might be valuable for developers or administrators to target. For example, server resources, server features (e.g., SSD drive), geolocation, etc. They can be created as follows:

**Definition 4.6.4.** *Labels are created using arbitrary long alphanumeric text, for both keys and values, separated by colon sign. For example `os:linux`, `arch:arm`, `model:rpi`, `cpu:2`, `memory:16GB`, `disk:300GB`.*

Following all things presented above, we can now give a formal description for the low-level health-check communication protocol (see Figure 4.7). The global protocol  $G_1$  (given below) conforms the informal description given at page 94: **node** connects **nodes** with `health_check` message and a payload of type  $T_1$  required by the system to properly register node into the system.

Based on the received information, **nodes** **either** connect **state** with `active` message, informing the node status alongside payload typed with  $T_2$  (containing informations required by the system to properly register active health-check sender), and then also connect **log** with the same message, **or** directly connect **log** informing the node is *free*.

$$G_1 = \text{node} \rightarrow \text{nodes:health\_check}(T_1). \\ \left\{ \begin{array}{l} \text{nodes} \rightarrow \text{state:active}(T_2). \text{nodes} \rightarrow \text{log:used}(T_2). \text{end} \\ \text{nodes} \rightarrow \text{log:free}(T_2). \text{end} \end{array} \right.$$

Notice that in  $G_1$  we indeed have a choice of **nodes** sending either to **state** or to **log**. Such communication patterns are impossible to be modeled using just standard MPST approaches. Also, notice that **state** will be introduced into the session only when receiving from **nodes**. Hence, if the session after the first ping from **node** to **nodes** proceeds with the second branch (i.e., connecting **nodes** with **log**) then

`state` is not considered as stuck, as it would be in standard MPST, such as, e.g., [43], but rather idle.

Projecting global type  $G_1$  onto participants `node`, `nodes`, `state` and `log` we can then get local types as follows:

$$S_{\text{node}} = \text{nodes}!!\text{health\_check}(T_1).\text{end}$$

$$S_{\text{nodes}} = \text{node}??\text{health\_check}(T_1). \\ + \begin{cases} \text{state}!!\text{active}(T_2).\text{log}!!\text{used}(T_2).\text{end} \\ \text{log}!!\text{free}(T_2).\text{end} \end{cases}$$

$$S_{\text{state}} = \text{nodes}??\text{active}(T_2).\text{end}$$

$$S_{\text{log}} = + \begin{cases} \text{nodes}??\text{used}(T_2).\text{end} \\ \text{nodes}??\text{free}(T_2).\text{end} \end{cases}$$

where, for instance, type  $S_{\text{nodes}}$  specifies `nodes` can receive the ping message from `node`, after which it will dynamically introduce either `state` or `log` into the session, where in the former case it also connects `log` (but now with message *free*).

### 4.6.3 Cluster formation protocol

Another communication protocol that the system relies on, appears in the cluster formation process, where users can form new clusters dynamically. Two different actions are distinguished:

- (1) The first action is user-system communication. Here user sends query parameters to the system to obtain a list of available nodes that satisfy specified query parameters;

- (2) The second action is a little bit more complicated than the previous one, and it starts when the user sends a message to the system with a new assembly specification. In this setting, the system involves participants: User, Queue, Scheduler, State, Nodes, Log, and NodesPool. These participants need to cooperate to successfully form new clusters, regions, or topologies dynamically, adhering to the scenario shown in Figure 4.8.

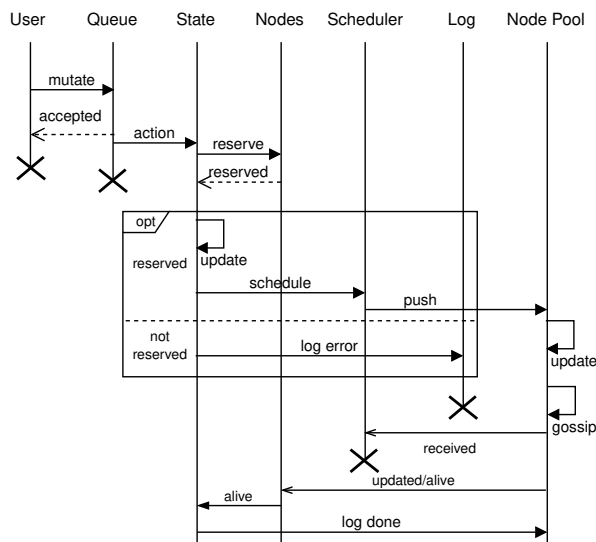


Figure 4.8: Low level cluster formation communication protocol diagram.

These participants included in Figure 4.8 follow the cluster formation protocol that is informally described below:

- (1) **User** query Nodes service, based on some predefined query parameters. User sends a new creation message to Queue. User either gets response *ok* if the message is accepted or *error* if the message cannot be accepted due to missing rights or other issues. This operation is called *mutation*;

- (2) **Queue** accepts a user message, and passes it to State. Messages are handled in FIFO (First In, First Out) order. The queue prevents system congestion, with received messages. The queue will also test if the specified message has already been handled before over idempotency check;
- (3) **State** accepts mutation messages from Queue, and tries to store new information about the cluster, region, or topology. If Nodes can reserve all desired nodes, the system will store new user desired specification and send a message to Scheduler to physically carry on the creation of clusters with desired nodes;
- (4) **Nodes** accept messages from State. It will reserve desired nodes, if possible, otherwise, it will send an error message to Log service. On a health-check message, it will either just store node information or if a node is used in some cluster, inform that cluster that the node is alive;
- (5) **Scheduler** waits for a message sent from State, and physically carry on cluster formation by pushing cluster formation messages to the chosen nodes;
- (6) **Log** contains records of all operations. Users can query this service to see if their tasks are finished or some problems occurred;
- (7) **Nodes Pool** represents the set of  $n$  free nodes that will accept mutation messages. When a node receives the message, it will follow some predefined steps:
  - (i) start gossip protocol to inform other nodes from the mutation message that they should form a cluster;
  - (ii) when cluster formation is done, send an event to Scheduler and Nodes that node is alive and can receive messages. The cluster formation is done, when all nodes have complete list of nodes that should form the cluster;

If a user wants to get a list of free nodes in the system, he must create a query using *the selector*, which is the set of key-value pairs where he can describe what type of nodes he desires. Algorithm 2 describes steps that are required to perform a proper node lookup based on a received selector value.

---

**Algorithm 2:** Nodes lookup

---

```

input : query
1 Initialize: nodes  $\leftarrow []$ 
2 foreach  $node \in freeNodes()$  do
3   if  $len(node.labels) == len(query) \wedge node.haveAll(query)$ 
4     then
5       nodes.append(node)
6   end
7 end
8 return nodes

```

---

We start with the empty selector  $Q = \emptyset$ , in which we append key-value pairs. Hence, when a user submits a set of  $p$  key-value pairs we have that

$$Q_{new} = Q_{old} \cup \bigcup_{i=1}^p \{(k_i, v_i)\} \quad (4.8)$$

Once the user submits query selector to the system with desired attributes, for every server in the set  $S$ , two things need to be checked:

- (1) the cardinality of the  $i_{th}$  server's set of labels and the query selector are identical in size

$$|s_i[L]| = |Q|, \text{ and} \quad (4.9)$$

- (2) every key-value pair from query set  $Q$  is present in the  $i_{th}$  server's labels set  $s_i[L]$ , hence the following predicate must yield true:

$$P(Q, s_i) = \left( \forall (k, v) \in Q \exists (k_j, v_j) \in s_i[L] \text{ such that } k = k_j \wedge v \leq v_j \right) \quad (4.10)$$

The  $i_{th}$  server from the server set  $S$  will be present in the result set  $R$ , iff both rules are satisfied so we have:

$$R = \{s_i \mid |s_i[L]| = |Q| \wedge P(Q, s_i), i \in \{1, \dots, n\}\} \quad (4.11)$$

If the result set  $R$  is not the empty set, we then reserve nodes for configurable time so that other users cannot see, and try to use them. Finally, reserved nodes with message data  $md$  are added to the task queue set:

$$TQ_{new} = TQ_{old} \cup \{(R, md)\}. \quad (4.12)$$

When the task comes to execution, the task queue will send messages to every node that is specified. Algorithm 3 describes the steps required for cluster formation.

---

**Algorithm 3:** Clustering formation message

---

**input :** request, config  
1 nodes  $\leftarrow$  searchFreeNodes(data.query)  
2 reserveNodes(nodes, config.time)  
3 pushMsgToQueue(nodes, data)  
4 key  $\leftarrow$  saveTopologyLogicState(data)  
5 watchForNodesACK(key)

---

Users are free to override existing node labels with their labels or keep predefined ones when including nodes in the cluster. If the node is free, or the user did not change the node labels on cluster formation, the system will use default labels like node geo-location, resources, operating system, architecture, etc.

When the node receives the cluster formation message, he will automatically pick and contact a configurable subset of nodes  $R_g \subset R$ , and start the gossip protocol, propagating pieces of information about nodes in the cluster (e.g, new, alive, suspected, dead, etc.). When

every node inside the newly formed cluster has a complete set of nodes  $R$  obtained through gossiping, the cluster formation process is over. Topology, region, or cluster formation should be done descriptively using YAML, or similar formats.

In the Algorithm 4 describes the required steps after nodes receive a cluster formation message.

---

**Algorithm 4:** Node reaction to clustering message

---

```

input : event
1 switch event.type do
2   case formationMessage do
3     updateId(event.topology, event.region, event.cluster)
4     newState  $\leftarrow$  updateState(event.labels, event.name)
5     sendReceived(newState)
6     nodes  $\leftarrow$  pickGossipNodes(event.nodes)
7     startGossip(nodes)
8   end
9 end

```

---

In the following, a low-level cluster formation communication protocol (see Figure 2.5) is described. We are using the same extension of MPSTs [42] used for the health-check protocol.

Global protocol  $G_2$  (given below) conforms the informal description of the cluster formation protocol given on page 100. The protocol starts with **user** connecting **state** by message *query* and a payload typed with  $T_1$  that contains user query data, and then **state** forwards the message by connecting **nodes**. Then, the protocol possibly enters into a loop, specified with  $\mu t$ , depending on the later choices. Further, **nodes** replies a response *resp* to **state**, that, in turn, forwards the message to **user**. The payload of the message is typed with  $T_2$  that has response data, based on a given query. At this point, **user** sends to **state** one of three possible messages:

- (1) *mutate*, and the mutation process, described with global protocol  $G'$ , starts;



- (2) *quit*, in which case the protocol terminates; or,
- (3) *query* – this means the process of querying starts again, the query message is forwarded to **nodes** and the protocol loops, returning to the point marked with  $\mu t$ .

The third branch is the only one in which protocol loops. Also, we can notice that **user** – **state** and **state** – **nodes** are connected before specifying recursion. Hence, even after several recursion calls, these connections will be unique. So it is not required to disconnect them before looping.

$$\begin{aligned}
 G_2 = & \text{user} \rightarrow \text{state:query}(T_1).\text{state} \rightarrow \text{nodes:query}(T_1). \\
 & \mu t.\text{nodes} \rightarrow \text{state:resp}(T_2).\text{state} \rightarrow \text{user:resp}(T_2). \\
 & \begin{cases} \text{user} \rightarrow \text{state:mutate}().G' \\ \text{user} \rightarrow \text{state:quit}().\text{end} \\ \text{user} \rightarrow \text{state:query}(T_1).\text{state} \rightarrow \text{nodes:query}(T_1).t \end{cases}
 \end{aligned}$$

The mutate protocol  $G'$ , activated in the first branch in  $G_1$ , starts with **user** sending **create** message to **state**, specifying also information about new user desired state typed with  $T_3$ , and **state** replies back with *ok*. Then, **state** sends *ids* of the nodes to be reserved (specified in the payload typed with  $T_4$ ) to **nodes**, that, in turn sends one of the two possible messages to **state**:

- (i) *rsrvd*, denoting all nodes are reserved and the protocol proceeds as prescribed with  $G''$ , or
- (ii) *error*, with error message in the payload, informing there has been unsuccessful reservation of nodes, in which case **state** connects **log** reporting the error and the protocol terminates.

$$\begin{aligned}
 G' = & \text{user} \rightarrow \text{state:create}(T_3).\text{state} \rightarrow \text{user:ok}(). \\
 & \text{state} \rightarrow \text{nodes:ids}(T_4). \\
 & \begin{cases} \text{nodes} \rightarrow \text{state:rsrvd}().G'' \\ \text{nodes} \rightarrow \text{state:err}(\text{String}).\text{state} \rightarrow \text{log:err}(\text{String}).\text{end} \end{cases}
 \end{aligned}$$

Finally, in  $G''$  **state** connects **sched** (Scheduler) with message *ids* and the payload that contains other data imported for mutation to be completed (typed with  $T_5$ ). Then, **sched** connects **pool** (Nodes Pool) with *update* specified with  $T_6$ , after which **pool** replies back with *ok*, and connects to **nodes** sending new id's *nids* typed with  $T_4$  ( that contains successfully reserved user desired nodes). Now **nodes** notifies **state** the action was successful, that in turn connects **log** with the same message, and the protocol terminates.

$$\begin{aligned}
 G'' = & \text{state} \rightarrow \text{sched:ids}(T_5). \text{sched} \rightarrow \text{pool:update}(T_6). \\
 & \text{pool} \rightarrow \text{sched:ok}(). \\
 & \text{pool} \rightarrow \text{nodes:nids}(T_4). \text{nodes} \rightarrow \text{state:succ}(). \\
 & \text{state} \rightarrow \text{log:succ}(). \text{end}
 \end{aligned}$$

We may now obtain the projections of global type  $G_2$  onto the participants **user**, **state**, **nodes**, **log**, **pool**, and **sched**:

$$\begin{aligned}
 S_{\text{user}} = & \text{state!!query}(T_1). \mu \mathbf{t}. \text{state?resp}(T_2). \\
 & + \left\{ \begin{array}{l} \text{state!mutate}(). \text{state!create}(T_3). \text{state?ok}(). \text{end} \\ \text{state!quit}(). \text{end} \\ \text{state!query}(T_1). \mathbf{t} \end{array} \right.
 \end{aligned}$$

$$\begin{aligned}
 S_{\text{state}} = & \text{user??query}(T_1). \text{nodes!!query}(T_1). \mu \mathbf{t}. \text{nodes?resp}(T_2). \\
 & \text{user!resp}(T_2). \\
 & + \left\{ \begin{array}{l} \text{user?mutate}(). \text{user?create}(T_3). \text{user!ok}(). \text{nodes!ids}(T_4). S' \\ \text{user?quit}(). \text{end} \\ \text{user?query}(T_1). \text{nodes!query}(T_1). \mathbf{t} \end{array} \right.
 \end{aligned}$$

where

$$S' = + \left\{ \begin{array}{l} \text{nodes?rsrvd}(). \text{sched!!ids}(T_5). \text{nodes?succ}(). \text{log!!succ}(). \text{end} \\ \text{nodes?err}(\text{String}). \text{log!!err}(\text{String}). \text{end} \end{array} \right.$$

$$\begin{aligned}
 S_{\text{nodes}} = & \text{state}??\text{query}(T_1).\mu\text{t.state!resp}(T_2). \\
 & + \left\{ \begin{array}{l} \text{state?ids}(T_4).+ \left\{ \begin{array}{l} \text{state!rsrvd}().\text{end} \\ \text{state!err}(\text{String}).\text{poll}??\text{nids}(T_4). \\ \text{state!succ}().\text{end} \end{array} \right. \\ \text{state?query}(T_1).\text{t} \end{array} \right.
 \end{aligned}$$

$$S_{\text{log}} = + \left\{ \begin{array}{l} \text{state}??\text{succ}().\text{end} \\ \text{state}??\text{err}(\text{String}).\text{end} \end{array} \right.$$

$$S_{\text{pool}} = \text{sched}??\text{update}(T_6).\text{sched!ok}().\text{nodes!!nids}(T_4).\text{end}$$

$$S_{\text{sched}} = \text{state}??\text{ids}(T_5).\text{pool!!update}(T_6).\text{pool?ok}().\text{end}$$

For instance, type  $S_{\text{sched}}$  specifies that participant **sched** gets included in the session only after receiving from **state** message *ids*, then **sched** connects **pool** with *update* message, after which it expects to receive *ok* message and finally terminates.

We remark that global type  $G_2$  could also be modeled directly by using standard MPST models (such as [43]). However, in such models, the projection of  $G_2$  onto, for instance, participant **sched** would be undefined (cf. [42]). Since we follow the approach of [42] with explicit connections, projection of  $G_2$  onto **sched** is indeed defined as  $S_{\text{sched}}$ .

#### 4.6.4 Idempotency check protocol

Mutate operation should be *atomic*, *immutable*, and *idempotent*. The user can specify the same topology details but in a different order, for example. We must ensure that the new cluster formation protocol should **not** be initiated, if the user changes order of regions, clusters, nodes, or labels in one or more node(s).

If *mutate* operation fails, **i.e.**, the same infrastructure already exists, the user should get a message that infrastructure is **already** formed, but the new formation protocol would **not** be initiated. If the user changes the number of labels per node, nodes per cluster, or similar, a new protocol **should** be initiated.

But since we have a different scenario than standard write to storage, and we do not specify steps on how operations should be done, to implement idempotency correctly we have to do it a little bit differently. First of all, we must use proper structure, and ensure that operation done over that structure is idempotent.

The idempotent structure  $I_S$  can be represented as a tuple of *topology name* and *data set* like:

$$I_S = (Name, DataSet) \quad (4.13)$$

*Name* could be used for faster lookup, while *DataSet* represents the set of data, because most set operations are idempotent, as described in SEC (see page 17). *DataSet* could be represented in two ways:

- (1) **Flat keyspace**, with this option all data could be part of the same set, and distinguishment could be achieved using *prefix* identity, for example: `region1_cluster1_node1_labels`, `region1_cluster1_node1_name` etc.;
- (2) **Hierarhical keyspace**, with this option we can create nested data-structures of elements, for example set of regions, where every region is a set of clusters, etc. We can go deep as long as we want, but we **must not** violate idempotency throughout the hierarchy — every (sub)structure **must** be idempotent. So we can use a set of sets or *N-ary trees*.

If we have cached idempotency data for user requests, and if a user tries to send the same request again, we can then test idempotency using set operation **intersection** because the intersection is an idempotent operation following the next proof.

*Proof.* Intersection of two sets  $x$  and  $y$   $x \cap y$  is an idempotent operation, because  $x \cap x$  is always equal to  $x$ . This means that the idempotency law 2.3  $\forall x, x \cap x = x$  is always *true*.  $\square$

If we have stored the cluster formation protocol request in some *idempotent storage*, and we receive a new request with the same name, **then** we can take the intersection of two sets. If we get the same set, that action is already done because of the definition 4.6.4. Otherwise, the request represents the new action, and new cluster formation protocol.

This can be made a little bit faster, by choosing the proper storage structure. If we first do (the same) topology, the name is already present in the idempotency store, and this lookup will spare us unnecessary comparison on sets. For example, data structures like *Hash tables* store element as pair of *key – value* and offers time and space complexity for lookups  $\mathcal{O}(1)$ , *on average* [177], and use Bloom filter [178] for key filtering. We can go even a little bit further and use CRDTs and SEC to store and replicate data required for the idempotency test. Figure 4.9 shows zoomed view in the *State* participant from figure 4.8, and idempotency check communication.

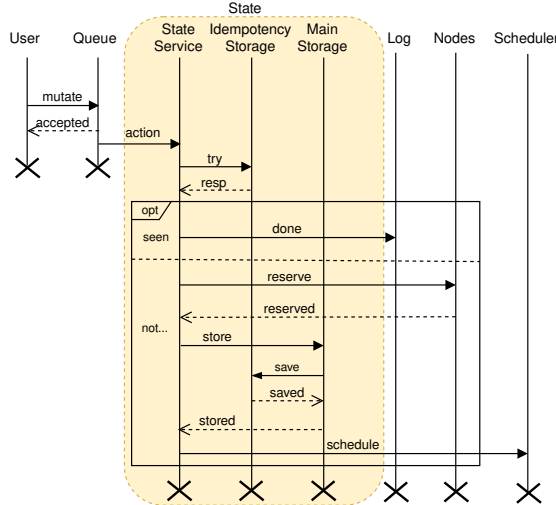


Figure 4.9: Low level view of idempotency check communication.

The participants follow the communication that we now describe informally:

- (1) **User** sends a list request to State service (same as [4.6.3](#));
- (2) **Queue** accepts the list request and the query local state based on the user selector. If a detailed view is required, the state gets metrics data from Nodes service (same as [4.6.3](#));
- (3) **State service** is a wrapper around system main storage. Interacts with main storage in order to create new topologies, regions or clusters, or to get data from the main storages about the same entities.
- (4) **Idempotency storage** contains idempotent set for all **already** formed topologies, regions and clusters.
- (5) **Main storage** contains records about the desired state for all formed topologies, regions, and clusters.
- (6) **Log** contains records of operations. Users can query this service to see if their tasks are finished or have any problems (same as [4.6.3](#));
- (7) **Nodes** accept messages from State. If possible, it will reserve desired nodes, otherwise, it will send an error message to Log service. On a health-check message, if a node is used in some cluster, it will inform that the node is alive (same as [4.6.3](#));
- (8) **Scheduler** waits for a message sent from State, and pushes cluster formation messages to desired nodes (same as [4.6.3](#));

When testing idempotency, we must have in mind that stored structure could be: (1) flat keypace, and (2) hierarchical keypace. The both options are valid, as long as the **structure is idempotent**, and we can do **idempotent operations** over that structure. For example,

*set* as a data structure and *intersection* are good candidates. It is worth noticing that the flat keypace is a special case of the hierarchical keypace. The algorithm that will test structure idempotency must be able to test both options.

Algorithm 5 describes steps required to test if the operation is done before.

---

**Algorithm 5:** Mutate idempotency check

---

```

// children are sorted in both trees
1 function equals(stored, topology):
2   | sChildren  $\leftarrow$  stored.children
3   | tChildren  $\leftarrow$  topology.children
4   | if topology == NULL  $\vee$  tChildren.size  $\neq$  sChildren.size
5   |    $\vee \neg$  stored.isIdempotent(topology.data) then
6   |     | return false
7   |   while sChildren.hasNext  $\wedge$  tChildren.hasNext do
8   |     |   if  $\neg$  equals(sChildren.next, tChildren.next) then
9   |     |     | return false
10  |   return true
    input : request
11 stored  $\leftarrow$  storage.findRequest(request.payload.name)
12 if stored  $\neq$  NULL then
13   | return false
14 else
15   | topology  $\leftarrow$  request.payload.topology
16   | return  $\neg$  equals(stored, topology)
17 end

```

---

A formal description of the idempotency check protocol (see Figure 4.9) by using [42] is presented next. The global protocol  $G_3$  (given below) conforms the informal description given at page 110. At this point, the user can choose one of two possible messages:

- (1) *quit*, in which case the protocol terminates; or,

- (2) *mutate*, and the mutation process, described with global protocol  $G'$ , starts;

$$G_3 = \begin{cases} \text{user} \rightarrow \text{service:quit}().\text{end} \\ \text{user} \rightarrow \text{service:create}(T_3).\text{state} \rightarrow \text{user:ok}().G' \end{cases}$$

The mutate protocol  $G'$ , activated in the first branch in  $G_3$ , starts with **user** sending **create** message to **state**, specifying also information about the new user desired state typed with  $T_3$ , and **state** replies back with *ok*. This process is the same as cluster formation protocol (see Section (1)). Now we start specific communication protocol for idempotency check so **service** sends payload  $T_3$  to **istorage** to test if this request is *seen* before. The **istorage** responds with payload  $T_6$  to **service**, and based on *Boolean* response **service** can do one of two things:

- (1) **service** sends message to **log** and this process terminates; or,  
 (2) **service** sends *ids* of the nodes to be reserved (specified in the payload typed with  $T_4$ ) to **nodes**;

same as cluster formation protocol (see Section (1)). For simplification, we can assume that all nodes are reserved, and now idempotency data store global protocol  $G''$ , starts;

$$\begin{aligned} G' &= \text{service} \rightarrow \text{istorage:try}(T_3). \\ &\quad \text{istorage} \rightarrow \text{service:resp}(\text{Boolean}). \\ &\quad \begin{cases} \text{service} \rightarrow \text{log:done}(\text{String}).\text{end} \\ \text{service} \rightarrow \text{nodes:ids}(T_4).\text{nodes} \rightarrow \text{service:rsrvd}().G'' \end{cases} \end{aligned}$$

The idempotency store protocol  $G''$ , starts with **service** sending mutation payload  $T_3$  to **mstorage**. Then **mstorage** sends the same payload to **istorage**. When data is saved for future testing, **istorage** responds



back to `mstorage`, and finally `mstorage` responds back to `service`. At this point user payload  $T_3$  is stored in both main storage and idempotency storage for future testing. Protocol continues with `state` connects `sched` (Scheduler) with message `ids` and the payload that contains other data imported for mutation to be completed (typed with  $T_5$ ), and the rest of cluster formation protocol may continue.

$$\begin{aligned} G'' = & \text{service} \rightarrow \text{mstorage:store}(T_3).\text{mstorage} \rightarrow \text{istorage:save}(T_3). \\ & \text{istorage} \rightarrow \text{mstorage:saved}().\text{mstorage} \rightarrow \text{service:stored}(). \\ & \text{service} \rightarrow \text{sched:ids}(T_5).\text{end} \end{aligned}$$

We may now obtain the projections of global type  $G_3$  onto the participants `user`, `service`, `istorage`, `mstorage`, `log`, `nodes`: and `sched`:

$$S_{\text{user}} = + \begin{cases} \text{service!quit}().\text{end} \\ \text{service!create}(T_3).\text{service?ok}().\text{end} \end{cases}$$

$$S_{\text{service}} = + \begin{cases} \text{user?quit}().\text{end} \\ \text{user?create}(T_3).\text{user!ok}().S' \end{cases}$$

where

$$\begin{aligned} S' = & \text{istorage!!try}(T_3).\text{istorage?resp}(\text{Boolean}). \\ & + \begin{cases} \text{log!!done}(\text{String}).\text{end} \\ \text{nodes!!ids}(T_4).\text{nodes?rsrvd}().S'' \end{cases} \end{aligned}$$

$$\begin{aligned} S_{\text{mstorage}} = & \text{service??store}(T_3).\text{istorage!!save}(T_3). \\ & \text{istorage?saved}().\text{service!stored}().\text{end} \end{aligned}$$

$$\begin{aligned} S'' = & \text{mstorage!!store}(T_3).\text{mstorage?stored}(). \\ & \text{sched!!ids}(T_5).\text{end} \end{aligned}$$

```
Sstorage = service??try(T3).service!resp(Boolean).
           mstorage??save(T3).mstorage!saved().end
```

```
Slog = service??done(String).end
```

```
Snodes = service??ids(T4).service!rsrvd().end
```

Similarly, as for  $G_2$ , we remark  $G_3$  could also be modeled using standard MPST (e.g., [43]), but again the projection types would be undefined while following the approach of [42] with explicit connections, all valid projections have been obtained.

#### 4.6.5 List detail protocol

The final communication protocol in our system appears in the information retrieval process. Using labels, the user can specify what part of the system he wants to retrieve, namely, on formed topologies. This protocol could be useful, for example, if the user wants to visualize his topologies, regions, or clusters on some dashboard and monitor for some changes, alerts, etc.

This operation may span over multiple services to retrieve complete informations about the clusters and/or nodes. For this purpose some of the distributed queries methods (see Section 2.2.2.1) can be used. For some patterns api composition may be more suited, while for others CQRS may be the better solution to optimize queries.

This protocol comes with two available options:

- (1) **global view** of the system – all topologies the user manages. This will return just basic information about regions and clusters and their utilization;

- (2) **specific clusters** details – in-depth details for specified clusters like resources utilization over time (using stored metrics information), node information, configuration data, and running or stopped services.

It is important to note, that similarly to the query operation (defined previously), both rules (4.9) and (4.10) **must** be satisfied for information to be presented in the response. The user can specify one additional information in the list request, and that is whether or not the user wants a detailed view or not. If such information is presented in the request, the user will get a detailed view back.

Figure 4.10 shows a low-level view of the list operation protocol, where users can get details about the formed system. This setting involves the next participants: User, State, Nodes, and Log.

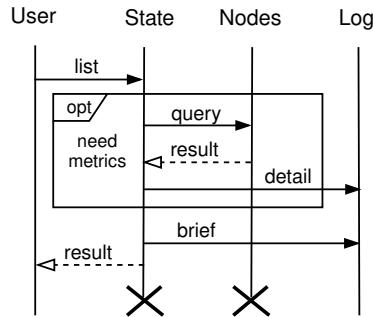


Figure 4.10: Low level view of list operation communication.

Now we can informally describe the roles of the participants in the protocol:, shown in figure 4.10

- (1) **User** sends a list request to State;
- (2) **State** accepts the list request and the query local state based on the user selector. If a detailed view is requested, the state gets metrics data from Nodes, and return details back to user;

- (3) **Nodes** contain node metrics data, and if required, it may send this data to State;
- (4) **Log** contains records of all operations. Users can query this service.

Algorithm 6 describes steps that are required after the state receives a list message.

---

**Algorithm 6:** List of current state of the system

---

```

input : request
1 Initialize: data  $\leftarrow []$ 
2 foreach (topology, isDetail)  $\in$  userData(request.query) do
3   if isDetail then
4     | data.append(topology.collectData())
5   else
6     | data.append(topology.data())
7   end
8 end
9 return data

```

---

A formal description of the list communication protocol (see Figure 2.11) by using [42] is presented next. Global type  $G_4$  (given below) starts with **user** connecting **state** with one of the two possible messages: (1) *list*, specifying a request for a detailed view, where sort  $T_1$  identifies which parts of the system the user wants to view in details, after which **state** connects **nodes** with *query* message, with a payload of sort  $T_2$  containing specification of which nodes need to show their metrics data, and then protocol proceeds as prescribed with  $G'$ ; (2) *list\**, specifies no need for a detailed view is specified, where a payload of sort  $T_4$  denotes user specified parts of the system the user wants to view, but without greater details. In the latter case protocol follows global type  $G''$ .

$$G_4 = \begin{cases} \text{user} \rightarrow \text{state} : \text{list}(T_1). \text{state} \rightarrow \text{nodes} : \text{query}(T_2). G' \\ \text{user} \rightarrow \text{state} : \text{list}^*(T_4). G'' \end{cases}$$

Global type  $G'$  starts with **nodes** replying to **state** *result* message and a payload identifying parts of the system user wants to see in greater detail typed with  $T_3$ . Then, **state** connects **log** with **details** and also sends **result** to **user**, and finally terminates. In  $G''$ , **state** also connects **log** with *brief* and a payload typed with  $T_5$  identifying parts of the system the user wants to see without greater detail. Then, **state** replies to **user** with **result** message, and the protocol terminates.

$$G' = \text{nodes} \rightarrow \text{state}:\text{result}(T_3).\text{state} \rightarrow \text{log}:\text{detail}(T_3). \\ \text{state} \rightarrow \text{user}:\text{result}(T_3).\text{end}$$

$$G'' = \text{state} \rightarrow \text{log}:\text{brief}(T_5).\text{state} \rightarrow \text{user}:\text{result}(T_5).\text{end}$$

Same as for the health-check and the cluster formation protocols, here we also present the projections of global type  $G_4$ , modeling the list protocol, onto participants **user**, **state**, **nodes**, and **log**:

$$S_{\text{user}} = + \begin{cases} \text{state}!!\text{list}(T_1).\text{state}?\text{result}(T_3).\text{end} \\ \text{state}!!\text{list}^*(T_4).\text{state}?\text{result}(T_5).\text{end} \end{cases}$$

$$S_{\text{state}} = + \begin{cases} \text{user}??\text{list}(T_1).\text{nodes}!!\text{query}(T_2).S' \\ \text{user}??\text{list}^*(T_4).\text{log}!!\text{brief}(T_5).\text{user}!\text{result}(T_5).\text{end} \end{cases}$$

where

$$S' = \text{nodes}?\text{result}(T_3).\text{log}!!\text{detail}(T_3).\text{user}!\text{result}(T_3).\text{end}$$

$$S_{\text{nodes}} = \text{state}??\text{query}(T_2).\text{state}!\text{result}(T_3).\text{end}$$

$$S_{\text{log}} = + \begin{cases} \text{state}??\text{detail}(T_3).\text{end} \\ \text{state}??\text{brief}(T_5).\text{end} \end{cases}$$

For instance, type  $S_{\log}$  specifies `log` gets included in the session only after receiving from `state`, either message *detail*, or message *brief*, and then terminates.

Similarly as for  $G_2$  and  $G_3$  we remark  $G_4$  could also be modeled using standard MPST (e.g., [43]), but again the projection types would be undefined, while following the approach of [42] with explicit connections, all valid projections have been obtained.

## 4.7 Long-lived transactions in micro clouds

To operate  $\mu$ Cs properly, such a system needs to be scalable. Architecture composed of loosely coupled services can be a way to go, because of all benefits such systems offer (see Section 2.2.2). But still, we must be aware of all problems that come along with them.

That being said, one of the problems we have to deal with are transactions that appear in such distributed system — distributed transactions (see Section 2.5.1). Because of the nature of the system *sagas* seems like a better pattern to use (see Section 2.5.2).

When a user submits new *cluster creation message* (see Section 4.6.3), the message will be accepted by the system, and the new task with *PENDING* state will be registered. If for whatever reason (e.g., no available resources or nodes, etc.), the system cannot proceed further, that task is terminated. It goes to *FAILED* state, concluding the transaction. Otherwise, no errors occurred, and the system can proceed with the cluster formation protocol, changing the task state to *IN PROGRESS*.

In this state, the system needs to do several things: **(1)** save newly formed cluster information, **(2)** prepare metrics service, **(3)** add watchers for the cluster nodes health-check, etc. This spans multiple services, creating a few sub-transactions in the process. The *task state* prevents the users to apply other tasks (e.g., configurations, actions, autoscaling, etc.) on a not yet fully formed cluster.

Separation on sub-transactions allow us to invoke the rollback mechanism, if the process yields any errors. Other strategy would be to start some of the retry strategies in order to fix the occurred issue.

If no errors occur during the process, the *cluster formation transaction* finishes, and the task state is changed to *CREATED*. Otherwise, the *cluster formation transaction* ends, *without creating the cluster*. The task state will be changed to *FAILED* state, concluding the transaction.

Figure 4.11 shows state diagram for cluster formation message.

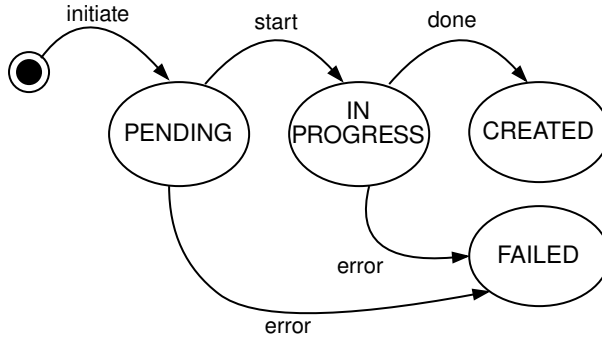


Figure 4.11: State diagram for cluster formation message

Cluster formation protocol will most likely span over multiple services, and strategies like sagas can be useful to implement transactions and rollbacks properly.

### 4.7.1 Garbage collection in micro clouds

When the user submits the task to the system, and if, for whatever reason that task fails, there is the option to immediately remove that task, and possibly all task connected items. Or, the system should have a specific process in the background, which only purpose is garbage collection to restore unused and fragmented system resources.

The second option seems more acceptable because one item can have a graph of items that are connected to that specific item. This

graph can be complex, and immediate deletion may slow down the entire system. In the background, the garbage collection process can mark unused items and delete them when the time is right.

Users should also be able to delete certain items with their dependencies as well, invoking *cascading deletion*. Dependencies may be deleted immediately, or marked as *orphans*, and deleted by the garbage collection process.

It makes sense that users can influence garbage collection by submitting their policies when and how often to do this operation, to delete only items or container images as well, and so on.

## 4.8 System observability

Observability is an important part of any real-world computing system (see page 35). It is an important part of any platform that strives to be offered as a service to large groups of different users.

At every moment, we must know what is happening in the system, where bottlenecks are, and how we can resolve problems. This is of significant importance in DS, where state and calculations are scattered across multiple nodes, clusters and even DCs.

In the  $\mu$ Cs environment that involves multiple layers (see page 77), this technique is even more important to give us insight what is happening in the system, and how we can optimize different layers of the system. On the other hand, users who use such a system need to know what is happening with their services.

Because of these two really important parameters, observability must be implemented on two levels:

- (1) **System level** contains data that is generated by the system. This data should be available to the administrators and providers of the system **only**. Operations people in the team (eg. DevOps or SREs), and developers cannot see it;



- (2) **User level** contains information about user requests that **only** operations people in the team (eg. DevOps or SREs) should be able see it. This type of data should not be visible to the system providers for privacy reasons.

Every service should log details about its usage and calls as well as traces how requests are going, regardless of whether it is a platform or user. Log data should be stored outside the service and the virtualization tool (eg. using sidecar pattern with containers [179]), while pieces of information will be sent to centralized log storage. Sending intervals should be able to be changed and adapted for both levels.

Log storage could be searched to see the general state of the system, and pieces of information about user requests and the state of their requests.

### 4.9 Access pattern

In section 4.3, we already discussed system access patterns from the applications point of view, using streams and topics. In this section, we are going to venture into the dashboard and system access described in 2.3.

To access the system, requests are going over the *master process*. This process is not responsible for any sort of synchronization, agreement, or something like that, but just to show dashboards and various system details since the cloud is available anywhere in the world, while  $\mu$ C should serve the local population.

One question that could that require an answer is, what if some cloud provider offering  $\mu$ C functionality, or running our master process, goes down, should the rest of the  $\mu$ Cs go to *read-only* mode and only accept **read** requests?

All communication is not exclusive to go over that master process, but if the user is nearby to  $\mu$ C he/she should be able to initiate commands directly to clusters, regions, and topologies. However, for

dashboards and full pieces of information about his  $\mu$ Cs, cloud would be a better solution, because of more available resources.

If some cloud provider goes down for whatever reason, we should foresee this, and to resolve it, we could use multiple cloud providers using *Multi-Cloud Computing* [59, 60] so that one cloud provider is a master process and many others are backup in case the whole cloud provider goes down.

We inevitably have some state synchronization here, and we could rely on SEC and CRDTs to do synchronization without some expensive coordination between providers.

This is not generally a problem if applications that are running in  $\mu$ Cs are not dependent on some process, location, sensors or services, etc. If this is the case, then we can connect  $\mu$ Cs in sort of *P2P network* and give them some logic to just route request to the proper location on the globe.

This might not be that fast, since even light is affected by the distance, but we would be able to issue requests to cluster in a different part of the world.

## 4.10 Auto scaling micro clouds

Configurable model structure bring any additional benefit to proposed  $\mu$ Cs model — infrastructure **auto scaling**.

The purpose of this process is to automatically scale up or down existing clusters, regions even topologies, adapting the number of nodes in the cluster, clusters per region, and even regions per topology to match newly workload increase.

The users can create a scaling plan via scaling policies, where they should be able to define: **(1)** how resources should be added (e.g., grow clusters, regions, or topologies) and at which rate, **(2)** what specific type of resources are desired (e.g., preferred values for CPU, disk, memory etc.), and **(3)** boundaries for both resource additions and resource types.

The third item is increasingly important, because it directly affect the costs of the applications running inside  $\mu$ Cs. As such, it should never be overlooked.

Based on the given policies the autoscaler dynamically adjusts the amount of resources according to user demands [180]. This elasticity allows accommodation of fluctuations in the workloads, that vary over time.

Auto scaling strategy could be applied to both applications and infrastructure, and  $\mu$ C provider should decide should they depend on each other, or should they be regulated separately. The proposed model could accomodate both options.

### 4.11 User data flow in micro clouds

One of the benefits of having  $\mu$ C in near proximity is that user data can be processed and stored locally, and users are able to control to which extent data can be shared with the cloud [20] (see section 1.2).

User data would be relatively easy to replicate inside a region using CRDTs for example (see 2.1.1), to eliminate expensive coordination [166]. On top of this, distributed file system can be built, distributed databases, processing frameworks specially designed for  $\mu$ C environment.

But since  $\mu$ Cs storage capacity is less than traditional cloud, they should store most recently used data [162]. The storage capacity is determined by the  $\mu$ C size given by the equation 4.1.

**Data locality** in  $\mu$ Cs creates two challenges, that must be addressed:

- (1) If a user leaves a place, is data stationary, or follows him using some pattern? The decision should be left to users,  $\mu$ C providers, and developers to decide. The proposed model allows the creation of different **data plans** that could be represented via **different policies** for every individual user or group of users. If data

needs to follow the user, as a data transfer medium the cloud could be used as a backbone, to transfer data to the  $\mu\text{C}$  in user proximity. The process is similar to the content delivery networks on the edge [27]. In such a scenario, we must pay attention to the network pressure and transfer only the requested data. As an alternative option, we can serve data directly from the cloud and do not transfer anything;

- (2) How long the user data should be stored in the  $\mu\text{C}$ , if we assume that locally existing  $\mu\text{C}$  will serve many users? The proposed model does not restrict how long the data should be stored and it is **policy-based**. Depending on the  $\mu\text{C}$  size given by equation 4.1,  $\mu\text{C}$  providers may determine **different policies** — Time to live (TTL) [181] and offer to users. This process is somewhat similar to the leases mechanism in cache systems [182].

The data retention and size in the  $\mu\text{C}$  are decided and optimized with 4.1 and 4.2, for the single user or group of users. This gives developers, and  $\mu\text{C}$  providers more room to optimize their policies offered to the users, while users may choose different options based on their needs [162].

## 4.12 Extendability

The proposed system allow extension, and could be used for addition of various configurations and artifacts to *nodes*, *clusters*, *regions* and even *topologies* —  $\mu\text{Cs}$ , all done remotely.

As the direct implication of *reconcile pattern* (see page 54), extension is as easy as adding the new worker to the system. This new worker, will deal with a specific resource type when that kind of file is submitted to the system.

On the other hand, a *reconcile loop* (see page 54) will provide the observed resource informations that this worker can act on. This mean

that specific worker (controller) will deal with just one resource type at any given point in time knowing what is specified, and what is returned from reconcile loop.

When registering new worker, a user should provide *unique* name of the worker so that reconcile loop can send informations to the right place. As a unique name we can use for example worker name and resource type name.

Figure 4.12 show high overview of adding a new worker to the system, with reconcile loop.

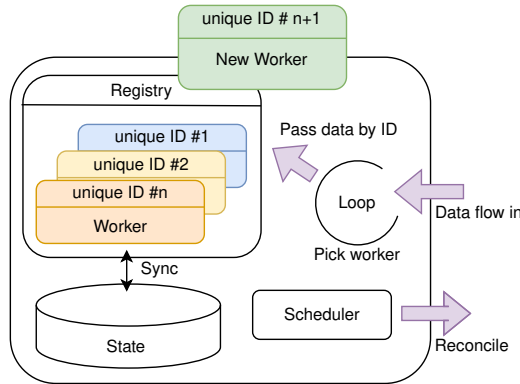


Figure 4.12: New worker addition to the system, with reconcile loop.

### 4.13 Repercussion

The model presented in this chapter, has four possible model repercussions:

- (1) **Stand alone**, the proposed model can serve as a base layer for future ECC as a service implementation. On top of it, we can implement other services and features like scheduling, storage, applications, management, monitoring, etc. As such it could be a viable option in the CC;

- (2) **Integration**, the proposed model could be integrated with existing systems like Kubernetes, OpenShift, or cloud provider infrastructure since they all operate over the cluster. This is possible, with some small infrastructure changes and adaptations because – the communication should be implemented via standard interfaces like HTTP and JSON, the integrations should be relatively easy to achieve. The proposed model could be used as a geo-distributed description and/or an organization tool;
- (3) **Combination**, this approach can be done over multi-cloud principles. Some cloud tasks could be offloaded to the nearest  $\mu C$ ;
- (4) **Enhancement**, big data tools, especially lambda architecture (see Section 2.2.1) could be enhanced using the proposed model creating **lambda++** architecture. Here we can reduce cost even more, by doing preprocessing and filtering of data closer to their source. After that we can send processed data to stream processing part of traditional lambda architecture, allowing it to process even more data than before.

# Chapter 5

## Proof of concept

In this chapter, more details are going to be given about the framework implemented based on the formal model and architecture specification given in the previous chapter.

Section 5.1 discusses framework architecture, and system implementation details, and framework limits. Section 5.2 presents implementation details about framework operations. Section 6.1 describes few possible scenarios and applications that could utilize  $\mu$ Cs platform. In Section 5.3 presents results of our experiments.

### 5.1 Platform implementation

In this section, we are going to introduce an implemented proof-of-concept framework based on the model proposed in the previous chapter 4. The framework is called **Constellations (c12s)**<sup>1</sup> because it is strongly influenced by nature and the neverending number of galaxies that the universe is (not only) composed of. Similarly, we are trying to create a universe of clusters that will serve users to help them with

---

<sup>1</sup><https://github.com/c12s>

their day-to-day tasks. The framework is **open-source**, and it is implemented using the microservice architecture with services that have distinct role and purpose to the entire system. These services are:

- **Gateway**, its purpose is to export services feature to the rest of the world. Gateway is designed as a REST service, accepting *JSON* style messages, so that various clients can communicate to the system. When the request arrives at the gateway, if the request is valid it will pass the request to the rest of the system. It communicates to the rest of the services to check if the user exists, if he/she has proper rights for actions he/she sends, and if not, returns the proper message and does not propagate it to the rest of the system;
- **Authentication & Authorization**, the sole purpose of this service is to store users and their credentials. This service will validate does user exists in the system, and does he have certain rights to perform some specific operation. Users that often use the system will be stored in the cache layer of the service so that on the next request his/her actions are done faster. Users that do not use the system that often will not be stored in the cache until the first use. After that, if the user does not use the system for some time, he/she will expire from the cache;
- **Queues**, the purpose of this service is to prevent huge request load to the system and to accept more user requests. When the user submits any **mutation** operation – an operation that changes the state of the system, these operations will be put in the queue. User can create their queues, to prevent long lines for specific tasks. For example, users can create queues for specific tasks, and use them only for those tasks, while other queues could be general-purpose queues. On system start, every user will start with one queue — **default**. When doing mutations on different parts of the system, the user can specify in metadata which queue he wants the task to go to. This service implements



a token bucket rate-limiting algorithm [183] to prevent congestion of the system;

- **Nodes**, this service stores and maintains pieces of information about registered nodes in the system. All node hardware and software details will be stored in this here. This service is also responsible for storing metrics data, accept health-check requests from nodes, and inform the rest of the system that the used node is alive;
- **State**, is the heart of the system. This service stores all information about architecture, clusters, regions, and topologies. When a new cluster/region/topology is created, this service will setup watchers for nodes, so that if the node does not send the health-check signal for some time, that node will be declared dead. This is important so that at any moment we must know the state of the clusters and their utilization. This service as well will cache frequently used nodes data, so that on the next request node lookup is faster, since we can have a huge number of nodes, topologies, clusters and pieces of information about them. To prevent data loss, this service will first store a copy of the operation before attempting any mutation of the system;
- **Log**, is responsible for storing all log and trace data from every service. Here a user can check whether all jobs are done, whether there is some error and possibly why the error happened to resolve it or fix it for the next time. From a user point of view, this is **read only** service and from a system view, this is **write only** service;
- **Command push**, the purpose of this service is to push commands and operations to the nodes user desired. This service implements a token bucket rate-limiting algorithm [183] to prevent constant data push to the nodes. Similar to the state service,

this service will store a copy of operation information locally before attempting any push to the nodes. This information will be deleted, once the operation reaches all decided nodes and all of the responses with the acknowledge (ACK) message;

All services are highly customizable and all have their configuration file that could be changed and adopted. All these services are easy to extend to support new operations and pieces of information about nodes, regions, clusters, topologies, and latter on applications, configurations, namespaces.

The system operates with four commands, where three of them follow formal models described in the previous chapter. The last command is a simple command to list logs for every user.

Regarding 4.9, the framework operates as a master process is running in Kubernetes, and through that process, all commands are issued. Future applications can communicate with nodes as stated in 4.3. They do not require communication with the master process at all, they can communicate with the formed cluster using topics and streams.

In the rest of this section, we will see details about all operations, as well as used technologies to implement the whole system. Possible applications and future directions for application development will also be described. System architecture is shown in Figure 5.1.

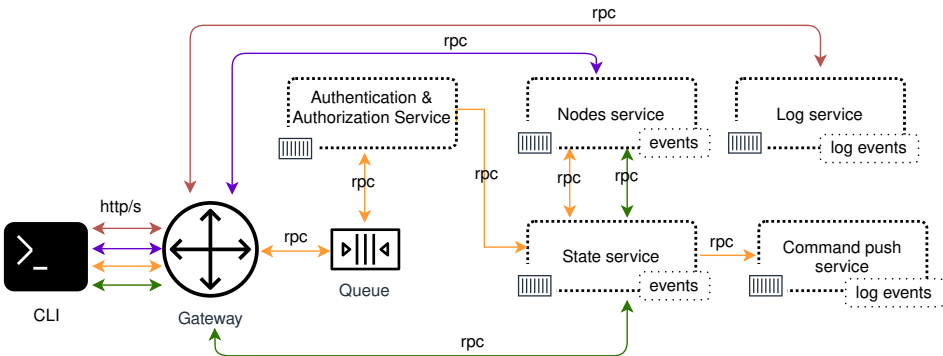


Figure 5.1: Proof of concept implemented system.

### 5.1.1 Technologies

All services are implemented using the Go<sup>2</sup> programming language, because of its well-known tooling, support for developing system software, web-based applications, small binaries, but also, great concurrency model, and ability to build binaries for almost any architecture without any code changes. All services rely on Go implicit *interface* implementation mechanism. Every technology or component used in the system can be swapped for some other, as long as that component implements *interface* fully. The framework is developed in such a way that is relatively easy to extend, or switch and use different components and technologies.

As the main storage layer for our system, we used etcd<sup>3</sup>, a popular open-source key-value database, that shows good performance, and it is mostly used for configuration data. Metrics data are stored in the open-source time-series database InfluxDB<sup>4</sup>.

Communication between microservices is implemented in RPC manner using gRPC<sup>5</sup>, and Protobuf<sup>6</sup> as a message definition. gRPC and Protobuf are open-source tools designed by Google to be scalable, interoperable, and available for general purposes. Communication between nodes and the system is carried out using NATS, an open-source messaging system. Health checking and action push to nodes are implemented over NATS<sup>7</sup> in a publish-subscribe manner.

Caching layer for every service is implemented using Redis<sup>8</sup> key-value, in-memory database. It is important to notice, that all concrete

---

<sup>2</sup><https://golang.org/>

<sup>3</sup><https://etcd.io/>

<sup>4</sup><https://www.influxdata.com/>

<sup>5</sup><https://grpc.io/>

<sup>6</sup><https://developers.google.com/protocol-buffers>

<sup>7</sup><https://nats.io/>

<sup>8</sup><https://redis.io/>

tools that are used, are easily swappable for some other as long as they implement a proper interface.

All communication with the outside world is done in a REST manner using JSON encoded messages over HTTP. To communicate with the platform, we have developed a simple command-line interface (CLI) application also using the Go programming language that sends JSON encoded messages over HTTP to the system.

Every service is packed in a container, and for this purpose Docker<sup>9</sup> containers are used. To achieve automatic orchestration, and self-heal, and up-time, all services that are packed in containers, are running inside Kubernetes.

Every service will log details about its usage and calls, as well as traces how requests are going. Log data is stored outside the service and container, and pieces of information will be sent to centralized log storage on every  $t$  seconds specified by the user. Sending intervals could be changed and adapted using the configuration file for every service independently.

Log data will be stored in the two levels:

- (1) **System level**, this data is generated by the system and could be viewed by administrators of the system **only**. Operations people in the team (eg. DevOps or SREs), and developers cannot see it, but providers can;
- (2) **User level** that stores information about user requests that **only** users can see. This type of data will not be visible to the developers of the system, and only users that created these logs will be able to see them.

Log storage could be searched to see the general state of the system, and pieces of information about user requests and the state of their requests.

---

<sup>9</sup><https://www.docker.com/>

### 5.1.2 Node daemon

Every node runs a simple daemon program implemented using the Go programming language, as an actor system (see Section 2.11.1). The actor system is developed for this purpose. When a message arrives, the proper actor will react based on the message type, or discard it if the type is not supported.

Extending such a system is rather easy because users need to simply write a new actor and logic that goes with them and register it to the system.

When the daemon starts, it will first read the configuration file to do the proper setup and then will contact the actor system to start all the actors.

System messages will be sent to the daemon, but it will not react to these messages. Daemon is not able to communicate with any actor directly. All communication goes through the actor system which is responsible to pass messages to the actors. The actor system at this point has only four existing actors:

- (1) **Cluster actor**, this actor reacts on messages about new cluster formation, but he is also responsible to contact rest of the system that message has arrived.
- (2) **Internal actor**, this actor reacts to messages from other actors to update the daemon state. For example, on new cluster creation, this actor will update daemon id, name, labels, etc.
- (3) **Health actor**, this actor will periodically send *health-check* data to the system about node state, utilization, etc. This actor will communicate to the rest of the hardware to get proper pieces of information, to collect logs from the node, and send all that data to the system.
- (4) **Gossip actor**, this actor will communicate with other peers in the group using SWIM protocol techniques.

The actor system will monitor these actors, and in case any of the crashes, the actor system will restart them.

Listing 5.1 show the actor system hierarchy of existing actors.

```

1  \StarSystem
2    +--\TopologyActor
3    +--\HealthcheckActor
4    +--\GossipActor
5    +--\InternalActor

```

Listing 5.1: Actor system hierarchy.

Before daemon starts, the user needs to specify some parameters for proper configuration like:

- (1) **Identifier** represents unique identifier of the node. When a node is not a part of some cluster, this can be whatever the user wants. Once the node is a part of some cluster, the identifier will be updated, and it is not advised to change it manually afterwards.;
- (2) **Name**, represent name of the node. This property also can be changed when a node is part of the cluster, otherwise, it can be whatever we want;
- (3) **Labels** represent the specific features of the node. Labels are used to query for free nodes, and there is no formal restriction of what they can or can not be. This property can be changed when the node is a part of some cluster. It is advised that as labels we put some specific features of the node that might be beneficial for the user who is looking for nodes to create new cluster(s);
- (4) **Health-check details**, here we have pieces of information to control the health-check mechanism. Since nodes communicate with the rest of the system via publish-subscribe events, we must specify the address of the rest of the system and the topic name, where we publish our pieces of information;

- (5) **System information** represents basic information for a node to know how to contact the rest of the system. We should specify the system address, so that node knows where to send messages, and where are messages are coming from. We can also specify the version of the system we are trying to contact. System version could be used to support **backward compatibility** if we have multiple API versions of the system running at the same time or some period.

Configuration can be done easily using the YAML configuration file.

Listing 5.2 shows simple YAML configuration file for daemon process.

```
1  star:
2    version: v1
3    nodeid: node1
4    name: noname
5    flusher: address
6    healthcheck:
7      address: address
8      topic: health
9      interval: 1m
10   labels:
11     os: linux
12     disk: flash
13     arch: arm
14     model: rpi
15     memory: 4GB
16     storage: 120GB
17     cpu: 1
18     cores: 4
```

Listing 5.2: Daemon configuration file

Based on the configuration file, the daemon will start a background health-check mechanism, and it will subscribe to the system, using an

identifier as a subscription topic. The background thread will contact the system repeatedly using a contact interval time, specified in the configuration file.

On every health-check, the node will send labels, names, IDs, and metrics to the system (e.g., CPU utilization, total, used, free ram or disk, etc.). The specified labels will be used when the user is querying for available nodes, while the node id will be used for reservation when forming a cluster.

At the first start of the daemon, when the node is free, the user can specify whatever node id he/she wants. Once, the node is a part of the cluster, the node id will be updated to match that. Node id update will happen on the cluster formation process.

### 5.1.3 Separation of concerns details

The implemented framework follow the clear SoC model, presented in Section 4.2. Since the presented model consists of three components, the implemented framework is deals only with **resources** (see page 78). Its job is to organize nodes into clusters, regions, and topologies, to make them communicate and expose their resources to the upper layer of SoC for utilization. The upper layer will run services on these resources, to collect data from data creators and process them as requested. It must have set up the infrastructure to do any processing or storage. This middle layer is the binding element between the layers.

### 5.1.4 Long-lived transactions details

Current implementation follows the sagas pattern (see page 44) for transactions, and isolation is added using semantic lock by adding states to the task. The task can be in one of following states: (1) *PENDING*, (2) *IN PROGRESS*, (3) *CREATED*, and (4) *FAILED*.

Until the task reaches *Created* state, no other operation can be done on that cluster configuration.



The cluster formation transaction is split into a few sub-transactions: (1) reserving nodes, (2) writing cluster information, and (3) sending cluster information to push service and ultimately sending pieces of information to the nodes.

The rollback mechanism is implemented using *choreography*, and communication between the various services is done asynchronously using message queues.

### 5.1.5 Garbage collection details

The current garbage collection process is pretty trivial since there is no complicated graph of items connected to the cluster information details.

The process will scan tasks with the state *FAILED* in the background, and it will remove those items. Related details to the cluster information, for example, node metrics, will expire automatically since there is a lease attached to them.

Another resource that will be automatically deleted is node data. If a node health-check message does not reach the system in some time, the node lease will expire and data will be removed automatically.

### 5.1.6 Limitations

The framework at the current state has some limitations that we **must address**. As shown in figure 5.1, for purpose of testing the system, and to give the users any way to interact with the system, a CLI application is implemented. This is a good option for initiating commands to the system. The problem with this approach is that showing logs, and topologies might be limiting the users, especially if they monitor and supervise multiple topologies with regions and clusters.

For monitoring, a tool with UI would be a better solution and it can show more details. For a small amount of data, when topologies are relatively small, CLI could be used. But for some real-world applications, desktop, and/or web-based UI will be a much better

solution, and CLI can be used for fast lookup on specific pieces of information about clusters and nodes, for example.

The current implementation of the queueing system is **limiting** because if we want to extend the system with new queues we need to shut down the whole service, and that is not the best solution. But since the goal of this thesis is not queue management, but  $\mu$ C formation, protocols, and formal modeling.

Since the goal of this thesis is not queue management, and the purpose of this thesis is not Human-computer interaction, but  $\mu$ C formation, protocols, and formal modeling, mentioned limitations should be the topic of the future work [7.3](#) section.

The current implementation do not include garbage collection of connected dependent items, and users cannot influence garbage collection decisions. This should be one of the topic of the future work [7.3](#) section as well.

## 5.2 Operations

In this section, we are going to describe all implemented operations in the framework and present specific details about every individual operation.

### 5.2.1 Query

The query operation is used to show all free nodes or filtered free nodes that are registered into the system, to the user (yellow arrows in [Figure 5.1](#)).

When a user wants to get information about free nodes, they need to submit a **selector** value which is composed of multiple **key-value** pairs. These key-value pairs can be any alphanumeric set of symbols for both keys and values. Based on that key-value pair the system will do the query of the free nodes.

The selector will be used as a search mechanism to compare the labels of every free node that exists in the system. The nodes that are satisfying the rules 4.11 defined in Section 4.6.3 will be present in the result.

This operation is done before the formation of new topologies, regions, or clusters — mutation of the system. The user first needs to get a list of free nodes, then he can choose nodes that are best suited for him and try to form topologies, regions, or clusters of them.

The querying process is a little bit changed from one presented in Section 4.6.3. The only change that is made is the addition of the *Gateway* service that will pass requests into the system and prevent overflow of requests. This change **does not** affect or validate the formal model presented before, since the added service does not interfere with the process of searching nodes or changes to the system.

Figure 5.2 shows a communication diagram for the query action, with the addition of the Gateway service.

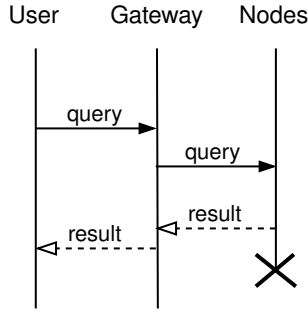


Figure 5.2: Low level communication protocol diagram of query operation.

Query operation is implemented using *api composition* pattern (see page 34).

### 5.2.2 Mutate

Mutate operation (orange arrows in Figure 5.1.) changes the system state by creating, editing, or deleting clusters, regions, and topologies. When a user wants to perform a mutation over the system, the desired state needs to be specified **declaratively** using a *YAML* file and submitted to the system. When the state is submitted, the system will do the rest of the job to ensure that the state desired by the user is reached. As a deployment strategy, we used *recreation* strategy (cf. section 4.5.1), for simplistic implementation reason.

The users specify which nodes are forming the cluster. Optionally, users can also specify labels and names on the node level, and retention period on the cluster level. The retention period is used to describe how long metrics are going to be kept. When the retention period expires, the metrics data will be deleted or moved to another location.

When forming a topology, users can assign a name and label to the entire topology. These parameters will be used when the user wants to query all topologies to get full information about regions, clusters, and nodes inside a topology.

Mutate operation is **immutable**, which mean that there will be no in-place changes to the existing state. If a user wants to do any update, he/she needs to specify a full new state that will replace the existing one. This operation is **atomic** as well, which mean that a whole new state must be able to replace the existing state. If this happens, the system will replace the old state with the new one. The main storage that stores configuration data is a key-value store implemented using an etcd database. The key that will be used to store the configuration data is the whole path of topology, regions, cluster, node, while the value represents the state desired by the user.

Listing 5.3 shows structure for key-value pair that is stored in the main database, where on top we can see a structure of the **key**, and below it we can see the structure of the **value**. This kind of structure is similar to OS file-system data organizations of files and folders.

```
1 \topology\region\cluster\node
```

```

2    +--\labels
3    +--\informations
4    +--\resources
5    +--\status

```

Listing 5.3: Structure of stored key-value element.

We store as well one additional information about labels **index** value. This index is used for a faster query of elements when doing labels comparison. To find elements we can query in a similar way like searching files and folder structure. The etcd in newer version **do not allow** hierarchical keyspace, but what they allow is **ranged** query by some **prefix**. This is useful as well because we can still get all regions in topology, clusters in the region, or nodes in the cluster if we know to which group they belong.

Mutate operation is not idempotent by nature, but the whole process behind it makes mutate an idempotent operation. Whether the user tries to create already existing infrastructure by changing the order of regions, clusters, or labels in the nodes, or if he for whatever reason did not get confirmation that infrastructure is created, an existing infrastructure will not be created. This is done in such a way, that *State* service (Figure 5.1) will keep the information set about created infrastructure. This information is implemented as a **flat keyspace** set, that have information. On every mutation request, the system will test if such a topology already exists. If such topology already exists, the user will get confirmation that his infrastructure is created. If such topology do not exist, a new cluster formation protocol 4.6.3 will be initiated. The mutation confirmation is followed by a **unique id**, which the user can use to query steps, logs, and traces that are done in process of working towards the state desired by the user. Logs service can give the user complete details about his task state using that unique id.

When creating topologies, the user is free to give whatever name he wants for every region, cluster, and node. The only restriction is that name should be the alphanumeric set of characters. Listing 5.4

shows an example of a user-defined state that forms the topology of one region with one cluster that contains three nodes, with a retention period of 24 hours. The whole topology will have the same set of labels, but *node3* redefines this rule by specifying its own.

```
1  constellations:
2    version: v1
3    kind: Topology
4    metadata:
5      taskName: default
6      queue: default
7    payload:
8      name: MyTopology
9      selector:
10       labels:
11         l1: v1
12         l2: v2
13         l3: v3
14     topology:
15       region1:
16         cluster1:
17           retention:
18             period: 24h
19           machineid1:
20             name: node1
21           machineid2:
22             name: node2
23           machineid3:
24             name: node3
25           selector:
26             labels:
27               os: linux
28               arch: arm
29               model: rpi
```

```
30                                storage: 100GB
31                                memory: 1GB
```

Listing 5.4: Example of mutate file using YAML.

After all, nodes that should form a cluster, acknowledge cluster formation message, they will inform the system that the message is received, and they will start the process of cluster formation. This process is done by using SWIM [32], a Gossip style protocol. When every node has a complete list of its peers that should be in the cluster, the process of cluster formation is done.

On the next health-check message, every node will send its **id** that is telling the system that he is part of some cluster. This kind of messages will be used in the *State* service to validate that cluster is alive and well, or that some nodes (or all), are dead or down if **id** is not received.

Gossip style protocols (like SWIM) could be used in the future to propagate information in the cluster, without explicitly ping every node in the cluster.

### 5.2.3 Queueing

When doing mutation, users can target a specific system queue, by adding a metadata part in the configuration file. With this ability, users can aim for specific queues just for the mutation to avoid long waiting times if other queues are full. Currently, the system does not have any limitations, restrictions, or logic that will specify which queues are used for what or give them special rules or permissions. This can be viewed as a “gentleman agreement”, that in one team, operations users can proclaim specific queues like *mutation* queues used maybe for specific topologies, and later on for other operations as well.

The queue service starts only when the *default* queue exists. Adding a new queue to the system is implemented using the configuration file, for convenience only.

Listing 5.5 shows an example of queue service with two additional queues with specifications of their parameter needed for token bucket operation [183]. Also, we can see configuration pieces of information for instrumentation of a single service, and the same configuration is implemented for all specified services shown in Figure 5.1.

```
1  blackhole:
2    db: address
3    queue:
4      myqueue1:
5        namespace: mynamespace
6        retry:
7          delay: linear
8          doubling: 1
9          limit: 5000
10       maxWorkers: 4
11       capacity: 4
12       tokens: 0
13       fillInterval:
14         interval: 1
15         rate: s
16     myqueue2:
17       namespace: default
18       retry:
19         delay: exp
20         doubling: 2
21         limit: 50000
22       maxWorkers: 5
23       capacity: 5
24       tokens: 0
25       fillInterval:
26         interval: 1
27         rate: s
28   instrument:
```



```
29      address: address
30      stopic: topic
31      collect: interval
32      location: location
```

Listing 5.5: Structure of stored key-value element.

### 5.2.4 List

The list command shows the current state of the system for the logged user (blue arrows in Figure 5.1.). Logged user is able **only** to see infrastructure he/she has created or maintains. All other infrastructures, created by other users, will not be visible to the users that are not creators or maintainers.

To view his/her infrastructures, the user can specify what part of the system he/she wants to see using a set of labels or list of key-value pairs, which will be used by the system to determine what the user wants to see. This process is similar to **selector** shown in the query 5.2.1 operation.

There are two levels of details that user can specify:

- (1) **Global view** of the syste, or all topologies he/she manages with just basic information like resource utilization, number of regions clusters and nodes;
- (2) **Detail view**, or details about a single topology (i.e., regions, clusters, and nodes in a single topology). Users can specify a more detailed view of a single cluster, meaning the users will get information about what resources the cluster has, but also resource utilization over time (using stored metrics information) and so on.

Both options can be useful if operations people need different details levels for different topologies. List operation is implemented using *api composition* pattern (cf. page 34).

### 5.2.5 Logs

The logs operation shows stored logs and traces to the user (purple arrows in Figure 5.1.). Same as previous operations, the user needs to be registered and logged in to be able to perform this action. With this action, the user can see the state of his/her operations and actions. The user can choose between two options for querying his/her logs:

- (1) **Get**, for this option a user needs to provide a unique task id that is given to the user when he/she creates a mutate operation. With this option, the user will get a full list of steps, traces, and logs collected over the time the system was setting up his/her desired state;
- (2) **List**, with this option user can specify **selector** using list of key-value pairs in a similar way to query 5.2.1 and mutate 5.2.2 to filter only parts of the traces that contain the same labels as selector does.

This action is implemented in some basic aspects, as this action can return a huge amount of data that require some better visualization than CLI.

## 5.3 Experimental Results

In our tests, we have used:

- 9 Raspberry Pi 3+ Model B with 1GB LPDDR2 SDRAM, 16GB SDCard storage, and 1.4GHz Cortex-A53 64-bit SoC running Raspbian Linux, a Debian-based operating system.
- 3 Beagle bone black devices with 512MB DDR3 RAM, 4GB 8-bit eMMC on-board flash storage, and 1GHz ARM Cortex-A8 running a version of Linux Debian operating system.

as test nodes for creating clusters, regions, and topologies.

### 5.3.1 Experiment

Using Go tooling, we were able to build daemon service without changes and disseminate on all nodes without problems.

We ran tests on different configurations and different nodes clusters using these nodes. All nodes were connected on the network, and experiments were conducted in a controlled environment. Nodes that should be a part of the same cluster were connected on same network.

Our experiment started with separating nodes into groups of **three**. This number is chosen because in DS, an odd number is used in cases when we need to reach some agreement and we need major majority like consensus, for example. This is not important for purpose of this thesis, we could pick any number, membership protocol does not makes a difference if there are three or four or eleven nodes in the cluster.

After nodes had been separated, we created a configuration file for every node, setting up default parameters for every property node daemon required 5.1.2. After all services were up and running, we turned the nodes on, and health-check protocol 4.6.2 started at upfront defined time, which we had set for test purposes at *1 minute* interval. Logs, resources and other node details were set to *15 seconds* interval, so between health-check intervals, daemon would collect resource information *four times* before sending it to the rest of the system.

For convenient testing, all nodes had the same set of labels, and as labels we chose OS name, OS version, architecture version, node name basic details about resources of the nodes.

After some time, we were able to see all nodes registered in the system. When all nodes had been sending health-check ping for some time, and we had a stable system, we issued a mutate operation creating clusters of nodes that are logically close ot each other, and we initiated cluster formation protocol 4.6.3. After the protocol was done, we ended up with four clusters as we intended. We tried to initiate the same command again to test idempotency check 4.6.4, and we got the message that clusters already existed.

To increase capacity, we extended clusters by creating new ones using *three clusters* with *four nodes*. We created new mutate file, and initiate new mutate command. After some time, we saw that one cluster was down and that we now had *three clusters* with *four-nodes*, as we intended. After successful creation of new clusters, we deleted down cluster.

The last test was to test health-check protocol once again - we disconnected one random node from the power supply, and since that node ping was missing, the system was able to detect which node was *down*. This concluded our experiment.

## 5.4 The existing solutions enhancement

The protocols defined in this thesis could serve as a base layer for the system developed from scratch. On top of such a solution, other services and features can be added like scheduling, storage, applications, management, monitoring, etc.

The protocols described in this thesis ensure proper node registration into the system, organization, and reorganization of node resources into clusters, regions, and topologies, bringing disposable  $\mu$ Cs closer to the users at the network edge, allowing that requests are served from the local  $\mu$ Cs.

The existing orchestrator engines (e.g., Kubernetes, Apache Mesos, Docker Swarm, etc.) operate one cluster level [91, 149, 150, 154, 184]. The single cluster could span over multiple availability zones in the cloud, minimizing the chance that a failure in one zone impairs services in other zones [184]. Kubernetes allow extension in terms of multi-cluster deployments [184]. In such a scenario, Kubernetes is handling these clusters as disposable — "treating **clusters** as cattle, not pets" (i.e., numerous servers/clusters built using automated tools designed for failure, where no servers/clusters are irreplaceable [185]).

The model proposed in this thesis goes one step further, allowing the creation of disposable  $\mu$ Cs. Such an extension allows infrastructure

optimization in more dimensions [186]. The users are allowed to build numerous  $\mu$ Cs designed for failure using automated tools where no  $\mu$ C is irreplaceable — "treating  $\mu$ Cs as cattle, not pets."

The existing solutions can integrate the model proposed in this thesis to serve as a node register. Such integration allows the registration of new nodes into the system, allowing the existing orchestration engines to use new nodes, and expand their available resource pool. In form of specification, the users can provide a list of which available EC nodes need to be part of some  $\mu$ C.

The system will communicate with the existing orchestrator agent to register/unregister them with the existing cluster. In such a scenario, we can hook on the existing orchestrator health-check mechanism informing our system that a node in some cluster is alive. Unused nodes still use the health-check protocol defined in this thesis informing the system that they are still available for utilization.

The proposed model preserves the node's topology allowing cloud providers and orchestrator engines to treat  $\mu$ Cs disposable, abstracting infrastructure to the level of software — infrastructure as software [131]. This approach benefits from the already available tools, principles, and techniques (e.g., reuse, testing, modeling, and evaluation). The already available tools can be used for the disposable  $\mu$ C infrastructure definitions.

The model developed in this thesis is not competing with the existing orchestrator tools. It is not orchestrating applications, but it is a free nodes register and  $\mu$ Cs infrastructure descriptor allowed to be offered as a stand-alone service bringing disposable  $\mu$ Cs model to the users.

The model allows integration with the existing orchestration tools, leveraging existing mechanisms and best practices. Cloud providers can create an operator for various orchestration engines. This allows them to offer dynamically created, disposable  $\mu$ Cs as a service to their users, using infrastructure as software principles, for already existing orchestration engines, that users are already familiar with. And if

cloud provider have their own orchestration engines, they can fine-tune integration between the model and their tools tightly integrating it into their their services.

# Chapter 6

## Model usability

This chapter discusses applicability of the proposed model. We are also going to give a case-study for the COVID-19 area traffic control example in city of Milan, Italy.

### 6.1 Applications

This research focuses on a platform with geo-distributed edge nodes that can be organized dynamically into the  $\mu$ DCs and regions based on the cloud model, but adapted for a different environment. This middle layer helps the power-hungry servers reduce traffic by serving the nearby population requests. Users are getting a new platform as a blank canvas, and there is no limitation in what applications they can develop. Integrating hardware and/or software, even more, connecting sensors and things around us and eventually an operating system that will be capable of running city/state infrastructure without human intervention. The system presented in this thesis is a stepping stone towards that idea.

The main advance of EC, when compared to the cloud-only approach, is the acceleration of the communication speed. The cloud

could bring huge latency, while EC originates from peer to peer systems [10], serving only local population needs, minimizing potentially huge round-trip time of the cloud. Furthermore, the presented model expands peer to peer systems into new directions and blends them with the cloud to allow novel human-centered applications.

If we imagine sensors are put on a specific group of users and/or places in the city/state and monitor them in real-time, we can process these streams of data directly close to where they are, where they are moving and going. We could monitor air pollution for example, and make decisions in real-time to suggest users not to walk in some area, especially if they have some known respiratory problem.

Geo-distributed nodes represent a great idea to do any kind of monitoring and processing, especially for natural phenomena and do alert as soon as probes, sensors, or other things detect even the slightest changes. Applications like self-driving cars, delivery drones, or power balancing in electric grids require real-time processing for proper decision making or any other application that future developers may develop. Content delivery networks, content sharing could be implemented to serve content to the users faster than going over the cloud, since  $\mu$ Cs should serve the local population that is nearby.

## 6.2 Area traffic control example

Let us consider a use case that can benefit from our model. In the context of the recent COVID-19 outbreak, we can examine the city of Milan, Italy. Into nine municipalities, numbered from 1 to 9 the city is split. Let us follow the natural subdivision having Milan *topology* where municipalities have one or more *regions*. Population density implemented applications and needs dictate the number of *clusters* per region serving the population nearby.

If city subdivision changes in the future, reorganization of *regions* and *clusters* is easy to be done dynamically, using *cluster formation protocol*. A prerequisite for this to be done: there must be EC nodes



deployed on the territory, and nodes are connected to the system using *health check protocol*.

During the COVID-19 outbreak, an increased amount of ambulance vehicles and medical personnel had to be routed to hospitals as fast as possible. Let us consider that the city is using some platform for supervised area traffic control. Utilizing the principles from the presented model, we can target ambulance vehicles, giving them a higher priority, compared to regular vehicles.

In such a scenario,  $\mu$ Cs can run three kinds of *frontend services*, specifically tailored for this application: **(1)** a service that follows the ambulance vehicles, **(2)** a service that will regulate the traffic light control, and **(3)** GPS routing service.

Suddenly increased number of ambulance vehicles causes an increased need for resources (e.g., decisions that require more processing power) at the *frontend service* regulating the traffic routing and light control.  $\mu$ C *clusters* allow resource rearrange, or even a dedicated *cluster* just for this purpose. If more resources are required, *regions* can be changed as well, and finally, the whole *topology* can be repurposed or merged with a city nearby to support an increased number of requests.

We can monitor patient health in real-time [187, 188, 189] with extremely short response time (below a millisecond), using internet-based application control [190] and transfer data to healthcare platform [191, 192]. On patient's arrival, health workers already have crucial information, while robotic systems can help in diagnosis and screening [193]. Such a platform offering telemedicine in cooperation with area traffic control increases patient survival chances. At the same time, reduce the hospital spending on unnecessary tests.

For the research purpose, the depersonalized data can be transferred to the *backend service* for deeper analysis, behavior modeling, etc. This should be regulated by a trustworthy body. The proposed model would be helpful for researchers [194], giving them valuable insight into the outbreak in real-time.

The same applications strategy could be reused by others or adjusted

to best suit their needs. Such a service may exist only during the outbreak. When the situation is under control, the service could be dismissed. If fewer resources are required, or the pandemic is over, we can rearrange resources according to different needs. We can reuse the strategy on the next pandemic spike.

Figure 6.1 depicts previously described example, through conceptual architecture model.

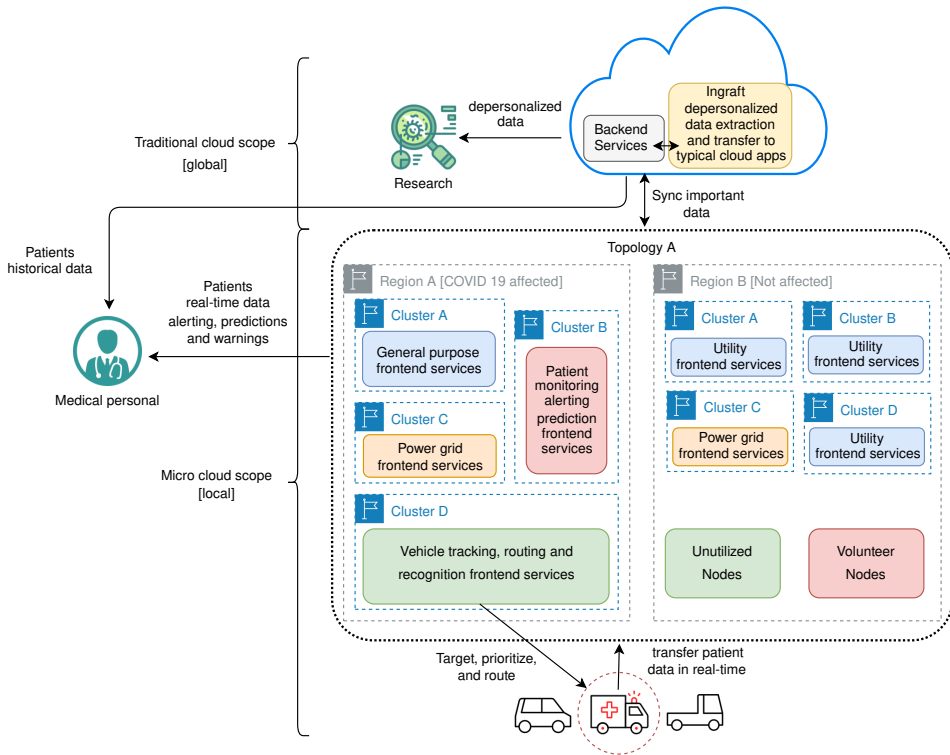


Figure 6.1: Conceptual architecture model for COVID-19 area traffic control example

# Chapter 7

## Conclusion

This chapter gives the summary of contributions for this thesis in Section 7.1. Section 7.2 presents limitations of our model, while Section 7.3 presents future work.

### 7.1 Contributions of the thesis

This thesis presents a possible solution to how to organize geo-distributed EC nodes into  $\mu$ Cs that will be able to serve the requests of the nearby population. We have introduced an extension of  $\mu$ DCs based on proven abstractions from cloud computing like zones and regions but adapted for different usage scenarios.

These easy to understand, yet powerful abstractions with slight adaptations, allowed us to cover any arbitrary vast geographic area, and yield a more available and reliable system forming the  $\mu$ C model. These abstractions are easy to organize and reorganize, and  $\mu$ C size is determined by the population needs descriptively.

We have also presented the cloud to ECC mapping, showing differences between two architecture models. Furthermore, we have given a formal model of the system and its protocols used to form such a system, with clear SoC and native application model for future  $\mu$ Cs

infrastructure and service development. The thesis also presents a proof of concept implementation and discusses integration into existing solutions, limitation of such a system with a few applications that could be used in.

Chapter 1 presents the motivation for this thesis with problem are, hypothesis combined with goals and research questions this thesis is built on.

A short introduction to the topic of distributed system is given in Chapter 2, with a short introduction to the topic of distributed systems, with a focus only on the areas that are important for the understanding of this thesis.

It is shown what distributed systems are. We present problems these systems create, and why they are so hard to implement and maintain. It also presents a few distributed computing applications, that we can use to employ nodes in the distributed system. Further on it is shown what scalability is and why it is important for distributed systems, with few organizational ideas like peer-to-peer and membership that protocol that is important in a distributed environment for various reasons.

Virtualization techniques are then described that can be used to pack and deploy applications and infrastructure, how to implement various deployment techniques especially in the CC environment, but also the difference between DS and few models that are usually confused as distributed like parallel computing and decentralized computing.

Chapter 3 shows related work done by various other researchers or, companies, focused again only on things that are related to this thesis.

We show different platform options, where people change the existing solutions like Kubernetes or OpenStack to made them work in areas like edge computing and mobile computing. Implementations of a few new platforms to use volunteer nodes to do some computation and storage on them with drop computing and systems like Nebula are shown further.

It is show how nodes can be organized to split some geographic area

into zones, and show how  $\mu$ DCs can help CC to serve requests from the local population. Different offloading techniques are used today, how to offload tasks from mobile devices closer to fog or edge nodes, but also various application models that could harness these offloading techniques and nodes organizations.

In the end, the position of this thesis, compared to other similar researches, is presented.

Chapter 4 presents the main contribution of this thesis. We dissect all important aspects that we need to have to help CC with latency issues, Big Data with huge volumes of data especially in the age of mobile devices and IoT.

Our model is based around  $\mu$ DCs that are zonally organized, that will serve local population and population nearby. We present a model that is based on CC but adapted for different scenarios and use cases.

We show how we can dynamically form new clusters, regions, and topologies and how we can use them in the new age of mobile devices and IoT. These newly formed systems will have clear the SoC, adopted from existing research to three-tier architecture. The formed model will serve as a pre-processing layer, firewall, or privacy layers, and it is adjustable in various dimensions.

The presented model can be formed in such a way to serve larger or smaller geographic areas, where users are located. This is a matter of agreement and a matter of choice. We present how developers can use this new infrastructure and what possible models of applications could exist, and how operations can deploy developed services onto existing infrastructure.

In the end, we present the repercussions of this model, and how it can be used as an integral part of existing systems to serve as topology storage or as a new model on top of which new subsystems and applications can be developed. We present protocols for the creation of such a system and model them using asynchronous session types. The system follows a formal model, and it is easy to extend.

Finally, we give limitations of such a system and things we must be aware of.

Chapter 5 presents an implemented framework that is based on knowledge compiled from previous chapters. We also present in detail operations that could be done in the framework, where it fits in the presented SoC model.

It further presents the results of our experiments in a controlled environment, what, the limitations of the framework at the current stage, and possible applications, and where this model could be used and beneficial.

The thesis is concluded in this last chapter with what is done, what can be done in the future in form of future work.

## 7.2 Limitations

The model proposed in this thesis has some limitations that we must be aware of. When talking about small-scale servers and  $\mu$ Cs, we must be aware of a few things.

- (1) Not all organizations will be able to deploy  $\mu$ Cs, due to the high initial investments required [12]. We can rely on government authorities, large cloud providers, or other big companies to build the initial infrastructure for their own needs, and lease it to others similar to the cloud. The general public can use them, similarly to the cloud – pay as you go, model;
- (2) There is no guarantee that existing public cloud providers will allow nodes that are not built, resigned, or deployed by them. If we are building a private cloud, then we can make a different decision. One way to resolve this issue is that the whole platform becomes open-source so that public cloud providers can engage in the development, and eventually use them as a solution;

- (3) These small-scale servers must be out of reach of people and protected in some way so that not everyone has access to them. Some degree of physical security must be implemented;
- (4) The places where these small scale servers will be deployed must have a stable internet connection, and the ability to integrate SDN or other similar technologies, so that complex network topologies could be implemented properly;
- (5) These servers can have some open architecture or could be custom built by other providers. In both cases, they must be able to satisfy rules that are presented in Section 4.2;
- (6) Splitting the processing into two parts and the possibility that users can be responsible for  $\mu$ Cs may raise some legal concerns. Either to develop interesting applications, use them as a firewall or simply use them as a privacy level for data, there must be a legal agreement that might not be that easy to achieve.

### 7.2.1 Discussion

The specialized infrastructure in models like [195, 26, 188, 187, 189] is required to solve a single problem. The model proposed in this thesis is more oriented towards a wider specter of applications, without the need for specialized hardware or software. Users should build applications, similarly as they build them for the cloud. Even existing applications models (e.g., microservices, serverless functions) could be transferred from cloud to  $\mu$ C.

One advantage is that specialized models are developed and optimized for a specific use case taking the maximum out of existing hardware and software. Compared to the presented model (in terms of speed), in some situations, they might outperform the proposed model. On the other hand, the proposed model offers much more freedom

for development (in terms of agility and applicability). This development freedom gives the users a new platform for creating interesting human-centered applications spanning over CC and  $\mu$ Cs [67].

Specialized platforms usually require special types of applications, while the proposed model does not limit users (in terms of development tools and techniques), as long as their application could be virtualized in some way (e.g., using virtual machines, containers, or unikernels). The developers may reuse existing knowledge to develop their applications.

The proposed model allows the organization of storage and processing resources according to priority if a catastrophic event occurs (e.g., COVID-19 see section 6.2). Providers can organize their resources and manage their digital infrastructure, where developers make applications that will help its citizens.

## 7.3 Future work

The work on the  $\mu$ Cs is at an early stage and leaves many open questions. As a part of our current and future work, we are planning to extend the proposed model in different directions. Future work might be separated into three options:

- (1) Features that operations people (eg. DevOps and SREs) users can benefit from;
- (2) Features developers might benefit from;
- (3) Infrastructure features, that both previous groups can benefit from.

For the first group, the first thing that should be implemented is remote cluster management, using configurations, security credentials, and actions over nodes in one or multiple clusters. On formed clusters, the users should be able to do remote configurations that nodes and/or



applications can use and set up the data without going from node to node.

The system should also be extended with namespaces for usage in environments with many users in multiple teams – multi-tenant environments. Namespaces provide the separation on virtual clusters, running on the same physical hardware. Speaking about multi-tenancy, we are also planning to implement role-based access control integrate with authentication and authorization services, with the addition of controlling different users with quotas, using rate limiting and resource limiting.

Also one possible approach would be to extend resources dynamically using infrastructure auto-scaling mechanism, so that system can grow and shrink as needed.

We are also planning to implement a full architecture and applications monitoring, alerting, and reporting that would be helpful to any administrator of such a system. We might also consider rethinking networking and making network isolation so that once formed topologies can communicate within themselves, and a possibility to specify strategies of communication with other topologies.

Queueing system mentioned in section [5.2.3](#) should be extended so that users and/or operations people can easily add new queues and possibly assign a role for them.

We should extend the access pattern so that users can issue commands to  $\mu$ Cs directly instead of going over the cloud master process only. Here we need to implement synchronization in multi-cloud deployment.

For the last part in this operations section, we should also think about continuous integration, deployment, and delivery of services onto the infrastructure, and as well as various UI dashboards that can be customized to present different aspects of the system.

Another direction for future work is the implementation which developers could benefit from. The first thing that should be implemented

here is the complete application framework so that users can start developing services that can do something. We should also implement a framework and maybe domain-specific language for the use case where users just want to pre-process the data before sending it to the cloud, in a more convenient way than writing the whole application.

Users can develop their applications with different models:

- (1) **mPaaS**, where the platform is doing all the management and offers a simple interface for developers to deploy their applications;
- (2) **mCaaS**, if users require more control over resources requirements, deployment and orchestration decisions;
- (3) **mSaaS**, users can develop their solutions only using  $\mu$ Cs, but this is not advised at the moment;

The second would be file system and database APIs that users can use to store their data. We should also provide an interface for extensions so that others can create their databases following different models from which developers can benefit, but also integrating existing ones.

The last part of the future work be extending the current system with tunable replicating strategies for the data, in case that any part of the topology fails for whatever reason, data would not be lost. Furthermore, we should provide tunable CC synchronization models that could be used.

We should implement a scheduling system for user-developed applications so that we can put applications into the formed architecture. And last but not least, we are planning to add several security layers to protect a system in general from malicious users.

We should investigate compression methods to reduce data stored and sent via the network. These tests should be conducted on ARM devices with existing methods, or maybe we can create a ground for new compression methods and techniques.

It is stated in section 4.13 of this thesis that this model could be integrated into existing solutions. Our efforts should go as well on

integrating this system with existing solutions so that they can benefit from this hierarchical and geo-distributed nodes organization in the same way or almost the same way as the stand-alone solution would.



# Bibliography

- [1] M. Chiang, T. Zhang, [Fog and iot: An overview of research opportunities](#), IEEE Internet Things J. 3 (6) (2016) 854–864. doi:[10.1109/JIOT.2016.2584538](#).  
URL <https://doi.org/10.1109/JIOT.2016.2584538>
- [2] H. S. Gunawi, M. Hao, R. O. Suminto, A. Laksono, A. D. Satria, J. Adityatama, K. J. Eliazar, [Why does the cloud stop computing? lessons from hundreds of service outages](#), in: M. K. Aguilera, B. Cooper, Y. Diao (Eds.), Proceedings of the Seventh ACM Symposium on Cloud Computing, Santa Clara, CA, USA, October 5-7, 2016, ACM, 2016, pp. 1–16. doi:[10.1145/2987550.2987583](#).  
URL <https://doi.org/10.1145/2987550.2987583>
- [3] W. Vogels, A head in the clouds the power of infrastructure as a service, in: Proceedings of the 1st Workshop on Cloud Computing and Applications, 2008.
- [4] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica, M. Zaharia, [Above the clouds: A berkeley view of cloud computing](#), Tech. Rep. UCB/EECS-2009-28, EECS Department, University of California, Berkeley (Feb 2009).  
URL <http://www.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-28.html>

- [5] Q. Zhang, L. Cheng, R. Boutaba, [Cloud computing: state-of-the-art and research challenges](#), J. Internet Serv. Appl. 1 (1) (2010) 7–18. doi:[10.1007/s13174-010-0007-6](#).  
URL <https://doi.org/10.1007/s13174-010-0007-6>
- [6] M. D. de Assunção, A. D. S. Veith, R. Buyya, [Distributed data stream processing and edge computing: A survey on resource elasticity and future directions](#), J. Netw. Comput. Appl. 103 (2018) 1–17. doi:[10.1016/j.jnca.2017.12.001](#).  
URL <https://doi.org/10.1016/j.jnca.2017.12.001>
- [7] S. K. A. Hossain, M. A. Rahman, M. A. Hossain, [Edge computing framework for enabling situation awareness in iot based smart city](#), J. Parallel Distributed Comput. 122 (2018) 226–237. doi:[10.1016/j.jpdc.2018.08.009](#).  
URL <https://doi.org/10.1016/j.jpdc.2018.08.009>
- [8] J. Cao, Q. Zhang, W. Shi, [Edge Computing: A Primer](#), Springer Briefs in Computer Science, Springer, 2018. doi:[10.1007/978-3-030-02083-5](#).  
URL <https://doi.org/10.1007/978-3-030-02083-5>
- [9] M. F. Bari, R. Boutaba, R. P. Esteves, L. Z. Granville, M. Podlesny, M. G. Rabbani, Q. Zhang, M. F. Zhani, [Data center network virtualization: A survey](#), IEEE Commun. Surv. Tutorials 15 (2) (2013) 909–928. doi:[10.1109/SURV.2012.090512.00043](#).  
URL <https://doi.org/10.1109/SURV.2012.090512.00043>
- [10] P. G. López, A. Montresor, D. H. J. Epema, A. Datta, T. Higashino, A. Iamnitchi, M. P. Barcellos, P. Felber, E. Rivière, [Edge-centric computing: Vision and challenges](#), Comput. Commun. Rev. 45 (5) (2015) 37–42. doi:[10.1145/2831347.2831354](#).  
URL <https://doi.org/10.1145/2831347.2831354>
- [11] M. B. A. Karim, B. I. Ismail, M. Wong, E. M. Goortani, S. Setapa, L. J. Yuan, H. Ong, [Extending cloud resources to the edge:](#)

## BIBLIOGRAPHY

---

- Possible scenarios, challenges, and experiments, in: International Conference on Cloud Computing Research and Innovations, ICC-CRI 2016, Singapore, Singapore, May 4-5, 2016, IEEE Computer Society, 2016, pp. 78–85. doi:10.1109/ICCCRI.2016.20.  
URL <https://doi.org/10.1109/ICCCRI.2016.20>
- [12] S. A. Monsalve, F. G. Carballeira, A. Calderón, [A heterogeneous mobile cloud computing model for hybrid clouds](#), Future Gener. Comput. Syst. 87 (2018) 651–666. doi:10.1016/j.future.2018.04.005.  
URL <https://doi.org/10.1016/j.future.2018.04.005>
- [13] M. Satyanarayanan, [The emergence of edge computing](#), Computer 50 (1) (2017) 30–39. doi:10.1109/MC.2017.9.  
URL <https://doi.org/10.1109/MC.2017.9>
- [14] A. Rotem-Gal-Oz, [Fallacies of distributed computing explained](#), Doctor Dobbs Journal (01 2008).  
URL [https://www.researchgate.net/publication/322500050\\_Fallacies\\_of\\_Distributed\\_Computing\\_Explained](https://www.researchgate.net/publication/322500050_Fallacies_of_Distributed_Computing_Explained)
- [15] H. Ning, Y. Li, F. Shi, L. T. Yang, [Heterogeneous edge computing open platforms and tools for internet of things](#), Future Gener. Comput. Syst. 106 (2020) 67–76. doi:10.1016/j.future.2019.12.036.  
URL <https://doi.org/10.1016/j.future.2019.12.036>
- [16] F. Bonomi, R. A. Milito, P. Natarajan, J. Zhu, [Fog computing: A platform for internet of things and analytics](#), in: N. Bessis, C. Dobre (Eds.), Big Data and Internet of Things: A Roadmap for Smart Environments, Vol. 546 of Studies in Computational Intelligence, Springer, 2014, pp. 169–186. doi:10.1007/978-3-319-05029-4\_7.  
URL [https://doi.org/10.1007/978-3-319-05029-4\\_7](https://doi.org/10.1007/978-3-319-05029-4_7)

- [17] S. Wang, X. Zhang, Y. Zhang, L. Wang, J. Yang, W. Wang, [A survey on mobile edge networks: Convergence of computing, caching and communications](#), *IEEE Access* 5 (2017) 6757–6779. doi:[10.1109/ACCESS.2017.2685434](#).  
URL <https://doi.org/10.1109/ACCESS.2017.2685434>
- [18] A. Khune, S. Pasricha, [Mobile network-aware middleware framework for cloud offloading: Using reinforcement learning to make reward-based decisions in smartphone applications](#), *IEEE Consumer Electron. Mag.* 8 (1) (2019) 42–48. doi:[10.1109/MCE.2018.2867972](#).  
URL <https://doi.org/10.1109/MCE.2018.2867972>
- [19] M. Chen, Y. Hao, Y. Li, C. Lai, D. Wu, [On the computation offloading at ad hoc cloudlet: architecture and service modes](#), *IEEE Commun. Mag.* 53 (6-Supplement) (2015) 18–24. doi:[10.1109/MCOM.2015.7120041](#).  
URL <https://doi.org/10.1109/MCOM.2015.7120041>
- [20] M. Satyanarayanan, G. Klas, M. D. Silva, S. Mangiante, [The seminal role of edge-native applications](#), in: E. Bertino, C. K. Chang, P. Chen, E. Damiani, M. Goul, K. Oyama (Eds.), 3rd IEEE International Conference on Edge Computing, EDGE 2019, Milan, Italy, July 8-13, 2019, IEEE, 2019, pp. 33–40. doi:[10.1109/EDGE.2019.00022](#).  
URL <https://doi.org/10.1109/EDGE.2019.00022>
- [21] R. V. Aroca, L. M. G. Gonçalves, [Towards green data centers: A comparison of x86 and ARM architectures power efficiency](#), *J. Parallel Distributed Comput.* 72 (12) (2012) 1770–1780. doi:[10.1016/j.jpdc.2012.08.005](#).  
URL <https://doi.org/10.1016/j.jpdc.2012.08.005>
- [22] M. Simic, M. Stojkov, G. Sladic, B. Milosavljević, [Edge computing system for large-scale distributed sensing systems](#), in: *Edge*



## BIBLIOGRAPHY

---

- computing system for large-scale distributed sensing systems, 2018.
- [23] M. Ryden, K. Oh, A. Chandra, J. B. Weissman, [Nebula: Distributed edge cloud for data intensive computing](#), in: 2014 IEEE International Conference on Cloud Engineering, Boston, MA, USA, March 11-14, 2014, IEEE Computer Society, 2014, pp. 57–66. doi:[10.1109/IC2E.2014.34](#).  
URL <https://doi.org/10.1109/IC2E.2014.34>
- [24] A. G. Greenberg, J. R. Hamilton, D. A. Maltz, P. Patel, [The cost of a cloud: research problems in data center networks](#), Comput. Commun. Rev. 39 (1) (2009) 68–73. doi:[10.1145/1496091.1496103](#).  
URL <https://doi.org/10.1145/1496091.1496103>
- [25] N. Abbas, Y. Zhang, A. Taherkordi, T. Skeie, [Mobile edge computing: A survey](#), IEEE Internet Things J. 5 (1) (2018) 450–465. doi:[10.1109/JIOT.2017.2750180](#).  
URL <https://doi.org/10.1109/JIOT.2017.2750180>
- [26] H. Guo, L. Rui, Z. Gao, [A zone-based content pre-caching strategy in vehicular edge networks](#), Future Gener. Comput. Syst. 106 (2020) 22–33. doi:[10.1016/j.future.2019.12.050](#).  
URL <https://doi.org/10.1016/j.future.2019.12.050>
- [27] I. Kurniawan, H. Febiansyah, J. Kwon, Cost-Effective Content Delivery Networks Using Clouds and Nano Data Centers, Vol. 280, 2014, pp. 417–424. doi:[10.1007/978-3-642-41671-2\\_53](#).
- [28] F. R. de Souza, C. C. Miers, A. Fiorese, M. D. de Assunção, G. P. Koslovski, [QVIA-SDN: towards qos-aware virtual infrastructure allocation on sdn-based clouds](#), J. Grid Comput. 17 (3) (2019) 447–472. doi:[10.1007/s10723-019-09479-x](#).  
URL <https://doi.org/10.1007/s10723-019-09479-x>

- [29] J. R. Hamilton, [An architecture for modular data centers](#), in: CIDR 2007, Third Biennial Conference on Innovative Data Systems Research, Asilomar, CA, USA, January 7-10, 2007, Online Proceedings, [www.cidrdb.org](http://cidrdb.org), 2007, pp. 306–313.  
URL <http://cidrdb.org/cidr2007/papers/cidr07p35.pdf>
- [30] K. Sonbol, Ö. Özkasap, I. Al-Oqily, M. Aloqaily, [Edgekv: Decentralized, scalable, and consistent storage for the edge](#), J. Parallel Distributed Comput. 144 (2020) 28–40. doi:[10.1016/j.jpdc.2020.05.009](https://doi.org/10.1016/j.jpdc.2020.05.009).  
URL <https://doi.org/10.1016/j.jpdc.2020.05.009>
- [31] J. Wang, D. Crawl, I. Altintas, W. Li, [Big data applications using workflows for data parallel computing](#), Comput. Sci. Eng. 16 (4) (2014) 11–21. doi:[10.1109/MCSE.2014.50](https://doi.org/10.1109/MCSE.2014.50).  
URL <https://doi.org/10.1109/MCSE.2014.50>
- [32] A. Das, I. Gupta, A. Motivala, [SWIM: scalable weakly-consistent infection-style process group membership protocol](#), in: 2002 International Conference on Dependable Systems and Networks (DSN 2002), 23-26 June 2002, Bethesda, MD, USA, Proceedings, IEEE Computer Society, 2002, pp. 303–312. doi:[10.1109/DSN.2002.1028914](https://doi.org/10.1109/DSN.2002.1028914).  
URL <https://doi.org/10.1109/DSN.2002.1028914>
- [33] C. Li, J. Bai, Y. Chen, Y. Luo, [Resource and replica management strategy for optimizing financial cost and user experience in edge cloud computing system](#), Inf. Sci. 516 (2020) 33–55. doi:[10.1016/j.ins.2019.12.049](https://doi.org/10.1016/j.ins.2019.12.049).  
URL <https://doi.org/10.1016/j.ins.2019.12.049>
- [34] E. Cau, M. Corici, P. Bellavista, L. Foschini, G. Carella, A. Edmonds, T. M. Bohnert, [Efficient exploitation of mobile edge computing for virtualized 5g in EPC architectures](#), in: 4th IEEE International Conference on Mobile Cloud Computing, Services,

- and Engineering, MobileCloud 2016, Oxford, United Kingdom, March 29 - April 1, 2016, IEEE Computer Society, 2016, pp. 100–109. doi:[10.1109/MobileCloud.2016.24](https://doi.org/10.1109/MobileCloud.2016.24).  
URL <https://doi.org/10.1109/MobileCloud.2016.24>
- [35] W. Yu, C.-L. Ignat, [Conflict-Free Replicated Relations for Multi-Synchronous Database Management at Edge](#), in: IEEE International Conference on Smart Data Services, 2020 IEEE World Congress on Services, Beijing, China, 2020.  
URL <https://hal.inria.fr/hal-02983557>
- [36] Y. Duan, G. Fu, N. Zhou, X. Sun, N. C. Narendra, B. Hu, [Everything as a service \(xaas\) on the cloud: Origins, current and future trends](#), in: C. Pu, A. Mohindra (Eds.), 8th IEEE International Conference on Cloud Computing, CLOUD 2015, New York City, NY, USA, June 27 - July 2, 2015, IEEE Computer Society, 2015, pp. 621–628. doi:[10.1109/CLOUD.2015.88](https://doi.org/10.1109/CLOUD.2015.88).  
URL <https://doi.org/10.1109/CLOUD.2015.88>
- [37] J. Gubbi, R. Buyya, S. Marusic, M. Palaniswami, [Internet of things \(iot\): A vision, architectural elements, and future directions](#), Future Gener. Comput. Syst. 29 (7) (2013) 1645–1660. doi:[10.1016/j.future.2013.01.010](https://doi.org/10.1016/j.future.2013.01.010).  
URL <https://doi.org/10.1016/j.future.2013.01.010>
- [38] C. Jiang, X. Cheng, H. Gao, X. Zhou, J. Wan, [Toward computation offloading in edge computing: A survey](#), IEEE Access 7 (2019) 131543–131558. doi:[10.1109/ACCESS.2019.2938660](https://doi.org/10.1109/ACCESS.2019.2938660).  
URL <https://doi.org/10.1109/ACCESS.2019.2938660>
- [39] X. Jin, S. Chun, J. Jung, K. Lee, [Iot service selection based on physical service model and absolute dominance relationship](#), in: 7th IEEE International Conference on Service-Oriented Computing and Applications, SOCA 2014, Matsue, Japan, November 17-19, 2014, IEEE Computer Society, 2014, pp. 65–72.

- [doi:10.1109/SOCA.2014.24](https://doi.org/10.1109/SOCA.2014.24).  
URL <https://doi.org/10.1109/SOCA.2014.24>
- [40] M. Satyanarayanan, P. Bahl, R. Cáceres, N. Davies, [The case for vm-based cloudlets in mobile computing](#), IEEE Pervasive Comput. 8 (4) (2009) 14–23. [doi:10.1109/MPRV.2009.82](https://doi.org/10.1109/MPRV.2009.82).  
URL <https://doi.org/10.1109/MPRV.2009.82>
- [41] Y. Yao, B. Xiao, W. Wang, G. Yang, X. Zhou, Z. Peng, [Real-time cache-aided route planning based on mobile edge computing](#), IEEE Wirel. Commun. 27 (5) (2020) 155–161. [doi:10.1109/MWC.001.1900559](https://doi.org/10.1109/MWC.001.1900559).  
URL <https://doi.org/10.1109/MWC.001.1900559>
- [42] R. Hu, N. Yoshida, [Explicit connection actions in multiparty session types](#), in: M. Huisman, J. Rubin (Eds.), Fundamental Approaches to Software Engineering - 20th International Conference, FASE 2017, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2017, Uppsala, Sweden, April 22-29, 2017, Proceedings, Vol. 10202 of Lecture Notes in Computer Science, Springer, 2017, pp. 116–133. [doi:10.1007/978-3-662-54494-5\\_7](https://doi.org/10.1007/978-3-662-54494-5_7).  
URL [https://doi.org/10.1007/978-3-662-54494-5\\_7](https://doi.org/10.1007/978-3-662-54494-5_7)
- [43] K. Honda, N. Yoshida, M. Carbone, [Multiparty asynchronous session types](#), in: G. C. Necula, P. Wadler (Eds.), Proceedings of the 35th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2008, San Francisco, California, USA, January 7-12, 2008, ACM, 2008, pp. 273–284. [doi:10.1145/1328438.1328472](https://doi.org/10.1145/1328438.1328472).  
URL <https://doi.org/10.1145/1328438.1328472>
- [44] H. El-Sayed, S. Sankar, M. Prasad, D. Puthal, A. Gupta, M. Mohanty, C. Lin, [Edge of things: The big picture on the](#)

- integration of edge, iot and the cloud in a distributed computing environment, *IEEE Access* 6 (2018) 1706–1717. doi:[10.1109/ACCESS.2017.2780087](https://doi.org/10.1109/ACCESS.2017.2780087).  
URL <https://doi.org/10.1109/ACCESS.2017.2780087>
- [45] Y. Jararweh, A. Doulat, O. AlQudah, E. Ahmed, M. Al-Ayyoub, E. Benkhelifa, [The future of mobile cloud computing: Integrating cloudlets and mobile edge computing](#), in: 23rd International Conference on Telecommunications, ICT 2016, Thessaloniki, Greece, May 16-18, 2016, IEEE, 2016, pp. 1–5. doi:[10.1109/ICT.2016.7500486](https://doi.org/10.1109/ICT.2016.7500486).  
URL <https://doi.org/10.1109/ICT.2016.7500486>
- [46] S. Miloš, P. Ivan, D. Jovana, S. Goran, M. Branko, Towards edge computing as a service: Dynamic formation of the micro data-centers, *IEEE Access* 9 (2021) 114468–114484. doi:[10.1109/ACCESS.2021.3104475](https://doi.org/10.1109/ACCESS.2021.3104475).
- [47] M. Hirsch, C. Mateos, A. Zunino, [Augmenting computing capabilities at the edge by jointly exploiting mobile devices: A survey](#), *Future Gener. Comput. Syst.* 88 (2018) 644–662. doi:[10.1016/j.future.2018.06.005](https://doi.org/10.1016/j.future.2018.06.005).  
URL <https://doi.org/10.1016/j.future.2018.06.005>
- [48] M. van Steen, A. S. Tanenbaum, [A brief introduction to distributed systems](#), *Computing* 98 (10) (2016) 967–1009. doi:[10.1007/s00607-016-0508-7](https://doi.org/10.1007/s00607-016-0508-7).  
URL <https://doi.org/10.1007/s00607-016-0508-7>
- [49] A. S. Tanenbaum, M. van Steen, *Distributed systems - principles and paradigms*, 2nd Edition, Pearson Education, 2007.
- [50] A. B. Bondi, [Characteristics of scalability and their impact on performance](#), in: Second International Workshop on Software and Performance, WOSP 2000, Ottawa, Canada, September 17-20,

- 2000, ACM, 2000, pp. 195–203. doi:10.1145/350391.350432.  
URL <https://doi.org/10.1145/350391.350432>
- [51] J. L. Gustafson, [Moore’s Law](#), Springer US, Boston, MA, 2011, pp. 1177–1184. doi:10.1007/978-0-387-09766-4\_81.  
URL [https://doi.org/10.1007/978-0-387-09766-4\\_81](https://doi.org/10.1007/978-0-387-09766-4_81)
- [52] R. Buyya, High Performance Cluster Computing: Architectures and Systems, Prentice Hall PTR, USA, 1999.
- [53] E. A. Brewer, [Towards robust distributed systems.](#), in: Symposium on Principles of Distributed Computing (PODC), 2000.  
URL <http://www.cs.berkeley.edu/~brewer/cs262b-2004/PODC-keynote.pdf>
- [54] S. Gilbert, N. A. Lynch, [Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services](#), SIGACT News 33 (2) (2002) 51–59. doi:10.1145/564585.564601.  
URL <https://doi.org/10.1145/564585.564601>
- [55] J. Gray, D. P. Siewiorek, [High-availability computer systems](#), Computer 24 (9) (1991) 39–48. doi:10.1109/2.84898.  
URL <https://doi.org/10.1109/2.84898>
- [56] M. Shapiro, N. M. Preguiça, C. Baquero, M. Zawirski, [Conflict-free replicated data types](#), in: X. Défago, F. Petit, V. Villain (Eds.), Stabilization, Safety, and Security of Distributed Systems - 13th International Symposium, SSS 2011, Grenoble, France, October 10-12, 2011. Proceedings, Vol. 6976 of Lecture Notes in Computer Science, Springer, 2011, pp. 386–400. doi:10.1007/978-3-642-24550-3\_29.  
URL [https://doi.org/10.1007/978-3-642-24550-3\\_29](https://doi.org/10.1007/978-3-642-24550-3_29)
- [57] P. M. Mell, T. Grance, Sp 800-145. the nist definition of cloud computing, Tech. rep., Gaithersburg, MD, USA (2011).

- [58] D. Cohen, T. Talpey, A. Kanevsky, U. Cummings, M. Krause, R. Recio, D. Crupnicoff, L. Dickman, P. Grun, [Remote direct memory over the converged enhanced ethernet fabric: Evaluating the options](#), in: K. Bergman, R. Brightwell, F. Petrini, H. Bubba (Eds.), 17th IEEE Symposium on High Performance Interconnects, HOTI 2009, New York, New York, USA, August 25-27, 2009, IEEE Computer Society, 2009, pp. 123–130. [doi:10.1109/HOTI.2009.23](#).  
URL <https://doi.org/10.1109/HOTI.2009.23>
- [59] J. Hong, T. Dreibholz, J. A. Schenkel, J. A. Hu, [An overview of multi-cloud computing](#), in: L. Barolli, M. Takizawa, F. Xhafa, T. Enokido (Eds.), Web, Artificial Intelligence and Network Applications - Proceedings of the Workshops of the 33rd International Conference on Advanced Information Networking and Applications, AINA Workshops 2019, Matsue, Japan, March 27-29, 2019, Vol. 927 of Advances in Intelligent Systems and Computing, Springer, 2019, pp. 1055–1068. [doi:10.1007/978-3-030-15035-8\\_103](#).  
URL [https://doi.org/10.1007/978-3-030-15035-8\\_103](https://doi.org/10.1007/978-3-030-15035-8_103)
- [60] D. Ardagna, [Cloud and multi-cloud computing: Current challenges and future applications](#), in: M. A. Babar, H. Paik, M. Chetlur, M. Bauer, A. M. Sharifloo (Eds.), 7th IEEE/ACM International Workshop on Principles of Engineering Service-Oriented and Cloud Systems, PESOS 2015, Florence, Italy, May 23, 2015, IEEE Computer Society, 2015, pp. 1–2. [doi:10.1109/PESOS.2015.8](#).  
URL <https://doi.org/10.1109/PESOS.2015.8>
- [61] I. Stoica, S. Shenker, [From cloud computing to sky computing](#), in: S. Angel, B. Kasikci, E. Kohler (Eds.), HotOS '21: Workshop on Hot Topics in Operating Systems, Ann Arbor, Michigan, USA, June, 1-3, 2021, ACM, 2021, pp. 26–32. [doi:10.1145/3458336](#).

3465301.  
URL <https://doi.org/10.1145/3458336.3465301>
- [62] L. G., A. L., [Multicloud Architecture Migration and Security](#), O'Reilly Media, 2019.  
URL <https://www.oreilly.com/library/view/multicloud-architecture-migration/9781492050407/>
- [63] A. Dadgar, J. Phillips, J. Currey, [Lifeguard: Local health awareness for more accurate failure detection](#), in: 48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops, DSN Workshops 2018, Luxembourg, June 25-28, 2018, IEEE Computer Society, 2018, pp. 22–25. doi:10.1109/DSN-W.2018.00017.  
URL <http://doi.ieeecomputersociety.org/10.1109/DSN-W.2018.00017>
- [64] N. Fernando, S. W. Loke, J. W. Rahayu, [Mobile cloud computing: A survey](#), Future Gener. Comput. Syst. 29 (1) (2013) 84–106. doi:10.1016/j.future.2012.05.023.  
URL <https://doi.org/10.1016/j.future.2012.05.023>
- [65] L. Lin, X. Liao, H. Jin, P. Li, [Computation offloading toward edge computing](#), Proc. IEEE 107 (8) (2019) 1584–1607. doi:10.1109/JPROC.2019.2922285.  
URL <https://doi.org/10.1109/JPROC.2019.2922285>
- [66] M. Villari, M. Fazio, S. Dustdar, O. Rana, D. N. Jha, R. Ranjan, [Osmosis: The osmotic computing platform for microelements in the cloud, edge, and internet of things](#), Computer 52 (8) (2019) 14–26. doi:10.1109/MC.2018.2888767.  
URL <https://doi.org/10.1109/MC.2018.2888767>
- [67] M. Villari, A. Celesti, M. Fazio, [Towards osmotic computing: Looking at basic principles and technologies](#), in: L. Barolli, O. Terzo (Eds.), Complex, Intelligent, and Software Intensive



- Systems - Proceedings of the 11th International Conference on Complex, Intelligent, and Software Intensive Systems (CISIS-2017), Torino, Italy, July 10-12, 2017, Vol. 611 of Advances in Intelligent Systems and Computing, Springer, 2017, pp. 906–915. doi:[10.1007/978-3-319-61566-0\\_86](https://doi.org/10.1007/978-3-319-61566-0_86).  
URL [https://doi.org/10.1007/978-3-319-61566-0\\_86](https://doi.org/10.1007/978-3-319-61566-0_86)
- [68] R. de Vera Jr., [Review of: Distributed systems: An algorithmic approach \(2nd edition\) by su Kumar ghosh](#), SIGACT News 47 (4) (2016) 13–14. doi:[10.1145/3023855.3023860](https://doi.org/10.1145/3023855.3023860).  
URL <https://doi.org/10.1145/3023855.3023860>
- [69] G. Andrews, [, parallel, and distributed programming](#), Addison-Wesley, 2000.  
URL <http://books.google.com.br/books?id=npRQAAAAAMAAJ>
- [70] D. Fisher, R. DeLine, M. Czerwinski, S. M. Drucker, [Interactions with big data analytics](#), Interactions 19 (3) (2012) 50–59. doi:[10.1145/2168931.2168943](https://doi.org/10.1145/2168931.2168943).  
URL <https://doi.org/10.1145/2168931.2168943>
- [71] C.-W. Tsai, C.-F. Lai, H.-C. Chao, A. V. Vasilakos, [Big data analytics: a survey](#), Journal of Big Data 2 (1) (2015) 21. doi:[10.1186/s40537-015-0030-3](https://doi.org/10.1186/s40537-015-0030-3).  
URL <https://doi.org/10.1186/s40537-015-0030-3>
- [72] Z. Guo, G. C. Fox, M. Zhou, [Investigation of data locality in mapreduce](#), in: 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, CCGrid 2012, Ottawa, Canada, May 13-16, 2012, IEEE Computer Society, 2012, pp. 419–426. doi:[10.1109/CCGrid.2012.42](https://doi.org/10.1109/CCGrid.2012.42).  
URL <https://doi.org/10.1109/CCGrid.2012.42>
- [73] P. G. Sarigiannidis, T. Lagkas, K. Rantos, P. Bellavista, [The big data era in iot-enabled smart farming: Re-defining systems, tools, and techniques](#), Comput. Networks 168 (2020). doi:[10.1016/j.comnet.2020.107388](https://doi.org/10.1016/j.comnet.2020.107388).

- 1016/j.comnet.2019.107043.  
URL <https://doi.org/10.1016/j.comnet.2019.107043>
- [74] R. Patgiri, A. Ahmed, [Big data: The v's of the game changer paradigm](#), in: J. Chen, L. T. Yang (Eds.), 18th IEEE International Conference on High Performance Computing and Communications; 14th IEEE International Conference on Smart City; 2nd IEEE International Conference on Data Science and Systems, HPCC/SmartCity/DSS 2016, Sydney, Australia, December 12-14, 2016, IEEE Computer Society, 2016, pp. 17–24. [doi:10.1109/HPCC-SmartCity-DSS.2016.0014](https://doi.org/10.1109/HPCC-SmartCity-DSS.2016.0014).  
URL <https://doi.org/10.1109/HPCC-SmartCity-DSS.2016.0014>
- [75] Y. Kumar, Lambda architecture - realtime data processing, Ph.D. thesis (01 2020). [doi:10.13140/RG.2.2.19091.84004](https://doi.org/10.13140/RG.2.2.19091.84004).
- [76] M. Kiran, P. Murphy, I. Monga, J. Dugan, S. S. Baveja, [Lambda architecture for cost-effective batch and speed big data processing](#), in: 2015 IEEE International Conference on Big Data, Big Data 2015, Santa Clara, CA, USA, October 29 - November 1, 2015, IEEE Computer Society, 2015, pp. 2785–2792. [doi:10.1109/BigData.2015.7364082](https://doi.org/10.1109/BigData.2015.7364082).  
URL <https://doi.org/10.1109/BigData.2015.7364082>
- [77] D. Medjedovic, E. Tahirovic, I. Dedovic, [Algorithms and Data Structures for Massive Datasets](#), Manning, 2022.  
URL <https://www.manning.com/books/algorithms-and-data-structures-for-massive-datasets>
- [78] T. R. Rao, P. Mitra, R. Bhatt, A. Goswami, [The big data system, components, tools, and technologies: a survey](#), Knowl. Inf. Syst. 60 (3) (2019) 1165–1245. [doi:10.1007/s10115-018-1248-0](https://doi.org/10.1007/s10115-018-1248-0).  
URL <https://doi.org/10.1007/s10115-018-1248-0>

## BIBLIOGRAPHY

---

- [79] J. E. Marynowski, A. O. Santin, A. R. Pimentel, [Method for testing the fault tolerance of mapreduce frameworks](#), *Comput. Networks* 86 (2015) 1–13. doi:[10.1016/j.comnet.2015.04.009](#). URL <https://doi.org/10.1016/j.comnet.2015.04.009>
- [80] N. Dragoni, S. Giallorenzo, A. Lluch-Lafuente, M. Mazzara, F. Montesi, R. Mustafin, L. Safina, [Microservices: yesterday, today, and tomorrow](#), *CoRR* abs/1606.04036 (2016). arXiv: [1606.04036](#). URL <http://arxiv.org/abs/1606.04036>
- [81] C. Pautasso, O. Zimmermann, M. Amundsen, J. Lewis, N. M. Josuttis, [Microservices in practice, part 1: Reality check and service design](#), *IEEE Softw.* 34 (1) (2017) 91–98. doi:[10.1109/MS.2017.24](#). URL <https://doi.org/10.1109/MS.2017.24>
- [82] S. Li, H. Zhang, Z. Jia, Z. Li, C. Zhang, J. Li, Q. Gao, J. Ge, Z. Shan, [A dataflow-driven approach to identifying microservices from monolithic applications](#), *J. Syst. Softw.* 157 (2019). doi: [10.1016/j.jss.2019.07.008](#). URL <https://doi.org/10.1016/j.jss.2019.07.008>
- [83] L. Krause, [Microservices: Patterns and Applications: Designing Fine-Grained Services by Applying Patterns](#), Lucas Krause, 2015. URL <https://books.google.rs/books?id=dd5-rgEACAAJ>
- [84] O. Al-Debagy, P. Martinek, [A comparative review of microservices and monolithic architectures](#), *CoRR* abs/1905.07997 (2019). arXiv: [1905.07997](#). URL <http://arxiv.org/abs/1905.07997>
- [85] N. Kratzke, P. Quint, [Understanding cloud-native applications after 10 years of cloud computing - A systematic mapping study](#), *J. Syst. Softw.* 126 (2017) 1–16. doi:[10.1016/j.jss.2017.01](#).

001.

URL <https://doi.org/10.1016/j.jss.2017.01.001>

- [86] G. Adzic, R. Chatley, [Serverless computing: economic and architectural impact](#), in: E. Bodden, W. Schäfer, A. van Deursen, A. Zisman (Eds.), Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2017, Paderborn, Germany, September 4-8, 2017, ACM, 2017, pp. 884–889. [doi:10.1145/3106237.3117767](#).  
URL <https://doi.org/10.1145/3106237.3117767>
- [87] W. Li, Y. Lemieux, J. Gao, Z. Zhao, Y. Han, [Service mesh: Challenges, state of the art, and future research opportunities](#), in: 13th IEEE International Conference on Service-Oriented System Engineering, SOSE 2019, San Francisco, CA, USA, April 4-9, 2019, IEEE, 2019. [doi:10.1109/SOSE.2019.00026](#).  
URL <https://doi.org/10.1109/SOSE.2019.00026>
- [88] P. Adamczyk, P. H. Smith, R. E. Johnson, M. Hafiz, [REST and web services: In theory and in practice](#), in: E. Wilde, C. Pautasso (Eds.), REST: From Research to Practice, Springer, 2011, pp. 35–57. [doi:10.1007/978-1-4419-8303-9\\_2](#).  
URL [https://doi.org/10.1007/978-1-4419-8303-9\\_2](https://doi.org/10.1007/978-1-4419-8303-9_2)
- [89] A. Balalaie, A. Heydarnoori, P. Jamshidi, Microservices architecture enables devops: Migration to a cloud-native architecture, IEEE Software 33 (3) (2016) 42–52. [doi:10.1109/MS.2016.64](#).
- [90] W. Felter, A. Ferreira, R. Rajamony, J. Rubio, [An updated performance comparison of virtual machines and linux containers](#), in: 2015 IEEE International Symposium on Performance Analysis of Systems and Software, ISPASS 2015, Philadelphia, PA, USA, March 29-31, 2015, IEEE Computer Society, 2015, pp. 171–172. [doi:10.1109/ISPASS.2015.7095802](#).  
URL <https://doi.org/10.1109/ISPASS.2015.7095802>

## BIBLIOGRAPHY

---

- [91] B. Burns, B. Grant, D. Oppenheimer, E. A. Brewer, J. Wilkes, [Borg, omega, and kubernetes](#), Commun. ACM 59 (5) (2016) 50–57. doi:[10.1145/2890784](#).  
URL <https://doi.org/10.1145/2890784>
- [92] J. Soldani, D. A. Tamburri, W. van den Heuvel, [The pains and gains of microservices: A systematic grey literature review](#), J. Syst. Softw. 146 (2018) 215–232. doi:[10.1016/j.jss.2018.09.082](#).  
URL <https://doi.org/10.1016/j.jss.2018.09.082>
- [93] G. Grätzer, B. Davey, R. Freese, B. Ganter, M. Greferath, P. Jipsen, H. Priestley, H. Rose, E. Schmidt, S. Schmidt, et al., [General Lattice Theory: Second edition](#), Birkhäuser Basel, 2002.  
URL <https://books.google.rs/books?id=SoGLVCpu0z0C>
- [94] C. Richardson, [Microservices Patterns: With examples in Java](#), Manning Publications, 2018.  
URL <https://books.google.rs/books?id=UeK1swEACAAJ>
- [95] Z. Long, Improvement and implementation of a high performance cqrs architecture, in: 2017 International Conference on Robots Intelligent System (ICRIS), 2017, pp. 170–173. doi:[10.1109/ICRIS.2017.49](#).
- [96] A. Akbulut, H. G. Perros, Performance analysis of microservice design patterns, IEEE Internet Computing 23 (6) (2019) 19–27. doi:[10.1109/MIC.2019.2951094](#).
- [97] C. J. Fidge, [Fundamentals of distributed system observation](#), IEEE Softw. 13 (6) (1996) 77–83. doi:[10.1109/52.542297](#).  
URL <https://doi.org/10.1109/52.542297>
- [98] J. Joyce, G. Lomow, K. Slind, B. W. Unger, [Monitoring distributed systems](#), ACM Trans. Comput. Syst. 5 (2) (1987) 121–150. doi:[10.1145/13677.22723](#).  
URL <https://doi.org/10.1145/13677.22723>

- [99] I. Beschastnikh, P. Wang, Y. Brun, M. D. Ernst, [Debugging distributed systems](#), Commun. ACM 59 (8) (2016) 32–37. doi: [10.1145/2909480](#).  
URL <https://doi.org/10.1145/2909480>
- [100] D. S. Daniels, A. Z. Spector, D. S. Thompson, [Distributed logging for transaction processing](#), in: U. Dayal, I. L. Traiger (Eds.), Proceedings of the Association for Computing Machinery Special Interest Group on Management of Data 1987 Annual Conference, San Francisco, CA, USA, May 27-29, 1987, ACM Press, 1987, pp. 82–96. doi: [10.1145/38713.38728](#).  
URL <https://doi.org/10.1145/38713.38728>
- [101] B. H. Sigelman, L. A. Barroso, M. Burrows, P. Stephenson, M. Plakal, D. Beaver, S. Jaspán, C. Shanbhag, [Dapper, a large-scale distributed systems tracing infrastructure](#), Tech. rep., Google, Inc. (2010).  
URL <https://research.google.com/archive/papers/dapper-2010-1.pdf>
- [102] R. Schollmeier, [A definition of peer-to-peer networking for the classification of peer-to-peer architectures and applications](#), in: R. L. Graham, N. Shahmehri (Eds.), 1st International Conference on Peer-to-Peer Computing (P2P 2001), 27-29 August 2001, Linköping, Sweden, IEEE Computer Society, 2001, pp. 101–102. doi: [10.1109/P2P.2001.990434](#).  
URL <https://doi.org/10.1109/P2P.2001.990434>
- [103] H. M. N. D. Bandara, A. P. Jayasumana, [Collaborative applications over peer-to-peer systems-challenges and solutions](#), Peer Peer Netw. Appl. 6 (3) (2013) 257–276. doi: [10.1007/s12083-012-0157-3](#).  
URL <https://doi.org/10.1007/s12083-012-0157-3>
- [104] M. Kamel, C. M. Scoglio, T. Easton, [Optimal topology design for overlay networks](#), in: I. F. Akyildiz, R. Sivakumar, E. Ekici,

## BIBLIOGRAPHY

---

- J. C. de Oliveira, J. McNair (Eds.), NETWORKING 2007. Ad Hoc and Sensor Networks, Wireless Networks, Next Generation Internet, 6th International IFIP-TC6 Networking Conference, Atlanta, GA, USA, May 14-18, 2007, Proceedings, Vol. 4479 of Lecture Notes in Computer Science, Springer, 2007, pp. 714–725. doi:[10.1007/978-3-540-72606-7\\_61](https://doi.org/10.1007/978-3-540-72606-7_61).  
URL [https://doi.org/10.1007/978-3-540-72606-7\\_61](https://doi.org/10.1007/978-3-540-72606-7_61)
- [105] I. Filali, F. Bongiovanni, F. Huet, F. Baude, [A survey of structured P2P systems for RDF data storage and retrieval](#), Trans. Large Scale Data Knowl. Centered Syst. 3 (2011) 20–55. doi:[10.1007/978-3-642-23074-5\\_2](https://doi.org/10.1007/978-3-642-23074-5_2).  
URL [https://doi.org/10.1007/978-3-642-23074-5\\_2](https://doi.org/10.1007/978-3-642-23074-5_2)
- [106] I. Stoica, R. T. Morris, D. R. Karger, M. F. Kaashoek, H. Balakrishnan, [Chord: A scalable peer-to-peer lookup service for internet applications](#), in: R. L. Cruz, G. Varghese (Eds.), Proceedings of the ACM SIGCOMM 2001 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication, August 27-31, 2001, San Diego, CA, USA, ACM, 2001, pp. 149–160. doi:[10.1145/383059.383071](https://doi.org/10.1145/383059.383071).  
URL <https://doi.org/10.1145/383059.383071>
- [107] N. Leavitt, [Will nosql databases live up to their promise?](#), Computer 43 (2) (2010) 12–14. doi:[10.1109/MC.2010.58](https://doi.org/10.1109/MC.2010.58).  
URL <https://doi.org/10.1109/MC.2010.58>
- [108] Q. H. Vu, M. Lupu, B. C. Ooi, [Peer-to-Peer Computing - Principles and Applications](#), Springer, 2010. doi:[10.1007/978-3-642-03514-2](https://doi.org/10.1007/978-3-642-03514-2).  
URL <https://doi.org/10.1007/978-3-642-03514-2>
- [109] E. Korach, S. Kutten, S. Moran, [A modular technique for the design of efficient distributed leader finding algorithms](#), ACM Trans. Program. Lang. Syst. 12 (1) (1990) 84–101. doi:[10.1145/](https://doi.org/10.1145/)

- 77606.77610.  
URL <https://doi.org/10.1145/77606.77610>
- [110] G. S. Almási, A. Gottlieb, Highly parallel computing (2. ed.), Addison-Wesley, 1994.
- [111] S. Leible, S. Schlager, M. Schubotz, B. Gipp, [A review on blockchain technology and blockchain projects fostering open science](#), Frontiers Blockchain 2 (2019) 16. doi:10.3389/fbloc.2019.00016.  
URL <https://doi.org/10.3389/fbloc.2019.00016>
- [112] G. Weikum, G. Vossen, Transactional Information Systems: Theory, Algorithms, and the Practice of Concurrency Control and Recovery, Morgan Kaufmann, 2002.
- [113] J. Gray, A. Reuter, Transaction Processing: Concepts and Techniques, Morgan Kaufmann, 1993.
- [114] L. Frank, T. U. Zahle, Semantic acid properties in multidatabases using remote procedure calls and update propagations, Softw. Pract. Exper. 28 (1) (1998) 77–98.
- [115] P. Helland, [Life beyond distributed transactions: an apostate's opinion](#), in: Third Biennial Conference on Innovative Data Systems Research, CIDR 2007, Asilomar, CA, USA, January 7-10, 2007, Online Proceedings, [www.cidrdb.org](http://cidrdb.org), 2007, pp. 132–141.  
URL <http://cidrdb.org/cidr2007/papers/cidr07p15.pdf>
- [116] H. Garcia-Molina, K. Salem, [Sagas](#), in: U. Dayal, I. L. Traiger (Eds.), Proceedings of the Association for Computing Machinery Special Interest Group on Management of Data 1987 Annual Conference, San Francisco, CA, USA, May 27-29, 1987, ACM Press, 1987, pp. 249–259. doi:10.1145/38713.38742.  
URL <https://doi.org/10.1145/38713.38742>



## BIBLIOGRAPHY

---

- [117] R. E. Jones, R. D. Lins, Garbage collection - algorithms for automatic dynamic memory management, Wiley, 1996.
- [118] J. McCarthy, [Recursive functions of symbolic expressions and their computation by machine, part I](#), Commun. ACM 3 (4) (1960) 184–195. doi:[10.1145/367177.367199](#).  
URL <https://doi.org/10.1145/367177.367199>
- [119] S. Crosby, D. Brown, [The virtualization reality](#), ACM Queue 4 (10) (2006) 34–41. doi:[10.1145/1189276.1189289](#).  
URL <https://doi.org/10.1145/1189276.1189289>
- [120] S. Sharma, Y. Park, Virtualization: A review and future directions executive overview, American Journal of Information Technology 1 (2011) 1–37.
- [121] E. E. Absalom, S. M. Buhari, S. B. Junaidu, [Virtual machine allocation in cloud computing environment](#), Int. J. Cloud Appl. Comput. 3 (2) (2013) 47–60. doi:[10.4018/ijcac.2013040105](#).  
URL <https://doi.org/10.4018/ijcac.2013040105>
- [122] C. Yang, K. Huang, W. C. Chu, F. Leu, S. Wang, [Implementation of cloud iaas for virtualization with live migration](#), in: J. J. Park, H. R. Arabnia, C. Kim, W. Shi, J. Gil (Eds.), Grid and Pervasive Computing - 8th International Conference, GPC 2013 and Colocated Workshops, Seoul, Korea, May 9-11, 2013. Proceedings, Vol. 7861 of Lecture Notes in Computer Science, Springer, 2013, pp. 199–207. doi:[10.1007/978-3-642-38027-3\\_21](#).  
URL [https://doi.org/10.1007/978-3-642-38027-3\\_21](https://doi.org/10.1007/978-3-642-38027-3_21)
- [123] K.-T. Seo, H.-S. Hwang, I. Moon, O.-Y. Kwon, B. jun Kim, Performance comparison analysis of linux container and virtual machine for building cloud, 2014.
- [124] R. Pavlicek, [Unikernels: Beyond Containers to the Next Generation of Cloud](#), O'Reilly Media, 2016.  
URL <https://books.google.rs/books?id=qfDXuQEACAAJ>

- [125] T. Goethals, M. Sebrechts, A. Atrey, B. Volckaert, F. D. Turck, [Unikernels vs containers: An in-depth benchmarking study in the context of microservice applications](#), in: 8th IEEE International Symposium on Cloud and Service Computing, SC2 2018, Paris, France, November 18-21, 2018, IEEE, 2018, pp. 1–8. doi:10.1109/SC2.2018.00008.  
URL <https://doi.org/10.1109/SC2.2018.00008>
- [126] R. Pavlicek, The next generation cloud: Unleashing the power of the unikernel, USENIX Association, Washington, D.C., 2015.
- [127] M. Plauth, L. Feinbube, A. Polze, [A performance survey of lightweight virtualization techniques](#), in: F. D. Paoli, S. Schulte, E. B. Johnsen (Eds.), Service-Oriented and Cloud Computing - 6th IFIP WG 2.14 European Conference, ESOC 2017, Oslo, Norway, September 27-29, 2017, Proceedings, Vol. 10465 of Lecture Notes in Computer Science, Springer, 2017, pp. 34–48. doi:10.1007/978-3-319-67262-5\_3.  
URL [https://doi.org/10.1007/978-3-319-67262-5\\_3](https://doi.org/10.1007/978-3-319-67262-5_3)
- [128] P. Helland, [Immutability changes everything](#), Commun. ACM 59 (1) (2016) 64–70. doi:10.1145/2844112.  
URL <https://doi.org/10.1145/2844112>
- [129] M. Perry, [The Art of Immutable Architecture: Theory and Practice of Data Management in Distributed Systems](#), Apress, 2020.  
URL <https://books.google.rs/books?id=Ea9tzQEACAAJ>
- [130] P. C. de Guzmán, F. Gorostiaga, C. Sánchez, [i2kit: A tool for immutable infrastructure deployments based on lightweight virtual machines specialized to run containers](#), CoRR abs/1802.10375 (2018). arXiv:1802.10375.  
URL <http://arxiv.org/abs/1802.10375>

## BIBLIOGRAPHY

---

- [131] B. Fitzgerald, N. Forsgren, K.-J. Stol, J. Humble, B. Doody, Infrastructure is software too!, SSRN Electronic Journal (01 2015). [doi:10.2139/ssrn.2681904](https://doi.org/10.2139/ssrn.2681904).
- [132] L. J. Osterweil, [Software processes are software too, revisited: An invited talk on the most influential paper of icse 9](#), in: Proceedings of the 19th International Conference on Software Engineering, ICSE '97, Association for Computing Machinery, New York, NY, USA, 1997, p. 540–548. [doi:10.1145/253228.253440](https://doi.org/10.1145/253228.253440).  
URL <https://doi.org/10.1145/253228.253440>
- [133] A. R. and Rezvan Mahdavi-Hezaveh and Laurie A. Williams, [A systematic mapping study of infrastructure as code research](#), Inf. Softw. Technol. 108 (2019) 65–77. [doi:10.1016/j.infsof.2018.12.004](https://doi.org/10.1016/j.infsof.2018.12.004).  
URL <https://doi.org/10.1016/j.infsof.2018.12.004>
- [134] A. Wittig, M. Wittig, [Amazon Web Services in Action](#), Manning Publications, 2018.  
URL <https://books.google.rs/books?id=-LRotAEACAAJ>
- [135] M. Artac, T. Borovsak, E. D. Nitto, M. Guerriero, D. A. Tamburri, [Devops: introducing infrastructure-as-code](#), in: S. Uchitel, A. Orso, M. P. Robillard (Eds.), Proceedings of the 39th International Conference on Software Engineering, ICSE 2017, Buenos Aires, Argentina, May 20-28, 2017 - Companion Volume, IEEE Computer Society, 2017, pp. 497–498. [doi:10.1109/ICSE-C.2017.162](https://doi.org/10.1109/ICSE-C.2017.162).  
URL <https://doi.org/10.1109/ICSE-C.2017.162>
- [136] M. Luksa, [Kubernetes in Action](#), Manning Publications, 2018.  
URL <https://books.google.rs/books?id=8bE5MQAACAAJ>
- [137] L. A. F. Leite, C. Rocha, F. Kon, D. S. Milojicic, P. Meirelles, [A survey of devops concepts and challenges](#), ACM Comput. Surv.

- 52 (6) (2020) 127:1–127:35. doi:10.1145/3359981.  
URL <https://doi.org/10.1145/3359981>
- [138] R. Jabbari, N. B. Ali, K. Petersen, B. Tanveer, [What is devops?: A systematic mapping study on definitions and practices](#), in: Proceedings of the Scientific Workshop Proceedings of XP2016, Edinburgh, Scotland, UK, May 24, 2016, ACM, 2016, p. 12. doi:10.1145/2962695.2962707.  
URL <https://doi.org/10.1145/2962695.2962707>
- [139] B. Beyer, C. Jones, J. Petoff, N. Murphy, [Site Reliability Engineering: How Google Runs Production Systems](#), O'Reilly Media, Incorporated, 2016.  
URL <https://books.google.rs/books?id=81UrjwEACAAJ>
- [140] C. Jones, T. Underwood, S. Nukala, [Hiring site reliability engineers](#), login Usenix Mag. 40 (3) (2015).  
URL <https://www.usenix.org/publications/login/june15/hiring-site-reliability-engineers>
- [141] B. Beyer, N. R. Murphy, L. Fong-Jones, T. Underwood, L. Nolan, D. K. Rensin, [How sre relates to devops](#) (2018).  
URL <https://www.safaribooksonline.com/library/view/how-sre-relates/9781492030645/>
- [142] R. Pike, [Concurrency is not parallelism](#), waza conference (2013).  
URL <https://blog.golang.org/waza-talk>
- [143] C. A. R. Hoare, [Communicating sequential processes](#), Commun. ACM 21 (8) (1978) 666–677. doi:10.1145/359576.359585.  
URL <https://doi.org/10.1145/359576.359585>
- [144] C. Hewitt, [Actor model for discretionary, adaptive concurrency](#), CoRR abs/1008.1459 (2010). arXiv:1008.1459.  
URL <http://arxiv.org/abs/1008.1459>

## BIBLIOGRAPHY

---

- [145] A. C. Baktir, A. Ozgovde, C. Ersoy, [How can edge computing benefit from software-defined networking: A survey, use cases, and future directions](#), *IEEE Commun. Surv. Tutorials* 19 (4) (2017) 2359–2391. doi:[10.1109/COMST.2017.2717482](#).  
URL <https://doi.org/10.1109/COMST.2017.2717482>
- [146] J. Sherwin, C. J. Sreenan, [Software-defined networking for data centre network management: A survey](#), *CoRR* abs/2106.10014 (2021). [arXiv:2106.10014](#).  
URL <https://arxiv.org/abs/2106.10014>
- [147] R. Ciobanu, C. Negru, F. Pop, C. Dobre, C. X. Mavromoustakis, G. Mastorakis, [Drop computing: Ad-hoc dynamic collaborative computing](#), *Future Gener. Comput. Syst.* 92 (2019) 889–899. doi:[10.1016/j.future.2017.11.044](#).  
URL <https://doi.org/10.1016/j.future.2017.11.044>
- [148] Y. Shao, C. Li, Z. Fu, L. Jia, Y. Luo, [Cost-effective replication management and scheduling in edge computing](#), *J. Netw. Comput. Appl.* 129 (2019) 46–61. doi:[10.1016/j.jnca.2019.01.001](#).  
URL <https://doi.org/10.1016/j.jnca.2019.01.001>
- [149] A. Verma, L. Pedrosa, M. Korupolu, D. Oppenheimer, E. Tune, J. Wilkes, [Large-scale cluster management at google with borg](#), in: L. Réveillère, T. Harris, M. Herlihy (Eds.), *Proceedings of the Tenth European Conference on Computer Systems, EuroSys 2015, Bordeaux, France, April 21-24, 2015*, ACM, 2015, pp. 18:1–18:17. doi:[10.1145/2741948.2741964](#).  
URL <https://doi.org/10.1145/2741948.2741964>
- [150] F. Rossi, V. Cardellini, F. L. Presti, M. Nardelli, [Geo-distributed efficient deployment of containers with kubernetes](#), *Comput. Commun.* 159 (2020) 161–174. doi:[10.1016/j.comcom.2020.04.061](#).  
URL <https://doi.org/10.1016/j.comcom.2020.04.061>

- [151] A. Lèbre, J. Pastor, A. Simonet, F. Desprez, [Revising openstack to operate fog/edge computing infrastructures](#), in: 2017 IEEE International Conference on Cloud Engineering, IC2E 2017, Vancouver, BC, Canada, April 4-7, 2017, IEEE Computer Society, 2017, pp. 138–148. doi:[10.1109/IC2E.2017.35](#). URL <https://doi.org/10.1109/IC2E.2017.35>
- [152] H. Sami, A. Mourad, [Dynamic on-demand fog formation offering on-the-fly iot service deployment](#), IEEE Trans. Netw. Serv. Manag. 17 (2) (2020) 1026–1039. doi:[10.1109/TNSM.2019.2963643](#). URL <https://doi.org/10.1109/TNSM.2019.2963643>
- [153] A. Kurniawan, [Learning AWS IoT: Effectively manage connected devices on the AWS cloud using services such as AWS Greengrass, AWS button, predictive analytics and machine learning](#), Packt Publishing, 2018. URL <https://books.google.rs/books?id=7NRJDwAAQBAJ>
- [154] Linux Foundation, KubeEdge, <https://kubedge.io/> (accessed November 7, 2020).
- [155] General Electric, GE. Predix, <https://www.ge.com/digital/iiot-platform/> (accessed November 7, 2020).
- [156] Y. Mao, J. Zhang, K. B. Letaief, [Dynamic computation offloading for mobile-edge computing with energy harvesting devices](#), IEEE J. Sel. Areas Commun. 34 (12) (2016) 3590–3605. doi:[10.1109/JSAC.2016.2611964](#). URL <https://doi.org/10.1109/JSAC.2016.2611964>
- [157] C. Shi, K. Habak, P. Pandurangan, M. H. Ammar, M. Naik, E. W. Zegura, [COSMOS: computation offloading as a service for mobile devices](#), in: J. Wu, X. Cheng, X. Li, S. Sarkar (Eds.), The Fifteenth ACM International Symposium on Mobile Ad Hoc Networking and Computing, MobiHoc’14, Philadelphia, PA,

## BIBLIOGRAPHY

---

- USA, August 11-14, 2014, ACM, 2014, pp. 287–296. doi:[10.1145/2632951.2632958](https://doi.org/10.1145/2632951.2632958).  
URL <https://doi.org/10.1145/2632951.2632958>
- [158] Z. Wang, Z. Zhao, G. Min, X. Huang, Q. Ni, R. Wang, [User mobility aware task assignment for mobile edge computing](#), *Future Gener. Comput. Syst.* 85 (2018) 1–8. doi:[10.1016/j.future.2018.02.014](https://doi.org/10.1016/j.future.2018.02.014).  
URL <https://doi.org/10.1016/j.future.2018.02.014>
- [159] B. Yee, D. Sehr, G. Dardyk, J. B. Chen, R. Muth, T. Ormandy, S. Okasaka, N. Narula, N. Fullagar, [Native client: a sandbox for portable, untrusted x86 native code](#), *Commun. ACM* 53 (1) (2010) 91–99. doi:[10.1145/1629175.1629203](https://doi.org/10.1145/1629175.1629203).  
URL <https://doi.org/10.1145/1629175.1629203>
- [160] M. Beck, M. Werner, S. Feld, T. Schimper, [Mobile edge computing: A taxonomy](#), in: *The Sixth International Conference on Advances in Future Internet (AFIN 2014)*, 2014.  
URL [https://www.researchgate.net/publication/267448582\\_Mobile\\_Edge\\_Computing\\_A\\_Taxonomy](https://www.researchgate.net/publication/267448582_Mobile_Edge_Computing_A_Taxonomy)
- [161] R.-A. Cherrueau, M. Delavergne, A. Lebre, [Geo-Distribute Cloud Applications at the Edge](#), in: *EURO-PAR 2021 - 27th International European Conference on Parallel and Distributed Computing*, Lisbon, Portugal, 2021, pp. 1–14.  
URL <https://hal.inria.fr/hal-03212421>
- [162] S. Miloš, S. Goran, Z. Miroslav, M. Branko, [Infrastructure as software in micro clouds at the edge](#), *Sensors* 21 (21) (2021). doi:[10.3390/s21217001](https://doi.org/10.3390/s21217001).  
URL <https://www.mdpi.com/1424-8220/21/21/7001>
- [163] M. Bertier, F. Desprez, G. Fedak, A. Lebre, A.-C. Orgerie, J. Pastor, F. Quesnel, J. Rouzaud-Cornabas, C. Tedeschi, [Beyond The](#)

- Clouds, How Should Next Generation Utility Computing Infrastructures Be Designed?, in: Z. Mahmood (Ed.), Cloud Computing: Challenges, Limitations and R and D Solutions, Springer, 2014.  
URL <https://hal.inria.fr/hal-01067888>
- [164] N. Hayashibara, X. Défago, R. Yared, T. Katayama, [The  \$\Phi\$  accrual failure detector](#), in: 23rd International Symposium on Reliable Distributed Systems (SRDS 2004), 18-20 October 2004, Florianopolis, Brazil, IEEE Computer Society, 2004, pp. 66–78. [doi:10.1109/RELDIS.2004.1353004](https://doi.org/10.1109/RELDIS.2004.1353004).  
URL <https://doi.org/10.1109/RELDIS.2004.1353004>
- [165] B. Satzger, A. Pietzowski, W. Trumler, T. Ungerer, [Variations and evaluations of an adaptive accrual failure detector to enable self-healing properties in distributed systems](#), in: P. Lukowicz, L. Thiele, G. Tröster (Eds.), Architecture of Computing Systems - ARCS 2007, 20th International Conference, Zurich, Switzerland, March 12-15, 2007, Proceedings, Vol. 4415 of Lecture Notes in Computer Science, Springer, 2007, pp. 171–184. [doi:10.1007/978-3-540-71270-1\\_13](https://doi.org/10.1007/978-3-540-71270-1_13).  
URL [https://doi.org/10.1007/978-3-540-71270-1\\_13](https://doi.org/10.1007/978-3-540-71270-1_13)
- [166] M. Simic, M. Stojkov, G. Sladic, B. Milosavljević, Crdts as replication strategy in large-scale edge distributed system: An overview, in: CRDTs as replication strategy in large-scale edge distributed system: An overview, 2020.
- [167] A. Farshin, A. Roozbeh, G. Q. M. Jr., D. Kostic, [Make the most out of last level cache in intel processors](#), in: G. Candea, R. van Renesse, C. Fetzer (Eds.), Proceedings of the Fourteenth EuroSys Conference 2019, Dresden, Germany, March 25-28, 2019, ACM, 2019, pp. 8:1–8:17. [doi:10.1145/3302424.3303977](https://doi.org/10.1145/3302424.3303977).  
URL <https://doi.org/10.1145/3302424.3303977>



## BIBLIOGRAPHY

---

- [168] D. Gannon, R. S. Barga, N. Sundaresan, [Cloud-native applications](#), IEEE Cloud Comput. 4 (5) (2017) 16–21. doi: [10.1109/MCC.2017.4250939](#).  
URL <https://doi.org/10.1109/MCC.2017.4250939>
- [169] J. J. M. M. Rutten, [Behavioural differential equations: a coinductive calculus of streams, automata, and power series](#), Theor. Comput. Sci. 308 (1-3) (2003) 1–53. doi: [10.1016/S0304-3975\(02\)00895-2](#).  
URL [https://doi.org/10.1016/S0304-3975\(02\)00895-2](https://doi.org/10.1016/S0304-3975(02)00895-2)
- [170] A. S. Abdelhamid, A. R. Mahmood, A. Daghistani, W. G. Aref, [Prompt: Dynamic data-partitioning for distributed micro-batch stream processing systems](#), in: D. Maier, R. Pottinger, A. Doan, W. Tan, A. Alawini, H. Q. Ngo (Eds.), Proceedings of the 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14-19, 2020, ACM, 2020, pp. 2455–2469. doi: [10.1145/3318464.3389713](#).  
URL <https://doi.org/10.1145/3318464.3389713>
- [171] J. hung Ding, C. jung Lin, P. hao Chang, C. hao Tsang, W. chung Hsu, Y. ching Chung, Armvisor: System virtualization for arm, in: In Proceedings of the Ottawa Linux Symposium (OLS, 2012, pp. 93–107.
- [172] M. Simić, M. Stojkov, G. Sladic, B. Milosavljević, M. Zarić, On container usability in large-scale edge distributed system, in: On container usability in large-scale edge distributed system, 2019.
- [173] S. Kuenzer, V. Badoiu, H. Lefeuvre, S. Santhanam, A. Jung, G. Gain, C. Soldani, C. Lupu, S. Teodorescu, C. Raducanu, C. Banu, L. Mathy, R. Deaconescu, C. Raiciu, F. Huici, [Unikraft: Fast, specialized unikernels the easy way](#), CoRR abs/2104.12721 (2021). arXiv: [2104.12721](#).  
URL <https://arxiv.org/abs/2104.12721>

- [174] Y. Yu, P. Manolios, L. Lamport, [Model checking tla<sup>+</sup> specifications](#), in: L. Pierre, T. Kropf (Eds.), *Correct Hardware Design and Verification Methods*, 10th IFIP WG10.5 Advanced Research Working Conference, CHARME '99, Bad Herrenalb, Germany, September 27-29, 1999, *Proceedings*, Vol. 1703 of *Lecture Notes in Computer Science*, Springer, 1999, pp. 54–66. doi:[10.1007/3-540-48153-2\\_6](#).  
URL [https://doi.org/10.1007/3-540-48153-2\\_6](https://doi.org/10.1007/3-540-48153-2_6)
- [175] M. Hennessy, *A distributed Pi-calculus*, Cambridge University Press, 2007.
- [176] J. Upadhyay, [Review of: Distributed computing through combinatorial topology by maurice herlihy and dmitry kozlov and sergio rajsbaum](#), *SIGACT News* 47 (2) (2016) 18–20. doi:[10.1145/2951860.2951867](#).  
URL <https://doi.org/10.1145/2951860.2951867>
- [177] T. H. Cormen, C. E. Leiserson, R. L. Rivest, C. Stein, [Introduction to Algorithms](#), 3rd Edition, MIT Press, 2009.  
URL <http://mitpress.mit.edu/books/introduction-algorithms>
- [178] S. Tarkoma, C. E. Rothenberg, E. Lagerspetz, [Theory and practice of bloom filters for distributed systems](#), *IEEE Commun. Surv. Tutorials* 14 (1) (2012) 131–155. doi:[10.1109/SURV.2011.031611.00024](#).  
URL <https://doi.org/10.1109/SURV.2011.031611.00024>
- [179] B. Burns, D. Oppenheimer, [Design patterns for container-based distributed systems](#), in: A. Clements, T. Condie (Eds.), 8th *USENIX Workshop on Hot Topics in Cloud Computing*, HotCloud 2016, Denver, CO, USA, June 20-21, 2016, *USENIX Association*, 2016.  
URL <https://www.usenix.org/conference/hotcloud16/workshop-program/presentation/burns>

## BIBLIOGRAPHY

---

- [180] M. A. Tamiru, J. Tordsson, E. Elmroth, G. Pierre, [An experimental evaluation of the kubernetes cluster autoscaler in the cloud](#), in: 12th IEEE International Conference on Cloud Computing Technology and Science, CloudCom 2020, Bangkok, Thailand, December 14-17, 2020, IEEE, 2020, pp. 17–24. doi: [10.1109/CloudCom49646.2020.00002](https://doi.org/10.1109/CloudCom49646.2020.00002).  
URL <https://doi.org/10.1109/CloudCom49646.2020.00002>
- [181] E. Cohen, E. Halperin, H. Kaplan, [Performance aspects of distributed caches using ttl-based consistency](#), Theor. Comput. Sci. 331 (1) (2005) 73–96. doi: [10.1016/j.tcs.2004.09.033](https://doi.org/10.1016/j.tcs.2004.09.033).  
URL <https://doi.org/10.1016/j.tcs.2004.09.033>
- [182] C. G. Gray, D. R. Cheriton, [Leases: An efficient fault-tolerant mechanism for distributed file cache consistency](#), in: G. R. Andrews (Ed.), Proceedings of the Twelfth ACM Symposium on Operating System Principles, SOSP 1989, The Wigwam, Litchfield Park, Arizona, USA, December 3-6, 1989, ACM, 1989, pp. 202–210. doi: [10.1145/74850.74870](https://doi.org/10.1145/74850.74870).  
URL <https://doi.org/10.1145/74850.74870>
- [183] K. Mathews, C. Kramer, R. Gotzhein, [Token bucket based traffic shaping and monitoring for wlan-based control systems](#), in: 28th IEEE Annual International Symposium on Personal, Indoor, and Mobile Radio Communications, PIMRC 2017, Montreal, QC, Canada, October 8-13, 2017, IEEE, 2017, pp. 1–7. doi: [10.1109/PIMRC.2017.8292201](https://doi.org/10.1109/PIMRC.2017.8292201).  
URL <https://doi.org/10.1109/PIMRC.2017.8292201>
- [184] Kubernetes.io, Running in multiple zones, <https://kubernetes.io/docs/setup/best-practices/multiple-zones/> (accessed November 7, 2020).
- [185] P. Andrade, T. Bell, J. Eldik, G. Mccance, B. Panzer-Steindel, M. Santos, S. and, U. Schwickerath, Review of cern data centre

- infrastructure, Journal of Physics Conference Series 396 (2012) 2002–. [doi:10.1088/1742-6596/396/4/042002](https://doi.org/10.1088/1742-6596/396/4/042002).
- [186] A. Forestiero, C. Mastroianni, M. Meo, G. Papuzzo, M. Sheikhalishahi, [Hierarchical approach for green workload management in distributed data centers](#), in: Euro-Par 2014: Parallel Processing Workshops - Euro-Par 2014 International Workshops, Porto, Portugal, August 25-26, 2014, Revised Selected Papers, Part I, Vol. 8805 of Lecture Notes in Computer Science, Springer, 2014, pp. 323–334. [doi:10.1007/978-3-319-14325-5\\_28](https://doi.org/10.1007/978-3-319-14325-5_28).  
URL [https://doi.org/10.1007/978-3-319-14325-5\\_28](https://doi.org/10.1007/978-3-319-14325-5_28)
- [187] M. Al-Khafajiy, T. Baker, C. Chalmers, M. Asim, H. Kolivand, M. Fahim, A. Waraich, [Remote health monitoring of elderly through wearable sensors](#), Multim. Tools Appl. 78 (17) (2019) 24681–24706. [doi:10.1007/s11042-018-7134-7](https://doi.org/10.1007/s11042-018-7134-7).  
URL <https://doi.org/10.1007/s11042-018-7134-7>
- [188] Y. J. Jeon, S. J. Kang, [Wearable sleepcare kit: Analysis and prevention of sleep apnea symptoms in real-time](#), IEEE Access 7 (2019) 60634–60649. [doi:10.1109/ACCESS.2019.2913849](https://doi.org/10.1109/ACCESS.2019.2913849).  
URL <https://doi.org/10.1109/ACCESS.2019.2913849>
- [189] G. Chiarini, P. Ray, S. Akter, C. Masella, A. Ganz, [mhealth technologies for chronic diseases and elders: A systematic review](#), IEEE J. Sel. Areas Commun. 31 (9-Supplement) (2013) 6–18. [doi:10.1109/JSAC.2013.SUP.0513001](https://doi.org/10.1109/JSAC.2013.SUP.0513001).  
URL <https://doi.org/10.1109/JSAC.2013.SUP.0513001>
- [190] F. Fitzek, S. Li, S. Speidel, T. Strufe, M. Simsek, M. Reisslein, [Tactile Internet: With Human-in-the-Loop](#), Elsevier Science, 2021.  
URL <https://books.google.rs/books?id=es4hEAAAQBAJ>
- [191] A. A. Omar, M. Z. A. Bhuiyan, A. Basu, S. Kiyomoto, M. S. Rahman, [Privacy-friendly platform for healthcare data in cloud](#)

## BIBLIOGRAPHY

---

- based on blockchain environment, *Future Gener. Comput. Syst.* 95 (2019) 511–521. doi:[10.1016/j.future.2018.12.044](https://doi.org/10.1016/j.future.2018.12.044).  
URL <https://doi.org/10.1016/j.future.2018.12.044>
- [192] M. Simic, G. Sladic, B. Milosavljević, A case study iot and blockchain powered healthcare, in: *A Case Study IoT and Blockchain powered Healthcare*, 2017.
- [193] Y. Shen, D. Guo, F. Long, L. A. Mateos, H. Ding, Z. Xiu, R. B. Hellman, A. King, S. Chen, C. Zhang, H. Tan, *Robots under COVID-19 pandemic: A comprehensive survey*, *IEEE Access* 9 (2021) 1590–1615. doi:[10.1109/ACCESS.2020.3045792](https://doi.org/10.1109/ACCESS.2020.3045792).  
URL <https://doi.org/10.1109/ACCESS.2020.3045792>
- [194] V. Ivančević, M. Knežević, M. Simić, I. Luković, D. Mandić, Public healthcare and epidemiology with dr warehouse, *International Journal on Advances in Software* 6 (3-4) (2013) 329–342.
- [195] E. Baccarelli, P. G. V. Naranjo, M. Scarpiniti, M. Shojafar, J. H. Abawajy, *Fog of everything: Energy-efficient networked computing architectures, research challenges, and a case study*, *IEEE Access* 5 (2017) 9882–9910. doi:[10.1109/ACCESS.2017.2702013](https://doi.org/10.1109/ACCESS.2017.2702013).  
URL <https://doi.org/10.1109/ACCESS.2017.2702013>

## BIBLIOGRAPHY

---

# Biography

The work in this thesis is synthesis of few individual parts:

- (1) The experience acquired on a university on the topic of software engineering,
- (2) The research conducted as part of the PhD studies, covering various aspects of the distributed systems,
- (3) The work done in collaboration with prominent software vendors,
- (4) The collaboration with researchers from different research areas.

Miloš Simić is a Ph.D. student and teaching assistant within the Department of Computing and Control, Faculty of Technical Sciences, University of Novi Sad since 2015. He received his B.Sc. degree in 2014, and M.Sc. degree in 2015, all in Computer Science from the University of Novi Sad, Faculty of Technical Sciences. He is owner of two team awards: **(1)** Best paper award (*academia*), and **(2)** ThinkX in the category Community and Social Impact (*industry*).

Over the years, Miloš worked with various prominent software vendors, in different fields. This allowed him combining the different skillsets developed over the years, to focus his expertise towards designing and implementing distributed and non-distributed software systems, for various usages. His research interests include: **(1)** distributed systems, **(2)** (multi) cloud computing, **(3)** edge computing, **(4)** big data and

(5) service oriented architectures and microservices.

As part his Ph.D., Miloš have studied the different distributed systems techniques, combined with various software engineering methodologies and practices, covering both standard-defined processes and industry-proven methods, to resolve and answer such complicated questions that are part of this thesis. Working with different software vendors, combined with traditional academic research, helped Miloš to clear his PhD vision, and guid him to the work that is described in this thesis.

Trough colaboration with people from different research areas, this theis gain forml description and formal model that is important leverage, to describe and validate such complicated system that is described in this thesis. In 2021, the work published in the *IEEE Access* journal (paper (cf. journal paper (1))), was presented to the eminent professors and colleagues from the *Imperial College London*, as an invited lecture.



# List of publications

This thesis is the result of years of research and development, and it is based on the previously presented papers at conferences in journals.

- Journal publications

- (1) M. Simić, I. Prokić, J. Dedeić, G. Sladić and B. Milosavljević, "Towards Edge Computing as a Service: Dynamic Formation of the Micro Data-Centers," in IEEE Access, vol. 9, pp. 114468-114484, 2021, doi: 10.1109/ACCESS.2021.3104475.
- (2) Simić, M.; Sladić, G.; Zarić, M.; Markoski, B. Infrastructure as Software in Micro Clouds at the Edge. Sensors 2021, 21, 7001. <https://doi.org/10.3390/s21217001>
- (3) Ivančević V., Knežević M., Simić M., Luković I., Mandić D.: Public Healthcare and Epidemiology with Dr Warehouse, International Journal on Advances in Software, 2013, Vol. 6, No. 3-4, pp. 329-342, ISSN 1942-2628

- Conference papers

- (1) Simić, M., Stojkov, M., Sladić, G., Milosavljević, B. CRDTs as replication strategy in large-scale edge distributed system: An overview. In: Zdravković, M., Konjović, Z., Trajanović, M. (Eds.) ICIST 2020 Proceedings Vol.1, pp.46-50, 2020, ISBN 978-86-85525-24-7.

- (2) Simić M., Stojkov M., Sladić G., Milosavljević B., Zarić M.: On container usability in large-scale edge distributed system, 9. International Conference on Information Science and Technology (ICIST), Kopaonik: Society for Information Systems and Computer Networks, 10-13 March, 2019, pp. 97-101, ISBN pp.97-101, 2019.
- (3) Simić M., Stojkov M., Sladić G., Milosavljević B.: Edge computing system for large-scale distributed sensing systems, 8. International Conference on Information Science and Technology (ICIST), Kopaonik: Society for Information Systems and Computer Networks, 11-14 March, 2018, pp. 36-39, ISBN 978-86-85525-22-3.
- (4) Simić M., Sladić G., Milosavljević B.: A Case Study IoT and Blockchain Powered Healthcare , 8. PSU-UNS International Conference on Engineering and Technology - ICET, Novi Sad: University of Novi Sad, Faculty of Technical Sciences , 8-10 June, 2017, pp. 1-4, ISBN 978-86-7892-934-2.
- (5) Simić M.: Clover: Property Graph based metadata management service, 6. International Conference on Information Science and Technology (ICIST), Kopaonik: Society for Information Systems and Computer Networks, 29-2 February, 2016, pp. 81-85, ISBN 978-86-85525-18-6.
- (6) Ivančević V., Knežević M., Simić M., Luković I., Mandić D.: Public Healthcare and Epidemiology with Dr Warehouse, International Journal on Advances in Software, 2013, Vol. 6, No. 3-4, pp. 329-342, ISSN 1942-2628

*Овај Образац чини саставни део докторске дисертације, односно докторског уметничког пројекта који се брани на Универзитету у Новом Саду. Попуњен Образац укоричити иза текста докторске дисертације, односно докторског уметничког пројекта.*

## План третмана података

<b>Назив пројекта/истраживања</b>
Динамичко формирање дистрибуираног микро окружења у рачунарству у облаку / Dynamic Formation of the Distributed Micro Clouds
<b>Назив институције/институција у оквиру којих се спроводи истраживање</b>
а) Факултет техничких наука, Универзитет у Новом Саду
<b>Назив програма у оквиру ког се реализује истраживање</b>
Рачунарство и аутоматика – докторска дисертација
<b>1. Опис података</b>
<p>1.1 Врста студије</p> <p><i>Укратко описати тип студије у оквиру које се подаци прикупљају</i></p> <p><b>Докторска дисертација</b></p> <p>1.2 Врсте података</p> <p>а) квантитативни</p> <p>б) <b>квалитативни</b></p> <p>1.3. Начин прикупљања података</p> <p>а) анкете, упитници, тестови</p> <p>б) клиничке процене, медицински записи, електронски здравствени записи</p> <p>в) генотипови: навести врсту _____</p> <p>г) административни подаци: навести врсту _____</p> <p>д) узорци ткива: навести врсту _____</p> <p>ђ) снимци, фотографије: навести врсту _____</p> <p>е) текст, навести врсту <u><b>Актуелна литература у области истраживања</b></u></p> <p>ж) мапа, навести врсту _____</p> <p>з) остало: описати _____</p>

### 1.3 Формат података, употребљене скале, количина података

#### 1.3.1 Употребљени софтвер и формат датотеке:

- a) Excel фајл, датотека \_\_\_\_\_
- b) SPSS фајл, датотека \_\_\_\_\_
- c) PDF фајл, датотека \_\_\_\_\_
- d) Текст фајл, датотека \_\_\_\_\_
- e) JPG фајл, датотека \_\_\_\_\_
- f) Остало, датотека \_\_\_\_\_

#### 1.3.2. Број записа (код квантитативних података)

- a) број варијабли \_\_\_\_\_
- б) број мерења (испитаника, процена, снимака и сл.) \_\_\_\_\_

#### 1.3.3. Поновљена мерења

a) да

**б) не**

Уколико је одговор да, одговорити на следећа питања:

- a) временски размак између поновљених мера је \_\_\_\_\_
- б) варијабле које се више пута мере односе се на \_\_\_\_\_
- в) нове верзије фајлова који садрже поновљена мерења су именоване као \_\_\_\_\_

Напомене: \_\_\_\_\_

*Да ли формати и софтвер омогућавају дељење и дугорочну валидност података?*

**a) Да**

**б) Не**

*Ако је одговор не, образложити* \_\_\_\_\_

\_\_\_\_\_

## 2. Прикупљање података

### 2.1 Методологија за прикупљање/генерисање података

#### 2.1.1. У оквиру ког истраживачког нацрта су подаци прикупљени?

- a) експеримент, навести тип \_\_\_\_\_
- б) корелационо истраживање, навести тип \_\_\_\_\_
- ц) анализа текста, навести тип Анализа доступне литературе** \_\_\_\_\_

д) остало, навести шта \_\_\_\_\_

2.1.2 Навести врсте мерних инструмената или стандарде података специфичних за одређену научну дисциплину (ако постоје).

---

---

## 2.2 Квалитет података и стандарди

### 2.2.1. Третман недостајућих података

а) Да ли матрица садржи недостајуће податке? Да **Не**

Ако је одговор да, одговорити на следећа питања:

а) Колики је број недостајућих података? \_\_\_\_\_

б) Да ли се кориснику матрице препоручује замена недостајућих података? Да Не

в) Ако је одговор да, навести сугестије за третман замене недостајућих података

---

### 2.2.2. На који начин је контролисан квалитет података? Описати

---

---

### 2.2.3. На који начин је извршена контрола уноса података у матрицу?

---

---

## 3. Третман података и пратећа документација

### 3.1. Третман и чување података

3.1.1. Подаци ће бити депоновани у \_\_\_\_\_ репозиторијум.

3.1.2. URL адреса \_\_\_\_\_

3.1.3. DOI \_\_\_\_\_

3.1.4. Да ли ће подаци бити у отвореном приступу?

а) Да

б) Да, али после ембарга који ће трајати до \_\_\_\_\_

в) Не

Ако је одговор не, навести разлог \_\_\_\_\_

3.1.5. Подаци неће бити депоновани у репозиторијум, али ће бити чувани.

Образложење

---

### 3.2 Метаподаци и документација података

3.2.1. Који стандард за метаподатке ће бити примењен? \_\_\_\_\_

3.2.1. Навести метаподатке на основу којих су подаци депоновани у репозиторијум.

\_\_\_\_\_  
\_\_\_\_\_  
*Ако је потребно, навести методе које се користе за преузимање података, аналитичке и процедуралне информације, њихово кодирање, детаљне описе варијабли, записа итд.*

### 3.3 Стратегија и стандарди за чување података

3.3.1. До ког периода ће подаци бити чувани у репозиторијуму? \_\_\_\_\_

3.3.2. Да ли ће подаци бити депоновани под шифром? Да Не

3.3.3. Да ли ће шифра бити доступна одређеном кругу истраживача? Да Не

3.3.4. Да ли се подаци морају уклонити из отвореног приступа после извесног времена?

Да Не

Образложити

## 4. Безбедност података и заштита поверљивих информација

Овај одељак МОРА бити попуњен ако ваши подаци укључују личне податке који се односе на учеснике у истраживању. За друга истраживања треба такође размотрити заштиту и сигурност података.

### 4.1 Формални стандарди за сигурност информација/података

Истраживачи који спроводе испитивања с људима морају да се придржавају Закона о заштити података о личности ([https://www.paragraf.rs/propisi/zakon\\_o\\_zastiti\\_podataka\\_o\\_licnosti.html](https://www.paragraf.rs/propisi/zakon_o_zastiti_podataka_o_licnosti.html)) и одговарајућег институционалног кодекса о академском интегритету.

4.1.2. Да ли је истраживање одобрено од стране етичке комисије? Да **Не**

Ако је одговор Да, навести датум и назив етичке комисије која је одобрила истраживање

\_\_\_\_\_  
4.1.2. Да ли подаци укључују личне податке учесника у истраживању? Да **Не**

Ако је одговор да, наведите на који начин сте осигурали поверљивост и сигурност информација везаних за испитанике:

- а) Подаци нису у отвореном приступу
- б) Подаци су анонимизирани
- ц) Остало, навести шта

## 5. Доступност података

### 5.1. Подаци ће бити

#### а) јавно доступни

б) доступни само уском кругу истраживача у одређеној научној области

ц) затворени

Ако су подаци доступни само уском кругу истраживача, навести под којим условима могу да их користе:

Ако су подаци доступни само уском кругу истраживача, навести на који начин могу приступити подацима:

### 5.4. Навести лиценцу под којом ће прикупљени подаци бити архивирани.

## 6. Улоге и одговорност

### 6.1. Навести име и презиме и мејл адресу власника (аутора) података

Милош Симић, [milos.simic@uns.ac.rs](mailto:milos.simic@uns.ac.rs)

### 6.2. Навести име и презиме и мејл адресу особе која одржава матрицу с подацима

### 6.3. Навести име и презиме и мејл адресу особе која омогућује приступ подацима другим истраживачима

