# A NETWORK SCIENCE BASED APPROACH FOR OPTIMAL MICROSERVICE GOVERNANCE

2020-021

| Student Name | Student ID |
|---|---|
| Leader: Saranga S.A. G | IT17016230 |
| Member 2: De Silva N. | IT17006880 |
| Member 3: L.S. Jayasinghe | IT17012966 |
| Member 4: M.V. Lakshitha | IT17410250 |

**Supervisor**
Dr. Dharshana Kasthurirathna

# Introduction

- **Why Micro-Services ?**

- **Current trends and practices in Micro-Services**

- **Why Kubernetes ?**

- **Problem definition**

# Problems

Constant monitoring of metrics through APM tools. E.g.-Resource utilization process[1-2]

Difficult to understand why a problem occurs even though there is knowledge that a problem has occurred[5]

Have to make use of multiple monitoring sources to make effective decisions[3-5]

In short the research problem that our research aims to fulfil can simply be described as follows

"In deploying Microservices through Kubernetes, there is no efficient and effective way for developers to evaluate and monitor the effectiveness and viability of a microservice deployment and identify possible performance bottlenecks. Furthermore, developers are not able to optimize their deployment such that they can make the optimal use of their deployed microservices in the cluster."

# Research Questions

Can current orchestration tools successfully monitor relationships among microservices ?

How effective is the proposed optimization model compared to current optimization models ?

What is the impact of pod load prediction in determining an optimal deployment strategy ?

What is the importance and the impact of identifying the resiliency of microservices in deployments ?

# Main Objective

To model a network science-based approach to govern microservice deployments through evaluation and analysis of metrics gathered, and ultimately produce a proposed model which aids to optimize microservice deployments.

To increase the efficiency of microservices deployments by applying the metrics used in network analysis, such as centrality and resilience measures and link predictions on identified dependency measurements.

To develop an improved auto-scaling policy for a deployment, based on load prediction

To development of a business intelligence dashboard to evaluate performance and monitor microservice deployments.

To identify key factors which lead to performance reduction in microservice deployments and come up with an optimal deployment strategy.

Sub Objectives

# Requirements

## Functional Requirements

- User should be able to view optimal deployment strategy for a given cluster

- User should be able to view possible deployment suggestions in order to optimize the deployment

- System should analyze the resiliency of the identified microservices in the cluster

- Load prediction should be performed on the identified microservices

- Load predictions should be made based for a maximum of 24hrs ahead

- Auto scaler should be configured to proactively auto scale based on predicted load

- Optimal deployment strategy for a particular deployment should be determined using load prediction, dependency measurements and resiliency evaluation

# Requirements

**Non Functional Requirements**

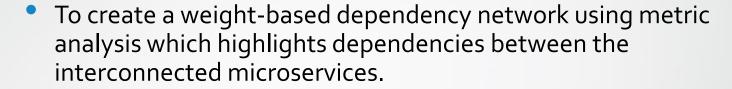- Usability
- Reliability
- Availability
- Interoperability
- Performance

# Expected Outcomes

- To create a weight-based dependency network using metric analysis which highlights dependencies between the interconnected microservices.

- To create an improved auto-scaling policy integrating load-prediction analysis.

- To improve the performance of microservice deployments through resiliency analysis.

- To increase the performance in microservice deployments through the combination of the above-mentioned approaches.

# Commercialization

Initial plan is develop model into an open source tool.

In the future, with the addition of newer features, we plan on making this tool a freemium tool.
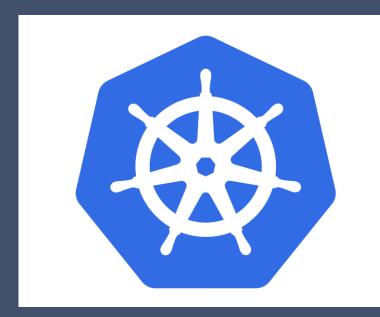
# A network science based approach for the generation of dependency map and a service mesh based on the network traffic and user behavior on a Kubernetes cluster

Saranga S.A.G
IT17016230

# Introduction

- **How Kubernetes works ?**

- **What is a service mesh ?**
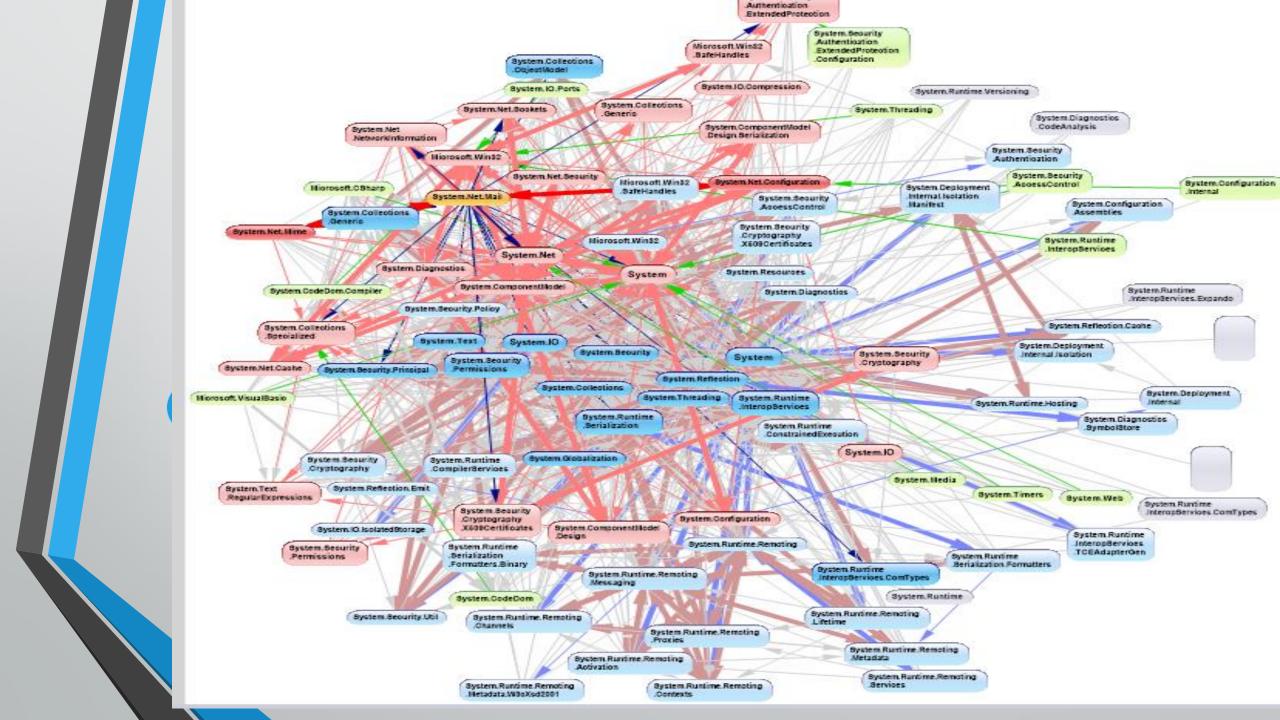
- **Service Gateway**

- **Problem definition**

# Related Research
# and
# Research Gap

- [1]Istio: Modernize digital applications with microservices management using the istio service mesh

- [4 Service Dependency Based Dynamic Load Balancing Algorithm for Container Clusters

- Deploying Microservice Based Applications with Kubernetes: Experiments and Lessons Learned

# A data science based approach for an improved auto-scaling policy in Kubernetes based on load prediction

De Silva N.

IT17006880

# Introduction

- **What is autoscaling ?**

- **Current Autoscaling tools by Kubernetes**

- **What is Kubernetes HPA and how it works ?**

- **Problems in the current implementation of HPA**

- **Why we focused about issues in HPA ?**

# Problems faced in the current implementation of HPA

I. Response delay caused by the time required for Pod initialization[1,7]

II. Response delay causes increases the user's request response time, resulting drop in service quality [1,7]

III. Ineffective scaling caused due to the under allocation of resources in the configuration process[1,7]

# Related Research

- [1] Uses a combination of ARIMA model and Empirical Mode Decomposition (EMD) in order to predict resource usage

  - Highlights current issues such as the response delay in Kubernetes

  - However, the use of EMD model has some drawbacks explained in [2]

- [3] Proposes a resource prediction algorithm called CRUPA

  - The CRUPA algorithm is based on a time series analysis model (ARIMA) combined with docker container techniques.

  - However this approach does integrate with existing auto-scaling tools provided by Kubernetes as well as does not make use of machine learning tools in order to perform the prediction process.

- [4-6] compares various models for short term load forecasting and their prediction accuracy

# Objective

- To develop an improved auto-scaling policy for a deployment, based on load prediction

# Research Questions

- What is the impact of load prediction in coming up with an optimal deployment strategy ?

- How accurately can we predict utilization values ?

- What is the impact of load prediction on the identified microservices in the overall performance of the cluster ?

# Steps in the development of load prediction algorithm

Obtain the list microservices that show high level of dependency based on the proposed dependency map

Identify the relevant pods in which the identified microservices are deployed

Obtain the CPU utilization metrics of the deployed pods

Convert the metrics obtained into a time series

# Steps in the development of load prediction algorithm cont.



Plot time series to identify the relevant time series components

Convert time series to a stationary dataset

Develop algorithm using prediction model to forecast future CPU utilization values for given pods

Forecasted CPU utilization values will be forwarded to optimization algorithm and used to configure auto-scaling of the cluster

# Expected Outcome

To create an improved auto-scaling policy integrating load-prediction analysis which solves the current issues found in the HPA in Kubernetes

# A network science based approach for evaluation of resiliency among microservices through the use of Chaos Engineering

Lakshitha M.V
IT17410250

# Introduction

- What is resilience?

- What are the ways to analyze resilience?

- What is Chaos Engineering?

- What are the tools used in Chaos Engineering?

- Importance of analyzing resilience

Currently the dependency network is not considered when measuring the importance of a service in the deployment

# Related Work

- [1]Gremlin: Systematic Resilience Testing of Microservices
    - how to use Gremlin to express common failure scenarios and how developers of an enterprise application were able to discover previously unknown bugs in their failure-handling code without modifying the application.

- [4]Chaos Monkey: Increasing SDN Reliability through Systematic Network Destruction

# Objectives

- Evaluate the resilience of microservices using Chaos Engineering to provide data to generate an optimal deployment strategy

# Research questions

- What are the impact of microservices on each other?

- How does the dependency network affect on the deployment of microservices?

- How to evaluate resiliency of microservices without interrupting customers?

# Chaos Engineering

Chaos Monkey - randomly kills a microservice

Chaos Gorilla - kills entire availability zone

Chaos Kong - kill whole reigon

Latency Monkey - Introduces communication delays to simulate degradation or outages in a network.

Doctor Monkey - Performs health checks, by monitoring performance metrics such as CPU load to detect unhealthy instances

Failure Injection Testing - designed to give developers a "blast radius" rather than unmanaged chaos

LATENCY MONKEY

other frameworks

http://techblog.netflix.com/2014/10/fit-failure-injection-testing.html
http://www.infoq.com/presentations/failure-as-a-service-netflix

Chaos Gorilla

NETFLIX                                              @atseitlin

CHAOS KONG

because regions fail

http://techblog.netflix.com/2015/09/chaos-engineering-upgraded.html

Steps in designing a Chaos Experiment

Pick a hypothesis

Choose the scope of the experiment

Identify the metrics you're going to watch

Analyze the results

Run the experiment

Notify the organization

Increase the scope

Automate

# Stimulating failovers to determine the highly dependent nodes



[7]How Netflix does failovers in 7 minutes flat

# Expected Outcome

- To generate a detailed evaluation report on the resiliency of the microservices

# A network science based approach for generating optimal deployment strategy by using dependency, load predictions and resiliency measures among microservices in Kubernetes.

L.S Jayasinghe
IT17016230

# Availability and Network latency

- When we deploy microservices, we deploy them in many nodes. because we need to increase availability.

- In the Kubernetes cluster, we have master node and salve nodes, these nodes can be located in the multiple places, may be another country.

- This problem causes delays between nodes. This is called network latency[2,4].

- **That is affected to overall response time[1]**

- E.g.: if network latency is high. Overall response time will be high.

# Network Latency between different regions (Azure Cloud)



- Lowest latency in December:
Australia central1 region – Australia central2 region(1ms)[5]

- Highest latency in December:
Australia central-South Africa(392ms)[5]

# How to avoid from high network latency

All the Microservices can be deployed in one Big Node.

Micro-services with higher communication rate can be deployed in nearby nodes.

# One Big Node

- If all the microservices are deployed in one big Node, we can decrease network latency.

- But then the application will become monolithic. Therefore availability can be reduced and it is one of the main features of Micro-Service architecture

# Decrease Overall network latency (scenario 1)

- If the couple of nodes are located in very far, the communication delay time is very low, also if couple of node are located near, communication is very fast.

- Best for multi-region problem

# Scenario 2

- Microservice A and C communicate more than microservice A and B.

# Related work

There is one research found on optimal deployment strategy[8], it is only concerned about ,

- resource power
- replication
- dependency map

Our research,

- Load prediction
- Resiliency

# Objective

- To identify key factors which lead to performance reduction in microservice deployments and come up with an optimal deployment strategy.

# Proposed Solution

To solved these previous scenario we develop an optimal deployment algorithm. It ensures the system has optimal availability and low network latency. It controls replication of the microservices and increases overall system performance.

To make this optimization algorithm we should be concerned about,

- Resource power
- Dependency map
- Future loads
- resiliency
- Network latency
- Node resource power
- Replication of microservices

# Methodology

GET THE METRIC VALUES FROM DEPENDENCY MAP, LOAD PREDICTIONS, RESILIENCY

ANALYSIS THE INFORMATION USING OPTIMAL DEPLOYMENT ALGORITHM[6]

CREATE A NODE MAP(NETWORK LATENCY MAP:SOCKPERF (LINUX))

SHOW NEXT 24 HOURS DEPLOYMENT MAP AND FUTURE SUGGESTION IN THE UI

# Path for the creation of optimal deployment strategy...

- Multi object optimization algorithm

- Identify shortest path between node(Dijkstra algorithm)[7]

- Rank all the metrics

- Identify correct deployment place using the algorithm

# Expected Outcomes

- To generate an optimal deployment strategy for a given cluster based on load prediction, dependency analysis and resilience evaluation.

- Next 24 hours deployment map

# Automated Deployment

YAML

Questions ?

# Thank you!

# References - main

- [1]"DataDog/the-monitor", GitHub, 2020. [Online]. Available: https://github.com/DataDog/the-monitor/blob/master/kubernetes/monitoring-kubernetes-performance-metrics.md

- [2]"Resource Usage Monitoring in Kubernetes", Kubernetes.io, 2020. [Online]. Available: https://kubernetes.io/blog/2015/05/resource-usage-monitoring-kubernetes/

- [3]"Monitoring Kubernetes (Part 1): The Challenges and Data Sources - DZone Performance", dzone.com, 2020. [Online]. Available: https://dzone.com/articles/monitoring-kubernetes-part-1-the-challenges-and-da.

- [4]"Challenges of Monitoring and Troubleshooting in Kubernetes Environments | Sumo Logic", Sumo Logic, 2020. [Online]. Available: https://www.sumologic.com/blog/troubleshooting-kubernetes/.

- [5]D. Snyder, "How to Overcome Kubernetes Monitoring Challenges", OverOps Blog, 2020. [Online]. Available: https://blog.overops.com/how-to-overcome-monitoring-challenges-with-kubernetes/.

- [6]M. Vizard, "Running Kubernetes at Scale Top 2020 Challenge - Container Journal", Container Journal, 2020. [Online]. Available: https://containerjournal.com/topics/container-management/running-kubernetes-at-scale-top-2020-challenge.

# References - main

- [7]"Kubernetes: The Challenge of Deploying & Maintaining", Techolution, 2020. [Online]. Available: https://techolution.com/kubernetes-challenges/.

- [8]"Kubernetes Challenges for Enterprises", Network Computing, 2020. [Online]. Available: https://www.networkcomputing.com/data-centers/kubernetes-challenges-enterprises.

- [9]"Monitoring Kubernetes (Part 1): The Challenges and Data Sources - DZone Performance", dzone.com, 2020. [Online]. Available: https://dzone.com/articles/monitoring-kubernetes-part-1-the-challenges-and-da.

- [10]Heinrich, Robert & van Hoorn, André & Knoche, Holger & Li, Fei & Lwakatare, Lucy Ellen & Pahl, Claus & Schulte, Stefan & Wettinger, Johannes. (2017). Performance Engineering for Microservices: Research Challenges and Directions.

- [11]Fazio, Maria & Celesti, Antonio & Ranjan, R. & Liu, Chang & Chen, Lydia & Villari, Massimo. (2016). Open Issues in Scheduling Microservices in the Cloud. IEEE Cloud Computing.

# References-Gihan

- [1] J. Zhang, R. Ren, C. Huang, X. Fei, W. Qun and H. Cai, "Service Dependency Based Dynamic Load Balancing Algorithm for Container Clusters," *2018 IEEE 15th International Conference on e-Business Engineering (ICEBE)*, Xi'an, 2018, pp. 70-77. doi: 10.1109/ICEBE.2018.00021 http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8592632&isnumber=8592611.

- [2] M. Gawel and K. Zielinski, "Analysis and Evaluation of Kubernetes Based NFV Management and Orchestration," *2019 IEEE 12th International Conference on Cloud Computing (CLOUD)*, Milan, Italy, 2019, pp. 511-513. http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8814580&isnumber=8814488

- [3] B. Carter, "Growing Software Applications by Incremental Development of Heterogeneous Micro-Applications Using Cellular Regeneration Concepts," *2015 Annual Global Online Conference on Information and Computer Technology (GOCICT)*, Louisville, KY, 2015, pp. 1-5. http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=7545087&isnumber=7545077

- [4] Preeth E N, F. J. P. Mulerickal, B. Paul and Y. Sastri, "Evaluation of Docker containers based on hardware utilization," *2015 International Conference on Control Communication & Computing India (ICCC)*, Trivandrum, 2015, pp. 697-700. L: http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=7432984&isnumber=7432856.

# References-Nishitha

- [1] Zhao, A., Huang, Q., Huang, Y., Zou, L., Chen, Z., & Song, J. "*Research on Resource Prediction Model Based on Kubernetes Container Auto-scaling Technology." IOP Conference Series: Materials Science and Engineering,2019.*

- [2]Yamin Wang, L. Wu and Shouxiang Wang, "Challenges in applying the empirical mode decomposition based hybrid algorithm for forecasting renewable wind/solar in practical cases," *2016 IEEE Power and Energy Society General Meeting (PESGM)*, Boston, MA, 2016, pp. 1-5

- [3]Y. Meng, R. Rao, X. Zhang and P. Hong, "CRUPA: A container resource utilization prediction algorithm for auto-scaling based on time series analysis," *2016 International Conference on Progress in Informatics and Computing (PIC)*, Shanghai, 2016, pp. 468-472

- [4] Nguyen, H., & Hansen, C. K. "*Short-term electricity load forecasting with Time Series Analysis." 2017 IEEE International Conference on Prognostics and Health Management (ICPHM).*

- [5] Bozkurt, Ö. Ö., Biricik, G., & Tayşi, Z. C. "*Artificial neural network and SARIMA based models for power load forecasting in Turkish electricity market".* 2017.

- [6] Vagropoulos, S. I., Chouliaras, G. I., Kardakos, E. G., Simoglou, C. K., & Bakirtzis, A. G. " *Comparison of SARIMAX, SARIMA, modified SARIMA and ANN-based models for short-term PV generation forecasting." 2016 IEEE International Energy Conference (ENERGYCON).*

- [7]Kim, Won-Yong & Lee, Jin-Seop & Huh, Eui-Nam. (2017). Study on proactive auto scaling for instance through the prediction of network traffic on the container environment.

# References-Lakshitha

- [1]"Systematic Resilience Testing of Microservices", 2016 IEEE 36th International Conference on Distributed Computing Systems. Available: https://ieeexplore.ieee.org/abstract/document/7536505[Accessed: 27- Jan- 2020]

- [2]"Chaos Engineering" Basiri, A., Behnam, N., de Rooij, R., Hochstein, L., Kosewski, L., Reynolds, J., & Rosenthal, C. (2016). . IEEE Software, 33(3), 35-Available: https://ieeexplore.ieee.org/abstract/document/7436642[Accessed: 28- Jan- 2020]

- [3]"Optimal and Automated Deployment for Microservices",FASE 2019: Fundamental Approaches to Software Engineering pp 351-368 Available:https://link.springer.com/chapter/10.1007/978-3-030-16722-6_21

- [4]"Chaos Monkey: Increasing SDN Reliability through Systematic Network Destruction"  Available : https://dl.acm.org/doi/pdf/10.1145/2785956.2790038 [Accessed: 25- Jan- 2020]

- [5]"Improving Kubernetes Resiliency with Chaos Engineering" Available:https://medium.com/faun/failures-are-inevitable-even-a-strongest-platform-with-concrete-operations-infrastructure-can-7d0c01643oc6 [Accessed: 20-Jan-2020]

- [6]"Chaos Engineering by Casey Rosenthal, Lorin Hochstein, Aaron Blohowiak, Nora Jones, and Ali Basiri Available: http://channyblog.s3-ap-northeast-2.amazonaws.com/data/channy/2018/01/18023151/chaos-engineering  [Accessed: 23-Jan-2020]

- [7]"How Netflix does failovers in 7 minutes flat" Available:https://opensource.com/article/18/4/how-netflix-does-failovers-7-minutes-flat  [Accessed: 23-Jan-2020]

# References - Sanjaya

- [1]"Analysis of HTTP Performance Problems", W3.org, 2020. [Online]. Available: https://www.w3.org/Protocols/HTTP-NG/http-prob. [Accessed: 27- Jan- 2020].

- [2]"What is cloud service latency? - Definition from WhatIs.com", WhatIs.com, 2020. [Online]. Available: https://whatis.techtarget.com/definition/cloud-service-latency. [Accessed: 27- Jan- 2020]

- [3]"What Is Latency and How to Reduce It - KeyCDN Support", KeyCDN, 2020. [Online]. Available: https://www.keycdn.com/support/what-is-latency. [Accessed: 27- Jan- 2020].

- [4]"Test Azure virtual machine network latency in an Azure virtual network", Docs.microsoft.com, 2020. [Online]. Available: https://docs.microsoft.com/en-us/azure/virtual-network/virtual-network-test-latency. [Accessed: 27- Jan- 2020].

- [5]2020. [Online]. Available: https://docs.microsoft.com/en-us/azure/networking/azure-network-latency https://searchcloudcomputing.techtarget.com/feature/Reducing-network-latency-means-focusing-on-location-location-location. [Accessed: 27- Jan- 2020].

- [6]"Multi-objective optimization", En.wikipedia.org, 2020. [Online]. Available: https://en.wikipedia.org/wiki/Multi-objective_optimization. [Accessed: 27- Jan- 2020].

- [7]"Dijkstra's algorithm", En.wikipedia.org, 2020. [Online]. Available: https://en.wikipedia.org/wiki/Dijkstra%27s_algorithm. [Accessed: 27- Jan- 2020].

- [8]M. Bravetti, S. Giallorenzo, J. Mauro, I. Talevi and G. Zavattaro, "Optimal and Automated Deployment for Microservices", Fundamental Approaches to Software Engineering, pp. 351-368, 2019. Available: 10.1007/978-3-030-16722-6_21 [Accessed 27 January 2020].