5G Network service resilience estimation under short time scale traffic variation

> Chair Risk and Resilience of Complex Systems Annual Scientific Seminar

Presenter: Rui Li (starting date: 02/11/2020) Supervisors: Bertrand Decocq (Orange Innovation) Anne Barros, Yiping Fang, Zhiguo Zeng (CentraleSupélec) 17 November 2022 CentraleSupélec Université

orange

Fondation CentraleSupélec





Main challenges for 5G system

5G system and resilience





5G system and resilience





Resilience: maintain an acceptable level of service facing various incidents.

5G system and resilience



5G networks are facing risks from RAN, TN, CN. Different vertexical industry may require different SLA.



5G and Resilience(continued)



5G system and resilience

5G networks are facing risks from RAN, TN, CN. Different vertical industry may require different SLA.



Different risks of a 5G network

5G system and resilience

□ External risks:

- <u>Traffic change</u>, attacks...
- □ Failure on different layer:



- Damage of switches, servers or physical links or failure on virtual switch, VNF, virtual link...
- □ Failure propagation
 - Direct: from lower level to a higher level; Indirect: e.g. a rerouting action may congest other entities





8

Traffic change, one of the most frequent risks

5G system and resilience







Attack



Simultaneous work (maintenance/upgrade)

[1] Xu F, Li Y, Wang H, Zhang P, Jin D. Understanding Mobile Traffic Patterns of Large Scale Cellular Towers in Urban Environment. *IEEE/ACM Transactions* on *Networking*. 2017; 25,(2):1147-1161. doi: 10.1109/TNET.2016.2623950.

[2] Pintér G, Felde I. Analyzing the Behavior and Financial Status of Soccer Fans from a Mobile Phone Network Perspective: Euro 2016, a Case Study. Information. 2021; 12(11):468. doi:10.3390/info12110468



Petri Net-based model



5G network is a dynamic system. It has many dynamic behaviors. Scalability: change the capacity according to current load. Healing: repair automatically the failed component (virtual layer). Mobility: connection (path, anchor point...) evolves with time.



Network modeling based on Petri Net





5-tuple: (P, T, F, W, M_0)

- P is a finite set of places
- T is a finite set of transitions



- F is a finite set of arcs with $F \subseteq (P \times T) \cup (T \times P)$ connecting places with transitions
- − *W* is a multiset of arcs $(P \times T) \cup (T \times P) \rightarrow \mathbb{N}$, it assigns the weight
- $-M_0$ is the initial marking of the Petri net graph

Classical Petri Net is not directly applicable to telecommunication systems.

Some extensions:

- Stochastic Petri Nets --failure process
- Colored Petri Nets --token selection
- Timed Petri Nets --delay



- We consider a multi-tuple, Timed Stochastic Colored Petri Net, TSCPN=(P, T, F, W, M₀, C, E, I, R, D)

Service delivery model Petri Net

Petri Net-based model



Service is delivered by following an acyclic service function chain consisting of Virtual Network Functions (VNFs); Each VNF is containerized and contains subfunctions in form of containers/pods.



Petri Net-based model





Congestion--queueing model (continued)

Petri Net-based model



rejected



Queueing model





14

Congestion--auto-scaling

Petri Net-based model

MS4 instances To cope with network congestion: *Reduce traffic*: 5G network level management cuts the traffic *Increase capacity*: in 5G MANO, Kubernetes launches auto-scaling Kubernetes, the container Scaling out Scaling in management system, collects Ó Scaling in Scaling out HPA MS4 metrics (e.g., CPU usage...) values MS3 instance Resource HPA MS1 intermittently. • After analyzing these metrics, HPA MS3 Scaling out Kubernetes knows if the system ΉPA MS2 MS1 instances Scaling in Ó (in level of VNF components) is Scaling in Scaling out Micro-service 1 overloaded or idle. Then it Available instances: 3 pods decides to create new instances Available resources: 22% or remove the existing ones. Kubernetes HPA notices that the micro-service is going to overloaded MS2 instances



Kubernetes Auto-scaling out

Kubernetes, the container management system, collects metrics (CPU usage, memory usage ...) values with a certain frequency.

After analyzing these metrics, Kubernetes knows if the system (in level of VNF components) is overloaded or too idle. Then based on this, it decides to create new instances or remove the existing ones. Desired CPU usage := 50%



15 The suto-scaling process modeled by Petri Net

Threshold UP:= 60% Threshold_DOWN := 30% # Current CPU usage is 78 % If Metrics > Threshold UP: New_scale := Current CPU usage / Desired **CPU usage** [0.78/0.5*3 = 4.67 pods → 5 pods] Deploy(New_scale) # The network instances will be doubled

Petri Net-based model



Kubernetes Auto-scaling in

Kubernetes, the container management system, collects metrics (CPU usage, memory usage ...) values with a certain frequency.

After analyzing these metrics, Kubernetes knows if the system (in level of VNF components) is overloaded or too idle. Then based on this, it decides to create new instances or remove the existing ones. Desired CPU usage := 50%



The suto-scaling process modeled by Petri Net

Threshold UP:= 60% Threshold_DOWN := 30% # Current CPU usage is 22 % Elseif Metrics < Threshold DOWN: New_scale := Current CPU usage / Desired **CPU usage** $[0.22/0.5*3 = 1.33 \text{ pods} \rightarrow 2 \text{ pods}]$ Deploy(New_scale) # The network instances will be halved

Petri Net-based model





Network service & Simulation

Latency

10 ms

50 ms

requirement



18

Resilience: maintain an acceptable level of service facing various incidents.

Reliability^[1] in the context of network layer packet transmissions: percentage value of the packets successfully delivered to a given system entity within the time constraint required by the targeted service out of all the packets transmitted. *Not rejected Within delay limit*

Network service & Simulation



[1]Third Generation Partnership Project (3GPP), "Management and orchestration; 5G performance measurements (Release 16)," 3GPP TS 22.261 version 16.14.1 Release 16, July 2022.



Network service & Simulation



From 0 to 18 seconds: preparation/ Launch PDU sessions arrival rate: 150 packet/PDU·s for service 1 or 75 packets/PDU·s for service 2

From 18 to 60 seconds: Traffic change with different patterns

Entropy LV = 0.0108 Entropy SV = 0.0207 Entropy S1 = 0.1019 Entropy S2 = 0.3676



Network traffic variation (continued)

3 strategies for scaling

Nothing

• Wait until traffic variation stops

Autoscaling – intermittence = 5 s

- Take actions based on the traffic load
- Check every 5 seconds

Autoscaling – intermittence = 5 s, windows = 15 s

- Take actions based on the decision made during last 15 seconds
- Check every 5 seconds

Scaling	Scaling	Scaling	Scaling
out	out	in	out
T-15	T-10	T-5	

Network service & Simulation

Network traffic variation – traffic pattern 1



Simulation result of traffic pattern 1

Strategy		Average latency (ms)	Resilience loss (second)	Resource cost (CPU·s)
Long varia	ation			
	Serv. 1	$10.437(\pm 0.014)$	$24.082(\pm 0.300)$	12000
NO AS	Serv. 2	18.863 ± 0.014)	$1.460(\pm 0.044)$	
Denie AC	Serv. 1	$9.742(\pm 0.003)$	$1.476(\pm 0.037)$	20067(170)
Basic AS	Serv. 2	$18.195(\pm 0.005)$	$0.073(\pm 0.005)$	$20067(\pm 70)$
Win AS	Serv. 1	$9.889(\pm 0.007)$	$6.584(\pm 0.141)$	$17024(\pm 42)$
wiii. AS	Serv. 2	$18.337(\pm 0.008)$	$0.445(\pm 0.026)$	$1/934(\pm 42)$

No scaling

6.002



Network service & Simulation

Window based auto-scaling



Long term traffic interruption

Network service & Simulation

Service 1 Reliability evolution

Service 2 Reliability evolution



Network traffic variation – traffic pattern 2



Simulation result of traffic pattern 2

Strategy		Average latency (ms)	Resilience loss (second)	Resource cost (CPU·s)
Short vari	ation			
No AS	Serv. 1	$9.786(\pm 0.004)$	$2.867(\pm 0.070)$	12000
NO AS	Serv. 2	$18.244(\pm 0.006)$	$0.172(\pm 0.011)$	12000
Basic AS	Serv. 1	$9.738(\pm 0.003)$	$1.404(\pm 0.033)$	$15364(\pm 55)$
Dasic AS	Serv. 2	$18.194(\pm 0.005)$	$0.070(\pm 0.005)$	15504(±55)
Win AS	Serv. 1	$9.785(\pm 0.004)$	$2.830(\pm 0.070)$	12275(±16)
wiii. AS	Serv. 2	$18.240(\pm 0.005)$	$0.164(\pm 0.010)$	$12273(\pm 10)$



Network service & Simulation

Basic auto-scaling



Window based auto-scaling



Short term traffic interruption

Network service & Simulation

Service 1 Reliability evolution



Service 2 Reliability evolution

Network traffic variation – traffic pattern 3



Simulation result of traffic pattern 3

Strategy		Average latency (ms)	Resilience loss (second)	Resource cost (CPU·s)
Sinusoidal	1			
No AS	Serv. 1	$9.832(\pm 0.004)$	$4.640(\pm 0.092)$	12000
NO AS	Serv. 2	$18.281(\pm 0.006)$	$0.247(\pm 0.011)$	12000
Dagia AS	Serv. 1	$9.793(\pm 0.004)$	$3.676(\pm 0.092)$	18220(174)
Dasic AS	Serv. 2	$18.252(\pm 0.005)$	$0.310(\pm 0.015)$	18320(±74)
Win AC	Serv. 1	$9.780(\pm 0.004)$	$3.394(\pm 0.088)$	12275(1.29)
wiii. AS	Serv. 2	$18.233(\pm 0.005)$	$0.182(\pm 0.010)$	$15575(\pm 28)$

No scaling Processing Transmission Waiting 0.0075 0.0050 Basic auto-scaling 8.016 Processing Transmission 8.014 Waiting 0.012 0.010 0.005 0.000 0.004 0.012

Network service & Simulation

Window based auto-scaling





Network service & Simulation

• • • •

Service 1 Reliability evolution



Service 2 Reliability evolution

Network traffic variation – traffic pattern 4



Network service & Simulation



Sinusoidal traffic interruption 2

Network service & Simulation

Service 1 Reliability evolution



Service 2 Reliability evolution





Conclusions

- A Petri Net-based model is constructed to represent the dynamics of 5G networks;
- Different management mechanisms are compared;
- 5G network service performance is evaluated from a resilience point of view.

Prospects

- Al-based management can be more usefull when the time series Entropy increases;
- More precise parameters to be collected from Orange experts / equipment suppliers;
- We will not stop in the data plane packet transfer but also a complete 5G architecture with signaling part.





Thank You

orange[™]

Orange Innovation

Rui Li rui.li@orange.com Bertrand Decocq bertrand.decocq@orange.com

5

CentraleSupélec





CentraleSupélec Université Paris-Saclay

Anne Barros anne.barros@centralesupelec.fr Yiping Fang yiping.fang@centralesupelec.fr Zhiguo Zeng zhiguo.zeng@centralesupelec.fr