# Availability Prediction of Telecommunication Application Servers Deployed on Cloud

Attila Hilt[1*], Gábor Járó[1], István Bakos[2]

## Abstract

*Availability and reliability considerations are discussed in this paper with special focus on 'cloud based' Mobile Switching and Telecommunication Application Servers (MSS and TAS). Before the extensive deployment of cloud based telecommunication networks, the essential question shall be answered: Will cloud technology ensure the 'carrier grade' requirements that are well established and proven on 'legacy telecommunications' hardware? This paper shows the possible redundancy principles and a simulation method to predict availability for 'cloudified' mobile communication network elements. As a calculation example, Nokia AS on 'telco-cloud' is presented, that combines several redundancy principles such as full protection (2N), standby and load sharing.*

## Keywords

*Availability, reliability, redundancy principles, simulation, telco-cloud, core networks, mobile networks, Mobile Switching Server, Telecommunication Application Server*

[1] Product Architecture Group, Mobile Broadband
Nokia Networks,
H-1092 Budapest, Köztelek u. 6., Hungary

[2] Institute of Mathematics, Faculty of Natural Sciences
Budapest University of Technology and Economics
H-1111 Budapest, Hungary

* Corresponding author, e-mail: attila.hilt@nokia.com

## 1 Introduction

Recently, the accelerating demand for quicker launch of new mobile services as well as the continuous increase of both the number of subscribers and their traffic turned the interest towards Cloud technology. On one hand, information technology (IT) introduces standardized hardware (HW) scenario for the telecommunication applications. On the other hand, 'telco-cloud' offers possibilities for more flexible and economic resource allocations, e.g. scaling (Fig. 1). These benefits are widely investigated recently [1-3].
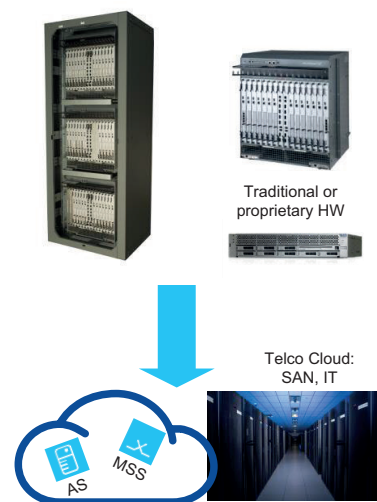


**Fig. 1** 'Cloudification' of telecommunication network elements

In a typical mobile network operator's (MNO) landscape [4] there is a wide variety of Network Elements (NE) based on different vendor specific hardware and software (SW) components as shown in Fig. 2. They are required to fulfill the various subscriber services [4]. Frequent introduction of new services as well as network (NW) maintenance e.g. performing SW upgrades or capacity expansions require careful preparation, planning and application specific HW.

In telecommunications networks the very strict service availability requirements are usually referred to as 'five nines' or A=99.999% availability. In practice, 'five nines' means maximum 316 seconds unplanned NE downtime in a year. Similarly,
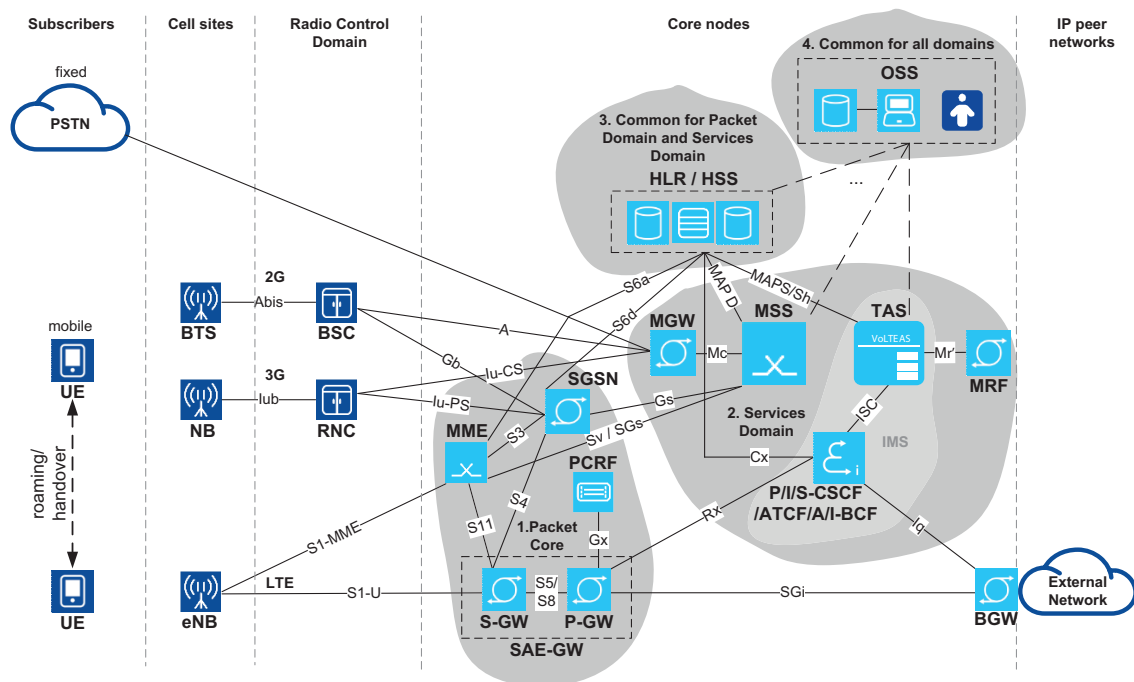
**Fig. 2** Core Network Elements (e.g. TAS and MSS) in a mobile network

'six nines' means maximum 32 seconds of outage over a year. Telecommunication HW (and SW) is specially designed to support these very strict requirements. On the other hand, IT HW components -even that of high quality- are not specially designed for telecom applications. Operators' target is to achieve at least as good capacity, performance and availability on the Cloud as it is provided nowadays on traditional, legacy or proprietary HW.

The rest of the paper is organized as follows. Part 2 summarizes availability and reliability definitions with special focus on their application in telecommunication networks. In Part 3 the basic redundancy principles and their calculations are shown. Part 4 discusses availability on network level. Finally, in Part 5 a calculation example is shown. The simulation results show that availability of cloud based telecommunication network elements can reach that of legacy ones deployed on traditional telecommunication HW.

## 2 Availability and Reliability definitions

Availability definitions and their explanations are summarized in this chapter. Please note that for some of the terms alternative definitions exist in the literature.

**Availability predictions** are mathematical calculations and models made for the different Network Elements to predict their availability performance. Predictions indicate the expected field performance only approximately. Prediction methods are useful when field data is not yet existing or very scarce. This is typically the case in the design phase of new products or when a new technology is introduced.

**Availability targets** in telecommunication networks are usually given in terms of "nines" as summarized in Table 1. Different levels of availability shall be distinguished:

- component level,
- Unit level,
- Network Element (NE) level,
- interconnection (e.g. IP backbone),
- Network (NW) level and
- service level availability.

**Table 1** Availability percentages and corresponding yearly Mean Down Time (MDT) values

| "nines" | A | yearly | |
| | % | MDT | unit |
| --- | --- | --- | --- |
| "2 nines" | 99.0 | 3.65 | days |
| "3 nines" | 99.9 | 8.76 | hours |
| "4 nines" | 99.99 | 52.56 | min |
| "5 nines" | 99.999 | 5.26 | min |
| "6 nines" | 99.9999 | 31.54 | sec |

Units are composed of components (combination of HW components and SW building blocks). Due to the strict availability requirements, troubleshooting cannot go down to component level. In case of fault, the entire unit is replaced as soon as possible, to minimize any possible outage time.

NEs are composed of units. NWs are composed of NEs and their interconnections. Service level availability has a wider sense than NW availability. Service should be granted to the subscribers in an end-to-end (e2e) manner when and where they would like to benefit the provided services.

A NW may operate with excellent availability, however where a geographical region is not covered, there the service

is not available. Another example is the interconnection of networks. E.g. when the subscriber's home NW is completely available but the visited network has an outage, then the subscriber cannot roam. On the contrary, when subscribers can make 2G or 3G calls, then they are not so sensitive for the lack or outage of 4G services in case of simple voice calls.

The very strict availability target required for telecommunication NEs is typically 5 or 6 nines, which is often approximated by (1).

$$A \cong \frac{MTBF}{MTBF + MTTR} \qquad (1)$$

Please note that in a well designed system, the outage of a unit does not automatically result in the outage or in the availability degradation of the entire network element. Similarly, the outage of a network element shall not automatically result in any availability degradation of the entire network.

On the contrary, proper NE and NW designs shall tolerate planned maintenance breaks. Planned maintenance breaks are used e.g. to check regularly, maintain or replace field replaceable units (FRU). Regular SW updates and upgrades are also preferably scheduled into planned maintenance windows.

**Estimation** is used instead of availability prediction when sufficient field data exists. Estimations correspond to actual measurement of failures. Estimated Mean Time Between Failures (MTBF) can be calculated based on the observation of similar NEs, usually after several field deployments. The longer time period is used for the observation and the larger population of similar NEs is observed the more accurate MTBF estimation can be reached.
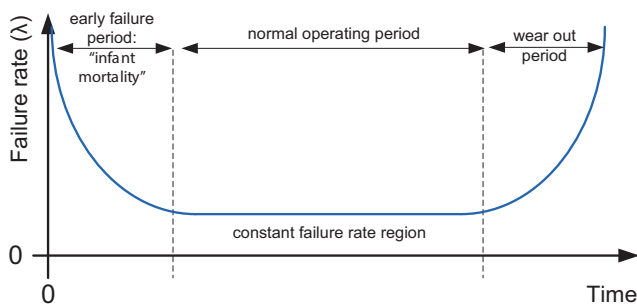


**Fig. 3** Typical failure rate curve over lifetime

**Failure and Fault:** A failure means any non-intended deviation of the system's behavior, defect or malfunction of HW and SW maintained. Fault may result in the loss of operational capabilities of the network element or the loss of redundancy in case of a redundant configuration.

**Failure Rate:** Failure Rate ($\lambda$) represents the number of failures likely to occur over a period of time. The failure rate of units (assembled from large number of components) is constant over the life expectancy. The constant failure rate period falls between the initial 'infant mortality" and the final "wear out"

phases of the lifetime. Figure 3 shows the well-known "bathtub" curve [5-7,14].

**FPMH and FIT:** Failure Rate is defined in units of Failures per Million Hours (FPMH) or failures per billion hours (FITs). If a unit has a failure rate of 1 FPMH, that unit is likely to fail once in one million hours. The failure rate is often measured in FITs (2):

$$FIT = \frac{number\ of\ failures}{10^9\ h} \qquad (2)$$

**MTBF:** Mean Time Between Failures (MTBF) is the expectation of the operating time duration between two consecutive failures of a repairable item. For field data MTBF is calculated as the total operating lifetime divided by the number of failures. MTBF is measured in hours or years. The following relationships apply between failure rate and MTBF:

$$MTBF(hours) = \frac{1}{\lambda} \qquad (3)$$

$$MTBF(years) = \frac{MTBF(hours)}{24 \cdot 365} = \frac{1}{\lambda \cdot 8760} \qquad (4)$$

**MTTF:** Mean Time to Failure is the mean proper operation time until the first failure (Fig. 4). Mainly used for the characterization of non-repairable or non-replaceable items, MTTF is a basic measure of reliability. As a statistical value, MTTF shall be preferably measured over a long period of time and with a large number of units.
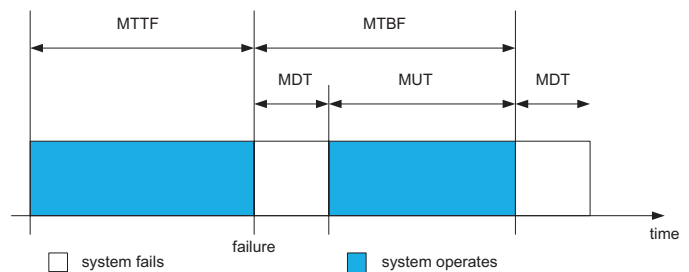


**Fig. 4** MTBF, MTTF, MDT and MUT on time scale

**MDT**: Mean Down Time (MDT) is the expectation of the time interval during which a unit or NE is in down state and cannot perform its function. MDT is the average time that a system is non-operational (Fig. 4). MDT includes all downtime associated with repair, corrective and preventive maintenance, self-imposed downtime, and any logistics or administrative delays. Please note that the down time of an individual unit (or more units) does not automatically result in down time of the entire NE element. Similarly, the down time of a single NE does not result in automatically a down time of the entire NW.

The addition of logistic delay times distinguishes MDT from MTTR, which includes only downtime specifically attributable to repairs. In order to minimize MDT, operators shall have

proper HW spare part management. In practice it means spare items of the different replaceable HW units stored in warehouse (in the amount according to the needs on NW level and calculated statistically). To reduce MDT, travel to the sites shall be also minimized. In practice this requires the possibility of remote NE access and remote SW management.

**MUT:** Mean Up Time is defined as the continuous operational time of the NE or the system without any down time (Fig. 4) [5]. MUT can be approximated with MTBF when MTTR is in the order of a few hours only and MTTF is in the order of several thousand hours. It is straightforward that system availability can be defined as MUT divided by the total operational time, the sum of MUT and MDT (5):

$$A = \frac{MUT}{MUT + MDT} \qquad (5)$$

**MTTR:** Mean Time to Repair (also known as Mean Time to Recovery) is the expectation of the time interval during which a unit or NE is down due to a failure that is under reparation. MTTR represents the average time required to repair a failed component or device. As seen in (1) MTTR affects availability. If it takes a long time to recover a system from a failure, the system will have a low availability. High availability can be achieved only if MTBF is very large compared to MTTR:

$$MTBF \gg MTTR \qquad (6)$$

The time that service persons take to acquire parts or modules, test equipment, and travel to the site is sometimes included in MTTR, but sometimes counted separately. MTTR generally does not include lead time for parts not readily available or other administrative and logistic downtimes. In our definition logistic delays shall be excluded from MTTR in order to achieve the required high availability.

**Reliability Block Diagram**: The Reliability Block Diagram (RBD) allows the graphical representation how the components of a system are reliability-wise connected. In most cases within a system, independence can be assumed across the components. Meaning, the failure of component A does not directly affect the failure of component B. Please note that the RBD shall not be equivalent or similar to the physical or logical setup or block diagram of the system.

It is worth to mention, that some parts of the full system are often omitted in the RBD. The parts that do not belong to the "functional" or "mission critical" ones of the system shall not be calculated in the availability figures. These parts can be for example displaying, statistical, logging or reporting subsystems, functions or units that help operators to supervise the entire system. Even though, the outages of these functions or units are inconvenient, they do not deteriorate the main function of the system it is designed for.

**Unavailability** is the complement of availability (7)-(9). It is the probability that the unit or NE cannot perform its function even though the required resources and normal operating conditions are provided.

$$U = 1 - A \qquad (7)$$

$$U = 1 - \frac{MUT}{MUT + MDT} = \frac{MDT}{MUT + MDT} \qquad (8)$$

$$U \cong 1 - \frac{MTBF}{MTBF + MTTR} = \frac{MTTR}{MTBF + MTTR} \qquad (9)$$

## 3 Availability and redundancy principles

A system composed of functional units in chain becomes unavailable, if any of the chained units fail (Fig. 5). Non-functional units are not considered as part of the chain, due to that fact that the unavailability of any non-functional unit does not deteriorate the desired function of the entire system. Thus the availability of end users' services is not affected.
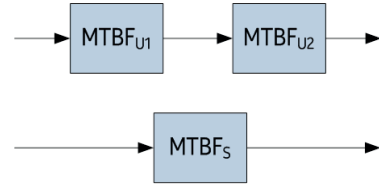


**Fig. 5** System MTBF of two units in series

Resulting availability (10), unavailability (11), system failure rate $\lambda_S$ (12) and $MTBF_S$ (13) of two units in series are written as:

$$A_S = A_{U1} \cdot A_{U2} = (1 - U_{U1}) \cdot (1 - U_{U2}) \qquad (10)$$

$$U_S = U_{U1} + U_{U2} - U_{U1} \cdot U_{U2} \qquad (11)$$

$$\lambda_S = \lambda_{U1} + \lambda_{U2} \qquad (12)$$

$$MTBF_S = \frac{1}{\lambda_S} = \frac{1}{\lambda_{U1} + \lambda_{U2}} = \frac{MTBF_{U1} \cdot MTBF_{U2}}{MTBF_{U1} + MTBF_{U2}} \qquad (13)$$

In case of two identical units in the chain the availability, unavailability, failure rate $\lambda_S$ and $MTBF_S$ are simplified to:

$$A_S = A_{U1} \cdot A_{U2}\big|_{U1=U2=U} = A_U^2 \qquad (14)$$

$$U_S = U_{U1} + U_{U2} - U_{U1} \cdot U_{U2}\big|_{U1=U2=U} = 2U_U - U_U^2 \qquad (15)$$

$$\lambda_S = \lambda_{U1} + \lambda_{U2}\big|_{U1=U2=U} = 2\lambda_U \qquad (16)$$

$$MTBF_S = \frac{1}{\lambda_S} = \frac{1}{2\lambda_U} = \frac{MTBF_U}{2} \qquad (17)$$

The system composed of several units in a chain becomes unavailable if any of the units (or any combination of two or more units) fails.
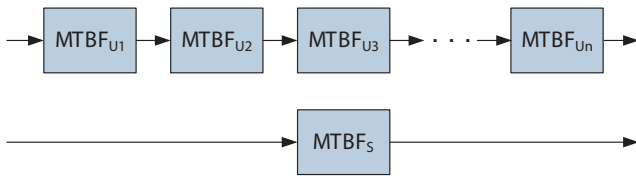


Fig. 6 Availability calculation and its model for n serial units

System failure rate $\lambda_S$ (18) and system $MTBF_S$ (19) for n units in chain (Fig. 6) are written as:

$$\lambda_S = \lambda_{U1} + \lambda_{U2} + \lambda_{U3} + \cdots + \lambda_{Un} = \sum_{i=1}^{n} \lambda_{Ui} \qquad (18)$$

$$MTBF_S = \frac{1}{\sum_{i=1}^{n} \lambda_{Ui}} = \frac{1}{\lambda_S} \qquad (19)$$

Supposing that all the n units are identical in the chain, the system failure rate $\lambda_S$ (18) and system $MTBF_S$ (19) formulas are simplified to:

$$\lambda_{U1} = \lambda_{U2} = \lambda_{U3} = \cdots = \lambda_{Un} = \lambda_U \qquad (20)$$

$$\lambda_S = n \cdot \lambda_U \qquad (21)$$

$$MTBF_S = \frac{1}{\lambda_S} = \frac{1}{n \cdot \lambda_U} = \frac{MTBF_U}{n} \qquad (22)$$

As it is seen in (22), the longer the chain composed of identical units is, the smaller the system $MTBF_S$ is. The system availability $A_S$ of n units forming a chain is:

$$A_S = A_{U1} \cdot A_{U2} \cdot A_{U3} \cdot \ldots \cdot A_{Un} \qquad (23)$$

Equation (23) shows the well-known fact that the less available unit within the chain determines the overall system availability. In case of n identical units (23) simplifies to:

$$A_S = A_{U1} \cdot A_{U2} \cdot A_{U3} \cdot \ldots \cdot A_{Un} \big|_{U1=U2=\ldots=Un=U} = A_U^{\,n} \quad (24)$$

Similarly to (7) the system unavailability is:

$$U_S = 1 - A_S = 1 - A_U^{\,n} \qquad (25)$$

The system unavailability $U_S$ of the entire chain can be approximated as the sum of the units' unavailability figures (upper bound). Equation (26) is valid in case of highly available units working in the chain (e.g. A=99.9% or better), where the products of the corresponding very small unavailability figures are falling into negligible orders of magnitude (e.g. $U_{Ui} \cdot U_{Uj}$ = 0.1%·0.1% = 0.001·0.001 ≈ 0):

$$U_S = U_{U1} + U_{U2} + U_{U3} + \cdots + U_{Un} - U_{U1} \cdot U_{U2} - \cdots$$
$$\leq \sum_{i=1}^{n} U_{Ui} \qquad (26)$$

Parallel protection of units significantly increases the availability of the network element (Fig. 7). It is very unlikely that both units fail the same time. Naturally, it is assumed that the failed unit is repaired or replaced as soon as possible to restore the back-up [6, 7].

Availability $A_P$ (27) and unavailability $U_P$ (28) of two units operating in parallel (supposing ideal switching between them in case of failure) is written as:

$$A_P = A_{U1} + A_{U2} - A_{U1} \cdot A_{U2} \qquad (27)$$

$$U_P = U_{U1} \cdot U_{U2} \big|_{U1=U2=U} = U_U^{\,2} \qquad (28)$$

$$\lambda_P = \lambda_{U1} \cdot \lambda_{U2} \cdot \left( MTTR_{U1} + MTTR_{U2} \right) \qquad (29)$$

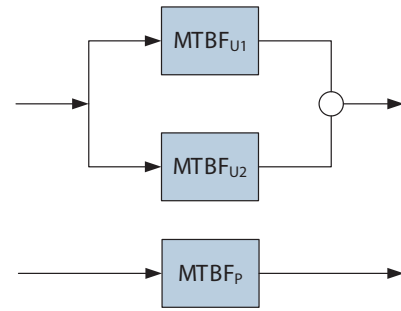$$MTBF_P = \frac{1}{\lambda_P} = \frac{MTBF_{U1} \cdot MTBF_{U2}}{MTTR_{U1} + MTTR_{U2}} \qquad (30)$$



Fig. 7 System MTBF of two units in parallel

The resulting availability for two uniform parallel units is:

$$A_P = 1 - U_U^{\,2} = 1 - \left(1 - A_U\right)^2 = A_U \left(2 - A_U\right) \qquad (31)$$

$$A_U < 1 \quad \rightarrow \quad 2 - A_U > 1 \quad \rightarrow \quad A_P > A_U \qquad (32)$$

As it is seen in Eqs. (31) and in (32), the combined availability $A_P$ of two parallel units is -in practice- always higher than the availability of the individual units. For two identical parallel units the system failure rate $\lambda_P$ and system $MTBF_P$ Eqs. (29), (30) are simply written as (33), (34) [8-10]:

$$\lambda_P = 2 \cdot MTTR_U \cdot \lambda_U^{\,2} \qquad (33)$$

$$MTBF_P = \frac{MTBF_U^{\,2}}{2 \cdot MTTR_U} \qquad (34)$$

The reliability block diagram of a complex system can be assembled from the basic (serial and parallel) reliability building blocks.

Figure 8 shows n units operating in parallel. Redundancy method is either N+X, where N working units are supported by X spare units or load sharing (LS).
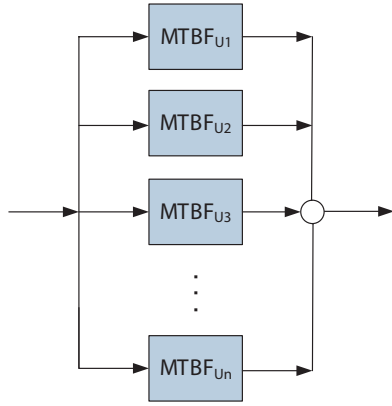
**Fig. 8** *n* units in parallel

Notation *m*/*n* is also often used, where the simultaneous outage of *m* units (out of the total number of *n* units forming the system) leads to a failure. Table 2 summarizes the different redundancy methods discussed.

**Table 2** Redundancy principles

| Notation | Explanation |
| --- | --- |
| N | *N* working units, no redundant unit, that means no protection. |
| n | total number of units, including working (*N*) and spare (*X*) units. *n=N+X* |
| m | number of simultaneously failed units (that may lead to the system failure) |
| 2N | full redundancy of N working units and N (hot or cold) spare units |
| X | number of spare units beside the working units |
| N+1 | N working units + one (hot or cold standby) spare unit. Only one unit may fail from the N working ones and the spare takes over its role. |
| N+X | N working units + X (hot or cold) spare units |
| SN+ | Load Sharing mode without any redundant unit: in case of any unit outage the remaining units carry as much traffic as they can handle. The group of load sharing units shall have enough spare capacity to bear a unit failure. |
| RN | Load Sharing in Recovery Groups (RG). One RG consists of several Functional Units that are dedicated to the same function. FUs cannot be allocated on the same physical resource (blade) if they belong to the same RG. But similar FUs belonging to different RGs can share the same blade. |
| Pooling | Grouping of several similar servers to increase network level availability. Pooling supports planned outages, geo-redundancy, etc. |

It is worth to mention that different combinations of the above protection methods are also possible. For example on Cloud, the HW blades may have N+1 protection while functional units can have either 2N, N+1 or RN depending on the function they provide. Affinity rules may ensure that 2N protected FUs are allocated onto different physical blades.

Overall Mean Time Between Failures $MTBF_P$ of n parallel units can be calculated according to (35), where m denotes the number of failed units [10].

$$MTBF_{P,m/n} = \frac{MTBF^m}{\dfrac{n!}{(n-m)!(m-1)!} \cdot MTTR^{m-1}} \qquad (35)$$

It is worth to mention that Eq. (35) in the simple case of *m* = *n* = 2 (full redundancy) gives back (34).

Figure 9 plots calculated $MTBF_P$ results as a function of n, the total number of units (sum of working and spare). Parameter of the curves is the increasing number of spare units (or NEs in a NW composed of parallel NEs).
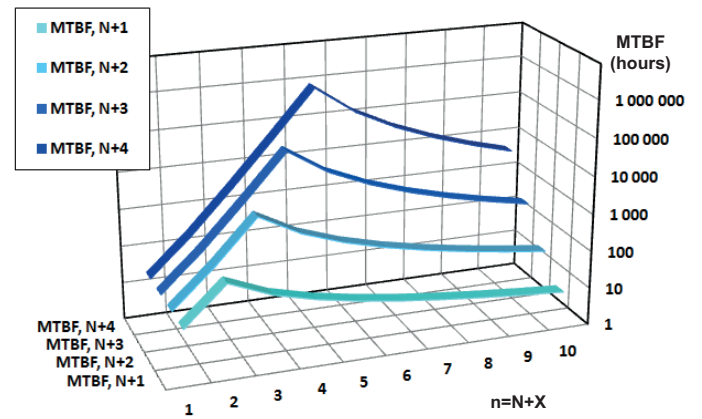


**Fig. 9** N+X redundancy, n=N+X=1,…,10

Please note that in the above models the interconnection between the units (represented by arrows in Fig. 5-8) have been assumed to be ideal. In a real network, however, this is not true, even though in most of the cases interconnection can be neglected as it has higher reliability than that of the NEs. The interconnections bring their own contribution to the network level availability that is discussed in the following part.

## 4 Network level availability

In a multinode network with only one spare, the system can tolerate the outage of only one single NE (Fig. 10). Outage may happen due to either interconnection (IP backbone composed of switches, routers, cables, optical fibers, connectors etc.) or NE failure. The simultaneous failure of two (or more) nodes would seriously overload or take down the network. To increase the overall network availability $A_{NW}$, sufficient amount of spares (NEs or units) shall be available.

Obviously, the higher the number *n* of the NEs is, the higher the number of the possible failures (36) is as it is discussed in [11].

$$U_{NW} = \frac{n \cdot (n-1)}{2} \cdot (1 - A_{NE})^2 = \frac{n \cdot (n-1)}{2} \cdot U_{NE}^2 \qquad (36)$$

**Table 3** Example of NE load values that tolerate any individual NE outage within the pool

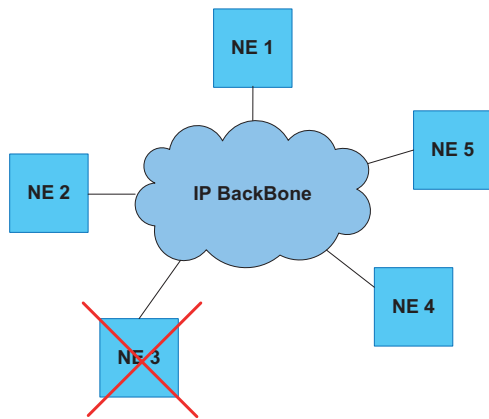| NE (server) | Weight Factor | Load [%] | with one server out of use (fault or planned maintenance) | | | |
|---|---|---|---|---|---|---|
| | | | NE1 out | NE2 out | NE3 out | NE4 out |
| NE1 | 100 | 100/250 = 40 % | out of use 0 % | 100/(100+50+30) = 56 % | 100/(100+70+30) = 50 % | 100/(100+70+50) = 45 % |
| NE2 | 70 | 70/250 = 28 % | 70/(70+50+30) = 47 % | out of use 0 % | 70/(100+70+30) = 35 % | 70/(100+70+50) = 32 % |
| NE3 | 50 | 50/250 = 20 % | 50/(70+50+30) = 33 % | 50/(100+50+30) = 28 % | out of use 0 % | 50/(100+70+50) = 23 % |
| NE4 | 30 | 30/250 = 12 % | 30/(70+50+30) = 20 % | 30/(100+50+30) = 17 % | 30/(100+70+30) = 15 % | out of use 0 % |
| total | 250 | 100 % | 100 % | 100 % | 100 % | 100 % |



**Fig. 10** Multiple NEs with one redundant NE

In a properly dimensioned and maintained network, the outage of one single NE shall not significantly decrease the overall network availability and thus the service availability (Fig. 11). This can be achieved with proper dimensioning of the server pool and the individual NE loads (Table 3).

Similarly, the bandwidth of the interconnections shall be planned with sufficient reserves [12]. All the paths remaining active after an outage situation shall be capable to tolerate the possible load increase due to the outage of any NE or its interconnection. Outages are either planned maintenance breaks or unplanned outages due to disaster situation (e.g. longer power supply outage, earthquake, flood etc.).
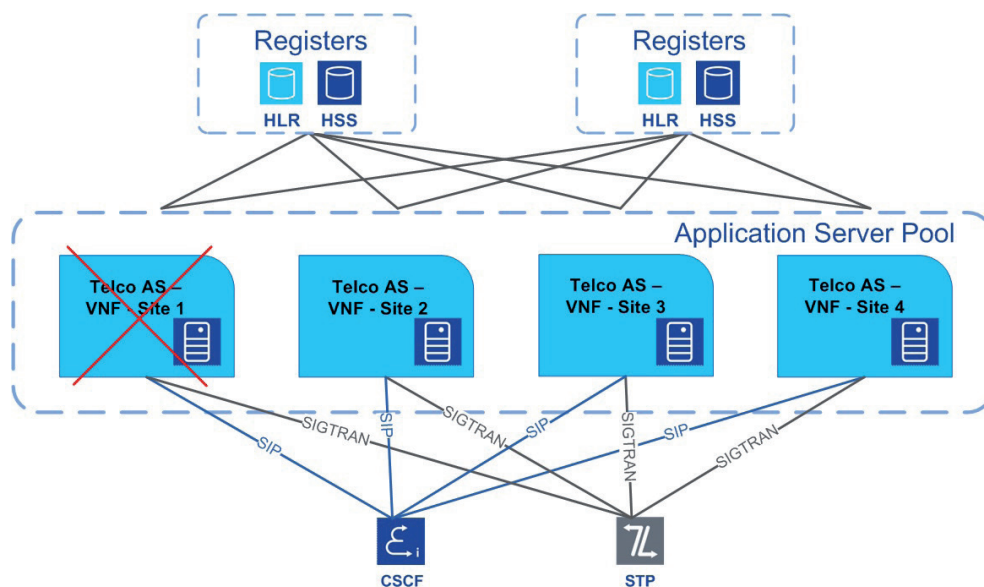


**Fig. 11** Pooling concept of telecommunication application servers
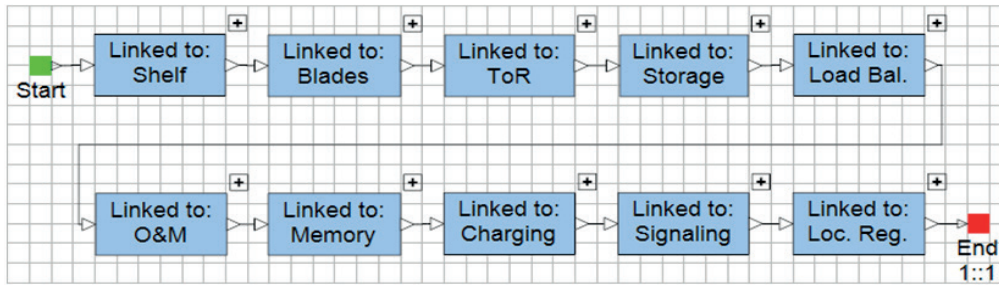
Fig. 12 Reliability block diagram example of a Telecommunication Application Server

## 5 Calculation example and simulation results

Availability predictions were calculated using Windchill (Relex) tool [13]. The reliability block diagram (Fig. 12) contains all the (critical) functional units of the NE. All the critical cloud HW building blocks are 2N redundant (e.g. EoR, ToR, bay switches and power supply units, (Fig. 13)). Storage system employs RAID 10 [14-16]. Functional units are either 2N or N+1 protected or using load sharing (e.g. signalling units) [17].

Affinity rules ensure that the virtual machines (VMs) of protected functional units are separated onto physically different HW blades. In this way, a single computer blade failure cannot cause simultaneous outage of working and standby units (of the same function, unless they belong to different RGs). Non-functional units (e.g. statistical units) are not involved in the availability calculations due to the fact that these units do not have any influence on call (or service) handling. Naturally they are involved in the dimensioning and load calculations,

as non- functional units also have their own physical resources such as VMs on the blades. Similarly to functional units, non-functional units are also consuming CPU, memory and storage.
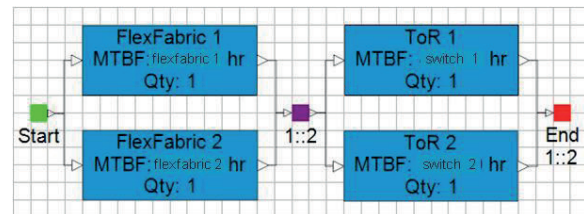


Fig. 13 Modeling the fully redundant (2N) IP switches

Monte-Carlo method provided the simulation results shown in Fig. 14. In this example, the total downtime was 58 seconds (0.016145 hours over one year (8760 h) as displayed in the figure). The predicted availability of the NE on Cloud is almost 'six nines'. (The simulation tool displays the calculated values rounded up to six decimal digits.)
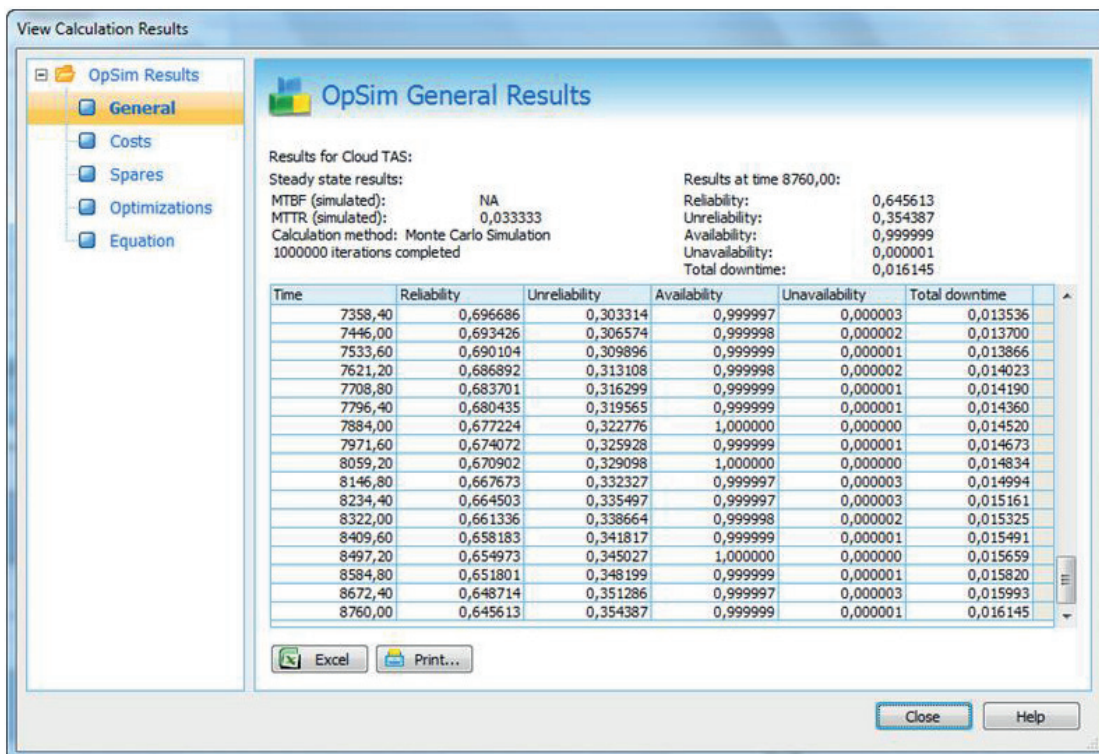


Fig. 14 Simulation results using Windchill [13] tool

## 6 Conclusions

Overall availability of telecommunication networks depends on NE, interconnection and NW level redundancy methods. Simulation results predict that 'telco-grade' availability can be achieved on cloud based core network elements (e.g. AS or MSS) of mobile networks. Critical HW and SW functional units shall be redundant. As we can see in Fig. 15-16, full protection and load sharing are more efficient than other methods, especially with increasing number of parallel nodes or units.
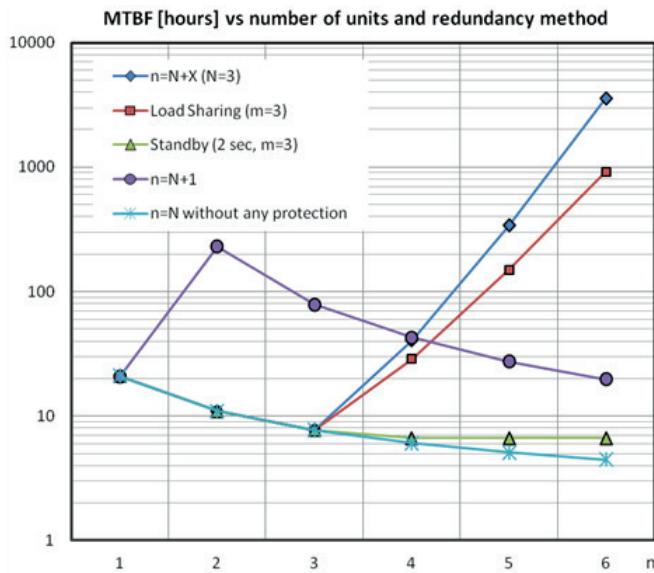


**Fig. 15** System MTBF$_S$ hours up to 6 parallel units (example of MTBF$_U$ = 20 hours, MTTR$_U$ = 1 hour)



**Fig. 16** System availability up to 6 parallel units (example of MTBF$_U$ = 20 hours, MTTR$_U$ = 1 hour). Comparison of non-protected and different redundancy systems.

Obviously, full protection (2N) is the most powerful, but 2N redundancy has the biggest footprint. Furthermore the increasing number of units results in increasing number of interconnections that brings additional possible failures.

Therefore, 2N is recommended only for 'mission critical' items (e.g. the IP switches). Interconnections must avoid single point of failure (SPoF) already in design phase (both HW and SW). In case of multiple items (e.g. load balancers or signalling units), load sharing gives an optimal trade-off because of its effectiveness and relatively smaller footprint (e.g. compared to the 2N). Furthermore, load sharing efficiently supports the dynamic scaling of the functions (functional units) on Cloud [3,18]. Pooling concept helps on NW level to overcome individual NE outages.

### List of Abbreviations

| | |
|---|---|
| A | Availability |
| A$_{NE}$ | Availability of a Network Element |
| A$_S$ | System Availability |
| A$_U$ | Availability of a unit |
| AS | Application Server |
| BSC | Base Station Controller |
| BGW | Border Gateway |
| BTS | Base Transceiver Station |
| CPU | Central Processing Unit |
| CSCF | Call Session Control Function |
| eNB | evolved Node B (LTE base station) |
| EoR | End of Row (switch) |
| e2e | end-to-end |
| FIT | Failures in Time |
| FPMH | Failure Per Million Hour |
| FRU | Field Replaceable Unit |
| FU | Functional Unit |
| HLR | Home Location Register |
| HSS | Home Subscriber Server |
| HW | Hardware |
| IMS | IP Multimedia Subsystem |
| IP | Internet Protocol |
| IT | Information Technology |
| λ | Failure Rate |
| λ$_S$ | System Failure Rate |

| | |
|---|---|
| $\lambda_U$ | Unit Failure Rate |
| LB | Load Balancing |
| LS | Load Sharing |
| LTE | Long Term Evolution |
| MDT | Mean Down Time |
| MGW | Media Gateway |
| MME | Mobility Management Entity |
| MNO | Mobile Network Operator |
| MRF | Media Resource Function |
| MSC | Mobile Switching Center |
| MSS | MSC Server |
| MTBF | Mean Time Between Failures |
| $MTBF_{NE}$ | MTBF of a Network Element |
| $MTBF_S$ | System MTBF |
| $MTBF_U$ | MTBF of a unit |
| MTTF | Mean Time to Failure |
| MTTR | Mean Time to Repair |
| $MTTR_U$ | MTTR of a unit |
| MUT | Mean Up Time |
| NB | Node B (3G base station) |
| NE | Network Element |
| NW | Network |
| O&M | Operation and Maintenance |
| OSS | Operations Support Systems |
| PCRF | Policy and Charging Rules Function |
| PSTN | Public Switched Telephone Network |
| P-GW | Packet Data Network Gateway |
| RAID | Redundant Array of Independent Disks |
| RBD | Reliability Block Diagram |
| RG | Recovery Group |
| RNC | Radio Network Controller |
| SAN | Storage Area Network |
| SGSN | Serving GPRS Support Node |
| SAE | System Architecture Evolution |
| S-GW | Serving Gateway |
| SIGTRAN | Signalling Transport (protocol) |
| SIP | Session Initiation Protocol |
| SPoF | Single Point of Failure |
| STP | Signalling Transfer Point |
| SW | Software |
| TAS | Telecommunication Application Server |
| ToR | Top of Rack switch |
| U | Unavailability |
| $U_{NE}$ | Unavailability of a Network Element |
| $U_S$ | System Unavailability |
| $U_U$ | Unavailability of a unit |
| UE | User Equipment (e.g. mobile phone) |
| VM | Virtual Machine |
| VNF | Virtual Network Function |
| 2G | 2nd generation mobile telephone technology |
| 3G | 3rd generation mobile telecommunications technology |

## References

[1] Csatári, G., László, T. "NSN Mobile Core Network Elements in Cloud, A proof of concept demo." In: *IE*EE International Conference on Communications Workshops (ICC), 2013, pp. 251-255, Budapest, Hungary, 9-13 June 2013. DOI: 10.1109/ICCW.2013.6649238

[2] Rotter, Cs., Farkas, L., Nyíri, G., Csatári, G., Jánosi, L., Springer, R. "Using Linux Containers in Telecom Applications.", accepted at Innovations in Clouds, Internet and Networks, ICIN 2016.

[3] Bakos, I., Bódog, Gy., Hilt, A., Jánosi, L., Járó, G. "Resource and call management optimization of TAS/MSS in Cloud environment." In: Infokom'2014 Conference, Hungary, Oct. 2014. (in Hungarian)

[4] 3GPP, TS 23.002, 3rd Generation Partnership Project, Technical Specification Group Services and System Aspects, Network architecture, Release 13, V13.3.0, 09. 2015.

[5] Lemaire, M. "Dependability of MV and HV protection devices." Cahier Technique Merlin Gerin n°175, pp.1-16, ECT 175, Aug. 1995. URL: http://www.schneider-electric.co.uk/documents/technical-publications/en/shared/electrical-engineering/dependability-availability-safety/high-voltage-plus-1kv/ect175.pdf

[6] Kleyner, A., O'Connor, P. "*Practical Reliability Engineering*." 5th edition, Wiley, New York, NY, 2011.

[7] Ayers, M. L. "*Telecommunications System Reliability Engineering, Theory and Practice*." Wiley, New York, NY, 2012. DOI: 10.1002/9781118423165

[8] Gnanasivam, P. "*Telecommunication Switching and Networks*." 2nd edition, New Age International Ltd., Publishers, 2006.

[9] Viswanathan, T., Bhatnagar, M. "*Telecommunication Switching Systems and Networks*." 2nd edition, PHI Learning Private Limited, Delhi-110092, 2015.

[10] Lin, D. L. "Reliability Characteristics for Two Sub- systems in Series or Parallel or *n* Subsystems in *m* out of *n* Arrangement." Technical report, Aurora Consulting Engineering LLC, 2006.

[11] Highleyman, W. H. "Calculating Availability – Redundant Systems." Sombers Associates Inc., [Online]. Available from: www.availability-digest.com, 2006. [Accessed: 26th January 2016]

[12] Nokia "Network Resilience in CS Core and VoLTE and Integrated IMS Services." DN0631359, 2015.

[13] Relex Software Corporation "Reliability: A Practitioner's Guide." [Online.] Available from: www.relexsoftware.com, 2003. [Accessed: 26th January 2016]

[14] Shooman, M. L. "*Reliability of Computer Systems and Networks, Fault Tolerance, Analysis and Design*." Wiley, New York, NY, 2002.

[15] Malhotra, M., Trivedi, K. S. "Reliability Analysis of Redundant Arrays of Inexpensive Disks." *Journal on Parallel and Distributed Computing*. 17 (1-2), pp. 146-151. 1993. DOI: 10.1006/jpdc.1993.1013

[16] Marcus, E., Stern, H. "*Blueprints for high availability*." 2nd edition, Indianapolis, John Wiley & Sons, 2003.

[17] Birolini, A. "*Reliability Engineering, Theory and Practice*." 3rd edition, Springer, 1999. DOI: 10.1007/978-3-662-03792-8

[18] Bakos, I., Bódog, Gy., Hilt, A., Jánosi, L., Járó, G. "Optimized resource management in core network element on Cloud based environment." Patent PCT/EP2014/075539, Nov. 2014.