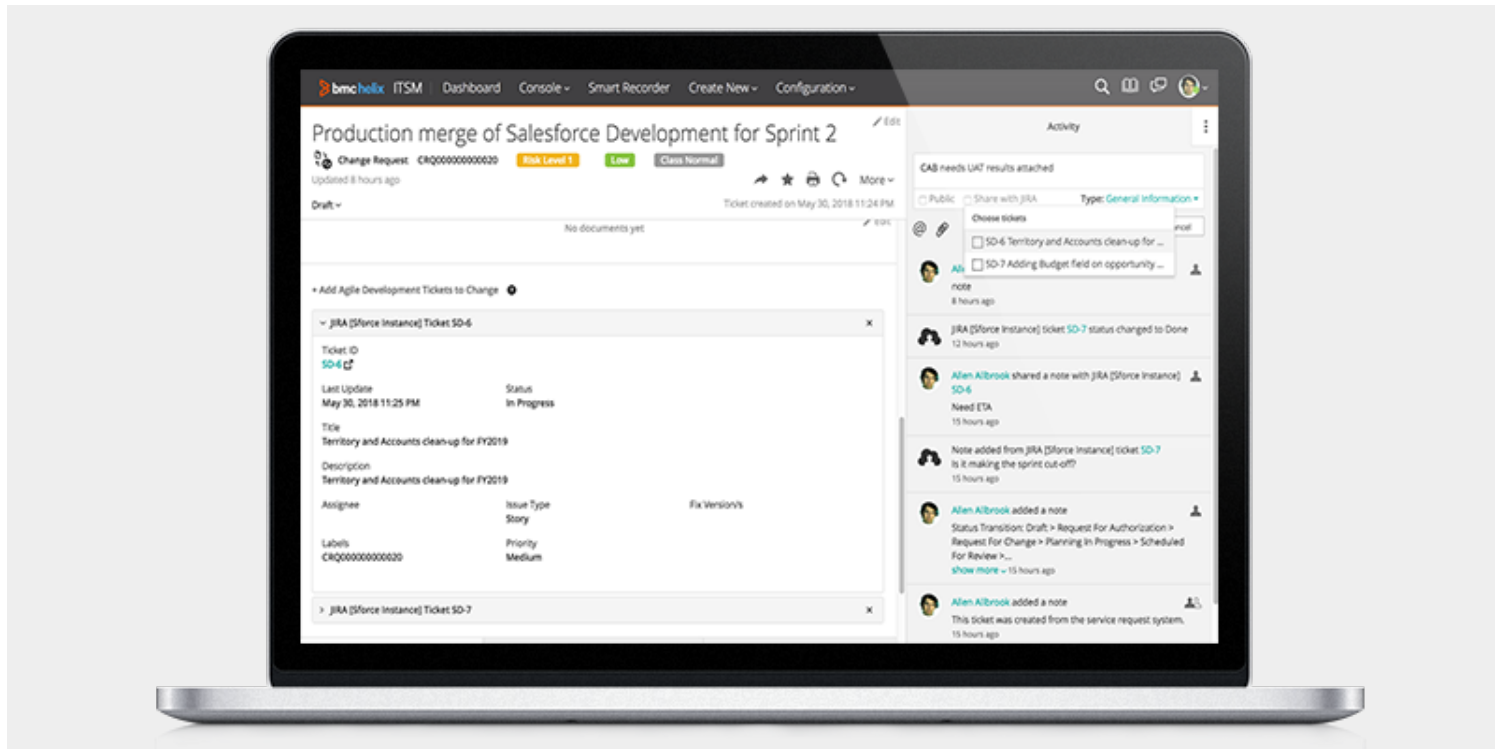# N, N+1, N+2, 2N, 2N+1, 2N+2, 3N/2 REDUNDANCY EXPLAINED



For an IT service to be readily used for critical use cases, the service provider must employ appropriate measures designed to enhance the dependability of the service.

Dependability refers to the trustworthiness of an IT system, which allows reliance to be justifiably placed on the IT services delivered by those systems. Introducing redundancy into IT systems is one approach to reduce the risk of service interruptions, protect against anomalous behavior and establish adequate fault tolerance within the system.

Redundancy refers to the system design principle that involves duplication of components constituting the IT systems to achieve this goal.

## N-Modular Redundancy

This duplication may be applied to the hardware, information, software as well as time that govern the operation of an IT system. Various configurations of redundant system design may be used based on the associated risk, cost, performance and management complexity impact. These configurations take various forms, such as N, N+1, N+2, 2N, 2N+1, 2N+2, 3N/2, among others. These multiple levels of redundancy topologies are described as N-Modular Redundancy (NMR):

- **N** refers to the bare minimum number of independent components required to successfully perform the intended operation. For instance, a data center server may require one power supply operating at specific power ratings to keep the server running at optimal availability conditions.

- **N+X** refers to a redundant system that contains **X** number of spare components to act as an independent backup when the appropriate component fails to operate as intended. For instance, N number of servers are running in a data center to provide the necessary computing power reliably to a specific number of users. Should one of the servers fail, an additional server is available to fill the gap in delivering the same computing capacity as originally intended. N+1 means that only one backup component is available to complementing the N original component(s). N+2 means that two backup components are available to ensure additional resilience.
- **YN** refers to the number of times the capacity is available to replace the entire set of original components. Consider the case of YN redundancy for a 10-server set up at a data center, which is also the maximum computing capacity required. With the 2N redundancy, an additional 10 servers of the same performance specifications and capacity are available to replace a few, or all of the original servers as required. Similarly, higher orders of redundancy such as 3N (e.g. 10 original servers plus 20 additional servers), 4N (10 original servers plus 30 additional servers) and so on may be introduced based on level of system reliability required.
- **YN+X** refers to the combination of the above two topologies. For instance, a 2N+1 redundancy topology for a 10-server data center set up means that at all times, twice the sever capacity (20 servers) is available, plus one additional server to operate as a backup if required. However, only a maximum N number of original servers are intended to operate under normal operating conditions without any server failure.
- **AN/B** refers to the shared redundancy topology, where **A** amount of backup capacity is available for total **B** amount of load or original components. For instance, consider a 3N/2 redundancy topology for the backup power supply of a server environment. The three available power backup systems will be available for every two server loads. Each of the three power supplies operate at a maximum of 67 percent at any time, or the inverse of the 3N/2 ratio. For a 4N/3 redundancy topology of the same UPS power backup, 4 UPS components will serve every 3 server machines, and each UPS will operate at a maximum of 75 percent capacity. It's interesting to note that 3N/2 may be similar to 2N+1 in a 2-server environment in terms of redundancy but differentiate in terms of operating cost and complexity.

The redundant components provide additional reliability based on different redundant topologies. However, the spare components may or may not be identical to the original components in terms of capacity. The redundancy system may offer Active, Passive, Load Sharing or Standby configuration. Active redundancy means that the redundant component is operating simultaneously to the original component, but the output is only used when the original component fails. The Passive component is switched on only after the original component fails. The Standby redundancy component fills in the availability gap temporarily until the startup of an original or Active component takes place. An additional load sharing redundancy may be applied to offer partial redundancy in meeting the necessary resilience goals.

It is important to note that the redundancy topologies and configurations are not defined consistently across the literature available on the redundancy theory. Authors may use different categories (such as N+1, 2N+1) to describe the redundancy capacity of spare components based on their Active, Passive, Standalone or Load Sharing capability.

The choice of various redundancy topologies depends on the resilience requirements of the system, the probability of failure of components individually and collectively, the cost and complexity of operating the redundant system. Data centers offering high dependability will use higher orders of

redundancy topologies, including complex configurations to optimize load distribution at higher operating cost.

From a customer perspective, the dependability of an IT service is determined by its reliability and availability, which data center companies may guarantee using different redundancy topologies. The relations between reliability and redundancy largely depends upon the failure rates of each individual component within the system. For instance, components with low failure rates may require simple redundancy topologies and few spare components to guarantee high availability, while components with high failure rates may require complex redundancy topologies to guarantee high reliability of the service.

According to the 2016 Ponemon Institute research report, Cost of Data Center Outages, the average cost of data center outage increased by 38 percent from $500,000 to $740,000 between the years 2010 to 2016. The analyzed cost centers include the direct, indirect and opportunity cost. However, the opportunity cost varies significantly between the customers, forcing some of them to pursue highly available services at affordable cost so the overall IT service is considered as dependable based on all decision factors. For service providers, a range of N, N+1, N+2, 2N, 2N+1, 2N+2, 3N/2 and similar redundancy topologies are used to meet these dependability requirements.