



# Next-Gen Resilience for Nonstop Business- Critical Apps

MAY 2018

COMMISSIONED BY





## About this paper

A Pathfinder paper navigates decision-makers through the issues surrounding a specific technology or business case, explores the business value of adoption, and recommends the range of considerations and concrete next steps in the decision-making process.

## About 451 Research

451 Research is a preeminent information technology research and advisory company. With a core focus on technology innovation and market disruption, we provide essential insight for leaders of the digital economy. More than 100 analysts and consultants deliver that insight via syndicated research, advisory services and live events to over 1,000 client organizations in North America, Europe and around the world. Founded in 2000 and headquartered in New York, 451 Research is a division of The 451 Group.

© 2018 451 Research, LLC and/or its Affiliates. All Rights Reserved. Reproduction and distribution of this publication, in whole or in part, in any form without prior written permission is forbidden. The terms of use regarding distribution, both internally and externally, shall be governed by the terms laid out in your Service Agreement with 451 Research and/or its Affiliates. The information contained herein has been obtained from sources believed to be reliable. 451 Research disclaims all warranties as to the accuracy, completeness or adequacy of such information. Although 451 Research may discuss legal issues related to the information technology business, 451 Research does not provide legal advice or services and their research should not be construed or used as such.

451 Research shall have no liability for errors, omissions or inadequacies in the information contained herein or for interpretations thereof. The reader assumes sole responsibility for the selection of these materials to achieve its intended results. The opinions expressed herein are subject to change without notice.

### NEW YORK

1411 Broadway  
New York NY 10018  
+1 212 505 3030

### SAN FRANCISCO

140 Geary Street  
San Francisco, CA 94108  
+1 415 989 1555

### LONDON

Paxton House  
(Ground floor)  
30, Artillery Lane  
London, E1 7LS, UK  
P +44 (0) 207 426 1050

### BOSTON

75-101 Federal Street  
5th Floor  
Boston, MA 02110  
Phone: +1 617.598.7200  
Fax: +1 617.357.7495

## Executive Summary

We expect the need for resiliency and the way it is achieved at both the datacenter and application levels to change significantly as the IT world gradually moves to a more cloud-based, hybrid and distributed architecture. Advances in hybrid and cloud computing, containers and virtualization, DevOps, replication, distributed databases, and global traffic management are coinciding with huge investments in datacenter and network capacity, and a move toward greater use of software to intelligently manage workloads, traffic and resiliency. All of these are combining to create an architectural shift: from single-site vertical resiliency to distributed, replicated resiliency.

This shift offers enormous promise for enterprise CIOs and has the potential to disrupt the infrastructure ecosystem of suppliers and many service provider operators. Studies by 451 Research suggest a more active, distributed and automated approach to managing infrastructures will become much more widely used. Indeed, the proliferation of distributed stateless cloud applications is beginning to make the traditional site-based approach to resiliency look increasingly limited – even if it remains a necessary component.

Negotiating a clear path ahead in this area presents challenges for operators, service providers and CIOs alike. For most businesses, service availability is critical and non-negotiable, so any move to newer, cheaper or more dynamic technologies or architectures must be undertaken with great care. The prize is better resiliency, with risk distributed and spread, and at lower cost of ownership; the price is more complexity, a possible loss of transparency and control, and a transitional period that could involve considerable costs, complexity and dependency on unproven – or less proven – IT.

### THE RESILIENCY IMPERATIVE

While single-site availability remains a critical concern, it is no longer sufficient to approach resiliency as based on a single datacenter. Nor is it about a single application, not least because applications are becoming increasingly failure-blind. Today, resiliency is primarily about the entire infrastructure or architecture. Additionally, it has to be concerned with recovery, as well as failure prevention. In the medium to long term, we believe the trend is for traditional disaster recovery (DR) and business continuity to become part of overall resiliency planning.

As systems evolve to have more complexity and interdependency, failure and recovery will often be less binary (on/off) than in the past. Systems can be degraded or be missing some services/components or critical data while still functioning, unaffected, in other areas. This ambiguity about what failure means will undoubtedly lead to confusion and ‘spinning’ by some operators.

It has occasionally been said that the move toward distributed and cloud architectures/business models represents a trend toward a more patient, gradual availability model, where failures are tolerated in exchange for cheaper services. This is not the case – underlying tradeoffs in data availability do not make businesses more tolerant of failures or slowdowns. In fact, the exact opposite is the case. In our 2017 Uptime Institute survey, only 8% of respondents said their management is less concerned about outages than it was a year ago. Incidents and failures appear to resonate outward, and have an ever-greater impact.

Resiliency is all too obvious when it is absent. Some recent datacenter and core infrastructure failures have led to national and global headlines and major financial and reputational losses (see Figure 1). Many of these incidents show that failures of equipment or processes quickly escalate, and frequently involve multiple facilities and IT systems.

#### DEFINITION

Resiliency describes the extent to which a system, digital infrastructure or application architecture is able to maintain its intended service levels, with minimal or no impact on users or business objectives, in spite of planned and unplanned disruptions. It also describes the ability of a system, infrastructure or application to recover full business operations after a disruption or disaster has occurred.

# PATHFINDER REPORT: NEXT-GEN RESILIENCE FOR NONSTOP BUSINESS-CRITICAL APPS

Figure 1: Selected recent major datacenter incidents/outages

COMPANY/DATACENTER	DATE	INCIDENT	IMPACT
AWS	Q1 2017	Many S3 and related AWS services at US East suffered partial or complete outage for several hours.	Hundreds of third-party services affected. Financial and reputational damage.
British Airways	Q2 2017	Failures and operator error in London datacenter led to loss of most operational systems for many hours.	Flights delayed or cancelled for several days. Financial (£80m) and reputational loss.
OVH	Q4 2017	Two incidents – the initial cause an external cooling system leak – crashed server and storage arrays in Paris.	More than 50,000 websites down for more than 24 hours.
US Customs & Border Protection	Q1 2018	Two-hour outage affected international arrival clearance at Atlanta, Denver, JFK, Miami and San Francisco airports.	Severe passenger delays. Followed similar (four-hour) outage in January 2017.
Microsoft Outlook	Q3 2018	Office 365, Outlook and Exchange down in APAC, EMEA and some of the US for up to a day. Incidents in January and April were attributed to attempted security upgrades.	Loss of customer confidence at a time when Microsoft is pushing cloud services.
TSB Bank	Q3 2018	Problematic systems migration at TSB bank left nearly two million customers unable to access their accounts. Inadequate testing was blamed.	Severe reputational and financial damage. Customers disrupted for a week or more.

Source: 451 Research

## How is Resiliency Changing?

This system-wide view represents an important change. A de facto design principle of almost all datacenters today, and much of IT, is a primary focus on physical or infrastructure resiliency – the ability to continue running in spite of maintenance, power or equipment failures through redundancy of equipment and power distribution. But over the next decade, we expect that for a sizable but as yet undetermined number of enterprises and operators, resiliency and redundancy at the individual datacenter level will, in whole or part, be complemented or replaced by resiliency at the IT level. This is not necessarily a resiliency tactic or strategy, but an inevitable move, resulting from the fact that applications themselves are becoming more distributed.

Achieving distributed resiliency will not happen uniformly. At best, it will involve the rapid handover of workloads to alternative venues in the event of failures, under automated software control, supported by reliable networks. This has been variously called ‘software-level resiliency,’ ‘network-level resiliency’ or ‘cloud-level resiliency.’ We prefer the term ‘distributed resiliency.’ Figure 2 shows some of the key enabling technologies for this approach.

Although adoption of these technologies will have a big impact on datacenters and datacenter design, they have wider implications as a key part of digital transformation for all companies. If resiliency can be fully and safely achieved based on distributed, lightweight systems and datacenters, then companies will be able to rewrite the way they plan, build and spend on IT and digital services.

**Figure 2: Distributed resiliency: enabling components/technologies**

<b>KEY ENABLING TECHNOLOGY/INFRASTRUCTURE</b>	<b>WHY?</b>
Datcenter capacity in multiple locations	There must be sufficient capacity to absorb workloads or move them around.
Homogenous off-the-shelf hardware	Servers and storage must be homogenous across sites to pick up loads or replicate data/applications.
Virtualization/cloud platform/containers	Applications will usually be portable and often composable/stateless.
Load-balancing software	Platform software is needed to balance loads and move work if certain sites lose or lack capacity.
GTM/domain name management	Traffic must be switched seamlessly if there are problems at any sites.
Application synchronization/systems of engagement for caching transactions	Users should have a consistent experience even if back-end systems are not working fully or at all.
Software-defined network management	Network tools should be able to reconfigure the network dynamically if required.
Distributed databases/not-only-SQL databases	Databases that are capable of working across multiple locations will provide integrity and transaction management in the event of node losses.
Cloud orchestration/management software	Most applications will run in cloud platforms, often multiple clouds. Orchestration software must ensure interworking.
Storage recovery/management/DR/resiliency management tools	Storage management systems and related tools must ensure recovery and management according to service levels across multiple sites.
Network and facility management tools/resiliency management tools	Tools for managing facilities and networks, including at the physical level, will help automation and rapid configuration/recovery.

Source: 451 Research

## Potential Benefits of Distributed Resiliency

To some degree, the general move to distributed applications will require a new approach to ensuring resiliency, regardless of whether this is an intentional strategic/architectural move. CIOs and service providers will find that, over time, they must tactically address weaknesses in their infrastructure that cause incidents and loss of service.

However, many of the benefits associated with a deliberate move to distributed resiliency are compelling. The long-term vision is that resiliency ultimately becomes autonomic – managing itself, shifting loads and traffic across geographies according to need, replicating data, and optimizing for performance and economics with little intervention. The short- and medium-term promise is that, after a transition, CIOs will have a more reliable, agile infrastructure that costs less and can support modern distributed applications far better. The key benefits, some theoretical and some already manifesting themselves, are as follows.

- **Availability:** As long as at least three datacenters are involved and the networking has sufficient capacity and variety, then extremely high availability should exist. Failure at two or more datacenters is extremely unlikely, and even less so across a larger number of datacenters or multiple regions.
- **Efficiency:** Cloud resiliency (whether public clouds are involved or not) means that as soon as more than two datacenters are used, spare active capacity is spread among a number of them. If utilization is pushed to high levels, then the more datacenters involved, the more efficient the operation becomes.
- **Agility:** A distributed approach means that applications can be developed to run at any of the sites, or across all – and appropriate infrastructure investments can also be made with greater flexibility.
- **Costs:** A move to cloud-based resiliency will ultimately mean a move to an all-active distributed datacenter model. This means both the processes and resources required for effective DR may no longer be needed, or at least may be reduced and redefined, with governance and policies becoming more important. Cloud-based DR services, now multiplying, represent an effective partial solution, avoiding the costs of unused capacity.

- **Reduced redundancy:** Because loads are running in multiple datacenters, or traffic can be rapidly switched to any data-center, the need for high levels of redundancy in power, cooling and IT at each site may, in theory and in some situations, be carefully reduced, potentially saving significant sums.
- **Testing:** One of the problems with traditional resiliency and DR strategies is that testing is extremely difficult and sometimes risky because tests must be carried out on live systems, sometimes at a datacenter-wide level. Distributed resiliency uses soft failovers, with 'hot swappable' software and cloud services. Many foreseeable failures can be easily and repeatedly tested in different ways.
- **Applications:** Almost by definition, single-site or traditional resiliency approaches cannot fully support applications that have components in multiple locations across networks. When failures occur, local applications will be missing core components. By using replication of components and multiple sites, a distributed strategy may make such incidents less likely.

## Types of Distributed Architecture

As we have seen, differing business requirements, including legacy investments, will influence the degree to which newer, distributed systems and databases can be used; similarly, the business requirements and the design of the existing systems will, to some extent, point toward certain resiliency architectures. We see the following models being used for resiliency, with the cloud-based models being markedly different from the earlier ones.



**Single-site availability** is the traditional setup, with high levels of redundancy at the infrastructure level, including facilities and basic IT. With sufficient redundancy and planned design, operations can continue in spite of planned (concurrent maintainability), and in some cases unplanned, facilities failure. At the IT level, resilience is further ensured by internal replication (typically through the use of clusters), so that loads may be replicated elsewhere and data, applications and configurations backed up to an off-site DR location.



**Limited site resiliency** describes two or more lower-tier datacenters within a campus, region or zone using a dedicated network to achieve a higher level of availability than is possible at any individual site, typically within synchronous replication distance (i.e., close enough to each other and to customers so that they are always synchronized). This distance will depend on the applications, but is usually less than 50 miles. This setup can be used to support either synchronous (fault-tolerant automated failover to the second site) or asynchronous (a second copy of applications, data and files is kept at the second site to pick up the load) replication.



**Distributed site resiliency** describes two or more independent sites, in or out of region or globally distributed (cloud or not), using shared internet/VPN networks to provide resiliency through multiple asynchronously connected instances. This can produce very high availability but can result in some (usually minor) loss of integrity between instances if outages occur. This is the architecture that underpins most DR services, and especially the modern cloud iteration, DR as a service. Improved network capacity, software tools, database synchronization protocols and, critically, homogenous IT infrastructure running virtualized workloads have made this option far more practical.



**Cloud-based resiliency** is provided by distributing virtualized applications, instances and/or containers with associated data across multiple datacenters, using middleware, orchestration and distributed databases, under the control of a comprehensive and distributed system. These systems will enable service or design choices to be made between, for example, absolute database integrity or immediate availability. Effectively, cloud-based resiliency moves the resiliency up to the IT level. Any facility resilience achieved through redundancy provides added security but may not prove essential. It does, however, assume that there is sufficient capacity in place, including the network, which is critical if loads are shifted from place to place. Developers do not need to concern themselves with location or infrastructure – this architecture is primarily suited for stateless or cloud-native applications.

## Which Architecture is Best?

Clearly, each type of resiliency architecture described above fulfills different purposes and has a different profile in terms of objectives, cost, level of availability and technical maturity. Cloud-based resiliency is the newest, and currently the most expensive; it may provide good total cost of ownership, but effectively can only be achieved at scale and with considerable capital. Each type is not mutually exclusive, at least at the facilities level.

For CIOs setting out to develop appropriate resiliency strategies, this is a challenging period because engineering control is being eroded, to be replaced with a more nuanced and strategic approach where good assessments are needed. With cloud services and architectures now part of the mix, or even the totality, CIOs must determine which type (or types) of resiliency is most appropriate for each type of application and data, based on business needs and technical risk, and then architect the best combination of IT infrastructure. This will span datacenter resiliency, applications, databases and networking, and must take into account organizational structure, processes, tools and automation. From all this, the organization must then deliver comprehensive and consistent applications that meet and exceed customer expectations for service availability and resiliency.

How do most enterprises achieve resiliency today? The traditional approach – a redundant datacenter infrastructure supporting a more complex IT architecture, with a backup and recovery plan using separate or remote non-live ('passive') facilities – still remains the dominant model, but many companies are using hybrid or newer methods.

## Application Considerations

When it comes to resiliency, many new cloud applications are no longer directly dependent on core systems sited at key datacenters because both the data and the application can be easily replicated in multiple places. The great majority of these applications are stateless, meaning the tight link between underlying systems and sessions that characterized earlier generations of IT is now broken. This new software design paradigm offers great advantages in agility, speed of development, reuse, portability and hardware costs, and is rapidly becoming the dominant way to develop new applications.

Understanding the implications of this shift is important when it comes to planning a resiliency strategy for a single application, or for the digital infrastructure as a whole. Effectively, distributed/cloud-based resiliency moves the resiliency up to the IT level – but not to the application level. By separating and virtualizing the application layer from the underlying infrastructure, it becomes possible for data and applications to be stored and processed in multiple places, spread across few or many sites or zones.

This is not always possible. In 'legacy' or traditional architectures, applications are hosted on single computers or systems, making calls to disks, databases or external systems when they need data. Many of these applications cannot be easily (or ever) transitioned to the newer model. Our research suggests that enterprises are finding the cost and the technical difficulty of changing from older to newer architectures to be very high and sometimes prohibitive. It is, therefore, reasonable to assume that for the next decade or so, most enterprises will have mixed or hybrid architectures. Figure 3 briefly describes the architecture suitable for each application type.

# PATHFINDER REPORT: NEXT-GEN RESILIENCE FOR NONSTOP BUSINESS-CRITICAL APPS

Figure 3: Application types and associated resiliency architectures

APPLICATION TYPE	SUITABLE RESILIENCY ARCHITECTURE (DATACENTER/IT)	MOST LIKELY VENUE
Transactional	Latency and availability are critical; most suited to single/linked site.	Single/linked site
Traditional	Limited by design to single sites or single virtual instances; rewrite needed to take advantage of cloud.	Single/linked site
Traditional HA	Limited by design to single sites or single virtual instances; rewrite needed to take advantage of cloud.	Single/linked site
Distributed HA	Multiple, virtualized geo-clusters replicating data and load balancing across large distances to protect against regional disasters.	Single/linked site (metro distances)
Cloud-optimized	Able to run on public/private cloud but requires modification/configuration to take advantage of distributed resiliency.	Single/linked or distributed resiliency
Cloud-native	Will run on cloud architectures (only) and requires trivial configuration for resiliency.	Distributed resiliency or cloud-based resiliency

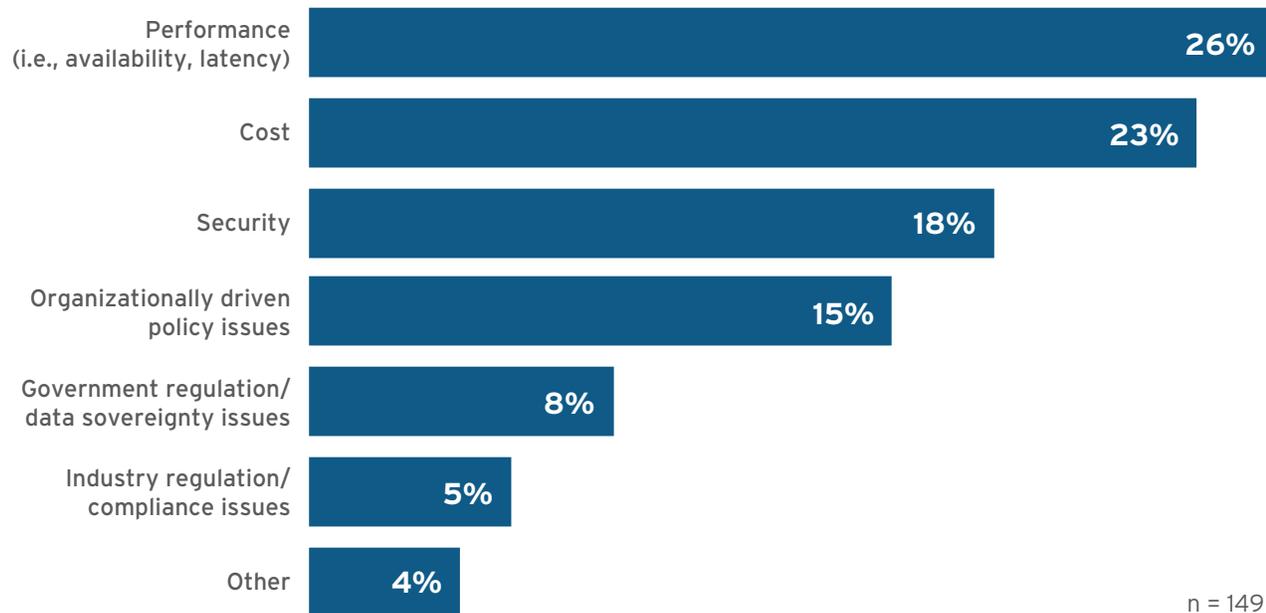
Source: 451 Research

## Are Enterprises Ready?

According to the most recent 451 Research Voice of the Enterprise: Cloud, Hosting and Managed Services survey (see Figure 4), many enterprise customers are still wary of the maturity of cloud-based availability and resilience – to the extent that performance/availability was cited as the top reason for migrating services previously moved to the public cloud back to private cloud and non-cloud on-premises systems.

Figure 4: Drivers for migrating from public cloud to private cloud or non-cloud environment

Q. What was the primary driver for migrating workloads from a public cloud to a private cloud or non-cloud environment?



Source: 451 Research Voice of the Enterprise: Cloud, Hosting and Managed Services, Workloads and Key Projects, May 2018

According to one enterprise IT customer participating in the survey, the need to consider the implications of distributed resiliency is top of mind today, even though the organization isn't yet ready to implement it: "We are not writing cloud-native apps yet. So, we had a huge debate over where to hide disaster recovery, full tolerance, and high availability, under what layer ... We agreed that it should be a commodity service, in the sense that applications should be potentially, soon enough, written to utilize cloud services but with the mindset of a platform."

Another enterprise IT executive voiced a different view, saying that the company remains "very worried" about maintaining uptime and spends a lot of time and resources maintaining the high level of internal resilience. For that organization, a critical decision factor when choosing a cloud service is making sure that the cloud service provider has a stable, resilient platform for the company to use.

### Summary and Recommendations

The CIO's approach to resiliency, availability and recovery is set to change in the next decade, bringing all operators more in line with large-scale cloud providers. For many with mission-critical operations, proven approaches based on tight process control, redundancy and even over-provisioning will still be needed in places, but at the CIO level, there will increasingly be a need for an overall strategy that is less binary and more nuanced. This approach will involve trading risks and costs, and sometimes accepting or supplementing what service providers can offer.

The types of resiliency described here will form the building blocks of resilient architectures, increasingly freeing up the applications teams (and cloud customers) to disregard most issues around availability, locality, site facilities and redundancy. The biggest mistake any CIO or operator can make is to assume that, in the new world of IT, the resiliency strategy can simply be outsourced.

The use of public cloud services presents both benefits and challenges. There is no doubt that big public cloud providers can offer a very high level of resiliency/availability. However, achieving this resiliency is only straightforward if the developers or customers are producing cloud-native applications – and even then, configuration, some extra services and a high degree of trust (probably involving provider lock-in) will be needed. For most of the rest, moving to a distributed, cloud-based or hybrid architecture is a positive step with many benefits, but resiliency will require careful attention and work; it is decidedly not trivial.

Enterprises undertaking digital transformation that do not want to find themselves in the headlines should conduct a detailed examination of their application and infrastructure resiliency at the outset of any big project – and regularly after that – as a matter of corporate due diligence.

## NONSTOP IT FOR BUSINESS-CRITICAL APPLICATIONS

Whether it is the planned downtime for system or hardware updates, or unplanned downtime from failures or even natural disasters, your business success counts on perpetual, predictable and nonstop IT. If your applications rely on precision timing and responsiveness, your business success depends on maintaining that predictability. Providing more uptime and reliability in your enterprise or data center is essential in meeting the demands of business-critical workloads. The SUSE product portfolio is designed to help:

- *Maximize service availability for nonstop IT*
- *Raise the bar on responsiveness and precision timing with a real time operating system*
- *Maintain business-critical continuity by applying kernel fixes on the fly without interrupting service*

## PRODUCT PORTFOLIO

- *SUSE Linux Enterprise Server (SLES)*
- *SUSE Linux Enterprise High Availability Extension (HA)*
- *Geo Clustering for SUSE Linux Enterprise High Availability Extension (GEO)*
- *SUSE Linux Enterprise Real Time (RT)*
- *SUSE Linux Enterprise Live Patching*
- *SUSE Global Services*

### 1. SUSE Linux Enterprise Server

SUSE Linux Enterprise Server is a world-class, secure open source server operating system, built to power physical, virtual and cloud-based mission-critical workloads. The operating system further raises the bar in helping organizations to accelerate innovation, enhance system reliability, meet tough security requirements and adapt to new technologies. The operating system optimizes your infrastructure with support for the latest hardware environments for ARM System on Chip, Intel, AMD, SAP HANA, Z Systems and NVM Express over Fabrics.

Key Benefits:

- Create and Support Agile IT Infrastructure
- Deploy Mission-Critical Services
- Continuously Improve IT Infrastructure

### 2. SUSE Linux Enterprise High Availability Extension

Maximize your service availability and virtually eliminate downtime. SUSE Linux High Availability Extension provides mature, industry-leading open-source high-availability clustering technologies that are easy to set up and use. It can be deployed in physical and/or virtual environments, and can cluster physical servers, virtual servers, or any combination of the two to suit your business' needs.

Key Benefits:

- Manage Pacemaker HA clusters using High Availability Web Konsole (Hawk)
- Easy and fast setup
- High performance Oracle Cluster File System 2 (OCFS2)
- Feature rich Global File System 2
- A mature, fifth-generation Pacemaker-based HA solution
- HAProxy support to complement Linux virtual server load balancer
- Update to the latest Relax & Recover (ReaR) version
- Support for EMC NetWorker connector & btrfs file system included

### 3. Geo Clustering for SUSE Linux Enterprise High Availability Extension

In the event of a regional disaster, the power can go down and your mission-critical workloads can fail. This is not acceptable. Geo Clustering for SUSE Linux Enterprise High Availability Extension provides rules-based failover for automatic and manual transfer of a workload to another cluster outside of the affected area. Your mission-critical workloads are transferred away from the affected region to continue to run.

Key Benefits:

- Protect workloads across globally distributed data centers
- Use multi-tenancy to manage geo clusters per business needs
- Achieve maximum protection for workloads that can never have unplanned downtime
- Comply with industrial regulations with manual failover option
- Protect mission-critical applications from regional disasters
- Deploy physical and virtual Linux clusters across data centers

### 4. SUSE Linux Enterprise Real Time

If your business can respond more quickly to new information and changing market conditions, you have a distinct advantage over those that cannot. Running your time-sensitive mission-critical applications using SUSE Linux Enterprise Real Time reduces process dispatch latencies, and gives you the time advantage you need to increase profits or avoid further financial losses, ahead of your competitors.

Key Benefits:

- Pre-emptible real-time kernel
- Ability to assign high-priority processes
- Greater predictability to complete critical processes on time, every time
- Increased reliability
- Lower infrastructure costs
- Tracing & debugging tools that help you analyze & identify bottlenecks in mission-critical applications

### 5. SUSE Linux Enterprise Live Patching

What if you could keep up with compliance and security, while at the same time increasing your service availability? The only way to pull that off would be to apply kernel fixes without interrupting your service. That's exactly what SUSE Linux Enterprise Live Patching lets you do. Plus, the subscription gives you access to fixes for any SUSE Linux Enterprise Server 12 maintenance Linux kernel released in the last 12 months.

Key Benefits:

- Increase service availability
- Reduce planned or unplanned downtime

### 6. SUSE Global Services

Nonstop IT means having the resources to proactively maintain and monitor your systems. The SUSE Global Services team is dedicated to keeping your mission-critical applications running without disruption. Made up of IT consultants, premium support engineers and customer success managers with deep product and technical experience, our services team is committed to your business success.

Key Benefits:

- Evaluate your existing processes and infrastructure
- Leverage our broad knowledge of best practices
- Work hand-in-hand with your team to provide business continuity