

How to Achieve and Measure the Benefits of Fault Tolerant Production Infrastructures

Emmanouil Serrelis, Nikos Alexandris
*Department of Informatics,
University of Piraeus,
80 Karaoli & Dimitriou,
18534, Piraeus, Greece
serrelis@unipi.gr, alexandr@unipi.gr*

Abstract

Disaster Recovery Infrastructures, which have become a common aspect of all major IT infrastructures, could transform to Fault Tolerant Infrastructures in order to increase productivity, effectiveness and availability. This paper suggests a methodology for the transformation of High Availability Systems on which Disaster Recovery Infrastructures are based, to Fault Tolerant Production Infrastructures and establishes some Key Performance Indicators (KPIs) as a means to measure the effectiveness of the approach, adopting the principles of the Information Technology Infrastructure Library (ITIL) framework to a cost cutting, ecological and security-aware environment.

Keywords: *Fault Tolerance, Disaster Recovery, Availability, Change Management, Key Performance Indicators, ITIL*

1. Introduction

Events, such as the recent global economic crisis, have stressed the need to reduce the expenses in every aspect of the effected organizations, including IT expenses. Additionally, Green IT has become equally important being nowadays a major strategic objective. As stated in "Gartner's 10 Strategic Trends for 2009" *"For IT, green is everything, and that includes anything that can help cut the energy bill and reduce fuel use."* [1].

Towards this direction, IT should consider, among others, the change of its existing infrastructures and services. Some of the most eminent expenses of IT-related infrastructures are the costs related to Disaster Recovery Infrastructures. How could these infrastructures be optimized in both operational and financial terms? What could be the effect of

transforming a disaster recovery infrastructure to a more cost-effective infrastructure? These are the basic questions this paper addresses.

In a "disaster avoidance" rather than a "disaster recovery" approach, the high availability solutions aim to proactively protect business continuity by monitoring the key business functions and mission critical applications that are predetermined as business priorities. In a situation where an IT component fails, it can be dealt (manually or automatically) well before its failure impacts the business. Designing IT system components with the ability to remain operational in the event of a failure has an additional benefit; that is to increase IT efficiency through continuous architecture. Moreover, such "preserve and protect" measures can facilitate maintenance projects when malfunctioning or low performing components can be upgraded or repaired during a planned downtime.

The current paper, based on [2], suggests the transformation of the existing "cold-standby" Disaster Recovery Infrastructures, based on High Availability Systems, to Fault Tolerant Production Infrastructures, presenting the various differences of the two approaches. Using the theory of change management adding the necessary technical aspects, a specialized transformation strategy is proposed. The results of both the transformation methodology and the the adoption of Fault Tolerant Production Infrastructures are examined using the concepts of Information Technology Infrastructure Library (ITIL) framework and especially through the use of Key Performance Indicators (KPIs).

The remainder of this paper is organized as follows. Section 2 gives an overview of the technical and terms related to availability. Section 3 highlights the benefits of migrating to a Fault Tolerant Infrastructure, whereas Section 4 introduces some basic principles of transition strategies. The proposed transformation strategy is presented in Section 5 and the measurement methodology of the related benefits is addressed in

Section 6. This approach is critically evaluated in Section 7 and the paper concludes in Section 8.

2. Availability terms

Currently, there are several approaches for developing Disaster Recovery Infrastructures. All of them aim to protect organizations' most valuable assets: Data and Services. In order to be able to understand the differences between them, it is essential to define several terms related to availability.

Production infrastructures should be distinguished from Backup (or Disaster Recovery) infrastructures. The **Production infrastructures** aim to serve all daily services of the organization, whereas **Backup infrastructures** operate only if a disaster occurs. This fact would classify Backup infrastructures as a rather luxury solution which would justify their existence only in the case of an extreme, disastrous event.

Availability is the proportion of time that an application can be used for productive work, measured against the time that it must be functional. The time that the application must be functional or available to users is called "mission time," which may be quite different than 7 days per week - 24 hours per day (24x7) or 5 days per week - 8 hours per day (usual working hours) [3].

There are two factors that determine application availability. The first is the reliability of the components that comprise the application: namely, how often any of the consisting components fail. The second is how long it takes for the application to be restored once a failure has occurred. The components that comprise an application minimally include the server hardware, operating system and the application itself. The application may also be dependent on data storage devices, network access devices, databases, file systems, and other hardware and software components. The amount of time it takes to bring up an application after a failure depends on what it was that caused the application to fail. This time period is called Recovery Time Objective (RTO). If the application itself failed, recovery may be as simple as restarting the application on the same system. If, on the other hand, the application has failed due to a hardware failure, recovery can take a significantly longer time since it could involve [3]:

- Notifying the service provider of the failure
- Waiting for the arrival of the service technician
- Determining what component failed
- Replacing the failed component
- Rebooting the operating system
- Recovering the file system
- Recovering the database

- Restarting the networking software
- Restarting the application

Fault tolerance differs from high availability by providing additional resources that allow an application to "ride through" a failure without interruption [3]. Many of the high-availability solutions on the market today actually provide fault tolerance for a particular application component. Disk mirroring, where there are two disk drives with identical copies of the data, is an example of a fault-tolerant component. If one of the disk drives fails, there is another copy of the data that is instantly available so that the application can continue execution. However, once such a failure occurs, the system becomes vulnerable to the failure of the single remaining disk drive, which now has the only copy of the data and represents a single point of failure. Action should be taken as soon as possible to create a mirror of the remaining disk drive. However, this process may have a negative impact on system performance, depending on where the processing to re-mirror the drive takes place.

A fully fault-tolerant solution requires that all the resources that the application is dependent on are replicated, including the application process itself. This requires an independent processor (not part of the same – probably – symmetrical multiprocessing system) and a copy of the memory that the application uses. In the worst-case failure scenario, one in which the processor or memory fails, the replicated version of the application continues to execute. Other failures simply require the application to use alternate resources (disk drives, disk adapters, communications devices). As a result of this complete hardware and process replication, fault-tolerant systems are significantly more expensive than highly available systems. A fault-tolerant system would be used in a situation where no downtime can be tolerated at all, such as an air-traffic-control system, an emergency-response system or financial trading systems (during trade hours).

In evaluating a fault-tolerant system, particular attention should be paid to the repair process. While the system may be capable of proper operation through a failure, to ensure that a subsequent failure will not bring the system down, the failed component must be immediately repaired.

Load balancing is a technique (usually performed by special software or hardware mechanisms called load balancers) to distribute work between many computers, processes, hard disks or other resources in order to get optimal resource utilization and decrease computing time. It is also the ability to make several servers participate in the same service, performing the same tasks or supporting the same service [4]. Load

balancing can also be used to increase the capacity of a server farm beyond that of a single server.

This technique is seen as complementary of fault tolerant services since it frequently provides the ability to maintain unaffected services during a certain predefined number of simultaneous failures [5]. Also, traditional implementations of fault-tolerant platforms often involve duplicate proprietary hardware and software with complex binding and mechanisms. This causes higher implementation costs and longer periods of inactivity, which could not make such solutions attractive to short term investment and productivity management. The challenge is to provide fault tolerant infrastructures that would contribute to the daily business operations as well as to the failure or disaster situations.

3. The benefits of fault tolerant approach

Having examined the background information of availability, fault tolerance as well as the load balancing techniques, it is necessary to present the benefits of migrating to a Fault Tolerant Infrastructure that could be used for production purposes as well.

One of the fundamental advantages of High Availability Systems that are based on load balancing techniques is the protection of systems operation. In addition to that their presence can vastly improve the overall performance. *"Capacity on demand, load balancing, offline maintenance capacity and zero-point backup windows are all examples of the added value [that] a continuous architecture can produce".* [6]

And where there is added value, there could be

Return On Investment (ROI). Still, the quantification of ROI in that situation is not straightforward. Increases in efficiency - unless they result tangible savings like staff reductions or other avoided bottom line expenses - are often elusive to measurement. Nevertheless, they should be examined for any possible ROI contribution.

As the frequency of planned downtime is rapidly escalating due to the increasing number of applications development and the corresponding increase in upgrades and patches, the need to compress the downtime as much as possible has become even more pressing. For some companies, downtimes or even slow downs of 5-10 minutes can have a substantial affect on revenues. [7]

Other sources, such as [8], have shown that there can be a parallel use of a segment of a disaster recovery infrastructure in order to tackle extreme attacks. Expanding this idea, the benefits described above are multiplied when segments of the primary site work together their equivalent ones in the disaster recovery infrastructure, in load balancing mode.

It is, therefore, evident that organizations have more that one reasons to transform their existing implementations and select the fault tolerant production infrastructure solution. The very same tools that are used for high availability such as clustering, volume management and load balancing, can automate key procedures that would decrease the length of the downtime window as well as the cost of downtime administration. Savings from these types of value-adding features are very real and can help reduce or eliminate costs associated with planned downtime. In

	High Availability	Fault Tolerance	Fault Tolerant Production Infrastructures
Purpose / Impact	To enable faster recovery of lost data and stalled business operations in the event of a disaster.	To proactively avoid some types of disasters before they occur.	Proactively avoid most types of disasters before they occur. Increase the productivity of the organisation's IT infrastructures.
Cost	Tangible IT investment.	Tangible IT investment. ROI can be measured in most of the cases.	Tangible IT investment with measurable ROI.
Benefits	Faster time to recovery, lower lost revenues/productivity, reduced recovery costs.	Reduced probability of disaster occurrence, improved operational efficiencies, reduced planned downtime windows and costs.	Minimal probability of disaster occurrence, improved operational efficiencies, reduced planned downtime windows and costs.
Return On Investment (ROI)	Soft since the benefits are only realized in the event of a disaster.	- Reduction of disaster probability is soft. - Reduced planned downtime generates real savings in IT costs through automation of procedures that can reduce the need for IT resources, eliminate human error and save in lost business from shorter downtime.	- Reduction of disaster probability is soft. - Improved operational productivity can have direct impact on revenues and expenses and could contribute tangible cash through sales and savings. - Reduced planned downtime generates real savings in IT costs through automation of procedures that can reduce the need for IT resources, eliminate human error and save in lost business from shorter downtime. These are hard, tangible benefits driven by avoided revenue loss and reduced operational expenses.

Table 1 – Comparison of Availability Solutions

addition, the outage window is compressed so that business functions can continue with little or no interruption. Table 1 concentrates the above points.

As far as the environmental requirements and international directives are concerned, Fault Tolerant Production Infrastructures could greatly contribute to Green IT by reducing the power demands needed for operation to one site distributing power (as well as the related CO₂ emissions) to multiple geographically dispersed IT sites which were consuming power anyway.

4. Transition strategies principles

Before presenting any transition strategies, some basic transition questions should be asked:

- How the transition should be planned and implemented?
- Which parts of the organization should be integrated into the fault tolerant production infrastructure solution?
- Who should be involved in the transition project?
- What this transition will cost in terms of money? Will the final outcome worth the transition costs?
- When is the right time to perform such a transition?

4.1 The change management theory

*“The concept of **change management** describes a structured approach to transitions in individuals, teams, organizations and societies that moves the target from a current state to a desired state”* [9]. This is exactly the situation one deals when transforming one IT Infrastructure to another, so it is considered very useful to see which points are suitable and applicable the situations presented in previous paragraphs. There are several theories regarding change management. The most popular ones are presented in [10] and [11]. However, as [12] points out the first question of someone diagnosing a problem is “what changed?” With a change management process in place, that question is far easier to answer. Change management is a process made up of people, software, and procedures. When properly followed the process results in many benefits including increased staff efficiency, reduced server and network device downtime and reduced Mean Time To Recover (MTTR). Change management can also bring about positive impacts on security, providing trusted audit data and increased control over ad-hoc changes, all of which lead to reduced IT costs.

Change management is critical for maintaining highly reliable systems that meet the defined service levels of the organization. To this end, best practice organizations are pushing all changes back into the build and test phases such that only rare emergency changes are actually performed on production systems. The whole network device change process must become formalized and incorporate security, testing, and documentation. The organization must ensure that appropriate preventive, detective and corrective controls exist in order to meet the challenges of regulatory frameworks, such as SOX, as well as to improve operational efficiencies.

Forrester’s “Best Practices For Infrastructure Change Management: Regain Control of Runaway IT Infrastructures,”[13] boldly states *“In IT, change is an engine of progress, as well as a source of doom... While application software change control is a relatively mature process, many organizations implement infrastructure change manually, relying primarily on the IT staff’s knowledge and expertise. This ad hoc process is nearing its limits in today’s complex environment, where the risks inherent to changes multiply”*.

Automating the change management process means addressing the six steps in an effective change management process:

1. A change is requested
2. Requested changes are reviewed, the impact assessed, and resources estimated and assigned
3. Changes are either approved or rejected
4. If approved, changes are developed and tested in a preproduction environment
5. Changes are implemented into production
6. Changes are verified and reconciled by someone else in the organization

The last step is the critical missing piece in most organizations. In order to effectively manage change, it is needed to complete the change process circle. This can be done by conducting a final verification that the requested change was implemented properly, verifying that change was implemented on all target systems and finally to have the ability to see if the change control process was bypassed. Without this step, the change management loop remains open ended, and it impossible to tell the difference between authorized, successful changes and unauthorized (or unsuccessful) changes.

The results are in and the experts agree that reducing service outages from human error through automated processes provide IT savings and a more efficient business. Eighty percent of IT budgets is used to maintain the status quo. By implementing

enforceable change management process, IT gains control of the infrastructure. By gaining visibility in what changed, IT closes the loop on change management and improves availability, improves audit performance, and lowers IT operational costs.

4.2 The technical experience

This section includes industrial experience as referenced in [14] and [15]. In today's IT infrastructures, applications are interrelated and integrated with others more than ever. At the same time, shared infrastructure elements are more common, while managing a maintenance window for each application can be exceedingly complex. However, a common maintenance window for infrastructure activity can be beneficial.

The technical experience of the current status, as highlighted above, has taught some basic lessons. The first one is that an organization should always aim to reduce unplanned downtime, since it costs on money and reputation.

The second lesson comes from [13] which states that *"80% or more of unplanned downtime is the result of People and Processes, not hardware or O/S failures"*. This means that this percentage is caused things like data corruption, application failures, software failures, errors in configurations, scheduling errors, operator errors, delayed batch jobs etc. So, in order to deal with these causes of downtime, an organization should provide funds and time in people (Proper staffing and training), problem management, event management, job scheduling, test and time recovery scenarios (in the form of production readiness reviews), Application and capacity planning and last, change management which is the area that is discussed in this paper.

The third lesson is more technology-related and mentions that an organization should minimize single points of failure, take care of environmental, facilities and network threats, make use of load balancers, redundant dispatchers, replication, cloning, RAID technologies, such as mirroring, striping and hot swap availability. Additionally an organization should plan to operate using High Availability, or even better Fault Tolerant solutions with clustering and auto fail over capabilities.

In order to implement infrastructures that could deal with the issues above and make use of the technologies mentioned, the organization should understand the application architecture and constraints as well as to understand all application dependencies and interrelationships to needed components, whereas they

should reduce any batch interference (delays, lockups etc.).

Furthermore, they should manage other planned changes, by developing suitable infrastructure and facility work and performing appropriate hardware, operating system, database, application changes and upgrades. Another need that should be covered is the need for proper infrastructure test environments. Within this framework the organization should aim to common maintenance windows, expecting increased coordination as well as staff overhead.

Taking all these lessons into account an organization should try to follow the following rules within the plans for implementing Fault Tolerant Production Infrastructures. Firstly, they must integrate application availability in their design, since this can be hardly be improved in later phases. Secondly, there should be a well planned transaction queuing as well as a highly optimized batch processing. The third rule is to set the requirements for scheduling and availability early in the design phase. Fourthly, an organization should choose to serve only business-critical functions with high-cost Fault Tolerant infrastructures, having in mind that these kind of infrastructures cost about 3.5 times as much as a standard infrastructure. [16]

5. The proposed transformation strategy

Taking into account all above sources, the **proposed transition strategy** combines the change management theory and the technical experience. There are seven phases to complete the transformation from the High Availability Standby Systems to the Fault Tolerant Production Infrastructures. These are:

- Phase 1: Definition of the transformation scope
- Phase 2: Categorization of System groups
- Phase 3: Application Analysis
- Phase 4: Process Analysis
- Phase 5: Cost Analysis
- Phase 6: Business Decision
- Phase 7: Execution of Transformation

5.1 Definition of the transformation scope

As can it be easily understood, a problem well defined is a problem that can be solved more easily. During the first phase of the transformation, the organization should decide which systems are candidates to transform. Thus, for each application area, it should be determined what the transformation scope is, with the correct user representative(s). At the same time, the schedule goal and the availability goal should also be agreed. Since it is more costly to re-change any infrastructures, it is important to determine and design

schedule and availability up front, just like any other application functional requirement.

5.2 Categorization of System groups

The second phase aims to categorize the system groups. For example an organization could distinguish between Business Support Systems, Operational Support Systems, Self Service / e-Commerce, Management Support Systems. This categorization will give the organization a rough idea on how these systems should be implemented in terms of availability, enriching the decisions taken in the first phase.

5.3 Application Analysis

During the third phase, the organization should understand each application's architecture, special constraints, "release tolerance" and flexibility to change. Additionally, the applications dependencies on other applications and components should be gathered, along with architecture diagrams and data flows. Finally, decisions on the whether the applications' modification for Fault Tolerance should be in-house or outsourced should be made.

5.4 Process Analysis

In this phase questions such what is the current Standing Operating and Testing Procedure should be answered with respect to technology. The current availability of each function/application should also be identified. Furthermore, what can the organization expect with existing budget. In order to answer these questions more easily, metrics related to availability, efficiency and performance have to be established. The Final of this phase is to identify root causes of unplanned downtime.

5.5 Cost Analysis

The most important phases are phases 5 and 6. This is where is actual decision on whether the transformation should be executed or not is taken. In the cost analysis phase, the basic question that the executive level will pose is what improvements can the organization make from existing budget. In order to answer this properly, the organization has to consider to invest in the right areas to expand schedule and availability. Additionally, the organization has to know costs to expand schedule beyond baseline to meet goals as well as the costs to increase availability beyond

baseline to meet goals. At this phase involvement from all areas of the organization should be encouraged.

5.6 Business Decision

This is the last phase before the actual execution of the transformation. During this phase, the organization should develop a consistent approach to weigh the business benefits against the cost, while maintaining focus on the business problem, which is to increase the availability and the usability of its systems. Towards that decision, a Steering Committee or the business owners of the applications need to determine the business need. Since it is difficult to cost and plan for applications individually an accurate categorization would be very useful. At all times, the decision committee should be aware of the transformation sponsor capabilities and wills that would also be effected by any potential future expenses that a Fault Tolerant Production Infrastructure may imply.

5.7 Transformation Execution

The final phase of the suggested transformation strategy is actual execution of the transformation. In order to achieve this, the organization, and especially the people involved and affected, should be committed to the project. A detailed and realistic definition of the resources in terms of people and budget is necessary. Another very important issue is to define the owner of the new infrastructure in order to establish a common communication point that could manage, adjust, develop, document the transformation plan, with goals, activities, responsibilities, dates, etc. Finally, the organization should measure the actual benefits against the initial goal, for use in future or parallel transformation projects.

6. Measuring the benefits

The application of the approach, as described above, has been demonstrated in past [17], resulting a rather successful outcome. However, as pointed out in that attempt, in order to provide more concrete evidence of the applicability of the methodology, some formal metrics of the methodology should be established. These metrics should enable consistent measurement of resources, time and cost.

Towards that direction, the Information Technology Infrastructure Library (ITIL) framework has been examined for suitability. ITIL is a globally accepted set of best practices used for the management of IT environments. In order to improve the level of IT services provided in an organization, the ITIL

framework suggests the adoption and combination of methodologies, tools, metrics and roles.

As it can be understood, the adaption of ITILs metrics to serve the needs of the transformation methodology described above, would strongly support the applicability of the methodology. Additionally, such an adaptation could also provide a known interface for people who are familiar with the ITIL framework as well as ITIL's measuring tools.

The following paragraph presents the foundations for the application of ITIL-based metrics to the transformation approach.

6.1 Key Performance Indicators (KPIs)

In ITIL terminology, KPIs are *"financial and non-financial metrics which help organizations to define and measure progress toward organizational goals"* [18]. KPIs main goal is to review the current state of an organization and provide the basis for the prescription of a course of improving actions. In order to obtain a more solid view of the organization's state, KPIs should be monitored in real-time, a process otherwise known as Business Activity Monitoring (BAM). Common uses of KPIs include the measurement of intangible benefits or values such as leadership development, engagement level, service delivery, and satisfaction rates. Being able to grasp such aspects, managers typically tie KPIs to organization's strategic management.

The selected KPIs may differ depending on the nature of the organization and the organization's objectives. In any case, their proper usage could assist an organization to measure progress towards their organizational goals, especially goals which include difficult to quantify knowledge-based processes.

Any KPI is a part of a "measurable objective" which is made up of a direction, KPI, benchmark, target and time frame. For example: "Increase Average Storage Utilization per Server from 20% to 60% by the end of the year 2010". In this case, "Average Storage Utilization per Server" is the KPI.

KPIs should not be confused with Marketing-related Critical Success Factors. For the example above, a critical success factor would be something that needs to be in place to achieve that objective; for the previous example, a file archiving software tool.

Performance indicators should also differ from business drivers & aims (or goals). A financial institution may consider the "increase rate of deposits" as a Key Performance Indicator which might help the institution understand its position in the market, whereas a telecommunications company could consider the "percentage of successful call attempts from its customers" as a potential Market-related KPI.

Nevertheless it is necessary for an organization to at least identify its KPIs. The basic rules for identifying KPIs are:

- To have pre-defined business processes.
- To have clear requirements for the aims and the performance for the business processes.
- To have a measurement that could quantify and qualify its results and compare these with the previously set goals.
- To examine the variances and adjust any processes or resources needed to achieve short-term goals.

The definition of any KPI should apply all of the following characteristics:

- Specific, so that it should not be confusing with other KPIs
- Measurable, so that it should be feasible to measure it or calculate using a specific measurement unit
- Achievable, so that it should be easy to obtain the necessary information
- Relevant, so that it should directly connect to the business objective
- Time-bound, so that it should take into account time constraints, in order to be able to tackle any issues related to time depended results and filter them.

Key Performance Indicators in practical and strategic development terms are objectives to be targeted that will add value to the business.

Having seen how KPIs are defined and used within a generic organization, it is now possible to use these principles within an IT infrastructure environment, where some more specific KPIs could be defined. [19] These KPIs will be used to represent the benefits from the adoption of Fault Tolerant Production Infrastructures as well as the benefits from the usage of the suggested transformation approach.

6.2 Generic KPIs for the adoption of Fault Tolerant Production Infrastructures

Although the transformation approach is quite specific as far as the transformation steps are concerned, each IT infrastructure involves different modules, processes and systems. Thus, the KPIs chosen to be presented in this paper could only be considered as a first, generic, set of KPIs. Additional KPIs can and should be considered in order to match the specific needs of an organization. However, the total number of the KPIs used to measure the success of the adoption of Fault Tolerant Production

Infrastructures should not be too large since this may affect the performance levels of the infrastructure.

The generic set of KPIs for measuring the success of the adoption of Fault Tolerant Production Infrastructures could be divided into the technical and business related KPIs. These two KPI categories are not directly related to each other. They aim to point up different aspects of the Fault Tolerant Production Infrastructures and measure the technical and business benefits of its adoption. It should be made clear that there is no need for conciliation, combination or synchronization between the results of these two categories.

Nevertheless, the measurement of the following KPIs should take place before and after the transformation, so that the comparison can confirm the benefits of the Fault Tolerant Production Infrastructures.

Technical KPIs

1. **Usable Storage in IT Site(s):** This KPI is the storage that can be used to store data in an IT site after any technical overhead, such as RAID configurations of storage boxes. The IT site can be Primary, Secondary, Other, Disaster Recovery or any IT Site. The measuring unit is in MBs or GBs.

2. **Average utilization of total Processing Power capacity Average in IT Site(s):** This KPI is the average percentage of utilization of processing power of all systems in a specific (or all) IT Site(s) during the measurement period. The IT site can be Primary, Secondary, Other, Disaster Recovery or any IT Site. The measuring unit of utilization is a percentage. The measuring unit of processing power is Million Instructions per Second.

3. **Average network throughput between servers and clients:** The term "Throughput" refers to the performance of data transmission, and is measured by characters actually transmitted or received during a certain period of time. Throughput is usually measured in bps (bits per second). A better (higher) throughput to the clients could signify the existence of a better infrastructure.

4. **Average Disks I/O in central storage in IT Site(s):** This particular KPI reveals the average percentage of storage disks utilization of all systems in a specific (or all) IT Site(s) during the measurement period. The IT site can be Primary, Secondary, Other, Disaster Recovery or any IT Site. The measuring unit is bps (bits per second).

5. **Average Memory utilization in IT Site(s):** The memory utilization expose the average percentage of memory utilization of all systems in a specific (or all) IT Site(s) during the measurement period. The IT site

can be Primary, Secondary, Other, Disaster Recovery or any IT Site. The measuring unit is a percentage.

6. **IT Site power usage effectiveness:** This KPI is calculated by dividing the total power usage of an IT Site by the power usage of IT equipment (computer, storage, and network equipment as well as switches, monitors, and workstations to control the IT Site). The IT site can be Primary, Secondary, Other, Disaster Recovery or any IT Site.

7. **Systems Footprint in IT Site(s):** The footprint represents the physical area that the systems occupy and is measured in square meters or square feet. A change in the measurement of this KPI would support the benefits earned by the adoption of the Fault Tolerant Production Infrastructures. Similarly to other KPIs the IT site can be Primary, Secondary, Other, Disaster Recovery or any IT Site.

8. **% of production servers located in Primary / Secondary IT Site(s):** This is one of the most profound benefits of Fault Tolerant Production Infrastructures. It appears as a percentage of production servers located in a particular IT Site (Primary or Secondary) over the total number of servers in all IT Sites.

Business KPIs

1. **Planned Downtime of offered business services:** Planned downtime is downtime of any business service caused by scheduled for system or application maintenance. It is measured in minutes or hours per year.

2. **Unplanned Downtime of offered business services:** This is the amount of downtime of any business service arising from reasons other than maintenance. It is measured in minutes or hours per year.

3. **Recovery time of business critical services:** This KPI presumes that there has been decided which are the business critical services. The recovery time is the duration of time within which the business critical services can be restored after a disaster in order to avoid unacceptable business consequences. It is measured in minutes, hours or days.

4. **Operational Expenses of IT Division:** The Operational Expenses, measured in any currency, are the yearly running costs of any organization, or in parts of the organization, such as IT Division. A decrease in these could signify a better usage or management of the existing resources.

5. **Capital Expenses of IT Division:** Opposite to the previous KPI, Capital Expenses are the one-off costs of products and non-consumable parts. It is measured in any currency and could relate to the financial benefits

of the adoption of Fault Tolerant Production Infrastructures.

6. Cost of Recovery of new business services: This is a very important KPI since it could depict the low cost expansion capabilities of Fault Tolerant Production Infrastructures. The measuring unit is any currency.

7. Satisfaction rate by IT staff (System owners): This is a qualitative measurement of the satisfaction of the IT staff. The staff's satisfaction rate could be based on periodic surveys of employees after a reasonable period of infrastructure maturity time such as 6 months. The maturity time could minimize non in type negative reactions, caused by staff's natural resistance to change [20]. This KPI is measured as a percentage of positive reactions.

8. Percentage of satisfaction by Business staff (Business owners): In the same way as in the previous KPI, this measurement is related to the satisfaction of the Business Staff which may (or may not) have a different opinion on the benefits of the implemented infrastructure. This KPI is measured as a percentage of positive reactions.

9. Average frequency of updates of disaster recovery plans: This KPI should portray staff awareness on the updating the Disaster Recovery plans. Since the Fault Tolerant Production Infrastructures amplify the role of Disaster Recovery IT Sites, it should be expected that this update frequency should be increased. It is measured in days.

10. % of growth of IT budget: An unusual growth of the IT budget may entail some form of relation to the new Infrastructure architecture.

11. New Systems Procurement rate (as % of existing systems): This KPI should confirm that the procurement of new systems should be less frequent since extra resources and capacities would be freed up (mainly in the Primary IT Site) after such a transformation.

12. Average time to provision new systems: This is the average time needed to provide a new system to an application or system owner. The time starts counting when the request is send and ends when the system is handed over. The measurement time is in minutes, hours or days. A more dynamic infrastructure, as Fault Tolerant Production Infrastructures aim to be, should decrease that time.

13. Average time to provision new business services: This KPI differs from the previous one for the reason that it also includes processes and people needed to provide the new business services. It is measured in minutes, hours or days.

6.3 KPIs for the usage of the suggested transformation methodology

The KPIs that could be used for measuring the benefits from using the suggested approach are less dependent on the IT infrastructure and business services of the organization that has chosen to use this approach, than the KPIs described in the previous paragraph. Again, the number of the selected KPIs should be limited to a level that would not effect the actual progress and effectiveness of the methodology.

Since the core of the transformation approach is a change management set of processes, the consequent KPIs for measuring the success of the suggested transformation methodology are solely related to project management metrics. The measurement of the following KPIs should take place during the transformation, and be compared to similar projects that have been (or will use) different transformation methodologies. These projects may also originate from outside the implementing organization.

Project Management KPIs

1. Return on the transformation process

Investment: This KPI illustrates the main idea behind this paper. It is a predominantly hard KPI to measure since the actual Return cannot be directly calculated. However all other KPIs mentioned in this paragraph could be used as input to its calculation. It could be measured in any currency or in time units such as days, weeks or months. When measured in time units the Return represents the time gained by using the proposed transformation methodology.

2. Total Time of transformation process: This is the time period the transformation project runs and includes all seven phases of the proposed transition strategy described in paragraph 5. It is measured in days, weeks or months.

3. Utilization rate of human resource for project purposes: This is the percentage of the time that a worker will dedicate to the transformation project in relation to its total time. It is similar to Full-time equivalent (FTE) which is a way to measure a worker's involvement in a project and is used by many organizations worldwide.

4. Downtime of Business Services due to transformation project: Some of the transformation phases described before could effect the operation of some Business Services and thus their availability. Less production outage time for each Business Service means less lost income of the organization and more value for the transformation methodology. It is measured in minutes, hours or days.

5. Total cost of transformation project: This is very important since the transformation project should be significantly less expensive than the expected earnings. It is measured in any currency.

6. Number of people involved in transformation project: This is also important in order to be able to appreciate the staffing needs of the project.

7. Percentage of administrative activities related to the transformation project: This is a project management quality KPI. It presents the number of administrative activities for the transformation project in relation to the total activities of the project which also include implementation activities. It is measured as a percentage.

8. Budgeted Cost of Work Scheduled: This is the sum of the budgets of the activities that were planned or scheduled to be completed, otherwise known as "planned value". It is measured in any currency.

9. Budgeted Cost of Work Performed: This KPI, measured in any currency, is the planned or scheduled cost of activities that were completed, also known as "earned value". It is measured in any currency.

10. Actual Cost of Work Performed: It is the sum of actual costs of activities that are completed. It is measured in any currency.

11. Schedule Performance Index: This is calculated by the use of the previous KPIs. It is the division of "Budgeted Cost of Work Performed" by the "Budgeted Cost of Work Scheduled".

12. Cost Performance Index: This is calculated by combining two of the previous KPIs. It is the "Budgeted Cost of Work Performed" divided by the "Actual cost of Work Performed".

12. Cost Schedule Index: This is the "Cost Performance Index" multiplied by the "Schedule Performance Index". The Cost Schedule Index measures the likelihood of recovery for any project that is late and/or over budget. The closer the index is to 1, the more likely the project's can be recovered from it's deviation to the original baseline. This can be useful for any organization that would decide to apply the proposed transformation methodology.

13. % of time coordinating project: This is an efficiency related KPI for the methodology and is represented as a percentage of time (in man hours) used to coordinate project relative to over the total time used to implement (and coordinate) the project.

14. % of milestones missed: Percentage of milestones recorded in all processes and phases as missed.

15. Number of incidents due to transformation project: In theory the transformation like any other planned change should not cause any incidents. However, a more practical evaluation of the methodology should also measure also the number of incidents caused by the methodology in relation to the

total number of incidents. In any case the changes should not cause more than a upper percentage of all the incidents.

16. Average rework per phase after implementation of each phase: This is a significant measurement of the quality of the Analysis and Design processes that is methodology involves. If the rework per phase is low then this could be an indication that the methodology is providing a solid foundation for the transformation to Fault Tolerant Production Infrastructures. It is measured in man-days, man-weeks or man-months.

7. Evaluation and Future Improvements

The suggested transformation is not a new idea. However the actual application of the change management theory to the specific transformation tasks which are based more on practical Management experience than Information Technology theory is a new addition to the arsenal of an IT manager.

The presented benefits range from low-level technical benefits to high level financial benefits as well as the contribution to Green IT Infrastructures.

The theory has been supported by establishing some ITIL-based metrics (KPIs) in order to challenge and prove its applicability. These metrics aim to measure, test and evaluate practically the proposal in a formal, accurate and consistent manner.

Using KPIs, such as the ones proposed, the IT managers are able to evaluate the outcomes of the transformation of High Availability Systems to Fault Tolerant Production Infrastructures. Furthermore, the transformation methodology itself can also be evaluated in terms of technical and business benefits.

As a more general remark, it should be pointed out that the transformation methodology can also be seen as part of ITIL's set of concepts and policies for managing change in infrastructures and services.

Following that way of thinking, future research should include a more thorough analysis of the relationship with ITIL's change management process, aligning the transformation of High Availability Systems to Fault Tolerant Production Infrastructures with the related ITIL's core components, namely the Service Strategy, Service Design and Service Transition.

A more detailed analysis of the selected KPIs and their usage can also offer more information on the prospective users of the methodology. Usage-related factors could include supplementary information on the measuring time, measuring frequency as well as number of measuring repetitions.

There are also other extensions to the proposed methodology, which can enforce the relationship with

ITIL's practices. These extensions might engage the use of Balanced Scorecard as well as the adjustment of IT organization Service Catalog. The Balanced Scorecard suggests that an organization should be viewed from four different perspectives (the Learning & Growth Perspective, the Business Process Perspective, the Customer Perspective and the Financial Perspective). Additionally, the Balanced Scorecard suggests the development of some other metrics, the collection of related data and the proper analysis of the perspectives' relations. The Service Catalog is a list of services that an organization provides to its employees or customers. Each service within the catalog may well include:

- A description of the service
- Timeframes or service level agreement for fulfilling the service
- Who is entitled to request/view the service
- Costs (if any)
- How to fulfill the service

Yet, this paper is to be perceived as a packaged proposal that includes a proposition for the target infrastructure, the transformation methodology as well as the metrics for the efficiency of both the infrastructure and the methodology.

8. Conclusion

This paper has suggested the transformation of the existing "cold-standby" Disaster Recovery Infrastructures, based on High Availability Systems, to Fault Tolerant Production Infrastructures. The various differences of the two approaches have been presented and there are clear indications that organizations can benefit from transforming their existing implementations and selecting the Fault Tolerant Production Infrastructure solution.

The business and technical results of the transformation as well as the effectiveness of the methodology are measured through the use of relevant KPIs using the ITIL framework. Savings from these types of value-adding features vary from case to case but the use of this methodology makes possible the reduction or elimination of costs associated with planned (and most unplanned) downtime. In addition, the outage window is compressed so that business functions can continue with little or no interruption. The suggested transformation strategy has used the theory of change management adding several technical aspects. This transformation strategy can be considered as a strong support tool in order to make the transformation less costly, less time consuming as well

as to effectively integrate people, systems and procedures.

9. References

- [1] P. Thibodeau, "Virtualization Tops Gartner's 10 Strategic Technologies for 2009", *Computerworld*, <http://www.cio.com/article/print/454906>
- [2] Em. Serrelis, N. Alexandris, "From High Availability Systems to Fault Tolerant Infrastructures", *IEEE Computer Society Press*, ICNS, Athens, 2007
- [3] "High Availability: A perspective", *Gartner Research*, ID Number: DPRO-90193
- [4] W. Tarreau, "Making applications scalable with Load Balancing", http://1wt.eu/articles/2006_lb/index.html
- [5] "Highly Available Embedded Computer Platforms Become Reality", *International Engineering Consortium*, http://www.iec.org/online/tutorials/ha_embed/topic01.html
- [6] K. Miller, "Don't Recover-Failover", *DM Direct*, Oct 2004
- [7] S. Atwood, "Planned Downtime", *DM Direct*, Veritas Software, October 2004
- [8] Em. Serrelis, N. Alexandris, "Disaster Recovery Sites as a Tool of Managing Extreme Attacks", *IEEE Computer Society Press*, ICISP, Cap Esterel, 2006.
- [9] "Change Management", *Wikipedia*, http://en.wikipedia.org/wiki/Change_management
- [10] J. Martin, "Organisational Behaviour", *Thomson Business Press*, 1998, ISBN 1-86152-180-4, pages 575-600
- [11] L. J. Mullins, "Management and Organisational Behaviour", 5th Edition, *Pitman Publishing*, 1999, ISBN 0-273-63552-2, pages 821-830
- [12] "Five basic principles of Change Management", <http://www.teamtechnology.co.uk/changemanagement.html>
- [13] J. P. Garbani, "Best Practices For Infrastructure Change Management: Regain Control Of Runaway IT Infrastructures", *Forrester Research*, 25-3-2004, <http://www.forrester.com/Research/Document/Excerpt/0,7211,34048,00.html>
- [14] "Microsoft Operations Framework 4.0", <http://www.microsoft.com/technet/solutionaccelerators/cits/mo/smf/smfchgmg.aspx>
- [15] "Change Management (ITSM)", *Wikipedia*, http://en.wikipedia.org/wiki/Change_Management_%28ITSM%29
- [16] D. Scott, Y. Natis, "Building Continuous Availability Into E-Applications", *GartnerGroup*, COM-12-1325, 29/9/2000
- [17] Em. Serrelis, N. Alexandris, "Fault Tolerant Production Infrastructures in Practice", *IEEE Computer Society Press*, PIMRC, Athens, 2007
- [18] F. John Reh, "Key Performance Indicators – What are Key Performance Indicators or KPI", *About.com*, <http://management.about.com/cs/generalmanagement/a/keyperfindic.htm>
- [19] "KPI Library", <http://kpilibrary.com/>
- [20] A. J. Schuler, "Overcoming Resistance to Change: Top Ten Reasons for Change Resistance", http://www.schulersolutions.com/resistance_to_change.html