

Deployment and Operations for Software Engineers 2nd Ed

Chapter 15—Disaster recovery



Outline

Disaster recovery plan

RTO and RPO and Tiers of systems

Primary and secondary data centers

Data Management

Software management

Failover



Disaster recovery plans

- A disaster is an event such as a flood, earthquake, hurricane, or tornado, that renders an entire data center or availability zone unusable.
- “Business continuity” refers to the complete set of concerns when a disaster occurs including people and customers
- “Disaster recovery” refers to information technology concerns.



Risk mitigation

- All business continuity planning is a risk-mitigation activity.
- Risk can be quantified: The risk of an event is the probability of the event occurring multiplied by the loss incurred if the event does occur.
- Insurance is an appropriate analogy.
 - Having no insurance can be catastrophic if a disaster event occurs.
 - Having too much insurance is a financial drain



Goal of a disaster recovery plan

- Restore normal operations as soon as possible.
 - Within cost constraints
 - 10% of operations costs is a disaster planning target for many organizations.
- Targets of recovery must be specified



Outline

Disaster recovery plan

RTO and RPO and Tiers of systems

Primary and secondary data centers

Data Management

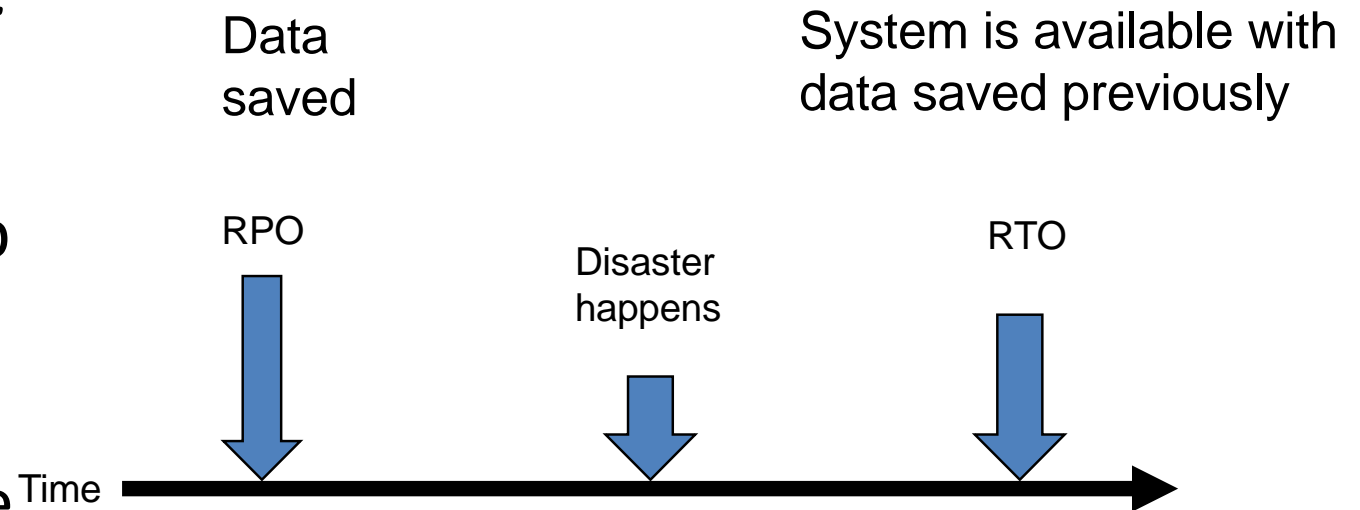
Software management

Failover



RPO and RTO

- *RPO—Recovery point objective.* the time interval between backups. Data stored in the database between last backup and disaster will be lost.
- *RTO—Recovery time objective.* The maximum amount of time that the system can be down before customers can access it again.





Prioritizing systems

- Your organization may have hundreds if not thousands of different systems.
- Not every system is of equal importance in terms of disaster recovery.
- The first step in disaster recovery planning is to prioritize your systems.
- For convenience, priorities are set in terms of tiers.
- All systems in a given tier have the same disaster plan.



Tiers

- A common model divides all systems in your organization into four tiers.
 - Tier 1 will have, nominally, a 15-minute RPO and RTO, These are the mission critical systems.
 - Tier 2 will have two hours. These are important support systems,
 - Tier 3 will have four hours, These are less important support systems.
 - Tier 4 will have 24 hours. These are everything else.



Discussion questions

1. Identify two Tier 1 systems for your organization
2. Identify two Tier 4 systems for your organization.



Outline

Disaster recovery plan

RTO and RPO and Tiers of systems

Primary and secondary data centers

Data Management

Software management

Failover



Data centers

- Your systems execute on hardware inside a data center.
- A data center is a physical location with
 - roughly 100,000 computers
 - backup power source (for some period),
 - physical security and access controls,
 - fire-suppression system,
 - air-conditioning system.



Secondary data center

- If a disaster occurs, your organization's only way to restore operation of software systems is to run them in a secondary data center,
- This secondary data center should be geographically located far enough away from your primary data center so that both data centers will not be affected by a single disaster.
- To restore a system to service (failover) at the secondary data center you need
 - all the software that the system comprises
 - valid data for the system



Disaster recovery strategy

- A disaster recovery strategy has three parts:
 - identify a secondary computing facility,
 - provide all needed software
 - provide data.
- Each part changes at a different rate.
 - The selection of a secondary computing facility might be reviewed once or twice a year.
 - Your software could change many times every day.
 - The data for your system will change nearly constantly



Warm and hot second locations

- Warm secondary location has space, power, and cooling. Depending on your cloud provider, this could be achieved by reserving space in a different availability zone or by expanding into a different zone if a disaster occurs.
 - No software is installed in a warm location. It just provides an available data center.
- Hot secondary location has all the features of a warm location, with the addition that the most current versions of software systems and infrastructure services are loaded and executing.
 - A hot location will have software but not have the most current system data.



Mirrored secondary location

- Mirrored locations have identical software and data at two or more data centers.
- In this configuration, the RTO and RPO can be truly “0.”
- This approach requires that the system be designed from the start to support mirrored data.



Discussion questions

- What are the costs in having a mirrored secondary data center?
- Can you use different secondary data centers for different tiers? Why would you do that?



Outline

Disaster recovery plan

RTO and RPO and Tiers of systems

Primary and secondary data centers

Data Management

Software management

Failover



Tier 2-4 Data Management

- Two issues exist in management of tier 2-4 data
 - Backup frequency
 - Storage media and recovery time



Back up frequencies

- When a disaster occurs, you must assume that all the data at the affected data center is lost. To have access to that data for recovery, you must have made a backup copy of that data.
- The frequency of the backups is determined by the RPO.
 - Tier 4 systems should be backed up once a day,
 - Tier 3 systems backed up every 4 hours
 - Tier 2 systems every two hours.



Storage media options

- Online storage to disk
 - Storage at your secondary location must always be available and accessible over the network from your primary data center
 - Replication is a one-way copy from disk storage in one data center to disk storage in another data center.
- Off line storage to tape
 - Replication is copying to tape, delivering the tape to the secondary data center, and copying to disk at the secondary data center



Why would you use tape?

- Two reasons for using tape
 - If the secondary data center location is not always available
 - If your data set is too large to be sent in a timely fashion over the internet.



Tier 1 data management

- To achieve Tier 1 RPO, you must keep an up-to-date copy of the system data at the secondary data center location.
- The replication can be
 - bidirectional (often called a *master-master* configuration), where either data set can be updated and the update is propagated to the other copy of the data set.
 - The one way (often called a *master-follower* configuration) where one data set is updated and the other maintains a copy.



Unreplicated data

- Not all data must be replicated.
 - *Session data*. In this case, a user would be required to log in again in the event of a disaster. Whether this is acceptable for your system is a business decision not a technical one.
 - *Infrequently changed data*. Items such as the static portion of web pages, e-shopping data, videos, and pictures are changed infrequently if at all.



Big data

- Big data is any data set that is too large to be backed up.
- Big data is split into chunks called *shards*.
- Each shard is replicated to create several copies that are distributed across multiple locations,
- Distributed coordination mechanisms are built into the database system keep all copies of a shard consistent across all locations.



Discussion questions

1. What is an incremental backup and when is it used?
2. Banks are moving from private clouds to public clouds. How does that affect their disaster recovery plans?



Outline

Disaster recovery plan

RTO and RPO and Tiers of systems

Primary and secondary data centers

Data Management

Software management

Failover



Software at the secondary location

- Tiers 2-4.
 - Use a configuration management system to keep the software up to date.
 - It must be identical in version/patch number to the primary. Dependencies must also be identical
- Tier 1
 - Updated as a portion of the normal deployment process. Thus, when the primary is changed, the secondary is also changed.
 - should execute silently—it should not allow any output to affect the behavior of the primary location or modify the backup database.



Licenses and keys

- If your system, infrastructure services, or development tools use any commercially licensed software, the software in the secondary data center must have the appropriate license keys.



Discussion questions

1. What are the software costs associated with keeping a mirrored or hot secondary site?
2. What are the database considerations for maintaining consistency between the primary and secondary data centers?



Outline

Disaster recovery plan

RTO and RPO and Tiers of systems

Primary and secondary data centers

Data Management

Software management

Failover



Failover

- Failover is the transfer of activity to a secondary data center.
- Three activities
 - trigger the switch to the secondary location,
 - activate the secondary location and restore data and software at the secondary location,
 - resume operation at the secondary location.



Manual failover

- Trigger
 - During failover, unmirrored systems are unavailable
 - Even mirrored systems may experience some disruptions. E.g. to unmirrored data
 - Consequently, a manual failover is triggered by a human
- Activating the secondary location should be scripted.
 - Fewer errors
 - Allows for one button failover
- Resuming operation is done by changing DNS setting.
 - Until this completes, user requests will be sent to your (now failed) primary data center.
 - If possible, a message “temporarily out of service” should be posted



Automatic failover

- Used for systems that have a very short RTO, where the time needed for human decision and action is too long.
- Assumes mirrored secondary location
- Requires monitor to trigger failover For these systems,
- Possible false positives
 - The data center did not fail, just a VM in the primary data center.
 - The VMs in the primary data center were slow to respond.
 - The network between the two data centers failed or was congested.
- A false positive may cause database inconsistency



Testing the failover process

- Concerns
 - You do not want to interrupt service to your clients.
 - The test fails and the production database becomes corrupted.
 - Your clients receive messages from the secondary location. While this is the desired behavior if there was a real disaster, these messages should not escape during a test.
- If possible, test during scheduled down time.
- Back up database prior to test



Discussion questions

- Suppose it is not possible to schedule down time for failover testing. How can you perform a test?
- What personnel should be available during a failover test?