# *Creating an effective EC2 Backup Strategy*

It's 11:30pm on a Friday night and you're ready to settle down, maybe watch some Leno. You check your email for the last time and realize to your horror that your EC2 cluster has just stopped responding. Your first thought, maybe my Internet connection is down, but that's not it. Maybe my Apache is down, no I can't ssh in. Maybe Amazon is down, no, some of my machines are responding. Lastly, you check the forums and there has been a hardware failure on several EC2 machines and all your data is gone.

If this sounds familiar, this discussion is for you. While I can't help you restore what is already lost, I'd like to give an overview of several different approaches to backing up and recovering your data on EC2 to ensure it never happens again.

The afore mentioned senario is exactly what happened to Enomaly earlier this year and lead us to create ElasticDrive, a continuous data protection server application for Amazon S3. You can learn more about ElasticDrive at http://www.elasticdrive.com. Please feel free try our Public AMI or VMware appliances.

## S3 Based Backups
Backing up your data to S3 is probably going to be the easiest and most cost effective solution for most EC2 users. S3 provides virtually limitless storage at a relatively low cost. The amount of data you want to backup is an important factor when planning the perfect strategy. Other variables such as the amount of lists/puts will also become a key factor. File systems typically work on a block level meaning that data is written and read on a frequent basis.

Creating a S3 based file system can be the simiplest solution but may also be the most costly if you're writing heavy amounts of data. Dumping files may be the cheapest but may require a lot of work upfront. Each have their own pros and cons. In building ElasticDrive we decided that a file system approach offered the least up-front configuration and setup for the typical consumer. Rather then having to re-design and develop an application to take advantage of the S3 API, all our customers need to do is install ElasticDrive and point their existing applications to either  the new file mount or configure a hard drive mirror using RAID. We also felt that giving our customers the choice of several remote storage solutions allowed for the most flexibility now and in the future - after all technologies change very quickly.

One of the typical starting points for using S3 is a periodic dump of data. For example, an SQL dump can happen on an hourly or daily basis. This is a simple yet effective way to backup key data and using various S3 applications such as ftp-based transfers. A popular solution has been jungledisk for FTP based S3 data transfer.

For more intensive or highly dynamic applications like databases or websites, a more broad data mirroring may be more efficient. This is a use case we at Enomaly have the most experience with using our ElasticDrive.

## Continuous Data Protection
We created ElasticDrive for the purposes of continuous data protection. Wikipedia defines continuous data protection (CDP) also called continuous backup, as backup of computer data by automatically saving a copy of every change made to that data, essentially capturing every version of the data that the user saves. It allows the user or administrator to restore data to any point in time. This becomes very important if you lose an instance in the middle of a transaction and you want to roll back to a point before the transaction had started, such an hour or two in the past. You could even roll back days, weeks or years if needed - any point in time before the corruption occurred.

ElasticDrive is a service that captures data changes to a separate storage location. There are multiple methods for capturing the continuous changes that serve different needs. CDP-based solutions can provide fine granularities of restorable objects ranging from crash-consistent images to logical objects such as files, mail boxes, messages, database files and logs.

For Enomaly the simplest solution is to run ElasticDrive in the form of a mirrored RAID drive (**Redundant Array of Independent Drives**) where all data is automatically written to a S3 backed virtual drive. In order to restore all our previous data we simply need start a new AMI with ElasticDrive automatically mounted at startup, which prepares us in case of such an emergency. If an instance is lost, a new AMI can be launched and downtime is kept to a minium. If you assume the AMI's are going to be lost, you could even create a monitor to automatically repair or rebuild lost instances with little to no human intervention. This is also very handy when using both local and remote server resources inside and outside of AWS.

### Virtual Tape Library & Off-site / Remote Backups
Another option is to create a virtual tape drive where data is backed up to look and act like a tape similar to a traditional enterprise but instead to a virtual storage environment.

There are literally dozens of commerical and open source applications dedicated to tape based backup. One such application is Zmanda which allows for an easy to manage environment with a web based graphical interface and various other features. Amanda provides the unique capability of writing backups to tape and disk simultaneously. The very same data can be available online at EC2 for quick restores from a local disk and off-site (low cost dedicated server host or S3) for disaster recovery and long-term retention. http://amanda.zmanda.com/

Using traditional dedicated hosting providers is also a popular trend among EC2 users and may be a solid approach for long term low cost data protection. A number of third party vendors offer API based S3 alternatives. For instance AOL xdrive offers 5GB for free.

### Network File Sharing & Virtual SAN
Creating a virtual storage area network is yet another option. This allows for either a master / slave type environment to be configured or even a hybrid where both a S3 backed drive and distributed storage engine can be used.

A recommendation for this type of solution is **dcache** is ideally suited for use with.The goal of this open source project is to provide a system for storing and retrieving huge amounts of data, distributed among a large number of heterogenous server nodes, under a single virtual filesystem tree with a variety of standard access methods. Depending on the Persistency Model, dCache provides methods for exchanging data with a variety of storage systems as well as space management, pool attraction, dataset replication, hot spot determination and recovery from disk or node failures. Connected to a distributed storage system, the cache simulates unlimited direct access storage space. Dataexchanges to and from the underlying HSM are performed automatically and invisibly to the user. Filesystem namespace operations can be performed through a standard nfs interface. http://www.dcache.org/

### Distributed / Replicated File System
Another option is to use a Google style replicated file system. One such system is GlusterFS. It enables a clustered file-system capable of scaling to several peta-bytes. It aggregates various storage bricks over Infiniband RDMA or TCP/IP interconnect into one large parallel network file system. GlusterFS effectively allows users to create a google style cluster with data distributed amount multiple EC2 nodes. This is great when working on large EC2 clusters as the chance of all cluster nodes failing at the same time is slim so additional S3 storage may not be needed. Although for most users I would still suggest some kind of secondary backup regardless of the amount of EC2 redundancy.

My concern with glusterFS is its lack of data security so be careful.

### MySQL Backup and Recovery
Most modern web applications depend heavily on a database. The most frequently used database for Enomaly is MYSQL. Having an effective database backup and recovery plan can save you time and money. There are a number of third party applications devoted to helping manage SQL backups ranging from ruby or java libraries to full turnkey solutions such as Zmanda's MYSQL backup.

The Zmanda offering is free and extremely user friendly. Features include; Schedule full and incremental backups of your MySQL database. Start immediate backup or postpone scheduled backups based on your needs. Choose to do more flexible logical or faster raw backups of your database.

Perform backup that is the best match for your storage engine and your MySQL configuration. Backup your remote MySQL database through a firewall. Configure on-the-fly compression and/or encryption of your MySQL backups to meet your storage and security needs.  Get e-mail notification about the status of your backups and receive MySQL backup reports via RSS feed.  Monitor and browse your backups. Define retention policies and delete backups that have expired. Recover a database easily to any point in time or to any particular transaction, e.g. just before a user made an error. Parse binary logs to search and filter MySQL logs for operational and security reasons http://www.zmanda.com/backup-mysql.html

In conclusion, there are various security measures you can take to make sure you don't lose important data. The best bet is to save your information in more then one place and plan for the worst.

You can learn more about ElasticDrive at http://www.elasticdrive.com

**About Reuven Cohen**
Reuven Cohen is the Founder and Chief Technologist at Enomaly Inc, a Toronto based open source technology firm. Reuven has extensive experience working with emerging enterprise technology and has developed in excess of 500 websites for companies including John Hancock, Intel, Alliance Atlantis, 20th Century Fox, Best Buy and Business Objects.

To Learn more about Enomaly, please visit http://www.enomaly.com or http://www.enomalylabs.com