# Large Lab Exercise Cloud Computing: Video Converter on AWS

13-11-2012

*Author:*
N. Singh
N.R.N. Timal
n.singh-2@student.tudelft.nl
n.r.n.timal@student.tudelft.nl

*Support cast:*
D.H.J. Epema
B.I. Ghit
A.Iosup
d.h.j.epema@tudelft.nl
b.i.ghit@tudelft.nl
a.iosup@tudelft.nl

**Abstract**

# 1 Introduction

This document is constructed in response to WantCloud BVs interest in moving a new application to the cloud. To this end the CTO of WantCloud wants us to create a (relatively) simple application to identify what types of trade-offs exists when deploying an application on the cloud.

We chose Amazon Web Services (AWS)[1] as our cloud provider because it is one of the major companies providing resources to external parties to date. If WantCloud BV will follow through with their wishes to deploy an application in the cloud most likely AWS will be one to seriously consider. As part of AWSs policy it offers new users of the cloud a free usage tier[2]. For our system we solely used the accommodations provided by the free usage tier, in particular the free use of Elastic Cloud Computing(EC2) micro instances and Elastic Load Balancer.

We created an application which allows users to convert a video file to a desired format using a client-server model[3]. Although this may seem to be a simple application its workflow can contain very demanding operations. That is client(/user) needs to upload a video file to the server, the server in turn processes the uploaded file and converts and returns the converted file to the client. The client then stores the converted file to its local file system. We therefore find this application to be fitting for assessing the tradeoffs inherent in cloud based applications.

Our cloud based system will implement the following. Automation, which means that a servlet will automatically be assigned to a EC2 micro instance. Auto-scaling, that is depending on the amount of client requests addition (when CPU usage of existing instances are high) or removal (when CPU usage of an instance is below a low threshold and can be compensated by an existing instance) EC2 micro instances. Load balancing, equally dividing the work between instances using Round Robin scheduling[4]. Reliability is achieved attaching by every micro instance to a persistent storage volume. And last but not least monitoring which will provide a comprehensive view of e.g. the amount of resources used by each live micro instance. The implementation of the aforementioned features will be accomplished by interfacing AWS with Pythons boto framework. The code of this framework can be found at this location (http://code.google.com/p/boto/). This interface extends among other components of AWS to EC2 and Simple Storage Service(S3) which will be used by our cloud based system.

The rest of the report is organized as follows. Section 2 gives a more in-depth view of the application that is implemented. Section 3 discusses how the features of the cloud based system are implemented (automation, auto-scaling, load balancing, reliability and monitoring). The experimental setup and the experiments conducted are discussed in Section 4. Any limitations of our cloud based system based on the experimental results from Section 4 appear in Section 5. Finally the paper concludes in Section 6.

---

[1]http://aws.amazon.com/
[2]http://aws.amazon.com/free/
[3]http://en.wikipedia.org/wiki/Clientserver_model
[4]http://en.wikipedia.org/wiki/Roundrobin_scheduling

# 2  Application

We have constructed a video converter tool in Java[5] which allows an user to specify what file needs to be converted and to what format it will be converted. The conversion from one format to another is done in the cloud. Using a client-server model, the client uploads a file and additionally provides the output format. Whereas the server, which is provided by the cloud, deals with converting the uploaded file. When the conversion is done the converted file is written to the local file system of the client. We have used the Jersey[6]Framework to create a servlet which is constructed using a REST[7] based approach. This servlet contains the algorithm for converting a video and is deployed on a Tomcat server in the cloud. The algorithm uses several methods from the Xuggle[8] framework. This framework is free and open-source and is used for audio/video manipulation. As for the client we used Jersey Client Framework which provides a comprehensive interface for communicating with a server.

## 2.1  Requirements

In this paragraph we will elaborate what requirements are essential for our application. The requirements can be subdivided into two categories: *functional* and *nonfunctional requirements*. A fitting *functional requirement* for our application is that it needs to convert to and from all commonly known codecs/video formats e.g. Windows Media Video(WMV)[9] and QuickTime MOV[10]. As a *nonfunctional requirement* the servlet needs to be deployed on the cloud, to achieve this goal we will be using the facilities provided by AWS. In particular Amazon Simple Storage (S3) for storing the WAR file created by the Jersey Servlet and Amazon Elastic Compute Cloud(EC2). An EC2 instance will be comprised of our servlets computational logic, that is once the servlet is launched on an EC2 instance, the instance will be able to execute a video conversion.

---

[5]http://www.java.com
[6]http://jersey.java.net/
[7]http://en.wikipedia.org/wiki/Representational_state_transfer
[8]http://www.xuggle.com/
[9]http://en.wikipedia.org/wiki/Windows_Media_Video
[10]http://en.wikipedia.org/wiki/QuickTime

# 3  System Design

The online video converter developed for WantCloudBV is deployed using Amazon Web Services (AWS). The video converter is build as a client-server model where the server part resides in the cloud and client is launched on the local machine. The application is developed using Java and we use AWS Java SDK to deploy our servers .war file to AWS Elastic Beanstalk Server (EBS) by uploading the application to Amazon S3. The Elastic Beanstalk environment is setup as Apache Tomcat Server 6 and deployed with a keypair to help in SSH-ing to the instances attached to EBS, the URL of this EBS environment is used by the applications client to send request to the server. The launch of the EBS is associated with the booting of a t1.small instance(http://aws.amazon.com/ec2/instance-types). Once the server is deployed and running (state of instance), video converter client can be launched to make requests to the server.

## 3.1  System Objectives

For the video converter application we designed an IaaS based system with the desired properties namely automation, elasticity, performance, reliability, and monitoring. We have used boto(a python interface to Amazon Web Services) to write a python script which implements all the above stated features.

### 3.1.1  Automation

This python script[11] implements all the features of the IaaS based system without any human intervention. When the server is ready we run this script and then proceed with video conversion requests. The script checks for all the running instances, starts monitoring on each of these instances, automates the scaling of instances based on CPU Utilisation metric, and balances the load on these instances. All the monitoring statistics are gathered in a log file.

### 3.1.2  Reliability and Durability

The application launches with minimum one instance, and initially attaches an ephemeral EBS storage volume to it, which means if the instance terminates the attached storage volume will be deleted as well. In our system(script) we first find out all the instances and attach each one of them with a unique EBS storage volume which is permanent and is not deleted with termination of an instance. This same volume is then attached to the newly launched instance, thus providing reliability between instance restarts and re-boots.

### 3.1.3  Monitoring

In the first prototype of our monitoring module, every instance was monitored by first SSH=ing into the instance and running the script which calculated the cpu usage and network I/O metric and stored the results in a file on the instance. This script was run as a daemon on the system. The log file was retrieved from every instance and stored as unique files corresponding to every instance. This approach was however rejected since the command used for calculating cpu usage did not give the correct and reliable results. Find the discussion here [12] [13] If you run a command like top from within an EC2 instance, you will not see an accurate CPU utilization metric for your instance. Your EC2 instance runs inside a virtual machine. The only way to capture an accurate metric for CPU utilization is to get that data from the virtual machine parent. Our CloudWatch metrics capture data directly from the parent of the virtual machine so you can get metrics for your instances. According to the suggestions, we were forced to use Amazons CloudWatch for monitoring our instances. The final design queries the CloudWatch for metrics and stores the

---

[11]https://github.com/nidhisingh88/CC_script
[12]https://forums.aws.amazon.com/thread.jspa?threadID=95288
[13]http://www.axibase.com/cloud/2010/07/22/ec2-monitoring-the-case-of-stolen-cpu/

results in a local file as log, the queries are launched at a gap of 60 seconds. The log file contains the metric with timestamp and instance id.

### 3.1.4   Load Balancing

The system script looks for the running EC2 instances and attaches them to a elastic load balancer which uses round robin technique to balance load between the instances. All newly launched instances are registered with the load balancer and all terminating instances are deregistered.

### 3.1.5   Security

Our system uses security groups to implement security of the instances. A new security group is created with traffic allowed only on ports 80(HTTP) and 22(SSH), all the running and newly spawned instances are associated with this security group.

### 3.1.6   AutoScale

The appication starts with one instance and our IaaS system is capable of autoscaling the number of instances. Every running instance is checked for the average CPU utilization in the past one minute, this value is checked against a minimum and a maximum threshold which are user configurable in the system script. If the CPU utilization exceeds the maximum configured threshold then a new instance is spawned which is a clone of the previous instance. The newly launched instance is configured for monitoring, a new permanent EBS volume is attached to it and it is also attached to the current load balancer. If the CPU usage of the instance is below the minimum configured then the instance is terminated and detached from the load balancer. Our system design also permits to specify the maximum number of instances that should be spawned and never terminates all the instances, there will always be one running instance. Autoscaling is automated except the configuration of initial threshold parameters. Once the system script is started it automatically checks for autoscaling implementation every 10 seconds.

# 4 Experimental Results

## 4.1 Experimental Setup

Our application is running on EC2 instances. For all our experiments we used t1.micro instances but found out that it cannot service more than 10 requests and cannot cope with bursts of traffic[14]. We were using t1.micro as it is part of AWS free tier group. So if the company anticipates say 100 users at a time, then they need 10 t1.micro instances and scale up or down accordingly. Hence, we use the next instance type m1.small for our application. The auto-scaling system implementation will only spawn m1.small instances as new instances. For testing each module of the system design, we developed simple python scripts which use threads. We use a 105KB file that will be converted to the desired format. The IaaS system also runs on a python script which uses boto to interface with the AWS.

## 4.2 Experiments

---

[14]http://docs.amazonwebservices.com/AWSEC2/latest/UserGuide/concepts_micro_instances.html

# 5   Discussion

# 6 Conclusion

| Metric | Hour |
|---|---|
| *total-time* | 133 |
| *think-time* | 8 |
| *dev-time* | 65 |
| *xp-time* | 32 |
| *analysis-time* | 12 |
| *write-time* | 12 |
| *wasted-time* | 4 |

Table 1: Time spend on project

| Metric | Hour |
|---|---|
| *total-time* 14 | |
| *dev-time* | 6 |
| *setup-time* 8 | |

Table 2: Time spend on experiment 1

| Metric | Hour |
|---|---|
| *total-time* | 7 |
| *dev-time* | 0.5 |
| *setup-time* | 2 |

Table 3: Time spend on experiment 2

| Metric | Hour |
|---|---|
| *total-time* | 11 |
| *dev-time* | 6 |
| *setup-time* | 5 |

Table 4: Time spend on experiment 3

# 7 Appendix A

In this appendix we evaluate how much time we have spent on realizing this project. For this we have defined 7 metrics. *Total-time* , determines the amount of time for completing this project. *Think-time*, reflects the time it took too figure out how to solve WantCloud BVs problem. *Dev-time*, rates the amount of time it took to code and setup the system. *Xp-time* is the time taken by conducting experiments. Evaluating the results gathered from the experiments is reflected with the metric *analysis-time*. The time taken by writing this report is represented by *write-time*. The time it took for every activity that cannot be accounted for in the aforementioned 6 metrics and which does facilitate progress of the project is reflected by the *wasted-time*.

Furthermore we define three metrics for each experiment. *Total-time* is the time it took to conduct the experiments. The amount of time to create the experiment is reflected by *dev-time*. *Setup-time* accounts for making the environment suitable to conduct the experiment.

Table 1 presents the time spend on the project. Whereas Table 2, Table 3 and Table 4 represent the time spend on the individual experiment, respectively Experiment 1, Experiment 2 and Experiment 3 conducted in Section 4. Taken into account that at the start of the project we lacked sufficient knowledge of Python the overall (*total-time*) of Table 1 is higher than anticipated. When conducting experiment 1 we ran in to the limititation that the micro instances have. We wrongfully assumed that there was something wrong with our cloud based system. Which in turn resulted in debugging the system and caused the excessive amount of time spend on Experiment 1.