

COMPX527 Assignment 1 Report

Glenn Cumming

Department of Computer Science
University of Waikato
Hamilton, New Zealand
glenn@hif.nz

Mitchell Grout

Department of Computer Science
University of Waikato
Hamilton, New Zealand
mjg44@students.waikato.ac.nz

Shufen Li

Department of Computer Science
University of Waikato
Hamilton, New Zealand
sl302@students.waikato.ac.nz

YingJun Huang

Department of Computer Science
University of Waikato
Hamilton, New Zealand
yh320@students.waikato.ac.nz

Abstract—Abstract Text

Index Terms—AWS, COMPX527

I. SOLUTION SUMMARY

Solution Summary

II. MOTIVATION

We decided on object detection in the cloud in order to learn and appreciate the challenges of providing a stateless, scalable service that could be used for a variety of other projects. Examples included

- Contextual threat analysis, such as humans detected in an area which should only include livestock.
- Numerical analysis, such as the number of a species of wildlife in an area over time.
- Absence detection, identifying objects that should be in an image but are not

III. PROPOSED SOLUTION

In order to

IV. SOLUTION ARCHITECTURE

Solution Architecture

A. The Image Processing Server

The project uses the Retinanet based Object Detection developed by Keras ¹. The web service was developed in Python using a Flask server ².

B. The Load Balancer

The object detection being inherently slow forces the use of load balancing and scaling techniques. Using the AWS Elastic Load Balancer allowed us to balance the HTTP requests over two or more EC2 servers running the Flask servers. We could create as many Flask servers as we wished, although of course this lent itself to wasted compute resources. Though we had great success using AWS AutoScaling Groups, which allows

a lower and upper limit of EC2 instances running our Flask server to be defined, this functionality was not available in the student accounts.

C. The Web User Interface

In order to have a better demo, it was decided to go ahead and produce a web site interface that would allow uploading

D. The Web API

The intended standard way to interact with the Flask Object Detection Cluster is to access it via HTTP calls using the url `http://<load_balancer_fqdn>/detect`. As the system is intended for the use of non-private images no SSL was implemented on the load balancer, though it is supported. The code submitted for this project includes *simple_upload.py* and *forked_upload.py*, which are Python 3 examples for using the HTTP API.

V. DEVELOPMENT

A. Technology

a) *Retinanet*: Retinanet is an open source object detection neural net. This technology was chosen for the following reasons.

- Open Source
- Python 3 for rapid development
- Specific instructions using it with the dataset available on the Registry of Open Data on AWS ³

b) *Terraform*: Terraform was given in as an example of the provisioning service to use in this assignment. AWS CloudFormation was another option suggested. The decision to go with Terraform rather than CloudFormation was based on the following

- CloudFormation is proprietary and is specific to AWS cloud offerings. Though this assignment is on a small and temporary scale, we still do not want to use a technology that would create vendor lock-in

¹<https://github.com/fizyr/keras-retinanet>

²<https://pypi.org/project/Flask/>

³<https://registry.opendata.aws/fast-ai-coco/>

- Terraform supports provisioning many different platforms, both open source and proprietary, such as Azure, Google Cloud, Kubernetes and OpenStack ⁴. Developing experience in deployments with Terraform therefore was deemed to be more useful.

c) *Ansible*: Though Terraform is capable of running commands post-install in order to install services and other software needed for our cluster, our cluster used Ansible playbooks instead. Though it meant learning another technology, we deemed a good use of our time since:

- Terraform can only run scripts, which would have to be created.
- Ansible's language is very flexible and is created specifically for the purpose of deployment.
- This is a recommended approach by HashiCorp, the developers of Terraform. ⁵

d) *Flask*: The decision to develop and deploy a service using the Keras-ResNet Python library immediately suggested to us we should again use Python to provide the web service. Our team had developers with experience in both Tornado and Flask; Flask was used as the web service development fell on the person with Flask experience. Both solutions would have been suitable.

B. Comparison with existing object detection solutions

Amazon Rekognition

VI. SECURITY ASSESSMENT

Security Assessment

A. Data Security

The Flask Object Detection Cluster was specifically created to be stateless and public. Interception of HTTP requests and results can be easily intercepted by third parties with access to the networks between the client and the service.

B. Access Security

Access Security

C. Network Security

Network Security

D. Vulnerability Assessment

E. Monitoring

AWS EC2 instance monitoring was enabled during the Terraform creation of the cluster. Future development of the cluster would likely use additional monitoring of performance, costs and uptimes using dedicated EC2 instances, and such existing monitoring systems as Icinga, Prometheus and Munin. Uptime monitoring would necessarily be run externally, such as on internal machines or another cloud providers offering.

VII. ACTUAL AWS EXPENDITURE

TODO

⁴<https://aws.amazon.com/cloudformation/>

⁵<https://www.hashicorp.com/resources/ansible-terraform-better-together>

VIII. FUTURE IMPROVEMENTS

a) *Packaging*: Deployment scripts were chosen for deployment and running the services. However, if the service was to be developed further we would be using `stdeb` ⁶ for packaging for the Ubuntu servers, and potentially set up a Personal Package Archive ⁷. This would ensure further development and upgrades by making it part of the standard apt package management system.

b) *Service Management*: Currently the starting and stopping of the service is done via control scripts. These would be integrated in to systemd service management ⁸.

c) *HTTP*: The service is stateless, public and insecure. If a layer of security is desired to protect the data in transit, then HTTPS can be added to the Elastic Load Balancer via the AWS Certificate Manager ⁹. The simple reason for not using this was cost: we would have to pay for extra services when our budget was tight.

IX. TEAM MEMBERS CONTRIBUTIONS

X. ASSIGNMENT REQUIREMENTS COMPLETION

XI. STANDARDS

This report was created in LaTeX using the IEEEtran document class provided by IEEE template available at ¹⁰, and generated into PDF by Gnome LaTeX ¹¹.

⁶<https://pypi.org/project/stdeb/>

⁷<https://help.ubuntu.com/community/PPA>

⁸<https://wiki.debian.org/systemd>

⁹<https://docs.aws.amazon.com/elasticloadbalancing/latest/classic/ssl-server-cert.html>

¹⁰<https://www.ieee.org/conferences/publishing/templates.html>

¹¹<https://wiki.gnome.org/Apps/GNOME-LaTeX>