# Consolidated Review of

# *Next Stop, the Cloud: Understanding Modern Web Service Deployment in EC2 and Azure*

## 1. Strengths:

This paper uses a diversity of data-sets including DNS subdomains, packet captures from a University campus network to characterize IaaS cloud usage patterns, and employs active measurements to evaluate the impact of Web service deployment strategies such as region and availability zones usage on the performance and fault tolerance of Web services. A first study looking at how IaaS clouds are used by web services.

Overall I liked this paper. It seems a nice contribution to the space. While we have seen all sorts of cloud papers I am not aware of others that tackle the problems this paper does. I.e., tries to determine the footprint of these clouds within popular web services. The paper is sound and quite methodical in execution. Papers like this may not exactly point out problems or things we should try to fix. However, they are valuable in that they inform our mental models of how the modern network is put together. So, I found this to be interesting and informative.

## 2. Weaknesses

The data set from traffic capture is limited. The traffic patterns of cloud-using services presented such as protocols and traffic volumes are skewed, since the data only includes network traffic between the clouds and a single campus network.

Some questions regarding how generalizable some of the results are. The methodology is relatively straightforward and the descriptions do not state how the limitations of the data impact the conclusions drawn (more details are discussed below).

I found the bits about resilience to be somewhat tenuous.

- ❖ There is some suggestion of the brittleness of concentrating functionality in one place. But, it's just a suggestion.
- ❖ We have no notion of how often this happens.
- ❖ There is no notion considered that the clouds may have schemes to fail over to a different zone / region if a large outage did happen.
- ❖ There is no notion of how likely it is that these places actually will lose all power.
- ❖ And, there is no notion of the cost of spreading a web service to more than one or two regions / zones. This all may just be simple money vs. likelihood of problems issue. But, since the downsides are not considered it looks like more of a no-brainer than it really might be.
- ❖ In general I think the numbers may make the situation appear worse than it may be.

## 3. Comments

The paper presents two datasets to analyze the usage of Amazon EC2 and Azure. One dataset is a list of cloud-using subdomains of Alexa's top 1 million list, which is generated by DNS queries. This list is combined with university packet trace. From this data source the paper examines a magnitude of aspects of web services deployed in EC2 and Azure, e.g. deployment patterns, protocol distribution and traffic patterns. In addition the paper studies the distribution of the deployed services over different areas and

resulting implications. Overall this paper is quite nice and gives a good summary of the state of the art. However, there are also some limitations that you may want to make rather explicit.

This paper uses one-week packet capture from one campus network to characterize traffic patterns of cloud-using services. This paper needs to emphasize the limitation of such a small set of data sets in characterizing traffic patterns of cloud tenants. Thousands of tech-savvy users in a campus network do not have "typical" usage patterns of Web services including those deployed in IaaS clouds.

Re the packet captures: Any notion of how much measurement-based loss there is in these? The traces are essentially uncalibrated. And, given the full payloads it seems that one could readily get some idea about the measurement loss using the techniques from Nechaev's 2009 IMC paper on calibration.

How representative is your university setup as compared with the overall population? One way to study this is to compare the popularity of the local services with the Alexa list...

This paper could continue to correlate traffic capture data with usage patterns of current Web service deployment by cloud tenants. For example, how much traffic was sent to these subdomains with only a single region or a single availability zone?

How are the services going to evolve and how will that impact your conclusions?

This paper uncovers an interesting and surprising fact that a large fraction of subdomains use only a single availability zone or region. However, this paper needs to realize that not all subdomains have the equal values to a Web site. For example, some subdomains might attract a large number of eyeballs and bring higher revenues, while some subdomains are quiet. It would be interesting if this paper could provide more details on the popularity of these subdomains using data-set from traffic capture.

This paper should separate the actual tables, e.g., Figure 2, from figures. In page 12, the final throughput could be expressed as "download_file_size/the_download_time" or other ways, instead of "download file size/the download time".

Table 2 is really hard to understand. It took me far too long to grok it. I.e., where things are supposed to sum to 100 and where they are not. Table 3 is even worse. This is space savings run amok. This should be three tables or something I think and yet (I'm guessing) to save real estate they have all been rammed together. And, it makes for a table that is very confusing. And, its also wrong---the EC2 total and the Azure total for #Domains does not add up to the "Total" in either raw count or percentage. What a mess ...

Clarifications are needed to justify how Figure (Table) 3 is able to definitively tell whether domains use EC2 only or Azure only. This seems to assume that all the subdomains can be surveyed, which are only possible if zone transfer requests were successful (which is clearly not the case as discussed in 2.1). The other

approach of identifying sub-domains in a brute-force manner is unlikely to exhaustively identify all possible subdomains.

The measurement and data analysis methodology rely on several assumptions, e.g., the IP address range advertised by EC2 and Azure are relatively complete. Similar such assumptions need to be clearly stated. Another related issue with any conclusions drawn from this measurement study is that the data used (e.g., packet traces, DNS data) have limitations (just like any other data). The paper doesn't discuss such limitations and the implications of the limitations on the analysis and conclusions drawn from the data. Some stats on the number of distinct IPs in the packet capture traces would be helpful to give some perspective of the representativeness of the data and conclusiveness of the data presented such as Figure/Table 5.

In Section 3.3, where traffic patterns are described, if the content is using HTTPS, it would not be possible to know the content type. Such limitation needs to be explicitly stated. The claim that very few zones are used needs to be better substantiated. It could be affected by the locations of Planetlab nodes used for probing that are not diverse enough in their geographic coverage.

In 4.1 I **THINK** this is talking about DNS probing. But, I am not entirely sure this data is not coming from the packet traces. The methodology here is under-specified. Please better explain what data you are looking at. A bunch of the CDFs are weird because the lines suggest continuous data, but the data is actually discrete. E.g., figure 8. There should be steps on the plot and not lines between points that suggest interpolations of (say) 1.5 virtual machines. It is also weird that tables are labeled as "figures". I guess that is just style, but it's a weird style if ya ask me! :) Nice to see that some of the data will be made available.

## 4. Summary from PC Discussion

**Strengths:**

❖ Interesting and the first study on how web services are deployed using EC2 and Azure
❖ Insight on how tenants should place their services for improved performance and fault tolerance

**Weaknesses:**

❖ The data set is limited (from a single campus network)
❖ Limited analysis to correlate usage patterns and service deployment within cloud, should weigh based on the popularity of subdomains.
❖ The claim on failure resilience (due to use of a single availability zone/region) is based on limited evidence, e.g., without considering other issues such as built-in fault

tolerance support of the cloud services, without considering the cost issues.

Overall, we all agree that this paper is interesting in terms of topic selection and has some nice measurement observations.

## 5. Authors' Response

As noted by reviews, the completeness of our Alexa subdomains dataset, and correspondingly our identification of cloud-using (sub)domains and characterization of tenant's deployment posture, may be limited by several factors: gaps in the list of IP address ranges published by EC2 and Azure, the use of a brute-force method to identify subdomains, and a lack of sufficient geo-diversity in performing DNS queries. In Section 2.1, we added text to acknowledge the former two limits, and we added a figure depicting the locations of the PlanetLab nodes used for performing DNS queries to address the later concern.

We also acknowledge the limits of relying on a packet capture from a single vantage point. We added text to Section 2.1 to make this explicit and provide a few more details about our capture (e.g., loss rate). Moreover, we use the packet capture only to extract information not attainable via DNS probing or active measurements---namely, protocol usage (Section 3.1), popularity estimates based on traffic volume and flow counts (Section 3.2), and traffic patterns (Section 3.3)---and we compare against findings from our Alexa subdomains dataset and other studies (e.g., DeepField's study) whenever possible. In the future, we plan to capture packet traces from additional vantage points.

To address confusion regarding the tables depicting protocol usage and (sub)domains' usage of EC2, Azure, and other infrastructure, we have split the former table into two, reorganized rows in the later table, and added additional text in Sections 3.1 and 3.2. In particular, we noted that the "EC2 total" and "Azure total" domain percentages in Table 2 sum to more than 100% because a small fraction of domains (0.7%) use both EC2 and Azure.

Finally, we extended our region usage analysis in Section 4.2 in two ways. First, we added a new "analysis of subdomain deployment vs. customer location" to answer the question of whether subdomains are deployed near their customers, revealing that a considerable fraction of web services are not deployed near their customers. Second, we noted several non-performance-related factors that may influence tenant's choice of number of regions.