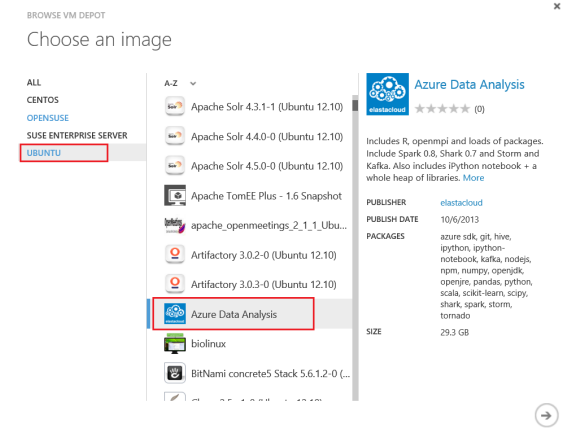
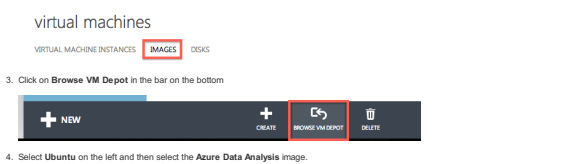


Processing GPS data with Storm and Kafka on Windows Azure Data Science Core

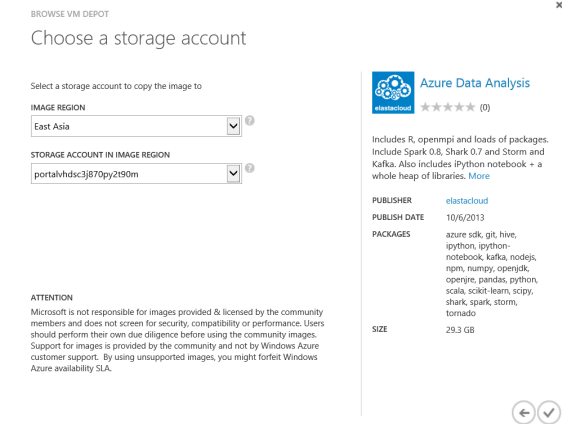
In this example, we'll show you how to deploy a Storm topology in Windows Azure that reads its data from the Kafka messaging system. We'll use a Kafka client application written in Java to send GPS coordinates from anywhere to the Kafka cluster. Our Storm topology will translate those coordinates into JSON objects, use GeoJSON to identify the country those coordinates belong to, and then keep a running count of how many times that a coordinate lands in a country. For persistence, the running count is stored in a Windows Azure Table Storage service, and the topology periodically dumps a compressed block of coordinates to a Windows Azure Blob Storage service. The topology also writes data to Redis for use by other services, such as the web application we use to display the data in real time. The web app is written in Node.js and uses Socket.IO and the express web application framework to read the data from Redis and display it via D3.js.

Use Azure Management Portal to Create a Windows Azure Data Analysis VM

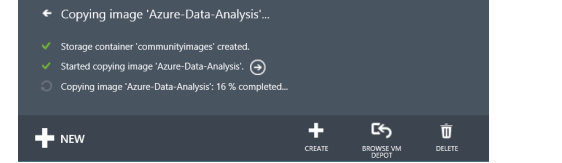
- 1. Log in to the Windows Azure Management Portal.
- 2. Click on the **Virtual Machines** tab and click on **Images** near the top of the screen.



- 5. Choose the **Image Region** that your storage account is in (i.e. the region you created your affinity group in) from the drop down box, then select your storage account from the drop down box.



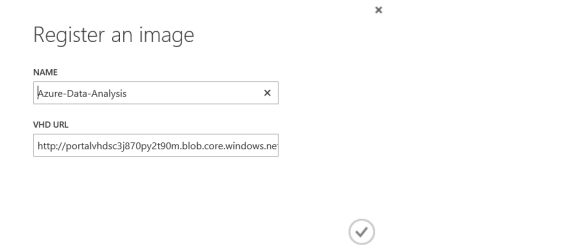
- 6. Click the check mark button to continue and wait for the disk image to be downloaded from the VM Depot to your storage account. You can click on **Details** in the status bar to see the transfer progress.



- 7. Once the image has copied you'll need to register it. Select the image and click **Register** in the bar on the bottom.



- 8. Enter a name for the image, click the checkmark button, and wait for registration to complete.



- 9. Stay on the Virtual Machines tab and click on **Virtual Machine Instances** near the top of the screen.

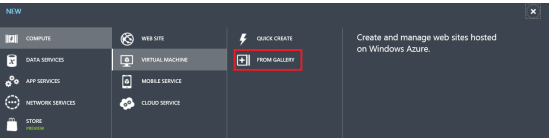
virtual machines

VIRTUAL MACHINE INSTANCES

IMAGES

DISKS

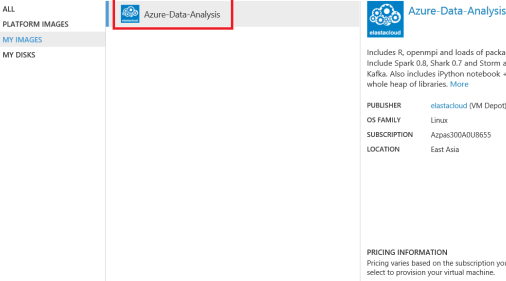
10. Click on **New** in the bottom bar and select **From Gallery**.



11. Select **My Images** on the left and then select the Azure-Data-Analysis you just registered. Go to the next page.

CREATE A VIRTUAL MACHINE

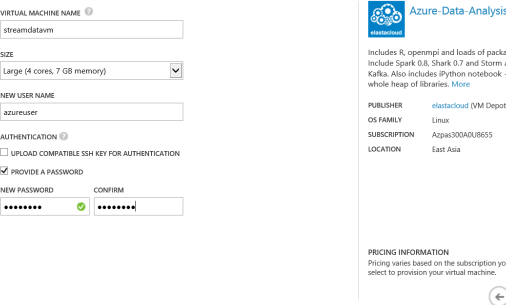
Virtual machine image selection



12. Enter the virtual machine name, select the **Large** machine size from the drop down list, enter a new user name, check the **Provide a Password** box and enter the new user password. Go to the next page.

CREATE A VIRTUAL MACHINE

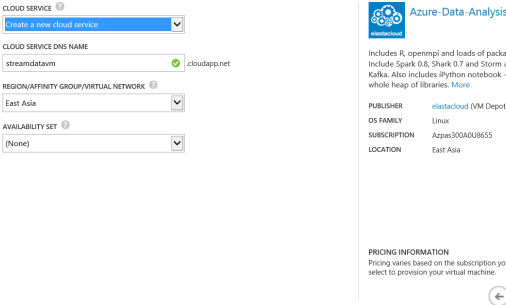
Virtual machine configuration



13. Enter a name for the new cloud service configuration and select your affinity group from the drop down box. Go to the next page.

CREATE A VIRTUAL MACHINE

Virtual machine configuration

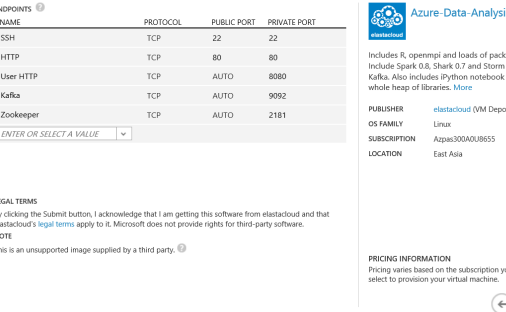


14. We'll need to add several endpoints for the VM. To add an endpoint, enter a name for the endpoint like "HTTP" or "Kafka" in the field under **Name** in the Endpoints table. Add the following TCP endpoints:

1. HTTP: Port 80
2. User HTTP: Port 8080
3. Kafka: Port 9092
4. Zookeeper: 2181 Your endpoints should look like this:

CREATE A VIRTUAL MACHINE

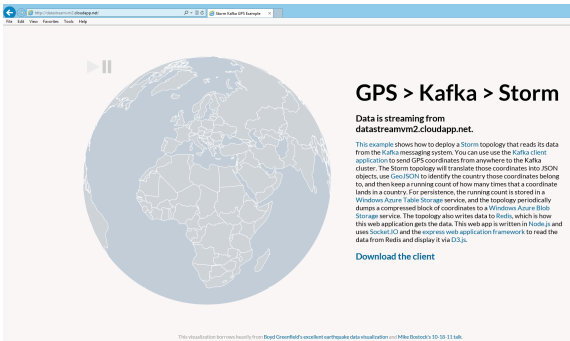
Virtual machine configuration



You should see output like this:

```
info - socket.io started
Listening on port 88
```

2. Open your web browser and go to the cloud service DNS name you specified when you created the VM. You should see the web app:



Nothing interesting is happening because we have not yet started the data stream. Leave this browser window open so you can see the effects of the following commands.

1. Open a second SSH connection to the VM (open PuTTY on Windows or a new terminal on Linux or OS X) and launch the Storm topology:

```
cd $HOME/gpskafkadeemo
java -cp $(lein classpath) storm.example.KafkaGpsTopology
```

You should see a lot of output as the topology starts. Once it's up and running the output will look like:

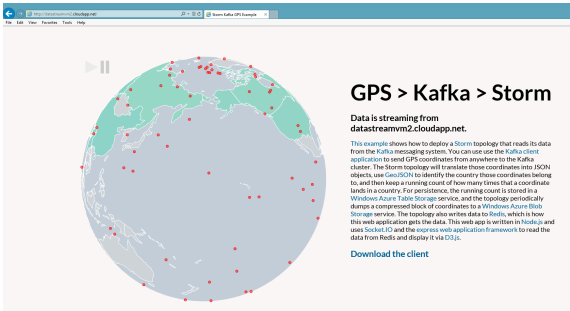
```
4105 [Thread-25] INFO storm.kafka.PartitionManager - Starting Kafka 127.0.0.1:0 from offset 21853734106 [Thread-25] INFO
backtype.storm.daemon.executor - Opened spout(spout(6) 4105 [Thread-25] INFO backtype.storm.daemon.executor - Activating spout
spout(6) 4224 [Thread-25] INFO storm.kafka.PartitionManager - Committing offset for 127.0.0.1:9092:0 4224 [Thread-25] INFO
storm.kafka.PartitionManager - Committed offset for 127.0.0.1:9092:0 6241 [Thread-25] INFO storm.kafka.PartitionManager - Committing offset
for 127.0.0.1:9092:0 6242 [Thread-25] INFO storm.kafka.PartitionManager - Committed offset for 127.0.0.1:9092:0
```

You won't see any change your browser because although the topology has been started there is no data being sent to Kafka for the Storm topology to process.

1. Open a third SSH connection to the VM and start the Kafka client:

```
cd $HOME/gpskafkadeemo/kafka-gps-client
java -cp $(lein classpath) kafka.example.KafkaGpsDataProducer localhost
```

The client generates random GPS coordinates and sends them to Kafka. Go back to your web browser and you'll see GPS coordinates being plotted on the globe. Countries will change color according to the frequency of "hits".



The client publishes data on the "gps" topic. **localhost** on the command line means we are connecting to Zookeeper on localhost to get connected to the Kafka server. You can specify the connection string as any of:

- zookeeper_host
- zookeeper_host:port
- brokerid:kafkaHost:kafka_port

- You can run this Kafka client from any machine with Java. To run from your local workstation, download the client JAR file from the web application page and run it as follows:

```
java -cp kafka-gps-client-0.0.1-SNAPSHOT-standalone.jar kafka.example.KafkaGpsDataProducer 0:ds-11nford.cloudapp.net:9092
```

In this case we've bypassed Zookeeper and connected directly to the Kafka server on port 9092. The 0: at the start of the connection address indicates that we wish to connect to the Kafka broker with ID 0.

Copyright 2013 Microsoft Corporation. All rights reserved. Except where otherwise noted, these materials are licensed under the terms of the Apache License, Version 2.0. You may use it according to the license as is most appropriate for your project on a case-by-case basis. The terms of this license can be found in <http://www.apache.org/licenses/LICENSE-2.0>.