

13) I started the load test for the web server with 100 users and 5 users/second, there was no failed requests and the scale-out alarm was triggered so the number of instances was scaled to 3. Then I changed to 100 users and 5 users/second. There appeared some failures and the failure rate was around 20% ~ 30%, as shown in figure1. The scale-out alarm is still shown in the in-alarm state and number of instances grew slowly as the alarm was triggered and it hasn't reached the maximum number of instances yet. Since I set the "And then wait" parameter to 300 seconds, the number of instances grew slowly at around 5 minutes per new instance, as shown in figure 2.

After it reached the maximum number of instances 10, I stopped the locust request. I expected the scale-in alarm would be triggered soon, but the two web server alarms were both showing "insufficient data" state for several minutes, after which I realized I still needed to make requests. So I set the locust load test with 1 user and 1 user/second and the scale-in alarm was triggered. The number of instances was decreasing at a fast rate as I set the "And then wait" parameter to 0, it only took the time to terminate an instance. When the number of instances was reduced to 5, it stopped decreasing as the target response time was ≥ 0.01 and scale-in alarm was in "OK" state. After around half an hour as the target response time was below 0.01 and the scale-in alarm was triggered and the number of instances decreased to 2 in the end.

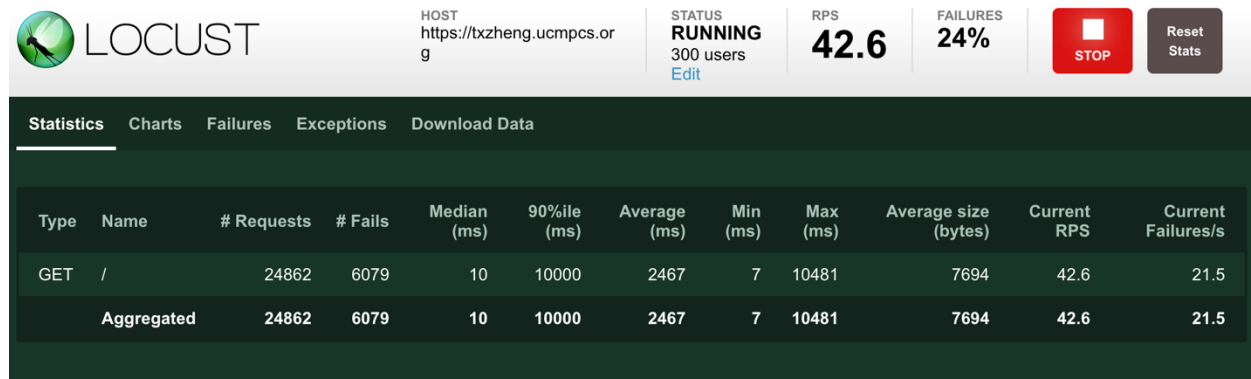
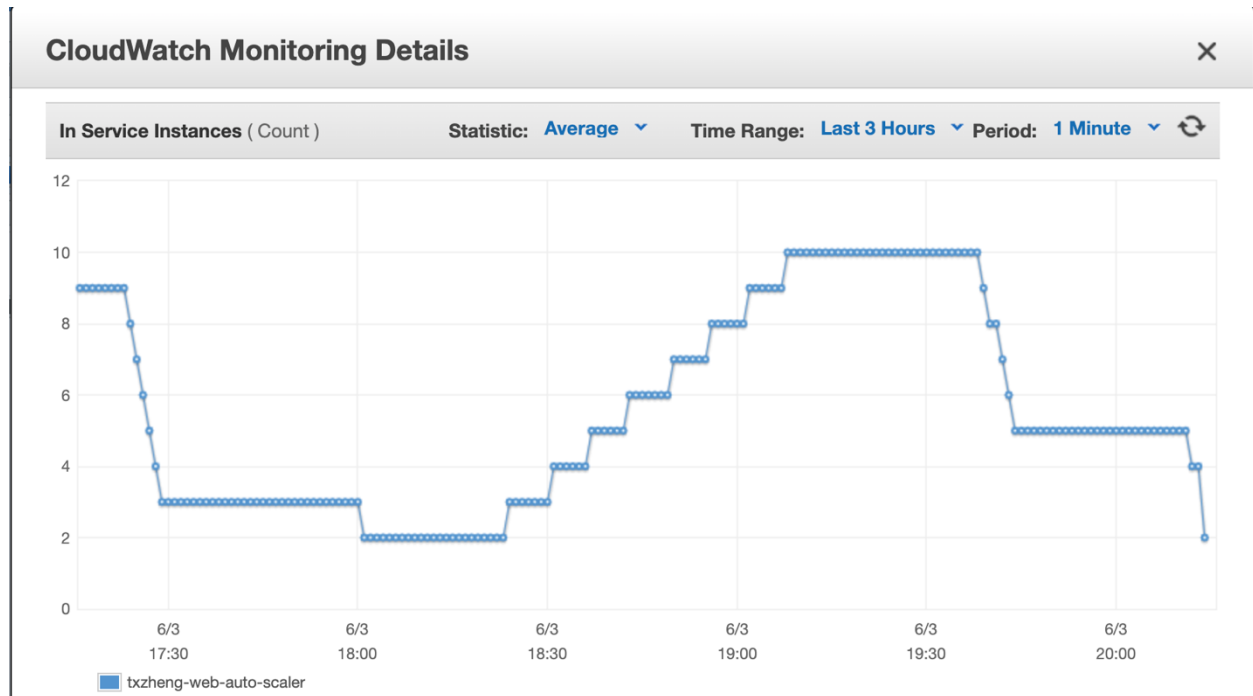


Figure 1



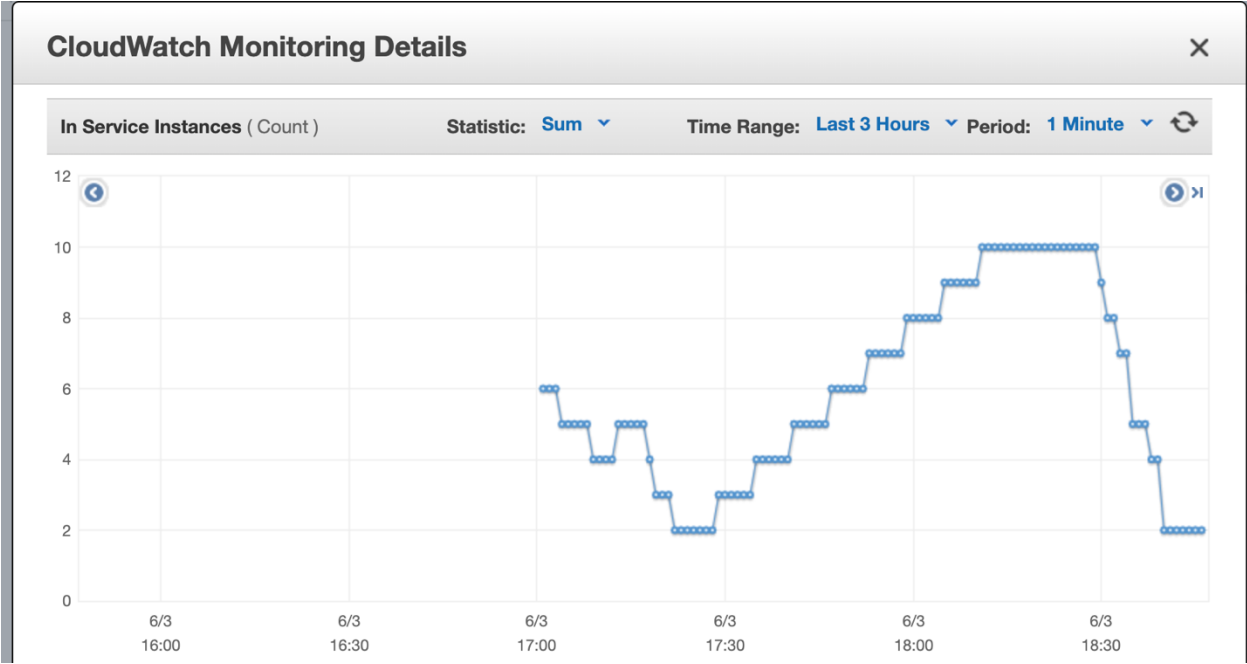


Figure 3