# Azure status history

Product:          Region:          Date:
[All]             [All]            [Most recent]

## May 2019

### 5/22   RCA - Service Management Operations - West Europe

**Summary of impact:** Between 15:10 and 21:00 UTC on 22 May 2019, a subset of customers in West Europe experienced intermittent service management delays or failures for resources hosted in this region. Impacted services included Azure Databricks, Azure Backup, Cloud Shell, HDInsight, and Virtual Machines.

Between 20:20 and 21:50 UTC on 22 May 2019, during the deployment of the permanent fix, a small subset of customers using Virtual Machines and Azure Databricks experienced increased latency or timeout failures when attempting service management operations in West Europe.

**Root Cause:** The issue was attributed to performance degradation in the Regional Network Manager (RNM) component of Azure software stack. The RNM component, called the partition manager, is a stateful service and has multiple replicas. This component saw an increase in latency due to a build-up of replicas being created for the services. Prolonged operational delays triggered a RNM bug which caused the primary replica to re-build on two occasions. This caused service management operations to fail while one of the other replicas was taking on the primary role.

**Mitigation:** Engineers identified the RNM bug and applied a hotfix to the region which helped resolve network operation failures. During the outage, an auto-validation process on RNM nodes was performing scans which slowed the replica buildout for impacted nodes. Engineers terminated the scans to improve performance. The network operation job queues began to drain and latency returned to normal.

We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Develop and roll out dedicated partition for large customers  [in progress]
- Implement automatic throttling to balance load  [in progress]
- Root cause and fix the cause for high commit latency [in progress]

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/3KB-5FZ

### 5/13   RCA - Network Connectivity - Increased Latency

**Summary of impact:** Between 05:00 and 15:33 UTC on 13 May 2019, a subset of customers in North America and Europe may have experienced intermittent connectivity issues when accessing some Azure services.

**Root cause and mitigation:** The impact was the result of inconsistent data replication in a networking infrastructure service. This resulted in unexpected throttling of network traffic to our name resolution servers. Once the issue was detected, engineers mitigated it by updating the configuration of the affected network infrastructure service to override the effect of this data inconsistency. Simultaneously, engineers performed operations to repair the data inconsistency.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Improve our monitoring to detect data inconsistencies similar to the one that caused this issue.
- Improvements in the system to help ensure such inconsistencies do not occur in the future.

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/7CC-H9G

### 5/7   SQL Services - West Europe

**Summary of impact:** Between 10:53 and 12:48 UTC on 07 May 2019, a subset of customers using SQL Database, SQL Data Warehouse, Azure Database for PostgreSQL, Azure Database for MySQL, Azure Database for MariaDB, in West Europe may have experienced issues performing service management operations – such as create, update, rename and delete- for resources hosted in this region.
In addition, customers may have been unable to see their list of databases using SSMS. However as this was a Service Management issue, these databases would not have been impacted (despite not being visible through SSMS).

**Preliminary root cause:** Engineers identified a back-end database service responsible for processing service management requests in the region became unhealthy preventing the requests from completing.

**Mitigation:** Engineers performed a manual restart of the impacting back-end service, which restored its capacity to process requests, mitigating the issue.

**Next steps:** Engineers will continue to investigate to establish the full root cause and prevent future occurrences. Stay informed about Azure service issues by creating custom service health alerts: https://aka.ms/ash-videos for video tutorials and https://aka.ms/ash-alerts for how-to documentation

### 5/2   RCA - Network Connectivity - DNS Resolution

**Summary of impact:** Between 19:29 and 22:35 UTC on 02 May 2019, customers may have experienced connectivity issues with Microsoft cloud services including Azure, Microsoft 365, Dynamics 365 and Azure DevOps. Most services were recovered by 21:40 UTC with the remaining recovered by 22:35 UTC.

**Root cause:** As part of planned maintenance activity, Microsoft engineers executed a configuration change to update one of the name servers for DNS zones used to reach several Microsoft services, including Azure Storage and Azure SQL Database. A failure in the change process resulted in one of the four name servers' records for these zones to point to a DNS server having blank zone data and returning negative responses. The result was that approximately 25% of the queries for domains used by these services (such as database.windows.net) produced incorrect results, and reachability to these services was degraded. Consequently, multiple other Azure and Microsoft services that depend upon these core services were also impacted to varying degrees.

**More details:** This incident resulted from the coincidence of two separate errors. Either error by itself would not have been non-impacting:

1) Microsoft engineers executed a name server delegation change to update one name server for several Microsoft zones including Azure Storage and Azure SQL Database. Each of these zones has four name servers for redundancy, and the update was made to only one name server during this maintenance. A misconfiguration in the operance of the automation being used to make the change resulted in an incorrect delegation for the name server under maintenance.
2) As an artifact of automation from prior maintenance, empty zone files existed on servers that were not the intended target of the assigned delegation. This by itself was not a problem as these name servers were not serving the zones in question.

Due to the configuration error in the change automation in this instance, the name server delegation made during the maintenance targeted a name server that had an empty copy of the zones. As a result, this name server replied with negative (nxdomain) answers to all queries in the zones. Since only one out of the four name server's records for the zones was incorrect, approximately one in four queries for the impacted zones would have received an incorrect negative response.

DNS resolvers may cache negative responses for some period of time (negative caching), so even though anomalous configuration was promptly fixed, customers continued to be impacted by the change for varying lengths of time.

**Mitigation:** To mitigate the issue, Microsoft engineers corrected the delegation issue by reverting the name server value to the previous settings. Engineers verified that all responses were then correct, and the DNS resolvers began returning correct results within 5 minutes. Some applications and services that accessed the incorrect values and cached the results may have experienced longer restoration times until the expiration of the incorrect cached information.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

1) Additional checks in the code that performs nameserver update in unrelated changes  [in progress].
2) Pre-execution modeling to accurately predict the outcome of the change and detect potential problems before execution  [in progress].
3) Improve per-zone, per-nameserver monitors to immediately detect changes that cause one nameserver's drift from the others  [in progress].
4) Improve DNS namespace design to better allow staged rollouts of changes with lower incremental impact  [in progress].

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/R5SSC-5RZ

### 5/2   RCA - Azure Map

**Summary of impact:** Between 04:25 and 11:00 UTC on 02 May 2019, a subset of customers using Azure Maps may have experienced 500 errors when attempting to make calls to Azure Maps Rest APIs.

**Root cause and mitigation:** Engineers were notified by internal monitoring that connectivity between internal components/services required by Azure Maps were disrupted. This lead to the failure to fulfill incoming customer requests. Upon investigation, engineers found that the authentication application was inadvertently removed in error during regular maintenance operations. Restoration of authentication application was performed manually, which led to the mitigation of this incident.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Moving the authentication applications to a more protected repository, and setting up fail safes to prevent future outages.
- Augmenting monitoring for authentication issues and improving troubleshooting guides for faster response.
- Longer term, engineers are planning to move to a different authentication platform which will provide better insights earlier.

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/DYD-J8Q

### 5/1   Issue signing in to https://shell.azure.com

**Summary of impact:** Between 18:00 UTC on 30 Apr 2019 and 17:20 UTC on 01 May 2019, customers may have experienced issues signing in to https://shell.azure.com

During this time, customers were able to access Cloud Shell through the Azure portal at https://portal.azure.com

**Preliminary root cause:** Engineers identified a mis-match between a configuration file which had been recently updated and its corresponding code in shell.azure.com

**Mitigation:** The Cloud Shell team developed, tested, and rolled out a new hotfix which addressed and corrected the issue.

**Next steps:** Engineers will continue to investigate to establish the full root cause and prevent future occurrences.

## April 2019

### 4/19   RCA - Availability degradation for Azure DevOps

**Summary of impact:** Between 17:00 and 23:20 UTC on 19 Apr 2019, a subset of customers experienced issues connecting to Azure DevOps. These issues primarily affected customers physically located on the East Coast and those whose organizations are located on the East Coast.

**Root cause:** During a planned maintenance event for Azure Front Door (AFD), a configuration change caused network traffic to be incorrectly advertised. The AFD ring impacted by this maintenance hosted Azure DevOps and other Microsoft internal tenants. This may have resulted in timeouts and 500 errors for customers of Azure DevOps.
The maintenance event started at 3:30 UTC, which started dropping around 5-10% of requests. When the environment severely degraded at 14:44 UTC, engineering observed the major impact start. The maintenance event was on a ToR (Top of Rack) switch. The standard operating procedure is to take the environment offline by removing edge machines. By design, the MUX stopped advertising BGP (Border Gateway Protocol) routes and traffic is routed through these MUX. Within this environment one of the MUX Load Balancers was in an unhealthy state but the BGP session between the load balancer and the TOR was still active. Consequently, the MUX was still active in the environment and the TOR was advertising traffic incorrectly.

**Mitigation:** The first impact window was mitigated by withdrawing the invalid route so that traffic would be routed correctly. The recurrence was caused by the maintenance process resetting the configuration back to the previous state, publishing an invalid route. The 2nd mitigation was re-applying the change again.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- reviewing and implementing more stringent measures for when we take environments offline for maintenance events.

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/WCMY-3QQ

### 4/16   RCA - Networking Degradation - Australia Southeast / Australia East

**Summary of impact:** Between 07:12 and 08:02 UTC on 16 Apr 2019, a subset of customers with resources in Australia Southeast / Australia East may have experienced difficulties connecting to Azure endpoints, which in-turn may have caused errors when accessing Microsoft Services in the impacted regions.

**Root cause:** Microsoft received automated notification alerts that the Australia East and Australia Southeast regions were experiencing degraded network availability from a select number of Internet Service Providers (ISPs). During this time, a subset of network paths used for the select number of ISPs, this manifested in traffic not reaching the destinations within the Australia East and Australia Southeast regions. The issue stemmed from a routing anomaly due to an erroneous advertisement of prefixes received via an ExpressRoute circuit to an Internet Exchange (IX).

**Mitigation:** Microsoft disabled the incorrect ExpressRoute peering. The IX also identified a high amount of traffic and automatically mitigated by bringing down the peering with the IX. Once the peering were brought down by Microsoft and the IX, availability was restored to Australia East and Australia Southeast regions.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Engage Internet Service Providers to add additional policies/protections to Internet facing routing infrastructure to block future routing anomalies [Complete]
- Add additional automated route mitigation steps within the Azure platform to reduce mitigation time [In Progress]
- Investigate further route optimizations in the Azure/Microsoft ecosystem to inherently block future routing anomalies [In Progress]

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/DPVY-1FG

### 4/12   RCA - Cognitive Services

**Summary of impact:** Between 02:50 and 11:30 UTC on 12 Apr 2019 a subset of customers using Cognitive Services including Computer Vision, Face and Text Analytics in West Europe and/or West Central US may have experienced 500-level response codes, high latency and/or timeouts when accessing to resources hosted in this region.

**Root cause:** Engineers determined a recent deployment introduced a software regression, manifesting in increased latency across two regions.

**Mitigation:** The issue was not detected in pre-deployment testing, however, once manually detected, engineers proceeded to roll-back the recent deployment task to mitigate the issue.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Improve pre-deployment tests to catch this kind of issue in future [In Progress]
- Improve monitoring to more closely represent production traffic patterns [In Progress]

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/QKRI-7FG

### 4/9   Virtual Machines - North Central US

**Summary of impact:** Between 21:34 on 9 Apr 2019 and 01:20 UTC on 10 Apr 2019, a subset of customers using Virtual Machines in North Central US may have experienced connection failures when trying to access some Virtual Machines hosted in the region. These Virtual Machines may have also restarted unexpectedly. Some residual impact was detected, impacting a small subset of recovered Virtual Machine connectivity into the underlying disk storage.

**Root cause:** Azure Storage team made a configuration change on 9 April 2019 at 21:30 UTC to our back-end infrastructure in North Central US to improve performance and latency consistency for Azure Disks running inside Azure Virtual Machines. This change was designed to be transparent to customers. It was enabled following our normal deployment process, first to our test environment, and lower impact scale units before being rolled out to the North Central US region. However, this region hit bugs which impacted customer VM availability. Due to a bug, VM hosts were able to establish session with the storage back end, but hit issues when trying to receive/send data from/to storage scale unit. This situation was designed to be handled with fallback to our existing data path, but an additional bug led to failure in the fallback path and resulted in in VM reboots.

**Mitigation:** The system automatically recovered. Some of the customer VMs which didn't auto recover, needed an additional recovery step.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- We have paused further deployment of this configuration change until the underlying bugs are fixed [complete].
- Fix bugs that caused the background operation to have customer-facing impact [in progress].
- Additional validation prior to cause for the scenario that caused the bugs to be missed in test environment [in progress].

## March 2019

### 3/29   RCA - SQL Database

**Summary of impact:** Between 11:45 and 22:05 UTC on 29 Mar 2019, a subset of customers may have experienced the following:

- Difficulties connecting to SQL Database resources in the East US, UK South, and West US 2 regions
- Difficulties connecting to Service Bus and Event Hubs resources in the East US and UK South regions
- Failures when attempting service management operations for App Service resources in the UK South and East US regions
- Failures when attempting service management operations for Azure Search resources in the East US region

**Root cause:** Azure SQL DB supports VNET service endpoints for connectivity specific databases to specific VNETs. A component used in this functionality, called the virtual network plugin, runs on each VM used by Azure SQL DB, and is involved at VM restart or restarts. A deployment of the virtual network plugin was rolling out worldwide. Deployments in Azure follow the Safe Deployment Practice (SDP), which aims to ensure deployment related incidents do not occur in many regions at the same time. SDP achieves this in part by limiting the rate of deployment for any one change. Thus the role of the incident this particular deployment had already successfully occurred across multiple regions and for multiple days such that the deployment had progressed to simultaneously roll out the change to several regions at once. This deployment was using a VM instant capability, which occurs without bringing or running workloads on those VMs.

On 5 capacity units across 3 regions, an error in the plugin lead process caused the VM to fail to restart. The virtual network plugin is configured as 'required to start', so absence of it prevents key VNET service endpoint functionality from being used on that VM. The error led to repeated restart attempts causing the VMs to continuously cycle. This occurred on enough VMs across those 5 capacity units that there were not enough resources available to provide placement for all databases in those units causing those database became unavailable. The plugin error was specific to the hardware types and configurations on the impacted capacity units.

The 5 capacity units affected included some of the databases used by Service Bus, Event Hub and App Services in those regions which led in-turn to the impact to those services. An impacted database if Azure SQL has the global service management state for Azure IoT Hub, hence the broad impact to that service.

**Mitigation:** Impacted database using the Azure SQL DB AutoDR capability were failed over to resources in other regions. Some impacted databases were moved to healthy capacity within the region. Full recovery occurred when sufficient affected VMs were manually rebooted on the impacted capacity units. This brought enough healthy capacity online for all databases to become available.

**Next steps:**

We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Fix the error in deployment, which led to continuous recycling on the specific hardware types and configurations [in progress]
- Repair deployment block system - if stopped of deployment in each capacity unit before widely, but not soon enough [in progress]
- Investigate why faster deployment block automation - it detected correlated impact at region level, but would have identified each impacted capacity unit separately [in progress]
- Improve resiliency for IoT Hub [in progress]

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/FHDR-7VZ

### 3/28   RCA - Data Lake Storage / Data Lake Analytics

**Summary of impact:** Between 22:10 on 28 Mar 2019 and 03:23 UTC on 29 Mar 2019, a subset of customers using Data Lake Storage and/or Data Lake Analytics may have experienced impact in these regions:

- East US 2 experienced impact from 23:40 UTC on 28 Mar to 03:23 UTC on 29 Mar 2019.
- West Europe and Japan East experienced impact from 22:10 to 23:50 UTC on 28 Mar 2019.

Impact symptoms would have been the same for all regions:

- Customers using Azure Data Lake Storage may have experienced difficulties accessing Data Lake Storage accounts hosted in the region. In addition, data ingress or egress operations may have timed out or failed.
- Customers using Azure Data Lake Analytics may have seen U-SQL job failures.

**Root cause:**

**Background:** ADLS (Gen1 uses a microservice to manage the metadata related to placement of data. This is a partitioned microservice where each partition serves a subset of the metadata. Each partition is served by a fault tolerant group of servers. Load across various partitions is managed by an entity called the partition.config - this a master file who's held information about an instance of the microservice; a per region file is generated by a tool. (This tool is applied to all config files, not just partition.config.) Load is balancing across multiple servers in the overall load across servers within a fault tolerant group. A large cause of metadata is served by the region. Currently, these load balancing actions on multiple partition.config files which must without impact to running workloads on those VMs.

All loads and config's are hosted across microservice-related deployments are staged and controlled such that deployment goes to a few machines in a region before moving to the next region. A software component called recharger for hosting the service microds and issuing, and resizing errors, which will stop a deployment when the first node is not on, and revert those that went. When the deployment is completed, this recharger wiped on those metadata state pages. Moving to next region requires success deployment in the current region AND approval of the engineer.

**What happened:** Some of the microservice instances across different regions needed balancing of load to continue to provide best experience and availability. An engineer made changes to partition.config files for the identified regions and triggered deployment using the process described above. After observing success in a canary region, the engineer approved deployment in additional regions. After deployment completed successfully, the engineer received alerts in two regions - East and West Europe.

Investigation revealed a syntax error in the partition.config. The tool which generates this file caused a syntax error, which generated the specific partition.config file and failed to generate a valid config for the specific partition.config file. This had the caused the error problem for the metadata which the partition.config was applied. This config file change was rolled out across region, and the missing partition.config would cause file to crash. The deployment in the canary region and other regions succeeded because they were different partition.config file.

**Mitigation:** The engineer reverted the bad syntax error in the partition.config file. The new version of partition.config fixed the error, mitigating those four regions as HS stopped crashing. BUT as the result of ADLS crashing which would not have detected faster if each capacity unit was not restoring the service availability.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Mandatory test run automatically at submit time, that sanity-checks partition.config. This test would catch both the syntax error and the logic error.
- Hardening the config deployment tooling, so that it has built-in delays between regions instead of manual approvals.
- Enhance the watchdogs so that they catch more and cause deployments to fail to automatically and revert.
- Enhance microservice logic to deal more gracefully with errors in partition.config.
- Fix the tool that generates per region config file for the issue that caused it to delete the region's partition.config file; instead have it raise an error to fail the deployment.
- Move partition.config to a data folder with separate file for each region, so that an error in one region doesn't affect other regions.

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/WFT3-3MS

### 3/27   RCA - Service Management Failures - West Europe

**Summary of impact:** Between approximately 15:20 UTC on 27 Mar 2019 and 17:30 UTC on 28 Mar 2019, a subset of customers have received failure notifications when performing service management operations such as create, update, delete, scale and restart for resources hosted in the West Europe region.

**Root cause and mitigation:**

**Root Cause:** Regional Network Manager (RNM) is a core component of the network control plane in Azure. RNM is an infrastructure service that works with another component called the Network Service provider (NSP) to orchestrate the network control plane and drive the networking goal states on host machines. Days leading up to the incident, peak load in RNM's partition manager sub-component had been increasing steadily due to organic growth and load spikes. In anticipation of this, the engineering team had prepared a code improvement to the lock acquisition logic to enhance the efficiency of queue draining operations in the RNM. On the day of the incident, before this change could be deployed, the load increased sharply, concentrating on a few subscriptions. This pushed RNM to a tipping point. The load caused queued times to rise, resulting in failures. As the load increased concurrently in a few subscriptions, leading to lock contentions where the thread waiting on the other, causing a slow-down of work items in queues resulting in failures in service management operations. The gateway component in RNM started to aggressively self-healing and the failure back in to the queue as retries, leading to a compounding effect. Highler in the stack such as NRP and Compute Resource Provider (CRP) further aggravated with its retries.

**Mitigation:** To mitigate the situation and restore RNM to its standard operating levels, the retries had to be stopped. A hotfix to stop the gateway component in RNM from adding retry jobs to the queue was successfully applied. In addition, the less-reproducible deployments were improving plans for the partition manager sub-component were suspended. The engineer took actions to reduce peak load by putting temporary throttling restrictions on a few subscriptions that were generating peak load were identified from analysis of the RNM, and throttling applied. As load decreased steadily, and the load returned to standard operating levels. Finally the originally planned partition manager optimization code change was rolled out to adjust the handling of the RNM, bringing RNM back to standard operating levels and providing the ability to take higher loads and improving its performance.

**Next steps:** We sincerely apologize for the impact to affected customers. We are continuously taking steps to improve the Microsoft Azure Platform and our processes to help ensure such incidents do not occur in the future. In this case, this includes (but is not limited to):

- Improve lock contention in RNM
- Improve RNM performance and scale capacity with enough headroom
- Mechanisms to throttle individual before RNM service hits tipping point

**Provide feedback:** Please help us improve the Azure customer communications experience by taking our survey: https://aka.ms/QK0N-3BG