# AIOps for a Cloud Object Storage Service

Anna Levin, Shelly Garion, Elliot K. Kolodner, Dean H. Lorenz, Katherine Barabash IBM Research  Haifa
lanna,shelly,kolodner,dean,kathy@il.ibm.com

Mike Kugler, Niall McShane IBM Cloud and Cognitive Software Mike.Kugler@ibm.com, nmcshane@us.ibm.com

platforms are intrinsically integrated into the IT system they help to operate. Other solutions, mostly vendor products and services, are more generic. These often cover only a part of the AIOps stack, such as data collection, or specific analytical methods and algorithms, or data integration in a data lake, etc. Such point solutions are hard to integrate into existing operational production systems in a meaningful way. We have addressed this gap by creating a useful AIOps system around an existing cloud-scale service, namely an object storage service. Cloud-scale object storage may contain trillions of objects, serve massive amounts of simultaneous access requests, and generate tons of operational data. Our goal is extracting useful business insights from the operational data collected dynamically from a production service. However, the sheer magnitude of data, as well as variability of formats and data sources, and speed of data generation, make it difficult to turn abstract theories into practice.

In this paper, we share our experience and lessons learned creating a data driven AIOps platform for a production cloud object storage service at IBM. Our contributions are: (1) describing service operations in production and the available operational data, (2) presenting the AIOps solution we have built for gaining operational insights, and (3) showing the operational pain points our solution helps to resolve.

## II. Problem Definition

We are developing AIOps capabilities for IBM Cloud Object Storage (COS) so that its operation will be data driven and automated. We collect several types of operational data from object storage, and do several different analyses on it. IBM COS encrypts and disperses the objects stored in it using erasure coding across multiple geographic locations [9]. Access to objects is over HTTP using a REST API. IBM COS has a two tier architecture: (1) front-end servers, called *Accesser*® nodes, that receive the REST commands and then orchestrate their execution across (2) storage-rich back-end servers, called *Slicestor*® nodes, that store the data.

*1) Operational data:* There are many types of IBM COS operational data, e.g., logs and metrics. We describe two kinds that we use in our analyses.

The first are access logs [10]. These are in JSON format. These logs contain an entry for each operation invoked on the object storage. The entry contains a wealth of information regarding the operation; this includes items such as the operation type, the bucket and object names on which the operation works, the HTTP return code, the start and end times of the

operation, and various latency statistics (total latency, time spent waiting for the client, time spent waiting for the storage back end, time for authentication and authorization, etc.). The access logs are generated on the Accesser nodes.

The second are connectivity logs. These are structured JSON logs produced once a minute on each IBM COS server. They provide information about the connectivity of the servers across all the offerings of IBM COS.

*2) Goals and challenges:* Our aim is to detect, predict and prevent failures and performance slowdowns that could impact users. Cloud object storage is engineered with redundancy, so that it continues working despite failures of individual components. Thus, when a failure occurs it is masked by the redundancy, which makes it difficult for the operators to discover. Our goal is to discover these failures and pinpoint their cause through analysis of the data.

In turn, the analysis of the operations data also poses several challenges. The first is that the schema of the log data is dynamic. Each log record is JSON, so its fields are labeled, but each individual record might have a different subset of the fields, e.g., depending on the operation type, the authentication mechanism, or whether the object data is encrypted. The second is the scale of the data to be processed. As described in Section III-B, we converted the JSON data to Apache Parquet, partitioned it and added derived fields in order to save space and speed subsequent processing. Poor choices, e.g., for partitioning, or producing Parquet objects that are too small, could lead to the costly need to reconvert and reprocess the data. Also, programs tested on small samples of data might fail after long running times when processing very large amounts of data, e.g., when encountering a new log record instance.

*3) Our Approach:* We apply the following stages:

**Ingestion (Sec. III-A):** Ingestion of historical data is from batch files, e.g., Elasticsearch dump files, as well as of streaming data, e.g., from Apache Kafka.

**Curation (Sec. III-B):** Data curation, cleaning and preparation for analytics is crucial, and requires a huge investment. In data warehousing this stage is called Extract, Load and Transfer (ELT).

**Features (Sec. III-C):** Features need to be generated in order to use statistical and machine learning algorithms. We present a "smart groupBy" method to generate features efficiently in parallel at scale using Apache Spark.

**Model (Sec. IV-A):** We use both statistics and machine learning to create our models and analyze incoming data.

**Causality (Sec. IV-B):** We do feature isolation to detect the root cause of a problem, such as a failing component.

**Visuals (Sec. IV-C):** Finally we present the root cause of the problem to the operator through reports, dashboards (e.g., Grafana), or by notifications (e.g., Slack).

## III. DATA PROCESSING

Our data processing flow was implemented in two ways. Initially the data samples were copied to a local storage system for easier exploration, see Fig. 1. However, moving the data far from its collection became unfeasible for the real

logs sizes we faced. Thus, we now do the data processing flow in the cloud, as shown in Fig. 2.

Fig. 1 shows the initial data processing flow. The IBM COS access logs in JSON format are collected in Elasticsearch [11] (ES) and historical dumps of ES are stored in an IBM COS bucket. In order to build a non-intrusive exploration pipeline, we copied the ES dumps to a local storage system: splitting files into chunks to optimize parallel copying and verifying correctness with MD5. The copied chunks were stored in HDFS [12] and converted to Parquet [13] to ease the analytics. For analytics we used Spark together with Zeppelin notebooks [14], a common platform for interactive data analytics.

Despite the advantages of having data close to the data scientist for efficient development and easy data manipulation, the initial solution was not scalable enough for our needs. As a result, we now do the processing in the cloud, as depicted in Fig. 2. In this second pipeline, the data processing is integrated with the collection stage and the logs are stored in IBM COS in Parquet format. The data is consumed for analytics by Apache Spark [15] applications running on the IBM Analytics Engine [16]. For interactive development, we use Jupyter notebooks [17] provided by IBM Watson Studio [18]. In the remainder of this section we describe our experience with the processing flows presented above.

### A. Ingestion

One of the common ways to ingest log files into a data lake is a Logstash processing pipeline [19] that ingests data from multiple sources simultaneously, transforms it, and then sends it to Elasticsearch - an open-source, RESTful, distributed search engine. We started with this ingestion approach in Fig 1. IBM COS operational logs were collected from the Accesser nodes, sent over Apache Kafka and stored in ES through Logstash. The ability to test development directly against ES before working with the full data volume provided fast initial time to value because of the easy integration with Spark and with visualization and reporting tools, e.g. Kibana.

We started the development of the pipeline by ingesting the access logs through ES to a local HDFS system, as depicted in Fig. 1. The log data was moved to HDFS, close to the processing, and after cleaning and format adjustment, the data was consumed for analytics. This method of ingestion provided an isolated environment for exploration independent from the production pipeline. This approach worked well for exploring relatively small samples; however, as the volume of the log data grew, it became infeasible to move the logs to HDFS.

Therefore, we decided to switch to a pipeline based on cloud services shown in Fig. 2. In this approach we collect and stream the logs directly to IBM COS, storing them close to the collection points and organizing them in Parquet, and also move the processing close to the data creation. Note that while this cloud-based approach is closer to a full data-lake pipeline, it is still a development environment. This approach is useful during development as it provides a sandbox with a

Fig. 1. Pipeline for processing operational logs



Fig. 2. Cloud-based pipeline

"cached" snapshot of the data. As demonstrated in the next section, this approach allows exploration of analytic methods in a repeatable manner without having to re-run the entire data pipeline.

*B. Curation*

It is not enough to collect the data to make it useful and meaningful. The ingested data must be cleaned, formatted and organized in order to make it useful for data scientists engaged in data discovery and analysis. The process of managing data, archiving and representing is called data curation. Data curators collect data from diverse sources and integrate it into repositories that are many times more valuable than the independent parts. Curation also ensures data quality and makes machine learning (ML) more effective.

It is hard to overestimate the importance of data cleaning and sanity checking. Our sanity checks include row counts, timestamps consistency checks and statistical data validation. In addition to counts, it is important to validate that statistical characteristics are preserved, e.g., the mean and median values of the important numerical fields. Further data cleaning and standardization require a unifying format for numerical data to allow numerical operations, and filling in missing data with appropriate values, e.g., null values of the appropriate type. Another important aspect of data cleaning is validating timestamps consistency and correctness, needed to ensure that no data is missing, ignored or duplicated. Validating timestamp consistency is a complex task in a heterogeneous system. Our data set is collected from different systems using multiple formats across many timezones. In order to validate timestamp consistency and enable data analysis, we unified the format and converted all timestamps to UTC.

In the data lake pipeline, curation prepares the data for consumption by analytic tools. Curation transforms the raw data into a structured view, annotates it using metadata, combines current inputs with historical data, and integrates previously aggregated data. During development, we chose to use a staging transformation approach. Staging the data in Parquet format greatly speeds up the debugging process while developing new analytic methods; it reduces the time wasted waiting for computations to complete. Moreover, staging simplifies repeatable experiments. We experimented with several choices before deciding on the right staging transformations and format to achieve the best performance.

Apache Parquet is a columnar data format, which supports compression and encoding schemes. The schema is embedded in the data itself, so it is a self-describing format. These features make it efficient to store data in Parquet as opposed to row-oriented schemes like CSV and TSV. Fig. 3 shows the

results of our experiments on the effectiveness of storing our log samples in Parquet. The figure shows for our case, when no explicit options are set and Spark uses default snappy compression, storing logs in Parquet rather than leaving them in JSON format in a form of compressed ES indexes, reduced the size by more than 10 times.

In addition, the columnar nature of Parquet enables efficient querying since most queries involve few columns and can be enhanced further by partitioning based on the values of one or more columns. It is crucial to know the common queries in order to benefit from by data partitioning. For example, among other fields our schema involves dates, locations and customer information. Since our typical analysis was done on a daily basis, we decided to partition data by date and location, storing it in an order dictated by dates and location columns. When the query is done for specific day and location only small sub-set of data is accessed. This order of partitioning served us well for daily performance logs analysis, but proved to be inefficient when per customer analysis was requested. In order to access specific customer data, we needed to access whole data set regardless of its date and location, losing all the advantages of partitioning.

Partitioning not only has benefits for performance, it also helps managing the data. Spark supports programmatic partitioning and can discover partitions automatically if the stored data is already partitioned. To store partitioned data in the most efficient way we performed experiments with Spark to determine the best way to write the data. The following write options were evaluated: (1) Each day-location gets its own dataset. (2) Hive-style partitioning using partition_by.write in *append* mode. (3) Hive-style partitioning explicitly writing each virtual directory in *overwrite* mode. Our experiments showed that option 1 is the slowest one. Options 2 and 3 are similar in terms of performance and twice as fast as option 1. The latter two options store the data and access it in a similar way; however, there is a difference between options 2 and 3 when a write job fails. If Spark job fails in option 2 in the middle of data writing, it cannot be easily reverted as the written data becomes part of a big data set with a single flag for successful writes. As a result retrying the write job might create duplicate data chunks. On the other hand, with option 3, if the write job fails writing a virtual directory, the partial results can be easily cleaned by rerunning the job, overwriting the specific directory. Therefore, we chose hive-style partitioning, explicitly writing a virtual directories in *overwrite* mode, as the fastest and safest option.

## C. Feature extraction

In large scale data sets, such as operational log data (e.g. [20]), a big challenge is preparing the raw data for use with statistical and ML methods. In the language of ML the prepared data is called "features". Moreover, for many statistical and ML methods, such as anomaly detection, choosing and generating the right features is much more important than the actual algorithm. For example, with the proper features basic outlier detection methods can find anomalies (see the discussion in [3]).

We use Spark [15], [21], which supports the Map/Reduce programming model, to generate the features efficiently in one parallel pass over the huge amount of log data. For the development of algorithms and code we use Jupyter notebooks, and for running the code in batch in production we use "spark-submit" to the IBM Analytics Engine [16]. We describe our "smart groupBy" method to generate the features efficiently, assuming that the input data is in a table format (e.g., Spark DataFrame) and that the processing is done using an SQL API (e.g., Spark-SQL).

**Map Step** We enrich the input data with additional columns:

*1) Computing ranges and buckets*: by ranges of timestamps (e.g., day/hour/minute), sizes, and other numerical components.

*2) Reducing a large set of values to a smaller subset (often called data cleaning)*: by parsing the original values and taking only dominant substrings, or taking only the most popular subset of values, etc.

*3) Computing a function over several columns*: by choosing a derived value based on values appearing in the input columns.

*4) Doing derivations and differences*: by computing a derived column $C$, adding a new column $C_{shift}$, which is a copy of $C$ shifted in one cell, and then computing the derived column $C_{derived} = C - C_{shift}$.

*5) Computing statistical functions*: counts, mean, standard deviation, ranks, percentiles, etc.

The one pass pipeline combines these transformations into a single map step.

**Reduce Step** We perform one SQL "group by" operation on many columns at once, including the columns generated in the map step, to produce all the features in parallel.

In order to choose the specific ranges, buckets, and subsets in steps 1 and 2, we use prior domain knowledge or learn from a small sample of the large data set. The statistical functions in step 5 may also be used to check the reliability of the generated features, for example, checking whether the number of items per feature is statistically meaningful. Since Spark is lazy [15], the entire calculation is done on-demand in the "Reduce step" to obtain the features. We perform the ML algorithm or statistical analysis on these features.

## IV. ANALYTICS AND INSIGHTS

In order to gain insight we perform statistical analysis and create ML models using the extracted features. Our analysis is done using Python packages, such as pandas [22] for general statistics, and scikitlearn [23] for ML.

## A. Analytical model: statistical or machine learning (ML)

In our use-cases we are able to provide an efficient solution using basic statistical and ML methods, and did not need deep-learning methods, which are less efficient and require more resources, especially when implemented at large scale.

*Statistical approach.* Computing statistical invariants of the large-scale dataset, such as counts, mean, standard deviation, histograms, median and other percentiles, is extremely useful to provide an overview of the operations data and perform basic reliability and sanity checks. This calculation can also be done efficiently at large scale using the "smart groupBy" method described above. We perform hierarchical aggregation of various metrics in order to detect and focus on a failed component, as can be seen in Fig. 6, to find the highest possible level of aggregation/hierarchy, representing the most significant problem currently occurring in the system. However, in cases where there are too many metrics and signals, it is not possible to manually handle too many graphs and alerts, so it is better to use automated ML tools.

*Machine learning approach.* There are various *anomaly detection* methods, for an extensive survey see [24]. The basic approach is based on calculating the *z-score* of a single metric, namely, the signed fractional number of standard deviations by which the value $x$ of an observation differs from the metric's mean value. In addition to a univariate approach, there exist multivariate anomaly detection methods, e.g., [25], which reduce false alerts. We use a multivariate anomaly detection algorithm on the features generate using the "smart groupBy" method described in Sec. III-C. These features include aggregations of various latency metrics of several components in the system, and our algorithm is in the spirit of Ng's algorithm in [26]. Fig. 4 shows a graph of our aggregated anomaly score and computed threshold ; Fig. 5 focuses on the anomaly inside the rectangle.

## B. Root cause and problem isolation

After completing the statistical and ML analysis, we perform feature isolation in order to find the root cause and detect the failed component. Figs. 5 and 6 demonstrate two equivalent points of view of the same failure that occurred in our system at the same time (from time $T_{start}$ to $T_{end}$).

Fig. 6 shows the *connectivity* point of view. We observe that certain application components were disconnected from the other components from time $T_{start}$ to $T_{end}$. In order to detect the most relevant problematic component for an alert, we reduced the problem of determining the failing components to a hierarchical flow problem. This approach allows us to pinpoint the problem at the most relevant level of aggregation/hierarchy, and notify the operations team of the most significant problem currently occurring in the system.

Fig. 5 shows the *performance* view. Our multivariate anomaly detection tool on the latency metrics shows a high peak at time $T_{start}$ that stays above the anomaly threshold for the same period from $T_{start}$ to $T_{end}$. In this case, we isolate the problem and focus on the misbehaving component by indicating the top features with the highest z-scores during
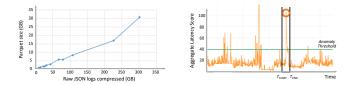
Fig. 3.   Size ratio Parquet vs. Raw JSON Logs Index Compressed



Fig. 4.   Anomaly score of the aggregated latency



Fig. 5.   Anomaly score of the aggregated latency – zoom in



Fig. 6.   Connectivity point of view: disconnected component

this anomalous period, and checking how many of them come from the same component.

### C. Visualizations and dashboards

In order to provide alerts and reports to the operation team, one can use static reporting tools that mainly present periodical graphs and textual reports. Another option is presenting semi-static tools, which allow interactive exploration of the pre-calculated results and support drill down and zoom-in at problematic points, such as Grafana [27]. Moreover, such semi-static tools allow to incorporate all the produced graphs into a few dashboards. However, such tools do not cover all types of required insights and graphs. For example, in our use-cases a connectivity matrix heatmap turned-out to be the most useful tool for presenting failed component in the context of overall system.

In addition to the graphical tools and dashboards, when an action of an operator is needed, it is necessary to provide a direct alert or real-time notification via tools like slack [28]. For example, our automated tool provides a slack notification to the operator, containing the identity of the specific application or network component experiencing issue. In addition to identifying the specific troubled component, the context of the issue is provided including exact time of the event, its geographical location, an estimation of the problem severity and the list of the additional components affected by this failure. It is important to provide the context of the event in order to assist the operator to discover the root cause of the problem and act quickly to restore normal system behavior.

## V. CONCLUSIONS

We have presented an AIOps solution that provides insights useful for the operation of the IBM COS service. The solution is now in the process of being globally deployed across multiple service offerings in IBM Cloud and will be further refined, optimized, and extended, e.g., to work with more cloud services and for cross-service operations. To conclude the paper we share a concise summary of the most important lessons and best practices that can illuminate the path to success for others who pursue similar goals.

*Know your tools.* There are multiple applicable tools and methods and it is very important to have good understanding of the suitability, efficiency, and compatibility of these methods in the context of a particular operational challenge.

*Keep it simple.* This timeless wisdom is your best friend when developing AIOps. The appeal of advanced analytics, including machine and deep learning, is so great that many fall in the trap of needless over complication.

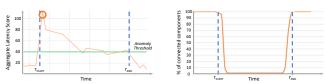*Iterate and refine as you go.* In an ideal world, one assumes sufficient and reliable data sources and specific questions to guide AIOps exploration. In reality, however, the data sources often are not intended for analytical purposes and asking the right questions is a significant challenge. Use an iterative approach whereby insights are gained incrementally, using snapshots, historical data, and calibration.

*Feed and sustain the engagement.* Last but not least, do not underestimate the importance of mutual trust between IT system owners and the AIOps team. System owners are often overwhelmed by their day-to-day load and if not seeing immediately useful results, may disengage and lose interest. Great cross-team collaboration is key to putting in place an AIOps system alongside the production service in a reasonable amount of time.

## REFERENCES

[1] A. Lerner, "AIOps Platforms," https://blogs.gartner.com/andrew-lerner/2017/08/09/aiops-platforms/.
[2] B. H. Sigelman and et al., "Dapper, a large-scale distributed systems tracing infrastructure," Google, Inc., Tech. Rep., 2010.
[3] D. Goldberg and Y. Shan, "The importance of features for statistical anomaly detection," in *USENIX Workshop, HotCloud'15*.
[4] M. Chow and et al., "The mystery machine: End-to-end performance analysis of large-scale internet services," in *OSDI'14*.
[5] "Introducing atlas: Netflix's primary telemetry platform," https://medium.com/netflix-techblog/introducing-atlas-netflixs-primary-telemetry-platform-bd31f4d8ed9a.
[6] "Moogsoft," https://www.moogsoft.com/.
[7] "OpsRamp," https://www.opsramp.com/.
[8] "Zenoss," https://www.zenoss.com.
[9] J. K. Resch *et al.*, "Aont-rs: Blending security and performance in dispersed storage systems," in *FAST'11*.
[10] "IBM Cloud Object Storage System Logs," www.ibm.com/support/knowledgecenter/STXNRM_3.14.1/coss.doc.
[11] "Elasticsearch," https://www.elastic.co/products.
[12] "Apache Hadoop," https://hadoop.apache.org/.
[13] "Apache Parquet," https://parquet.apache.org/.
[14] "Apache Zeppelin," https://zeppelin.apache.org.
[15] M. Zaharia and et al., "Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing," in *NSDI'12*.
[16] "IBM Analytics Engine," https://console.bluemix.net/catalog/services/analytics-engine.
[17] "Project Jupyter," https://jupyter.org/.
[18] "IBM Watson Studio," https://console.bluemix.net/catalog/services/watson-studio.
[19] "Logstash," http://www.elastic.co/products/logstash.
[20] J. Barr, "Amazon s3 – two trillion objects, 1.1 million requests per second," AWS News Blog 2013.
[21] "Apache Spark," https://spark.apache.org/.
[22] "Pandas: Python Data Analysis Library," https://pandas.pydata.org/.
[23] "scikit-learn: Machine Learning in Python," https://scikit-learn.org/.
[24] V. Chandola and et al., "Anomaly detection: A survey," *ACM Computing Surveys*, Jul. 2009.
[25] H. Cheng and et al., "Detection and characterization of anomalies in multivariate time series," *SIAM Int. Conf. on Data Mining'09*.
[26] A. Ng, "Anomaly detection using the multivariate gaussian distribution," in *Machine Learning Yearning*, 2018.
[27] "Grafana," https://grafana.com/.
[28] "Slack," https://slack.com/.