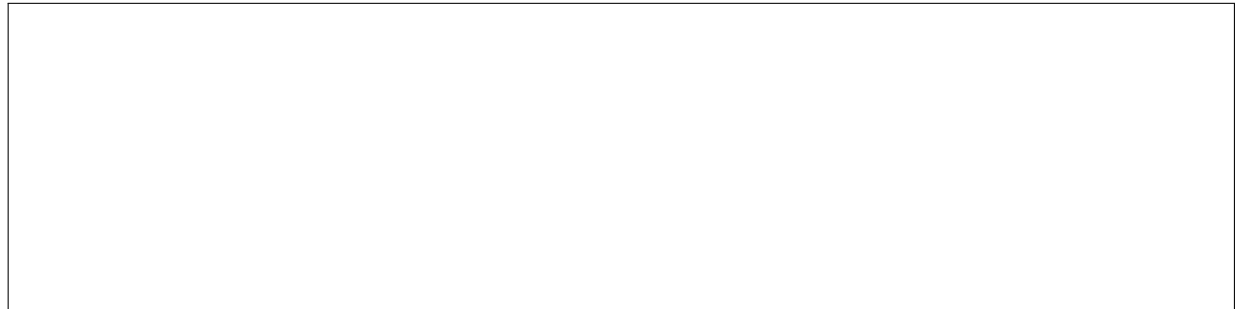


MAIN TITLE OF PAPER

Complete Research and Research in Progress



Abstract

The panorama of data is ever evolving, and big data has emerged to become one of the most hyped terms in the industry. Today, users are the perpetual producers of data that if gleaned and crunched, will yield game-changing patterns. This has introduced an important shift about the role of data in organizations and many strived to harness to power of this new material. Howbeit, institutionalizing data is not an easy task and requires the absorption of a great deal of complexity. According to various sources, it is estimated that only 13% of organizations succeeded in delivering on their data strategy. Among the root challenges, big data system development and data architecture are prominent. To this end, this study aims to facilitate data architecture and big data system development by applying well-established patterns of microservices architecture to big data systems. This objective is achieved by two systematic literature reviews, and infusion of results through thematic synthesis. The result of this work is a series of theories that explicates how microservices patterns could be useful for big data systems. These theories are then validated through a semi-structured interview with experts from the industry. The findings emerged from this study indicates that big data architecture can benefit from many principles and patterns of microservices architecture.

Keywords: big data, microservices, microservices patterns, big data architecture, data architecture, data engineering

1 Introduction

Today, we live in a world that produces data at an unprecedented rate. The attention toward these large volume of data has been growing rapidly and many strive to harness the advantages of this new material. Along these lines, academics and practitioners have considered means through which they can incorporate data-driven functions and explore patterns that were otherwise unknown. While the opportunities exist with big data (BD), there are many failed attempts. According to Davenport and Bean (2021) in 2022, only 26.5% of companies successfully become data-driven. Another survey by Databricks (2021) highlighted that only 13% of organizations succeeded in delivering on their data strategy. Among the challenges of adopting BD, data architecture, organizational culture, lack of talent, and rapid technology change are bold.

Therefore, there is an increasing need for more research on reducing the complexity involved with BD projects. One area with good potential is data architecture. Data architecture allows for a flexible and scalable BD system that can account for emerging requirements. One way to absorb the body of knowledge

available on data architecture, can be reference architectures (RAs). By presenting proven ways to solve common implementation challenges on an architectural level, RAs support the development of new systems by offering guidance and orientation (Ataei and Litchfield, 2022).

Another concept that has the potential to help with development of BD systems is the use of microservices (MS) architecture. MS architecture allows for division of complex applications into small, independent, and highly scalable parts and, therefore, increase maintainability and allows for a more flexible implementation (Richardson, 2022, p. 20). Nevertheless, design and development of MS is sophisticated, since heterogenous services have to interact with each other to achieve the overall goal of the system. One way to reduce that complexity is the use of patterns. Comparable to RAs, they are proven artifacts on how certain problems could be solved. In the realm of MS, there are numerous patterns that can be utilized, depending on the desired properties of the developed system. Despite the potential of RAs and MS architectures to solve some of complexities of BD development, to our knowledge, there is no study that properly bridge these two concepts.

To this end, this study aims to explore the application of MS patterns to BD system, in aspiration to solve some of the complexities of BD system development. For this purposes, the result of two distinct systematic literature review (SLR) are combined. The first SLR is done by Ataei and Litchfield (2020) to find all BD RAs available in the body of knowledge and to point out architectural constructus and limitations. The second SLR is conducted as part of this study to collect all MS patterns in the body of knowledge. The results of these SLRs are collected, captured and combined through thematic synthesis. As a result, various design theories are generated and validated through a semi-structured interview.

The contribution of the publication, is thereby threefold: 1) it generates numerous design theories and validate them through expert opinion, 2) it assembles an overview of relevant microservice patterns and, most importantly, 2) it creates a connection between the two SLRs to facilitate BD system development and architecture.

2 Related Work

To the best of our knowledge, there is no study in academia that has shared the same goal as our study. Laigner et al. (2020) applied an action research and reported on their experience on replacing a legacy BD system with a microservice based event-driven system. This study is not a systematic review and aims to create contextualized theory based on a specific experience. In another effort, Zhelev and Rozeva (2019) described why event-driven architectures could be a good alternative to monolithic architectures. This study does not follow any clear methodology, and seems to contribute only in terms of untested theory. Staegemann et al. (2021) examined the interplay between BD and MS by conducting a bibliometric review. This study aims to provide with a general picture of the topic, and does not aim to explore MS patterns and their relationship to BD systems. While the problem of BD system development has been approached through a RA that absorbs some of the concepts from MS as seen in Phi (Maamouri, Sfaxi, and Robbana, 2021) and Neomycelia (Ataei and Litchfield, 2021), there is no study that aimed to apply MS patterns to BD systems through a systematic methodology.

3 Methodology

Since the goal of this study is to map BD architectures and microservice patterns, it is consequently mandatory to get a comprehensive overview over both domains. For this purpose, it was decided to combine the results of two systematic literature reviews (SLR), one for each domain. Both SLRs are conducted following the guidelines presented in Kitchenham, Dyba, and Jorgensen (2004) and Page et al. (2021) on Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA). The former was used because of its clear instructions on critically appraising evidence for validity, impact and applicability in software engineering and the latter was used because it's a comprehensive and well-

established methodology for increasing systematicity, transparency, and prevention of bias. To synthesize our findings, thematic synthesis proposed by Cruzes and Dyba was applied (Cruzes and Dyba, 2011).

3.1 First Review

The first SLR, was a rigorous approach from scratch and was conducted in the following steps: 1) selecting data sources 2) developing a search strategy 3) developing inclusion and exclusion criteria, 4) developing the quality framework 5) pooling literature based on the search strategy, 6) removing duplicates, 7) scanning studies titles based on inclusion and exclusion criteria, 8) removing studies based on publication types, 9) scanning studies abstract and title based on inclusion and exclusion criteria, 10) assessing studies based on the quality framework (includes three phases), 11) extracting data from the remaining papers, 12) coding the extracted data, 13) creating themes out of codes, 14) presenting the results. These steps are not direct mappings to the following sub-sections. Some sub-sections include several of these steps.

3.1.1 Selecting data sources

To assure the comprehensiveness of the review, a broad set of scientific search engines and databases was queried. To increase the likelihood of finding all relevant contributions, it was decided to not discriminate between meta databases and publisher bound registers. Thus, both types were utilized. To achieve this, ACM Digital Library, AISel, IEEE Xplore, JSTOR, Science Direct, Scopus, Springer Link, and Wiley were included into the search process. For all of these, the initial keyword search was conducted on June 19, 2022, and there was no limitation to the considered publishing date.

3.1.2 Developing a search strategy

Since there are differences in the filters of the included search engines, it was not possible to always use the exact same search terms and settings. Nevertheless, the configurations for the search were kept as similar as possible. The exact keywords and search strategy used can be found at Ataei and Staegemann (2022c). These search terms are chosen because *patterns* are exactly what was sought for, *architectures* can contain such patterns, and *design* is often used as a synonym for architecture. Further, patterns can be seen as *building blocks*, therefore, the building blocks was also included.

3.1.3 Developing inclusion and exclusion criteria

Inspired by PRISMA checklist Tricco et al. (2018), Our inclusion and exclusion criteria are formulated as following:

Inclusion Criteria: 1) Primary and secondary studies between Jan 1st 2012 and June 19th 2022, 2) The focus of the study is on MS patterns, and MS architectural constructs, 3) Scholarly publications such as conference proceedings and journal papers.

Exclusion Criteria: 1) Studies that are not written in English, 2) Informal literature surveys without any clearly defined research questions or research process, 3) Duplicate reports of the same study (a conference and journal version of the same paper). In such cases, the conference paper was removed. 4) Complete duplicates (not just Updates) were also removed. 5) Short papers (less than 6 pages).

3.1.4 Developing the quality framework

Quality of the evidence collected as a result of this SLR has direct impact on the quality of the findings, making quality assessment an important undertaking. To address this, we developed a criteria made up of 7 elements. These criteria are informed by guidelines provided by Kitchenham Kitchenham, Dyba, and Jorgensen (2004) on empirical research in software engineering. These 7 criteria are discussed in table 1.

3.1.5 Pooling literature based on the search strategy

Overall, the keyword search yielded 3064 contributions. The total number of found publications per source as well as an overview of the further search process can be seen in Figure 1.

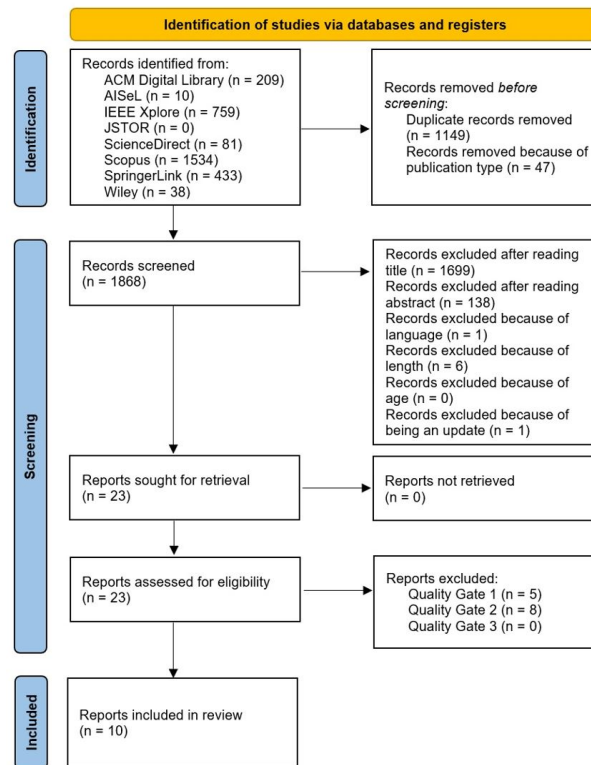


Figure 1. Overview of the search process

3.1.6 Evaluating papers based on inclusion and exclusion criteria

The remaining 1868 papers were filtered by title to evaluate their relevance to the concepts of microservice patterns or architectural constructs related to MS. For this purpose, the first two authors separately evaluated each entry. If both agreed, this verdict was honored. In case of disagreement, they discussed the title to come to a conclusion. In this phase, the first author initially included 113 papers and the second author 146. Of those, 41 were present in both sets and 1650 were excluded by both. This equates to an agreement rate of 90,5 percent (1691 of 1868 records) between the authors. After discussing the contributions with divergent evaluations, in total, 1699 of the 1868 papers were excluded, leaving 169 items for the next round.

The same approach was followed for abstracts. As a result, the first author evaluated 40 papers positively, and the second one 28. Both agreed on the inclusion of 22 papers and the exclusion of 123. This equates to an agreement rate of 85,8 percent (145 of 169 records) between the authors. In total of the 169 papers, 138 were removed and 31 were included in the next phase.

From there on, the papers that were not written in English (despite the abstract being in English), were published before the year 2012, and had a length of less than six pages were removed.

3.1.7 Evaluating papers based on the quality framework

After having filtered out the pooled studies based on inclusion and exclusion criteria, we initiated a deeper probing, by running the remaining studies against the quality framework. The filtering based on the quality

Table 1. The quality framework

Quality Gate	Criterion	Considered Aspect	Rating to pass
1	Minimum quality threshold	1) Does the study report empirical research or is it merely a 'lesson learnt' report based on expert opinion? 2) The objectives and aims of the study are clearly communicated, including the reasoning for why the study was undertaken? 3) Does the study provide with adequate information regarding the context in which the research was carried out?	5/6
2	Rigor	1) Is the research design appropriate to address the objectives of the research? 2) Is there any data collection method used and is it appropriate?	3/4
3	Credibility Relevance	1) Does the study report findings in a clear and unbiased manner? 2) Does the study provide value for practice or research?	3/4

criteria was divided into three differently focused phases, with each of them requiring the passing of a quality gate as portrayed in Table 1. In the first phase, the aim was to ensure that reports fulfill at least a desired minimum level of comprehensiveness. For this purpose, studies were evaluated for their content to see if they are actual research or just a mere report on some lessons or expert opinions. In addition, we checked if objectives, justification, aim and context of the studies are clearly communicated.

Authors independently rated the three aspects for all 23 remaining papers, giving one point respectively, if they deemed a criterion fulfilled and no point if they considered that aspect lacking. Consequently, for each aspect, zero to two points were achievable and for all aspects, six points were available per paper. For inclusion into the second phase, at least five out of six points were demanded to assure a sufficient base quality. This corresponds to having at least 75 percent of the points.

In total, the authors agreed on 51 of 69 evaluations, resulting in an agreement rate of 73,9 percent. The second phase was focused on rigor. In this phase, studies were judged based on their research design and the data collection methods. The general procedure with the first two authors independently evaluating the reports remained the same. For inclusion in the next phase, again, 75 percent of the obtainable point were needed (this time three out of four). In total, the authors agreed on 23 of 36 evaluations, resulting in an agreement rate of 63,9 percent. While this value is rather low, this is likely caused by the narrow margins for some decisions.

Once more, the papers with the highest score (this time two) were discussed before inclusion, to further counteract possible fuzziness in the individual evaluations. The remaining 10 papers went through the third and final phase. Here, the credibility of the reporting and the relevance of the findings were evaluated. The procedure was the same as the previous phases. However, this time, all of the remaining papers passed. In this last phase, the authors agreed on 14 of 20 evaluations, resulting in an agreement rate of exactly 70 percent.

All ten publications have been published in 2018 or later, with three of them being published in 2022, which shows the timeliness of the topic. Eight of the ten papers were found via Scopus, whereas the remaining two have been identified through IEEE Xplore.

3.1.8 Data synthesis

After selecting the quality papers, we embarked on the data synthesis process. For this phase we follow the guidelines of thematic synthesis discussed by Cruzes et al. Cruzes and Dyba (2011). To begin, we first extracted the following data from each paper: 1) findings, 2) research motivation, 3) author, 4) title,

5) research objectives, 6) research method, 7) year. We extracted these data through coding, using the software Nvivo. After that, we created two codes: 1) patterns, and 2) quality attributes, and coded the findings based on it. By the end of this process, various themes emerged.

3.2 Second Review

The second SLR is conducted by Ataei and Litchfield (2022) on available BD RAs in academia and industry. This study is comprehensive and covers various aspects of BD RAs such as limitations, and common architectural blocks. While the common architectural constructs was discussed in in this study, the focus was not on BD system requirements. Therefore, we assessed and coded the findings of this SLR to point out the system requirements of BD systems. This was achieved by creating a new code in Nvivo called: 'software and system requirements'. This was necessary as we needed to map patterns against requirements. This is further elaborated in the results section.

4 Results

In this section, we present with three integral elements: 1) BD software and system requirements, 2) MS patterns, 3) the mapping between the two and theories that emerged as a result.

4.1 Big Data Software and System Requirements

The results of our data synthesis emerged a few themes in regards to BD requirements. While we could find BD major building blocks and system requirements by coding the findings gained from Ataei and Litchfield (2022) work, our coding process did not include categorization and representation of these requirements. We also did not know what type of requirements is the most suitable to the goal of this study. To this end, we performed a lightweight literature review in the body of knowledge to realize three major elements: 1) the type of requirements that we need, 2) an approach to categorizing the requirements, 3) presentation of these requirements.

4.1.1 Type of requirements

System and software requirements come in different flavours and can range from a formal (mathematical) specifications to sketch on a napkin. There's been various attempts to defining and classifying software and system requirements. For the purposes of this study, we opted for a well-received approach discussed by P. A. Laplante (2017). In this approach, requirements are classified into three major types of 1) functional requirements, 2) non-functional requirements, and 3) domain requirements.

Our objective is to define the high-level requirements of BD systems, thus we do not fully explore 'non-functional' requirements. Majority of non-functional requirements are emerged from the particularities of an environment, such as a banking sector or healthcare, and do not correlate to our study. Our primary focus is one functional and domain requirements.

4.1.2 Categorizing requirements

After having filtered out the right type of requirement, we then sought for a rigorous and relevant method to categorize the requirements. For this purpose, we followed the well-established categorization method based on BD characteristics, that is the 5Vs. These 5Vs are velocity, veracity, volume, variety and value (Bughin, 2016; Rad and Ataei, 2017). We took inspiration from various studies such as the work of Nadal et al. (2017), and the requirements categories presented in NIST BD Public Working Group (Chang and Grady, 2019).

In addition, we studied the RAs developed for BD systems to understand general requirements. For this purpose, we used the SLR published by Ataei and Litchfield (2022). This SLR assessed the body of

evidence and presented with a comprehensive list of BD RAs. This study helped us realize the spectrum of BD RAs, how they are designed and the general set of requirements. By evaluating the design and requirement engineering required for BD RAs, we adjusted our initial categories of requirements and added security and privacy to it.

4.1.3 Present requirements

After knowing the type and category of requirements, We looked for a rigorous approach to present these requirements. There are numerous approaches used for software and system requirement representation including informal, semiformal and formal methods. For the purposes of this study, we opted for an informal method because it is a well established method in the industry and academia (Kassab, Neill, and P. Laplante, 2014). Our approach follows the guidelines explained in ISO/IEC/IEEE standard 29148 (2018) for representing functional requirements. We have also taken inspiration from Software Engineering Body of Knowledge (Abran et al., 2004). However, our requirement representation is organized in term of BD characteristics. These requirements are described in Table 2.

Table 2. *BD system requirements*

Volume	Vol-1) System needs to support asynchronous, streaming, and batch processing to collect data from centralized, distributed, and other sources, Vol-2) System needs to provide a scalable storage for massive data sets
Velocity	Vel-1) System needs to support slow, bursty, and high-throughput data transmission between data sources, Vel-2) System needs to stream data to data consumers in a timely manner, Vel-3) System needs to be able to ingest multiple, continuous, time varying data streams, Vel-4) System shall support fast search from streaming and processed data with high accuracy and relevancy, Vel-5) System should be able to process data in real-time or near real-time manner
Variety	Var-1) System needs to support data in various formats ranging from structured to semi-structured and unstructured data, Var-2) System needs to support aggregation, standardization, and normalization of data from disparate sources, Var-3) System shall support adaptations mechanisms for schema evolution, Var-4) System can provide mechanisms to automatically include new data sources
Value	Val-1) System needs to able to handle compute-intensive analytical processing and machine learning techniques, Val-2) System needs to support two types of analytical processing: batch and streaming, Val-3) System needs to support different output file formats for different purposes, Val-4) System needs to support streaming results to the consumers
Security & Privacy	SaP-1) System needs to protect and retain privacy and security of sensitive data, SaP-2) System needs to have access control, and multi-level, policy-driven authentication on protected data and processing nodes.
Veracity	Ver-1) System needs to support data quality curation including classification, pre-processing, format, reduction, and transformation, Ver-2) System needs to support data provenance including data life cycle management and long-term preservation.

4.2 Microservice Patterns

As a result of this SLR, our data synthesis yielded 50 MS patterns. These patterns are classified based on their function and the problem they solve. Our categories are inspired by the works of Richardson (2022). While we elaborate the patterns adopted for BD requirements in detail, the aim of our study is not to explain each MS pattern. Nevertheless, the name of all patterns found and their category can be found in Ataei and Staegemann (2022b). Additionally the definition of these patterns with examples can be found in the works of

5 Application of Microservices Design Patterns to Big Data Systems

In this section, we combine our findings from both SLRs, and present new theories on application of MS design patterns for BD systems. The patterns gleaned are established theories that are derived from actual problems in MS systems in practice, thus we do not aim to re-validate them in this study.

The main contribution of our work is to propose new theories and try to apply some of the well-known software engineering patterns to the realm of data engineering and in specific BD. Based on this, we map BD system requirements against a pattern and provide with reasoning on why such pattern might work for BD systems. We support our arguments by the means of modeling. We use Archimate (Lankhorst, 2013) as recommend in ISO/IEC/IEEE 42010 (Chaabane, Bouassida, and Jmaiel, 2017).

We posit that a pattern alone would not be significantly useful to a data engineering or a data architect, and propose that a collection of patterns in relation to current defacto standard of BD architectures is a better means of communication. To achieve this, we've portray patterns selected for each requirement in a reference architecture. We then justify the components and describe how patterns could address the requirement. These descriptions are presented as sub section, each describing one characteristic of BD systems.

5.0.1 Volume

To address the volume requirements of BD , and in specific for Vol-1 and Vol-2 we suggest the following patterns to be effective; 1) Gateway offloading, 2) API gateway, 3) External Configuration Store

5.0.1.1 Gateway Offloading and API Gateway

In a typical flow of data engineering, data goes from ingestion, to storage, to transformation and finally to serving. However there are various challenges to maintain this. One challenge is the realization of various data sources as described in Vol-1. Data comes in various formats from structured to semi-structured to unstructured, and system needs to handle different data through different interfaces. There is also streaming data that needs to be handled separately with different architectural constructs and data types. So some of the key engineering consideration for the ingestion process is that; 1) is the BD system ingesting data reliably? How frequently should data be ingested? In what volume the data typically arrives?

Given the challenges and particularities of data types, different nodes may be spawned to handle the volume of data. Another popular approach is the segregation of concerns by separating batch and streaming processing nodes. Given the requirement of horizontal scaling for BD systems, it is safe to assume that there is usually more then one node associated to ingesting data. This can be problematic as different nodes will need to account for security, privacy and regional policies, in addition to their main functionality.

This means that each node needs to reimplement the same interface for the cross-cutting concerns, which makes scalability and maintainability a daunting task. This also introduces unnecessary repetition of codes and can result in incompatible implmenetations and interfaces. To solve this problem, we explore the concept of gateway offloading and API gateway patterns. Offloading cross-cutting concerns to a single

architectural construct, achieves ‘separation of concerns’, improves security and performance. In addition, this implies processing and filtering incoming data through a well specified ingress controller.

Moreover, if data producers directly communicate with the processing nodes, they will have to update the endpoint address every now and on. This issue is exacerbated when the data producer tries to communicate to a service that is down. Given that lifecycle of a service in a typical distributed cloud environment is not deterministic and many container orchestration systems constantly recycle services to proactively address this issue, reliability and maintainability of the BD system can be compromised.

Additionally, the gateway can increase the system reliability and availability by doing a constant health check on services, and distribute traffic based on liveness probes. There are other benefits to these patterns such as weighted distribution, customized cache mechanism through specific HTTP headers, and consolidated security. This also means that if the gateway is down, service nodes won’t introduce bad state into the overall system. We have portrayed a simplistic representation of this pattern in Figure 2.

5.0.1.2 External Configuration Store

As discussed earlier, BD systems are made up of various services in order to achieve horizontal scalability. These services will have to communicate with each other in order to achieve the goal of the system. Thus, each one of them will require a set of runtime environmental configurations.

These configurations could be database network locations, feature flags, and third party credentials. Moreover, different stages of the data engineering may have different environments for different purposes, for instance, privacy engineers may require a completely different environment to achieve their requirements. Therefore, the challenge is the management of these configurations as the system scales, and enabling services to run in different environments without modification. To address this problem, we propose the external configuration store pattern.

By externalizing all nodes configuration to another service, each node can request its configuration from an external store on boot up. This can be achieved in Docker files through the CMD command, or could be written in Terraform codes for a Kubernetes pod. This pattern solves the challenges of handling large number of nodes in BD systems and provide with a scalable solution for handling configurations. This pattern is portrayed in Figure 2.

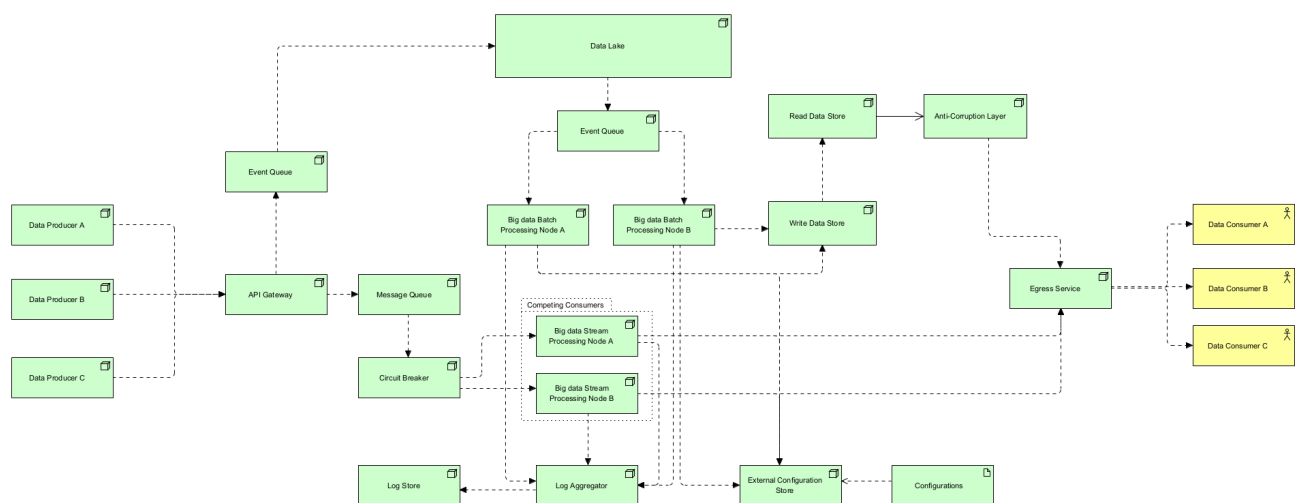


Figure 2. Big data reference architecture with microservices patterns

5.0.2 Velocity

Velocity is perhaps one of the most challenging aspects of the BD systems, which if not addressed well, can result in series of issues from system availability to massive losses and customer churn. To address some of the challenges associated with the velocity aspect of BD systems, we recommend the following patterns for the requirements Vel-1, Vel-2, Vel-3, and Vel-5; 1) Competing Consumers, 2) Circuit Breaker and 3) Log Aggregation.

5.0.2.1 Competing Consumers

BD doesn't imply only 'big' or a lot of data, it also implies the rate at which data can be ingested, stored and analyzed to produce insights. According to a recent MIT report in collaboration with Databricks, one of the main challenges of BD 'low-achievers' is the 'slow processing of large amounts of data' (Databricks, 2021). If the business desires to go data driven, it should be able to have an acceptable time-to-insight, as critical business decisions cannot wait for data engineering pipelines.

Achieving this in such a distributed setup as BD systems with so many moving parts, is a challenging task, but there are MS patterns that can be tailored to help with some of these challenges. Given the very contrived scenario of a BD system described in the previous section, at the very core, data needs to be ingested quickly, stored in a timely manner, micro-batch, batch, or stream processed, and lastly served to the consumers. So what happens if one node goes down or becomes unavailable? in a traditional Hadoop setup, if Mesos is utilized as the scheduler, the node will be restarted and will go through a lifecycle again. This means during this period of time, the node is unavailable, and any workload for stream processing has to wait, failing to achieve requirements Vel-2, Vel-3 and Vel-5. This issue is exacerbated if the system is designed and architected underlying monolithic pipeline architecture with point-to-point communications. One way to solve some of these issues, is to introduce an event driven communication as portrayed in the works of Ataei and Litchfield (2021), and try to increase fault tolerance and availability through competing consumers, circuit breaker, and log aggregation.

Underlying the event-driven approach, we can assume that nodes are sending each other events as a means of communication. This implies that node A can send an event to node B in a 'dispatch and forget' model on a certain topic. However this pattern introduces the same problem as the point-to-point REST communication style; if node B is down, then this will have a ripple effect on the whole system. To address this challenge, we can adopt the competing consumer pattern. Adopting this pattern means instead of one node listening on the topic, there will be a few nodes.

This can change the nature of the communication to asynchronous mode, and allow for a better fault tolerance, because if one node is down, the other nodes can listen to the event and handle it. In other terms, because now there are a few consumers listening on the events being dispatched on a certain topic, there's a competition of consumers, therefore the name 'competing consumers'. For instance, three stream processing consumer nodes can be spawned to listen on data streaming events being dispatched from the upstream. This pattern will help alleviate challenges in regards to Vel-2, Vel-3 and Vel-5.

5.0.2.2 Circuit Breaker

On the other hand, given the large number of nodes one can assume for any BD system, one can employ the circuit breaker pattern to signal the service unavailability. Circuit breakers can protect the overall integrity of data and processes by tripping and closing the incoming request to the service. This communicates effectively to the rest of the system that the node is unavailable, allowing engineers to handle such incidents gracefully. This pattern, mixed with competing consumers pattern can increase the overall availability and reliability of the system, and this is achieved by providing an even-driven asynchronous fault tolerance communication mechanisms among BD services. This allows system to be able to be resilient and responsive to bursty, high-throughput data as well as small, batch oriented data, addressing requirements Vel-1, Vel-4, and Vel-5.

5.0.2.3 Log Aggregator

Given that BD systems are comprising of many services, log aggregation can be implemented to shed lights on these services and their audit trail. Traditional single node logging does not work very well in distributed environments, as engineers are required to understand the whole flow of data from one end to another. To address this issue, log aggregation can be implemented, which usually comes with a unified interface that services communicate to and log their processes from. This interface then, does the necessary processes on the logs, and finally store the logs.

In addition, reliability engineers can configure alerts to be triggered underlying certain metrics. This increases teams agility to proactively resolve issues, which in turn increases reliability and availability which in turn addresses the velocity requirement of BD systems. While this design pattern does not directly affect any system requirements, it indirectly affects all of them. A simplistic presentation of this pattern is portrayed in Figure 2.

5.0.3 Variety

Variety, being another important aspect of BD, implies the range of different data types and the challenges of handling these data. As BD system grows, newer data structures emerge, and an effective BD system must be elastic enough to handle various data types. To address some of the challenges of this endeavour, we recommend the following patterns to address requirements Var-1, Var-3, Var-4; 1) API Gateway, 2) Gateway Offloading.

5.0.3.1 API Gateway and Gateway Offloading

We have previously discussed the benefits of API Gateway and Gateway Offloading, however in this section we aim to relate it more to BD system requirements Var-1, Var-3, and Var-4. Data engineers need to keep an open line of communication to data producers on changes that could break the data pipelines and analytics. Suppose that developer A changes a field in a schema of an object that may break a data pipeline or introduce a privacy threat. How can data engineers handle this scenario effectively?

To address this problem, API Gateway and Gateway Offloading can be used. API Gateway and Gateway Offloading could be good patterns to offload some of the light-weight processes that maybe associated to the data structure or the type of data. For instance, a light weight metadata check or data scrubbing can be achieved in the gateway. Moreover, as the data engineering loads increases, there will be more servers spawned for batch and stream processing. Gateways help with handling the traffic to these new servers and provide with zero down-time transitions.

However, Gateways themselves should not be taking a lot of responsibility and become a bottleneck to the system. Therefore, as nodes increase and requirements emerge, one might chose to opt for 'Backend for Frontend' pattern. We do not do any modeling for this section, as the high-level overview of API Gateway pattern is portrayed in Figure 2.

5.0.4 Value

Value is the nucleus of any BD endeavour. In fact all components of the system pursue the goal of realizing a value, that is the insight derived from the data. Howbeit, realizing these insights require absorption of great deal of complexity. To address some of these challenges, we propose the following patterns to address the requirements Val-1, Val-3, and Val-4; 1) Command and Query Responsibility Segregation (CQRS), 2) Anti-Corruption Layer.

5.0.4.1 Command and Query Responsibility Segregation

Suppose that there are various application that would like to query data in different ways and with different frequencies (Val-3, Val-4). Different consumers such as business analysts and machine learning engineers

have very different demands, and would therefore create different workloads for the BD systems. As the consumers grow, the application has to handle more object mappings and mutations to meet the consumers demands. This may result in complex validation logics, transformations, and serialization that can be write-heavy on the data storage. As a result, the data serving layer can end up with an overly complex logic that does too much.

Read and write workloads are different in nature, and this is something a data engineer should consider from the initial data modeling, to data storage, retrieval and serialization. And while the system may be more tolerant on the write side, it may have a requirement to provide reads in a timely manner (checking a fraudulent credit card in banking systems). Read and write representation of the data are often different and miss-matching and require a specific approach and modeling. For instance a snowflake schema maybe expensive for writes, but cheap for reads.

To address some of these challenges, we suggest CQRS pattern. CQRS separates the read from writes, using commands to update the data, and query to read data. This implies that the read and write databases can be physically segregated and consistency can be achieved through an event. To keep databases in sync, the write database can publish an event whenever an update occurs, and the read database can listen to it, retrieve the data, optimize for read and persist. This allows for elastic scaling of the read nodes, and increased query performance. This also allows for a read optimized data modeling tailored specifically for data consumers. Therefore, this pattern can potentially address the requirement Val-1, and Val-3. This pattern is portrated in Figure 2.

5.0.4.2 Anti-Corruption Layer

Another pattern that comes useful when handling large number of data consumers is the anti-corruption layer. Given that the number of consumers and producers can grow and data can be created and requested in different formats with different characteristics, the ingestion and serving layer may be coupled to these foreign domains and try to account for an abstraction that aims to encapsulate all the logic in regards to all the external nodes. As the system grows, this abstraction layer becomes harder to maintain, and its maintainability becomes more difficult.

One approach to solve this issue is anti-corruption layer. Anti-corruption layer is a node that can be placed between the serving layer and data consumers or data producer, isolating different systems and translating between domains. This eliminates all the complexity and coupling that could have been otherwise introduced to the ingestion layer or the serving layer. This also allows for nodes to follow the 'single responsibility' pattern. Anti-corruption layer can define strong interfaces and quickly serve new demands without affecting much of the serving node's abstraction. In another terms, it avoids corruption that may happen among systems, by separating them. This pattern can help with requirements Val-3 and Val-4. We have portrayed this pattern in Figure 2.

5.0.5 Security and Privacy

Security and privacy should be on top of mind for any BD system development, as these two aspects play an important role in the overall data strategy and architecture of the company. At the intersection of data evolution, regional policies, and company policies, there's a great deal of complexity. To this end, we propose the following pattern to address requirements SaP-1 and SaP-2; 1) Backend for Frontend (BFF)

5.0.5.1 Backend for Frontend

API gateway has been discussed in several sections in this study, however, in this section we are interested to see how it can improve security and privacy of BD systems. In terms of privacy, given the increasing load of data producers, and how they should be directed to the right processing node, how does one comply with regional policies such as GDPR? how do we ensure, for example, that data is anonymized and identifiable properties are omitted? one approach is to do this right in the API gateway. However as

data consumers grow and more data gets in, maintaining the privacy rules and applying them correctly to the dataset in the API gateway becomes more difficult. In addition, this can result in a bloated API gateway with many responsibilities, that can be a potential bottleneck to the system.

One approach to this problem can be the BFF pattern. By creating backends (services) for frontends (data producers), we can logically segregate API gateways for data that requires different level of privacy and security. This logical separation can include other factors such as quality of service (QoS), key accounts, and even the nature of the API (GraphQL or RPC). Implementing this pattern means that instead of trying to account for all privacy related concerns in one node (API gateway), we separate the concerns to a number of nodes that are each responsible for a specific requirement. This means, instead of creating a coupled, loosely abstracted implementation of privacy mechanisms, the system can benefit from hiding sensitive or unnecessary data in a logically separated node. This is also a great opportunity for data mutation, schema validation, and potentially interface change (receive REST, and return GraphQL).

On the other hand, from the security point of view, and in specific in relation to authorization and authentication, this pattern provides with a considerable advantage. BFF can be implemented to achieve token isolation, cookie termination, and a security gate before requests can reach to upstream servers. Other security procedures such as sanitization, data masking, tokenization, and obfuscation can be done in this layer as well. As these BFF servers are logically isolated for specific requirements, maintainability and scalability is increased. This addresses the requirements SaP-1 and SaP-2.

5.0.6 Veracity

Next to value, veracity is an integral component of any effective BD system. Veracity in general is about how truthful and reliable data is, and how signals can be separated from the noises. Data should conform with the expectations from the business, thus data quality should be engineered across the data lifecycle. According to Eryurek et al. (2021), data quality can be defined by three main characteristics 1) accuracy, 2) completeness, and 3) timeliness. Each of these characteristics posit a certain level of challenge to architecture and engineering of BD systems. To this, we propose the following patterns for addressing requirements Ver-1, and Ver-4; 1) Pipes and Filters, 2) Circuit breaker

5.0.6.1 Pipes and Filters

Suppose that there is a data processing node that is responsible for performing variety of data transformation and other processes with different level of complexities. As requirements emerge, newer approaches of processing may be required, and soon this node will turn into a big monolithic unit that aims to achieve too much. Furthermore, this node is likely to reduce the opportunities of optimization, refactoring, testing and reusing. In addition, as the business requirements emerge, the nature of some of these tasks may be different. Some processes may require a different metadata strategy that requires more computing resources, while others might not require such expensive resources. This is not elastic and can produce unwanted idle times.

One approach to this problem could be the pipes and filters pattern. By implementing pipes and filters, processing required for each stream can be separated into its own node (filter) that performs a single task. This is a well-established approach in unix-like operating systems. Following this approach allows for standardization of the format of the data and processing required for each step. This can help avoiding code duplication, and results in easier removal, replacement, augmentation and customization of data processing pipelines, addressing the requirements Ver-1 and Ver-4.

5.0.6.2 Circuit breaker

In an inherently distributed environment like BD, calls to different services may fail due to various issues such as timeouts, transient faults or service being unavailable. While these faults may be transient, this

can have a ripple effect on other services in the system, causing a cascading failure across several nodes. This affects system availability and reliability and can cause major losses to the business.

One solution to this problem can be the circuit breaker pattern. Circuit breaker is a pattern that prevents an application from repeatedly trying to access a service that is not available. This improves the fault tolerance among services, and signals the service unavailability. The requesting application can decide accordingly on how to handle the situation. In other terms, circuit breakers are like proxies for operations that might fail. This proxy is usually implemented as a state machine having the states close, open, and half-open. Having this proxy in place provides stability to the overall BD system, when the service of interest is recovering from an incident. This can indirectly help with Ver-4.

6 Validation

After the generation of the design theories, we sought for a suitable model of validation. This involved a thorough research in some of the well-established methods for validation such as single-case mechanism experiment, technical action research and focus groups Wieringa, 2014. For the purposes of this study we chose semi-structured interviews (SSIs), following the guidelines of Kallio et al. Kallio et al., 2016.

6.1 Interview design

Our SSI methodology is made up of four phases: 1) identifying the rationale for using semi-structured interviews, 2) formulating the preliminary semi-structured interview guide, 3) pilot testing the interview guide, 4) presenting the results of the interview. SSI are suitable for our study, because our conceptual framework is made up of architectural constructs that can benefit from in-depth probing and analysis. Our questions are categorized into main themes and follow-up questions, with main themes being progressing and logical, as recommended by Kallio et al., 2016. We pilot tested our interview guide using internal testing, which involved an evaluation of the preliminary interview guide with the members of the research team.

6.2 Sampling strategy

After having our interview guide designed, we used purposive sampling (Baltes and Ralph, 2022) to select experts. We chose purposive sampling because it allowed us to collect rich information by expert sampling. In addition, this approach enabled us to ensure representativeness and removed the need for a sampling frame. We also attempted 'heterogeneity sampling' by approaching candidates from various industries. We reached out to colleagues, our connections on ResearchGate and LinkedIn, and tried to look for experts with the titles 'data engineer', 'data architect', 'senior data engineer', 'solution architect', 'lead architect', and 'big data architect'. We also included founders of big data companies, or academics who have been working on BD systems. We interviewed 6 experts from various industries over a period of 3 months.

6.3 Data Synthesis

All the interviews have been done through the software Zoom. We saved all of the recordings, and then downloaded the automatically generated transcripts. Transcripts for each interview has been added to Nvivo and then codes are created. We created a code for each BD characteristics discussed in Section ?? . After the initial coding process, through axial coding, we created higher level codes. These higher level codes were connected to create themes.

6.4 Results

From the results of these interviews, we gathered a lot of insights and probed deep into our architectural constructs. Every interview involved in deep analysis of the design patterns with one question from the interviewee trying to understand the problem space and solution better. Our interviewees had at least 8 years of experience. Our interview guide is available at Ataei and Staegemann, 2022a.

6.4.1 Volume

For volume, we went through the theories elaborated in Section ???. All of the interviewees took the idea of API gateway and gateway offloading naturally, while we had to explore the 'external configuration store' a bit deeper. We used the idea of Kubernetes ingress to help with elaboration of API gateway. We used AWS load balancer example, and discussed the challenges of maintaining certificates and authentication. For the externalized configuration pattern we had to go a bit deeper and talk about a scenario in which the developer of the batch processing node may need to account for the development, trial and production workloads that have different DNS requirements, and configurations. We discussed how environment variables may vary, have development and trial environments may not need as much resources as the production, how ingress may vary.

The quick brown fox jumps right over the lazy dog. the quick brown fox jumps right over the lazy dog. the quick brown fox jumps right over the lazy dog. the quick brown fox jumps right over the lazy dog. the quick brown fox jumps right over the lazy dog. the quick brown fox jumps right over the lazy dog.

After explaining a scenario, interviewees agreed that external configuration store pattern can help with some of the challenges of data engineering. One interviewee mentioned that this can even be utilized for special privacy requirements. In addition, there's been discussion in regards to regional privacy and security requirements and how configuration can help derive them. Some interviewees discussed that this is a general pattern that any system can utilize to its benefit.

Another interviewee discussed how they are taking extensive measures to embark on a fully event-driven process, and how a lot of things that we theorize and modeled may sound easy to do, but daunting to implement. The interviewee explained how they are planning to store data in their AWS S3 initially and then having a Lambda function trigger to start the ETL process. The interviewee then explained how they need to obtain different configurations from different data providers, and how that can affect the data prepared for data consumers. He discussed an archetype in which the data would be stored in AWS S3, which triggers a Lambda function to initial an ETL process. Furthermore, he added how externalized configuration pattern could be implemented with DynamoDB and Lambda functions.

One of the interviewees from insurance and finance sector mentioned that scaling the gateway and corresponding nodes may not be as easy as it seems. He mentioned that during normal days there are hardly any claims, and while there's a special event, the storm comes. The interviewee mentioned that scaling forecast is usually based on historical data. Further, he mentioned that even the delay in auto-scaling groups in AWS can be problematic for them.

The same interviewee from the insurance sector discussed how centralizing configuration may sound like a good idea. Howbeit, he added that this approach makes him slightly nervous, because every service is unique in its own, and may require a specific configuration. He added, that as configurations store, the externalized configuration node can be bloated, taking so much responsibilities. He added that at times, his team had to reconfigure a service at the fly to prevent customer churn, and with this pattern he finds everything more complicated. At last, he added that in a multi-region operating companies, a centralized configuration store can really help with standardization and maintenance.

Another interiveee from the financial sector has affirmed us that gateway offloading and API gateways are pretty common patterns, and he has witnessed it in several major banks. The same candidate has elaborated that 'external configuration store' pattern is sometimes referred to as 'declarative configuration

management'. The candidate then continued to explain how this pattern can be witnessed in Kubernetes clusters through metadata objects, kube-system and Etcd.

6.4.2 Velocity

For Velocity, we first started by exploring an event-driven data engineering architecture, and then justified the idea of competing consumers. We then explored how competing consumers can fail, and how circuit breaker pattern can help. Finally we explore the idea of logging and how tail logging and distributed tracing can be achieved through it. An interviewee challenged the idea of competing consumer and stated that a leader election may be a better choice for a distributed setup as such. The interviewee also mentioned that circuit breakers could be implemented on the competing consumers themselves, but he could see the value of separating it to its own service. Of particular argument was the fact that circuit breaker's implementation may not be that complicated and a dedicated service for it can increase costs.

In another interview, interviewee asked about the amalgamation technique for the logs, and discussed how dimensionality of the logs can be challenging. We took both feedbacks of 'leader election' and 'more in-detailed logging approach' into consideration. We researched deeper into leader election, logging approaches, and distributed tracing. We found leader election a bit hard to justify, as it introduces a single point of failure, can potentially introduce faulty data as there's only one point of trust, and partial deployments are really hard to apply. We found that benefits of 'leader election' pattern to be outweighed by the complexity it introduces. In regards to logging, we found various approaches to distributed tracing and log merging, however these were mostly in-detailed micro approaches, which is not in the scope of our study.

In another interview, the interviewee discussed how circuit breakers may need to do load balancing as well. We then discussed how circuit breakers could be implemented in data processing nodes themselves, or in a side car. The interviewee then explained how they've created a system that resembles to the log aggregator pattern. The interviewee elaborated how the system has a graphical interface that captures errors from various ETL jobs.

An interviewee from the insurance sector discussed how log aggregator might be a good pattern, but it's not always great to add so many technologies to the stack. Then he added that each system may have a different logging library and interface and aggregating them may need an effective methodology. The interviewee described that logging is better be approached through several layers of abstraction. He described how it would be useful to have some easy to understand metrics on the surface level. He added there should not be a need for technical skills to read these metrics. Nevertheless, there should be detailed logs abstracted for more technical users.

Moreover, it was mentioned that only important pieces of information should be collected and presented, as most web servers such as Nginx create so many logs.

In an another interviewee, the candidate brought to light the challenges of time-sychroonization in log aggregator pattern. The candidate, who had a background in financial sector, discussed how handling logs from a large amount of can introduce a challenge of its own. He continued to describe how critical these logs can be during sensitive stream processing tasks, and how data can easily get into petabytes in the banks.

The candidate recommended to design services in a manner that promotes self-awareness. This is to prevent them from breaking silently, which makes debugging and issue resolution much longer. He added that this 'awareness' can be complementary to log aggregator, as services can reflect and dispatch an event in regards to the root cause of the failure. In addition, the candidate discussed the benefits of dynamically defining the level of logging.

He illustrated how dynamically setting the level of logging has been really helpful in his personal experience. The candidate then explore further on low-level technical details of implementing OpenTelemetry for different cases and with different level of logging.

6.4.3 Variety

For variety, we discussed common data types that need support, and how system may use parquet, JSON, or how unstructured data can introduce challenges. By this point, interviewees had a better grasp of our models and the gateway patterns, thus there wasn't much questions. An interviewee suggested the 'API composition' pattern and suggested that we may have various services that handle different data types, but the composition of these data may be necessary. The interviewee suggested that 'API composition' can occur at the egress level.

One interviewee provided details on how painful it has been for his team to onboard new data producers and how that dramatically slowed the project deadline. The interviewee added that data received from data producers hardly have the standards necessary, as these data are generated by third-party software that they have no control over. He explained how different versions of the same software create different schema and how this can sometimes break the data engineering pipelines. Then, the interviewee suggested off-loading more compute intensive checks to gateways. We discussed how that could result in a bloated architectural construct and both parties decided that BFF pattern is probably a better suit.

Another interviewee from insurance sector discussed how the rate of change is very low and most things are standard in insurance and finance sectors. For instance, he stated that if Avro is being used as the data format, the industry will be using the same format for the next 5 years. Additionally, the interviewee explained breaking that changes, specially schema changes are usually avoided. He added that in German insurance companies, almost everything is standard, and introducing any change would require large scale communication with all insurance companies which is an extensive measure.

One of the participants who had an experience with firmware development, depicted the challenges of working with Eletronic Data Interchange (EDI) formats. In his experience the data format has hardly changed despite the recent technological advancements, and that had introduced significant challenges to his team.

He then explained how gateway offloading could be useful to isolate this data format only to a group of specialized engineers. This meant that newer, less interested engineers could be working on different nodes concurrently without having to worry about introducing side effects to the pipelines. He mentioned that at times, there were very few people available who were well-aware of EDI. He explained how at the very least with the gateways, the data could be stored in a storage for later processing through special headers.

6.4.4 Value

For value, we discussed CQRS and anti-corruption layer. We first began by exploring the challenges of having to optimize for read and write loads. We discussed how it could be essential for the business to provide read queries in a timely manner, and how trying to model for both read and write queries may not be efficient. For instance, we explored snowflake schemas against star schemas, discussed a typical data analysis flow and provided with challenges. None of our interviewees have had an experience with CQRS in practice, but they were aware of the pattern.

An interviewee discussed how this pattern can be helpful in companies that have adopted domain-driven design, and how each bounded context can decide how the data should be modeled. Our interviewees shared the same idea that CQRS should only be applied when the needs arise and not proactively. This is due to the fact that implementing and getting CQRS right comes with complexity, and can dramatically increase cost. An interviewee suggested that CQRS is perhaps unnecessary in many cases and should be utilized only in special cases. The interviewee also suggested that a reporting database should suffice, and discussed other optimization strategies that could be applied to optimize for read and write without needing to implement CQRS. One example is using different access techniques for reads and writes.

These interviews shed some lights on how complex implementing CQRS can be, and we deduced that this pattern can introduce challenges and should be adopted when the benefits outweigh the challenges. We also received questions about event sourcing and if that could be applied, as CQRS is usually implemented

with event sourcing. However, we do not think that event sourcing can scale to account for big data systems, and the challenges of maintaining event logs can introduce risk to modifiability and performance of the overall system.

An interviewee discussed how they have implemented something similar to CQRS with Elasticsearch. Another interviewee mentioned how a lot of things are going on in their MySQL databases, and how during write-heavy times, database is locked and unresponsive. He added how waiting for database to become available again has been a pain point, and how their services have timed out on this. The interviewee explained how the stochastic nature of database locks, even harder for them to predict and tackle this issue. This interviewee found the idea behind CQRS relevant and effective in solving some of their problems.

Another interviewee discussed how they have successfully deployed CQRS into production and how it's been really effective for them. For instance, the interviewee discussed that Avro data format has been utilized for the write data store and how without this approach the cost of operation and infrastructure would have been doubled. He added that only a part of complexity is associated to bringing data to the platform and storing it in the write database. He discussed that different data consumers have different use cases, and not everyone would appreciate Parquet data format. He stated that some consumers are more interested in row-based data formats and need more aggregation.

Along the same lines, the interviewee depicted the fact that human side of things is just as complicated as the technical side. For instance, he gave us several examples in which the data consumer did not actually know what's really the most optimized format for his/her workload. This is due to the fact that some consumers are not technical stakeholders, and need to be accounted for. The interviewee continued describing how his team has to sometimes go to the data consumer directly and understand the usage patterns or algorithms run on the data. From there on, his team then would decide the best data format. Nevertheless, as stakeholders change and requirements emerge, there might be a need for doing this several times, which introduces constant challenge to data engineers.

Another interesting fact we learnt was that in financial and insurance sectors, it's not that unlikely for people to press a button on Friday and come back to get their data on Monday. He added that there are various Fortran and legacy Java applications that are widely used in practice, and are really un-optimized.

In another interview, the interviewee discussed the known issues for not applying CQRS to big data systems. The main argument was around the management of overall data volume and the stress that CQRS can introduce to storage media. The interviewee discussed how CQRS is challenging even in non data-intensive systems, and how BD can make it challenging. In addition, as discussed by the interviewee, the network and OS overheads introduces CQRS and microservices in general Sriraman and Wenisch, 2018 may not perform well in BD systems.

Furthermore, we explained the anti-corruption layer. We discussed how the consumer needs can emerge, and how coupling it all to the read service can affect the system modifiability negatively. This pattern was well-perceived by interviewees, however there were concerns about the anti-corruption layer itself getting bloated and introducing 'corruption'! However a system architect can tackle this by introducing several anti-corruption layers, or egress nodes that are each responsible for a class of data consumers.

An interviewee raised the concern that defining the scope of anti-corruption layer may be a challenge. In his experience, data scientists need 'all the data' available in the company, and that's been a challenge for his team in the past. He continued discussing that this pattern can be useful not only from decoupling perspective but from a security and governance point of view. We failed to realize this in our research. The interviewee discussed that at times his team has been asked to provide with a lot of data, and providing it could have cause major security issues. He added that defining these anti-corruption layers with clearly defined contracts between the consumers and the big data system canonical data can be an effective measure to govern what should be provided, and is a great opportunity to eliminate security risks.

6.4.5 Security and Privacy

For security and privacy, we started the discussion by exploring how different companies and regions may have different requirements, and how consuming data from data producers might be affected. We then discussed how having a single gateway to encompass all that logic can be daunting to scale. We then introduced the BFF pattern and elaborated that how each class of data consumers can be associated to a specific BFF. This pattern was well-received. An interviewee pointed out a potential of the access token pattern to be applied to the BFF. The interviewee elaborated that how having BFFs can help with cloud design and potential usage of private networks to increase security.

An interviewee discussed how data engineers are usually not well educated on security matters in his professional experience. He added how expensive it is to train engineers to a good level on security and privacy and even after that the company may not be able to retain them. The interviewee explained how IT giant companies like Google have the resources necessary to constantly account for emerging privacy and security requirements, while small to medium sized businesses are struggling. Finally he stated that following privacy and security standards is really costly for companies.

In another interview, the participant elaborated on how challenging it would be to have several ingresses into the system, and how BFF pattern may provide some stress on security and platform teams. While he admitted that performance and maintainability may be increase, he found challenges of controlling what comes into the system significant.

The interviewee added that going BFF requires substantial resources and may not be ideal for every company. From his perspective, BFF was only an unnecessary complexity, stating that his life is gonna be hard if he brought something like this into production.

Moreover, an interview discussed how encryption should be taken more seriously in BD world. He admitted that in his experience, most BD architects and data engineers were not in favour of data encryption. This was due to performance issues associated to encryption of large amount of data. The interviewee then further elaborated on issues that may arise if data is not encrypted. From his perspective, in today's world there is really no borders with data connectivity and one has to make sure that data is safe. He added that if the architecture depends on the perimeters, you need to make sure these perimeters are concretely defined, he then stated that 'there are no concretely defined perimeters'!

In his view, having access to data storage should not mean having access to data. He suggested hardware encryption to solve some of the performance challenges. The same interviewee pointed out the challenges of GDPR and privacy. He suggested that 'deleting data' is as important as storing it, and one should proactively look for opportunities to delete sensitive data.

6.4.6 Veracity

For veracity, we discussed the transformation process required for any data engineering pipeline, and how this pattern was tacitly incorporated to all of our models. We made an analogy to Linux philosophies and how different commands pipe into each other. This pattern was perceived to be one of the defacto patterns that many data engineering pipelines use. In addition, we discussed circuit breakers again, but this time for capturing transformation nodes that are unavailable. This section did not require a lot of deep probing and discussions.

An interviewee discussed how pipe and filters have been the key for them in production, and how it helped them scale and avoid data corruption. He added that without adopting such pattern, if something broke in a large transformation, you'd never know what went wrong, and you might be forced to rerun a process that takes 5 hours to complete. Furthermore, the interviewee depicted how data quality is becoming more and more important for his team and company. This is due to the fact that the interviewee works in the insurance sector, and data is used in deciding some of the claims.

The interviewee admitted that sometimes in the past, many years ago, they had to make difficult decisions because data did not possess the qualities necessary. Moreover, he added that separating transformation into their own service (filters), creates an opportunity for introducing data quality metrics for each

transformation, which can be used later to probe what has gone wrong and the team can recover from a corrupted data.

In another interview, the participant discussed how circuit breaker should be tied to the end of the data processing and not only to the beginning of it. He elaborated that the server might be healthy when the transformation starts, but that might not be the case when it's about to end, therefore corrupting the data. He added that this can introduce unnecessary reprocessing.

6.4.7 Other feedbacks, closing thoughts

The most experienced interviewee (14 years) suggested us to further break down our processing requirements into domains and then utilize gateway aggregate patterns to do 'data as a service'. This idea was driven by data mesh and data fabrics. All of our interviewees found the study interesting, and were eager to know more after the interview.

Another feedback was the idea of having an egress that encapsulate the anti-corruption layer and adds some more into it as well. The pattern 'backend for frontend' was well received, and our event driven thinking seemed to be very well accepted by the interviewees. By the result of this interview we realized that we have missed an architectural construct while discussing velocity requirements, which was the message queue. These interviews increased our confidence in our results and reasoning and have shed some lights on possible new patterns that we could employ.

We have received a lot of good insights into how else we could model and approach this problem. While some of these ideas are really interesting, due to time and resource constraint, we opted not to apply all suggestions for the purposes of this study. Lastly, some of our interviewees connected some of the patterns discussed to their own context of practice and helped us realized further improvements.

7 Discussion

The result of this study have provided us with two major findings; 1) the progress in the data engineering space seems to be uneven in comparison to software engineering, 2) MS pattern provide with a great potential for resolving some of the BD system development challenges. While there has been adoption of a few practices from software engineering into data engineer like DataOps, we posit that data engineering space can benefit from some of the well-established practices of software engineering.

Majority of the studies that we've analyzed to understand BD systems, seems to revolve around crunching and transforming data without much attention to data lifecycle management. This is bold when it comes to addressing major cross-cutting concerns of successful data engineering practice such as security, data quality, DataOps, data architecture, data interoperability, data versioning and testing. In fact, while we found a lot of mature approaches in MS and event driven architectures, we could not find many well-established patterns in the data engineering space. Based on this, we think that data architecture remains a significant challenge and requires more attention from both academia and industry.

8 Conclusion

With all the undeniable benefits of BD, the success rate of BD projects is still rare. One of the core challenges of adopting BD lies in data architecture and data engineering. While software engineers have adopted many well-established methods and technologies, data engineering and BD architectures don't seem to benefit a lot from these advancements.

The aim of this study was to explore the relationship and application of MS architecture to BD systems through two distinct SLR. The results derived from these SLRs presented us with interesting data on the potential of MS patterns for BD systems. Given the distributed nature of BD systems, MS architectures seems to be a natural fit to solve myriad of problems that comes with decentralization. Even though we created many design theories, modeled patterns against systems, and validated our theories, we believe

that our results could be further validated by an empirical study. We therefore posit that there is a need for more attention in the area of MS and event-driven architectures in relation to BD systems from both academia and industry.

References

- 29148, I. (2018). *ISO/IEC 29148:2018. Systems and software engineering — Life cycle processes — Requirements engineering*. Ed. by I. 29148. URL: <https://www.iso.org/standard/72089.html>.
- Abran, A., J. W. Moore, P. Bourque, R. Dupuis, and L. Tripp (2004). "Software engineering body of knowledge." *IEEE Computer Society, Angela Burgess*, 25.
- Ataei, P. and A. Litchfield (2020). "Big Data Reference Architectures, a systematic literature review." In: *Australasian Conference on Information Systems (ACIS) 2020*. AIS.
- (2021). "NeoMycelia: A software reference architecture for big data systems." In: *2021 28th Asia-Pacific Software Engineering Conference (APSEC)*. IEEE, pp. 452–462.
- (2022). "The state of big data reference architectures: a systematic literature review." *IEEE Access*.
- Ataei, P. and D. Staegemann (2022a). *Interview guide for the paper: Application of microservices patterns to big data systems*. URL: <https://anonymous.4open.science/r/SSI-Repo-F90E/SSI.pdf>.
- (2022b). *Microservices patterns classified for the Paper Titled: Application of Microservices Patterns to Big Data Systems*. URL: <https://anonymous.4open.science/r/MS-Patterns-5519/>.
- (2022c). *Systematic literature review search terms table for the Paper Titled: Application of Microservices Patterns to Big Data Systems*. URL: <https://anonymous.4open.science/r/SLR-Search-Terms-3147/>.
- Baltes, S. and P. Ralph (2022). "Sampling in software engineering research: A critical review and guidelines." *Empirical Software Engineering* 27 (4), 1–31.
- Bughin, J. (2016). "Big data, Big bang?" *Journal of Big Data* 3 (1), 2. ISSN: 2196-1115. DOI: 10.1186/s40537-015-0014-3.
- Chaabane, M., I. Bouassida, and M. Jmaiel (2017). "System of systems software architecture description using the ISO/IEC/IEEE 42010 standard." In: *Proceedings of the Symposium on Applied Computing*, pp. 1793–1798.
- Chang, W. L. and N. Grady (2019). *NIST Big Data Interoperability Framework: Volume 1, Definitions*. Ed. by W. L. Chang and N. Grady. URL: <https://doi.org/10.6028/NIST.SP.1500-1r2>.
- Cruzes, D. S. and T. Dyba (2011). "Recommended Steps for Thematic Synthesis in Software Engineering." In: *2011 International Symposium on Empirical Software Engineering and Measurement*. IEEE, pp. 275–284. ISBN: 978-1-4577-2203-5. DOI: 10.1109/ESEM.2011.36.
- Databricks, M. technology review insights in partnership with (2021). *Building a high-performance data organization*. Tech. rep. URL: <https://databricks.com/p/whitepaper/mit-technology-review-insights-report>.
- Davenport, T. H. and D. D. Bean (2021). *Big Data and AI Executive Survey 2021*. Technical Report. NewVantage Partners. URL: <https://www.newvantage.com/thoughtleadership>.
- Eryurek, E., U. Gilad, V. Lakshmanan, A. Kibunguchy-Grant, and J. Ashdown (2021). *Data Governance: The Definitive Guide*. "O'Reilly Media, Inc."
- Kallio, H., A.-M. Pietilä, M. Johnson, and M. Kangasniemi (2016). "Systematic methodological review: developing a framework for a qualitative semi-structured interview guide." *Journal of advanced nursing* 72 (12), 2954–2965.
- Kassab, M., C. Neill, and P. Laplante (2014). "State of practice in requirements engineering: contemporary data." *Innovations in Systems and Software Engineering* 10 (4), 235–241.
- Kitchenham, B. A., T. Dyba, and M. Jorgensen (2004). "Evidence-based software engineering." In: *Proceedings of the 26th International Conference on Software Engineering*. IEEE Comput. Soc, pp. 273–281. ISBN: 0-7695-2163-0. DOI: 10.1109/ICSE.2004.1317449.

- Laigner, R., M. Kalinowski, P. Diniz, L. Barros, C. Cassino, M. Lemos, D. Arruda, S. Lifschitz, and Y. Zhou (2020). "From a monolithic big data system to a microservices event-driven architecture." In: *2020 46th Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*. IEEE, pp. 213–220.
- Lankhorst, M. (2013). "A Language for Enterprise Modelling." In: *Enterprise Architecture at Work*. Springer, pp. 75–114.
- Laplante, P. A. (2017). *Requirements engineering for software and systems*. Auerbach Publications.
- Maamouri, A., L. Sfaxi, and R. Robbana (2021). "Phi: A Generic Microservices-Based Big Data Architecture." In: *European, Mediterranean, and Middle Eastern Conference on Information Systems*. Springer, pp. 3–16.
- Nadal, S., V. Herrero, O. Romero, A. Abelló, X. Franch, S. Vansummeren, and D. Valerio (2017). "A software reference architecture for semantic-aware Big Data systems." *Information and software technology* 90, 75–92.
- Page, M. J., D. Moher, P. M. Bossuyt, I. Boutron, T. C. Hoffmann, C. D. Mulrow, L. Shamseer, J. M. Tetzlaff, E. A. Akl, S. E. Brennan, R. Chou, J. Glanville, J. M. Grimshaw, A. Hróbjartsson, M. M. Lalu, T. Li, E. W. Loder, E. Mayo-Wilson, S. McDonald, L. A. McGuinness, L. A. Stewart, J. Thomas, A. C. Tricco, V. A. Welch, P. Whiting, and J. E. McKenzie (2021). "PRISMA 2020 explanation and elaboration: updated guidance and exemplars for reporting systematic reviews." *BMJ (Clinical research ed.)* 372, n160. DOI: 10.1136/bmj.n160.
- Rad, B. B. and P. Ataei (2017). "The big data ecosystem and its environs." *International Journal of Computer Science and Network Security (IJCSNS)* 17 (3), 38.
- Richardson, C. (2022). *A pattern language for microservices*. URL: <https://microservices.io/patterns/index.html>.
- Sriraman, A. and T. F. Wenisch (2018). "μ suite: a benchmark suite for microservices." In: *2018 IEEE International Symposium on Workload Characterization (IISWC)*. IEEE, pp. 1–12.
- Staegemann, D., M. Volk, A. Shakir, E. Lautenschläger, and K. Turowski (2021). "Examining the Interplay Between Big Data and Microservices—A Bibliometric Review." *Complex Systems Informatics and Modeling Quarterly* 27 (27), 87–118.
- Tricco, A. C., E. Lillie, W. Zarin, K. K. O'Brien, H. Colquhoun, D. Levac, D. Moher, M. D. Peters, T. Horsley, L. Weeks, et al. (2018). "PRISMA extension for scoping reviews (PRISMA-ScR): checklist and explanation." *Annals of internal medicine* 169 (7), 467–473.
- Wieringa, R. J. (2014). *Design science methodology for information systems and software engineering*. Springer.
- Zhelev, S. and A. Rozeva (2019). "Using microservices and event driven architecture for big data stream processing." In: *AIP Conference Proceedings*. Vol. 2172. 1. AIP Publishing LLC, p. 090010.