

Application of Microservices Patterns to Big Data Systems

1st Pouya Ataei

*School of Engineering, Computer and Mathematical Science
Auckland University of Technology
Auckland, New Zealand
pouya.ataei@aut.ac.nz*

2nd Daniel Staegemann

*Magdeburg Research and Competence Cluster VLBA
Otto-von-Guericke University Magdeburg
Magdeburg, Germany
daniel.staegemann@ovgu.de*

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. However, 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 approximately 70% of big data projects fail to deliver. Among the root causes of these failures, 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 explicate 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.

Index Terms—big data, microservices, microservices patterns, big data architecture, data architecture, data engineering,

I. 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, academicians 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, there are many failed attempts. According to New Vantage Partners report in 2022, only 26.5% of companies successfully become data-driven [1]. Another survey by Databricks highlighted that only 13% of organizations succeeded in delivering on their data strategy [2].

Therefore, there is an increasing need for more research on reducing the complexity involved with big data projects [3]. One area with good potential is data architecture. Data architecture allows for a flexible and scalable big data 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. A comprehensive overview of the existing big data reference architectures up until 2020 has been given in [4].

Another concept that has the potential to help with development of big data systems is the use of microservices (MS) architecture [5]. MS architecture allows for division of complex applications into small, independent, and highly scalable parts and, therefore, increase maintainability and allow for a more flexible implementation [6]. Nevertheless, design and development of MS is sophisticated, since heterogeneous 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.

While practitioners in the domain of MS architecture seem to benefit from many well-established practices, data engineering does not seem to be absorbing many of these concepts. Despite the potential of RAs and MS architectures to solve some of complexities of big data 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, two distinct systematic literature review (SLR) is conducted. The first SLR is an updated SLR on BD RAs [7], aiming to capture the years 2020-2022, and the second SLR is on microservices patterns. The result of these SLRs are then captured through thematic synthesis, and design theories are generated. The theories are then explained and finally, a semi-structured interview is conducted to validate these theories.

The contribution of the publication at hand is thereby threefold. It provides an updated synopsis of the existing big data reference architectures, it assembles an overview of relevant microservice patterns and, most importantly, it creates a connection between the two to facilitate big data system development and architecture.

II. 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. [8] applied an action research and reported on their experience on replacing a legacy BDS 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 et al [9] 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 [10] examined the interplay between big data and microservices by conducting a bibliometric review. This study aims to provide with a general picture of the topic, and does not aim to explore microservices patterns and their relationship to big data systems.

While the problem of big data system development has been approached through a RA that absorbs some of the concept from microservices as seen in Phi [11] and Neomycelia [12], there is no study that aimed to apply microservices patterns to big data systems through a systematic methodology.

III. METHODOLOGY

Since the goal of this study is to map big data architectures and microservice patterns, it is consequently mandatory to get a comprehensive overview over both domains. For this purpose, it was decided to conduct two systematic literature reviews (SLR), one for each domain.

Both SLRs are conducted following the guidelines presented in Kitchenham et al. [13] and Page et al. [14]. 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.

While, Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) provided us with a strong underpinning on how to conduct a SLR, it had many assumptions that arose from the healthcare and nursing domains, and did not correlate directly to our study. To overcome some of this limitations, we combined PRISMA with the guidelines of Kitchenham et al. to form a rigorous methodology.

To synthesize our findings, thematic synthesis proposed by Cruzes and Dyba was applied [15]. While it was initially planned to capture and present all identified microservice patterns using the Buschmann et al. template [16], this was later on omitted to save space and because a considerable number of them can already be found in the works of Richardson [17].

A. First Review

The first SLR was the less extensive one, since it is just an update to an already existing study. For this purpose, the SLR conducted by Ataei et al. [7], as a thematically fitting study on big data reference architectures was extended up unto the current date, providing us with the necessary overview. We followed the exact same methodology and covered the years

2020-2022. Our main objective for this SLR was to highlight the fundamental building blocks and requirements of big data systems.

While the common architectural constructs was discussed in [7], the study was not focused on requirements. We therefore extended the data synthesis by adding a new code: software and system requirements. This code had sub-codes each being named after one characteristics of big data. This was necessary as we needed to map patterns against requirements. This is further elaborated in the results section, subsection A.

B. Second Review

The second 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.

1) *Selecting data sources*: To assure the comprehensiveness of the review and following the recommendations of PRISMA-S [18], 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.

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 [19]. For those engines that yielded a high number of results, the scope was reduced by adding variations of "pattern", "architecture", "design", "building block", or "best practice" to appear in title, abstract or keywords.

If this could not be realized because of the interface, the most similar setting that is more lenient (therefore, potentially yielding more results) was chosen. This was the case for IEEE Xplore and SpringerLink. 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. Finally, the use of patterns is often highlighted as a best practice and hence, in reverse, papers that refer to best practices might also contain information regarding the use of patterns.

As it can be seen at [19], due to the specifics of their search masks, the searches in the ACM Digital Library (title, abstract, keywords could only be searched separately), JSTOR (no support of wildcards), and Science Direct (no support of wildcards) had to be split in several parts. Those were afterwards merged for each of them, and duplicates were removed.

3) *Developing inclusion and exclusion criteria:* Our inclusion and exclusion criteria is inspired by the PRISMA checklist [20] and the works of Ataei et al. [4], is 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 microservices patterns, and microservices 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).

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 those proposed by CASP for assessing the quality of qualitative research [21] and by guidelines provided by Kitchenham [22] on empirical research in software engineering. These 7 criteria are discussed in I.

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.

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 microservices. 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. With a difference that authors agreed to also allow themselves to look into the actual paper and not just the abstract, if they wanted to further explore certain aspects of the study to improve their judgement. 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

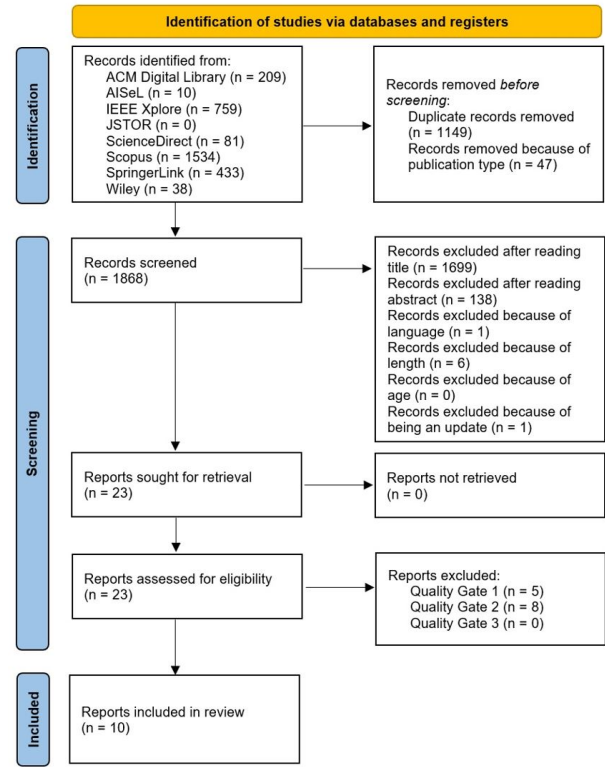


Fig. 1: Overview of the search process

agreement rate of 85,8 percent (145 of 169 records) between the authors. It was stipulated so if the agreement could not be reached, an arbitrator would be invited to the research. Yet, this was not necessary. In total of the 169 papers, 138 were removed and 31 were included in the next phase.

While conducting the first steps of our filter process, we encountered several hurdles that shall be highlighted to ensure transparency, especially since they can slightly affect the number of remaining entries after those initial phases. However, the final set of literature was not impacted and, therefore, those factors did not pose a threat to the studies validity. For once, since not all entries of the combined literature list specified a digital object identifier (DOI), the duplicate removal had to be conducted based on the publication title. Yet, in some rare cases, there were duplicates for which the spelling of the title was slightly altered (e.g., the two parts of a title were in one search engine separated by a hyphen and in another by a double colon), and which were, therefore, not detected in the initial duplicate removal phase. Instead, they were only identified during the scanning of the title.

Furthermore, in SpringerLink, conference papers are classified as book chapter, since conference proceedings are published as books. This makes them indistinguishable from real book chapters, when only looking at the metadata. Book chapters are, however, not part of the search's scope. Consequently, the removal of book chapters for SpringerLink could only be processed when inspecting the respective publications. To slightly reduce the effort, it was decided to only do this for

those publications that passed the filtering by title.

From this step, the papers that were not written in English (despite the abstract being in English) or for which an updated version exists were filtered out. The same would have applied to papers published before the year 2012, because in the previous years, the concept of microservices as the focus of this publication was not yet present. However, while there were three corresponding records in the initially obtained set of literature (from the years 2003, 2007, and 2010), those were already filtered out for other reasons by this stage. Along the lines, publications that had a length of less than six pages, without counting the reference section and acknowledgments, were also removed, since those can't provide the desired degree of comprehensiveness.

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 criteria was divided into three differently focused phases, with each of them requiring the passing of a quality gate as portrayed in I.

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. Further, the objectives, justification, and aim of the study shall be clearly communicated. Finally, also the context of the conducted research needed to be sufficiently described.

The first and second author 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.

TABLE I: 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	3.1 Credibility 3.2 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

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.

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. [15]. 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.

IV. RESULTS

In this section, we present with three integral elements: 1) the requirements that has emerged from the first SLR, 2)

A. Requirements Specification

The results of our data synthesis emerged a few themes in regards to BD requirements. While we could find BD major building blocks and requirements from the body of evidence, our SLR did not include categorization and representation of these requirements. To this end, we performed a lightweight literature review in the body of evidence to find a rigorous approach for to categorized and represent BD requirements.

Precursor to theorizing about the potential of microservices patterns for big data systems, we needed to define what are the requirements of these systems. System and software requirements come in different flavours and can range from a sketch on a napkin to formal (mathematical) specifications. Therefore, we first needed to identify what kind of requirements is the most suitable for the purposes of this study. To answer this question, we first explored the body of evidence to understand the current classification of software requirements.

There's been various attempts to defining and classifying software and system requirements. For instance, Sommerville ([23]) classified requirements into three levels of abstraction that are namely 1) user requirements, 2) system requirements and 3) design specifications. The author then mapped these requirements against user acceptance testing, integration testing and unit testing. While this could satisfy the requirements of this study, we opted for a more general framework provided by Laplante ([24]). In Laplante's approach, requirements are categorized into three categories of 1) functional requirements, 2) non-functional requirements, and 3) domain requirements.

Our objective is to define the high-level requirements of big data 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. Therefore, our primary focus is one functional and domain requirements and secondly on non-functional requirements.

After clarifying the type of requirements, we then explored the body of evidence to realize the general requirements of big data systems. By the result of this, we realized that the most discussed characteristics of big data systems are the popular 5Vs which are velocity, veracity, volume, Variety and Value ([25], [26], [27], [28], [29], [30]). Many researchers such as Nadal et al. ([31]) and Klein et al. (klein2016reference) have underpinned their reference architecture on these characteristics and requirements that goes with them.

On the other hand, in an extensive effort, NIST Big Data Public Working Group embarked on a large scale study to extract requirements from variety of application domains such as Healthcare, Life Sciences, Commercial, Energy, Government, and Defense. The result of this study was the formation of general requirements under seven categories. In another effort by Volk et al. ([32]), nine use cases for big data projects are identified by collecting theories and use cases from the literature and categorizing them using a hierarchical clustering algorithm. Bashari et al. ([33]) focused on the security and privacy requirements of big data systems, Yu et al. presented the modern components of big data systems ([34]), Eridaputra et al. ([35]) created a generic model for big data requirements using goal oriented approaches, and Al-jaroodi et al. ([36]) investigated general requirements to support big data software development.

We've also studied the reference architectures developed for big data systems to understand general requirements. In one study, Ataei et al. ([4]) assessed the body of evidence and presented with a comprehensive list of big data reference

architectures. This study helped us realize the spectrum of big data reference architectures, how they are designed and the general set of requirements.

By analyzing these studies and by evaluating the design and requirement engineering required for big data reference architectures, we created a set of high-level requirements based on big data characteristics. We have then 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's a well established method in the industry and academia ([37]).

Our approach follows the guidelines explained in ISO/IEC/IEEE standard 29148 ([38]) for representing functional requirements. Our requirement representation is organized in system modes, that is we explain the major components of the system and then describe the requirements. This approach is inspired by the requirement specification expressed for NASA WIRE (wide-field infrared explorer) system explained in [24]. We also taken inspiration from Software Engineering Body of Knowledge ([39]).

Taking all into consideration, we categorized our requirements based on the major characteristics of big data, that is value, variety, velocity, veracity, and volume ([40]), plus security and privacy ([33]). These requirements are described in table II.

B. Microservice Patterns

As a result of this SLR, 50 microservice patterns have been found. These patterns are then classified based on their function and the problem they solve. Each classification and its reasoning is depicted in table III.

V. 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 microservices design patterns for big data systems. The patterns gleaned, are established theories that are derived from actual problems in microservices systems in practice, thus we do not aim to re-validate them in this study. Moreover, we do not aim to validate the theories proposed in this study through an empirical 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, big data. Based on this, we map big data system requirements against a pattern and provide with reasoning on why such pattern might work for big data systems.

Big data systems and microservices architecture are both inherently distributed. While majority of current big data applications are designed underlying a monolithic data pipeline

Volume	Vol-1) System needs to support asynchronous, streaming, and batch processing to collect data from centralized, distributed, and cloud data sources, and sensors, instrument and other IOT devices 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 and computing clusters 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 graph, web, text, document, timed, spatial, multimedia, simulation, instrumental, and geo-spatial 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 be 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 such as descriptive analytics, predictive analytics, reporting and visualizations. 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. Ver-3) System needs to support data validation in two ways: automatic and human annotated. Ver-4) System should be able to handle data loss or corruption.

TABLE II: Big data system requirements

architecture, here, we propose microservices architecture for a domain-driven and decentralized big data architecture. We support our arguments by the means of modeling. We use Archimate ([41]) as recommend in ISO/IEC/IEEE 42010 ([42]).

We posit that a pattern alone would not be significantly useful to a data engineering or a data architect, and propose that collection of a pattern 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 big data systems.

A. Volume

There has been two requirements associated to the Volume aspect of big data systems which are about handling various data types (Vol-1) and providing with a scalable storage (Vol-2).

For Vol-1 and Vol-2 we suggest the following patterns to be effective; 1) External Configuration Store, 2) API gateway, 3) Gateway offloading

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 achieve this process. One challenge in this process 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) what are the typical use cases for the data being ingested ? is the big data system ingesting data reliably ? what is the next data destination ? How frequently should data be ingested ? In what volume the data typically arrives? Does streaming data need to be transformed before reaching the destination ?

Given the challenges and particularities of data types, different nodes maybe spawned to handle the volume of data as witnessed in big data reference architectures studied by Ataei et al ([4]). Another popular approach is the segregation of concerns by separating batch and streaming processing nodes. Given the requirement of horizontal scaling for big data systems, it is safe to assume that there is usually more

Category	Pattern
Data Management	Database per Service, Shared Database, Event Sourcing, Command and Query Responsibility Segregation
Platform and Infrastructure	Multiple service instances per host, External configuration store, Sidecar, Static content hosting, Computer resource consolidation
Communicational	API gateway, Anti-corruption layer, Self Registration, Service Discovery, Competing consumers, Pipes and filters, Priority queue, Ambassador, Gateway aggregate, Gateway offloading, Aggregator, Backend for Frontend, API Composition, Saga transaction management, Gateway routing, Leader election
Fault Tolerance	Circuit breaker, Bulkhead pattern
Observability	Log Aggregation Pattern

TABLE III: Microservices categorization

then one node associated to the data being ingested. This can be problematic as different nodes will need to account for security, privacy and overall regulations of the context, alongside, the software engineering demand that each may have.

This means that each node needs to reimplement the same interface for the aforementioned cross-cutting concerns, which makes scalability and maintainability of the big data system a daunting task. This also introduces unnecessary repetition of codes. To solve this problem, we explore the concept of gateway offloading and API gateway patterns. By offloading cross-cutting concerns that are shared across nodes to a single architectural construct, the API gateway in this case, not only we will achieve a separation of concerns and a good level of usability, but we increase security and performance, by processing and filtering incoming data through a well specified ingress.

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 service tries to communicate with a service that is down. Given that, the 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 big data system can be compromised. This scenario remains the same, and can be even worst if the company decides to have an on-premise data center.

Additionally, the gateway can increase the system reliability and availability by doing a constant health check on services, and distribute traffic based on healthy nodes. There is also an array of other benefits such as having a weighted distribution, and creating a special cache mechanism through specific HTTP headers. This also means that if the gateway is down, service nodes won't introduce bad data or state into the overall system. We have portrayed a very simplistic representation of this pattern in fig 2.

2) *External Configuration Store*: As discussed earlier, big data systems are made up of various nodes in order to achieve horizontal scalability. While these systems are logically separated to their own service, they 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 configuration to achieve their functionality.

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, also known as the 'externalized configuration 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 solve the challenges of handling large number of nodes in big data systems and provide with a scalable solution for handling configurations. This pattern is portrayed in fig 2.

B. Velocity

Velocity is perhaps one of the most challenging aspects of the big data 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 big data 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.

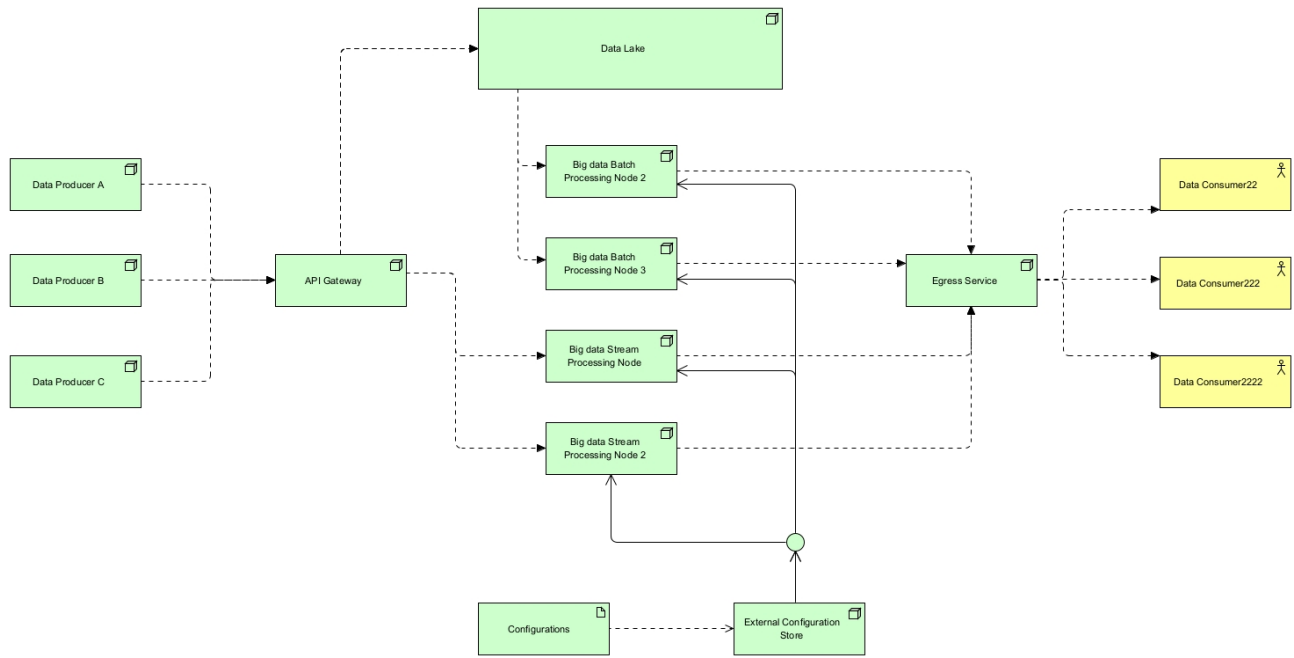


Fig. 2: Design patterns for volume requirement

1) *Competing Consumers*: Big data 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 big data 'low-achievers' is the 'slow processing of large amounts of data' ([2]). 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 big data systems with so many moving parts, is a challenging task, but there are microservices patterns that can be tailored to help with some of these challenges. Given the very contrived scenario of a big data 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 later 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 et al ([12]), 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 com-

munication. This implies that node A can send an event to node B in a 'dispatch and forget' fashion 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 (could be ingress or data producers). This pattern will help alleviate challenges in regards to Vel-2, Vel-3 and Vel-5.

2) *Circuit Breaker*: On the other hand, given the large number of nodes one can assume for any big data 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 event-driven asynchronous fault tolerance communication mechanisms among big data 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 Val-1, Val-4, and Val-5.

C. Log Aggregator

Given that big data 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 communicates to and log their processes. 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 issue, which in turn increases reliability and availability which in turn addresses the velocity requirement of big data systems. While this design pattern does not directly affect any system requirements, it indirectly affects all of them. A simplistic reference architecture of these patterns have been portrayed in fig 3

D. Variety

Variety, being another important aspect of big data, implies the range of different data types and the challenges of handling these data. As big data system grows, newer data structures emerge, and an effective big data 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.

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 big data 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 pipeline or introduce a privacy threat. How can data engineers handle this scenario effectively?

To address this problem, and to address big data requirements Var-1, Var-3, and Var-4, 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.

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 fig 2.

E. Value

Value is the nucleus of any big data 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, 3) Gateway Offloading

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 big data 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 serving layer can end up with an overly complex logic that does too much.

Read and write workloads are really different, and this is something a data engineer should consider from the initial data modeling, to data storage, retrieval and potential serialization (JSON to Parquet). 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). 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 this 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 and update its values. This allows for elastic scaling of the read nodes, and increased query performance, especially in big data systems that have got egress services sitting on an edge. Therefore, this pattern can potentially address the requirement Val-1, and Val-3.

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 is placed between the serving layer and data consumers, isolating different systems and

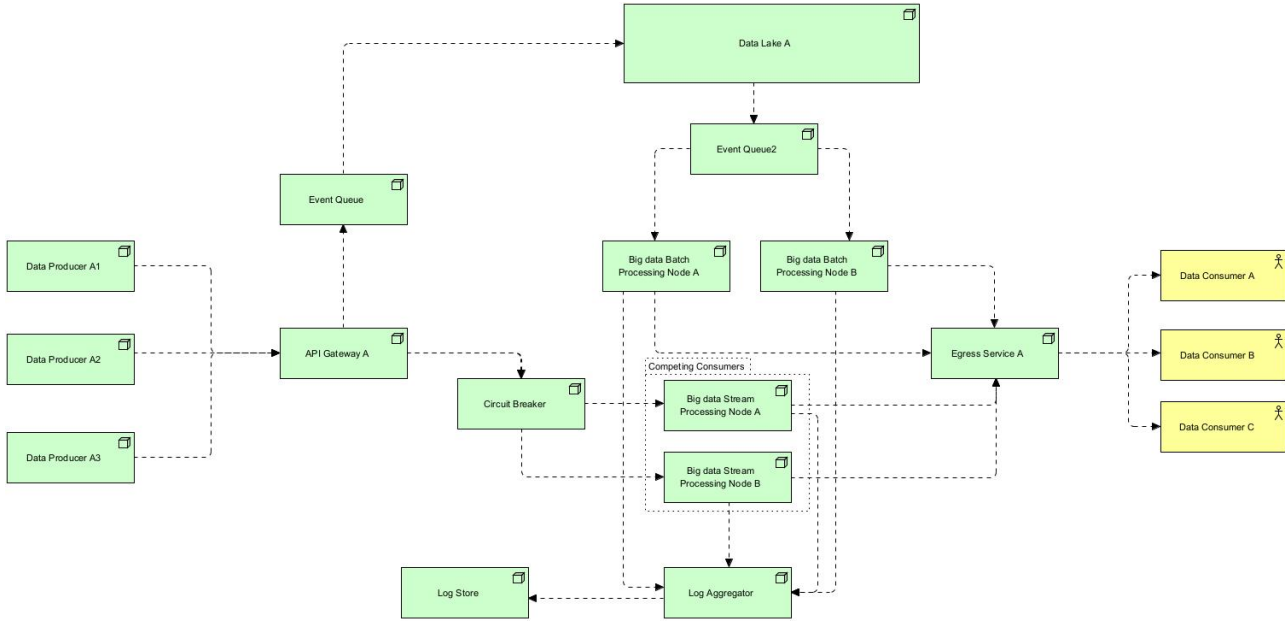


Fig. 3: Design patterns for velocity requirements

translating communication. This eliminates all the complexity and coupling that could have been otherwise introduced to the serving layer. This also allows for nodes to follow the 'single responsibility' pattern ([43]). 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 and CQRS in fig 4.

3) *Gateway Offloading*: We have previously discussed this pattern, but in this section we aim to relate more to the value requirements of BD. Given that the various nodes in the system may require services such as authentication, authorization, monitoring, logging, and throttling, it would become really difficult to address the value aspect of big data in a timely manner. To address these issues and to achieve Val-3 and Val-4, we recommend the gateway offloading pattern. Employing this pattern abstracts out cross-cutting services from each node and creates a unified interface that each node can utilize. This simplifies the development of new services for handling new data formats, and free data engineers from implementing features that requires special knowledge such as security and privacy.

F. Security and Privacy

Security and privacy should be on top of mind for any big data 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)

1) *BFF*: 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 big data 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 batch or stream processing nodes. However as data consumers grow and more data gets in, maintaining the privacy rules and applying them correctly to the dataset becomes more difficult. This becomes a perfect ground for mistake, and can potentially introduce legal issues to the company.

On approach to this problem can be the BFF pattern. By creating backends (services) for frontends (data producers), we can logically segregate system's ingress for data that requires different level of privacy and security. This logical separation can include other factors such as 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, 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 protocol 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

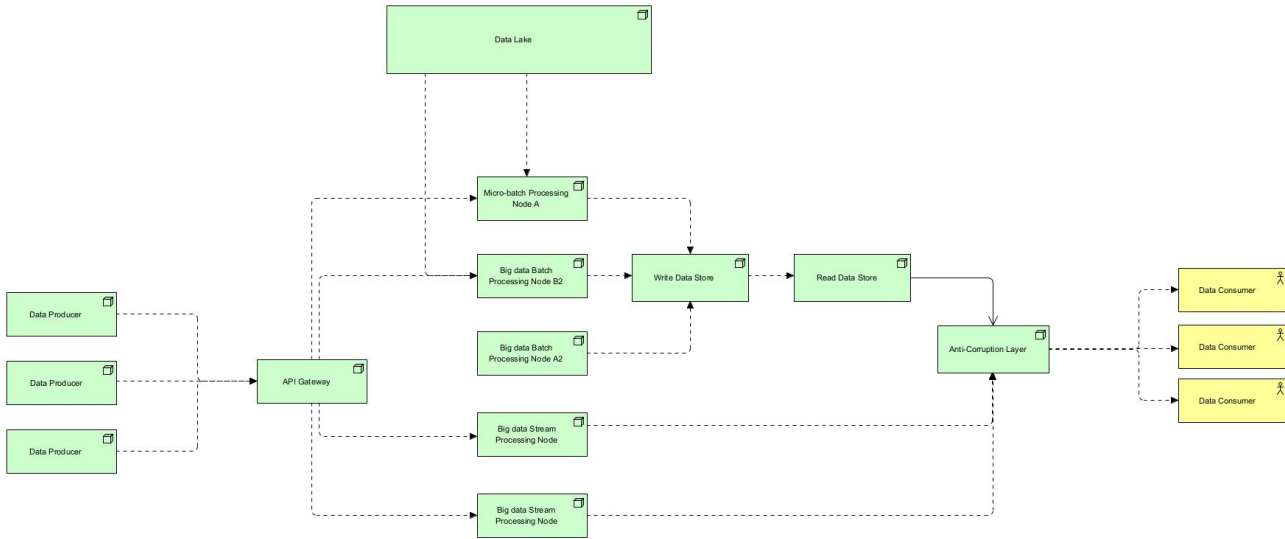


Fig. 4: Design patterns for value requirement

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.

G. Veracity

Next to value, veracity is an integral component of any effective big data 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 ([44]), 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 big data systems.

To this, we propose the following patterns for addressing requirements Ver-1, and Ver-4; 1) Pipes and Filters, 2) Circuit breaker

H. 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 resembles to the well-known pattern of 'single responsibility' ([43]). 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.

I. Circuit breaker

In an inherently distributed environment like big data, 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 big data system, when the service of interest is recovering from an incident. This can indirectly help with Ver-4.

VI. 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 ([45]). For the purposes of this study we chose semi-structured interviews (SSIs), following the guidelines of Adams ([46]) and Kallio et al. ([47]).

A. Methodology

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. As we examine an uncharted territory with a lot of potential, we posit that these interviews can post useful leads which we can pursue to further improve the theories of this study. We've formulated our SSI guide based on our research objective to achieve the richest possible data. Our guide is flexible, to increase our opportunity to explore new ideas, and allow for participant-orientation. Nevertheless, we do have some close-ended questions at the start of our interview which is a good starter, and also helps us with some statistics.

Our questions are categorized into main themes and follow-up questions, with main themes being progressing and logical, as recommended by Kallio et al. ([47]). Follow-up questions were utilized to direct the dialogue towards the subject of our study, and make things easier for candidate to understand. Some of these follow-up questions were improvised, as we did not aim to rigidly control the flow of the interview. After this, to ensure the rigour and relevance of the interview guide, we've conducted a pilot test. This step was necessary to make informed adjustments to the guide, and to improve quality of data collection.

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. We aimed to assume the role of the interviewee and gain insight into the limitations of our guide. This approach helped us capture some issues with the questions, and remove some questions that may be deemed eccentric. From there on, we presented the results as theories with some statistics attained from close-ended questions.

B. Results

From the results of these interviews, we gathered a lot of insights and probed deeper some of our architectural decisions. Almost 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 and held titles such as 'lead development architect' and 'solution architect'.

While some of our interviewee had more experience with big data and AI, some others were well-versed in microservices architecture. We first asked interviewees about their depth of understanding with microservices and big data, and then asked them if patterns discussed for each characteristics makes sense. We went through each pattern using our Archimate model, and explained and discussed why we've chosen it. We asked every interviewee if they can think of a pattern that we failed to consider. Our interview guide is available at [48].

We realized that API gateway and gateway offloading pattern is easily accepted as an effective pattern, while CQRS needed more reasoning and explanation. One interviewee had a concern about the write and read services being the single point of failure in the CQRS pattern and how those services can scale. Another interviewee wanted to know if we looked into event sourcing and if that could be applied. We've also received comments on the usage of priority queue for sensitive stream processing requirements. The most experienced interviewee (14 years) have 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. Almost 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. Some of our interviewees connected some of the patterns discussed to their own context of practice and helped us realized further improvements. By the result of this interview we realized that we have missed an architectural construct while discussing velocity requirements, which was a 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.

VII. DISCUSSION

This study was our attempt to explore two major areas, big data systems and microservices architecture. By adopting a rigorous methodology, we aimed at applying microservices patterns to big data systems. As a result of this we created a number of theories that can help shed lights on adoption of these patterns to big data systems in practice.

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) microservices pattern provide with a great potential for resolving some of the big data 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 big data 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 microservices and event driven architectures, we could not find many well-established patterns in the data engineering space. A part of it is due to adoption of microservices architecture in the industry, and in specifically by IT giants ([49]). While big data is not been successfully adopted as much. A survey by MIT technology review insights presented that only 13% of companies excel at delivering their big data strategy. Based on this, we think that data architecture remains a significant challenge and requires more attention from both academia and industry.

VIII. CONCLUSION

With all the undeniable benefits of big data, the success rate of big data projects is still rare. One of the core challenges of adopting big data lies in data architecture and data engineering. While software engineers has matured to go through cadence of architectures from monolithic to service-oriented and to microservices and event-driven architectures, data engineering and big data architectures don't seem to benefit a lot from these advancements.

The aim of this study was to explore the relationship and application of microservices architecture to big data systems through two distinct SLR. The results derived from these SLRs presented us with interesting data on the potential of microservices patterns for big data systems. Given the distributed nature of big data systems, microservices 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 microservices and event-driven architectures in relation to big data systems from both academia and industry.

REFERENCES

- [1] N. Partners, "Big data and ai executive survey 2021," 2022. [Online]. Available: <https://www.newvantage.com/thoughtleadership>
- [2] M. technology review insights in partnership with Databricks, "Building a high-performance data organization," 2021. [Online]. Available: <https://databricks.com/p/whitepaper/mit-technology-review-insights-report>
- [3] M. Volk, D. Staegemann, M. Pohl, and K. Turowski, "Challenging big data engineering: Positioning of current and future development," in *Proceedings of the 4th International Conference on Internet of Things, Big Data and Security*. SCITEPRESS - Science and Technology Publications, 2019, pp. 351–358.
- [4] P. Ataei and A. T. Litchfield, "Big data reference architectures, a systematic literature review," 2020.
- [5] A. Freymann, F. Maier, K. Schaefer, and T. Böhnelt, "Tackling the six fundamental challenges of big data in research projects by utilizing a scalable and modular architecture," in *Proceedings of the 5th International Conference on Internet of Things, Big Data and Security*. SCITEPRESS - Science and Technology Publications, 2020, pp. 249–256.
- [6] I. Nadareishvili, R. Mitra, M. McLarty, and M. Amundsen, *Microservice architecture: Aligning principles, practices, and culture*, first edition ed. Beijing and Boston and Farnham and Sebastopol and Tokyo: O'Reilly, 2016.
- [7] P. Ataei and A. Litchfield, "Big data reference architectures, a systematic literature review," in *Australasian Conference on Information Systems (ACIS) 2020*. AIS, 2020.
- [8] R. Laigner, M. Kalinowski, P. Diniz, L. Barros, C. Cassino, M. Lemos, D. Arruda, S. Lifschitz, and Y. Zhou, "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, 2020, pp. 213–220.
- [9] S. Zhelev and A. Rozeva, "Using microservices and event driven architecture for big data stream processing," in *AIP Conference Proceedings*, vol. 2172, no. 1. AIP Publishing LLC, 2019, p. 090010.
- [10] D. Staegemann, M. Volk, A. Shakir, E. Lautenschläger, and K. Turowski, "Examining the interplay between big data and microservices—a bibliometric review," *Complex Systems Informatics and Modeling Quarterly*, no. 27, pp. 87–118, 2021.
- [11] A. Maamouri, L. Sfaxi, and R. Robbana, "Phi: A generic microservices-based big data architecture," in *European, Mediterranean, and Middle Eastern Conference on Information Systems*. Springer, 2021, pp. 3–16.
- [12] P. Ataei and A. Litchfield, "Neomycelia: A software reference architecture for big data systems," in *2021 28th Asia-Pacific Software Engineering Conference (APSEC)*. IEEE, 2021, pp. 452–462.
- [13] B. A. Kitchenham, T. Dyba, and M. Jorgensen, "Evidence-based software engineering," in *Proceedings of the 26th International Conference on Software Engineering*. IEEE Comput. Soc, 2004, pp. 273–281.
- [14] M. J. Page, 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, "Prisma 2020 explanation and elaboration: updated guidance and exemplars for reporting systematic reviews," *BMJ (Clinical research ed.)*, vol. 372, p. n160, 2021.
- [15] D. S. Cruzes and T. Dyba, "Recommended steps for thematic synthesis in software engineering," in *2011 International Symposium on Empirical Software Engineering and Measurement*. IEEE, 2011, pp. 275–284.
- [16] F. Buschmann, R. Meunier, H. Rohnert, P. Sommerlad, and M. Stal, *Pattern-Oriented Software Architecture: A System of Patterns, Volume 1*. John Wiley & sons, 2008, vol. 1.
- [17] C. Richardson, "A pattern language for microservices," 2022. [Online]. Available: <https://microservices.io/patterns/index.html>
- [18] M. L. Rethlefsen, S. Kirtley, S. Waffenschmidt, A. P. Ayala, D. Moher, M. J. Page, and J. B. Koffel, "Prisma-s: an extension to the prisma statement for reporting literature searches in systematic reviews," *Systematic reviews*, vol. 10, no. 1, pp. 1–19, 2021.
- [19] P. Ataei and D. Staegemann, "Systematic literature review search terms table for the paper titled: Application of microservices patterns to big data systems," 2022. [Online]. Available: <https://anonymous.4open.science/r/SLR-Search-Terms-3147/>
- [20] A. C. Tricco, E. Lillie, W. Zarin, K. K. O'Brien, H. Colquhoun, D. Levac, D. Moher, M. D. Peters, T. Horsley, L. Weeks *et al.*, "Prisma extension for scoping reviews (prisma-scr): checklist and explanation," *Annals of internal medicine*, vol. 169, no. 7, pp. 467–473, 2018.
- [21] [Online]. Available: <https://casp-uk.net/casp-tools-checklists/>
- [22] B. A. Kitchenham, S. L. Pfleeger, L. M. Pickard, P. W. Jones, D. C. Hoaglin, K. El Emam, and J. Rosenberg, "Preliminary guidelines for empirical research in software engineering," *IEEE Transactions on software engineering*, vol. 28, no. 8, pp. 721–734, 2002.
- [23] I. Sommerville, *Software Engineering, 9/E*. Pearson Education India, 2011.
- [24] P. A. Laplante, *Requirements engineering for software and systems*. Auerbach Publications, 2017.
- [25] Y. Demchenko, C. De Laat, and P. Membrey, "Defining architecture components of the big data ecosystem," in *2014 International Confer-*

ence on Collaboration Technologies and Systems (CTS). IEEE, 2014, Conference Proceedings, pp. 104–112.

- [26] J. Bughin, “Big data, big bang?” *Journal of Big Data*, vol. 3, no. 1, p. 2, 2016.
- [27] M. Bahrami and M. Singhal, *The role of cloud computing architecture in big data*. Springer, 2015, pp. 275–295.
- [28] B. B. Rad and P. Ataei, “The big data ecosystem and its environs,” *International Journal of Computer Science and Network Security (IJCSNS)*, vol. 17, no. 3, p. 38, 2017.
- [29] N. Marz and J. Warren, *Big Data: Principles and best practices of scalable real-time data systems*. New York; Manning Publications Co., 2015.
- [30] H.-M. Chen, R. Kazman, and S. Haziye, “Agile big data analytics development: An architecture-centric approach,” in *2016 49th Hawaii International Conference on System Sciences (HICSS)*. IEEE, 2016, Conference Proceedings, pp. 5378–5387.
- [31] S. Nadal, V. Herrero, O. Romero, A. Abelló, X. Franch, S. Vansumeren, and D. Valerio, “A software reference architecture for semantic-aware big data systems,” *Information and software technology*, vol. 90, pp. 75–92, 2017.
- [32] M. Volk, D. Staegemann, I. Trifonova, S. Bosse, and K. Turowski, “Identifying similarities of big data projects—a use case driven approach,” *IEEE Access*, vol. 8, pp. 186 599–186 619, 2020.
- [33] B. Bashari Rad, N. Akbarzadeh, P. Ataei, and Y. Khakbiz, “Security and privacy challenges in big data era,” *International Journal of Control Theory and Applications*, vol. 9, no. 43, pp. 437–448, 2016.
- [34] J.-H. Yu and Z.-M. Zhou, “Components and development in big data system: A survey,” *Journal of Electronic Science and Technology*, vol. 17, no. 1, pp. 51–72, 2019.
- [35] H. Eridaputra, B. Hendradjaya, and W. D. Sunindyo, “Modeling the requirements for big data application using goal oriented approach,” in *2014 international conference on data and software engineering (ICODSE)*. IEEE, 2014, pp. 1–6.
- [36] J. Al-Jaroodi and N. Mohamed, “Characteristics and requirements of big data analytics applications,” in *2016 IEEE 2nd International Conference on Collaboration and Internet Computing (CIC)*. IEEE, 2016, pp. 426–432.
- [37] M. Kassab, C. Neill, and P. Laplante, “State of practice in requirements engineering: contemporary data,” *Innovations in Systems and Software Engineering*, vol. 10, no. 4, pp. 235–241, 2014.
- [38] I. 29148:2018, “Iso/iec 29148:2018,” 2018. [Online]. Available: <https://www.iso.org/standard/72089.html>
- [39] A. Abran, J. W. Moore, P. Bourque, R. Dupuis, and L. Tripp, “Software engineering body of knowledge,” *IEEE Computer Society, Angela Burgess*, p. 25, 2004.
- [40] B. B. Rada, P. Ataeib, Y. Khakbizc, and N. Akbarzadehd, “The hype of emerging technologies: Big data as a service,” 2017.
- [41] M. Lankhorst, “A language for enterprise modelling,” in *Enterprise Architecture at Work*. Springer, 2013, pp. 75–114.
- [42] M. Chaabane, I. Bouassida, and M. Jmaiel, “System of systems software architecture description using the iso/iec/ieee 42010 standard,” in *Proceedings of the Symposium on Applied Computing*, Conference Proceedings, pp. 1793–1798.
- [43] E. Gamma, R. Helm, R. Johnson, R. E. Johnson, J. Vlissides *et al.*, *Design patterns: elements of reusable object-oriented software*. Pearson Deutschland GmbH, 1995.
- [44] E. Eryurek, U. Gilad, V. Lakshmanan, A. Kibunguchy-Grant, and J. Ashdown, *Data Governance: The Definitive Guide*. ” O’Reilly Media, Inc.”, 2021.
- [45] R. J. Wieringa, *Design science methodology for information systems and software engineering*. Springer, 2014.
- [46] W. C. Adams, *Conducting Semi-Structured Interviews*. John Wiley & Sons, Ltd, 2015, ch. 19, pp. 492–505. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781119171386.ch19>
- [47] H. Kallio, A.-M. Pietilä, M. Johnson, and M. Kangasniemi, “Systematic methodological review: developing a framework for a qualitative semi-structured interview guide,” *Journal of advanced nursing*, vol. 72, no. 12, pp. 2954–2965, 2016.
- [48] P. Ataei and D. Staegemann, “Interview guide for the paper: Application of microservices patterns to big data systems,” 2022. [Online]. Available: <https://anonymous.4open.science/r/SSI-Repo-F90E/SSI.pdf>
- [49] “Organizations’ adoption level of microservices worldwide in 2021.” [Online]. Available:

<https://www.statista.com/statistics/1233937/microservices-adoption-level-organization/>