

# Conference Paper Title\*

\*Note: Sub-titles are not captured in Xplore and should not be used

1<sup>st</sup> Given Name Surname  
*dept. name of organization (of Aff.)*  
*name of organization (of Aff.)*  
City, Country  
email address or ORCID

2<sup>nd</sup> Given Name Surname  
*dept. name of organization (of Aff.)*  
*name of organization (of Aff.)*  
City, Country  
email address or ORCID

3<sup>rd</sup> Given Name Surname  
*dept. name of organization (of Aff.)*  
*name of organization (of Aff.)*  
City, Country  
email address or ORCID

4<sup>th</sup> Given Name Surname  
*dept. name of organization (of Aff.)*  
*name of organization (of Aff.)*  
City, Country  
email address or ORCID

5<sup>th</sup> Given Name Surname  
*dept. name of organization (of Aff.)*  
*name of organization (of Aff.)*  
City, Country  
email address or ORCID

6<sup>th</sup> Given Name Surname  
*dept. name of organization (of Aff.)*  
*name of organization (of Aff.)*  
City, Country  
email address or ORCID

**Abstract**—This document is a model and instructions for L<sup>A</sup>T<sub>E</sub>X. This and the IEEEtran.cls file define the components of your paper [title, text, heads, etc.]. \*CRITICAL: Do Not Use Symbols, Special Characters, Footnotes, or Math in Paper Title or Abstract.

**Index Terms**—component, formatting, style, styling, insert

## I. REQUIREMENTS SPECIFICATION

Precursor to theorizing about the potential of microservices patterns for big data systems, we need to define what we mean by big data systems and what are the requirements of these systems. System and software requirements come in different flavour and can range from a sketch on a napkin to formal (mathematical) specifications. Therefore, we first need 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 systems requirements. For instance, Sommerville ([?]) 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 ([?]). 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 seek to explore 'non-functional' requirements. Non-functional requirements are emerged from the particularities of an environment, such as a banking sector and do not correlate to our study. Therefore, the type of

requirements we are looking for is functional and domain requirements.

After clarifying the type of requirements, we then explored the body of evidence to realize the general requirements of big data systems. Indeed, the most discussed characteristics of big data systems are the popular 5Vs which are velocity, veracity, volume, Variety and Value ([?], [?], [?], [?], [?]). Many researchers such as Nadal et al. ([?]) have underpinned their artifact development on these characteristics and requirements that emerge from them.

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 and 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. ([?]), 9 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. ([?]) focused on the security and privacy requirements of big data systems, Yu et al. presented the modern components of big data systems [?], Eridaputra et al. ([?]) created a generic model for big data requirements using goal oriented approaches, and Al-jaroodi et al. ([?]) 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. ([?]) assessed the body of evidence and presented with a comprehensive list of big data reference architectures. This study helped us realized 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

Identify applicable funding agency here. If none, delete this.

a rigorous approach to present these requirements. There are numerous approaches used for 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 ([?]).

Our approach follows the guidelines explained in ISO/IEC/IEEE standard 29148 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 [?]. We also taken inspiration from Software Engineering Body of Knowledge Version ([?]).

Taking all into consideration, we categorized our requirements based on the major characteristics of big data, that is value, variety, velocity, veracity, and volume ([?]), plus . These requirements are as followings:

## II. MICROSERVICE PATTERNS

As a result of this SLR, 50 microservice patterns have bene found. These patterns are then classified based on their function and the problem they solve. Each classifications and it's reasoning is depicted in table ??.

- 1) Database per service ✓
- 2) Shared database ✓
- 3) Event sourcing ✓
- 4) Multiple service instances per host ✓
- 5) API gateway ✓
- 6) Self registration ✓
- 7) Service discovery ✓
- 8) Circuit breaker ✓
- 9) Bulkhead pattern ✓
- 10) Command and query responsibility segregation ✓
- 11) Competing consumers ✓
- 12) Pipes and filters ✓
- 13) Strangler ✓
- 14) Anti-corruption layer ✓
- 15) External configuration store ✓
- 16) Priority queue ✓
- 17) Log Aggregation ✓
- 18) Ambassador ✓
- 19) Sidecar ✓
- 20) Gateway aggregate ✓
- 21) Gateway offloading ✓
- 22) Aggregator ✓
- 23) Backend for Frontend ✓
- 24) API Composition ✓
- 25) Saga transaction management ✓
- 26) Gateway routing (duplicate) ✓
- 27) Static content hosting ✓
- 28) Computer resource consolidation ✓
- 29) Leader election ✓

## III. 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.

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 the application handling various data types (Vol-1) and the application providing with a scalable storage (Vol-2).

For Vol-1 we suggest the following patterns to be effective:

- 1) Database per service
- 2) Event sourcing
- 3) Command and query responsibility segregation
- 4) External configuration
- 5) API gateway
- 6) Service discovery
- 7) Self registration
- 8) Priority queue
- 9) Gateway offloading
- 10) Gateway aggregate
- 11) Leader election
- 12) Log aggregation pattern

Big data systems and microservices architecture are both inherently distributed. While majority of current big data applications are designed underlying a monolithic data pipeline 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 ([?]) as recommend in ISO/IEC/IEEE 42010 ([?]).

We do not suggest any pattern to be very effective, and propose that a collection of these patterns can provide the effectiveness of microservices in to big data systems. For this purpose

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

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

Requirement	Patterns	Reasoning
Vol-1	1) Database per Service 2) Event Sourcing 3) Command and Query Responsibility Segregation 4) External Configuration Store 5) API gateway 6) Anti-Corruption Layer 7) Service Discovery 8) Self Registration 9) Priority Queue 10) Gateway Offloading 11) Gateway Aggregate 12) Leader Election 13) Log Aggregation Pattern	Reasoning
Vol-2	1) Database per Service 2) Command and Query Responsibility Segregation	Reasoning
Vel-1	1) API Gateway 2) Service Discovery 3) Pipes and Filters 4) Gateway Routing 5) Leader Election 6) Circuit Breaker 7) Log Aggregation	Reasoning
Vel-2	1) Command and Query Responsibility Segregation 2) API gateway 3) Competing consumers 4) Gateway aggregate 5) Gateway Offloading 6) Leader Election 7) Gateway routing	Reasoning
Vel-3	1) Command and Query Responsibility Segregation 2) API gateway 3) Competing consumers 4) Gateway aggregate 5) Gateway Offloading 6) Leader Election 7) Gateway routing	Reasoning
Vel-3	1) API composition 2) API gateway 4) Gateway aggregate 5) Gateway Offloading 7) Gateway routing	Reasoning
Vel-4	1) API composition 2) API gateway 4) Gateway aggregate 5) Gateway Offloading 7) Gateway routing 8) Event Sourcing 9) Command and Query Responsibility Segregation	Reasoning
Vel-5	1) Leader Election 2) Log Aggregation Pattern	Reasoning
Var-1	None	Reasoning
Var-2	1) Database per Service 2) API Gateway 3) Gateway Routing	Reasoning
Var-3	None	Reasoning
Var-4	1) API Gateway 2) Gateway Routing 3) Gateway Offloading 4) Gateway Aggregate	Reasoning

Val-1	<ul style="list-style-type: none"> <li>1) Event Sourcing</li> <li>2) Command and Query Responsibility Segregation</li> <li>3) Priority Queue</li> <li>4) Leader Election</li> <li>5) Bulkhead Pattern</li> </ul>	Reasoning
Val-2	<ul style="list-style-type: none"> <li>1) Event Sourcing</li> <li>2) Command and Query Responsibility Segregation</li> <li>3) API gateway, Gateway aggregate</li> <li>4) Gateway offloading</li> <li>5) Gateway routing</li> <li>6) Priority queue</li> </ul>	Reasoning
Val-3	<ul style="list-style-type: none"> <li>1) API gateway</li> <li>2) Anti-corruption layer</li> <li>3) Service Discovery</li> <li>4) Gateway aggregate</li> <li>5) Gateway routing</li> <li>6) Backend for Frontend</li> </ul>	Reasoning
Val-4	<ul style="list-style-type: none"> <li>1) Event Sourcing</li> <li>2) Command and Query Responsibility Segregation</li> <li>3) Backend for Frontend</li> </ul>	Reasoning
SaP-1	<ul style="list-style-type: none"> <li>1) External Configuration Store</li> <li>2) API Gateway</li> <li>3) Gateway Aggregate</li> <li>4) Backend for Frontend</li> </ul>	Reasoning
SaP-2	<ul style="list-style-type: none"> <li>1) External Configuration Store</li> <li>2) API Gateway</li> <li>3) Gateway Aggregate</li> <li>4) Backend for Frontend</li> </ul>	Reasoning
Ver-1	<ul style="list-style-type: none"> <li>1) Pipes and filters</li> </ul>	Reasoning
Ver-2	None	Reasoning
Ver-3	None	Reasoning
Ver-4	<ul style="list-style-type: none"> <li>1) Circuit Breaker</li> </ul>	Reasoning