



# A taxonomy of service identification approaches for legacy software systems modernization

Manel Abdellatif<sup>a,\*</sup>, Anas Shatnawi<sup>b</sup>, Hamed Mili<sup>b</sup>, Naouel Moha<sup>c</sup>, Ghizlane El Boussaidi<sup>c</sup>, Geoffrey Hecht<sup>b</sup>, Jean Privat<sup>b</sup>, Yann-Gaël Guéhéneuc<sup>d</sup>

<sup>a</sup> Polytechnique Montréal, Montreal, Quebec, Canada

<sup>b</sup> Université du Québec à Montréal, Montreal, Quebec, Canada

<sup>c</sup> École de Technologie Supérieure, Montréal, Québec, Canada

<sup>d</sup> Concordia University, Montréal, Québec, Canada

## ARTICLE INFO

### Article history:

Received 8 April 2019

Received in revised form 15 October 2020

Accepted 16 November 2020

Available online 20 November 2020

### Keywords:

Service identification

Microservices

Taxonomy

Legacy system

Migration

## ABSTRACT

The success of modernizing legacy software systems to Service-Oriented Architecture (SOA) depends on Service Identification Approaches (SIAs), which identify reusable functionalities that could become services. The literature describes several SIAs. However, the selection of an identification approach that is suitable for a practitioner is difficult because it depends on several factors, including the goal of modernization, the available legacy artifacts, the organization's development process, the desired output, and the usability of the approach. Accordingly, to select a suitable service identification approach, a practitioner must have a comprehensive view of existing techniques.

We report a systematic literature review (SLR) that covers 41 SIAs based on software-systems analyses. Based on this SLR, we create a taxonomy of SIAs and build a multi-layer classification of existing identification approaches. We start from a high-level classification based on the used inputs, the applied processes, the given outputs, and the usability of the SIAs. We then divide each category into a fine-grained taxonomy that helps practitioners in selecting a suitable approach for identifying services in legacy software systems. We build our SLR based on our experience with legacy software modernization, on discussions and experiences working with industrial partners, and analyses of existing SIAs. We validate the correctness and the coverage of our review with industrial experts who modernize(d) legacy software systems to SOA. The results show that our classification conforms to the industrial experts' experiences. We also show that most of the studied SIAs are still at their infancy. Finally, we identify the main challenges that SIAs need to address, to improve their quality.

© 2020 Elsevier Inc. All rights reserved.

## 1. Introduction

The maintenance and migration of legacy software systems are central IT activities in many organizations in which these systems are mission-critical. These systems embed hidden knowledge that is still of significant values. They cannot be removed or replaced because they execute effectively and accurately critical and complex business logic. Yet, legacy software systems suffer from several drawbacks including high maintenance costs, scalability and portability problems, and so forth (Lewis et al., 2005). Thus, these systems should be migrated to more flexible and modern

architectures to retain their business values while decreasing their maintenance costs.

The migration of legacy software systems to a Service-Oriented Architecture (SOA) is one avenue for the modernization of these systems. SOA allows developing complex and inter-organizational systems by integrating and composing services that are reusable, distributed, relatively independent, and often heterogeneous (Erl, 2008). Also over the past few years, increasing efforts have been made to migrate legacy systems to microservices, which are, in a SOA architecture, any services having a single responsibility, running in their own processes, and communicating with lightweight mechanisms (Newman, 2015). In the following, we use the term “service” to cover any form/granularity of services, including microservices.

The migration of legacy software systems to SOA is difficult because it depends on many factors, e.g., the choice of the migration process, the service-identification approach, the desired quality characteristics of the generated services, the implementation

\* Corresponding author.

E-mail addresses: [manel.abdellatif@polymtl.ca](mailto:manel.abdellatif@polymtl.ca) (M. Abdellatif), [shatnawi.anas@uqam.ca](mailto:shatnawi.anas@uqam.ca) (A. Shatnawi), [mili.hafedh@uqam.ca](mailto:mili.hafedh@uqam.ca) (H. Mili), [naouel.moha@etsmtl.ca](mailto:naouel.moha@etsmtl.ca) (N. Moha), [ghizlane.elboussaidi@etsmtl.ca](mailto:ghizlane.elboussaidi@etsmtl.ca) (G.E. Boussaidi), [hecht.geoffrey@uqam.ca](mailto:hecht.geoffrey@uqam.ca) (G. Hecht), [privat.jean@uqam.ca](mailto:privat.jean@uqam.ca) (J. Privat), [yann-gael.gueheneuc@concordia.ca](mailto:yann-gael.gueheneuc@concordia.ca) (Y.-G. Guéhéneuc).

and integration of the services, etc., which we discuss in details later. Also, the modernization of legacy systems may have some side effects that could affect the expected or claimed benefits of the migration of legacy systems (Khadka et al., 2015; Wagner, 2014). Such side effects could be the decrease of the system's performance, users resistance to the new technology/system, the unexpected high cost of the modernization, the increasing time to finish the migration, etc.

An organization may adopt one of three strategies to migrate legacy software systems to SOA. It can migrate its legacy systems through a *top-down*, forward-engineering strategy by: (1) performing a high-level decomposition of its domain artifacts, (2) modeling the needed services that will take part of the targeted SOA, (3) implementing those services, and (4) implementing the process that orchestrates all these services.

An organization may also want to use a *bottom-up* strategy to re-engineer its legacy software systems to a service-oriented style by: (1) extracting all the dependencies of their legacy system, (2) mining the existing applications for reusable functionality that could qualify as services, (3) packaging these functions as services to enable their reuse and to delete their dependencies to the legacy infrastructures, and (4) rewriting some existing applications to use the newly-identified services.

An organization may also adopt a *hybrid* strategy and reuse its legacy artifacts by: (1) grouping the functions of the applications into coarse functional blocks, (2) mapping those functional blocks to available services while deleting their dependencies to the legacy infrastructure, and (3) implementing the process orchestrating these services.

*Service identification* is central to all aforementioned three migration strategies, and has been recognized by practitioners as the most challenging step of the overall migration process (Khadka et al., 2013; Abdellatif et al., 2018). The services identified through a Service Identification Approaches (SIAs) must meet a range of expectations regarding their capabilities, quality of service, efficiency of use, etc. Lewis et al. (2005), which we also discuss in details later. To the best of our knowledge, all bottom-up and hybrid SIAs focus solely on identifying services in legacy software systems, not in ensuring that they can be then called "*as identified*" by different clients immediately. Indeed, once services become available, multiple clients may call them simultaneously, which may and may not cause problems in the services themselves (because they store some states) or related databases (because they do not take into account multiple clients/tenants). The challenges of turning such legacy code into autonomous and self-contained services include dealing with multi-tenancy, data consistency and statefulness. The legacy code might have side effects that violate one or more service design principles. These challenges must be considered after identifying the services, as part of the whole migration process of legacy software systems (Furda et al., 2017).

Due to the importance of SIAs and their impact on the success of legacy migrations to SOA, the literature proposed several approaches for identifying services in legacy systems. The selection of a SIA that is suitable for some practitioners among all other SIAs is however difficult and depends on several factors, e.g., the available legacy artifacts, the process of analyzing these legacy artifacts, the available inputs, the desired outputs and the usability degree of the approach. As a result, practitioners need a comprehensive view of existing SIAs to select the identification approach fulfilling their needs.

In the following, we propose a systematic literature review (SLR) of published SIAs, with focusing on *bottom-up* and *hybrid* approaches that use existing software artifacts. We chose to focus on bottom-up and hybrid approaches because previous studies (Abdellatif et al., 2018; Rodriguez et al., 2013) and our own

preliminary study showed that companies often have only source code as most up-to-date source of information about their legacy software systems.

Based on this SLR, we also present a taxonomy of SIAs, i.e., a multi-layer classification of SIAs. This classification helps practitioners in selecting a suitable service identification approach that corresponds to their migration needs. We perform our SLR using our experience with legacy software modernization, discussions with industrial partners, and the analysis of 41 papers retained from a first set of 3246 papers. We validate the correctness and coverage of our SLR through a survey and one-to-one interviews with 45 industrial experts in legacy software-systems modernization. The results show that our taxonomy conforms to the industrial experts' experiences, with a precision of 99%, and a recall of 94%.

### 1.1. Research questions

Through our SLR, we study the SIAs following four dimensions: the used inputs, the applied processes, the resulting outputs, and the usability degree of the approaches. We set out to answer the following research questions:

- **RQ1: What are the inputs used by SIAs?** We aim to identify the different inputs used by SIAs that are based on software systems analyses. We aim to classify the targeted SIAs based on the artifacts used for the identification.
- **RQ2: What are the processes followed by SIAs?** We aim to describe the processes that underlie the service identification approaches reported in the literature. This entails gathering information about, (1) the techniques used to identify candidate services, (2) the desired quality metrics, (3) the direction of the identification, (4) the automation level, and (5) the type of analysis used.
- **RQ3: What are the outputs of SIAs?** We aim to report information about the generated outputs of service identification approaches in terms of the targeted service types.
- **RQ4: What is the usability of SIAs?** We aim to study the usability degree of service identification approaches in the literature based on the systems used to validate the results, the accuracy of the identification method (when reported), the tool support, and the quality of the reported identification results.

We answer these questions and conclude that the state-of-the-art SIAs are still at their infancy. This is due to four main reasons: (1) the lack of validation on real enterprise-scale systems; (2) the lack of tool support, (3) the lack of automation of SIAs, and (4) the lack of assessment of the quality of the identified services. The results also show that the proposed SIAs generally ignore the economic aspects of the identification phase such as the implementation and maintenance costs, the re-factoring costs, and time-to-market issues. We believe that more work should be done to automate state-of-the-art SIAs and consider enterprise-scale systems to validate the proposed approaches. We also believe that regardless of the sought quality attributes, SIAs should provide means to assess the quality of the identified services and consider economic aspects in their identification process.

### 1.2. Outline

The remainder of this paper is structured as follows. Section 2 describes our SLR methodology. Section 3 describes the inputs used by SIAs. Section 4 describes the processes that underlie the studied SIAs. Section 5 surveys the outputs of SIAs. Section 6 describes the usability level of these SIAs. Section 7 details the

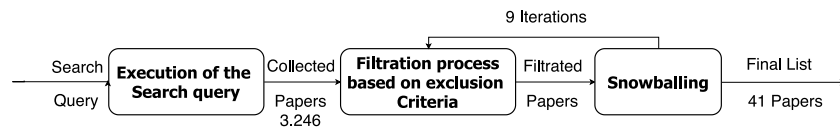


Fig. 1. Paper selection.

validation of our taxonomy. Section 8 synthesizes the comparison between the studied SIAs. Section 9 describes related work. Finally, Section 10 concludes our work.

## 2. Search methodology

In this section, we describe the design methodology of our systematic literature review as well as the mechanisms and data that we analyze to answer our research questions. We follow the procedures proposed by Kitchenham (2004) for performing systematic reviews.

Fig. 1 depicts our methodology. We first collected research papers based on search queries. We started by identifying relevant query terms based on our research questions and the context of our work: *service identification*, *SOA*, and *migration*. Then, for each keyword, we identified a set of related terms and synonyms using an online synonym finder tool<sup>1</sup> and defined the following query:

(service identification OR service mining OR service packaging) AND (migration OR modernization OR transformation OR re-engineering) AND (legacy OR existing systems OR Object-Oriented)

We executed this query in different scientific search engines, such as Google Scholar, ACM Digital Library, and IEEE Xplore Digital Library, Engineering Village, etc.

Our search queries returned a total of 3,246 unique references. We then filtered these references, first, based on their titles, second, based on their abstracts, and finally, based on their contents. Two of the authors manually and independently analyzed all the papers and then reconciled any differences through discussions. We excluded from our review papers meeting one of the following criteria:

- Papers not written in English.
- Papers not related to service identification.
- Papers about *top-down* SIAs.
- Papers that did *not* propose a technique or a methodology for service identification.
- Papers published before 2004 and after 2019.

Based on these exclusion criteria, we reduced the number of references and retain 26 papers that focus on SIAs that analyze software artifacts. We believe that our search string may not cover all query terms related to service identification (e.g., *microservices*, *decomposition*, *restructuring*, etc.) and thus we risk to miss important studies. To minimize these threats, we (1) included in our search string the most important keywords related to *service identification*, and (2) applied forward and backward snowballing (Wohlin, 2014; Felizardo et al., 2016) to minimize the risk of missing important papers. Forward snowballing refers to the use of the bibliographies of the papers to identify new papers that are referenced. Backward snowballing refers to the identification of new papers citing the papers being considered. We iterated the backward and forward snowballing and apply for each candidate paper our exclusion criteria. We stopped the iteration process when we have found no new candidate paper. We performed a total of nine iterations and added 15 papers. We thus obtained 41 papers that describe different SIAs, presented in Table 1.

## 3. RQ1: What are the inputs used by SIAs?

Using suitable inputs for service identification is crucial to the quality of the identified services and thus the migration process (Zadeh et al., 2012). When it comes to legacy systems, not all software-related artifacts (e.g., use cases, business process models, activity diagrams, etc.) are always available. Consequently, as depicted in Table 1, many SIAs in the literature relied on different types of inputs. When considering bottom-up and hybrid approaches, they all use source code or related models, as well as other types of input. We classify the inputs into three main categories: (1) executable models of the systems, (2) non-executable models of the systems, and (3) domain artifacts. We discuss them in turn, below.

### 3.1. Executable models

Executable models of the systems include source code and database schemas and test cases.

#### 3.1.1. Source code

*"If the map and the terrain disagree, trust the terrain".*

—Swiss Army Aphorism

With legacy systems, documentation (the map) is often missing or out of date. The source code (the terrain) becomes the only reliable source of information about the system. Source code is the most commonly used software artifact by the existing SIAs, due to its availability. SIAs that use source code as input identify business capabilities of the existing legacy systems and expose them as reusable services. Such SIAs rely on reverse and re-engineering processing to (1) extract dependencies between program elements such as variables, functions, modules/classes, etc.; (2) recover other kinds of information such as data flow diagrams, use cases, business process models, state machine diagrams, etc.; (3) map the source code to other artifacts such as business process models, use cases and database schemas, to complete the system map; and, usually, (4) apply clustering techniques to extract reusable services.

For legacy object-oriented systems, some SIAs rely on the relationships among classes to analyze the system structure and identify highly cohesive and loosely coupled reusable parts that could be exposed as services. For example, Adjoyan et al. (2014) relied on the analysis of dependencies between the classes of legacy object-oriented software systems. They proposed a fitness function that takes into account the type of relationship between the classes and assigns a score for each relationship. They then applied an agglomerative clustering technique to group classes into candidate services. Aversano et al. (2008) mined candidate services from the analysis of legacy source-code. They applied reverse-engineering techniques to extract UML diagrams of systems and analyze the signatures of related methods to identify candidate services.

Other SIAs identify services by analyzing the source code of non-object-oriented software systems. For example Rodriguez et al. (2013) reported the analysis of a large legacy system in an Argentinian government agency written in COBOL and running on

<sup>1</sup> <https://www.synonym-finder.com/>.

**Table 1**

Inputs of service identification approaches (**SC** for Source Code, **DB** for Database, **LogT** for Log Traces, **UAI** for User Application Interaction, **BPM** for Business Process Model, **UC** for Use Case, **AD** for Activity diagram, **DFD** for Data Flow Diagram, **SMD** for State Machine Diagram, **Ont** for Ontology, **Hu.Exp** for Human Expertise, **Doc** for documentation)

Method	Ex. Rep. of the Soft.			Non Ex. Rep. of the Soft.								Domain artifacts		
	SC	DB	TEST	Runtime artifacts		Model artifacts						Ont	Hu.Exp	Doc
				LogT	UAI	BPM	UC	AD	DFD	SMD				
Service identification based on quality metrics ( <a href="#">Adjoyan et al., 2014</a> )	x													
A spanning tree based approach to identifying web services ( <a href="#">Jain et al., 2004</a> )	x						x	x						
Generating a REST Service Layer from a Legacy System ( <a href="#">Rodríguez-Echeverría et al., 2014</a> )	x													
A service identification framework for legacy system migration into SOA ( <a href="#">Alahmari et al., 2010</a> )	x					x		x						
Reusing existing object-oriented code as web services in a SOA ( <a href="#">Sneed et al., 2013</a> )	x					x								
Mining candidate web services from legacy code ( <a href="#">Aversano et al., 2008</a> )	x													x
From objects to services: toward a stepwise migration approach for Java applications ( <a href="#">Marchetto and Ricca, 2009</a> )	x		x			x								x
Migrating interactive legacy systems to web services ( <a href="#">Canfora et al., 2006</a> )					x		x			x				
MDCSIM: A method and a tool to identify services ( <a href="#">Huergo et al., 2014a</a> )	x					x				x		x		
Reverse engineering relational databases to identify and specify basic Web services with respect to service oriented computing ( <a href="#">Baghdadi, 2006</a> )		x												
Identifying services in procedural programs for migrating legacy system to service oriented architecture ( <a href="#">Nakamur et al., 2012</a> )	x								x					x
A service-oriented analysis and design approach based on data flow diagram ( <a href="#">Zhao et al., 2009</a> )	x								x					
Service discovery using a semantic algorithm in a SOA modernization process from legacy web applications ( <a href="#">Sosa-Sánchez et al., 2014</a> )	x					x						x	x	
Incubating services in legacy systems for architectural migration ( <a href="#">Zhang and Yang, 2004</a> )	x				x	x								x
Migrating to web services: A research framework ( <a href="#">Sneed, 2007</a> )	x													x
Service identification and packaging in service oriented re-engineering ( <a href="#">Zhang et al., 2005</a> )	x					x								
A wrapping approach and tool for migrating legacy components to web services ( <a href="#">Chenghao et al., 2010</a> )	x													
Extracting reusable object-oriented legacy code segments with combined formal concept analysis and slicing techniques for service integration ( <a href="#">Zhang et al., 2006</a> )	x											x		

(continued on next page)

**Table 1** (continued).

Method	Ex. Rep. of the Soft.			Non Ex. Rep. of the Soft.						Domain artifacts		
	SC	DB	TEST	Runtime artifacts		Model artifacts				Ont	Hu.Exp	Doc
				LogT	UAI	BPM	UC	AD	DFD	SMD		
Using dynamic analysis and clustering for implementing services by reusing legacy code (Fuhr et al., 2011)	x			x		x					x	
Service mining from legacy database applications (Saha, 2015)		x										
An approach for mining services in database oriented applications (Del Grosso et al., 2007)		x									x	
Using user interface design to enhance service identification (Mani et al., 2008)	x				x							
A method to identify services using master data and artifact-centric modeling approach (Huergo et al., 2014b)	x	x				x					x	
Multifaceted service identification: Process, requirement and data (Amiri et al., 2016)		x				x					x	
The service modeling process based on use case refactoring (Kim and Doh, 2007)	x						x				x	
Extracting reusable services from legacy object-oriented systems (Bao et al., 2010)	x		x	x			x				x	x
Locating services in legacy software: information retrieval techniques, ontology and FCA based approach (Djeloul, 2012)	x				x						x	
Microservices identification through interface analysis (Baresi et al., 2017)	x				x							
Functionality-Oriented microservice extraction based on execution trace clustering (Jin et al., 2018)	x		x	x	x						x	
Bottom-up and top-down cobol system migration to web services (Rodriguez et al., 2013)	x										x	x
Extraction of microservices from monolithic software architectures (Mazlami et al., 2017)	x											
Service Cutter: A systematic approach to service decomposition (Gysel et al., 2016)	x						x					
An approach to align business and IT perspectives during the SOA services identification (Souza et al., 2017)	x					x					x	
Discovering microservices in enterprise systems using a business object containment heuristic (De Alwis et al., 2018a)	x	x		x							x	
A heuristic approach to locate candidate web service in legacy software (Abdelkader et al., 2013)	x											

(continued on next page)

IBM mainframes. They analyzed the legacy source code to identify the transactions to be migrated to services. These transactions are then translated into Java code, which is easier to expose as Web services.

Although the identification of candidate services using source code analysis leads to reusable and fine grained services, a combination of this kind of input with other artifacts (e.g., business

processes, databases, etc.) can be used to identify services with more business values.

### 3.1.2. Databases

Architecturally, the database layer is important to manage the persistence of data. Database contents, schemas and transactions are the artifacts used by database-related SIAs (Baghdadi, 2006; Saha, 2015; Del Grosso et al., 2007). These approaches identify



**Table 1** (continued).

Method	Ex. Rep. of the Soft.			Non Ex. Rep. of the Soft.							Domain artifacts		
	SC	DB	TEST	Runtime artifacts		Model artifacts					Ont	Hu.Exp	Doc
				LogT	UAI	BPM	UC	AD	DFD	SMD			
Identifying microservices using functional decomposition (Tyszberowicz et al., 2018)	x						x						
Toward the understanding and evolution of monolithic applications as microservices (Escobar et al., 2016)	x												
From monolithic systems to microservices: A decomposition framework based on process mining (Taibi and Systä, 2019)	x			x		x							
Function-splitting heuristics for discovery of microservices in enterprise systems (De Alwis et al., 2018b)		x		x		x	x						
From a monolith to a microservices architecture: An approach based on transactional contexts (Nunes et al., 2019)	x												
Re-architecting OO software into microservices a quality-centered approach (Selmadji et al., 2018)	x												

data/entity services that provide access to, and management of, the persistent data of the systems (C.f. Section 5).

For example, Baghdadi (2006) identified entity services by extracting SQL statements from systems. They then re-factored these statements and added them to the specification of a list of candidate services using *CRUD operations patterns* (Create, Read, Update and Delete). Saha (2015) relied on identifying instances of database-access patterns (database related operations) to identify reusable services. Using specific quality metrics, they refined database-related operations and wrapped them into data/entity services. Interactions between the application to migrate and the database have been also used by Del Grosso et al. (2007) to identify pieces of functionalities that can be exported as services. They performed the identification using clustering techniques and formal concept analysis.

Although the identification of candidate services based on the study of database queries or schema leads to reusable and fine grained services—which can only be entity services (cf. Section 5), a forward-engineering process is needed to build more coarse-grained services, that combine these finer-grain services, into business services.

### 3.1.3. Test cases

A test case can be defined as a specification of the inputs, execution conditions, testing procedure, and expected output results that must be executed to achieve a testing objective, such as to verify compliance with a specific requirement.

We found only three SIAs that use test cases, among other inputs, to identify reusable services (Bao et al., 2010; Jin et al., 2018; Marchetto and Ricca, 2009). For example, Bao et al. (2010) use test cases as an intermediate input for service identification. They first analyzed the legacy system source code and manually identified candidate use cases that correspond to potential reusable services. Then, they derived test cases from these use cases and used them to drive the execution of legacy-software systems. They used dynamic analysis techniques to analyze the execution log traces and generate coarse-grained code segments for each candidate use case that corresponds to an identified service. Also, Jin et al. (2018) only used test cases to execute different paths of the system and generate the corresponding

log traces. They analyzed these log files to get all classes and method invocations of the system. They then applied a clustering algorithm to group high cohesive and loosely coupled group of classes that will be mapped into services.

As shown by Table 1, test cases are rarely used by SIAs. They are only used as an intermediate artifact to guide the service identification process, probably because test cases are seldom available, and when they are, they cover only a small portion of the system.

### 3.2. Non-executable models

We distinguish between two categories of non-executable models: runtime artifacts extracted during the execution of the systems, and non-executable models that describe the architecture of the systems. We discuss them below.

#### 3.2.1. Runtime artifacts

Runtime artifacts are extracted during the executions of the systems. They contain log traces and user-application interactions (e.g., user interfaces).

**Log traces.** Execution traces of legacy software systems depict the dynamic behavior of the systems. Six SIAs rely on log traces to extract sequence calls related to specific execution scenarios (Fuhr et al., 2011; Bao et al., 2010; Jin et al., 2018; De Alwis et al., 2018a; Taibi and Systä, 2019; De Alwis et al., 2018b). These approaches identify pieces of legacy code executed during a set of business processes (Fuhr et al., 2011) or use cases (Bao et al., 2010), which are usually identified manually by business analysts. Then, they suggest those pieces of code as potential implementations of services. For example, Fuhr et al. (2011) applied mapping techniques of legacy code to business processes. They used log trace analyses and clustering techniques. They cluster the classes identified in the log traces according to their usage during the business processes.

We note that SIAs do not rely *solely* on log traces to identify services; they usually combine them with other types of inputs such as business process models, use cases, or human expertise.

**User interactions.** User-interface inputs capture the relationship between users and the system's functionalities. User interfaces usually embody data requirements and workflows (Hix and Hartson, 1993). If the workflow model of a system is not available, knowledge extracted from its user interfaces is useful to recover its underlying business logic (Mani et al., 2008; Zhang and Yang, 2004).

We found five SIAs that analyze users' interactions with user interfaces to identify services (Baresi et al., 2017; Mani et al., 2008; Canfora et al., 2006; Zhang and Yang, 2004; Djeloul, 2012; Jin et al., 2018). For example, Mani et al. (2008) proposed an XML-based representation, Unified User Interface Design Specification (UUIDS), to describe user interfaces, including data bindings and navigation events. They use this representation to automate the analysis of user interfaces and retrieve useful information for candidate services requirements.

The analyses of user interactions help to retrieve navigational information through the operations performed by users. They also help to identify reusable tasks with high business values, which could become services. However, SIAs based on user interactions are hardly automated. Further, they require a model of the tasks, as input, which may not be readily available.

### 3.2.2. Model artifacts

*Model artifacts* abstract the structure and execution behavior of systems. They include business process models, use cases, activity diagrams, and state machine diagrams, which are discussed in turn below.

**Business Process Model (BPMs).** They describe sets of activities and tasks that accomplish an organizational goal (Weske, 2012). BPMs have been used extensively by SIAs because of their ability to describe the business logic of legacy software systems at a high-level of abstraction. Business processes can be modeled with the Business Process Model and Notation (BPMN) and executed through their corresponding Business Process Execution Language (BPEL). The decomposition of business processes is a common strategy to identify services (Amiri et al., 2016). Business process-driven SIAs usually decompose business processes into tasks. These tasks are then clustered and exposed as services.

For example, Alahmari et al. (2010) identified services based on analyzing business process models. These business process models are derived from questionnaires, interviews and available documentations that provide atomic business processes and entities on the one hand, and activity diagrams that provide primitive functionalities on the other hand. The activity diagrams are manually identified from UML class diagrams extracted from the legacy code using IBM Rational Rose. Different service granularity levels are distinguished, as they pertain to atomic business processes and entities. Related atomic processes and entities are grouped together within the same service candidates to maximize cohesion of candidate services and minimize coupling between them. Fuhr et al. (2011) relied on business process models to correlate classes of legacy object oriented systems. Each activity in the business process model is executed. The classes that are called during the execution of a task are considered to be related. The identification of services is based on a clustering technique where the similarity measurement is based on how many classes are used together in the activity executions.

In the context of service identification, BPMs help to understand and capture the broad functional domains of legacy systems and how they interact with each other. Furthermore, business process-driven approaches identify high-level candidate services (based on process and tasks activities). However, the major problem with relying on BPMs to identify services is that such models are not always available especially for legacy software systems.

**Use cases.** They help to identify, at a high-level of abstraction, the interactions between users and systems to achieve goals. Use cases depicts functional requirements as well as sequences of actions that can be used for service identification (Vemulapalli and Subramanian, 2009). We found seven SIAs that use such artifact (Jain et al., 2004; Canfora et al., 2006; Kim and Doh, 2007; Bao et al., 2010; Gysel et al., 2016; Tyszbewicz et al., 2018; De Alwis et al., 2018b).

For example, Bao et al. (2010) analyze of the relationships between use-case elements to identify reusable services. They consider independent use cases of object-oriented systems are *candidate* services. If a use case *A* extends a use case *B*, they consider *B* as a candidate service, whereas *A* is not. Further, if use case *A* specializes (inherits from) use case *B*, then *A* is considered as a candidate service, whereas *B* is not.

The main reasons for SIAs to rely on use cases is that they offer systematic and intuitive means of capturing functional requirements with a focus on value to the users. However, to the best of our knowledge, SIAs based on use cases are difficult to automate to the extent that they often rely on human expertise.

**Activity diagram.** They show interactions in systems as well as the different steps involved in executing tasks (Schmuller, 2004). Only two SIAs use activity diagrams to identify services (Jain et al., 2004; Alahmari et al., 2010). For example, Alahmari et al. (2010) extracted, from activity diagrams, useful information and transform them to BPMN using mapping rules. They then analyzed the business process models to extract reusable services. They used activity diagrams of legacy systems as input but concretely relied on analyzing the BPMNs to identify reusable services in the system.

None of the identified SIAs relied only on activity diagrams. Other types of inputs are usually used such as source code, BPMs, and use cases to complement the identification process of candidate services.

**Data flow diagram.** A *Data Flow Diagram* (DFD) is a graphical representation of *functional dependencies*, based on the analysis of data flows, between business functions or processes (Ambler, 2004). The main entities of a DFD are the (1) data stores storing data for later use, (2) external entities representing sources/destinations of the data, (3) processes manipulating the data, and (3) data flows. Only two SIAs use DFDs to identify reusable services (Nakamur et al., 2012; Zhao et al., 2009).

For example, Zhao et al. (2009) rely on DFDs to identify services. They start by elaborating DFDs based on the system source code analysis. They recommend to design new DFDs for coarse-grained processes and to delete from the diagrams the fine-grained ones. They map each process of the elaborated DFDs to a service. They finally recommend to design a composite service that will capture the operations provided by identified services and allow these operations to be invoked in a defined workflow structure.

DFDs can describe the business logics of a software system. However, they are not always available nor straightforward to generate from legacy systems. SIAs based on DFDs of ill-structured systems do not guarantee as well the identification of relevant services (Nakamur et al., 2012; Zhao et al., 2009). Further, DFDs cannot represent dynamic dependencies because they are only based on the source code of software systems.

**State machine diagram.** A *State Machine Diagram* (SMD) shows a dynamic view of a system and describes the different states that entities can have during their lifetimes (Aggarwal and Sabharwal, 2012). We found that only two SIAs use state machine diagrams as inputs (Canfora et al., 2006; Huergo et al., 2014a). Canfora et al. (2006) used these diagrams to model the interactions between

users and systems. [Huergo et al. \(2014a\)](#) used them to model the life-cycle of master data, defined as any information considered to play a key role in the operation of a business.

Although state machine diagrams are ideal for describing the behavior of a limited number of objects, they are not suitable for SIAs that are dealing with large systems due to the state-explosion problem. Further, they are seldom available, and are not easy to obtain from source code or documentation.

### 3.3. Domain artifacts

Domain artifacts provide knowledge about the application domain of the systems. They include software documentation, human expertise, and ontologies.

#### 3.3.1. Documentation

Software documentation describes and documents systems at different levels of abstraction ([Lethbridge et al., 2003](#)). Software documentation includes textual descriptions as well as diagrams and models, such as the ones discussed above. Software documentation can guide SIAs by reducing the search space for candidate services by describing key functionalities of the systems. Some SIAs rely on software documentation to better understand the system at hand, which helps to identify reusable services ([Rodriguez et al., 2013](#); [Aversano et al., 2008](#); [Nakamur et al., 2012](#); [Zhang and Yang, 2004](#); [Sneed, 2007](#); [Bao et al., 2010](#)). For example, [Aversano et al. \(2008\)](#) proposed a SIA that analyzes the Javadoc documentation of systems to calculate lexical similarity between the classes or methods of the systems; they then used that similarity to identify clusters of functionality that can map to services. [Rodriguez et al. \(2013\)](#) described an industrial case study in which the documentation of a COBOL system was used to understand the system and to identify business rules in the code.

As with many other inputs (e.g., business process models, log traces, use cases, etc.), software documentation is not always available, and often outdated or out of sync with the source code of legacy systems.

#### 3.3.2. Human expertise

Human expertise appears in different ways in SIAs. It has been used to fine tune the parameters of various service identification algorithms (see e.g. [Jain et al. \(2004\)](#)). It has also been used to define the business logic and translate it into business processes ([Alahmari et al., 2010](#); [Amiri et al., 2016](#); [Sosa et al., 2013](#)). It is also needed to analyze use cases and identify candidate services ([Bao et al., 2010](#)). Finally, human expertise is needed to define data flow diagrams of the system to then identify candidate services ([Nakamur et al., 2012](#); [Zhao et al., 2009](#); [Sneed, 2007](#)).

Human expertise in SIAs limits the automation of service identification approaches and it appears in most of SIAs at different steps of the identification process.

#### 3.3.3. Ontologies

An ontology is a structured set of terms representing the semantics of a domain, whether through metadata or elements of a knowledge domain ([Bechhofer, 2009](#)). Several SIAs use ontologies to identify services ([Del Grosso et al., 2007](#); [Amiri et al., 2016](#); [Djeloul, 2012](#); [Chen et al., 2009](#)).

For example, [Djeloul \(2012\)](#) proposed a WordNet-based technique to identify services. They built queries by analyzing users interfaces. They then used WordNet to expand the queries and identify pieces of code participating in services. They also used information-retrieval techniques, such as vector-space model and latent-semantic analysis, to map queries to the relevant code.

[Chen et al. \(2009\)](#) started by analyzing the source code of systems and used three types of ontologies: a domain concept ontology, a functionality ontology, and a software-component ontology. They used formal and relational concept analysis to map source code of legacy systems to the ontologies they specified to identify candidate services.

The major challenge of ontology-based SIAs lies in defining the proper ontologies for the system. Also, the high cost of developing ontologies in terms of time, effort and resources remain a well-known bottleneck in the ontology development process ([Zhao et al., 2008](#)). Finally, ontology-based SIAs are complex and require a lot of human expertise.

## 4. RQ2: What are the processes followed by SIAs?

A service-identification process applies one or more identification techniques (e.g. wrapping, clustering, formal concept analysis, etc.) that target a set of quality metrics (e.g. coupling, cohesion, granularity, etc.) based on a predefined identification direction (i.e. bottom-up, top-down or hybrid). Human expertise defines the automation degree of the process, based on specific analysis types (e.g. static, dynamic, lexical, etc.).

### 4.1. Techniques of SIAs

We classified techniques of SIAs into six types (see [Table 2](#)):

- **Wrapping:** A black-box identification technique that encapsulates the legacy system with a service layer without changing its implementation. The wrapper provides access to the legacy system through a service encapsulation layer that exposes only the functionalities desired by the software architect ([Canfora et al., 2006](#); [Sneed, 2006](#)).
- **Genetic Algorithm:** A meta-heuristic for solving optimization problems that is based on "natural selection". It relies on the calculation of a fitness function to reach an optimal (or near-optimal) solution. By definition, an optimal solution is a feasible solution where the fitness function reaches its maximum (or minimum) value ([Balabanović and Shoham, 1997](#)).
- **Formal concept analysis (FCA):** A method for data analysis where we derive implicit relationships between objects in a formal way. It is also considered as a principled way of grouping objects that have common properties ([Birkhoff, 1940](#)). To use FCA, we should first specify the context denoted by a triple  $C=(E, P, R)$  where  $E$  is a set of finite elements,  $P$  is a set of finite properties and  $R$  is a binary relation based on  $E$  and  $P$ . Also a *formal concept* is defined as a grouping of all the elements that share a common set of properties. A partial order could be defined on the formal concepts with *concept lattices* ([Gratzer, 2009](#)), which also offer a structured visualization of the concepts hierarchy.
- **Clustering:** It consists of classifying and partitioning data into clusters (also called groups, categories or partitions) that share common properties. These clusters are built based on the internal homogeneity of their elements and the external separation between them. In fact, elements in the same cluster should be similar to each other while elements in different clusters should not ([Xu and Wunsch, 2005](#)).
- **Custom heuristics:** Some authors proposed their own heuristic algorithms, instead of using predefined algorithms, to decompose legacy software into SOA.
- **General guidelines:** they refer to approaches that only propose best practices, lessons learned, or recommendations for service identification.

In the following, we describe and discuss the use of these techniques to identify services from legacy systems.



#### 4.1.1. Wrapping

Wrapping-based SIAs use this technique for encapsulating a legacy system (or subset thereof) with a service layer and exporting its functionalities without changing its implementation (Canfora et al., 2006). Seven SIAs use/propose wrapping techniques (Rodríguez et al., 2013; Rodríguez-Echeverría et al., 2014; Sneed et al., 2013; Canfora et al., 2006; Sneed, 2007; Chenghao et al., 2010; Saha, 2015). For example, Canfora et al. (2006) proposed a wrapping methodology to expose the interactive functionalities of systems as services. The wrapper acts as an interpreter of a Finite State Automaton (FSA) that describes the interaction model between the system interfaces and their users. Also, Sneed et al. (2013) proposed an automatic wrapping technique based on the analysis of the public method interfaces of object-oriented code. They transform the public method interfaces into a relational table. Then based on this table, they generate WSDL interfaces that describe the functionalities of web services. Finally, they generate from the definitions of WSDL service interfaces the corresponding BPEL scripts to manage the service, as well as the corresponding test script to test the service. Wrapping techniques do not require to understand fully the architectures/implementations of the legacy software systems. It avoids the decomposition of the systems into reusable services. However, the underlying systems still must be maintained and so still need legacy expertise.

#### 4.1.2. Genetic algorithms

We found only three SIAs that rely on Genetic Algorithms to identify services from legacy software systems (Jain et al., 2004; Abdelkader et al., 2013; Amiri et al., 2016). For example, Jain et al. (2004) used Genetic Algorithms to identify services in legacy source code. They proposed an identification technique that is based on spanning trees. They used these representations to provide developers with a set of possible solutions for the identification problem. They also used a multi-objective genetic algorithm to refine the initial set of service decompositions. The multi-objective Genetic Algorithm relied on a fitness function that takes into consideration a set of managerial goals (i.e., cost effectiveness, ease of assembly, customization, reusability, and maintainability) to get a near-optimal solution for the service identification problem. Abdelkader et al. (2013) proposed also a Genetic Algorithm-based SIA. However, they only take into consideration the functional cohesion of a set of legacy system modules.

Although Genetic Algorithm-based SIAs may yield near-optimal solutions of reusable services, these SIAs do not guarantee to obtain systematically the optimal services that (1) maximize (or minimize) the fitness function, and (2) are architecturally relevant for the identification problem. Also, the relevance of the identified services highly depend on the choice of the objectives/managerial goals of the identification.

#### 4.1.3. Formal concept analysis

SIAs based on formal concept analysis basically rely on ontologies and/or concept lattices (Gratzer, 2009) to identify services (Zhang et al., 2006; Del Grosso et al., 2007; Chen et al., 2009). These SIAs usually rely on concept lattices to order the identified formal concepts and/or to visualize these concepts as well as the specified ontologies—when used. For example, Zhang et al. (2006) used formal concept analysis and program slicing to identify services in object-oriented systems. They begin by mapping the program entities (classes, methods) into elements and properties, using documentation and human expertise. They then applied the Ganter algorithm (Ganter, 2010) to build the concept lattices. Finally, they visualized, interpreted and analyzed these concepts to get meaningful, useful, and reusable services.

Also, Del Grosso et al. (2007) identified database-related features to be exported as services. They started by collecting database queries, using the dynamic execution of the database oriented systems. They then performed an analysis of the queries fields (i.e., the SELECT and the FROM clauses) and constraints (i.e., the WHERE clauses). They built a formal context using the concept lattice technique (Wille, 1982). They used FCA to group related queries into concepts and map them to candidate services.

The big challenge of using FCA for service identification consists in well identifying the concepts related to the entities of legacy systems. A proper setting of the formal context and their entities is required to ensure proper identification of reusable services. Also, the lack of automation in setting the formal context of the system may hinder the use of FCA algorithms to identify services in enterprise-scale systems.

#### 4.1.4. Clustering

SIAs use clustering to group classes or functionalities in legacy systems and consider each group as a candidate service. In general, they combine clustering techniques and custom heuristics. SIAs based on clustering belong to either one of two categories: classes clustering (Adjoyan et al., 2014; Jain et al., 2004; Zhang and Yang, 2004; Fuhr et al., 2011; Amiri et al., 2016; Baresi et al., 2017; Jin et al., 2018; Mazlami et al., 2017; Gysel et al., 2016; De Alwis et al., 2018a; Tyszbrowicz et al., 2018; Escobar et al., 2016; De Alwis et al., 2018b; Nunes et al., 2019; Selmadji et al., 2018) or functionalities clustering techniques (Saha, 2015; Zhang et al., 2005). The main clustering techniques used in the literature are k-means (Jain, 2010; Fuhr et al., 2011) and hierarchical-agglomerative clustering (Murtagh and Legendre, 2014; Zhang and Yang, 2004).

For example, Zhang and Yang (2004) proposed an agglomerative hierarchical clustering technique to extract reusable services from object-oriented legacy code. They started by analyzing legacy source code to calculate the similarity between the source code entities. The similarity metric consider the relationship between classes (i.e., inheritance, association, etc.) as well as the semantic similarity between them according to their names. They finally express the results in a dendrogram, which presents a hierarchic view of several possible decompositions of the system into services. Also, Fuhr et al. (2011) used k-means clustering techniques to identify services according to their type. The similarity measurement is based on how many classes are used together in a targeted activity execution.

K-means clustering techniques are indeed straightforward to apply. However, their results in the context of service identification show below-average performance. On the other hand, SIAs based on hierarchical clustering techniques do not require to specify in advance the number of the needed clusters/services. However, a subjective choice of the cutting point level in the generated dendrogram is needed to get the final set of services. This could be problematic for enterprise-scale systems where the number of possibilities for cutting points could be important. The choice between K-means and hierarchical clustering depends on the application context where K-means could be a good option when practitioners already know the number of services to be identified. On the other hand, hierarchical clustering is good for the case of unknowing the number of services to be identified. In this case, the hierarchical clustering will partition the system into a number of services based on the inter and intra cluster scaling.

#### 4.1.5. Custom heuristics

Some SIAs use dedicated heuristics (Adjoyan et al., 2014; Jain et al., 2004; Nakamura et al., 2012; Zhao et al., 2009; Zhang and Yang, 2004; Fuhr et al., 2011; Mazlami et al., 2017; Gysel et al., 2016) to identify services from legacy systems. Heuristics techniques are usually used with clustering techniques and genetic

algorithms. They also rely on quality metrics to identify candidate services.

For example, [Adjoyan et al. \(2014\)](#) proposed a fitness function based on three characteristics of services: composability, self-containment, and functionality. They grouped classes from object-oriented legacy software systems using a hierarchical-agglomerative clustering algorithm, which groups classes using the value of the fitness function. Also, [Jain et al. \(2004\)](#) proposed a set of heuristics based on dynamic and static relationships among classes in object-oriented systems. Then, they used these heuristics with a multi-objective optimization algorithm to get sets of classes representing services.

Although the use of heuristics is common in SIAs, their main challenge consists in establishing reliable heuristics to guide the process of identifying reusable services.

#### 4.1.6. General guidelines

We found two works in the literature that propose only general guidelines for service identification ([Alahmari et al., 2010](#); [Sneed, 2007](#)).

For example, [Alahmari et al. \(2010\)](#) proposed to extract UML activity diagrams from legacy systems and perform a model-to-model transformation to obtain BPMN from the diagrams. They argued that having a well defined SOA migration meta-model is important to make the migration process effective. They recommended the use of ad-hoc metrics because they assist in deriving optimal services with suitable granularity. Also [Sneed \(2007\)](#) proposed several guidelines for discovering potential services, evaluating these services and extracting their code from legacy systems. They recommended the use of a highly customizable rule based decision making mechanisms to identify which portions of legacy code could be potential services. They also recommended the use of DFDs to analyze data flow of the identified portions of code and decide about its business value.

SIAs based on guidelines propose general ideas to extract services from legacy software systems. They are indeed difficult to validate and automate.

### 4.2. Quality of identified services

Achieving the desired level of quality is critical for service based architectures. As a result, some SIAs use/target some quality metrics/requirements to obtain high-quality candidate services.

#### 4.2.1. Quality requirements

We describe the quality requirements targeted by the studies SIAs as follows:

- **Reuse:** The ability of a service to participate in multiple service assemblies (compositions) ([Feuerlicht et al., 2007](#)). Better reusability should provide better return of investment (ROI) and shorter development times ([Alkkiomäki and Smolander, 2016](#)).
- **Maintainability:** Services should ease the effort to modify their implementation, to identify root causes of failures, to verify changes, etc. ([Pereplechikov et al. \(2007\)](#)).
- **Interoperability:** The ability of a service to communicate and be invoked by other systems/services implemented in different programming languages ([Erl, 2005](#)).
- **Self-containment:** A service should be completely self-contained to be deployed as a single unit, without depending on other services ([Adjoyan et al., 2014](#)).
- **Composability:** Services should be composable with one another to be reused and integrated as services that control other services or that provide functionalities to other services ([Sindhgatta et al., 2009](#)).

As we can see in [Table 3](#), a few SIAs consider quality requirements in their identification techniques. However, service reuse is the most considered requirement by these approaches. On the other hand, we notice that few studies consider the study of composability, self-containment, maintainability, and interoperability of the identified services. This could be because these quality requirements are (1) difficult to characterize and measure and (2) hardly provide useful insights to identify services.

#### 4.2.2. Quality metrics

We describe the quality metrics targeted by the studied SIAs as follows:

- **Coupling:** The dependencies among services should be minimized and the functionalities should be encapsulated to limit the impact of changes in one service to other services ([Pereplechikov et al., 2007](#)).
- **Cohesion:** Cohesion is a measure of the strength of the relationships among programming entities (e.g., classes, functions, etc.) implementing a service and the functionality provided by the service ([Abdelkader et al., 2013](#)).
- **Granularity:** An adequate granularity is a primary concern of SIAs. It can be adjusted to the scope of the functionality offered by the service ([Huerigo et al., 2014a](#)).
- **Total number of services:** SIAs must not have too many "small" services or not enough services ([Gysel et al., 2016](#)).

[Table 4](#) shows that state-of-the-art SIAs highly rely on the use of some specific quality metrics such as loose coupling, high cohesion, and granularity. However, these SIAs fail at providing a comprehensive quality model to assess and evaluate the quality of the identified services.

### 4.3. Directions of SIAs

SIAs can follow three directions: top-down, bottom-up, and hybrid.

- A top-down process starts with high-level artifacts, e.g., domain analysis or requirement characterization of systems to define their functionalities. They do not consider low-level artifacts to identify services. Hence, we do not consider these SIAs in our study.
- A bottom-up process starts with low-level artifacts to maximize code reuse and minimize changes. It extracts more abstract artifacts, e.g., architectures, which can be used to identify candidate services. It can also identify new services that fill implementation gaps or meet new requirements ([Bell, 2009](#)).
- A hybrid process combines a top-down and a bottom-up process. It uses both requirements and implementation artifacts to identify the candidate services.

As we focus in this SLR on SIAs that follow the bottom-up and hybrid direction, we report in [Table 5](#) the distribution of SIAs over these two directions. [Table 5](#) shows that there are almost equal numbers of bottom-up and hybrid SIAs in the literature. Finally we notice that bottom-up SIAs are more successful at delivering services in the short-term but they usually identify fine-grained services with limited reuse. Moreover, Hybrid SIAs tend to complement and reduce the limitations of bottom-up approaches by also considering requirements.

**Table 2**  
Targeted techniques of SIAs.

Technique	SI Method	Total
Wrapping	Rodriguez et al. (2013), Rodríguez-Echeverría et al. (2014), Sneed et al. (2013), Canfora et al. (2006), Sneed (2007), Chenghao et al. (2010) and Saha (2015)	7
Genetic algorithm	Jain et al. (2004), Abdelkader et al. (2013) and Amiri et al. (2016)	3
Formal concept analyses	Zhang et al. (2006), Del Grosso et al. (2007), Djeloul (2012) and Chen et al. (2009)	4
Clustering	Adjoyan et al. (2014), Jain et al. (2004), Zhang et al. (2005), Zhang and Yang (2004), Fuhr et al. (2011), Saha (2015), Amiri et al. (2016), Baresi et al. (2017), Jin et al. (2018), Mazlami et al. (2017), Gysel et al. (2016), De Alwis et al. (2018a), Tyszbrowicz et al. (2018), Escobar et al. (2016), De Alwis et al. (2018b), Nunes et al. (2019) and Selmadji et al. (2018),	17
Custom heuristics	Adjoyan et al. (2014), Jain et al. (2004), Rodríguez-Echeverría et al. (2014), Aversano et al. (2008), Marchetto and Ricca (2009), Huergo et al. (2014a), Nakamur et al. (2012), Zhao et al. (2009), Souza et al. (2017), Zhang and Yang (2004), Fuhr et al. (2011), Mani et al. (2008), Kim and Doh (2007), Bao et al. (2010), Chen et al. (2009), Jin et al. (2018), Mazlami et al. (2017), Gysel et al. (2016), Souza et al. (2017), De Alwis et al. (2018a), Taibi and Systä (2019), Escobar et al. (2016), De Alwis et al. (2018b), Nunes et al. (2019) and Selmadji et al. (2018)	25
General guidelines	Sneed (2007) and Alahmari et al. (2010)	2

**Table 3**  
Targeted quality requirements by SIAs.

Quality requirement	SI Method	Total
Reuse	Rodriguez et al. (2013), Marchetto and Ricca (2009), Zhang and Yang (2004), Sneed (2007), Amiri et al. (2016) and Kim and Doh (2007)	6
Maintainability	Zhang et al. (2005)	1
Interoperability	Sneed (2007)	1
Self-containment	Adjoyan et al. (2014) and Zhang and Yang (2004)	2
Composability	Adjoyan et al. (2014)	1

**Table 4**  
Targeted quality metrics by SIAs.

Quality metric	SI method	Total
Coupling	Adjoyan et al. (2014), Jain et al. (2004), Alahmari et al. (2010), Huergo et al. (2014a), Nakamur et al. (2012), Zhao et al. (2009), Zhang and Yang (2004), Zhang et al. (2005), Huergo et al. (2014b), Baresi et al. (2017), Tyszbrowicz et al. (2018), Escobar et al. (2016), Taibi and Systä (2019), De Alwis et al. (2018b), Nunes et al. (2019) and Selmadji et al. (2018)	16
Cohesion	Rodriguez et al. (2013), Adjoyan et al. (2014), Jain et al. (2004), Alahmari et al. (2010), Huergo et al. (2014a), Nakamur et al. (2012), Zhao et al. (2009), Zhang et al. (2005), Huergo et al. (2014b), Amiri et al. (2016), Baresi et al. (2017), Tyszbrowicz et al. (2018), Escobar et al. (2016), De Alwis et al. (2018b) and Selmadji et al. (2018)	15
Granularity	Rodriguez et al. (2013), Marchetto and Ricca (2009), Huergo et al. (2014a), Nakamur et al. (2012), Zhao et al. (2009), Zhang and Yang (2004), Zhang et al. (2005), Huergo et al. (2014b), Kim and Doh (2007), Baresi et al. (2017), Tyszbrowicz et al. (2018), Escobar et al. (2016) and Nunes et al. (2019)	13
Number of services	Adjoyan et al. (2014), Marchetto and Ricca (2009), Gysel et al. (2016) and Selmadji et al. (2018)	4

**Table 5**  
Identification process directions of service identification methods in the literature.

Direction	SI method	Total
Bottom-up	Rodriguez et al. (2013), Adjoyan et al. (2014), Jain et al. (2004), Rodríguez-Echeverría et al. (2014), Baghdadi (2006), Nakamur et al. (2012), Zhao et al. (2009), Sneed (2007), Chenghao et al. (2010), Zhang et al. (2006), Saha (2015), Del Grosso et al. (2007), Djeloul (2012), Baresi et al. (2017), Jin et al. (2018), Mazlami et al. (2017), Gysel et al. (2016), De Alwis et al. (2018a), Abdelkader et al. (2013), Escobar et al. (2016), Selmadji et al. (2018) and Chen et al. (2009)	22
Hybrid	Alahmari et al. (2010), Sneed et al. (2013), Aversano et al. (2008), Marchetto and Ricca (2009), Canfora et al. (2006), Huergo et al. (2014a), Sosa-Sánchez et al. (2014), Zhang and Yang (2004), Zhang et al. (2005), Fuhr et al. (2011), Mani et al. (2008), Amiri et al. (2016), Kim and Doh (2007), Bao et al. (2010), Souza et al. (2017), Tyszbrowicz et al. (2018), Taibi and Systä (2019), De Alwis et al. (2018b) and Nunes et al. (2019)	19

#### 4.3.1. Analyses types

SIAs may perform static, dynamic, lexical analyses, or some combination thereof to identify services.

- Static analysis is performed without executing a software system. Dependencies between classes are potential relationships, like method calls and access attributes. These dependencies are analyzed to identify strongly connected classes, for example, to identify services. Adjoyan et al. (2014), Rodríguez-Echeverría et al. (2014), Sneed et al. (2013), Aversano et al. (2008), Baghdadi (2006), Zhang and Yang (2004), Sneed (2007), Zhang et al. (2005) and Cheng-

hao et al. (2010) are examples of identification methods based only on static analysis. The main advantage of static analysis is that it depends only on the source code. It does not address polymorphism and dynamic binding.

- Dynamic analysis is performed by examining the software system at run time. Dependencies between software elements (e.g., class instantiations and accesses (Bao et al., 2010), function calls (Zhang et al., 2005; De Alwis et al., 2018b), relationships between database tables (De Alwis et al., 2018a), etc.) are collected during the program execution (Shatnawi et al., 2018). The execution is performed



based on a set of cases that covers the system functionalities, called execution scenarios.

- Lexical analysis techniques suppose that the similarity between the classes should be taken into account during service identification process. This analysis plays the main role in approaches that used features location and textual similarity techniques.

Table 6 shows that 76% of SIAs rely on static analysis, 39% on dynamic analysis, and 21% on lexical analysis. Finally we found that 38% rely on a combination of analyses to reduce the limitations of each individual analysis.

#### 4.4. Automation of SIAs

Automation is the degree to which a SIA needs human experts. We distinguish three levels of automation: manual, semi-automatic, and fully automatic.

- Manual SIAs depend entirely on human experts. They only provide general guidelines to experts to identify services without automating any step of the service identification process (Rodríguez-Echeverría et al., 2014; Marchetto and Ricca, 2009).
- Semi-automatic SIAs need human experts to perform some of their tasks. For example, Jain et al. (2004) proposed a SIA that require a human expert to provide objective functions and specify weights for each of them.
- Automatic SIA do not need any human intervention during the identification process. We did not find any approach in the literature that fully automates the identification of services in existing systems.

Table 7 shows that there is a lack of automation of SIAs: 88% of the SIAs are semi-automatic or manual.

### 5. RQ3: What are the outputs of SIAs?

In the following, we discuss the output of SIAs in terms of the target service architecture (service-based/microservice-based) and discuss the types of services considered by these approaches.

#### 5.1. Service architecture

Service identification approaches aim at identifying services that will be integrated in a SOA.

In the past few years, several SIAs have been interested in identifying microservices – a variant of the service-oriented architecture style – to migrate legacy systems to microservice-based systems (Baresi et al., 2017; Mazlami et al., 2017; Gysel et al., 2016; Tyszberowicz et al., 2018; Escobar et al., 2016; Taibi and Systä, 2019; De Alwis et al., 2018b; Nunes et al., 2019; Selmadji et al., 2018). For example, Escobar et al. (2016) proposed a microservice identification approach to migrate a monolithic Java Enterprise Edition (JEE) application to microservices. They performed a static analysis to cluster session and entity beans into microservices. They started by associating a cluster to each session bean. They grouped these clusters according to a *clustering threshold* that focuses on structural coupling and cohesion. The distance between clusters is calculated based on the number of shared entity beans.

Mazlami et al. (2017) proposed a microservices identification approach that relies on the analysis of data collected from a version control repository of a monolithic application. They also applied clustering and custom heuristics to extract loosely-coupled and high-cohesive set of classes that will be mapped to

microservices. Both semantic and logical coupling metrics were considered by their clustering algorithm. In particular, they combined three metrics to identify microservices: *semantic coupling* (to identify groups of classes that belong to the same domain), *single responsibility principle* (to analyze classes that change together in commits), and *contributor coupling* (to identify classes accessed by the same development team). All these metrics were combined and used by a clustering algorithm to identify groups of classes that belong to the same domain and could represent a microservice.

We notice that microservices identification approaches rely on clustering and custom heuristics to decompose the system into small services. Although the granularity is an important characteristic for qualifying microservices, none of the studied approaches provided a comprehensive model to evaluate whether microservices are identified with the right level of granularity. Also, the granularity difference between services and microservices is still neither well defined nor clearly discussed by the studied microservices identification approaches.

#### 5.2. Service types

We identified only four SIAs that identify specific types of services in existing systems (Alahmari et al., 2010; Marchetto and Ricca, 2009; Huergo et al., 2014a; Fuhr et al., 2011) and nine papers proposing service taxonomies (Alahmari et al., 2010; Marchetto and Ricca, 2009; Huergo et al., 2014a; Fuhr et al., 2011; Erl, 2007; Cohen, 2007; Krafzig et al., 2005; Ani and Baghdadi, 2015; Gu and Lago, 2010), that classify services with hierarchical-layered schemas to support the communication among stakeholders during the implementation of SOAs. These existing taxonomies offer several service types with different classification criteria (e.g., granularity Alahmari et al., 2010; Erl, 2007; Cohen, 2007; reuse Fuhr et al., 2011; Cohen, 2007; Ani and Baghdadi, 2015, etc.) and different names for the same service types. We studied these previous works and identified the following six service types that are generic and cover most of the existing service types. We validated our taxonomy through an industrial survey with practitioners (Abdellatif et al., 2018) that we detail in Section 7.

1. **Business-process services:** (Also called business service (Alahmari et al., 2010; Fuhr et al., 2011; Cohen, 2007; Ani and Baghdadi, 2015)), they correspond to business processes or use cases. These are services used by users. These services compose or use the enterprise-task, application-task, and entity services described in the following. Examples of business-process services include flight booking services, hotel booking services and sales order services.
2. **Enterprise-task services:** (Also called capabilities Cohen, 2007), they are of finer granularity than business-process services. They implement generic business functionalities reused across different applications. Examples of enterprise-task services include “online payment” and “tax calculation”.
3. **Application-task services:** (Also called task, activity or composite service (Alahmari et al., 2010; Marchetto and Ricca, 2009; Baghdadi, 2006; Ani and Baghdadi, 2015)), they provide functionalities specific to one application. They exist to support reuse within one application or to enable business-process services (Cohen, 2007). Examples of application-task services include quoting request and invoicing that take part in the sales order business process of a typical ERP system.



**Table 6**  
Analyses types of SIAs.

Analysis type	SI method	Total
Static analysis	Adjoyan et al. (2014), Jain et al. (2004), Rodríguez-Echeverría et al. (2014), Sneed et al. (2013), Aversano et al. (2008), Huergo et al. (2014a), Baghdadi (2006), Sneed (2007), Chenghao et al. (2010), Sosa-Sánchez et al. (2014), Zhang and Yang (2004), Zhang et al. (2005), Mani et al. (2008), Amiri et al. (2016), Kim and Doh (2007), Souza et al. (2017), Zhang et al. (2006), Saha (2015), Del Grosso et al. (2007), Djeloul (2012), Chen et al. (2009), Baresi et al. (2017), Rodríguez et al. (2013), Mazlami et al. (2017), Gysel et al. (2016), De Alwis et al. (2018a), Abdelkader et al. (2013), Tyszbrowicz et al. (2018), Escobar et al. (2016), Nunes et al. (2019) and Selmadji et al. (2018)	31
Dynamic analysis	Jain et al. (2004), Alahmari et al. (2010), Marchetto and Ricca (2009), Canfora et al. (2006), Huergo et al. (2014a), Nakamur et al. (2012), Zhao et al. (2009), Bao et al. (2010), Zhang et al. (2005), Fuhr et al. (2011), Mani et al. (2008), Jin et al. (2018), Mazlami et al. (2017), De Alwis et al. (2018a), Taibi and Systä (2019) and De Alwis et al. (2018b)	16
Lexical analysis	Aversano et al. (2008), Sosa-Sánchez et al. (2014), Zhang and Yang (2004), Zhang et al. (2006), Del Grosso et al. (2007), Djeloul (2012), Baresi et al. (2017) and Mazlami et al. (2017)	8

**Table 7**  
Automation of SIAs.

Analysis type	SI method	Total
Automatic	Baghdadi (2006), Saha (2015), Gysel et al. (2016), Souza et al. (2017) and Abdelkader et al. (2013)	5
Semi-automatic	Rodríguez et al. (2013), Adjoyan et al. (2014), Jain et al. (2004), Alahmari et al. (2010), Sneed et al. (2013), Aversano et al. (2008), Canfora et al. (2006), Huergo et al. (2014a), Nakamur et al. (2012), Zhao et al. (2009), Sosa-Sánchez et al. (2014), Zhang and Yang (2004), Sneed (2007), Zhang et al. (2005), Chenghao et al. (2010), Zhang et al. (2006), Fuhr et al. (2011), Del Grosso et al. (2007), Mani et al. (2008), Amiri et al. (2016), Bao et al. (2010), Djeloul (2012), Baresi et al. (2017), Jin et al. (2018), Mazlami et al. (2017), De Alwis et al. (2018a), Escobar et al. (2016), Taibi and Systä (2019), De Alwis et al. (2018b), Nunes et al. (2019), Selmadji et al. (2018) and Chen et al. (2009)	32
Manual	Rodríguez-Echeverría et al. (2014), Marchetto and Ricca (2009), Kim and Doh (2007) and Tyszbrowicz et al. (2018)	4

4. **Entity services:** (Also called information or data services (Alahmari et al., 2010; Erl, 2007; Ani and Baghdadi, 2015)), they provide access to and management of the persistent data of legacy software systems. They support actions on data (CRUD) and may have side-effects (i.e., they modify shared data). Examples of entity services include management services for clients, bank accounts, and products.
5. **Utility services:** They do not support directly the business-process services but provide some cross-cutting functionalities required by domain-specific services (Ani and Baghdadi, 2015; Fuhr et al., 2011; Huergo et al., 2014a). Examples of typical utility services include notification, logging, and authentication.
6. **Infrastructure services:** They allow users deploying and running SOA systems. They include services for communication routing, protocol conversion, message processing and transformation (Alahmari et al., 2010). They are sometimes provided by an Enterprise Service Bus (ESB). They are reused in more services than utility services. Examples of infrastructure services include publish-subscribe, message queues, and ESB.

Most of SIAs identify general services of SOA without specifying different service types, e.g., Adjoyan et al. (2014), Zhang et al. (2005) and Nakamura et al. (2009). Only a few approaches (Alahmari et al., 2010; Marchetto and Ricca, 2009; Huergo et al., 2014a; Fuhr et al., 2011) considered the identification of specific types of services in existing systems.

For example, Alahmari et al. (2010) identified services based on analyzing business process models. These business process models are derived from questionnaires, interviews and available documentations that provide atomic business processes and entities on the one hand, and activity diagrams that provide primitive functionalities on the other hand. The activity diagrams are manually identified from UML class diagrams extracted from the legacy code using IBM Rational Rose. Different service granularity are distinguished in relation to atomic business processes

and entities. Dependent atomic processes as well as the related entities are grouped together at the same service to maximize the cohesion and minimize the coupling. There is no details about how to identify the different service types. Fuhr et al. (2011) identified three types of services. These are business, entity and utility services. The services are identified from legacy codes based on a dynamic analysis technique. The authors relied on a business process model to identify correlation among classes. Each activity in the business process model is executed. Classes that have got called during the execution are considered as related. The identification of services is based on a clustering technique where the similarity measurement is based on how many classes are used together in the activity executions. The identified clusters are manually interpreted and mapped into the different service types. Classes used only for the implementation of one activity are grouped into a business service corresponding to this activity. Entity services are composed of clusters of classes that contribute to implement multiple activities but not all of them. A Cluster of classes that are used by all of the activities represent the implementation of utility services. A strong assumption regarding this approach is that business process model should be available to identify execution scenarios.

We notice that there is a lack of SIAs that are type-centric: only four SIAs focus on the identification of specific types of services from legacy systems. These approaches focus on identifying business (Alahmari et al., 2010; Huergo et al., 2014a; Fuhr et al., 2011), entity (Marchetto and Ricca, 2009; Huergo et al., 2014a) and utility services (Alahmari et al., 2010; Marchetto and Ricca, 2009; Huergo et al., 2014a; Fuhr et al., 2011). Also, none of the studied SIAs tried to identify enterprise-task or infrastructure services through the analysis of legacy systems. These type-centric SIAs do not distinguish in their service identification process between enterprise and application-task services as the scope of reuse of the identified services is not well studied or specified.

## 6. RQ4: What is the usability of SIAs?

Fig. 2 shows that we consider four elements to estimate the usability of SIAs: validation, accuracy, tool support, and result quality. We then introduce a measure of the usability of the SIAs based on these four elements and their values for each SIA.

### 6.1. Validation

Validation refers to the legacy software systems (if any) on which the SIA was applied. It can be industrial (e.g., real industrial systems), experimental (small, experimental systems), or none at all. We evaluate the usability of a SIA as follows. If the validation is performed on (1) industrial systems, it is “high”; (2) experimental systems, it is “medium”, else (3) it is “low”. We found that only 34% of SIAs were validated on real industrial systems, with most SIAs validated on experimental systems or not validated at all. This lack of industrial validation is a major threat to the applicability of SIAs.

### 6.2. Accuracy/precision

We assign “high”, “medium”, and “low” values to the accuracy/precision of SIAs. We assign “high” if it is greater than 80%, medium if it is between 50% and 79% in the SIA, and low if it is less than 50%.

Although the accuracy/precision of SIAs is important, we found that only few SIAs have reported accuracy/precision (as depicted in Table 8).

### 6.3. Tool support

Tool support refers to the tool(s) implementing a SIA and their maturity, if any.

We consider the tool support of a SIA as “high” if it is open-source or industry ready, “medium” if it is only a prototype, and “low” if there is little or no tool support.

### 6.4. Result quality

Result quality is an estimation of the quality of the identified candidate services and whether or not the authors detailed well their proposed SIA. It can be “high”, “medium”, or “low”.

### 6.5. Usability

We consider these four preceding elements to estimate the usability of SIAs. We assign to each SIA a usability degree (UD) as follows:

$$UD = \sum_{i=1}^4 Score_i$$

$Score_i \in \{high = 1, medium = 0, low = -1\}, \forall i \in \{1, \dots, 4\}$  and refers to validation, accuracy, tool support, and usability, respectively.

If  $UD \geq 1$ , then  $UD = high$ .

If  $UD = 0$ , then  $UD = medium$ .

If  $UD \leq -1$ , then  $UD = low$ .

We tried our best to consider the most important usability criteria and give a rational estimation of the usability degree of the studied SIAs. For example, as shown in Table 8, to calculate the usability of the SIA of Rodríguez-Echeverría et al. (2014), we studied the scores relative to tool support, validation, identification accuracy, and quality results of the approach. This study has a high tool support through the tool named MIGRARIA

(tool-support score is 1). It is validated on an experimental system (validation score is 0). There was no mention of the accuracy/precision of the approach and thus we did not consider associated scores for calculating the usability of the approach. Finally, based on our judgment of the whole approach, we estimated that this SIA has high quality results (quality result is 1). We added all these scores and obtain a usability score of two, which we qualified as a high usability degree.

Table 8 shows that 39% of SIAs have a high usability degree while 22% have medium usability, and 39% have low usability. These results show that the studied SIAs are still in their infancy, mainly due to (1) the lack of validation on industrial systems, (2) the lack of estimation of their accuracy/precision, (3) their lack of tool support, and (4) their lack of automation.

## 7. Taxonomy and validation

Fig. 2 shows the taxonomy resulting from our answers to the research questions. This taxonomy directly derive from the previous sections.

We believe that the validation of a taxonomy is difficult for several reasons. In fact, it is a tool for researchers and practitioners and, as such, it should be used to assess its strengths and limitations. Also, a taxonomy often cannot be compared against other ones, either because they do not exist or because they have different objectives. Consequently, to validate our taxonomy, we performed a survey with industrial experts.

### 7.1. Methodology

We conducted a survey with 45 industrial experts to validate our taxonomy and also obtain their informed opinions about legacy-to-SOA migration in general and service identification in particular (Abdellatif et al., 2018). We conducted this survey between October 2017 and March 2018 in five main steps:

#### 7.1.1. Preparation of the survey

We created a Web-based survey<sup>2</sup> using Google Forms. We built our survey on our taxonomy: the individual questions correspond to each composite node of the taxonomy and their possible answers correspond to the leave nodes.

Before releasing the survey, we performed a pilot study with six participants, three from academia and three from industry, and validated the relevance of the questions, their wording, the coverage of their answers, etc. The six participants went through the questions and suggested few minor changes.

The final survey contained six sections: (1) participants' professional and demographic data, (2) type of migrated systems and reasons for migration, (3) general information about SIAs (perception of importance, strategy, inputs, level of automation), (4) technical information about SIAs (techniques, quality metrics), (5) types of sought services, and (6) used tools and best practices.

For example, we asked the participants the following questions: “What information do/did you use to identify services?”, “What kind of identification techniques do/did you apply?”, “What are the types of the migrated services?”, etc. We provide a list of possible answers for each question and ask the participants to mention any other answer if he/she did not select any possibility from the provided list.

<sup>2</sup> <https://goo.gl/forms/EE31KeA7R7pUeTYI2>.

**Table 8**  
Usability of SIAs.

Method	Tool support	Validation	Accuracy/Precision	Quality of results	Usability
Service identification based on quality metrics (Adjoyan et al., 2014)	Prototype	Experimental	Medium	Medium	Medium
A spanning tree based approach to identifying web services (Jain et al., 2004)	MOGA-WSI	Industry	NA	High	High
Generating a REST Service Layer from a Legacy System (Rodríguez-Echeverría et al., 2014)	MIGRARIA	Experimental	NA	High	High
A service identification framework for legacy system migration into SOA (Alahmari et al., 2010)	Prototype	Experimental	NA	Low	Low
Reusing existing object-oriented code as web services in a SOA (Sneed et al., 2013)	Industry ready	Industry	NA	High	High
Mining candidate web services from legacy code (Aversano et al., 2008)	NA	Experimental	NA	Low	Low
From objects to services: toward a stepwise migration approach for Java applications (Marchetto and Ricca, 2009)	NA	Experimental	NA	Low	Low
Migrating interactive legacy systems to web services (Canfora et al., 2006)	NA	Case Study	NA	Medium	Low
MDCSIM: A method and a tool to identify services (Huergo et al., 2014a)	MDCSIM	Industry	NA	High	High
Reverse engineering relational databases to identify and specify basic Web services with respect to service oriented computing (Baghdadi, 2006)	CASE	Experimental	NA	Medium	High
Identifying services in procedural programs for migrating legacy system to service oriented architecture (Nakamur et al., 2012)	NA	Experimental	NA	Low	Low
A service-oriented analysis and design approach based on data flow diagram (Zhao et al., 2009)	SOAD	Experimental	NA	Low	Medium
Service discovery using a semantic algorithm in a SOA modernization process from legacy web applications (Sosa-Sánchez et al., 2014)	MigraSOA	Experimental	NA	Low	Medium
Incubating services in legacy systems for architectural migration (Zhang and Yang, 2004)	Prototype	Industry	NA	Low	Medium
Migrating to web services: A research framework (Sneed, 2007)	NA	No Validation	NA	Low	Low
Service Identification and Packaging in Service Oriented Reengineering (Zhang et al., 2005)	Prototype	Case Study	NA	Medium	Medium
A wrapping approach and tool for migrating legacy components to web services (Chenghao et al., 2010)	Prototype	Case Study	NA	Low	Low
Extracting reusable object-oriented legacy code segments with combined formal concept analysis and slicing techniques for service integration (Zhang et al., 2006)	Prototype	Experimental	NA	Low	Low
Using dynamic analysis and clustering for implementing services by reusing legacy code (Fuhr et al., 2011)	Prototype	Case Study	Medium	Low	Low
Service mining from legacy database applications (Saha, 2015)	Prototype	Industry	NA	High	High
An approach for mining services in database oriented applications (Del Grosso et al., 2007)	Prototype	Industry	High	High	High
Using user interface design to enhance service identification (Mani et al., 2008)	Prototype	Industry	NA	Medium	High
A method to identify services using master data and artifact-centric modeling approach (Huergo et al., 2014b)	NA	Experimental	NA	Low	Low
Multifaceted service identification: Process, requirement and data (Amiri et al., 2016)	Prototype	Experimental	High	Low	Medium
The service modeling process based on use case refactoring (Kim and Doh, 2007)	Prototype	Case Study	NA	Low	Low
Extracting reusable services from legacy object-oriented systems (Bao et al., 2010)	Prototype	Industry	NA	Medium	High

(continued on next page)

Table 8 (continued).

Method	Tool support	Validation	Accuracy/Precision	Quality of results	Usability
Locating services in legacy software: information retrieval techniques, ontology and FCA based approach (Djeloul, 2012)	Prototype	Case Study	NA	Low	Low
Microservices identification through interface analysis (Baresi et al., 2017)	NA	Case Study	NA	Low	Low
Extraction of microservices from monolithic software architectures (Mazlami et al., 2017)	Prototype	Industry	NA	High	High
Service Cutter: A systematic approach to service decomposition (Gysel et al., 2016)	ServiceCutter	Experimental	NA	High	High
Bottom-up and top-down cobol system migration to web services (Rodriguez et al., 2013)	Industry ready	Industry	NA	High	High
Functionality-Oriented microservice extraction based on execution trace clustering (Jin et al., 2018)	FOME	Experimental	NA	Low	Medium
An approach to align business and IT perspectives during the SOA services identification (Souza et al., 2017)	Prototype	Experimental	NA	Low	Low
Discovering microservices in enterprise systems using a business object containment heuristic (De Alwis et al., 2018a)	Prototype	Industry	NA	Medium	High
A heuristic approach to locate candidate web service in legacy software (Abdelkader et al., 2013)	Prototype	Experimental	NA	Low	Low
Identifying microservices using functional decomposition (Tysberowicz et al., 2018)	Prototype	Experimental	NA	Low	Low
Toward the understanding and evolution of monolithic applications as microservices (Escobar et al., 2016)	Prototype	Industry	NA	High	High
From monolithic systems to microservices: A decomposition framework based on process mining (Taibi and Systä, 2019)	Prototype	Industry	NA	High	High
Function-Splitting heuristics for discovery of microservices in enterprise systems (De Alwis et al., 2018b)	Prototype	Industry	NA	Medium	High
From a monolith to a microservices architecture: An approach based on transactional contexts (Nunes et al., 2019)	Prototype	Experimental	Medium	Medium	Medium
Re-architecting OO software into microservices a quality-centered approach (Selmadji et al., 2018)	Prototype	Experimental	NA	Medium	Medium

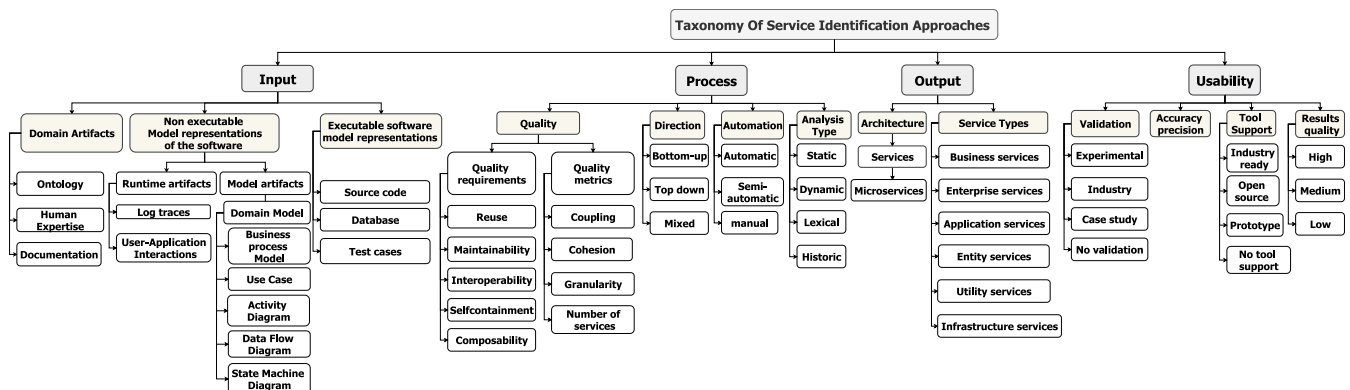


Fig. 2. Taxonomy of service identification approaches.

### 7.1.2. Selection of the participants

We targeted developers with an industrial experience in SOA migration. Identifying and soliciting such developers was challenging. We relied on (1) information about companies that offer modernization services, (2) online presentations and webinars made by legacy-to-SOA migration experts, and (3) search queries

on LinkedIn profiles: “legacy migration OR legacy modernization OR SOA architect OR SOA migration OR Cloud migration OR service migration OR service mining”. We attempted neither to be exhaustive in our search for participants nor to cover different strata of developers working on SOA migration. As such, our sample is a random sample of convenience.



Once we identified potential participants, we sent them invitations via e-mail, LinkedIn, Facebook, and Twitter. We chose *not* to solicit more than three participants from any given company (1) to have a diversity of companies and (2) to avoid overburdening a single company.

#### 7.1.3. Administering the survey

We invited 289 participants and recommended about 15 min to complete the survey. We asked potential participants to forward our invitations to colleagues with experience in SOA migration in their social/professional networks. The survey was completed by 47 participants, two of whom did not participate in SOA migration and whose answers we removed, for 45 complete answers.

#### 7.1.4. Validation of the survey

We assessed the reliability of the answers to the survey by searching for spurious/facetious answers, contradictory answers, etc. We also performed follow-up interviews with 24 of the 45 participants who agreed to such interviews by providing their e-mail addresses.

We interviewed 8 of these 24 participants. We used a two-pass method (Charmaz and Belgrave, 2012) to analyze our transcripts of the individual interviews.<sup>3</sup> We first performed a *thematic coding* to identify broad issues related to legacy-to-SOA migration in general and SIAs in particular. We then performed an *axial coding* to identify relationships among the identified issues. We identified major issues related to our taxonomy using *meta-codes*, which we then used to code manually the data of all the interviewees (Abdellatif et al., 2018).

#### 7.1.5. Validation of the taxonomy

We could then measure the precision and accuracy of our taxonomy (input, process, and output of SIAs) as:

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

where TP (*True Positive*) corresponds to a leaf/answer identified by both a participant and our taxonomy, FP (*False Positive*) corresponds to a leaf/answer identified in our taxonomy but not mentioned by any participant, and FN (*False Negative*) corresponds to an answer from at least one participant but not identified in our taxonomy.

We do not assess the usability of the SIAs for two reasons. First, we did not want to overburden the participants with the assessment of a very subjective metric and because we believe that such a measure deserves a complete, independent study (in future work).

### 7.2. Participants

We reached a total of 45 participants involved in legacy-to-SOA migration projects in different capacities: 50% were software architects, 23.7% were directors of technology, and 21% were software engineers. The remaining 5.3% were migration specialists, project managers, and CEOs.

The participants worked in different industries: 64% were in technology and telecommunication, 20% in banking and insurance, 12.8% in health, and 3.2% in education. In terms of experience, 78% had more than 10 years of experience, also reflected in their age distributions: 23% were less than 35 years old, 39% were between 36 and 45, 20.5% were between 46 and 55, and 17.5% were over 55.

**Table 9**

Validation results of the service-identification taxonomy.

Feature	Precision	Recall
Inputs	93%	100%
Techniques	100%	82%
Quality metrics	100%	73%
Direction	100%	100%
Automation	100%	100%
Analysis type	100%	100%
Service types	100%	100%

### 7.3. Validation results

Table 9 shows the validation results of our taxonomy with precision and recall: our taxonomy is conformed to the experts' experiences with a precision between 93% and 100% and a recall between 73% and 100%.

In particular, for inputs, precision is 93% and recall 100%. None of the participants mentioned the use of test scenarios to identify services. However, we kept this input in our taxonomy as it is used by two SIAs. In terms of the techniques, recall is 82%, which is acceptable, because some participants mentioned the use of in-house tools or manual identification. For quality metrics, recall is 73%, which is acceptable, because some participants mentioned the use of other economic quality metrics not considered in any of the SIAs: identification cost, adaptation effort, and time to market.

## 8. Discussions

In this section, we will discuss our observations about the studied SIAs in terms of the main nodes of our taxonomy: inputs, processes, outputs, and usability.

### 8.1. Inputs

SIAs rely on diverse types of inputs to identify services. We found that the most used inputs are source code and business-process models (BPMs). Combining multiple inputs is also common. The most used combination of inputs are also source code and BPMs (Sneed et al., 2013; Marchetto and Ricca, 2009; Zhang et al., 2005; Fuhr et al., 2011). Only 10 SIAs rely on a single input type (Adjoyan et al., 2014; Rodríguez-Echeverría et al., 2014; Baghdadi, 2006; Chenghao et al., 2010; Saha, 2015; Mazlami et al., 2017; Abdelkader et al., 2013; Escobar et al., 2016; Nunes et al., 2019; Selmadji et al., 2018), either source code again or databases.

### 8.2. Processes

Most SIAs rely on clustering and custom heuristics to identify services. The main challenge for these approaches is in using adequate heuristics to identify services.

The success of a SOA depends on the quality of the services. Services with low quality attributes may (1) affect reuse negatively and (2) compromise business agility and reduce return on investment (Huerger et al., 2014c). Quality attributes are therefore important to identify services. However, not all service quality requirements are considered by state-of-the-art SIAs. Moreover, regardless of the adopted quality requirements, SIAs should provide means to assess/control the quality of the candidate services. Also, there are many economic factors that SIAs should take into account. Such aspects could be the implementation and maintenance cost, the re-factoring cost of the system, and time-to-market. The economic aspects of the identified methods are widely ignored in the studied SIAs. We believe that more efforts should be done in SIAs to consider as well such economic aspects which play an important role to select the appropriate SIA for an organization.

<sup>3</sup> <https://goo.gl/ZYv2Ut> for sample transcripts.

### 8.3. Outputs

We noticed that microservices architectures have been gaining a lot of consideration in the past few years as we found many studies focusing on the identification of microservices in legacy systems. The applied identification techniques are quite similar to those used for identifying services. On the other hand, few SIAs focus on the identification of specific types of services. In particular we observed that these SIAs focus on identifying business, entity, and utility services but not enterprise/application-task and infrastructure services. Also, we noticed that these type-sensitive SIAs do not distinguish between enterprise and application task services as the scope of reuse of the identified services is not well specified/studied. We believe that the identification of services according to their types is a challenging problem because (1) we have to build a taxonomy that cover all service types, (2) define detection rules/signature for each service type, and (3) target the metrics or detection rules that are appropriate for each type. We believe that not all service types have distinct signatures as two different service types may leave similar or indistinguishable signatures in the code. The taxonomy of service types may not be representative of all existing service types. To mitigate this threat, we validated our taxonomy through an industrial survey with 45 practitioners who were involved in migration projects of legacy systems to SOA (Abdellatif et al., 2018). None of them mentioned the identification of new/other types of services.

### 8.4. Usability

We reported that 51% of the state-of-the-art SIAs have medium or low usability degree due to (1) their lack of validation on real industrial systems, (2) their lack of tool support, and (3) their lack of automation. In particular, most SIAs consider only small examples in their validation, also confirmed by some participants in our survey (Abdellatif et al., 2018). The participants reported that a problem exists in the knowledge transfer between academia and industry because of the lack of consideration of enterprise-scaled systems to validate the proposed SIAs in academia.

Finally, we believe that measuring the usability of a given SIA is quite difficult. Our proposed metric may partially measure the usability of a given SIA as we do not cover all possible usability-related aspects. However, we tried our best to consider the most important usability criteria such as the tool support, the quality of SIA results, the validation of the process and the accuracy/precision of the SIA. As a future work, we aim to empirically validate our proposed metric of usability with people from academia and industry to study its feasibility of quantifying/estimating the usability degree of a SIA.

### 9. Related work

Several systematic literature reviews and surveys on SIAs have been proposed in the literature. In the period from 2009 to 2019, ten surveys (Boerner and Goeken, 2009; Birkmeier et al., 2015; Gu and Lago, 2010; Cai et al., 2000; Vale et al., 2012; Taei Zadeh et al., 2012; Huergo et al., 2014c; Fritzsche et al., 2018; Bani-Ismael and Baghdadi, 2018a,b) on service identification were identified. Although these surveys had different goals, neither of them fully addressed all our research questions. Table 10 contains a summary and comparison between the most relevant surveys focusing on service identification in the literature.

For example, Boerner and Goeken (2009), only studied business-driven SIAs techniques and focused on their strategic and economic aspects. They stressed the consideration of economic aspects when identifying services based only on top

down approaches. Birkmeier and Overhage (2009) proposed a classification of SIAs between 1984 and 2008. This SLR is indeed old, does not fully addressed our research questions and does not cover recent SIAs. Cai et al. (2011) proposed another survey where they identified the most frequent activities in the state-of-the-art SIAs between 2004 and 2011. Then, Vale et al. (2012) made a comparison of SIAs and a list of recommendation of the most suitable SI technique according to stakeholders' needs in the Service-Oriented Product Line Engineering context. Bani-Ismael and Baghdadi (2018a,b) proposed two different surveys about service identification. In the first one they studied the evaluation frameworks for 24 state-of-the-art SIAs. In the second survey they only identified the challenges of 14 service identification approaches and their limitations. Both studies do not fully address our research questions as we do in our SLR.

Finally Fritzsche et al. (2018) provided a classification of refactoring approaches of monolithic applications to microservices. They studied 10 microservices identification approaches and provided a guide for decomposition approaches using microservices identification requirements.

Although there are several SLRs on service identification in the literature, none of these surveys fully addressed our research questions. Their focus differ deeply as we cover more in details state-of-the-art service identification approaches in terms of (1) the artifacts used by SIAs, (3) the processes of these approaches, (4) the outputs of these processes, and (5) the usability degree of these approaches. We also propose a taxonomy of SIAs and validate its correctness and coverage with industrial experts in legacy-to-SOA migration through surveys and one-on-one interviews.

### 10. Conclusion and future work

We presented in this paper a systematic literature review (SLR) on service identification approaches (SIAs) that use the artifacts to build legacy software systems as input. We studied the SIAs in terms of their inputs, their processes, their outputs, and their usability. We built our taxonomy on our experience with legacy software modernization, discussions with industrial partners, and the analysis of existing SIAs. We validated the correctness and the coverage of our taxonomy with industrial experts in legacy-to-SOA migration through surveys and one-on-one interviews. The validation results showed that our taxonomy is conformed to the industrial experts' experiences with 99% of precision and 94% of recall.

The results of our SLR show that the state-of-the-art SIAs are still at their infancy mainly due to (1) the lack of validation on real enterprise-scale systems; (2) the lack of tool support, and (3) the lack of automation of SIAs. The results also show that the proposed SIAs generally ignore the economic aspects of the identification phase as well as the identification by service type. Indeed despite of their importance in the migration process, only few SIAs consider the economic aspects of the service identification process such as the implementation and maintenance cost, the re-factoring cost of the system, and time-to-market. Also, most of the existing SIAs look for services based on their functional cohesion and low coupling with other parts of the applications, regardless of service types. Furthermore, we showed that the current trend of SIAs is the identification of microservices in existing systems. However, the applied identification techniques were very similar to those used for identifying services. The granularity border between services and microservices is still not well defined nor clearly discussed by these approaches. Finally, we found that most SIAs usually do not try to improve the quality attributes of the identified candidate services. We believe that regardless of the sought quality attributes, SIAs should provide means to assess the quality of the identified services. Also,

**Table 10**Systematic literature reviews of Service Identification in the literature (**A** for Addressed, **PA** for Partially Addressed, **NA** for Not Addressed)

SIA	Goal	Year of publication	Covered years	Included papers	RQ1	RQ2	RQ3	RQ4
Boerner and Goeken (2009)	Business-driven SI techniques comparison with the study of their strategic and economic aspects	2009	2005–2008	5	NA	PA	PA	A
Birkmeier et al. (2015)	Classification of service identification techniques	2009	1984*–2008	15	PA	A	PA	NA
Gu and Lago (2010)	Providing the basic elements of SI to help practitioners selecting the most suitable one basic on their needs	2010	2004–2009	30	A	A	A	NA
Cai et al. (2011)	Identify frequent used activities done in several SI research works	2011	2004–2011	41	PA	A	PA	NA
Vale et al. (2012)	Comparison of SI methods and recommendation of the most suitable SI technique according to stakeholders' needs in the Service-Oriented Product Line Engineering context	2012	2005–2012	32	PA	PA	PA	PA
Taei Zadeh et al. (2012)	Suitable inputs identification for SI methods in small and medium enterprise	2012	2002–2010	48	PA	PA	PA	NA
Huergo et al. (2014c)	Classification of SI methods	2014	2002–2013	105	PA	A	PA	NA
Bani-Ismael and Baghdadi (2018b)	Exploring existing evaluation frameworks for state-of-the-art SIAs	2018	2007–2016	23	PA	PA	NA	PA
Bani-Ismael and Baghdadi (2018a)	Identifying service identification challenges in service oriented architecture	2018	2005–2016	14	PA	NA	NA	NA
Fritzsch et al. (2018)	Classification of Refactoring Approaches of monolithic applications to microservices	2018	2015–2017	10	PA	PA	PA	NA
Our SLR	Focusing on bottom-up and hybrid SIAs based on the used input, the applied process, the generated output and the usability of the approach Reviewing SI from the point of view of researchers and practitioners interest	2020	2004–2019	41	A	A	A	A

we believe that more work should be done to automate the SIAs and consider enterprise-scaled systems to validate the proposed approaches.

As future work, we will generalize our survey and study top-down service identification approaches. We will study in detail SIAs that use some architecture-centric methods such as Architecture Tradeoff Analysis Method (ATAM) (Kazman et al., 1998), Attribute-Driven Design (ADD) (Nord et al., 2004), and Cost Benefit Analysis Method (CBAM) (Nord et al., 2003). These methods could assist a service identification approach by providing and evaluating architectural descriptions of the system to migrate. Also, we will study empirically the gap of the state of the practices of SIAs between academia and industry. We want to identify issues that the research community can address to ease knowledge transfer between academia and industry in the context of legacy-to-SOA migration. Finally, we believe that the identification of services according to their types is a challenging but interesting problem. As future work, we aim to propose a type-centric service identification approach that promote better reuse at the application, enterprise, and business levels.

## Acknowledgment

The authors would like to thank all the practitioners who replied to us. This work was supported by the Fonds de Recherche du Quebec Nature et Technologies (FRQNT) N2016-PR-191087.

## References

- Abdelkader, Mostefai, Malki, Mimoun, Benslimane, Sidi Mohamed, 2013. A heuristic approach to locate candidate web service in legacy software. *Int. J. Comput. Appl. Technol.* 47 (2–3), 152–161.
- Abdellatif, Manel, Hecht, Geoffrey, Mili, Hafedh, Elboussaidi, Ghizlane, Moha, Naouel, Shatnawi, Anas, Privat, Jean, Guéhéneuc, Yann-Gaël, 2018. State of the practice in service identification for soa migration in industry. In: *International Conference on Service-Oriented Computing*. Springer, pp. 634–650.
- Adjoyan, Seza, Seriai, Abdelhak-Djamel, Shatnawi, Anas, 2014. Service identification based on quality metrics - object-oriented legacy system migration towards SOA. In: *The 26th International Conference on Software Engineering and Knowledge Engineering*. Hyatt Regency, Vancouver, BC, Canada, pp. 1–6.
- Aggarwal, Manuj, Sabharwal, Sangeeta, 2012. Test case generation from UML state machine diagram: A survey. In: *Computer and Communication Technology (ICCT)*, 2012 Third International Conference on. IEEE, pp. 133–140.
- Alahmari, Saad, Zaluska, Ed, De Roure, David, 2010. A service identification framework for legacy system migration into SOA. In: *Services Computing (SCC)*, 2010 IEEE International Conference on. IEEE, pp. 614–617.
- Alkkio, Ville, Smolander, Kari, 2016. Anatomy of one service-oriented architecture implementation and reasons behind low service reuse. *Serv. Orient. Comput. Appl.* 10 (2), 207–220.
- Ambler, Scott W., 2004. *The Object Primer: Agile Model-Driven Development with UML 2.0*. Cambridge University Press.
- Amiri, Mohammad Javad, Parsa, Saeed, Lajevardi, Amir Mohammadzade, 2016. Multifaceted service identification: Process, requirement and data. *Comput. Sci. Inform. Syst.* 13 (2), 335–358.
- Ani, Bashar Al, Baghdadi, Youcef, 2015. A taxonomy-centred process for service engineering. *Int. J. Comput. Appl. Technol.* 52 (1), 1–17.
- Aversano, Lerina, Cerulo, Luigi, Palumbo, Ciro, 2008. Mining candidate web services from legacy code. In: *10th International Symposium on Web Site Evolution*. IEEE, pp. 37–40.
- Baghdadi, Youcef, 2006. Reverse engineering relational databases to identify and specify basic web services with respect to service oriented computing. *Inform. Syst. Front.* 8 (5), 395–410.
- Balabanović, Marko, Shoham, Yoav, 1997. Fab: content-based, collaborative recommendation. *Commun. ACM* 40 (3), 66–72.
- Bani-Ismael, Basel, Baghdadi, Youcef, 2018a. A literature review on service identification challenges in service oriented architecture. In: *International Conference on Knowledge Management in Organizations*. Springer, pp. 203–214.
- Bani-Ismael, Basel, Baghdadi, Youcef, 2018b. A survey of existing evaluation frameworks for service identification methods: towards a comprehensive evaluation framework. In: *International Conference on Knowledge Management in Organizations*. Springer, pp. 191–202.
- Bao, Liang, Yin, Chao, He, Weigang, Ge, Jun, Chen, Ping, 2010. Extracting reusable services from legacy object-oriented systems. In: *Software Maintenance (ICSM)*, 2010 IEEE International Conference on. IEEE, pp. 1–5.



- Baresi, Luciano, Garriga, Martin, De Renzis, Alan, 2017. Microservices identification through interface analysis. In: *European Conference on Service-Oriented and Cloud Computing*. Springer, pp. 19–33.
- Bechhofer, Sean, 2009. OWL: Web ontology language. In: *Encyclopedia of Database Systems*. Springer, pp. 2008–2009.
- Bell, Michael, 2009. SOA Modeling Patterns for Service Oriented Discovery and Analysis. John Wiley & Sons.
- Birkhoff, Garrett, 1940. *Lattice Theory*, Vol. 25. American Mathematical Soc..
- Birkmeier, Dominik, Klöckner, Sebastian, Overhage, Sven, 2015. A survey of service identification approaches—classification framework, state of the art, and comparison. *Enterp. Model. Inform. Syst. Arch.* 4 (2), 20–36.
- Birkmeier, Dominik, Overhage, Sven, 2009. On component identification approaches—classification, state of the art, and comparison. In: *Component-Based Software Engineering*. Springer, pp. 1–18.
- Boerner, René, Goeken, Matthias, 2009. Service identification in SOA governance literature review and implications for a new method. In: *Digital Ecosystems and Technologies*, 2009. DEST'09. 3rd IEEE International Conference on. IEEE, pp. 588–593.
- Cai, Simin, Liu, Yan, Wang, Xiaoping, 2011. A survey of service identification strategies. In: *Services Computing Conference (APSCC)*, 2011 IEEE Asia-Pacific. IEEE, pp. 464–470.
- Cai, Xia, Lyu, Michael R., Wong, Kam-Fai, Ko, Roy, 2000. Component-based software engineering: technologies, development frameworks, and quality assurance schemes. In: *Software Engineering Conference*, 2000. APSEC 2000. Proceedings. Seventh Asia-Pacific. IEEE, pp. 372–379.
- Canfora, Gerardo, Fasolino, Anna Rita, Frattolillo, Gianni, Tramontana, Porfirio, 2006. Migrating interactive legacy systems to web services. In: *Conference on Software Maintenance and Reengineering (CSMR'06)*. IEEE, pp. 10–pp.
- Charmaz, Kathy, Belgrave, Liska, 2012. Qualitative interviewing and grounded theory analysis. In: *The SAGE handbook of interview research*. pp. 347–365.
- Chen, Feng, Zhang, Zhuopeng, Li, Jianzhi, Kang, Jian, Yang, Hongji, 2009. Service identification via ontology mapping. In: *2009 33rd Annual IEEE International Computer Software and Applications Conference*, Vol. 1. IEEE, pp. 486–491.
- Chenghao, Guo, Min, Wang, Xiaoming, Zhou, 2010. A wrapping approach and tool for migrating legacy components to web services. In: *First International Conference on Networking and Distributed Computing (ICNDC)*, 2010. IEEE, pp. 94–98.
- Cohen, Shy, 2007. Ontology and taxonomy of services in a service-oriented architecture. *Arch. J.* 11 (11), 30–35.
- De Alwis, Adambarage Anuruddha Chathuranga, Barros, Alistair, Fidge, Colin, Polyvyanyy, Artem, 2018a. Discovering microservices in enterprise systems using a business object containment heuristic. In: *OTM Confederated International Conferences" on the Move To Meaningful Internet Systems"*. Springer, pp. 60–79.
- De Alwis, Adambarage Anuruddha Chathuranga, Barros, Alistair, Polyvyanyy, Artem, Fidge, Colin, 2018b. Function-splitting heuristics for discovery of microservices in enterprise systems. In: *International Conference on Service-Oriented Computing*. Springer, pp. 37–53.
- Del Grosso, Concettina, Di Penta, Massimiliano, de Guzman, Ignacio Garcia-Rodriguez, 2007. An approach for mining services in database oriented applications. In: *11th European Conference on Software Maintenance and Reengineering*, 2007. CSMR'07. IEEE, pp. 287–296.
- Djeloul, M.A.M., 2012. Locating services in legacy software: information retrieval techniques, ontology and FCA based approach. *WSEAS Trans. Comput.* 11 (1), 19–26, legacy software; information retrieval techniques; FCA based approach; Web services technology; WORDNET ontology; formal concepts analysis; source code;.
- Erl, Thomas, 2005. *Service-Oriented Architecture*, Vol. 8. Pearson Education Incorporated.
- Erl, Thomas, 2007. *SOA Principles of Service Design*. Prentice Hall PTR, Upper Saddle River, NJ, USA.
- Erl, Thomas, 2008. *SOA Design Patterns*. Pearson Education.
- Escobar, Daniel, Cárdenas, Diana, Amarillo, Rolando, Castro, Eddie, Garcés, Kelly, Parra, Carlos, Casallas, Rubby, 2016. Towards the understanding and evolution of monolithic applications as microservices. In: *2016 XLII Latin American Computing Conference (CLEI)*. IEEE, pp. 1–11.
- Felizardo, Katia Romero, Mendes, Emilia, Kalinowski, Marcos, Souza, Érica Ferreira, Vijaykumar, Nandamudi L, 2016. Using forward snowballing to update systematic reviews in software engineering. In: *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. ACM, p. 53.
- Feuerlicht, George, et al., 2007. Understanding service reusability. In: *International Conference Systems Integration*. Department of Information Technologies and Czech Society for Systems Integration.
- Fritzsche, Jonas, Bogner, Justus, Zimmermann, Alfred, Wagner, Stefan, 2018. From monolith to microservices: a classification of refactoring approaches. In: *International Workshop on Software Engineering Aspects of Continuous Development and New Paradigms of Software Production and Deployment*. Springer, pp. 128–141.
- Fuhr, Andreas, Horn, Tassilo, Riediger, Volker, 2011. Using dynamic analysis and clustering for implementing services by reusing legacy code. In: *Reverse Engineering (WCRE)*, 2011 18th Working Conference on. IEEE, pp. 275–279.
- Furda, Andrei, Fidge, Colin, Zimmermann, Olaf, Kelly, Wayne, Barros, Alistair, 2017. Migrating enterprise legacy source code to microservices: on multitenancy, statefulness, and data consistency. *IEEE Software* 35 (3), 63–72.
- Ganter, Bernhard, 2010. Two basic algorithms in concept analysis. *Formal Concept Anal.* 312–340.
- Gratzer, George, 2009. *Lattice Theory: First Concepts and Distributive Lattices*. Courier Corporation.
- Gu, Qing, Lago, Patricia, 2010. Service identification methods: a systematic literature review. In: *Towards a Service-Based Internet*. Springer, pp. 37–50.
- Gysel, Michael, Kölbener, Lukas, Giersche, Wolfgang, Zimmermann, Olaf, 2016. Service cutter: A systematic approach to service decomposition. In: *European Conference on Service-Oriented and Cloud Computing*. Springer, pp. 185–200.
- Hix, Deborah, Hartson, H. Rex, 1993. *Developing User Interfaces: Ensuring Usability Through Product & Process*. John Wiley & Sons, Inc..
- Huergo, Rosane S., Pires, Paulo F., Delicato, Flávia C., 2014a. Mdcsim: A method and a tool to identify services. *IT Converg. Pract.* 2 (4), 1–27.
- Huergo, Rosane S., Pires, Paulo F., Delicato, Flávia C., 2014b. A method to identify services using master data and artifact-centric modeling approach. In: *Proceedings of the 29th Annual ACM Symposium on Applied Computing*. ACM, pp. 1225–1230.
- Huergo, Rosane S., Pires, Paulo F., Delicato, Flávia C., Costa, Bruno, Cavalcante, Everton, Batista, Thais, 2014c. A systematic survey of service identification methods. *Serv. Orient. Comput. Appl.* 8 (3), 199–219.
- Jain, Anil K., 2010. Data clustering: 50 years beyond K-means. *Pattern Recognit. Lett.* 31 (8), 651–666.
- Jain, Hemant, Zhao, Huimin, Chinta, Nageswara R., 2004. A spanning tree based approach to identifying web services. *Int. J. Web Serv. Res.* 1 (1), 1.
- Jin, Wuxia, Liu, Ting, Zheng, Qinghua, Cui, Di, Cai, Yuanfang, 2018. Functionality-oriented microservice extraction based on execution trace clustering. In: *2018 IEEE International Conference on Web Services (ICWS)*. IEEE, pp. 211–218.
- Kazman, Rick, Klein, Mark, Barbacci, Mario, Longstaff, Tom, Lipson, Howard, Carriere, Jeremy, 1998. The architecture tradeoff analysis method. In: *Proceedings. Fourth IEEE International Conference on Engineering of Complex Computer Systems (Cat. No. 98EX193)*. IEEE, pp. 68–78.
- Khadka, Ravi, Saeidi, Amir, Jansen, Slinger, Hage, Jurriaan, 2013. A structured legacy to SOA migration process and its evaluation in practice. In: *Maintenance and Evolution of Service-Oriented and Cloud-Based Systems (MESOCA)*, 2013 IEEE 7th International Symposium on the. IEEE, pp. 2–11.
- Khadka, Ravi, Shrestha, Prajan, Klein, Bart, Saeidi, Amir, Hage, Jurriaan, Jansen, Slinger, van Dis, Edwin, Bruntink, Magiel, 2015. Does software modernization deliver what it aimed for? A post modernization analysis of five software modernization case studies. In: *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, pp. 477–486.
- Kim, Yookyong, Doh, Kyung-Goo, 2007. The service modeling process based on use case refactoring. In: *International Conference on Business Information Systems*. Springer, pp. 108–120.
- Kitchenham, Barbara, 2004. Procedures for performing systematic reviews. *Keele, UK, Keele University* 33 (2004), 1–26.
- Krafzig, Dirk, Banke, Karl, Slama, Dirk, 2005. *Enterprise SOA: Service-Oriented Architecture Best Practices*. Prentice Hall Professional.
- Lethbridge, Timothy C., Singer, Janice, Forward, Andrew, 2003. How software engineers use documentation: The state of the practice. *IEEE Softw.* 20 (6), 35–39.
- Lewis, Grace, Morris, Ed, O'Brien, Liam, Smith, Dennis, Wrage, Lutz, 2005. SMART: The Service-Oriented Migration and Reuse Technique. Technical report, DTIC Document.
- Mani, Senthil, Sinha, Vibha S, Sukaviriya, Noi, Ramachandra, Thejaswini, 2008. Using user interface design to enhance service identification. In: *Web Services*, 2008. ICWS'08. IEEE International Conference on. IEEE, pp. 78–87.
- Marchetto, Alessandro, Ricca, Filippo, 2009. From objects to services: toward a stepwise migration approach for java applications. *Int. J. Softw. Tools Technol. Transf.* 11 (6), 427.
- Mazlami, Genc, Cito, Jürgen, Leitner, Philipp, 2017. Extraction of microservices from monolithic software architectures. In: *2017 IEEE International Conference on Web Services (ICWS)*. IEEE, pp. 524–531.



- Murtagh, Fionn, Legendre, Pierre, 2014. Ward's hierarchical agglomerative clustering method: which algorithms implement ward's criterion?. *J. Classification* 31 (3), 274–295.
- Nakamura, Masahide, Igaki, Hiroshi, Kimura, Takahiro, Matsumoto, Kenichi, 2012. Identifying services in procedural programs for migrating legacy system to service oriented architecture. *Implementation Integr. Inform. Syst. Serv. Sector* 237.
- Nakamura, Masahide, Igaki, Hiroshi, Kimura, Takahiro, Matsumoto, Ken-ichi, 2009. Extracting service candidates from procedural programs based on process dependency analysis. In: *Services Computing Conference, 2009. APSCC 2009. IEEE Asia-Pacific. IEEE*, pp. 484–491.
- Newman, Sam, 2015. *Building Microservices: Designing Fine-Grained Systems*. "O'Reilly Media, Inc."
- Nord, Robert L., Barbacci, Mario R., Clements, Paul, Kazman, Rick, Klein, Mark, 2003. Integrating the Architecture Tradeoff Analysis Method (ATAM) with the cost benefit analysis method (CBAM). Technical report, Carnegie-Mellon Univ Pittsburgh Pa Software Engineering Inst.
- Nord, R.L., et al., 2004. Integrating the Quality Attribute Workshop (QAW) and the Attribute-Driven Design (ADD) Method. Inf. Téc. Technical report, CMU/SEI-2004-TN-017, Software Engineering Institute–Carnegie Mellon ....
- Nunes, Luís, Santos, Nuno, Silva, António Rito, 2019. From a monolith to a microservices architecture: An approach based on transactional contexts. In: *European Conference on Software Architecture*. Springer, pp. 37–52.
- Pereplechikov, Mikhail, Ryan, Caspar, Frampton, Keith, Tari, Zahir, 2007. Coupling metrics for predicting maintainability in service-oriented designs. In: *2007 Australian Software Engineering Conference (ASWEC'07). IEEE*, pp. 329–340.
- Rodriguez, Juan Manuel, Crasso, Marco, Mateos, Cristian, Zunino, Alejandro, Campo, Marcelo, 2013. Bottom-up and top-down cobol system migration to web services. *IEEE Internet Comput.* 17 (2), 44–51.
- Rodríguez-Echeverría, Roberto, Maclás, Fernando, Pavón, Víctor M., Conejero, José M., Sánchez-Figueroa, Fernando, 2014. Generating a REST service layer from a legacy system. In: *Information System Development*. Springer, pp. 433–444.
- Saha, Diptikalyan, 2015. Service mining from legacy database applications. In: *Web Services (ICWS), 2015 IEEE International Conference on. IEEE*, pp. 448–455.
- Schmuller, Joseph, 2004. *Sams Teach Yourself UML in 24 hours*. Sams publishing.
- Selmadji, Anfel, Seriai, Abdelhak-Djamel, Bouziane, Hinde Lilia, Dony, Christophe, Mahamane, Rahina Oumarou, 2018. Re-architecting OO software into microservices. In: *European Conference on Service-Oriented and Cloud Computing*. Springer, pp. 65–73.
- Shatnawi, Anas, Shatnawi, Hudhaifa, Saied, Mohamed Aymen, Shara, Zakarea Al, Sahrroui, Houari, Seriai, Abdelhak, 2018. Identifying software components from object-oriented APIs based on dynamic analysis. In: *Proceedings of the 26th Conference on Program Comprehension. ACM*, pp. 189–199.
- Sindhgatta, Renuka, Sengupta, Bikram, Ponnalagu, Karthikeyan, 2009. Measuring the quality of service oriented design. In: *Service-Oriented Computing*. Springer, pp. 485–499.
- Sneed, Harry M., 2006. Integrating legacy software into a service oriented architecture. In: *Software Maintenance and Reengineering, 2006. CSMR 2006. Proceedings of the 10th European Conference on. IEEE*, pp. 11–pp.
- Sneed, Harry, 2007. Migrating to Web services: A research framework, In: *Proceedings of the International*.
- Sneed, Harry M., Verhoef, Chris, Sneed, Stephan H., 2013. Reusing existing object-oriented code as web services in a SOA. In: *Maintenance and Evolution of Service-Oriented and Cloud-Based Systems (MESOCA), 2013 IEEE 7th International Symposium on the. IEEE*, pp. 31–39.
- Sosa, Encarna, Clemente, Pedro J., Conejero, José M., Rodríguez-Echeverría, Roberto, 2013. A model-driven process to modernize legacy web applications based on service oriented architectures. In: *2013 15th IEEE International Symposium on Web Systems Evolution (WSE). IEEE*, pp. 61–70.
- Sosa-Sánchez, Encarna, Clemente, Pedro J., Sánchez-Cabrera, Miguel, Conejero, José M., Rodríguez-Echeverría, Roberto, Sánchez-Figueroa, Fernando, 2014. Service discovery using a semantic algorithm in a SOA modernization process from legacy web applications. In: *Services (SERVICES), 2014 IEEE World Congress on. IEEE*, pp. 470–477.
- Souza, Eric, Moreira, Ana, De Faveri, Cristiano, 2017. An approach to align business and IT perspectives during the soa services identification. In: *2017 17th International Conference on Computational Science and Its Applications (ICCSA). IEEE*, pp. 1–7.
- Taei Zadeh, Ali, Mukhtar, Muriati, Sahran, Shahnorbanun, Khabbazi, Mahmood Reza, 2012. A systematic input selection for service identification in SMEs. *J. Appl. Sci.* 12 (12), 1232–1244.
- Taibi, Davide, Systä, Kari, 2019. From monolithic systems to microservices: A decomposition framework based on process mining. In: *8th International Conference on Cloud Computing and Services Science, CLOSER*.
- Tyszbrowicz, Shmuel, Heinrich, Robert, Liu, Bo, Liu, Zhiming, 2018. Identifying microservices using functional decomposition. In: *International Symposium on Dependable Software Engineering: Theories, Tools, and Applications*. Springer, pp. 50–65.
- Vale, Tassio, Figueiredo, Gustavo Bittencourt, de Almeida, Eduardo Santana, de Lemos Meira, Silvio Romero, 2012. A study on service identification methods for software product lines. In: *Proceedings of the 16th International Software Product Line Conference-Volume 2. ACM*, pp. 156–163.
- Vemulapalli, Anisha, Subramanian, Nary, 2009. Transforming functional requirements from UML into BPEL to efficiently develop SOA-based systems. In: *OTM Confederated International Conferences on the Move to Meaningful Internet Systems*. Springer, pp. 337–349.
- Wagner, Christian, 2014. *Model-Driven Software Migration: A Methodology: Reengineering, Recovery and Modernization of Legacy Systems*. Springer Science & Business Media.
- Weske, Mathias, 2012. Business process management architectures. In: *Business Process Management*. Springer, pp. 333–371.
- Wille, Rudolf, 1982. Restructuring lattice theory: an approach based on hierarchies of concepts. In: *Ordered Sets*. Springer, pp. 445–470.
- Wohlin, Claes, 2014. Guidelines for snowballing in systematic literature studies and a replication in software engineering. In: *Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering. ACM*, p. 38.
- Xu, Rui, Wunsch, D., 2005. Survey of clustering algorithms. *IEEE Trans. Neural Netw.* 16 (3), 645–678. <http://dx.doi.org/10.1109/TNN.2005.845141>.
- Zadeh, Ali Taei, Mukhtar, Muriati, Sahran, Shahnorbanun, Khabbazi, MR, 2012. A systematic input selection for service identification in SMEs. *J. Appl. Sci.* 12 (12), 1232.
- Zhang, Zhuopeng, Liu, Ruimin, Yang, Hongji, 2005. Service identification and packaging in service oriented reengineering. In: *SEKE*, 5, pp. 620–625.
- Zhang, Zhuopeng, Yang, Hongji, 2004. Incubating services in legacy systems for architectural migration. In: *11th Asia-Pacific Software Engineering Conference, 2004. IEEE*, pp. 196–203.
- Zhang, Zhuopeng, Yang, Hongji, Chu, William C., 2006. Extracting reusable object-oriented legacy code segments with combined formal concept analysis and slicing techniques for service integration. In: *2006 Sixth International Conference on Quality Software (QSIC'06). IEEE*, pp. 385–392.
- Zhao, Shuxin, Chang, Elizabeth, Dillon, Tharam, 2008. Knowledge extraction from web-based application source code: An approach to database reverse engineering for ontology development. In: *2008 IEEE International Conference on Information Reuse and Integration. IEEE*, pp. 153–159.
- Zhao, Yun, Si, Huayuan, Ni, Yulin, Qi, Hengnian, 2009. A service-oriented analysis and design approach based on data flow diagram. In: *International Conference on Computational Intelligence and Software Engineering CISE 2009. IEEE*, pp. 1–5.

**Manel Abdellatif**, is a Ph.D. student at Polytechnique Montreal. Her research focus is on service identification to support the migration of legacy object-oriented software systems to SOA. She holds a Master degree from Ecole de Technologie Supérieure in Canada and an engineering degree from Ecole Nationale d'Ingénieurs de Tunis.

Her research interests focus on automated software engineering including: reengineering and reverse engineering, empirical software engineering, software reuse, maintenance, evolution, program analysis and refactoring, data mining for software engineering.

**Anas Shatnawi** is a postdoc researcher at LIP6 lab of Sorbonne University, Paris, France. Before joining Sorbonne University, he was a postdoctoral researcher at LTA lab of Università degli Studi di Milano-Bicocca (UniMib) and at LATECE of Université du Québec à Montréal (UQAM), Canada. He obtained his Ph.D. degree in Computer Science from LIRMM at the Université de Montpellier, France. He worked under the supervision of Prof. Abdelhak Seriai, University of Montpellier and Prof. Houari Sahraoui, University of Montreal, Canada. He received his Master and bachelor's degrees in computer science respectively from Jordan University of Science and Technology and Irbid National University, Jordan. He also worked in industry for more than 4 years as software developer at Maysalward R&D, Jordan.

**Hafedh Mili** is a Professor at Université du Québec à Montréal. He obtained his Ph.D. in Computer Science (Artificial Intelligence) from the George Washington University in 1988.

His research interest focuses on software development from business models, representation and classification of business processes, representation and execution of flexible processes, pattern making, software re-engineering of legacy applications, customer experience management, the internet of things, and intelligent processes and infrastructure.

**Naouel Moha** is a Professor at École de Technologie Supérieure, Montreal, Canada since September 2020. She was before professor at Université du Québec à Montréal (UQAM), Canada. She received the master's degree in computer science from the University of Joseph Fourier, Grenoble, in 2002. She also received the Ph.D. degree, in 2008, from the University of Montreal (under Professor Yann-Gaël Guéhéneuc's supervision) and the University of Lille (under the supervision of Professor Laurence Duchien and Anne-Françoise Le Meur). Her research interests include software quality and evolution, in particular refactoring and the identification of patterns and antipatterns in different systems including service-based systems and mobile apps.

**Ghizlane El Boussaidi** is a Professor at École de Technologie Supérieure, Montreal, Canada since 2010. She obtained her Ph.D. in Software engineering from Université de Montréal, Canada. Her research interests include software architecture, design patterns and architectural styles, architecture recovery, modernization of legacy systems, model-driven engineering, domain-specific modeling, and the design and certification of safety critical systems. She also has over 15 years of industrial experience during which she contributed to the development and the implementation of various software systems.

**Geoffrey Hecht** is a Postdoctoral Research Fellow at Université du Québec à Montréal (UQAM) in Canada and at Departamento de Ciencias de la Computación, FCFM, Universidad de Chile. His research interests include mobile computing, cloud computing and legacy-to-SOA migration.

**Jean Privat** is a Professor at Université du Québec à Montréal (UQAM) in Canada. His research interests include specification and implementation of object languages Compilers, code generators, and virtual machines.

**Yann-Gaël Guéhéneuc** is a Professor at Concordia University where he leads the Ptidej team on evaluating and enhancing the quality of object-oriented programs by promoting the use of patterns, at the language-, design-, or architectural-levels. In 2009, he was awarded the NSERC Research Chair Tier II on Software Patterns and Patterns of Software. He holds a Ph.D. in software engineering from University of Nantes, France (under Professor Pierre Cointe's supervision) since 2003 and an Engineering Diploma from École des Mines of Nantes since 1998. His Ph.D. thesis was funded by Object Technology International, Inc. (now IBM OTI Labs.), where he worked in 1999 and 2000. His research interests are program understanding and program quality during development and maintenance, through the use and the identification of recurring patterns. He was the first to use explanation-based constraint programming in the context of software engineering to identify occurrences of patterns. He is interested also in empirical software engineering; he uses eye-trackers to understand and to develop theories about program comprehension. He has published many papers in international conferences and journals.