



# Cloud CBR: Towards Microservices Oriented Case-Based Reasoning

Ikechukwu Nkisi-Orji<sup>1</sup>  , Nirmalie Wiratunga<sup>1</sup> ,  
Chamath Palihawadana<sup>1</sup> , Juan A. Recio-García<sup>2</sup> , and David Corsar<sup>1</sup> 

<sup>1</sup> School of Computing Science and Digital Media,  
Robert Gordon University, Aberdeen AB10 7GJ, Scotland, UK  
{i.o.nkisi-orji,n.wiratunga,c.palihawadana,d.corsar1}@rgu.ac.uk

<sup>2</sup> Department of Software Engineering and Artificial Intelligence,  
Universidad Complutense de Madrid, Madrid, Spain  
jareciog@fdi.ucm.es

**Abstract.** CBR applications have been deployed in a wide range of sectors, from pharmaceuticals; to defence and aerospace to IoT and transportation, to poetry and music generation; for example. However, a majority of these have been built using monolithic architectures which impose size and complexity constraints. As such these applications have a barrier to adopting new technologies and remain prohibitively expensive in both time and cost because changes in frameworks or languages affect the application directly. To address this challenge, we introduce a distributed and highly scalable generic CBR system, CLOUD, which is based on a microservices architecture. This splits the application into a set of smaller, interconnected services that scale to meet varying demands. Experimental results show that our CLOUD implementation retrieves cases at a fairly consistent rate as the casebase grows by several orders of magnitude and was over 3,700 times faster than a comparable monolithic CBR system when retrieving from half a million cases. Microservices are cloud-native architectures and with the rapid increase in cloud-computing adoption, it is timely for the CBR community to have access to such a framework.

**Keywords:** Cloud CBR · Microservices · Elasticsearch · CBR framework

## 1 Introduction

Several case-based reasoning (CBR) development frameworks and toolkits have been introduced to the CBR community [13–15]. These have been extended for recommender systems [8] and textual CBR [12] and more recently for self-management systems [1]. However many of these CBR systems are mostly implemented with monolithic architectures such as desktop standalone applications,

---

J. A. Recio-García—Supported by the Spanish Committee of Economy and Competitiveness (TIN2017-87330-R).

© Springer Nature Switzerland AG 2020

I. Watson and R. Weber (Eds.): ICCBR 2020, LNAI 12311, pp. 129–143, 2020.

[https://doi.org/10.1007/978-3-030-58342-2\\_9](https://doi.org/10.1007/978-3-030-58342-2_9)

with heavy demands due to siloed in-memory batch processing. This is not compatible with recent software development trends, which are increasingly using REST APIs<sup>1</sup> for communication with cloud computing platforms.

Cloud computing is a term used to describe the use of remote hardware and software to deliver on-demand computing services through a network (usually the Internet). In the past, applications or programs were run from software downloaded on to a physical computer or server. In contrast cloud computing lets users access these applications through the internet. Implementing software applications in the cloud offer several benefits which include efficient/cost reduction, scalability, mobility, and disaster recovery. Distribution of CBR applications and cases enables, MapReduce type algorithms to exploit the parallelism opportunity that is to be had with pair-wise similarity computations [19]. Interestingly, CBR has also been applied to support cloud provisioning, whereby similar Amazon Web Services (AWS)<sup>2</sup> configurations are recommended given a characterisation of a user’s compute task [9]. This helps the user to make decisions about the types of cloud services for the given task. But having to monitor resource utilisation and change service requirements accordingly is a challenge which in turn has paved the way for microservice based architectures.

A CBR framework using a microservice based architecture provides (amongst other things) flexibility in both the technology being used (e.g., programming language) as well as dynamic scalability that can adapt to user application demands (e.g., spikes in casebase querying, seasonal effects). This is because, individual microservices are independently scaled and developed such that the overall system architecture is a scalable distributed application [6]. Importantly, the computation of services are stateless since they are automatically provisioned only when needed and then stopped when no longer required. This is particularly advantageous to CBR in situations where there is in-memory demand due to its inherent nature of being a lazy learner.

In this paper we discuss how the CBR cycle can be organised into multiple microservices and how service discovery is facilitated between these independent components using rest communications. A microservice is considered efficient when the system is loosely coupled and highly cohesive [10]. Identifying which functionalities within the CBR cycle should be decoupled and organising them into microservices is a key design challenge that we address in this paper. We do this by introducing, CLOUD<sup>3</sup>, a generic open-source CBR cloud-based microservice framework, and make the following key contributions:

- create a novel design using the microservice paradigm for CBR;
- introduce, CLOUD, an extensible open source microservice CBR framework<sup>4</sup>;
- evaluate the scalability of the retrieval phase on a recommender task; and

---

<sup>1</sup> An architectural style and approach for communication based on representational state transfer (REST) that is often used in web services development.

<sup>2</sup> <https://aws.amazon.com>.

<sup>3</sup> Cloud is “Cloud” in Scottish dialect.

<sup>4</sup> CLOUD CBR repository: <https://github.com/RGU-Computing/cloud>.

- identify areas of future development that are essential for the sustainability of CLOUD CBR.

Rest of the paper is organised as follows; in Sect. 2 we discuss existing frameworks, jCOLIBRI and *myCBR*. The design paradigm appears in Sect. 3 and the CLOUD implementation is discussed in Sect. 4. Results from a scalability experiment with half a million cases is presented in Sect. 5 followed by conclusions and future directions in Sect. 6.

## 2 Related Work in CBR Development Architectures

There are two well-established open-source frameworks for building CBR applications: *myCBR* and COLIBRI, though they follow different approaches and support different phases of the CBR application development.

*myCBR*<sup>5</sup> has been a tool for researchers and practitioners over the last ten years [16]. This framework is focused on the developing of a knowledge model for representing cases and computing similarity through the myCBR-workbench tool [2]. This knowledge model can be instantiated through the building blocks and functionality provided by the myCBR-SDK, that is a Java library following a classical monolithic software architecture. However, their authors have recently presented the *myCBR* Rest API which exposes the functionality of both myCBR-SDK and myCBR-workbench through a RESTful API [1]. Instead of forcing users to integrate their *myCBR* systems into a Java environment, this novel API enables users to model a CBR system using myCBR's workbench and then deploying the application as a web service. The goal is to make it easier to build, test, compare and deploy CBR applications.

COLIBRI, on the other hand, is focused on the development of a wide range of CBR applications [11]. As a platform, COLIBRI offers a well defined architecture for designing CBR systems, a reference implementation of that architecture: the jCOLIBRI framework [13], and several development tools that aid users in the implementation and sharing of new CBR systems and components. These tools have been integrated in the COLIBRI *Studio* development environment [14]. Both tools make up the COLIBRI platform following a two layer architecture. jCOLIBRI is the white-box layer of the architecture: a framework for developing CBR applications in Java. This framework represents the bottom layer of the platform. It includes most of the code required to implement a wide collection of CBR systems: Standard, Textual, Knowledge-Intensive, Data-Intensive, Recommender Systems, and Distributed CBR applications. It also includes evaluation, maintenance and casebase visualisation tools. All this functionality has established jCOLIBRI as a reference CBR framework with more than 35K downloads<sup>6</sup>. However, jCOLIBRI still follows the same monolithic Java architecture like *myCBR* and is not suitable for modern web environments.

<sup>5</sup> <http://mycbr-project.org>.

<sup>6</sup> <http://gaia.fdi.ucm.es/research/colibri/jcolibri>.

The need for both these platforms to evolve into web services architecture is clear. However, there are different approaches to implement this evolution. *myCBR* proposes wrapping its existing java components as web services. It is a straightforward option but has several drawbacks. Mostly, the wrapping of the existing java components does not allow to take advantage of the capabilities of cloud architectures regarding availability or scalability. The alternative option is to create a cloud-based CBR framework from scratch in order to exploit the features of modern cloud architectures. This is the option adopted by CLOUD, that can be considered as a re-implementation of the functionalities provided by the jCOLIBRI and *myCBR* frameworks, but instead of wrapping its existing java components, it redesigns entirely the CBR architecture for the cloud. In this manner, CLOUD adopts the CBR architecture defined in COLIBRI based on a pre/post-CBR-cycle to load/release required resources. CLOUD also reproduces the case structure representation based on a composite pattern, and the similarity computation through global/local similarity functions that both jCOLIBRI and *myCBR* implement.

In summary, our goal is to create a cloud architecture that is able to provide the same functionalities using familiar methods currently being used in jCOLIBRI and, thereafter, further integrate existing web services found in *myCBR*. As we will present in the following section, CLOUD re-implements jCOLIBRI's methods using modern web services technologies such as Elasticsearch or JSON-based communications that extend the existing capabilities of the framework regarding flexibility and data-intensive processing.

### 3 Microservices Design Paradigm for CBR

A microservice is an independent process which can carry out specific tasks in isolation [6]. These should be deployed, tested and scaled independently for a single functional responsibility; such as similarity, ranking, casebase editing, etc. Key to this architecture are the concept of serverless functions also referred to as *Function-as-a-Service* (FaaS) [3] - logic that is split into small code snippets and executed in a managed compute service. Well known examples include AWS Lambda and Google Cloud Functions<sup>7</sup>.

#### 3.1 Cloud Architecture

Fig. 1 shows a high-level overview of the system's design consisting of 3 main components: REST API; serverless functions (compute service); and data service. The core CBR tasks – retrieve, reuse, revise, retain – are implemented as serverless computing functions. Functions can interact with external applications (e.g., a dashboard) and internally with other functions through REST APIs. Decomposition of the CBR cycle into smaller functions provides flexibility to introduce similarity functions and deploy them independently. Such functions will also include relevant knowledge container provisions. The post-cycle

<sup>7</sup> <http://aws.amazon.com/lambda> and <http://cloud.google.com/functions>.

or maintenance tasks, like forgetting cases or recomputing footprint cases can be confined to the Retain service. The data service is used as the casebase which allows the serverless functions to query and retrieve. Data sources and connectors forming the pre-cycle communicate with the casebase once they are synced with the data service. Data sources can either be external or within the cloud platform which gives flexibility for the community to use existing data sources. An important distinction here with the pre-cycle is that it remains lean (as compared to jCOLIBRI, or *myCBR*); in that it does not involve loading cases into memory once cases are made persistent.

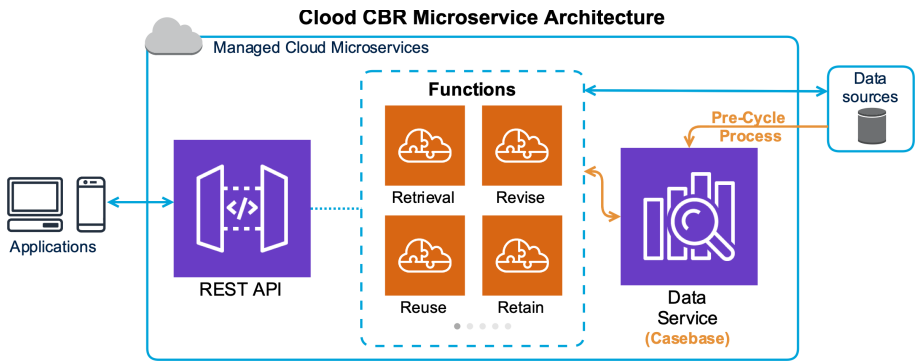


Fig. 1. Proposed CLOUD CBR architecture diagram

### 3.2 The Casebase

Popular CBR systems like jCOLIBRI keep the casebase in memory during operation. An in-memory casebase guarantees speed when interacting with the casebase but will incur massive costs to scale up for big data. Also, using the CBR system in a distributed manner can be problematic with in-memory casebase as memory is an expensive resource even on the cloud. In the serverless architecture, we maintain the casebase in the data service. The data service is a NoSQL full-text distributed search engine for all types of data. Elasticsearch and Solr are popular examples of such distributed, scalable open-source search tools for textual, numerical, geospatial, structured, and unstructured data. These tools provide a significant improvement regarding the representation of cases in previous CBR frameworks, because the case structure does not need to be fixed. Therefore, the cases in the casebase can have different attributes, and similarity metrics are applied according to each particular data types. Moreover, as these search tools are built on Apache Lucene, they are extendable, allowing users to write custom similarity metric scripts against a data index. Accordingly, the

type of operations that would normally occur in-memory can be done in the data store index which is usually file-based<sup>8</sup>.

### 3.3 Local Similarity

A subset of the serverless functions for the retrieve phase are used to generate similarity scripts to measure local similarity. These metric functions perform retrieval from the casebase at the attribute level. Each generated similarity function script depends on the data type of the attribute. Supported data types include string, numeric, boolean, date and object. Some similarity metrics, such as metrics to retrieve exact matches, are in-built in several distributed search engines that can be used for the data service. A suitable data service should enable the implementation of custom similarity metrics functions to support other local similarity functions that are used for CBR retrieval in the jCOLIBRI and *myCBR* frameworks.

### 3.4 Global Similarity

The global similarity function which aggregates local similarities determines the order in which cases are retrieved from the casebase and their ranking. Both a weighted and non-weighted form can be used to identify the nearest neighbours and is managed directly by the data service. Each local similarity function script is executed in the data service, in response to a single query, to obtain the global similarity as a sum. Custom scripts can be created as needed to vary the weights associated with different attributes. These weights can be dynamically modified for each retrieval task or alternatively remain static for all queries. The latter corresponds to learning an attribute weighting scheme that is used unchanged with every casebase query; whilst the former provides the opportunity to change attribute weights to suit the query context. The default global aggregation can be replaced with a custom aggregation script; whilst this does offer greater flexibility it will also incur greater computing memory when working with medium to large casebases since all the cases that are returned by the local functions will be held in memory (as with the monolithic organisation of jCOLIBRI and *myCBR*).

### 3.5 Implication for CBR Cycle

The major improvement over the architectures used by jCOLIBRI and *myCBR* is the lack of a two-layer persistence strategy. In previous frameworks there is a need to load cases into memory from a persistence media such as a database, text file, etc. However, the use of the CLOUD data service allows to manage cases directly from its internal data index.

Absence of the two-layer persistence strategy, has an immediate impact on the application structure because unlike previously where a *pre-cycle* step was

<sup>8</sup> Elasticsearch index store <http://www.elastic.co/guide/en/elasticsearch/reference/7.6/index-modules-store.html> (accessed May 14, 2020).

needed prior to the CBR cycle itself for loading cases into memory, this is no longer required. However CLOUD maintain the possibility of executing a pre-cycle (or its complementary post-cycle) in order to perform additional pre/post-processing of the data, if the CBR system requires it.

Another significant benefit of cloud-based technologies is concurrency, which directly creates the opportunity to execute CBR processes in parallel. This feature is quite limited in current frameworks and is also very relevant in order to parallelise time-consuming algorithms such as kNN or noise removal methods such as BBNR (Blame-based noise reduction), CRR (Conservative Redundancy Removal), RENN (Repeated Edited Nearest Neighbour), RC (Relative Cover), or ICF (Iterative Case Filtering) [5].

## 4 Cloud CBR System

CLOUD is implemented using python functions following the design paradigm presented in Sect. 3. These functions run on Amazon Web Services (AWS) Lambda, which is the serverless event-driven computing service of AWS. The casebase uses the AWS ES service and the client application is implemented with JavaScript and HTML using the AngularJS framework<sup>9</sup>. Using a test application provided by jCOLIBRI<sup>10</sup> we describe the CLOUD implementation (see Fig. 2) and discuss how CBR functionality is achieved with cloud capabilities. Services that are not core to CBR operations include Cognito which is used for authenticated access to the system and Cloud Watch which is used to collect and monitor event logs.

### 4.1 Casebase Using Elasticsearch

Elasticsearch (ES) is an open-source highly distributed and horizontally scalable full-text search engine with various capabilities built on Apache Lucene [7]. ES uses RESTful interfaces to manipulate its schema-free JSON document store and performs searches at very high speeds maintaining an index that is about 20% the size of the indexed documents [18]. Compared to traditional database management systems, the ES “index” is somewhat like the database table as queries are executed against the index. While there are several schema-free databases with search capability to choose from, we choose ES as the data service in our implementation because of its popularity and close integration with existing cloud service providers.

Although it is “schema-free”, ES internally generates a schema based on the field (attributes/columns) values of documents to be indexed. Relying on an ES-generated schema can be problematic in some cases. For example, a field for storing alphanumeric values can be designated as numeric by ES if the first documents to be indexed have numeric values only for that field. In order to

<sup>9</sup> <http://angularjs.org>.

<sup>10</sup> <http://gaia.fdi.ucm.es/research/colibri/jcolibri/doc/apidocs/es/ucm/fdi/gaia/jcolibri/test/test1/package-summary.html> (accessed May 14, 2020).

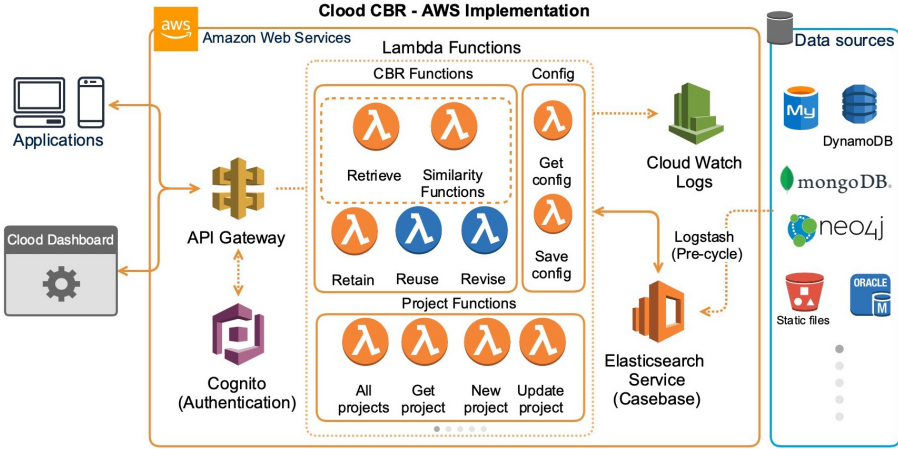


Fig. 2. CLOUD CBR implementation on AWS

avoid undesirable field properties, we create an explicit mapping which indicates the data type to be stored for each field in the casebase. The ES index “mapping” is comparable to the database schema as it describes the fields (columns) in the JSON documents along with their data types.

An explicit index mapping supports the specification of how a field’s values should be indexed and the local similarity metric to be used for retrieving the values of that field. Where possible, we delay specifying the local similarity function for a field until retrieval time for greater flexibility. This is because the index specification for a field cannot be modified once data is added to the index. With query script similarity functions supplied at retrieval time, the method of retrieval can be varied without having to modify the underlying index mapping. Introducing a new attribute to an existing casebase can be done by extending the index mappings with the new field. The structure of cases that do not have values for newly created fields will remain unchanged. CLOUD’s serverless functions interact with ES by HTTP requests and responses using a python Elasticsearch client, `elasticsearch-py`<sup>11</sup>. The casebase is a separate service which can be hosted anywhere with exposed API end-points further highlighting the distributed nature of CLOUD.

### 4.2 Cloud Similarity Functions

Table 1 shows the local similarity metric functions that are currently implemented on CLOUD, reproducing some relevant functions available in `jCOLIBRI` and `myCBR`. Although several similarity metrics are currently missing in CLOUD, the goal here is to demonstrate the potential of the framework and to encourage code contributions in the future. Each similarity metric is

<sup>11</sup> <http://elasticsearch-py.readthedocs.io/en/master> (accessed May 14, 2020).



implemented as a python function which generates and returns a *Painless* script<sup>12</sup> that can be executed on ES during retrieval operations. Painless is the scripting language that is specifically designed for writing inline and stored scripts on ES. Generated scripts for the local similarity of each case attribute are combined into a single multi-match query script at retrieval.

**Table 1.** CLOUD’s local similarity metrics

Data type	Similarity metric	Description
All	Equal	Similarity based on exact match
String	EqualIgnoreCase	Case-insensitive string matching
	BM25	TF-IDF similarity with TF normalisation based on Okapi BM25 ranking function
	Semantic USE	Similarity based on the similarity of vector representations
Numeric	Interval	Similarity of two numbers inside an interval
	INRECA	Similarity following the INRECA More is Better and Less is Better
	McSherry	Similarity following the McSherry More is Better and Less is Better
Enum	EnumDistance	Similarity of values based on their relative positions within an enumeration
Date	ClosestDate	Similarity depending on the extent two dates are to each other

McSherry, INRECA, Interval and EnumDistance are re-implementations of local similarity metrics found in jCOLIBRI. For textual CBR, we specifically implemented the Semantic local similarity metric (Semantic USE) for text content, using the Universal Sentence Encoder (USE) which embeds texts in a dense vector space of 512 dimensions [4]. This vector representation is generated using a lite version of USE based on the Transformer architecture<sup>13</sup> [17] and is stored as a dense vector field on ES. Textual retrieval follows the same process of generating the vector representation of a query string. Afterwards, the Semantic USE local similarity function measures the cosine similarity between query vectors and documents’ vectors to identify the most semantically similar content.

### 4.3 REST API

REST APIs are stateless in that the API server does not remember the state of its clients and every call to an end-point is independent of other calls. REST

<sup>12</sup> <http://www.elastic.co/guide/en/elasticsearch/painless/master/painless-guide.html> (accessed May 14, 2020).

<sup>13</sup> <http://github.com/tensorflow/tfjs-models/tree/master/universal-sentence-encoder> (accessed May 14, 2020).

API uses existing protocols such as HTTP for Web APIs. As a result, client applications do not need additional software to use the service. REST improves portability to different types of platforms since all interactions are completed through universally understood interfaces. With CLOUD, each REST API end-point is a serverless function. The replication of an end-point and the resources allocated to it vary to meet changing demands without affecting the other end-points. REST APIs are created and published using the API Gateway (see Fig. 2) and Table 2 summarises the major REST API end-points of CLOUD.

CLOUD is able to concurrently manage multiple CBR applications (use-cases) referred to as “project” in Table 2. The system’s capabilities can be easily extended by introducing new serverless functions (e.g., similarity functions, reuse functions, revise functions). Functions that will become part of the REST API are specified in a YAML file along with their access protocols.

**Table 2.** CLOUD’s REST API end-points

End-point	Request method	Description
/project	HTTP GET	Retrieves all the CBR projects
/project/{id}	HTTP GET	Retrieves a specific CBR project with specified id
/project	HTTP POST	Creates a new CBR project. The details of the project are included as a JSON object in the request body
/project/{id}	HTTP PUT	Updates the details of a CBR project. Modifications are included as a JSON object in the request body
/project/{id}	HTTP DELETE	Removes a CBR project with specified id
/case/{id}/list	HTTP POST	Bulk addition of cases to the casebase of the project with specified id. Cases are included in the request body as an array of objects
/retrieve	HTTP POST	Performs the retrieve task
/retain	HTTP POST	Performs the retain task
/config	HTTP GET	Retrieves the system configuration
/config	HTTP POST	Adds or updates the system configuration

#### 4.4 Cloud CBR Dashboard

Client applications can perform CBR operations through the RESTful API end-points of CLOUD. The CLOUD CBR client application is a light-weight HTML and JavaScript implementation that is able to manage multiple CBR projects through API calls. Figure 3 shows the interface for specifying the attributes of a project’s casebase. CLOUD system’s configuration provides guidance on allowed

operations when specifying attributes. For example, it indicates that the Interval local similarity metric only applies to numeric attributes. Once the attribute specifications are completed, CLOUD generates an index mapping for the case representation on ES.

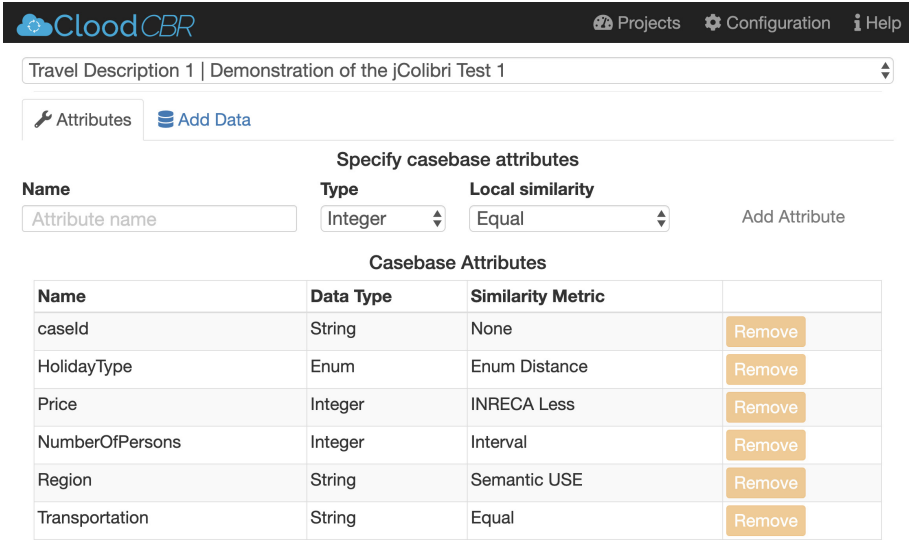


Fig. 3. Specifying attributes for a casebase.

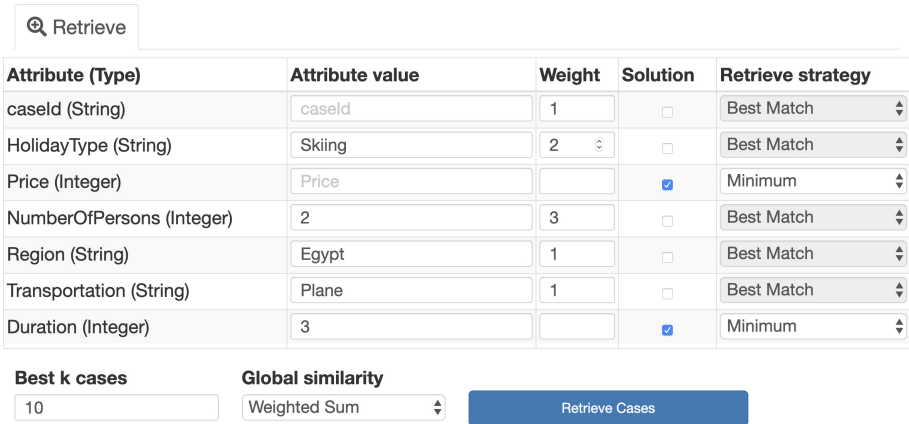


Fig. 4. Retrieve stage query specification.

Logstash is an open-source data processing pipeline from the ES stack for ingesting data into ES<sup>14</sup>. Using Logstash, cases can be added to a CLOUD's

<sup>14</sup> <http://www.elastic.co/guide/en/logstash/current/input-plugins.html> (accessed May 14, 2020).

casebase from multiple data sources including files (e.g., CSV file), databases with JDBC interfaces (e.g., MySQL), and NoSQL databases (e.g., MongoDB). However, we also include a file upload utility for adding cases from CSV files through a RESTful end-point and which should be sufficient for file sizes that will not overwhelm the Web browser.

The retrieve operation begins with specifying some attribute values along with weights for aggregating the local similarity measures. Attributes with known values become part of the problem space while attributes with unknown values form the solution space. Furthermore, a retrieve strategy can be specified per attribute as shown on the user interface in Fig. 4. For example, the Best match can be retrieved for one attribute while the Mean of the  $k$  best matches retrieved for another attribute. The  $k$  nearest neighbours to retrieve and the global similarity method can also be specified at the retrieve phase.

A reuse interface displays the retrieval results for reuse. The recommended case (candidate solution) mixes the user-supplied attribute values with the retrieved values for unknown attribute values. The  $k$  most similar cases to the query case are also presented for possible reuse. The reuse button against a retrieved case is used to make it the recommended case. The recommended case can be revised by adjusting it as required. Afterwards, the case can be retained by adding to the casebase.

## 5 Evaluation

A scalability test is conducted to evaluate CLOUD based CBR application, to examine how resource demands both on the casebase and the serverless CBR functions are met. We expect a fairly consistent compute performance for different CBR tasks across different project sizes (compared to a jCOLIBRI application). We focus on case retrieval for evaluation since it is the most commonly performed and time-consuming stage of the CBR cycle.

### 5.1 Experimental Setup and Dataset

Six CBR projects of increasing casebase sizes (10,  $10^2$ ,  $10^3$ ,  $10^4$ ,  $10^5$ , and 540,394) were created from a used cars dataset<sup>15</sup> (1.35 GB CSV file), and case retrieval efficiency compared with CLOUD and jCOLIBRI. A case has 25 attributes<sup>16</sup> describing the physical attributes of a car (e.g., colour), identification attributes (e.g., vehicle identification number), and location attributes (e.g., region, state, coordinates), and the listing price.

In the comparative study, 10 nearest neighbours (NN) are retrieved with Equal similarities (Table 1) using the following query.

<sup>15</sup> <http://www.kaggle.com/austinreese/craigslist-carstrucks-data/data> (accessed February 25, 2020).

<sup>16</sup> Dataset attributes are `id`, `url`, `region`, `region_url`, `price`, `year`, `manufacturer`, `model`, `condition`, `cylinders`, `fuel`, `odometer`, `title_status`, `transmission`, `vin`, `drive`, `size`, `type`, `paint_color`, `image_url`, `description`, `county`, `state`, `lat`, `long`.

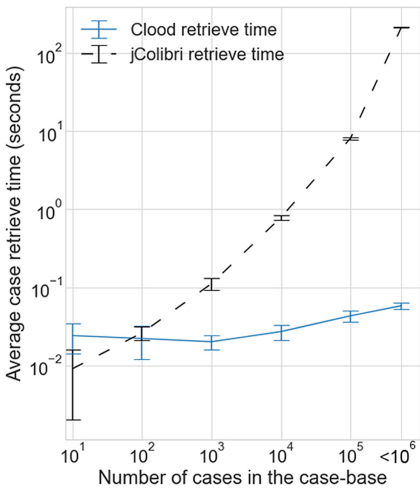
```
{ 'year': '2017', 'manufacturer': 'ford', 'model': 'focus',
  'condition': 'good', 'fuel': 'gas', 'title_status': 'clean',
  'transmission': 'automatic', 'drive': '4wd',
  'size': 'compact', 'paint_color': 'grey' }
```

Time taken by the Retrieval function (Retrieve time) is recorded which for CLOUD, consists of: the time spent to dynamically generate a query using the appropriate similarity functions for the query case, retrieve the 10 NN of the query case from the casebase, generate a recommended case for reuse using specified reuse strategy, and generate a response through the API. We do not include the time lapse between the client application and the API endpoints as that is very dependent on the network connection speed and client's platform resources. For jCOLIBRI, Retrieve time is measured in the cycle phase consisting of: the time spent to retrieve the similarity configuration, perform NN scoring over the cases (in-memory), and select the 10 best cases. jCOLIBRI was run on a Windows 10 PC having 6th generation Intel core i7 processor and 16 GB RAM with 2GiB Java heap size. CLOUD uses AWS Lambda functions for its operations while the casebase was hosted on a single cluster of the AWS ES Service with 2GiB and 1 vCPU (t2.small.elasticsearch instance).

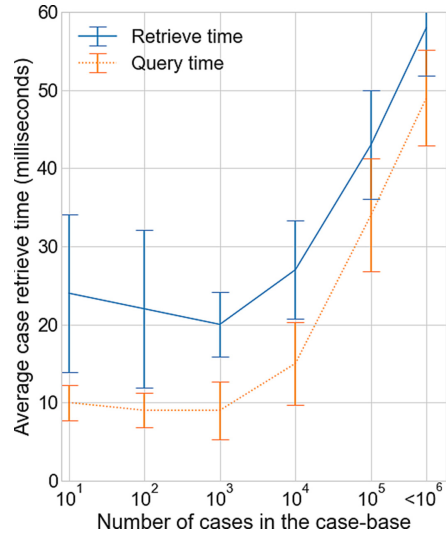
## 5.2 Results and Discussion

Figure 5 shows the average case retrieve times for CLOUD and jCOLIBRI on log scales with standard deviations as error bars. jCOLIBRI was marginally faster on the smallest casebase (10) but the superior performance of CLOUD is apparent with increasing casebase sizes. Similar case retrieval times were obtained by both systems at about casebase size of 100 cases; however at casebase size of 1,000, CLOUD was 5.5 times faster than jCOLIBRI and at casebase size of 540,394, CLOUD was 3,737 times faster than jCOLIBRI. Close examination of CLOUD's Retrieve time spent on the ES casebase when measured separately (Query time) shows to have increased due to time spent querying the casebase (see Fig. 6). We used the smallest AWS ES instance, and we expect Query time to improve when using an ES instance with improved resources. Also, several optimisation techniques can be employed to improve Query time. In the current implementation, we apply each local similarity function to the target attribute of every case in the casebase. An improvement can optimise the querying process such that it uses filters to reduce the number of similarity computations. For example, in the query above where the 'year' must match '2017', we can apply the year limit as a filter when matching 'manufacturer' so that it only searches for 'ford' in 2017 models.

The use of cloud services typically involves usage costs. The microservice architecture with pricing per run-time keeps the costs minimal. For example, running the CloudCBR system with core services (Lambda and ES and data transfer) costs 14 USD a month. In comparison, a similar monolith system hosted



**Fig. 5.** Case retrieve times as casebase size increases. Both axes are log scales.



**Fig. 6.** CLOUD retrieve times compared to the query times. X-axis is a log scale.

on a medium-sized AWS machine (t3a.medium and 50 GB storage) costs 22 USD a month (cost estimates as of April 2020).

## 6 Conclusion

We introduced CLOUD CBR, a novel microservices-oriented CBR framework which leverages the serverless architecture for CBR operations and a distributed data storage service (Elasticsearch) for CBR knowledge persistence. Implementation of the extensible CLOUD CBR framework is an ongoing open-source project. We demonstrated the robustness of CLOUD on a CBR project of half a million cases and showed how CLOUD is scalable for different project sizes. Ongoing work for the future sustainability of CLOUD include extending support for additional similarity and data types (e.g., *myCBR*'s table similarity); and include functions for reuse and revise, casebase maintenance and visualisation. Also, overcoming the performance bottleneck of handling the intermediate results of local similarity functions in-memory, when implementing custom global similarity functions, is beneficial for extending the capabilities of CLOUD. We intend to make CLOUD a Python library to reuse the CLOUD Elasticsearch similarity functions for the community and add seamless integration for deploying on more cloud providers.

## References

1. Bach, K., Mathisen, B.M., Jaiswal, A.: Demonstrating the myCBR rest API. In: Demo Session of the 27th International Conference on CBR (2019)

2. Bach, K., Sauer, C.S., Althoff, K., Roth-Berghofer, T.: Knowledge modelling with the open source tool myCBR. In: CEUR Workshop Proceedings, KESE@ECAI, vol. 1289. CEUR-WS.org (2014)
3. Castro, P., Ishakian, V., Muthusamy, V., Slominski, A.: Serverless programming (function as a service). In: 2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS), pp. 2658–2659. IEEE (2017)
4. Cer, D., et al.: Universal sentence encoder. arXiv preprint [arXiv:1803.11175](https://arxiv.org/abs/1803.11175) (2018)
5. Cummins, L., Bridge, D.: On dataset complexity for case base maintenance. In: Ram, A., Wiratunga, N. (eds.) ICCBR 2011. LNCS (LNAI), vol. 6880, pp. 47–61. Springer, Heidelberg (2011). [https://doi.org/10.1007/978-3-642-23291-6\\_6](https://doi.org/10.1007/978-3-642-23291-6_6)
6. Dragoni, N., et al.: Microservices: yesterday, today, and tomorrow. Present and Ulterior Software Engineering, pp. 195–216. Springer, Cham (2017). [https://doi.org/10.1007/978-3-319-67425-4\\_12](https://doi.org/10.1007/978-3-319-67425-4_12)
7. Gormley, C., Tong, Z.: Elasticsearch: The Definitive Guide: A Distributed Real-time Search and Analytics Engine. O’Reilly Media Inc., Sebastopol (2015)
8. Jorro-Aragoneses, J.L., Recio-Garcia, J.A., Diaz-Agudo, B., Jiménez-Díaz, G.: Recolibry-core: a component-based framework for building recommender systems. *Knowl.-Based Syst.* **182**, 104854 (2019)
9. Minor, M., Schulte-Zurhausen, E.: Towards process-oriented cloud management with case-based reasoning. In: Lamontagne, L., Plaza, E. (eds.) ICCBR 2014. LNCS (LNAI), vol. 8765, pp. 305–314. Springer, Cham (2014). [https://doi.org/10.1007/978-3-319-11209-1\\_22](https://doi.org/10.1007/978-3-319-11209-1_22)
10. Pahl, C., Jamshidi, P.: Microservices: a systematic mapping study. In: CLOSER (1), pp. 137–146 (2016)
11. Recio-García, J.A., Díaz-Agudo, B., González-Calero, P.A.: The COLIBRI platform: tools, features and working examples. In: Montani, S., Jain, L. (eds.) Successful CBR Applications-2, pp. 55–85. Springer, Heidelberg (2014). [https://doi.org/10.1007/978-3-642-38736-4\\_5](https://doi.org/10.1007/978-3-642-38736-4_5)
12. Recio, J.A., Díaz-Agudo, B., Gómez-Martín, M.A., Wiratunga, N.: Extending jCOLIBRI for textual CBR. In: Muñoz-Ávila, H., Ricci, F. (eds.) ICCBR 2005. LNCS (LNAI), vol. 3620, pp. 421–435. Springer, Heidelberg (2005). [https://doi.org/10.1007/11536406\\_33](https://doi.org/10.1007/11536406_33)
13. Recio-García, J.A., González-Calero, P.A., Díaz-Agudo, B.: jCOLIBRI2: a framework for building CBR systems. *Sci. Comput. Program.* **79**, 126–145 (2014)
14. Recio-Garcia, J.A., González-Calero, P.A., Diaz-Agudo, B.: Template-based design in COLIBRI studio. *Inf. Syst.* **40**, 168–178 (2014)
15. Roth-Berghofer, T., Recio-Garcia, J.A., Severing-Sauer, C., Althoff, K.D., Diaz-Agudo, B.: Building CBR applications with myCBR and COLIBRI. In: Proceedings of 17th UK Workshop on CBR, pp. 71–82. University of Brighton (2012)
16. Stahl, A., Roth-Berghofer, T.R.: Rapid prototyping of CBR applications with the open source tool myCBR. In: Althoff, K.-D., Bergmann, R., Minor, M., Hanft, A. (eds.) ECCBR 2008. LNCS (LNAI), vol. 5239, pp. 615–629. Springer, Heidelberg (2008). [https://doi.org/10.1007/978-3-540-85502-6\\_42](https://doi.org/10.1007/978-3-540-85502-6_42)
17. Vaswani, A., et al.: Attention is all you need. In: Advances in Neural Information Processing Systems, pp. 5998–6008 (2017)
18. Voit, A., Stankus, A., Magomedov, S., Ivanova, I.: Big data processing for full-text search and visualization with elasticsearch. *Int. J. Adv. Comput. Sci. Appl.* **8**(12), 18 (2017)
19. Zhong, Z., Xu, T., Wang, F., Tang, T.: Text CBR framework for fault diagnosis and prediction by cloud computing. *Math. Probl. Eng.* **2018**, 1–10 (2018)