Enhanced Clients for Data Stores and Cloud Services

Arun Iyengar, Fellow, IEEE

Abstract—Data stores and cloud services are typically accessed using a client-server paradigm wherein the client runs as part of an application process which is trying to access the data store or cloud service. This paper presents the design and implementation of enhanced clients for improving both the functionality and performance of applications accessing data stores or cloud services. Our enhanced clients can improve performance via multiple types of caches, encrypt data for providing confidentiality before sending information to a server, and compress data for reducing the size of data transfers. Our clients can perform data analysis to allow applications to more effectively use cloud services. They also provide both synchronous and asynchronous interfaces. An asynchronous interface allows an application program to access a data store or cloud service and continue execution before receiving a response which can significantly improve performance.

We present a Universal Data Store Manager (UDSM) which allows an application to access multiple different data stores and provides a common interface to each data store. The UDSM also can monitor the performance of different data stores. A workload generator allows users to easily determine and compare the performance of different data stores.

We also present NLU-SA, an application for performing natural language understanding and sentiment analysis on text documents. NLU-SA is implemented on top of our enhanced clients and integrates text analysis with Web searching. We present results from NLU-SA on sentiment on the Web towards major companies and countries. We also present a performance analysis of our enhanced clients.

Index Terms—data store client, database client, cloud service client, cognitive service client, Web service client, caching, client caching, sentiment analysis.

1 INTRODUCTION

Data stores are typically accessed using a client-server paradigm in which clients run as part of an application program and communicate with data store servers. Cloud services are accessed using a similar client-server paradigm. There is often little support provided for data store and cloud service clients. It is incumbent on application writers to provide features beyond the basics. This paper addresses the problem of providing enhanced clients for data stores and cloud services which improve both functionality and performance.

A broad range of data stores are currently available including SQL (relational) databases, NoSQL databases, caches, and file systems. An increasing number of data stores are available on the cloud and through open source software. There clearly is a need for software which can easily provide access to multiple data stores as well as compare their performance. One of the goals of this work is to address this need.

A second goal of this work is to improve data store performance. Latencies for accessing data stores are often high. Poor data store performance can present a critical bottleneck for users. For cloud-based data stores where the data is stored at a location which is distant from the application, the added latency for communications between the data store server and the applications further increases data store latencies [1], [2]. Techniques for improving data store performance such as caching are thus highly desirable.

A related issue is that there are benefits to keeping data sizes small; compression can be a key component for achieving this. Client-based compression is important since not all servers support compression. Even if servers have efficient compression capabilities, client-side compression can still improve performance by reducing the number of bytes that need to be transmitted between the client and server. In cloud environments, a data service might charge based on the size of objects sent to the server. Compressing data at the client before sending the data to the server can save clients money in this type of situation.

A third goal of this work is to help applications more effectively use cloud services which are not necessarily data stores, such as natural language processing services. Our work can improve both the functionality and performance of these types of applications.

A fourth goal of this work is to provide data confidentiality as it is critically important to many users and applications. Giving users the means to encrypt data may be essential either because a data store lacks encryption features or data is not securely transmitted between the application and data store server. For certain applications, encryption at the client is a necessity as the data store provider simply cannot be trusted to be secure. There have been many serious data breaches in recent years in which confidential data from hundreds of millions of people have been stolen. Some of the widely publicized data breaches have occurred at the US government, Equifax, Yahoo!, Anthem, Democratic National Committee, eBay, Home Depot,

Arun Iyengar is with IBM's T.J. Watson Research Center, Yorktown Heights, NY, 10598



Fig. 1: Data stores are typically accessed using clients. This work focuses on enhancing the functionality of clients.

JP Morgan Chase, Sony, and Target.

Data stores are often implemented using a client-server paradigm in which a client associated with an application program communicates with one or more servers using a protocol such as HTTP (Figure 1). Clients provide interfaces for application programs to access servers. Some but not all data stores are cloud-based. This paper focuses on providing multiple data store options, improving performance, and providing data confidentiality by enhancing data store clients. We also provide enhanced clients for cloud services which are not necessarily data stores. No changes are required to servers or cloud services. That way, our techniques can be used by a broad range of data stores and cloud services. Requiring changes to a server or cloud service would entail significantly higher implementation costs and would seriously limit the number of data stores and cloud services our techniques could be applied to.

We present architectures and implementations of clients which provide enhanced functionality such as caching, encryption, compression, asynchronous (nonblocking) interfaces, and performance monitoring. We also present a universal data store manager (UDSM) which gives application programs access to a broad range of data store options along with the enhanced functionality for each data store. Caching, encryption, compression, and asynchronous (nonblocking) interfaces are essential; users would benefit considerably if they become standard features of data store clients. Unfortunately, that is not the case today. Research in data stores often focuses on server features with inadequate attention being paid to client features.

The use of caching at the client for reducing latency is particularly important when data stores are remote from the applications accessing them. This is often the case when data is being stored in the cloud. The network latency for accessing data at geographically distant locations can be substantial [3]. Client caching can dramatically reduce latency in these situations. With the proliferation of cloud data stores that is now taking place, caching becomes increasingly important for improving performance.

Our enhanced data store clients and UDSM are architected in a modular way which allows a wide variety of data stores, caches, encryption algorithms, and compression algorithms. Widely used data stores such as Cloudant (built on top of CouchDB), OpenStack Object Storage, and Cassandra have existing clients implemented in commonly used programming languages. Popular language choices for data store clients are Java, Python, and Node.js (which is actually a JavaScript runtime built on Chrome's V8 JavaScript engine). These clients handle low level details such as communication with the server using an underlying protocol such as HTTP. That way, client applications can communicate with the data store server via method (or other type of subroutine) calls in the language in which the client is written. Examples of such clients include the Cloudant Java client [4], the Java library for OpenStack Storage (JOSS) [5], and the Java Driver for Apache Cassandra [6] (Figure 1).

The UDSM is built on top of existing data store clients. That way, we do not have to re-implement features which are already present in an existing client. The UDSM allows multiple clients for the same data store if this is desired.

It should be noted that client software for data stores is constantly evolving, and new clients frequently appear. The UDSM is designed to allow new clients for the same data store to replace older ones as the clients evolve over time.

A key feature of the UDSM is a common key-value interface which is implemented for each data store supported by the UDSM. If the UDSM is used, the application program will have access to both the common key-value interface for each data store as well as customized features of that data store that go beyond the key-value interface, such as SQL queries for a relational database. If an application uses the key-value interface, it can use any data store supported by the UDSM since all data stores implement the interface. Different data stores can be substituted for the key-value interface as needed.

The UDSM and enhanced clients provide a synchronous (blocking) interface to data stores and cloud services for which an application will block while waiting for a response to a request. They also provide an asynchronous (nonblocking) interface to data stores and cloud services wherein an application program can make a request and not wait for the request to return a response before continuing execution. The asynchronous interface is important for applications which do not need to wait for all data store or cloud service operations to complete before continuing execution and can often considerably reduce the completion time for such applications.

Most existing data store clients only provide a synchronous interface and do not offer asynchronous operations on the data store. A key advantage to our UDSM is that it provides an asynchronous interface to all data stores it supports, even if a data store does not provide a client with asynchronous operations on the data store.

The UDSM also provides monitoring capabilities as well as a workload generator which allows users to easily determine the performance of different data stores and compare them to pick the best option.

While caching can significantly improve performance, the optimal way to implement caching is not straightforward. There are multiple types of caches currently available with different performance trade-offs and features [7], [8], [9], [10], [11]. Our enhanced clients can make use of multiple caches to offer the best performance and functionality. We are not tied to a specific cache implementation. As we describe in Section 3, it is important to have implementations of both an in-process cache as well as a remote process cache like Redis [7] or memcached [8] as the two approaches are applicable to different scenarios and have different performance characteristics. We provide key additional features on top of the base caches such as expiration time management and the ability to encrypt or compress data before caching it.

The way in which a cache such as Redis is integrated with a data store to properly perform caching is key to achieving an effective caching solution. We have developed three types of integrated caching solutions on top of data stores. They vary in how closely the caching API calls are integrated with the data store client code. We discuss this further in Section 3.

Our paper makes the following key contributions:

- We present the design and implementation of enhanced data store clients which improve performance and security by providing integrated caching, encryption, and compression. We have written a library for implementing enhanced data store clients and made it available as open source software [12]. This is the first paper to present caching, encryption, and compression as being essential features for data store clients and to describe in detail how to implement these features in data store clients in a general way. Our enhanced data store clients are being used by IBM customers.
- We have generalized our enhanced clients to support applications using cloud services which are not necessarily data stores. An open source version of our enhanced clients for supporting cloud services such as IBM's Natural Language Understanding is available [13] and is being used by customers. We are not aware of any other software which provides similar functionality.
- We present NLU-SA, an application for performing natural language understanding and sentiment analysis on text documents. NLU-SA is implemented on top of our enhanced clients and integrates text analysis with Web searching. We also provide new results on sentiment for major companies and countries expressed on the Web which were collected by NLU-SA. We are not aware of any other software which provides the same functionality as NLU-SA. Furthermore, the results on sentiment analysis that we report from using NLU-SA are novel.
- We present the design and implementation of a universal data store manager (UDSM) which allows application programs to access multiple data stores. The UDSM allows data stores to be accessed with a common key-value interface or with an interface specific to a certain data store. The UDSM provides the ability to monitor the performance of different data stores as well as a workload generator which can be used to easily compare the performance of data stores. The UDSM also has features provided by our enhanced clients so that caching, encryption, and compression can be used to enhance performance and security for all of the data stores supported by the UDSM. The UDSM provides both a synchronous (blocking) and asynchronous (nonblocking) interface to all data stores it supports, even if a data store fails to provide a client supporting asynchronous op-



Fig. 2: Enhanced data store client.

erations on the data store. Asynchronous interfaces are another commonly ignored feature which should become a standard feature of data store clients. The UDSM is available as open source software [14] and is being used by IBM customers. We are not aware of any other software which provides the full functionality of our UDSM or is designed the same way.

- We present key issues that arise with client-side caching and offer multiple ways of implementing and integrating caches with data store clients.
- We present performance results from using our enhanced clients and UDSM. The results show the high latency that can occur with cloud-based data stores and the considerable latency reduction that can be achieved with caching. The latency from remote process caching can be a problem and is significantly higher than latency resulting from in-process caching. We also quantify the performance gains that can be achieved using asynchronous interfaces.

The remainder of this paper is structured as follows. Section 2 presents the overall design of our enhanced clients and UDSM. Section 3 presents some of the key issues with effectively implementing client-side caching. Section 4 presents NLU-SA, an application we have written on top of our enhanced clients which performs natural language understanding and sentiment analysis. Section 5 presents a performance evaluation of several aspects of our system. Section 6 presents related work. Finally, Section 7 concludes the paper.

2 DESIGN AND IMPLEMENTATION OF ENHANCED CLIENTS

2.1 Enhanced Clients for Data Stores

Our enhanced data store clients are built on top of a data store client library (DSCL) which handles features such as caching, encryption, compression, and delta encoding (Figure 2). For important features, there is an interface and multiple possible implementations. For example, there are multiple caching implementations which a data store client can choose from. A Java implementation of the DSCL is available as open source software [12]. A guide for using this DSCL is also available [15].

Commonly used data stores such as Cloudant, Open-Stack Object Storage, Cassandra, etc. have clients in widely used programming languages which are readily available as open source software [4], [5], [6]. These clients make it easy for application programs to access data stores since they can use convenient function/method calls of a programming language with properly typed arguments instead of dealing with low level details of a protocol for communicating with a data store.

While existing clients for data stores handle the standard operations defined on that data store, they generally do not provide enhanced features such as caching, encryption, compression, performance monitoring, and a workload generator to test the performance of a data store. Our DSCL and UDSM provide these enhanced features and are designed to work with a wide variety of existing data store clients.

Our DSCL can have various degrees of integration with an existing data store client. In a tight integration, the data store client is modified to make DSCL calls at important places in the client code. For example, DSCL calls could be inserted to first look for an object in a cache when a data store client queries a server for an object. DSCL API calls could also be inserted to update (or invalidate) an object in a cache when a data store client updates an object at a server. More complicated series of API calls to the DSCL could be made to achieve different levels of cache consistency. A similar approach can be used to enable a data store client to perform encryption, decryption, compression, and/or decompression transparently to an application.

Tight integration of a data store client with the DSCL requires source code modifications to the data store client. While this is something that software developers of the data store client should be able to, it is not something that a typical user of the client can be expected to do. However, tight integration is not required to use the DSCL with a data store client. The DSCL has explicit API calls to allow caching, encryption, and compression. Users can thus use the DSCL within an application independently of any data store. The advantage to a tight integration with a data store is that the user does not have to make explicit calls to the DSCL for enhanced features such as caching, encryption, and compression; the data store client handles these operations automatically. In a loosely coupled implementation in which the data store client is not modified with calls to the DSCL, the user has to make explicit DSCL calls to make use of enhanced features.

Optimal use of the DSCL is achieved with a tight coupling of the DSCL with data store clients along with exposing the DSCL API to give the application fine grained control over enhanced features. In some cases, it may be convenient for an application to make API calls to client methods which transparently make DSCL calls for enhanced features. In other cases, the application program may need to explicitly make DSCL calls to have precise control over enhanced features.

The size of data transfers between the client and server can sometimes be reduced by delta encoding in which a client communicates an update to a server as a difference between the current and previous version of a data object [16]. If this difference between the current and previous version of the object is smaller than the current object itself after compression, delta encoding results in fewer bytes communicated between the client and server. The use of



Fig. 3: We provide enhanced clients for a variety of cloud applications.

delta encoding for enhanced clients is presented in [17].

2.2 Enhanced Clients for Cloud Services

Our enhanced clients are of particular value for applications accessing cloud services. We provide enhanced clients for applications using cloud services in addition to data stores (Figure 3). For example, there are several services accessible over the Web related to artificial intelligence and cognitive computing, including those from IBM [18], Microsoft [19], Amazon [20], and Google [21]. These services are typically accessed over HTTP and often return data in a standard format such as JSON or XML. Our enhanced clients wrap HTTP calls to Web services in method, function or procedure calls in a programming language such as Java. That way, an application can invoke a cloud service using a method, function, or procedure instead of having to having to make HTTP requests in the exact format required by the Web service. Our clients also handle failures and timeouts resulting from HTTP requests.

We have developed a number of features specifically for cloud services which perform natural language processing (NLP) such as IBM's Natural Language Understanding service [22]. These services typically analyze single documents at a time. Application programs must then analyze results from each individual document. Our enhanced clients can perform analyses on multiple documents by sending requests for each individual document to one or more Web services and aggregating and analyzing the results from the Web services calls. We provide the ability to store results persistently at the client, which can be used to avoid redundant calls to a Web service.

Our enhanced clients integrate Web searching with NLP. It is possible to perform Web searches on a wide variety of topics and analyze the results from the Web searches using NLP techniques. The Web searches can be restricted to news stories.

An open source version of our enhanced clients for supporting cloud services such as IBM's Natural Language Understanding is available [13] and is being used by customers. In Section 4, we describe an application for natural language understanding and sentiment analysis which has been implemented using our enhanced clients.



Fig. 4: Universal data store manager.

2.3 Universal Data Store Manager

Existing data store clients typically only work for a single data store. This limitation is a key reason for developing our Universal Data Store Manager (UDSM) which allows application programs to access multiple data stores including file systems, SQL (relational) databases, NoSQL databases, and caches (Figure 4).

The UDSM provides a common key-value interface. Each data store implements the key-value interface. That way, it is easy for an application to switch from using one data store to another. The key-value interface exposed to application programs hides the details of how the interface is actually implemented by the underlying data store.

In some cases, a key-value interface is not sufficient. For example, a MySQL user may need to issue SQL queries to the underlying database. The UDSM allows the user to access native features of the underlying data store when needed. That way, applications can use the common keyvalue interface when appropriate as well as all other capabilities of the underlying data store when necessary.

A key feature of the UDSM is the ability to monitor different data stores for performance. Users can measure and compare the performance of different data stores. The UDSM collects both summary performance statistics such as average latency as well as detailed performance statistics such as past latency measurements taken over a period of time. It is often desirable to only collect latency measurements for recent requests. There is thus the capability to collect detailed data for recent requests while only retaining summary statistics for older data. Performance data can be stored persistently using any of the data stores supported by the UDSM.

The UDSM also provides a workload generator which can be used to generate requests to data stores in order to determine performance. The workload generator allows users to specify the data to be stored and retrieved. The workload generator automatically generates requests over a range of different request sizes specified by the user. The workload generator can synthetically generate data objects to be stored. Alternatively, users can provide their own data objects for performance tests either by placing the data in input files or writing a user-defined method to provide the data. The workload generator also determines read latencies when caching is being used for different hit rates specified by the user. Additionally, the workload generator also measures the overhead of encryption and compression.

The workload generator is ideal for easily comparing the performance of different data stores across a wide range of data sizes and cache hit rates. Performance will vary depending on the client, and the workload generator can easily run on any UDSM client. The workload generator was a critical component in generating the performance data in Section 5. Data from performance testing is stored in text files which can be easily imported into data visualization and analysis tools.

A Java implementation of the UDSM is available as open source software [14]. A guide for using this UDSM is also available [23]. This UDSM includes existing Java clients for data stores such as the Cloudant Java client [4] and the Java library for OpenStack Storage (JOSS) [5]. Other data store clients, such as the Java Driver for Apache Cassandra [6], could also be added to the UDSM. That way, applications have access to multiple data stores via the APIs in the existing clients. It is necessary to implement the UDSM key-value interface for each data store; this is done using the APIs provided by the existing data store clients. The UDSM allows SQL (relational) databases to be accessed via JDBC. The key-value interface for SQL databases can also be implemented using JDBC.

The common key-value interface serves a number of useful purposes. It hides the implementation details and allows multiple implementations of the key-value interface. This is important since different implementations may be appropriate for different scenarios. In some cases, it may be desirable to have a cloud-based key-value store. In other cases, it may be desirable to have a key-value store implemented by a local file system. In yet other cases, it may be desirable to have a cache implementation of the key-value interface such as Redis or memcached. Since the application accesses all implementations using the same key-value interface, it is easy to substitute different key-value store implementations within an application as needed without changing the source code.

Another advantage of the key-value interface is that important features such as asynchronous interfaces, performance monitoring, and workload generation can be performed on the key-value interface itself, automatically providing the feature for all data stores implementing the keyvalue interface. Once a data store implements the key-value interface, no additional work is required to automatically get an asynchronous interface, performance monitoring, or workload generation for the data store. In our Java implementation of the UDSM, applying important features to all data store implementations in the same code is achieved by defining a

public interface KeyValue<K,V> {

which each data store implements. The code which implements asynchronous interfaces, performance monitoring, and workload generation takes arguments of type KeyValue<K, V> rather than an implementation of KeyValue<K, V>. That way, the same code can be applied to each implementation of KeyValue<K, V>.

The UDSM provides data encryption and compression in a similar fashion as the DSCL. The DSCL can be used to pro-

vide integrated caching for any of the data stores supported by the UDSM. In addition, the key-value interface allows UDSM users to manually implement caching without using the DSCL. The key point is that via the key-value interface, any data store can serve as a cache or secondary repository for one of the other data stores functioning as the main data store. The user would make appropriate method calls via the key-value interface to maintain the contents of a data store functioning as a cache or secondary repository.

2.4 Asynchronous Interfaces to Data Stores and Cloud Services

Most interfaces to data stores and cloud services are synchronous (blocking). An application will access a data store or service via a method or function call and wait for the method or function to return before continuing execution. Performance can often be improved via asynchronous (nonblocking) interfaces wherein an application can access a data store (to store a data value, for example) or service and continue execution before the call to the interface returns. Our enhanced clients offer both synchronous and asynchronous interfaces to data stores and services.

The asynchronous interface allows applications to continue executing after a call to a method A to access a data store or service by using a separate thread for invoking method A. Since creating a new thread is expensive, our enhanced clients use thread pools in which a given number of threads are started up when the client is initiated and maintained throughout the lifetime of the client. A method invoked via an asynchronous interface is assigned to one of the existing threads in the thread pool which avoids the costly creation of new threads to handle asynchronous method calls. Users can specify the thread pool size via a configuration parameter.

The Java enhanced client implementation implements asynchronous calls using the ListenableFuture interface [24]. Java provides a Future interface which can be used to represent the result of an asynchronous computation. The Future interface provides methods to check if the computation corresponding to a future is complete, to wait for the computation to complete if it has not finished executing, and to retrieve the result of the computation after it has finished executing. The ListenableFuture extends the Future interface by giving users the ability to register callbacks which are code to be executed after the future completes execution. This feature is the key reason that we use ListenableFutures instead of only Futures for implementing asynchronous interfaces. In Section 5, we evaluate the performance of both synchronous and asynchronous interfaces and quantify the significant performance improvements that can be achieved using asynchronous interfaces.

3 CACHING

Our enhanced data store clients use three types of caching approaches. In the first approach, caching is closely integrated with a particular data store. Method calls to retrieve data from the data store can first check if the data are cached, and if so, return the data from the cache instead of the data store. Methods to store data in the data store can also update the cache. Techniques for keeping caches updated and consistent with the data store can additionally be implemented. This first caching approach is achieved by modifying the source code of a data store client to read, write, and maintain the cache as appropriate by making appropriate method calls to the DSCL.

We have used this first approach for implementing caching for Cloudant. The source code for this implementation is available from [25]. We have also used this approach for implementing caching for OpenStack Object Storage by enhancing the Java library for OpenStack Storage (JOSS) [5]. While this approach makes things easier for users by adding caching capabilities directly to data store API methods, it has the drawback of requiring changes to the data store client source code. In some cases, the client source code may be proprietary and not accessible. Even if the source code is available (e.g. via open source), properly modifying it to incorporate caching functionality can be time consuming and difficult.

The second approach for implementing caching is to provide the DSCL to users and allow them to implement their own customized caching solutions using the DSCL API. The DSCL provides convenient methods allowing applications to both query and modify caches. The DSCL also allows cached objects to have expiration times associated with them; later in this section, we will describe how expiration times are managed. It should be noted that if the first approach of having a fully integrated cache with a data store is used, it is still often necessary to allow applications to directly access and modify caches via the DSCL in order to offer higher levels of performance and data consistency. Hence, using a combination of the first and second caching approaches is often desirable.

The third approach for achieving caching is provided by the UDSM. The UDSM key-value interface is implemented for both main data stores as well as caches. If an application is using the key-value interface to access a data store, it is very easy for the application writer to store some of the key-value pairs in a cache provided by the UDSM. Both the main data store and cache share the same key-value interface. This approach, like the second approach, requires the application writer to explicitly manage the contents of caches. Furthermore, the UDSM lacks some of the caching features provided by the DSCL such as expiration time management. An advantage to the third approach is that any data store supported by the UDSM can function as a cache or secondary repository for another data store supported by the UDSM; this offers a wide variety of choices for caching or replicating data.

The DSCL also supports multiple different types of caches via a Cache interface which defines how an application interacts with caches. There are multiple implementations of the Cache interface which applications can choose from.

There are two types of caches. In-process caches store data within the process corresponding to the application [26]. That way, there is no interprocess communication required for storing the data. For our Java implementations of in-process caches, Java objects can directly be cached. Data serialization is not required. In order to reduce overhead when the object is cached, the object (or a reference to it) can be stored directly in the cache. One consequence of this approach is that changes to the object from the application will change the cached object itself. In order to prevent the value of a cached object from being modified by changes to the object being made in the application, a copy of the object can be made before the object is cached. This results in overhead for copying the object.

Another approach is to use a remote process cache [27]. In this approach, the cache runs in one or more processes separate from the application. A remote process cache can run on a separate node from the application as well. There is some overhead for communication with a remote process cache. In addition, data often has to be serialized before being cached. Therefore, remote process caches are generally slower than in-process caches (as shown in Section 5). However, remote process caches also have some advantages over in-process caches. A remote process cache can be shared by multiple clients, and this feature is often desirable. Remote process caches can often be scaled across multiple processes and nodes to handle high request rates and increase availability.

There are several caches that are available as open source software which can be used in conjunction with our DSCL. Redis [7] and memcached [8] are widely used remote process caches. They can be used for storing serialized data across a wide range of languages. Clients for Redis and memcached are available in Java, C, C++, Python, Javascript, PHP, and several other languages.

Examples of caches targeted at Java environments include the Guava cache [9], Ehcache [10], and OSCache [11]. A common approach is to use a data structure such as a HashMap or a ConcurrentHashMap with features for thread safety and cache replacement. Since there are several good open source alternatives available, it is probably better to use an existing cache implementation instead of writing another cache implementation unless the user has specialized requirements not handled by existing caches.

Our DSCL allows any of these caches to be plugged into its modular architecture. In order to use one of these caches, an implementation of the DSCL Cache interface needs to be implemented for the cache. We have implemented DSCL Cache interfaces for a number of caches including redis and the Guava cache.

The DSCL allows applications to assign (optional) expiration times to cached objects. After the expiration time for an object has elapsed, the cached object is obsolete and should not be returned to an application until the server has been contacted to either provide an updated version or verify that the expired object is still valid. Cache expiration times are managed by the DSCL and not by the underlying cache. There are a couple of reasons for this. Not all caches support expiration times. A cache which does not handle expiration times can still implement the DSCL Cache interface. In addition, for caches which support expiration times, objects whose expiration times have elapsed might be purged from the cache. We do not always want this to happen. After the expiration time for a cached object has elapsed, it does not necessarily mean that the object is obsolete. Therefore, the DSCL has the ability to keep around a cached object of whose expiration time has elapsed. If of is requested after its expiration time has passed, then the client



Fig. 5: Handling cache expiration times.

might have the ability to revalidate o1 in a manner similar to an HTTP GET request with an If-Modified-Since header. The basic idea is that the client sends a request to fetch o1 only if the server's version of o1 is different than the client's version. In order to determine if the client's version of o1 is obsolete, the client could send a timestamp, entity tag, or other information identifying the version of o1 stored at the client. If the server determines that the client has an obsolete version of o1, then the server will send a new version of o1 to the client. If the server determines that the client has a current version of o1, then the server will indicate that the version of o1 is current (Figure 5).

Using this approach, the client does not have to receive identical copies of objects whose expiration times have elapsed even though they are still current. This can save considerable bandwidth and improve performance. There is still latency for revalidating o1 with the server, however.

If caches become full, a cache replacement algorithm such as least recently used (LRU) or greedy-dual-size [28] can be used to determine which objects to retain in the cache.

Some caches such as redis have the ability to back up data in persistent storage (e.g. to a hard disk or solid-state disk). This allows data to be preserved in the event that a cache fails. It is also often desirable to store some data from a cache persistently before shutting down a cache process. That way, when the cache is restarted, it can quickly be brought to a warm state by reading in the data previously stored persistently.

The encryption capabilities of the DSCL can also be used in conjunction with caching. People often fail to recognize the security risks that can be exposed by caching. A cache may be storing confidential data for extended periods of time. That data can become a target for hackers in an insecure environment. Most caches do not encrypt the data they are storing, even though this is sometimes quite important.

Remote process caches also present security risks when an application is communicating with a cache over an unencrypted channel. A malicious party can steal the data being sent between the application and the cache. Too often, caches are designed with the assumption that they will be deployed in a trusted environment. This will not always be the case.

For these reasons, data should often be encrypted before it is cached. The DSCL provides the capability for doing so. There is some CPU overhead for encryption, so privacy needs need to be balanced against the need for fast execution.

The DSCL compression capabilities can also be used to reduce the size of cached objects, allowing more objects to be stored using the same amount of cache space. Once again, since compression entails CPU overhead, the space saved by compression needs to be balanced against the increase in CPU cycles resulting from compression and decompression.

4 USING ENHANCED CLIENTS FOR NATURAL LANGUAGE UNDERSTANDING AND SENTIMENT ANALYSIS

Our enhanced clients can be used by a wide variety of applications for improving both functionality and performance. In this section, we describe an application we have built using our enhanced clients which performs natural language understanding and sentiment analysis on text documents. Our application is known as NLU-SA based on the first letters of *natural language understanding* and *sentiment analysis*. We also present the results of using NLU-SA to analyze Web documents for sentiment over a period of time.

NLU-SA uses our enhanced clients to invoke cloudbased NLP services. The cloud-based NLP services will typically analyze a single text document at a time. NLU-SA aggregates multiple text documents, passes each document to a cloud-based NLP service, and aggregates and analyzes the results from the cloud-based NLP service. One of the key use cases we have implemented is analyzing the results from Web searches.

NLU-SA can summarize the contents of several documents. For example, Tables 1 and 2 illustrate the types of analyses that can be performed on search queries. NLU-SA invoked a Web search engine to obtain the top 50 URLs resulting from a query on "United States" on September 24, 2017 made from Yorktown Heights, New York. The corresponding Web documents were analyzed by NLU-SA sending each URL to IBM's Alchemy Language Web service. NLU-SA analyzes the results from each invocation of Alchemy Language to produce aggregate results across all documents. NLU-SA can also be used with other Web services performing base NLP functions instead of IBM's Alchemy Language.

Table 1 shows the top 10 disambiguated entities based on a relevance score. A disambiguated entity is an item such as a person, place, or organization that is present in the input text and has been properly identified [29]. The relevance score indicates how relevant an entity is to one or more documents. While the relevance score is correlated with how many times the entity appears in the documents, the correlation is not perfect, as can be seen in the tables. Entities also have a sentiment score associated with them. Sentiment scores are a quantification of attitudes, opinions, or feelings expressed in the text being analyzed to the entity [30]. Sentiment scores range from -1 to 1 with 0 being neutral and higher scores representing more favorable sentiment towards the entity.

Table 2 shows the top 10 concepts based on a relevance score.

We have used NLU-SA to study sentiment expressed in Web documents and how it changes over a period of

TABLE 1: Top 10 disambiguated entities sorted by relevance from Web documents on the United States.

Disambiguated Entity	Occurrences	Relevance	Sentiment
United States	688	21.223	-0.223
Facebook	21	3.366	0.470
Washington, D.C.	31	2.789	-0.003
Canada	41	2.324	0.040
Federal govt. of the US	34	2.311	-0.124
Donald Trump	20	2.154	-0.215
New York City	23	2.090	-0.184
United States Congress	51	2.082	-0.156
Hawaii	37	1.902	0.248
North America	20	1.899	0.047

TABLE 2: Top 10 concepts sorted by relevance from Web documents on the United States.

Concept	Occurrences	Relevance
United States	28	21.002
U.S. state	16	8.286
US President	14	7.991
Federal govt. of the US	9	5.429
Washington, D.C.	11	5.255
US Constitution	9	4.875
US Congress	7	3.883
New York City	7	3.783
Puerto Rico	6	3.357
World War II	5	3.269

several weeks. For the experimental results we are about to present, NLU-SA obtained several URLs from a query to a Web search engine for an entity such as *Canada*. Each of the URLs obtained from the Web search engine was then passed to IBM's Alchemy Language Web service which calculates the sentiment for the entity (such as *Canada*) by analyzing the Web document corresponding to the URL. NLU-SA determines a single sentiment value for the entity over all URLs returned by the search query by computing an average of the sentiment values returned by Alchemy Language weighted by the number of occurrences of the entity in each Web document.

Figure 6 shows the sentiment expressed for four major technology companies (referred to as Tech 1, Tech 2, Tech 3, and Tech 4) and a major financial company (referred to as Finance) in the first 50 documents obtained by a Web search on the company name. The results were obtained over a 32-day period starting on June 30 and ending on July 31 of 2017. There are clear differentiations between the companies with Tech 3 finishing first followed by Tech 2, Tech 1, Tech 4, and the financial company respectively.

Figure 7 shows the sentiment expressed for the same companies over the same time period in the first 50 news stories obtained from Web searches restricted to news stories for a query comprised of the company name. The financial company finishes lower than the other companies. For the technology companies, the differentiations in sentiment are less clear cut than in Figure 6. This is because news stories are ephemeral and constantly changing. By contrast, the documents obtained from a general Web search do not change as much from day to day.

Figure 8 shows the sentiment expressed for Canada, China, Germany, Russia, and the United States in the first



Fig. 6: Sentiment for companies in text documents returned from Web searches on the company name.



Fig. 7: Sentiment for companies in news stories from Web searches on the company name.

50 documents obtained by a Web search on the country name over the same 32-day period. One of the surprising results is the relatively favorable sentiment expressed towards China compared with the relatively unfavorable sentiment expressed towards the United States and Canada. The Web searches were made in Yorktown Heights, New York.

Figure 9 shows the sentiment expressed for the same countries over the same time period in the first 50 news stories obtained from Web searches restricted to news stories for a query comprised of the country name. Once again, the relatively favorable sentiment expressed towards China compared with other countries stands out.

NLU-SA can use different search engines. Figure 10 shows the sentiment expressed for the five companies in the first 50 documents obtained by searches on the company name on August 13, 2017. Search engine 1 is a major search engine which we also used for the other experimental results contained in this section. Search engine 2 is another major search engine which competes with Search engine 1. While

Sentiment Scores, Regular Search



Fig. 8: Sentiment for countries in text documents from Web searches on the country name.

Sentiment Scores, News Stories



Fig. 9: Sentiment for countries in news stories from Web searches on the country name.

the values differ somewhat when different search engines are used, the relative rankings of the companies do not change.

We performed sentiment analysis on documents obtained from Web searches on the two search engine names. Both search engines support the fact that Search engine 1 has more favorable sentiment on the Web than Search engine 2. In fact, Search engine 2 returns documents with a higher difference in sentiment in favor of Search engine 1 over Search engine 2 than the documents returned by Search engine 1. While this seems to suggest that neither search engine is biasing its search results to favor itself over the other search engine, more analysis would be needed before definitive statements along these lines can be made.

Figure 11 shows the sentiment expressed for Canada, China, Germany, Russia, and the United States in the first 50 documents obtained by both search engines on the country name on August 13, 2017. Germany finishes highest based on Search engine 1 documents, while China is highest based Comparison of Different Search Engines



Fig. 10: Sentiment for companies in text documents from Web searches on the company name using different search engines.



Fig. 11: Sentiment for countries in text documents from Web searches on the country name using different search engines.

on Search engine 2 documents. Russia is second based on Search engine 2 documents but fourth based on Search engine 1 documents.

Figures 12 - 15 show changes in sentiment over a longer period of time. Initial values were obtained from the top 50 documents obtained from Web searches on June 30, 2017. Final values were obtained from the top 50 documents obtained from Web searches on September 9, 2017.

5 PERFORMANCE EVALUATION

In this section, we present a performance evaluation of our enhanced clients. We use the UDSM to determine and compare read and write latencies that a typical client would see using several different types of data stores. We show the performance gains that our enhanced data store clients can achieve with caching. We compare the performance of asynchronous and synchronous interfaces. We also quantify

Start and End, Regular Search

Tech1 Tech2 Tech3 Tech4 Fin

Fig. 12: Change over time in sentiment for companies in text documents from Web searches on the company name.

Starting and Ending Sentiment Scores, News



Fig. 13: Change over time in sentiment for companies in news stories from Web searches on the company name.



Fig. 14: Change over time in sentiment for countries in text documents from Web searches on the country name.

Starting and Ending Sentiment Scores, News



Fig. 15: Change over time in sentiment for countries in news stories from Web searches on the country name.

the overheads resulting from encryption, decryption, compression, and decompression.

We test the following data stores by using the UDSM to send requests from its workload generator using the keyvalue interface:

- A file system on the client node accessed via standard Java method calls.
- A MySQL database [31] running on the client node accessed via JDBC.
- A commercial cloud data store provided by a major cloud computing company (referred to as Cloud Store 1).
- A second commercial cloud data store provided by a major cloud computing company (referred to as Cloud Store 2).
- A Redis instance running on the client node accessed via the Jedis client [32].

The Redis instance also acts as a remote process cache. A Guava cache [9] acts as an in-process cache.

Our enhanced clients and UDSM run on a 2.70GHz Intel i7-3740QM processor with 16 GB of RAM running a 64-bit version of Windows 7 Professional. Several of the performance graphs use log-log plots because both the xaxis (data size in bytes) and y-axis (time in milliseconds) values span a wide magnitude of numbers. Experiments were run multiple times. Data points are averaged over 4 runs of the same experiment. We did not find significant correlations between the types of data read and written and data store read and write latencies. There was often considerable variability in the read and write latency for the experimental runs using the same parameters.

Figure 16 shows the average time to read data as a function of data size. Cloud Store 1 and 2 show the highest latencies because they are cloud data stores geographically distant from the client. By contrast, the other data stores run on the same machine as the client. Another factor that could adversely affect performance for the cloud data stores, particularly in the case of Cloud Store 1, is that the requests coming from our UDSM might be competing for



Fig. 16: Read latencies for data stores.

server resources with computing tasks from other cloud users. This may be one of the reasons why Cloud Store 1 exhibited more variability in read latencies than any of the other data stores. Significant variability in cloud storage performance has been observed by others as well [2]. The performance numbers we observed for Cloud Store 1 and 2 are not atypical for cloud data stores, which is a key reason why techniques like client-side caching are essential for improving performance.

Redis offers lower read latencies than the file system for small objects. For objects 50 Kbytes and larger, however, the file system achieves lower latencies. Redis incurs overhead for interprocess communication between the client and server. There is also some overhead for serializing and deserializing data stored in Redis. The file system client might benefit from caching performed by the underlying file system.

Redis offers considerably lower read latencies than MySQL for objects up to 50 Kbytes. For larger objects, Redis offers only slightly better read performance, and the read latencies converge with increasing object size.

Figure 17 shows the average time to write data as a function of data size. Cloud Store 1 has the highest latency followed by Cloud Store 2; once again, this is because Cloud Store 1 and 2 are cloud-based data stores geographically distant from the client. MySQL has the highest write latencies for the local data stores. Redis has lower write latencies than the file system for objects of 10 Kbytes or smaller. Write latencies are similar for objects of 20-100 Kbytes, while the file system has lower write latencies than Redis for objects larger than 100 Kbytes.

Write latencies are higher than read latencies across the data stores; this is particularly apparent for MySQL for which writes involve costly commit operations. It is also very pronounced for larger objects with the cloud data stores and the file system. Write latencies for the file system and MySQL exhibit considerably more variation than read latencies.

Figures 18 - 26 show read latencies which can be achieved with the Guava in-process cache and Redis as a remote process cache for the cloud data stores. Multiple runs



Fig. 17: Write latencies for data stores.

were made to determine read latencies for each data store both without caching and with caching when the hit rate is 100%. From these numbers, the workload generator can extrapolate performance for different hit rates. Each graph contains 5 curves corresponding to no caching and caching with hit rates of 25%, 50%, 75%, and 100%.

The in-process cache is considerably faster than any of the data stores. Furthermore, read latencies do not increase with increasing object size because cache reads do not involve any copying or serialization of data. By contrast, Redis is considerably slower, as shown in Figure 26. Furthermore, read latencies increase with object size as cached data objects have to be transferred from a Redis instance to the client process and deserialized. An in-process cache is thus highly preferable from a performance standpoint. Of course, application needs may necessitate using a remote process cache like Redis instead of an in-process cache.

Figure 25 shows that for the file system, remote process caching via Redis is only advantageous for smaller objects; for larger objects, performance is better without using Redis. This is because the file system is faster than Redis for reading larger objects. Figure 23 shows a similar trend. While Redis will not result in worse performance than MySQL for larger objects, the performance gain may be too small to justify the added complexity.

One of the key features of our enhanced clients is the presence of both synchronous and asynchronous interfaces. Asynchronous interfaces are critically important for speeding up a wide range of applications using cloud services or data stores. For example, NLU-SA spends the vast majority of its time making calls to Web services for analyzing text documents. These calls can be made in parallel. A simple sequential implementation will consume a considerable amount of time.

We now quantify the performance improvements that can be achieved with asynchronous interfaces. Figure 27 shows the latency for analyzing text documents using IBM's AlchemyLanguage NLP Web service with a synchronous interface. The text documents analyzed were the top 100 Web documents from a Web search for *United States* on July 19, 2017 made from Yorktown Heights, New York. The first n URLs retrieved from the Web search were passed to Alchemy Language, where n is shown on the X-axis.

Figure 28 shows the latency for the same analysis of text documents using an asynchronous interface with different thread pool sizes. The graph shows that the asynchronous interface is considerably faster than the synchronous interface. As the number of documents analyzed increases, the asynchronous interface can improve performance by more than a factor of 10. This is due to the fact that the asynchronous interface can make several calls to Web services in parallel. Since Alchemy Language can satisfy many requests concurrently, the asynchronous interface can speed things up considerably. The thread pool size limits the number of concurrent requests sent by a client. For our data set, there was not much of a performance gain by using a thread pool size larger than 20.

Since Web services calls consume the vast majority of time spent by NLU-SA, these graphs are indicative of the overall speedup for NLU-SA which can be achieved using asynchronous interfaces.

Figure 29 shows the times that an enhanced data store client requires for encrypting and decrypting data using the Advanced Encryption Standard (AES) [33] and 128-bit keys. Since AES is a symmetric encryption algorithm, encryption and decryption times are similar.

Figure 30 shows the times that an enhanced data store client requires for compressing and decompressiong data using gzip [34]. Decompression times are roughly comparable with encryption and decryption times. However, compression overheads are several times higher.

6 RELATED WORK

Clients exist for a broad range of data stores. A small sample of clients for commonly used data stores includes the Java library for OpenStack Storage (JOSS) [5], the Jedis client for Redis [32] and the Java Driver for Apache Cassandra [6]. These clients do not have the capabilities that our enhanced clients provide. Amazon offers access to data stores such as DynamoDB via a software development kit (SDK) [35] which provides compression, encryption, and both synchronous and asynchronous APIs. Amazon's SDK does not provide integrated caching, performance monitoring, or workload generation; we are not aware of other clients besides our own which offer these features. Furthermore, we offer coordinated access to multiple types of data stores, a key feature which other systems lack. Microsoft's Azure SDK [36] also does not provide the full range of features that we provide.

Remote process caches which provide an API to allow applications to explicitly read and write data as well as to maintain data consistency were first introduced in [27], [37]. A key aspect of this work is that the caches were an essential component in serving dynamic Web data efficiently at several highly accessed Web sites. Memcached was developed several years later [8]. Design and performance aspects for in-process caches were first presented in [26].

There have also been a number of papers which have studied the performance of cloud storage systems. Dropbox, Microsoft SkyDrive (now OneDrive), Google Drive, Wuala (which has been discontinued), and Amazon Cloud drive

Cloud Store 1 Latency, In-Process Cache



Fig. 18: Cloud Store 1 read latencies with in-process caching.



Fig. 20: Cloud Store 2 read latencies with in-process caching.



Fig. 22: MySQL read latencies with in-process caching.

Cloud Store 1 Latency, Remote Cache



Fig. 19: Cloud Store 1 read latencies with remote process caching.

Cloud Store 2 Latency, Remote Cache



Fig. 21: Cloud Store 2 read latencies with remote process caching.



Fig. 23: MySQL read latencies with remote process caching.

14

File System Latency with In-Process Cache



Fig. 24: File system read latencies with in-process caching.



Fig. 26: Redis read latencies with in-process caching.



Fig. 27: Response time for synchronous interface.

File System Latency with Remote Cache



Fig. 25: File system read latencies with remote process caching.

Response Time with Multithreading







Fig. 29: Encryption and decryption latency.



Fig. 30: Compression and decompression latency.

are compared using a series of benchmarks in [1]. No storage service emerged from the study as being clearly the best one. An earlier paper by some of the same authors analyzes Dropbox [38]. A comparison of Box, Dropbox, and SugarSync is made in [2]. The study noted significant variability in service times which we have observed as well. Failure rates were less than 1%. Data synchronization traffic between users and cloud providers is studied in [39]. The authors find that much of the data synchronization traffic is not needed and could be eliminated by better data synchronization mechanisms. Another paper by some of the same authors proposes reducing data synchronization traffic by batched updates [40].

There have been several past works in sentiment analvsis [29], [41], [42]. However, we are not aware of any existing tools with the capabilities of NLU-SA. In addition, our empirical results on sentiment for major companies and countries are new. Methods for supporting data analytics applications which use cognitive services are presented in [43].

7 CONCLUSIONS

This paper has presented enhanced clients which improve both the functionality and performance of applications accessing data stores or cloud services. Our clients support caching, data compression, encryption, and provide both synchronous and asynchronous interfaces. We presented NLU-SA, an application for performing natural language understanding and sentiment analysis on text documents. NLU-SA is built on top of our enhanced clients. We presented results from NLU-SA on sentiment on the Web towards major companies and countries. We also presented a detailed performance analysis of our enhanced clients.

We have also presented a Universal Data Store Manager (UDSM) which allows applications to access multiple data stores. The UDSM provides common synchronous and asynchronous interfaces to each data store, performance monitoring, and a workload generator which can easily compare performance of different data stores from the perspective of the client. Software for implementing enhanced clients and

the UDSM are available as open source software and are being used by IBM customers.

Most existing data store clients only have basic functionality and lack a rich set of features. Users would significantly benefit if caching, encryption, compression, and asynchronous (nonblocking) interfaces become commonly supported in data store clients.

8 ACKNOWLEDGMENTS

Rich Ellis and Mike Rhodes provided assistance and support in developing client-side caching for Cloudant. German Attanasio Ruiz helped develop the Cognitive Client for IBM's Watson Developer Cloud.

REFERENCES

- [1] I. Drago, E. Bocchi, M. Mellia, H. Slatman, and A. Pras, "Benchmarking Personal Cloud Storage," in Proceedings of IMC '13, 2013, pp. 205–212.
- [2] R. Gracia-Tinedo, M. Artigas, A. Moreno-Martinez, C. Cotes, and P. Garcia-Lopez, "Actively Measuring Personal Cloud Storage," in Proceedings of the IEEE 6th International Conference on Cloud *Computing*, 2013, pp. 301–308. J. Dean and P. Norvig, "Latency numbers every programmer
- [3] should know," https://gist.github.com/jboner/2841832.
- Cloudant, "Cloudant java client," https://github.com/cloudant/ [4] java-cloudant.
- Javaswift, "JOSS: Java library for OpenStack Storage, aka Swift," [5] http://joss.javaswift.org/.
- DataStax, "Datastax java driver for apache cassandra," https:// [6] github.com/datastax/java-driver.
- [7] Redis, "Redis home page," http://redis.io/.
- B. Fitzpatrick, "Distributed Caching with Memcached," Linux [8] Journal, no. 124, p. 5, 2004.
- C. Decker, "Caches explained," https://github.com/google/ [9] guava/wiki/CachesExplained.
- [10] Ehcache, "Ehcache: Java's most widely-used cache," http://www. ehcache.org
- OSCache, "Oscache," https://java.net/projects/oscache. [11]
- [12] IBM, "Data Store Client Library," https://developer.ibm.com/ code/open/projects/data-store-client-library/.
- 'Cognitive Client (for IBM's Watson Developer [13] Cloud)," https://github.com/watson-developer-cloud/ cognitive-client-java.
- , "Universal Data Store Manager," https://developer.ibm. [14]com/code/open/projects/universal-data-store-manager/
- [15] A. Iyengar, "Enhanced Storage Clients," IBM Research Division, Yorktown Heights, NY, Tech. Rep. RC 25584 (WAT1512-042), December 2015.
- [16] F. Douglis and A. Iyengar, "Application-specific Delta-encoding via Resemblance Detection," in *Proceedings of the USENIX 2003* Annual Technical Conference, 2003, pp. 113–126.
- [17] A. Iyengar, "Providing Enhanced Functionality for Data Store Clients," in Proceedings of the IEEE 33rd International Conference on Data Engineering (ICDE 2017), April 2017. IBM, "Watson Developer Cloud," https://www.ibm.com/
- [18] IBM, watson/developercloud/
- Microsoft, "Microsoft Cognitive Services," [19] https://www. microsoft.com/cognitive-services.
- [20] Amazon, "Amazon AI," https://aws.amazon.com/amazon-ai/.
- [21] Google, "Cloud Natural Language API," https://cloud.google. com/natural-language/.
- [22] IBM, "Natural Language Understanding," https://www.ibm. com/watson/services/natural-language-understanding/.
- [23] A. Iyengar, "Universal Data Store Manager," IBM Research Division, Yorktown Heights, NY, Tech. Rep. RC 25607 (WAT1605-030), May 2016.
- "ListenableFutureExplained," [24] Google, https://github.com/ google/guava/wiki/ListenableFutureExplained.
- "Java cloudant cache, [25] https://github.com/ Cloudant, cloudant-labs/java-cloudant-cache.

- [26] A. Iyengar, "Design and Performance of a General-Purpose Software Cache," in Proceedings of the 18th IEEE International Performance, Computing, and Communications Conference, 1999, pp. 329-336.
- [27] A. Iyengar and J. Challenger, "Improving Web Server Performance by Caching Dynamic Data," in *Proceedings of the USENIX Sympo*sium on Internet Technologies and Systems, 1997, pp. 49-60.
- [28] P. Cao and S. Irani, "Cost-Aware WWW Proxy Caching Algorithms," in Proceedings of the USENIX Symposium on Internet Technologies and Systems, 1997, pp. 193-206.
- [29] S. Cucerzan, "Large-Scale Named Entity Disambiguation Based on Wikipedia Data," in Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL 2007), June 2007, pp. 708-716
- [30] B. Pang, L. Lee et al., "Opinion mining and sentiment analysis," Foundations and Trends® in Information Retrieval, vol. 2, no. 1-2, pp. 1-135, 2008.
- [31] MySQL, "MySQL home page," https://www.mysql.com/.
- [32] Jedis, "A blazingly small and sane redis java client," https:// github.com/xetorthio/jedis.
- [33] NIST, "Announcing the Advanced Encryption Standard (AES)," National Institute of Standards and Technology, Tech. Rep. Federal Information Standards Publication 197, November 2001.
- [34] Gzip, "The gzip home page," http://www.gzip.org/.[35] Amazon, "AWS SDK for Java," https://aws.amazon.com/ sdk-for-java/. [36] Microsoft, "Microsoft Azure Downloads," https://azure.
- microsoft.com/en-us/downloads/.
- [37] J. Challenger and A. Iyengar, "Distributed Cache Manager and API," IBM Research Division, Yorktown Heights, NY, Tech. Rep. RC 21004 (94070), 1997.
- [38] I. Drago, M. Mellia, M. Munafo, A. Sperotto, R. Sadre, and A. Pras, "Inside Dropbox: Understanding Personal Cloud Storage Services," in Proceedings of IMC '12, 2012, pp. 481-494.
- [39] Z. Li et al., "Towards Network-level Efficiency for Cloud Storage Services," in Proceedings of IMC '14, 2014, pp. 115–128.
- -, "Efficient Batched Synchronization in Dropbox-like Cloud [40] Storage Services," in Proceedings of Middleware 2013, 2013, pp. 307-327
- [41] B. Liu, "Sentiment analysis and opinion mining," Synthesis lectures on human language technologies, vol. 5, no. 1, pp. 1–167, 2012.
- [42] K. Ravi and V. Ravi, "A survey on opinion mining and sentiment analysis: tasks, approaches and applications," Knowledge-Based Systems, vol. 89, pp. 14-46, 2015.
- [43] A. Iyengar, "Supporting Data Analytics Applications Which Uti-lize Cognitive Services," in Proceedings of the 37th IEEE International Conference on Distributed Computing Systems (ICDCS 2017), June 2017.