

Fig. 1. The organization of the SecureNoSQL.

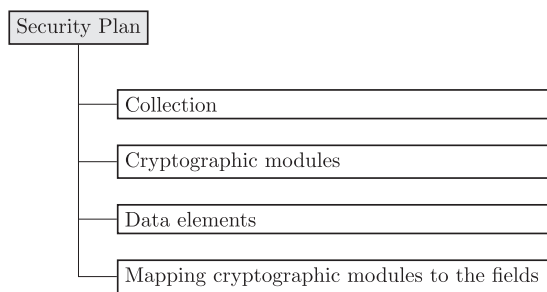


Fig. 2. The high level structure of the security plan.

5.1. Database security plan

The security plan identifies the mechanism to maintain the security of the data elements in a database. It also determines how to interpret queries issued by a specific application. The security plan has four sections, see Fig. 2, describing the security rules for the data elements and for meta-data such as the field-name (Key) and the collection name. These sections are the building blocks of the security plan showing how the rules are enforced. The sections and their roles are:

1. *Collection*: includes the name of a collection and a reference to the encryption module used to encrypt the name of the collection and the name of fields (metadata).
2. *Cryptographic modules*: lists the cryptographic modules for encrypting the fields of the database entries in the query.
3. *Data elements*: lists the properties of each data field including the data type; the data type determines the cryptographic modules to be applied to each field.
4. *Mapping cryptographic modules to the fields*: assigns the cryptographic modules to data fields; proxy uses this information to encrypt and decrypt the data elements.

5.1.1. Collection

A collection is defined as a group of NoSQL documents, the equivalent of relational database table, see Fig. 3. The name of

the collection must be encrypted. The listing 3b illustrates how to secure a sample collection using the description language.

The key-value pairs (KVP) are the primary data model for a NoSQL database. The *key* is used as an index to access the associated value of the data pointed by the reference *ref*. The *initialization vector* (IV) is a fixed-size, random input to the cryptographic module *encryption*. Additionally, a collection exists within a single database. Documents within a collection can have different fields. Typically, all documents in a collection are related with one another.

5.1.2. Cryptographic modules

The choice of a particular cryptosystem depends on the security policy of application. Multiple criteria for algorithm selection include: (i) the security against theoretical attacks; (ii) the cost of implementation; (iii) the performance; and (iv) whether the encryption and decryption can be parallelized. Other factors involved in the selection of an algorithm are the memory requirements and the integration in the overall system design.

The *Cryptographic modules* introduce all encryption modules and their parameters such as key, key-size, initialization vector and output-size. The structure of this section is shown in Fig. 4a complemented by the listing in Fig. 4b presenting the second section of security plan for the previous example.

Our proof of concept uses the parametric Order Preserving Encryption (OPE) and the Advanced Encryption Standard (AES) modules. The system is open-ended; users can add the cryptosystems best suited to the security requirements of their application. In our design the definitions of the cryptographic modules and of the pairs, encryption key and initialization value, are separated following the so-called *key separation principle* (Galiegue & Zyp, 2013). This security practice is based on the observations that users have long- and short-term security policies. The cryptographic modules are less likely to change while the key and the initialization value change frequently.

5.1.3. The data elements

The third section of security plan, the data elements and their properties are covered. Fig. 5 presents the structure and description of *Data element* section of *Security plan*. The listing

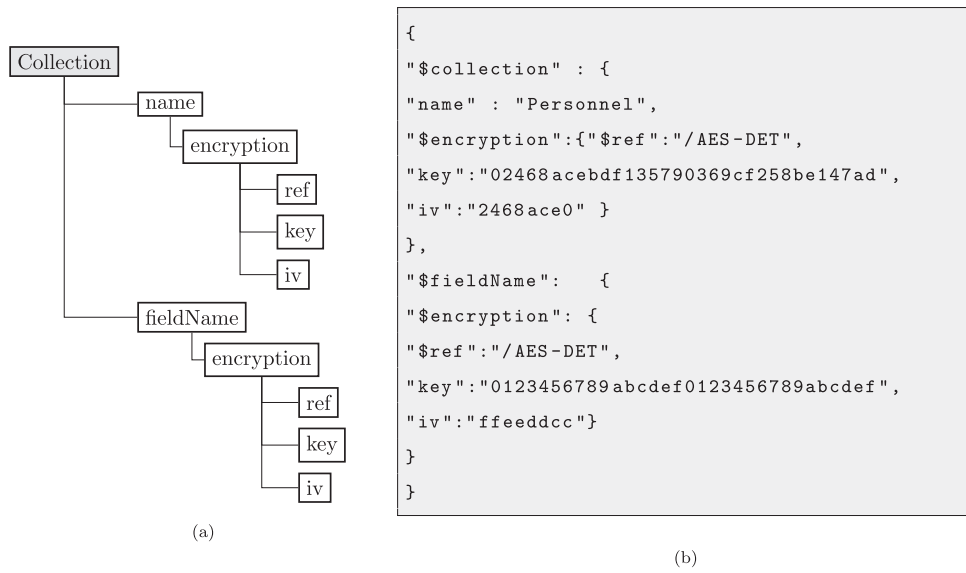


Fig. 3. The structure of a collection: (a) The chart outlines the structure of a collection containing the name of collection and name of all attributes which are considered as a meta-data, and should be protected with proper cryptographic module. (b) The description of a collection and security parameters in designed JSON based language. In this specific case the Advanced Encryption Standard in deterministic (AES-DET) mode with a 128-bit key and an *initialization vector* (IV) is assigned to encrypt the name of the collection and the fields name.

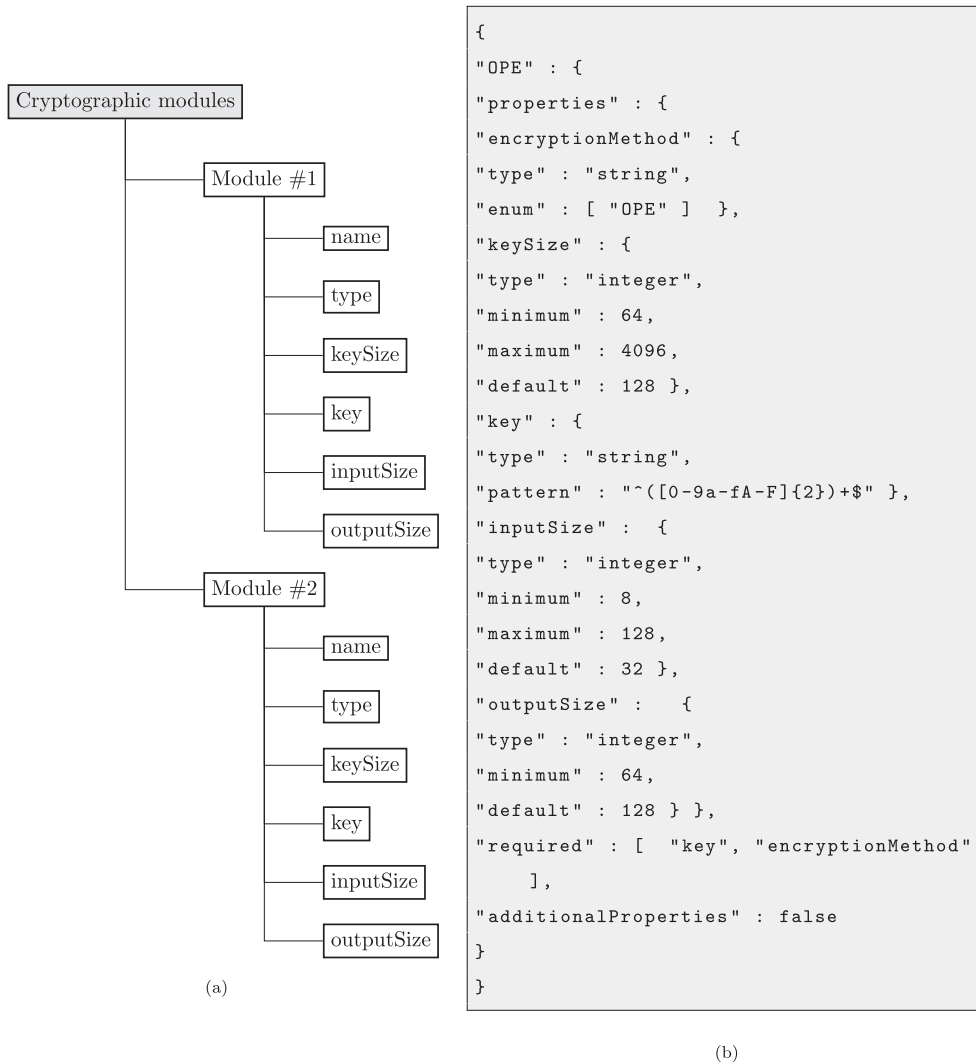


Fig. 4. The structure and function of Cryptographic modules: (a) The *Security Plan* with the second section, the cryptographic module, expanded. The attributes included for each module are: name, type, key size, key, input and output size. (b) The OPE encryption including the cryptosystems and their attributes. The proxy applies these modules using the key-value pairs (KVP).

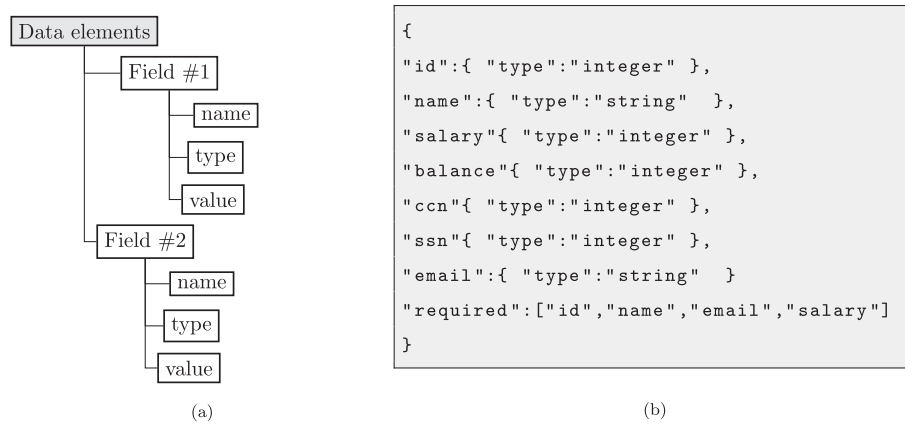


Fig. 5. Structure and description of *Data element*: (a) The chart outlines the structure of *Data elements* containing attributes of data elements such as name, type and value for of collection and name and then introduces security parameters for each data element. (b) The data element section of a sample database which is represented in designed notation. A data item has 7 fields: id, name, salary, balance, ccn, ssn, and email. The id, name, email, and salary are required fields.

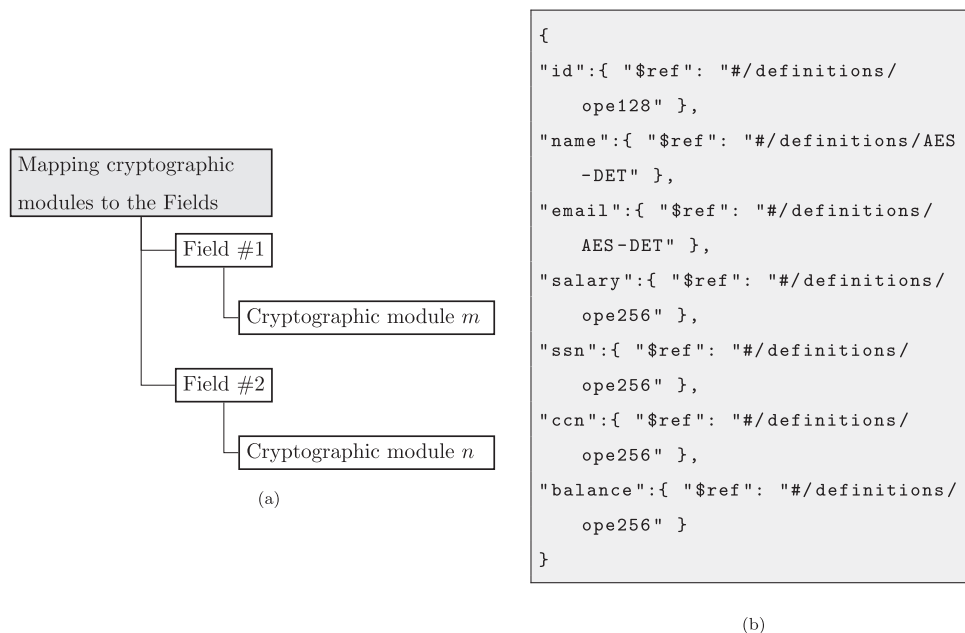


Fig. 6. The structure and description of *Mapping cryptographic modules to the Data element*: (a) *Security plan* with the fourth section expanded. This section establishes a correspondence between the data fields and the cryptographic modules used to encrypt and decrypt the data fields. (b) The mapping section of the schema for a sample database with 7 fields. For example, the *id* and the *name* will be encrypted with *OPE* 128 bit and *AES-DET*, respectively.

displayed in Fig. 5b displays data elements and its JSON description for previous example. To ensure the desired level of security the security plan should provide the description of all sensitive data elements of database in third section of security plan.

5.1.4. Mapping cryptographic modules to the fields

The last section of security plan specifies all cryptographic modules for all sensitive data fields. Fig. 6 and the listing presented in Fig. 6b shows the mapping of the cryptographic modules and the corresponding JSON format for a sample application.

The method presented in this paper can be easily extended to the other NoSQL data models discussed in Section 2. Fig. 7 shows how this extension from the key-value pair to the document store model can be carried out.

5.2. Query and data validation

The proxy validates the data and the query as a JSON-formatted input with the reference security plan. Then the proxy enforces the crypto-primitives and generates new query following the NoSQL query semantics. During this process the proxy applies to each field the cryptographic modules. Finally, the proxy forwards the newly encrypted query/data to the NoSQL database server. Fig. 8 depicts the schema validation process.

For better illustration, consider listings depicted in Fig. 9a as an input data; after running validation process the output is generated (see Fig. 9b). The output of validation process is a single file which contains descriptive information for data and meta-data in designed format and ready to execute on the *SecureNoSQL*.

The output of validation process is a single file containing descriptive information for data and meta-data expressed in the

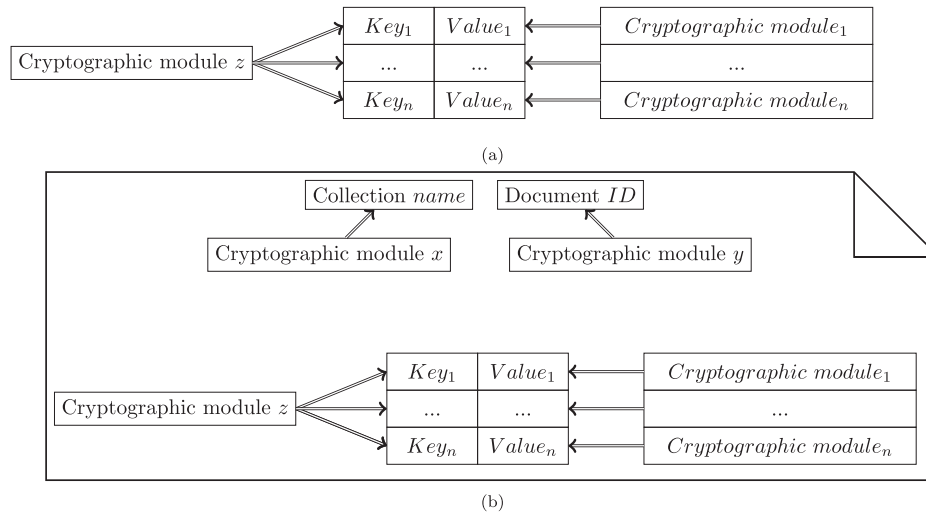


Fig. 7. SecureNoSQL applied to: (a) The key-value data model; Key_1, \dots, Key_n are all encrypted using the cryptographic module z while the corresponding values, $Value_1, \dots, Value_n$ are encrypted with cryptographic modules $1, 2, \dots, n$, respectively. (b) The document store data model; the meta-data such as collection name encrypted as well as attributes with assigned cryptographic modules.

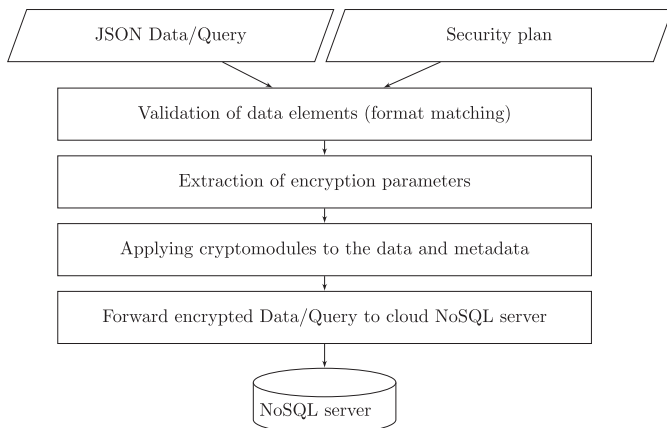


Fig. 8. The validation process of input data against security plan in the client side.

Table 1
The overhead of encryption for several encryption schemes.

Database	Plain	OPE64	OPE128	OPE256	OPE512
Size (MB)	170	430	508	662	1000

required format and ready to execute. The output of validation process for the example is illustrated in Fig. 9b. As it was noted earlier, the schema reflects the desired security level expressed by the security plan for the database. Table 1 shows the overhead for several parameters and crypto-primitives.

6. Processing queries on encrypted data

According to the proposed scheme, in order to process queries over encrypted data the queries should transfer to the encrypted version with respect to *security plan*, and this task is designed to be conducted in the proxy. The *security plan* discussed in Section 5, supplies the parameters of the cryptographic modules to be applied for the data elements involved in the query. Fig. 10 displays the processing and rewriting of a sample query.

For better understanding the query encryption, in Table 2 you can find some sample encrypted queries after enforcing *security plan*. As it can be seen, data elements and immediate values are

encrypted; however, the output is consistent with NoSQL semantics.

7. Integrity of data/query/response

Integrity and confidentiality are two critical components of data security. Integrity refers to the consistency of the outsourced data. The proposed integrity verification algorithm in SecureNoSQL guarantees the integrity of data/queries (see Algorithm 1 and Fig. 11). Data owner first applies encryption scheme on the documents, and then calculates Hashed Message Authentication Code (HMAC) for each one of encrypted documents. A hash value of any given document is a fixed length of 512 bit and data owner concatenates a unique document identifier (ID) with hash value and stores the results in efficient structure like HashTable which has constant look-up time $O(1)$. Next, data owner transfers the encrypted dataset to the cloud and sends HashTable containing hash values to the proxy. Once the proxy receives the query response from the server, it initiates the verification process to check the authenticity of the documents by recalculating the hash values. This process is illustrated in Fig. 11.

Algorithm 1. Document Integrity Verification Algorithm in the Proxy

Input: Plaintext query q from client application c_i

Output: Are the documents in the response authentic? Yes/No

$q_{enc} = \text{Encrypt}(q)$;

$q_{enc} \xrightarrow{\text{forward}}$ to cloud database server;

$R_{enc} \xleftarrow{\text{receive}}$ from cloud database server;

repeat

$H_d = \text{HMAC}(R_{enc}[i], key)$;
if ($\text{HashTable}[user_i] \neq H_d$) then
 return false

until (*There is a document in R_{enc}*);

return true;

In this configuration the data owners just trust the proxy (SecureNoSQL) and cloud servers are not trustworthy. Thus, a result of data integrity verification, all active attacks done by internal or external attacker will be detected by the proposed



Fig. 9. The security plan for the sample input: (a) The data element section of sample security plan. (b) The output of the JSON data validation for the sample database.

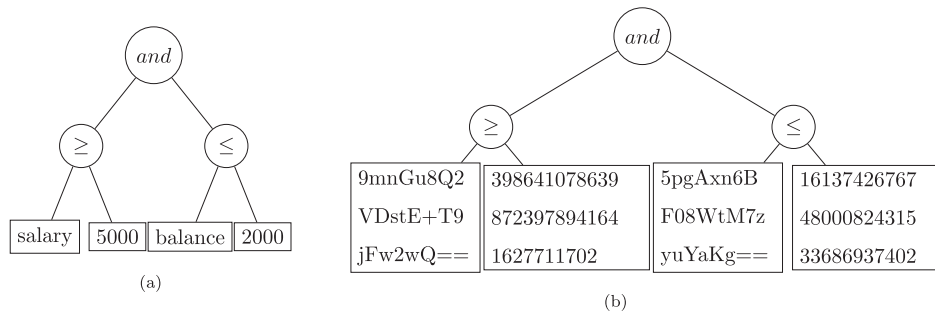


Fig. 10. The query $db.customers.find(\{salary:\{ \$gt:5000\}, balance:\{ \$lt:2000\}\})$ received from an application. (a) The parsing tree of the query. (b) The cryptographic modules applied to the data elements according to schema definition.

approach. The message authentication code (MAC) is created by using the keyed Hash Message Authentication Code (HMAC) as rephrased in Eq. (2).

$$H_{MAC}(K, document) = H((K \oplus okeyPad) || H((K \oplus ikeyPad) || document)) \quad (2)$$

Where:

H represents the hash function \oplus is the XOR operator

okeyPad is one-block-long outer pad **ikeyPad** is one-block-long inner key pad

Algorithm 2 presents the pseudo-code of the HMAC function for a block size of 64 bytes. The computed hash values with correspondent document's unique identifier can be stored in the form of key-value pair in a hash-table, thus allowing the proxy to carry the lookup in constant time during the verification process.

Algorithm 2. Keyed Hash Message Authentication Code (HMAC) generation

Input: Document: d, user key: k

Output: hash value

if $(length(key) > blocksize)$ then

key = hash(key);

if $(length(key) < blocksize)$ then

key = key || [0x00 * (blocksize - length(key))];

okeyPad = [0x5c * (blocksize)] \oplus k;

ikeyPad = [0x36 * (blocksize)] \oplus k;

return hash(okeyPad || hash(ikeyPad || d));

Table 2
Five sample queries and their corresponding encrypted version.

Q	Query	Encrypted query
1	db.customers.find({ssn:936136916})	db["k/IevnbanDMQHnkb9cRgUg=="].find({ "5pgAxn6BF08WtM7zyuYaKg==": 74172405478441908041711118833862143778})
2	db.customers.find({balance:{\$gte: 5084610},balance:{\$lte:9911843}})	db["k/IevnbanDMQHnkb9cRgUg=="].find({ "3iXpo2l8xZpW7J7TezFdeA==":{\$gte: 402982988013604629517872370128473753}, "3iXpo2l8xZpW7J7TezFdeA==" {\$lte: 785596355698717592780268633369454231}})
3	db.customers.aggregate({{\$group: {_id:null,minBalance:{\$min: "\$balance"}}}})	db["k/IevnbanDMQHnkb9cRgUg=="].aggregate({{\$group: {_id:null,EncMinBalance:{\$min: "\$3iXpo2l8xZpW7J7TezFdeA=="}}}})
4	db.customers.aggregate({{\$group: {_id:null,maxBalance:{\$max: "\$balance"}}}})	db["k/IevnbanDMQHnkb9cRgUg=="].aggregate({{\$group: {_id:null,EncmaxBalance:{\$max: "\$3iXpo2l8xZpW7J7TezFdeA=="}}}})
5	db.customers.find({\$or:{{Salary: {\$gt:516046}},{balance: {\$lt:285462}}}})	db["k/IevnbanDMQHnkb9cRgUg=="].find({\$or:{{ "9mnGu8Q2VDstE+T9jFw2wQ==":{\$gt: 40994186216785746613193244129885849}},{ "3iXpo2l8xZpW7J7TezFdeA==": {\$lt:22657430453144634679791167652174833}}}})

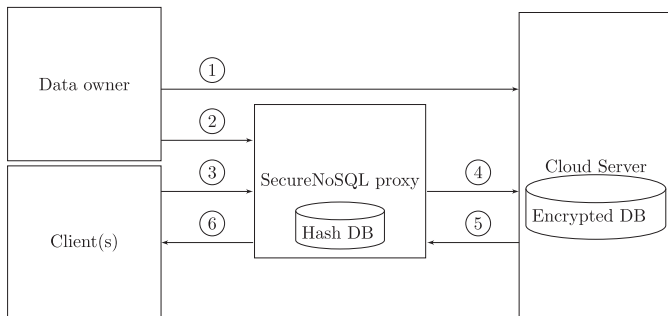


Fig. 11. (1) Data owner transfers the encrypted database to the cloud server. (2) Data owner sends the Hash database to proxy. (3) Clients send plain queries to the proxy. (4) The proxy translates queries to the encrypted version, and forwards them to the cloud server. (5) The cloud server returns the query response set. (6) The proxy runs a hash verification process on the query response set, and then based on the result either forwards to the decrypted response or reports integrity violation to the client.

8. Results and discussion

The response time of a query to an encrypted *NoSQL* database has several components:

1. the time to encode the query in *JSON* format;
2. the time to encrypt and decrypt the data;
3. the communication time to/from the server;
4. the database response time.

For our experiments we first created a sample database with one million records and then determined the overhead of searching an encrypted database. To do so we measured the database response time for queries when the records were unencrypted versus when

records were encrypted. Then, we measured the encryption and the decryption time for different sizes of the ciphertext. We wanted to isolate the different components of the response time dominated by the communication time.

The environment used for testing was set up on the Linux operating system. We chose *MongoDB* (Dede, Govindaraju, Gunter, Canon, & Ramakrishnan, 2013), classified as a *NoSQL* document store database 3.0.2. The random data generator in *JS*, *PHP*, and *MySQL* format was generated by using a tool (Keen, 2016) to generate a one million record plaintext data set. Each record had seven different data fields including *name*, *email*, *salary*, as shown in Listing 9b.

We applied OPE 64, 128, 256 and 512 bit to numeric data type, and the *AES-DET 128* bit for the string data type of the plaintext data set and generated four encrypted data sets of one million records each. Finally, we uploaded the five datasets and created five *MongoDB* databases, one with the unencrypted data, and four with the encrypted data. Once the *MongoDB* databases were created we run several types of queries including equality, greater than, less than, greater than or equal to, less than or equal to, and OR logical operations.

The experiments to measure the query time must be carefully designed. To construct average query processing time each experiment has to be carried out repeatedly. We noticed a significant reduction of database management response time after the first execution of a query, a sign that *MongoDB* is optimized and caches the results of the most recent queries. A solution is to disable the cache, or if this is not feasible, to clear the cache before repeating the query. Another important observation is that modern processors have a 64-bit architecture and are optimized for operations on 64-bit integers. For three of the five types of queries, *Q2* (Range query), *Q3* (equality), and *Q4* (logical), database response time is slightly shorter for the encrypted database than for the unencrypted one

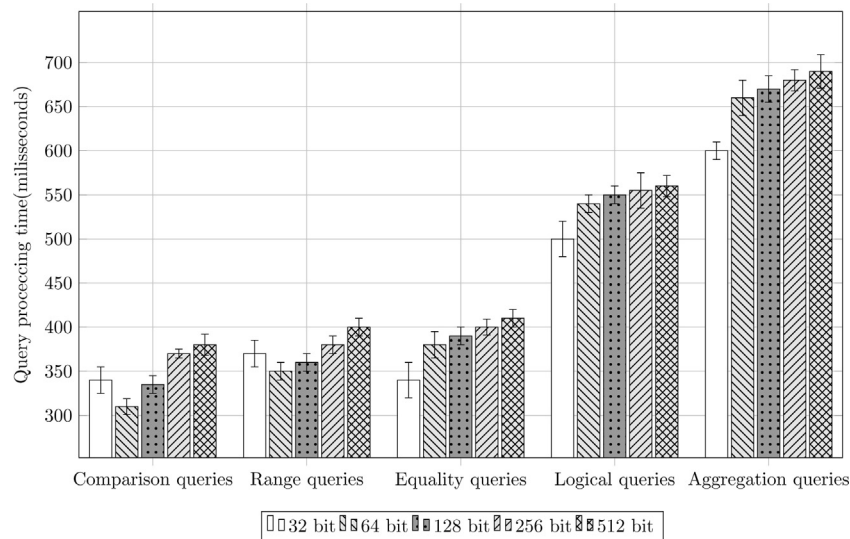


Fig. 12. The query processing time in milliseconds (ms) for the unencrypted database and for the encrypted databases when the 32-bit keys are encrypted as 64, 128, 256 and 512-bit integers.

Table 3

The query processing time in milliseconds (ms) for the plaintext and for the ciphertext. 32-bit plaintext integers are encrypted as 64, 128, 256 and 512-bit integers. The record count gives the number of records retrieved by each one of the five types of queries, Q1–Q5.

Query type	Number of matching record(s)	32-bit plaintext	64-bit ciphertext	128-bit ciphertext	256-bit ciphertext	512-bit ciphertext
Q1: Comparison	461,688	340	310	355	370	380
Q2: Equality	1	340	380	390	400	410
Q3: Range	991,225	370	350	360	380	400
Q4: Logical	551,380	500	540	550	555	560
Q5: Aggregation	1	600	660	670	680	690

when the keys are 32-bit integers. A plausible explanation for this is most likely related to the cache management.

The results reported in Table 3 and in Fig. 12 show the database response time for the five MongoDB experiments. Each query was carried out 100 times with disabled query cache and the average query response in milliseconds was calculated. We also measured the encryption and the decryption time and the results are reported in Fig. 13. The measurement process was automated, and it was running under the control of a script which generated the data and reported the processing time.

Our measurements show that the response time of the NoSQL database management system to encrypted data depends on the

type of the query. The shortest and longest database response times occur for Q1 (comparison) and Q5 (aggregated queries), respectively; for these two extremes the time for the unencrypted database was almost double, but the time for encrypted databases increases only by 70–80%. As expected, the query processing time for a given type of query increases, but only slightly, less than 5% when the key length increases from 64, to 128, 256, and 512 bit.

The OPE encryption time increases significantly with the size of the encryption space; it increases almost tenfold when the size of the encrypted output increases from 64-bit to 1024-bit and it is about 10 ms for 256-bit. The decryption time is considerably smaller; it increases only slightly from 0.11 ms to 0.17 when the size of the encrypted key increases from 64-bit to 1024 bit.

Secure proxy is an important element for the proposed architecture; therefore, the potential attacks that could affect the proxy, also should be taken into consideration. In general, two major possible attacks on proxy are Denial of Service (DoS) and unauthorized access. In DoS attack, the attacker sends so many network traffic to the proxy, that the system is not capable of processing within the expected time frame. Successful DoS attacks can turn the proxy to a bottleneck of the system. In unauthorized access attacks, attackers use a proxy to mask their connections while attacking the different targets.

Several solutions exist for improving the security of proxy against DoS attacks and reducing the consecutive impacts, including blocking the undesired packets or using multiple proxies with load balancers. Moreover, for prevention of unauthorized access attacks, it is required to use best fit authorization to access the proxy. User authentication based on group membership with different authorizations are the best practical solutions.

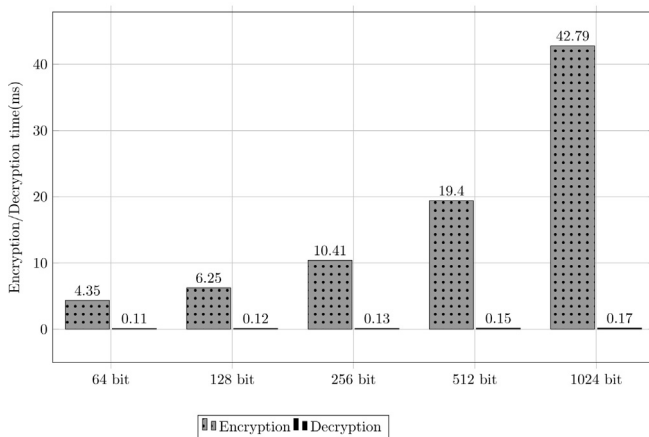


Fig. 13. Execution time of the OPE module when the key is encrypted as 64, 128, 256, 512, and 1024 bit.

