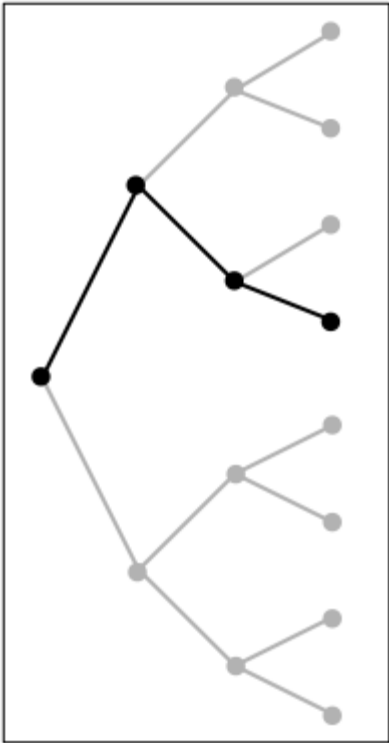


# A Method for Searching Encrypted Ranges Using Comparator Values Generated by Binary Trees

Justin Parr

November, 2022

JustinParrTech.com



J.P. TECH  
(-)

# Table of Contents

1 Abstract.....	3
2 Overview / Executive Summary.....	4
3 Encryption and Databases.....	5
3.1 Impact to Database Indexes.....	8
4 Comparators.....	9
4.1 Generating Comparators.....	9
4.1.1 Irregular Distribution Prevents Data Leakage.....	12
4.2 Tree Capacity and Scaling.....	14
4.3 Tree Balance, Refactoring, and Collisions.....	16
4.3.1 Refactoring.....	16
4.3.2 Collisions.....	17
4.4 Searching Using Comparators.....	19
4.4.1 Complex Queries.....	22
5 Conclusion.....	23



# 1 Abstract

A ‘comparator’ is a numerical value that is congruent to a data element’s ordinal value, but generated without disclosing any information.

One or two relational search terms can be passed to a database, and by converting the search terms to comparator values, range comparisons can be made on encrypted values based on comparator relationships to the underlying data.

Comparators are generated using a binary tree, which preserves the ordinal relationship of each data element, and can then be quantified as an integer value.



## 2 Overview / Executive Summary

Although it's widely agreed that field-level data encryption provides the best protection against data breaches, it also limits an application's ability to perform ranged searches, where an inequality comparison is performed against search terms.

Although there are existing strategies which use search trees, these depend on complex key-management schemes, or trees with fixed intervals.

Another common approach is to assign artificial search keys, but the general problem with this approach is that it can result in information disclosure as well as key collisions. If the keys are regularly-spaced, this can lead to an inference about the underlying data values, and if the spacing is too narrow, this can lead to insufficient search keys when presented with a large quantity of data values for given key interval.

In the scheme proposed herein, a binary tree is used to generate integer search keys, called comparators, that are non-sequential but maintain the same ordinal relationship as the underlying data. Because comparators have no fixed relationship to each other, they don't leak any information.

### 3 Encryption and Databases

Encryption scrambles data in a predictable way, such that it can be decrypted (or unscrambled) later. This is an important tool that helps protect against a data breach, where an attacker is able to access a database where sensitive data is stored in bulk.

In an encrypted state, sensitive information such as names, addresses, birth dates, credit card numbers, bank account numbers, and contact details are worthless when stolen by an attacker.

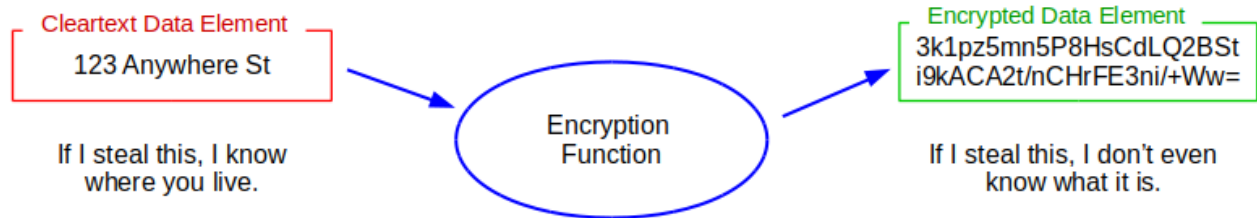


Illustration 1

Using OpenSSL to encrypt an address using the password “password” and no salt.

```
echo 123 Anywhere St|openssl aes256 -e -k password -nosalt | base64
```

To decrypt:

```
echo 3k1pz5mn5P8HsCdLQ2BSti9kACA2t/nChrFE3ni/+Ww=| base64 -d|openssl aes256 -d -k password -nosalt
```

When not encrypted (in a cleartext state), an attacker can use these details to commit fraud and identity theft, and data breaches account for billions of dollars lost to fraud every year.

Encryption can be implemented in a number of different ways, but experts generally agree that field-level encryption is the most secure.

Type of Encryption	Technology Layer	Implication
Field-Level	Application	Encryption is managed by the application, and data is only presented in cleartext when the application is using it.
Column-Level	Database	The Database Management System (DBMS) manages the encryption and decryption.  Data values are stored in an encrypted format, similar to field-level encryption, but data values are passed to the application in cleartext (although typically using an encrypted transmission channel)

Transparent Data Encryption (TDE)	Database	<p>The DBMS stores the field values in cleartext, but the database file itself is encrypted by the DBMS as it's written to disk.</p> <p>As with column-level encryption, TDE is transparent to the application.</p>
Encrypting File System (EFS) / Encrypted File System / Block-Level	Operating System	<p>The Operating System is configured to encrypt specific files, folders, or an entire partition.</p> <p>The encryption is managed by the OS, and is transparent to both the database and application.</p>
Storage-Level / Block-Level	Virtualization / Storage	<p>The Operating System runs in an environment where the hardware and / or storage are virtualized.</p> <p>Encryption is managed by the Storage Controller and / or Hypervisor, and is transparent to the OS, database, and application.</p>

As the encryption moves farther from the application, it becomes faster and less expensive because some other layer is handling the encryption and decryption.

However, each layer that presents the data in cleartext also presents an attacker with an opportunity to steal it. For example, if using block-level encryption and the OS is compromised, an attacker has access to the entire database file in cleartext, which could be exfiltrated in-tact, and dissected later.

Field-level encryption, while it has the most overhead, is also the most secure because only the application sees the data in cleartext.

Within the application, searching for a specific, encrypted value (equality search) is relatively straightforward:

- Because the encryption process is deterministic, a given cleartext data value will always result in the same encrypted version of that value.
- Assuming that the user enters a specific search value, the application encrypts it, and then performs a database search using the encrypted value.
- Again, because the encryption process is deterministic, if the encrypted value exists in the database, then the cleartext search value is the same as the cleartext version of the encrypted value found in the database.

However, a byproduct of the encryption process is that because numbers, dates, and strings are scrambled, they lose their ordinal relationship. This makes inequality searches impractical unless the application decrypts and compares every value against the search value.

For example:

- Searching ranges of encrypted birth dates, account numbers, or account balances
- Searching ranges of encrypted health data, such as blood pressure or heart rate
- Searching ranges of encrypted text data, such as names beginning with “J”

What would normally be a simple inequality search, such as “WHERE birthDate > 1/1/1990” is no longer feasible because every birthDate in the database is scrambled, and has no ordinal relationship to the search value. Therefore, each birthDate must be retrieved by the application, decrypted, and then compared to the search value, ‘1/1/1990’. The same is true of ranged text searches, for the same reason.

For clarity, we will differentiate the two following, unrelated terms:

- **Encryption Key:** Allows an encryption algorithm called a cipher to encrypt or decrypt data. The encryption process takes the data and key as input, and scrambles the data in a deterministic way, so that it can be unscrambled later using the appropriate key. If the wrong key is specified during the decryption process, the data remains scrambled.
- **Search Key / Database Key:** Allows a database engine to identify a specific row or rows based on a unique value. For example, a “people” table might have a unique key called “PersonID”. Even if there are multiple people with the same name, each would each have a unique PersonID. To tell the database to retrieve, update, or delete a specific person, you would specify the PersonID rather than the name.

Although there are existing strategies which use search trees, these depend on complex encryption key-management schemes, or trees with fixed intervals. For example, MRQED uses an n-dimensional lattice with n search trees, each with a fixed encryption key distribution. Other schemes use hidden vectors.

Another common approach is to assign artificial search keys, but the general problem with this approach is that it can result in information disclosure. As an example, given sequential search keys, and knowing some of the underlying associated values allows an attacker to deduce other values by induction. For example:

If we have encrypted street addresses, where k is the set of search keys,  $k_{10}=1230$  and  $k_{20}=1232$ , then it's reasonable to assume  $k_{15}=1231$ .

The other problem with this approach is that, despite the interval gap, there is no way to predict the relative ordinal value of new data added to the database, whose quantity might exceed the gap. For example:

If we have encrypted street addresses, where  $k_{10}=1230$  and  $k_{20}=1250$ , what happens if we add all of the street addresses between? There are 19 possible data values in the range, but only 10 search keys in the key interval, which leads to the possibility that this scheme could have to deal with 9 data values without search keys, or would have to re-allocate and shuffle search keys, etc.

### 3.1 Impact to Database Indexes

Each data table within a database may have one or more indexes, and the purpose of an index is to maintain a sorted, searchable list of pointers to the data. This can be used to speed up searches and sorts for frequently-used columns, especially when searching key fields that are used for joining tables.

Of course, the Database Management System (DBMS) doesn't *physically* rearrange the data – an index works like a linked list that maintains the proper sequence when the DBMS is asked to return data from that particular table.

In addition to facilitating equality searches (a search for a specific value), indexes speed up ranged searches because each index maintains a specified order for the affected data elements.

For example, creating an index on a column called 'date' would allow the DBMS to quickly return a list of dates within a specific range without having to scan each row of data (called a table scan) which is both expensive and slow. Instead, the DBMS consults the index, which already maintains a sort order for 'date', and can simply do a binary search within the index to find the top and bottom of the range, and then return all the rows in between.

Likewise, an application can request the most recent transaction by asking the database for the largest number from an integer column called 'transaction\_ID', and the DBMS finds this almost instantly by simply jumping to the end of the appropriate index.

However, if the indexed column is encrypted, this largely negates the value of the index because the values are no longer sorted properly.



## 4 Comparators

In the proposed scheme, every unique encrypted data value has a corresponding, unique, integer comparator value. The purpose of the comparator is to maintain the ordinal relationship of the underlying data.

Assuming that encryption is handled within the application, a user passes one or more search parameters to the application. The application generates a comparator for each search term, and then performs a regular database query using the comparator values. Because the comparators maintain an ordinal relationship to the data, this allows the database to perform both equality and ranged searches.

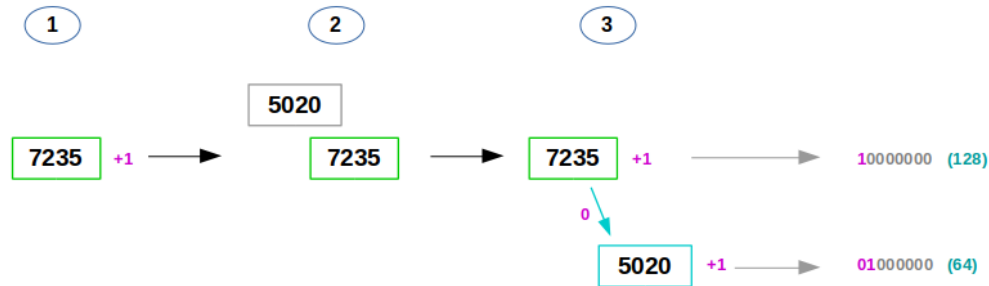
As a result of the query, the database returns a data set containing the relevant comparator values for each record, along with the corresponding encrypted data elements, which are then decrypted as needed by the application.

### 4.1 Generating Comparators

Comparators are created using a binary tree whose nodes are the underlying cleartext data values. At each step, the application retrieves an encrypted data element from the database, decrypts it, and performs a comparison.

- The application selects a data element to be the “pivot”, which is the first element in the tree. The pivot is assigned the address of “1”.
- As new elements are added, the binary tree is built in the usual fashion.
  - If the new value is greater than the decrypted node, the application descends to the right-hand child node.
  - If the new value is less than the decrypted node, the application descends to the left-hand child node.
  - If the new value matches a node, it shares the same address as the matching node.
  - If the application finishes descending in to the tree without finding a matching node, the new value gets added as a child of the last node it visited – either as the right-hand or left-hand child based on whether it is greater or less than the decrypted node.
- The address of each node is the path taken through the tree, plus “1”.
  - Starting at the pivot, the address is empty.
  - For each left-hand turn (value is less than the node), a “0” is added to the address.

- For each right-hand turn (value is greater than the node), a “1” is added to the address.
- A “1” is added to the end of the address.
- **The resulting comparator is the address, left-justified in an n-bit integer field.**



*Illustration 2*

*The initial process of building an example binary tree. Data values are shown in cleartext, but would be encrypted in the database and only visible to the application.*

In Illustration 2, we start with a pivot value, 7235 (1). The next value, 5020 is evaluated against the pivot (2). In (3), 5020 is less than 7235, so it gets added as the “left” child node. Each node’s address is the path through the tree, plus “1”.

- Our pivot, 7235, has address “1”, which is “” (blank) with “1” appended.
- Our next data value, 5020, is less than 7235, and gets added as the “left” child node. Tracing the path through the tree for 5020, we start at the pivot (current address=“”), go left (current address=“0”) and then append “1” (final address=“01”)
- If we left-justify both of these addresses in an 8-bit integer field, we get 128 (10000000) and 64 (01000000) respectively

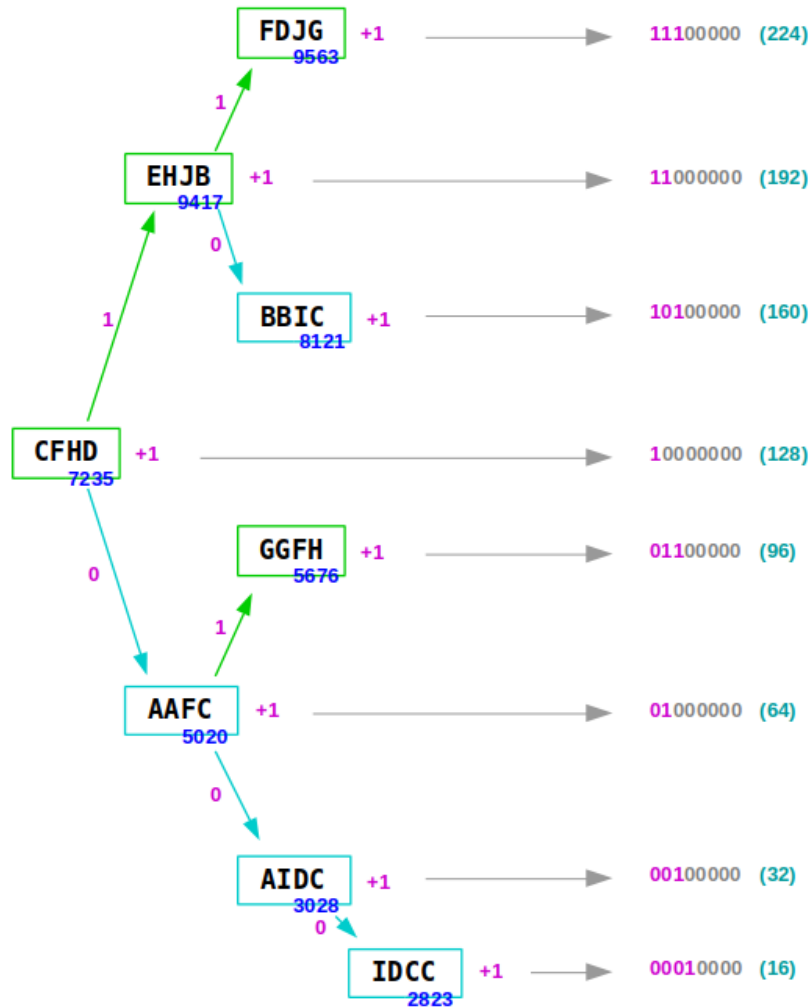


Illustration 3

After a few more nodes have been added, this is what the tree looks like.

In Illustration 3, the data values have been replaced with their encrypted counterparts. In reality, all database operations are conducted in an encrypted state.

Looking at the comparator values, the necessity of adding a “1” at the end of the address becomes clear. Without it, the left-hand three nodes would all have the same address:

Data Element	Tree Address	Without “1”	With “1”
5020	0	00000000 = 0	01000000 = 64
3028	00	00000000 = 0	00100000 = 32
2823	000	00000000 = 0	00010000 = 16

The trailing “1” ensures that each comparator is discreet.

### 4.1.1 Irregular Distribution Prevents Data Leakage

As we’ve seen, schemes that use a fixed search key spacing can leak data by allowing an attacker to analyze search keys, and then interpolate the underlying data values.

This is possible because the search keys maintain a linear relationship to each other, and to the underlying data.

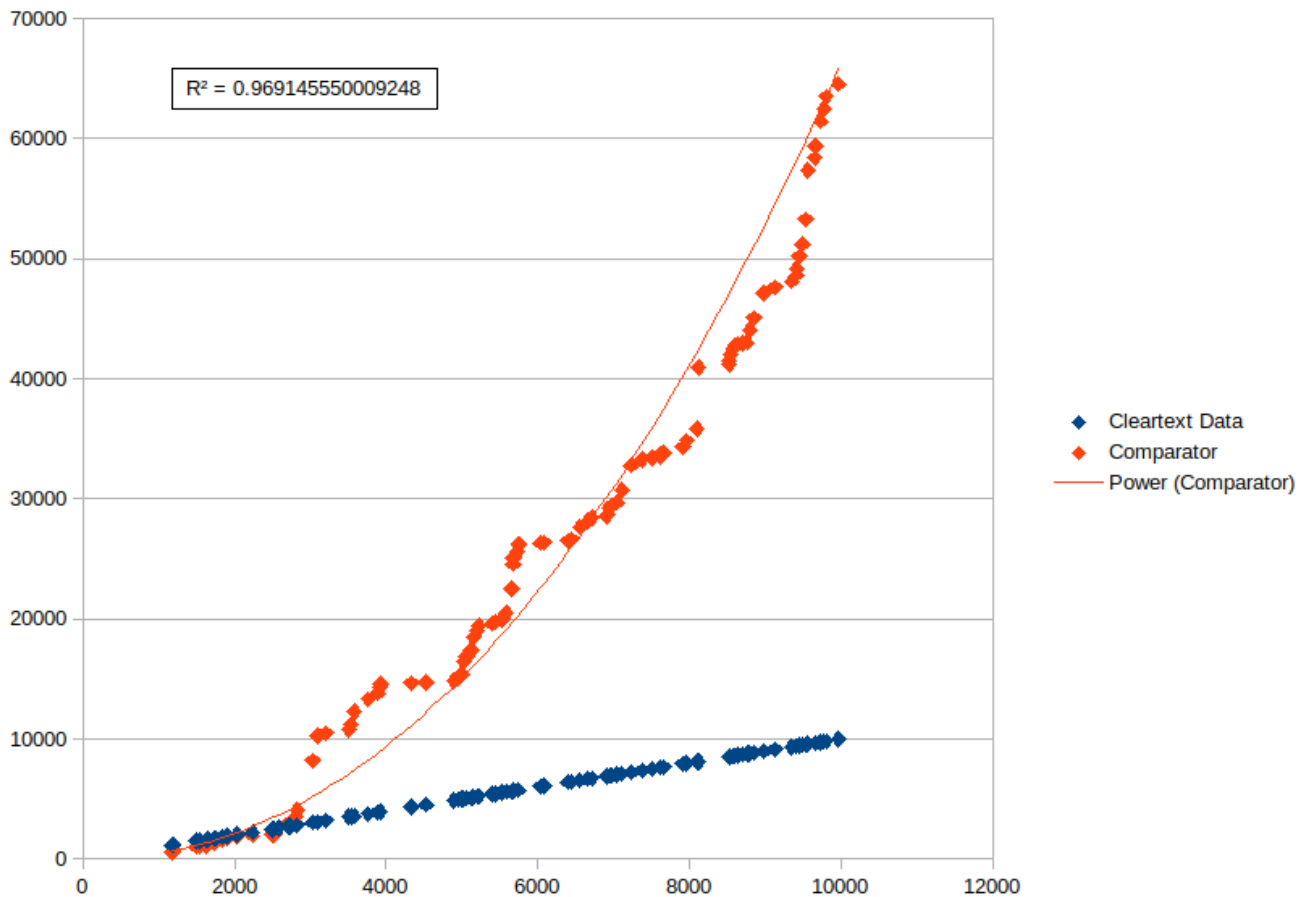


Illustration 4

For cleartext values  $N$ , the distribution of comparators progresses at approximately  $N^2$ .

However, just as with the effort involved in determining large prime factors of a composite number, determining individual data values based on comparators is polynomially-complex, because the distribution isn't consistent.

Comparators, despite being ordered, occur based on their relationship within the tree rather than at specific intervals, making them much harder to predict.

Although the distribution of comparators progresses at approximately  $n^2$ , any new comparator can appear anywhere between two existing ones. Rather than being exactly between them, it could be arbitrarily close to either one. And, despite this, the scheme always allows new comparators to be inserted between two others, regardless of how close they are.

This works because the tree address is left-justified within a bit field, making it possible to use subsequent bits. Despite being an integer, the comparator behaves similar to a decimal, where you can always create a new number that's between two others by simply adding another decimal digit.

For example, creating a number between 0.345 and 0.346 can be accomplished by adding a digit. All of the numbers between and including 0.3451 and 0.3459 are between the two. And for each pair, there are an infinite number of numbers between them.

The comparator scheme works in a similar fashion, for as many bits that exist within the integer bit field.

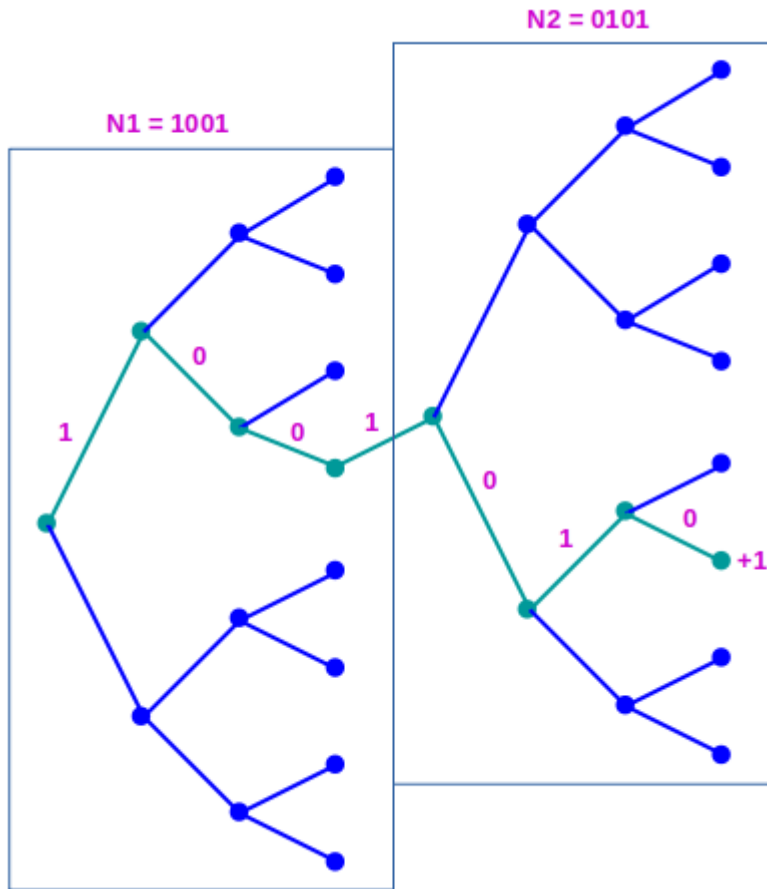
## 4.2 Tree Capacity and Scaling

If  $n$  is the number of data elements predicted to occur in the set, the size of the largest comparator is approximately  $n^2$ . Because the tree is binary,  $\log_2(n^2)+1$  bits are required – the extra bit accommodates the “1” address terminator.

Bit Field Size	Approximate Number of Elements
8 (7)	11
16 (15)	181
32 (31)	46,340
64 (63)	3,037,000,500

Counter-intuitively, the size of the comparator only depends on the number of elements, not the magnitude of the underlying data value. Therefore, even large data elements such as long strings can be represented concisely by integers.

Although either a 32-bit or 64-bit comparator would be suitable for most applications, the comparator could consist of multiple integers or a string of integers. Chaining an arbitrary number of integers would allow for a virtually unlimited tree size.



*Illustration 5*

*Here, two nibbles (4-bit field) are chained together, which is effectively a single 8-bit field.*

For example, a chained 128-bit (127) comparator would consist of 4 chained 32-bit integers or 2 chained 64-bit integers, and would be capable of representing about 13 quadrillion data elements.

## 4.3 Tree Balance, Refactoring, and Collisions

If a specific area of the tree becomes overpopulated, this results in excess depth, and could lead to a condition where new values can no longer be inserted.

### 4.3.1 Refactoring

One method used to address an unbalanced tree is refactoring.

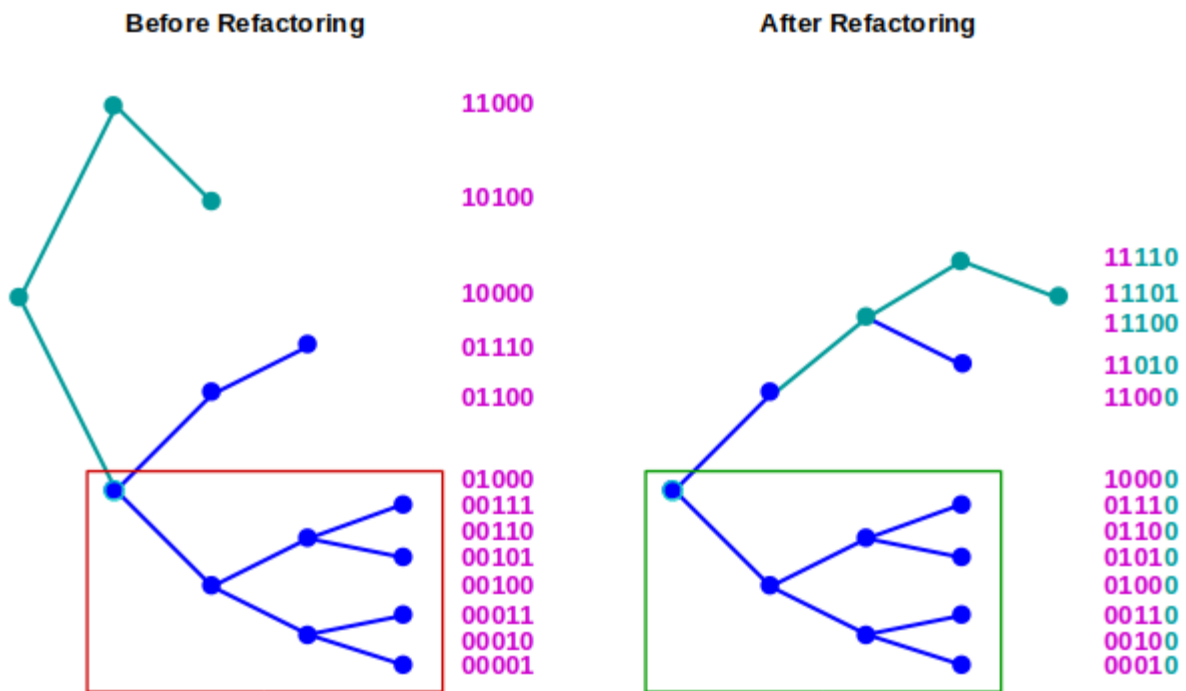


Illustration 6

On the left, the tree is unbalanced, resulting in crowding within the red area – in this example, there is no room to add new comparitors, which would exceed the bit (tree) depth.

On the right, the tree has been refactored by selecting a new pivot. Because each node in the crowded section is one position closer to the new pivot, this effectively doubles the capacity within that portion of the tree. The previous pivot and its child nodes have been re-grafted to the new pivot.

As there are many algorithms for accomplishing this, and all are well documented elsewhere, a discussion of the specifics will be excluded here. In general, these algorithms select a new pivot (root), and then re-graft orphaned nodes within the new sub-tree.



However, the end result of refactoring is that the tree is more balanced, where each node has approximately the same number of descendants on both child branches. Because the descendant nodes are better distributed, this results in a more shallow overall depth.

This approach is elegant, because refactoring can be performed using the old comparators without having to actually decrypt any of the nodes' data values. Refactoring results in new tree addresses, and therefore new comparator values for all data nodes, but the old ones are simply discarded at the end of the refactoring process. Because decryption is not required, the refactoring process can be conducted by the database engine, for example, within a stored procedure.

This is similar to the process for re-indexing a database table, because database indexes also use a tree structure. In fact, the final step would be to rebuild any database indexes using the new comparator values, so that the new database index is also balanced properly.

### 4.3.2 Collisions

Without refactoring, the other option is to accept **collisions**, where a single comparator could relate to multiple data values, even though the underlying data values are different.

In some cases, it may be acceptable to allow collisions, understanding that this affects how queries behave.

For example, given the following data and associated comparators, we can analyze how allowing collisions would affect each type of operation.

Comparator (Collisions)	Comparator (No Collisions)	Cleartext Data Values
1	1 2 3	1111 2222 3333
4	4 5	4444 5555
6	6	6666
7	7 8 9	7777 8888 9999

In the table below, “c” is used for the comparator, and “A” and “B” represent data values.

Operator	Without Collisions	With Collisions
Equality (A=B)	Only exact matches are returned.	Values within a small range of A may be returned.

	<p>For example, WHERE c=4 ONLY returns 4444.</p>	<p>For example, WHERE c=4 returns 4444 and 5555, but WHERE c=6 only returns 6666.</p> <p>If the application must make provisions for situations where the database returns multiple records.</p>
<p>Greater Than (A&gt;B) Less Than (A&lt;B)</p>	<p>Works as expected.</p> <p>For example, WHERE c&gt;1 returns 2222, 3333, 4444, and so on.</p>	<p>Could result in gaps.</p> <p>For example, WHERE c&gt;1 skips 2222 and 3333, and returns values starting at 4444</p> <p>To avoid this, the application must issue &gt;= (Greater or Equal) or &lt;= (Less or Equal) and then filter out equalities.</p> <p>Likewise, WHERE c&gt;=3 would revert to c&gt;=1, and would improperly include 1111 and 2222, despite the fact that the application doesn't expect these values.</p>
<p>BETWEEN</p>	<p>BETWEEN works like a pair of inequality comparisons:</p> <p>A BETWEEN B<sub>1</sub> AND B<sub>2</sub></p> <p>Is the same as:</p> <p>(A &gt;= B<sub>1</sub>) AND (A &lt;= B<sub>2</sub>)</p>	<p>Just as with &gt;= and &lt;=, some extra data values may be improperly included.</p>
<p>IN</p>	<p>IN works like multiple equality comparisons:</p> <p>A IN (B<sub>1</sub>, B<sub>2</sub>, B<sub>3</sub>)</p> <p>Is the same as:</p> <p>(A=B<sub>1</sub>) OR (A=B<sub>2</sub>) OR (A=B<sub>3</sub>)</p>	<p>Just as with =, some extraneous data values may be improperly included.</p>

Depending upon the nature of the application, collisions and their associated limitations may be acceptable, in order to avoid refactoring or integer chaining.

## 4.4 Searching Using Comparators

In the illustration below, the original data values are shown in blue for clarity, but in a real database, the cleartext data wouldn't be present.

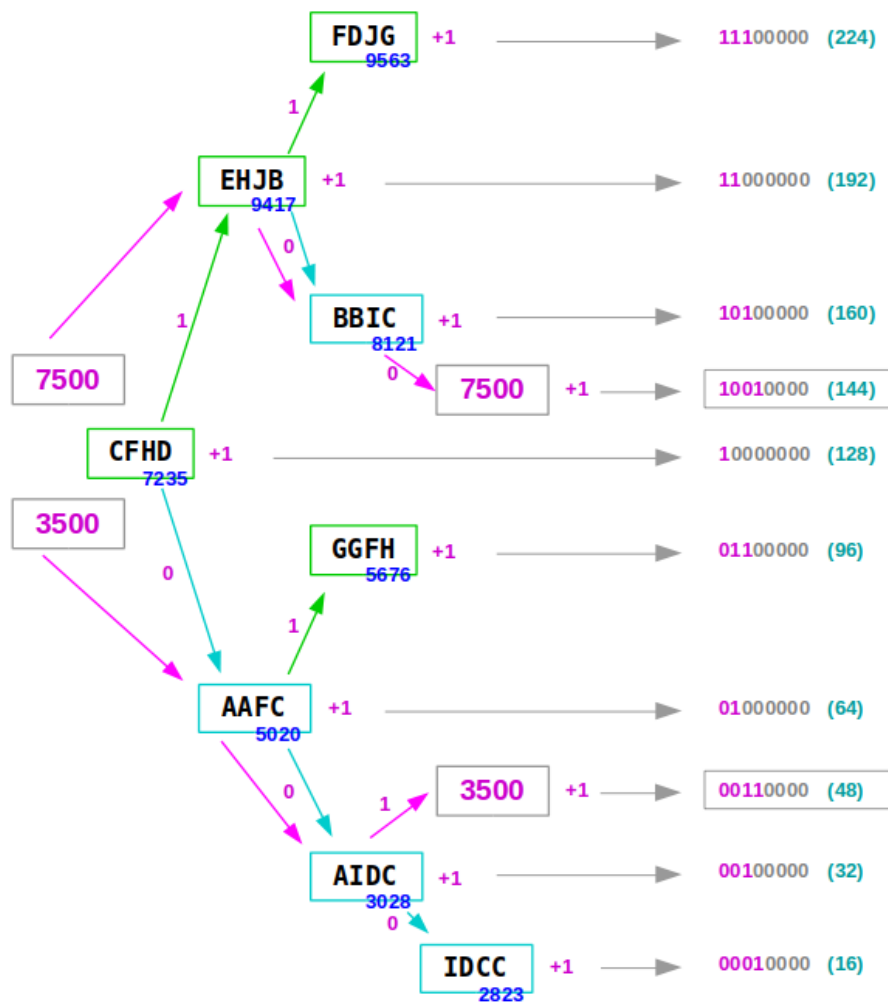


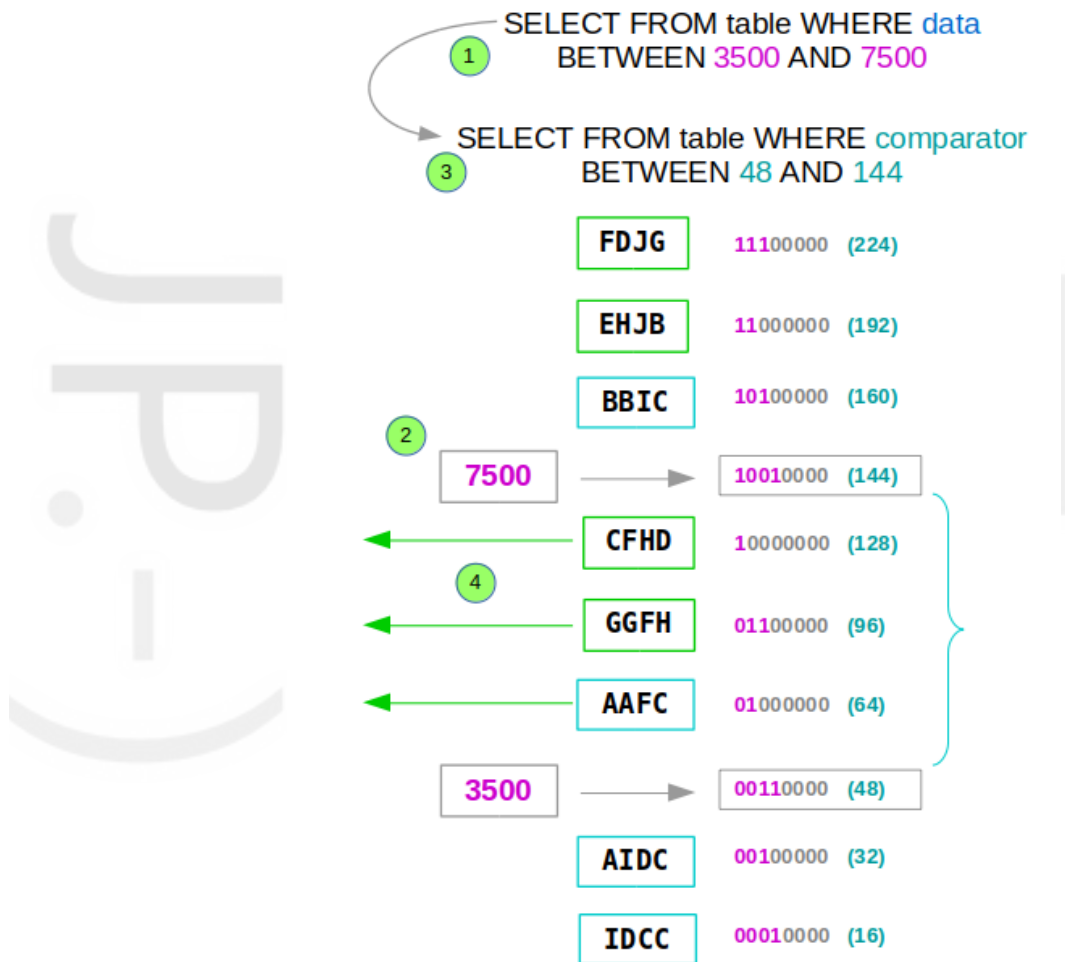
Illustration 7

Creating comparators from search terms follows the same process as adding a new node, but without permanence. When the search is complete, the search comparators are discarded.

In Illustration 7, we create comparators for two numeric parameters that are search terms for a BETWEEN operation in a WHERE clause.

The resulting search comparators, 48 and 144, maintain an ordinal relationship to the underlying data, even though we can't see the actual values.

Rather than having to decrypt every value in the database, we simply let the database perform a normal search on the comparators.



*Illustration 8*

*The application replaces database search terms with search comparators, and then the database executes the query.*

In Illustration 8, the application wants to run a query for data elements within a specific range. These could be confidential performance data for a new engine, or maybe protected health data.

1. The application needs to run a query using normal terms.
2. The application converts the search terms to search comparators.
3. The application writes a new query using search comparators.
4. The database executes the comparator query, and returns the resulting encrypted data elements to the application.

Once search values are converted to search comparators, the database query process behaves very similar to querying the underlying cleartext data.

However, the data returned to the application by the database could vary based on context.

Result Set	Use Cases
<p>A list of comparators, assuming that comparators have a UNIQUE constraint (no duplicates are allowed)</p>	<p>If the application needs to identify certain records without actually acting on them, the application can act upon the comparator as if it was the underlying data, without having to decrypt it.</p> <p>For example, perhaps you want to send birthday cards to everyone born this month – the application can act on all records whose birth date is within a range without decrypting the actual birth dates.</p> <p>Another example is updating a group of accounts whose account number falls within a certain range, without ever decrypting the account numbers.</p>
<p>A data set containing encrypted data values</p>	<p>If the application needs to act upon the the data directly, the database can return the encrypted values, which are then decrypted as needed by the application.</p> <p>For example, a healthcare application might look for a patient’s encrypted blood pressure readings that are within a certain range. Because the readings themselves have specific meaning, the database would return the encrypted values, which would then be decrypted by the application as needed.</p>
<p>The result of an aggregation function, such as MIN, MAX, or COUNT</p>	<p>Aggregation functions can be used to find the oldest account, or simply tally the number of customers within a certain age range.</p> <p>In another example, an application can answer a simple question such as “is the account holder a minor” by using</p>

	comparators. The application takes today's date, subtracts 18 years, and generates a comparator. If the comparator for a person's birth date is less than this value, they are a minor.
--	---

In addition to the above, and similar to data masking, comparators can be passed to downstream applications without compromising the underlying data.

#### 4.4.1 Complex Queries

An application can query multiple, protected parameters, assuming that each one has an equivalent comparator.

For example, an application might need to know which account holders are within a certain age range, and whose credit scores are within a certain range. Both of these are protected elements, and therefore encrypted within the database.

However, if both birth date and credit score each have comparators, a complex query allows searching within an overlapping range:

```
SELECT FROM accounts WHERE c_dob BETWEEN c1,1 AND c1,2 AND c_credit_score BETWEEN c2,1 AND c2,2
```

In another example, if an account holder's address is encoded with latitude and longitude, each encrypted, and each having comparators, then a complex query can be used to search for customers within a specified geographical area:

```
SELECT FROM accounts WHERE c_latitude BETWEEN a1 AND a2 AND c_longitude BETWEEN b1 AND b2
```

And, as with data masking, if key values have unique comparators, joins can be performed using only the comparators:

```
SELECT a.* FROM tbl1 a INNER JOIN tbl2 b ON a.c_account = b.c_account
```

This would come in to play, for example, if joining an "accounts" table to a "transactions" table without having to decrypt individual account values.

## 5 Conclusion

Field-level encryption offers the strongest level of encryption, but creates challenges when performing ranged searches.

Because comparators are left-justified within an integer bit field, they maintain an ordinal relationship to the underlying data without leaking information.

Comparators, can be used to perform ranged queries without having to decrypt the underlying data, because they have the same ordinal relationship as the underlying data.

Scaling can be accomplished using integer chaining as well as tree refactoring.

As with data masking, if comparators are suitably unique, the comparator itself can be substituted for the underlying cleartext data value.

Because comparators are ordered, databases can perform ranged searches with optimum performance, because the sort order is preserved. Only the application sees the data in cleartext.

Unlike MRQED and other ranged-search schemes, comparators don't rely on key generation mechanisms, key distribution schemes, nor any fixed relationship with the underlying data.