

# Privacy Protected Shared Person-Specific Genomic Sequence of Queries for DNA Databases

Parvez Khan<sup>1</sup> and Md Ateeq Ur Rahman<sup>2</sup>

<sup>1</sup>Research Scholar, Dept. of Computer Science & Engineering, SCET, Hyderabad

<sup>2</sup>Professor and Head, Dept. of Computer Science & Engineering, SCET, Hyderabad

**Abstract-** To support large-scale medicine analysis comes, organizations ought to share person-specific genomic sequences while not violating the privacy of their information subjects. within the past, organizations protected subjects' identities by removing identifiers, like name and Social Security number; but, recent investigations illustrate that deidentified genomic information are often "identified" to named people exploitation straightforward machine-driven ways. This paper addresses the matter of sharing person-specific genomic sequences while not violating the privacy of their knowledge subjects to support large-scale medical specialty analysis comes. The projected methodology builds on the framework projected by Kantarcioglu et al. however extends the leads to variety of how. One improvement is that our theme is settled, with zero likelihood of a wrong answer (as opposition an occasional probability). we have a tendency to conjointly give a replacement operational purpose within the coordinate system exchange, by providing a theme that's doubly as quick as theirs however uses doubly the cupboard space. this time is impelled by the actual fact that storage is cheaper than computation in current cloud computing evaluation plans. Moreover, our coding of the info makes it doable for U.S.A. to handle a richer set of queries than precise matching between the question and every sequence of the information, including: (i) enumeration the quantity of matches between the question images and a sequence; (ii) logical OR matches wherever a question image is allowed to match a set of the alphabet thereby creating it doable to handle (as a special case) a "not equal to" demand for a question symbol (e.g., "not a G"); (iii) support for the extended alphabet of ester base codes that encompasses ambiguities in DNA sequences (this happens on the DNA sequence facet rather than the question side); (iv) queries that specify the quantity of occurrences of every quite image within the specified sequence positions (e.g., 2 'A' and 4 'C' and one 'G' and 3 'T', occurring in any order within the question-specified sequence positions); (v) a threshold query whose answer is 'yes' if the quantity of matches exceeds a query-specified threshold (e.g., "7 or additional matches out of the fifteen query-specified positions"). (vi) For all question varieties we are able to hide the answers from the decrypting server, so solely the shopper learns the solution. (vii) all told cases, the shopper deterministically learns solely the query's answer, aside from question sort (v) wherever we have a tendency to quantify the (very small) applied mathematics discharge to the shopper of the particular count. during this paper, we have a tendency to

gift a unique cryptologic framework that permits organizations to support genomic data processing while not revealing the raw genomic sequences. Organizations contribute encrypted genomic sequence records into a centralized repository, wherever the administrator will perform queries, like frequency counts, while not decrypting the info. we have a tendency to measure the potency of our framework with existing databases of single ester polymorphism (SNP) sequences and demonstrate that the time required to complete count queries is possible for planet applications. for instance, our experiments indicate that a count question over forty SNPs in an exceedingly information of 5000 records are often completed in around thirty min with off-the-peg technology. we have a tendency to additional show that approximation methods are often applied to considerably speed up question execution times with least loss in accuracy. The framework are often enforced on prime of existing data and network technologies in medicine environments.

**Index Terms-** DNA Databases, Cloud Security, Secure Outsourcing, Cryptography, Genomics, Bioinformatics, Databases, Large-scale systems, Data privacy, Protection, Data security, Data mining.

## I. INTRODUCTION

In this paper, we tend to propose an alternate approach to genomic information privacy protection that's supported cryptography. Our model ensures that: 1) the info utility of protected records is resembling that achieved by deidentification and 2) the info privacy is resembling that achieved by information augmentation schemes. As an overview, our model works as follows. information holders John and electro-acoustic transducer transmit encrypted versions of their records to a 3rd party's information repository. The repository administrator executes queries on behalf of Charlie the researcher while not decrypting any of the records. The results of the question square measure then sent to a 3rd party World Health Organization decrypts the aggregation of the result and sends the answer to the somebody. This design incorporates two completely different third parties for security-related advantages. There is no chance to decipher the info unless each third parties collaborate. As a result, the utilization of multiple third parties ensures that there's no single purpose of knowledge compromise. Thus, if a hacker breaks into one amongst the third party's pc systems, the hacker cannot learn the sensitive data within the encrypted records.

Recognize that tho' the info remains encrypted in the least times, the results of queries themselves will violate privacy needs. for example, if the solution to Charlie's question is specified there's just one record with DNA sequence "AATCAATGAA" and juvenile Alzheimer's unwellness, then Charlie has unambiguously pinpointed associate individual's record. Thus, it's necessary for the third party to make sure that question results, or the combination of a series of question results issued by a researcher, don't allow the triangulation of associate individual's record. This method, referred to as question restriction, is important to make sure that our framework achieves identity protection; but, this subject has been studied extensively within the info security community, and thus, we tend to neglect the presentation of query restriction during this paper. the most contribution of our model is within the analysis of encrypted genomic information. To the most effective of our data, there's no ready-to-wear product or literature that may be applied to satisfy this component of the framework. As such, this paper focuses on the cryptologic protocols that square measure necessary to create and question encrypted genomic databases. additionally, we offer experimental validation in order that in our framework, queries will be answered with efficiency for planet medical specialty applications. DNA or polymer is that the medium of long-run storage and transmission of genetic info for all trendy living organisms. Human desoxyribonucleic acid information (DNA sequences inside the twenty three body pairs) area unit non-public and sensitive personal info. However, such information is essential for conducting medical specialty analysis and studies, for instance, diagnosing of pre-disposition to develop a particular unwellness, drug allergic reaction, or prediction of success rate in response to a particular treatment. Providing a in public offered desoxyribonucleic acid info for fostering analysis during this field is especially confronted by privacy considerations. Today, the verdant computation and storage capability of cloud services allows sensible hosting and sharing of desoxyribonucleic acid databases and economical process of genomic sequences, like playing sequence comparison, actual and approximate sequence search and numerous tests (diagnosis, identity, ancestry and paternity). what's missing is AN economical security layer that preserves the privacy of individuals' records and assigns the burden of question process to the cloud. Whereas anonymization techniques like de-identification, information augmentation, or info partitioning solve this downside partly, they're not comfortable as a result of in several cases, re-identification of persons is feasible. It follows that the desoxyribonucleic acid information should be protected, not simply unlinked from the corresponding persons.

In this paper, we have a tendency to take into account the framework planned in wherever the desoxyribonucleic acid records coming back from many hospitals area unit encrypted and keep at a knowledge storage website, and medical specialty researchers area unit able to submit

mixture investigation queries to the present website. investigation queries area unit significantly fascinating for applied math analysis.

This paper provides a brand new technique that addresses a bigger set of issues and provides a quicker question interval than the technique introduced in. Our approach is predicated on the very fact that, given current valuation plans at several cloud services suppliers, storage is cheaper than computing. Therefore, we have a tendency to favor storage over computing resources to optimize value. Moreover, from a user expertise purpose of read, interval is that the most tangible indicator of performance; thence it's natural to aim at reducing it. Our technique enhances the state of the art at each the abstract level and therefore the implementation level. additional concretely:

- At the abstract level, we offer a settled theme, with zero chance of a wrong answer (as hostile an occasional probability). this offers confidence to the users that they get actual results to any or all their queries, while not impacting security.

- we have a tendency to conjointly offer a brand new operative purpose within the reference frame trade-off, by giving a theme that's double as quick as theirs however uses double the space for storing. A variant of this theme uses only one.5 their space for storing at the expense of further latency.

Moreover, our secret writing of the information makes it doable for United States to handle a richer set of queries than actual matching between the question and every sequence of the info, including:

- i. investigation the quantity of matches between the question symbols and a sequence;

- ii. Logical OR matches wherever a question image is allowed to match a set of the alphabet thereby creating it doable to handle (as a special case) a "not equal to" demand for a question image (e.g., "not a G");

- iii. Support for the extended alphabet of ester base codes that encompasses ambiguities in desoxyribonucleic acid sequences (contrary to the previous item this happens on the desoxyribonucleic acid sequence aspect rather than the question side);

- iv. Queries that specify the quantity of occurrences of every reasonably image within the fixed sequence positions (e.g., 2 'A' and 4 'C' and one 'G' and 3 'T', occurring in any order within the query-specified sequence positions);

- v. A threshold question whose answer is 'yes' if the quantity of matches exceeds a query-specified threshold (e.g., "7 or additional matches out of the fifteen query-specified positions").

- vi. For all question sorts we are able to hide the answers from the decrypting server, in order that solely the consumer learns the solution.

- vii. altogether cases the consumer deterministically learns solely the query's answer, aside from question kind (v) wherever we have a tendency to quantify the (very small) applied math outpouring to the consumer of the particular count.

II. RELATED WORKS

In our model, hospitals World Health Organization have DNA sequences don't have the computing and process capabilities to method researchers' requests, so that they all store their DNA sequences at a server (which is additionally a lot of convenient to try and do queries across all hospitals). The shoppers, World Health Organization area unit usually researchers, question the server to get statistics on the incidence of a given subsequence within the pool of DNA sequences keep on the server. as a result of the sensitivity of DNA, of these operations have to be compelled to be performed securely: the goal of securing queries is creating each the shopper and also the server unaware of specifically that sequences match the question however solely knowing the aggregative results of the question (i.e., the count).To be a lot of precise the safety model is as follows:

- Hospitals need to safeguard the confidentiality of the DNA sequences that they own and no external party has the correct to access these DNA sequences for privacy reasons. Thus, alternative parties (be it the server or the clients) ought to solely work on encrypted sequences and ne'er have access to the DNA cleartext.
- The server is associate external repository of DNA sequences provided by the varied hospitals. The server is taken into account honest-but-curious by hospitals: hospitals trust him to perform the queries requested by shoppers however they are doing not need the server to access the DNA sequences in clear.
- shoppers area unit entities approved to perform queries on the database of encrypted DNA sequences. they're solely allowed to get statistics on the database: the quantity or proportion of sequences matching a given question. The queries don't seem to be confidential and area unit processed by the server, but the server shouldn't grasp the result of the queries.
- we have a tendency to assume that none of those entities conspire.

Additively homomorphic coding is appropriate for the aim of activity count statistics on encrypted knowledge. Paillier's homomorphic coding possesses the subsequent properties: (i) it is a public key theme, which implies coding may be performed by anyone World Health Organization is aware of the general public key, whereas coding will solely be done by the matching non-public key, identified solely to a trustworthy party. (ii) it's probabilistic. In alternative words, it's not possible for associate soul to inform whether or not 2 ciphertexts area unit encryptions of identical plaintext or not. (iii) It possesses the homomorphic properties

A. Existing System

Several works specialise in protective biometric computations over genomic sequence records within the context of secure multi-party computations (SMC). Secure outsourcing may be a specific case of SMC wherever a shopper with low resources (energy, memory, CPU)

requests the service of 1 or additional outsourcing agents with abundant resources. Secure outsourcing finds a true projection within the current business models due to the proliferation of cloud-based services.

B. Disadvantages

What is missing is Associate in Nursing economical security layer that preserves the privacy of individuals' records and assigns the burden of question process to the cloud.

III. PROPOSED SYSTEM

In this paper, we have a tendency to contemplate the framework planned in wherever the polymer records returning from many hospitals area unit encrypted and keep at a knowledge storage website, and medicine researchers area unit ready to submit mixture investigation queries to the present website. investigation queries area unit notably attention-grabbing for applied math analysis.

This paper provides a brand new methodology that addresses a bigger set of issues and provides a quicker question latent period than the technique introduced in . Our approach is predicated on the actual fact that, given current rating plans at several cloud services suppliers, storage is cheaper than computing. Therefore, we have a tendency to favor storage over computing resources to optimize value. Moreover, from a user expertise purpose of read, latent period is that the most tangible indicator of performance; thus it's natural to aim at reducing it. Our methodology enhances the state of the art at each the abstract level and therefore the implementation level. a lot of concretely:

STORAGE AND COMPUTATION SCHEMES:

- Quaternary storage:
- Quaternary query:

IV. SYSTEM ARCHITECTURE

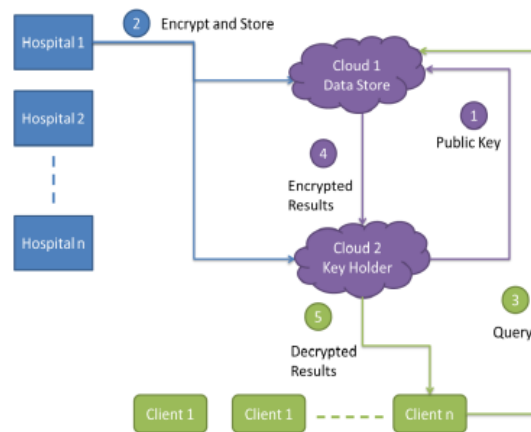


Fig.1: Framework of secure aggregate queries over encrypted DNA database

Security evaluation:

From a security perspective the framework that we have a tendency to use is comparable to [1]. each schemes square measure supported well-known security building blocks like

Paillier's cryptography, public key cryptography and isobilateral key cryptography. this {can be} why our security analysis focuses on the interconnection of those building blocks and what can get it wrong in our settings. To be additional precise:

□ The desoxyribonucleic acid sequences square measure continuously encrypted at Cloud1, therefore Cloud1 cannot access these sequences in clear. the sole entity that might decode them is Cloud2 that could be a trustworthy entity by the hospital (again with named Cloud2 to stress that it may be deployed within the cloud however it's terribly completely different from cloud1 therein it performs marginal operations (decryption oracle) and doesn't ought to store immense information, therefore it's a selected trustworthy platform within the cloud. The confidentiality of desoxyribonucleic acid sequences is so properly preserved.

□ Cloud1 additionally doesn't get ANy leak from the queries of shoppers as a result of he processes the queries in an encrypted (he cannot decode the outcome). The result's decrypted by Cloud2 that sends the result on to the shopper United Nations agency created the question through a secure channel because of the protection associations between Cloud2 and therefore the shoppers. This additionally implies that Cloud1 cannot act as a shopper and acquire the results of his own queries, unless he colludes with a true shopper, that is out of scope of our model.

□ shoppers solely acquire statistics on the amount of desoxyribonucleic acid sequences across all hospitals that match their question however they don't get the desoxyribonucleic acid sequences themselves. Learning the variety of sequences could be a leak that's acceptable in our model as this leak would happen even with the perfect model of a trustworthy entity doing all the process between hospitals and shoppers.

□ the sole entity that very has a position in our framework is Cloud2 because it owns the non-public key. we have a tendency to argue but that:

o Cloud2 could be a trustworthy entity

o Cloud2 doesn't have access to encrypted desoxyribonucleic acid sequences unless he colludes with Cloud1 or a Hospital

o Cloud2 sees the queries and therefore the outcome of the queries on every desoxyribonucleic acid sequence severally therefore doubtless he incorporates a higher leak than the shoppers. but to avoid this case, the shoppers will interact with Cloud1 to forestall Cloud2 from obtaining the outcomes on every sequence severally, as explained within the next paragraph.

#### A. Module Description:

In this project, we have these modules.

##### ➤ Conceptual level

At the conceptual level, we provide a deterministic scheme, with zero probability of a wrong answer (as opposed to a low probability).

This gives confidence to the users that they get exact results to all their queries, without impacting security

##### ➤ Operating point

The space-time tradeoff, by giving a scheme that is twice as fast as theirs but uses twice the storage space.

A variant of this scheme uses only 1.5 their storage space at the expense of additional latency

##### ➤ Data Encryption

This module is create key for provide security to the DNA data Based on well-known security building blocks like Paillier's encryption, public key encryption and symmetric key encryption.

##### ➤ Data Decryption

Queries that specify the number of occurrences of each kind of symbol in the specified sequence

A threshold query whose answer is 'yes' if the number of matches exceeds a query-specified threshold

For all query types we can hide the answers from the decrypting server, so that only the client learns the answer.

Counting the number of matches between the query symbols and a sequence Logical OR matches where a query symbol is allowed to match a subset of the alphabet thereby making it possible to Support for the extended alphabet of nucleotide base codes that encompasses ambiguities in DNA sequences

##### ➤ Set Query

This helps supporting ambiguity from the query side or biomedical researchers, important queries have often the form "How many records contain a diagnosis of Alzheimer disease and gene variant X?" Secure outsourcing of the database and allowing such type of queries without requiring the server to decrypt the data has been addressed.

The paper presents very practical results. For example, a count query over 40 records in a database of 5000 records takes 30 minutes. Our paper extends these results by proposing a variant storage and computation scheme

## V. CONCLUSION

In this paper, we've got revisited the challenge of sharing person-specific genomic sequences while not violating the privacy of their knowledge subjects so as to support large-scale medical specialty analysis comes. we've got used the framework planned by Kantarcioglu et al. [1] supported additive homomorphic secret writing, and 2 servers: one holding the keys and one storing the encrypted records. The planned technique offers 2 new operative points within the coordinate system exchange and handles new styles of queries that aren't supported in earlier work. moreover, the tactic provides support for extended alphabet of nucleotides that could be a sensible and important demand for medical specialty researchers.

Big knowledge analytics over genetic knowledge could be a sensible future work direction. There area unit fast recent advancements that address performance limitations of homomorphic secret writing techniques. we tend to hope that these advancements can cause additional sensible solutions within the future which will handle larger-scale

genetic science knowledge. It's price mentioning that our approach isn't restricted to a hard and fast homomorphic secret writing technique and thus, it'd be attainable to use and inherit the benefits of freshly developed ones. Even though cryptography may be a ground-breaking technology, de-identification of information is that the bottle neck for the protection of genomic sequences, thus we tend to providing security.

#### VI. REFERENCES

- [1]. M. Kantarcioglu, W. Jiang, Y. Liu, and B. Malin, "A cryptographic approach to securely share and query genomic sequences," *Inf. Technol. Biomed. IEEE Trans.*, vol. 12, no. 5, pp. 606–617, 2008.
- [2]. B. Malin and L. Sweeney, "How (not) to protect genomic data privacy in a distributed network: using trail re-identification to evaluate and design anonymity protection systems," *J. Biomed. Inform.*, vol. 37, no. 3, pp. 179–192, 2004.
- [3]. Z. Lin, A. B. Owen, and R. B. Altman, "Genomic research and human subject privacy," *Science (80-. )*, vol. 305, no. 5681, p. 183, 2004.
- [4]. A. E. Nergiz, C. Clifton, and Q. M. Malluhi, "Updating outsourced anatomized private databases," in *Proceedings of the 16th International Conference on Extending Database Technology*, 2013, pp. 179–190.
- [5]. L. Sweeney, A. Abu, and J. Winn, "Identifying Participants in the Personal Genome Project by Name," Available SSRN 2257732, 2013.
- [6]. E. Aguiar, Y. Zhang, and M. Blanton, "An Overview of Issues and Recent Developments in Cloud Computing and Storage Security," in *High Performance Cloud Auditing and Applications*, 2014, pp. 3–33.
- [7]. P. Bohannon, M. Jakobsson, and S. Srikwan, "Cryptographic Approaches to Privacy in Forensic DNA Databases," in *Public Key Cryptography*, vol. 1751, H. Imai and Y. Zheng, Eds. Springer Berlin Heidelberg, 2000, pp. 373–390.
- [8]. F. Esponda, E. S. Ackley, P. Helman, H. Jia, and S. Forrest, "Protecting data privacy through hard-to-reverse negative databases," *Int. J. Inf. Secur.*, vol. 6, no. 6, pp. 403–415, 2007.
- [9]. F. Bruekers, S. Katzenbeisser, K. Kursawe, and P. Tuyls, "Privacy-preserving matching of dna profiles," *IACR Cryptol. ePrint Arch.*, vol. 2008, p. 203, 2008.
- [10]. M. J. Atallah and J. Li, "Secure outsourcing of sequence comparisons," *Int. J. Inf. Secur.*, vol. 4, no. 4, pp. 277–287, Mar. 2005.
- [11]. M. Blanton, M. M. J. Atallah, K. B. K. Frikken, and Q. Malluhi, "Secure and Efficient Outsourcing of Sequence Comparisons," *Comput. Secur.* 2012, pp. 505–522, 2012.
- [12]. M. Franklin, M. Gondree, and P. Mohassel, "Communication-efficient private protocols for longest common subsequence," in *Topics in Cryptology--CT-RSA 2009*, Springer, 2009, pp. 265–278.