

Protecting Genomic Privacy in Medical Tests using Distributed Storage

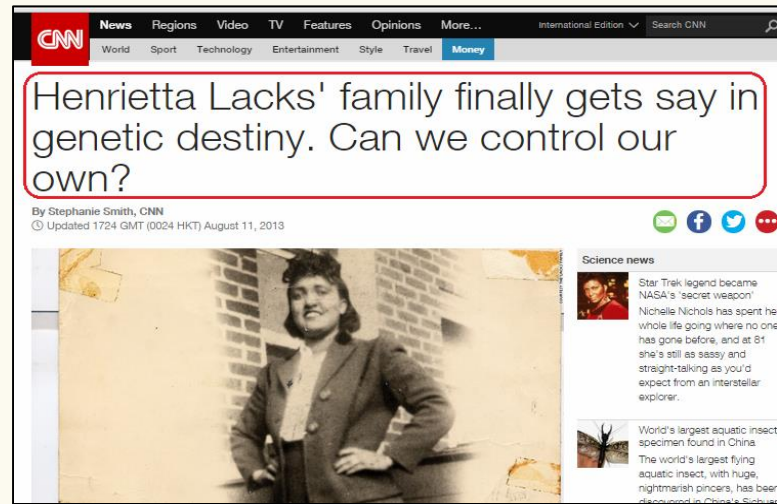
Sharmin Afrose (0905028), Maitraye Das (0905052)

Introduction

Individual's chances of diseases are largely associated with personal genetic variations. Hence, genomic data is significantly used in disease susceptibility tests and personalized medicine.

Privacy Threats

- Reveals trait, ancestry, vulnerability of diseases etc.
- Exposes relatives' genome.
- Initiates genomic discrimination in insurance, employment etc.



Thesis Goal

Privacy-preserved and precise computation of multiple disease risks using genomic and clinical data.

Novelty:

We offer substantial improvement over cryptography-based methods^{1,2,3} regarding queries for multiple diseases related to both the alleles of the same SNP.

Genomic Background

- Four nucleotides : A, C, G, T.
- SNP:** Difference of a single nucleotide between
 - Members of same species
 - Paired chromosomes of an individual.
- Each SNP carries two alleles; one from each parent.
- Both alleles can contain risks of different diseases.

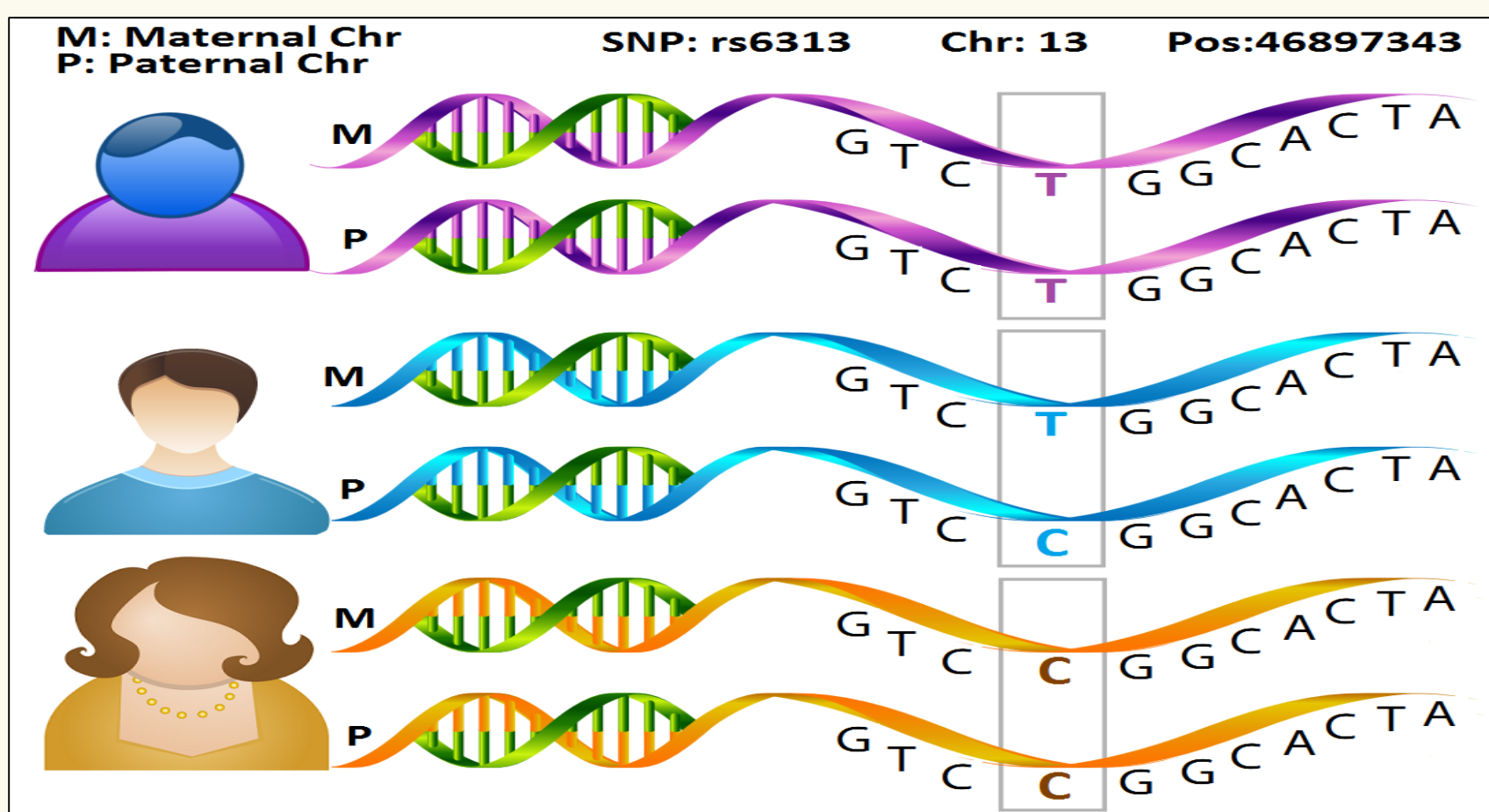


Fig 1: DNA fragments showing SNP rs6313

System Architecture

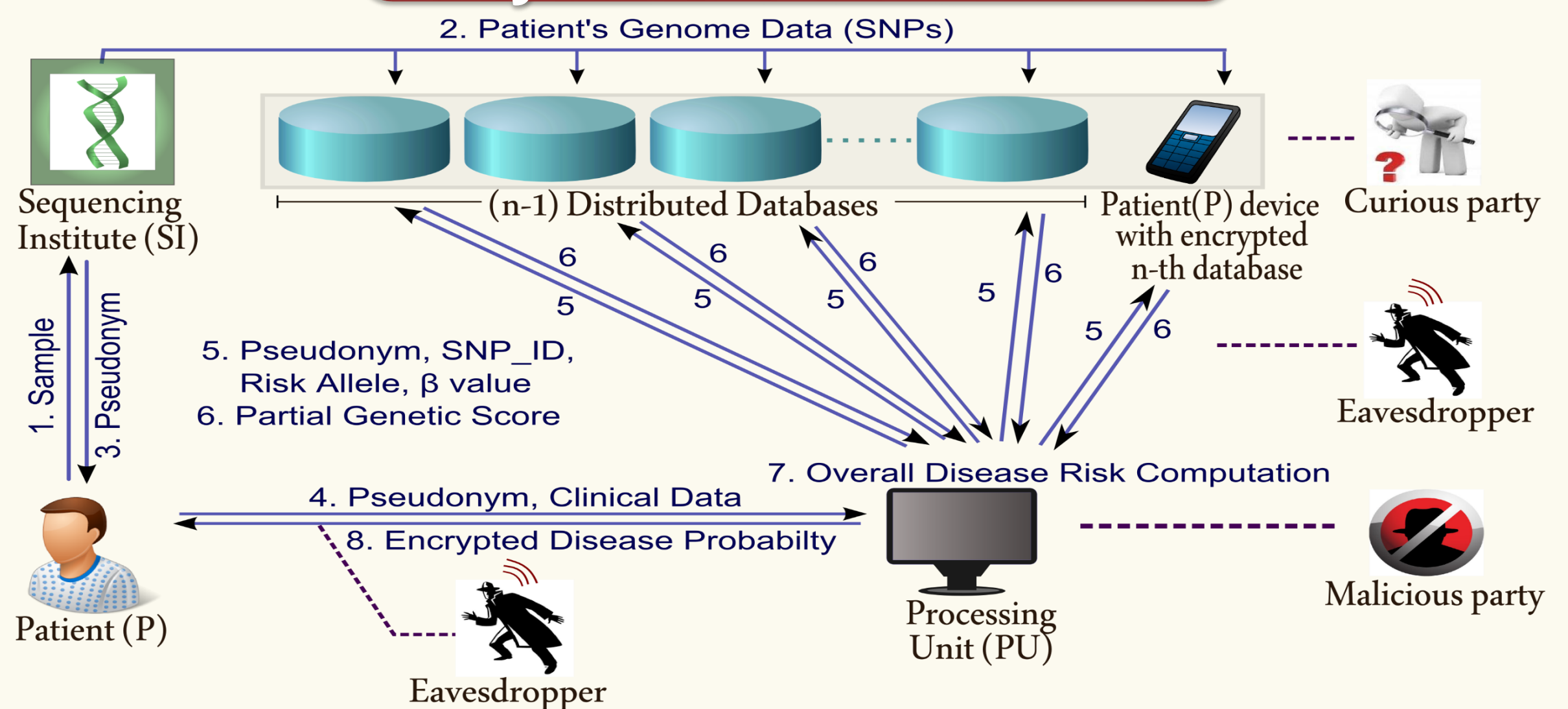


Fig 2: System architecture and threat model for disease risk computation

Our Approach

Genomic data is distributed in DBs such that only aggregated data reveals true contents. Dummy diseases are added in query using ℓ -diversity.

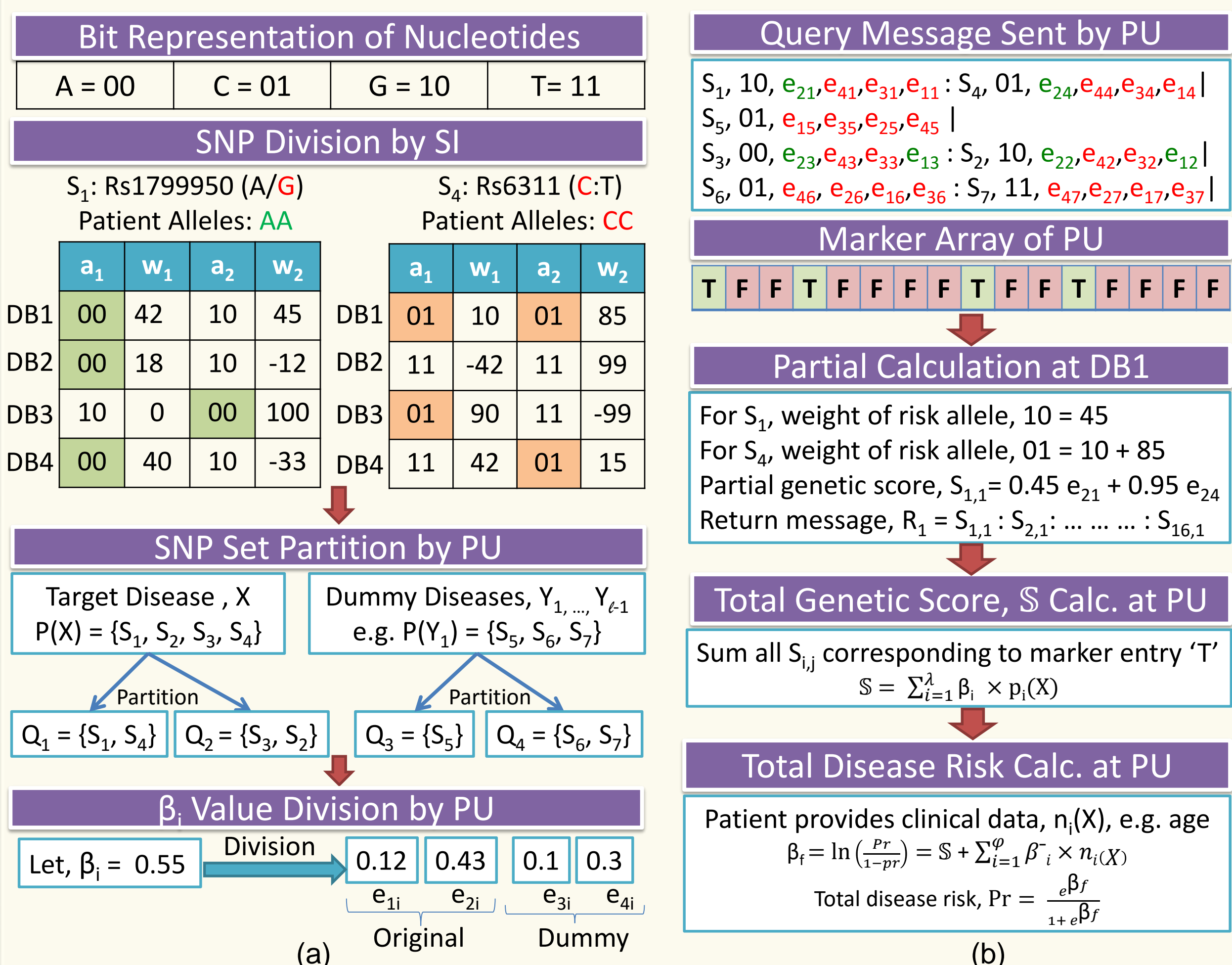


Fig 3: Flowchart of (a) query message generation (b) total disease risk computation

Implementation and Evaluation

Privacy Analysis

- SNP retrieval attack by
 - Semi-honest DDB
 - Semi-honest PU
 - Dishonest-but-covert PU
- Test inference attack by semi-honest DDB
- β value inference attack by semi-honest DDB

Dataset Size

- Over 0.3 million.

Assumption

- DDBs maintain protocols.

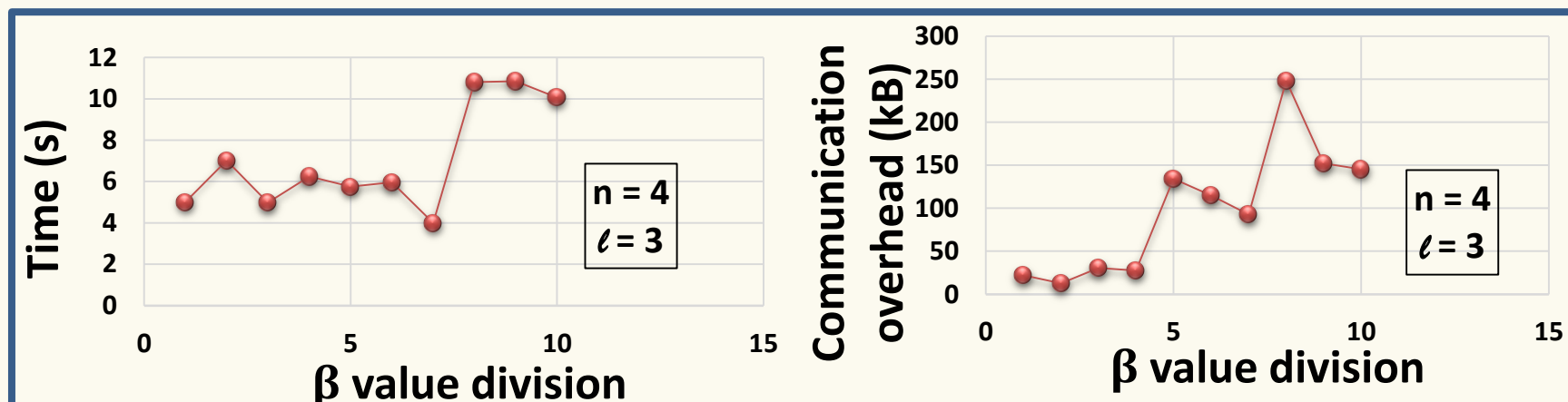


Fig 4: Increasing privacy level (β value division) has no particular effect on time complexity and communication overhead.

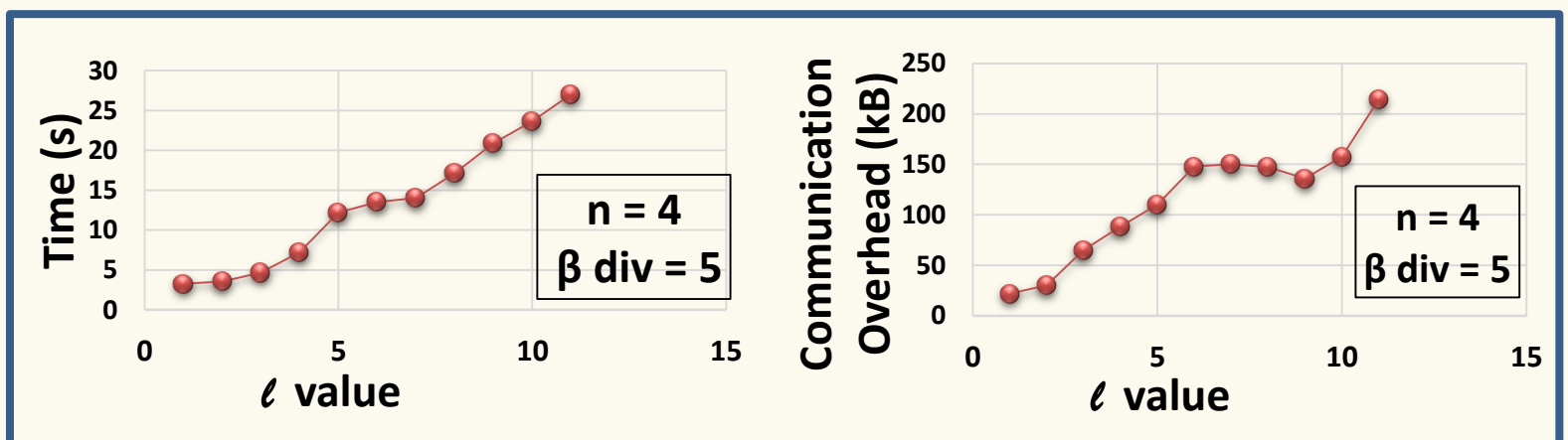


Fig 5: Increasing privacy level (ℓ value) has almost linear effect on time complexity and communication overhead.

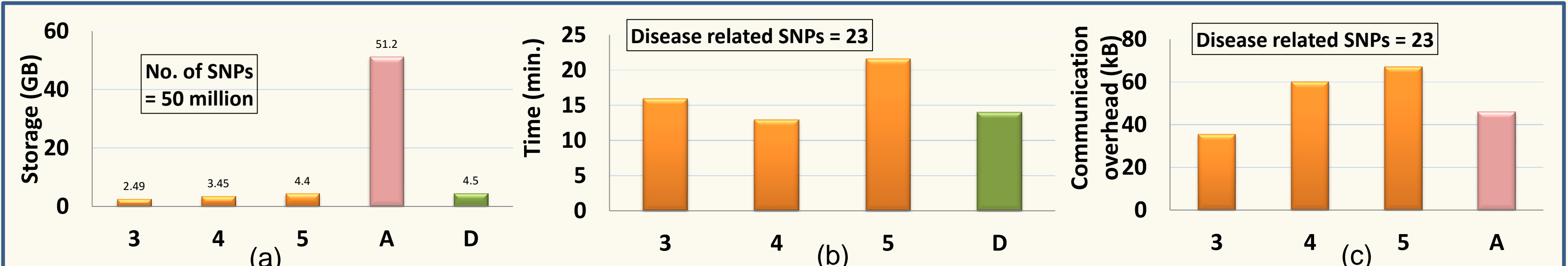


Fig 6: Performance of our approach compared to existing methods – A¹ and D² varying privacy level (no. of DDBs, n) for $\ell = 3$ and β division = 5 with respect to (a) storage (b) time complexity (c) communication overhead

Conclusion

Our proposed system preserves genomic privacy in medical tests using distributed storage and precisely computes risks for multiple diseases in all possible scenarios. Estimation of performance compared to existing methods over real dataset shows its practicality in real-life implementation.

References

- E. Ayday, J. L. Raisaro, P. J. McLaren, J. P. Hubaux, and J. Rougemont, "Protecting and evaluating genomic privacy in medical tests and personalized medicine", in WPES, 2013.
- G. Danezis and E. D. Cristofaro, "Fast and private genomic testing for disease susceptibility", in WPES, 2014.
- L. Barman, M.-T. Elgraini, J. L. Raisaro, J.-P. Hubaux, and E. Ayday, "Privacy threats and practical solutions for genetic risk tests", in GenoPri, 2015.
- <http://edition.cnn.com/2013/08/07/health/henrietta-lacks-genetic-destiny/>