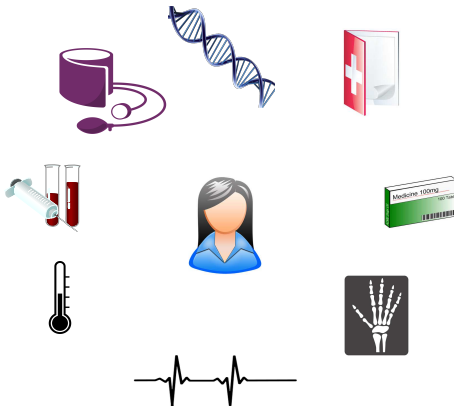# Security and Privacy of Genomic Data

**Sara Foresti**
Dipartimento di Informatica
Università degli Studi di Milano
sara.foresti@unimi.it

Challenges in Data-Driven Genomic Computing
Como, Villa del Grumello – March 8, 2019

# Genomic data (1)

- Large collections generated thanks to reduction of sequencing costs

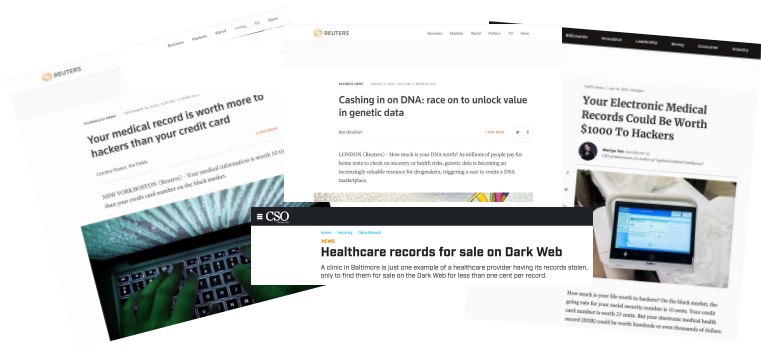- Highly related with personal and medical data

# Genomic data (2)

- Advantages for research
  - + data analysis for extracting valuable information
  - + sharing for collaborative computation

# Disclosure risks

- Considerable economic value

- Highly sensitive content

# Disclosure risks

- Considerable economic value

- Highly sensitive content

- High risk of exposure is case of attacks

# Security: a complex problem



Publication

Outsourcing

Sharing

Regulations

# Data protection – Publication

- Minimize release/exposure
  - correlation among different data sources
  - indirect exposure of sensitive information
  - de-identification $\neq$ anonymization
  - privacy vs utility

# Data protection – Outsourcing

- Encryption protects data confidentiality but
  - limits functionality
  - indirect exposure

# Data protection – Sharing

- Scientific research demands for data sharing
  - $+$ combine data collections owned by different subjects
  - $+$ enables collaboration
  - $-$ requires controlled data release

# Characterization of
# Data Protection Challenges

# Scientific and technical challenges



SECURITY PROPERTIES

SLA compliance

Integrity

Confidentiality

Data archival

Data retrieval / extraction

Data updates

ACCESS REQUIREMENTS

1 user - 1 provider

N users - * providers

* users - N providers

ARCHITECTURES

# Security properties



**Confidentiality**
- data externally stored
- users identities
- actions that users perform on the data



**Integrity**
- data externally stored
- computation and query results



**SLA compliance**
- assurance and certification

# Access requirements

**Data archival**
- upload/download
- protection of data in storage

**Data retrieval/extraction**
- support for fine-grained data retrieval and queries
- protection of computations and query results

**Data update**
- support for access retrieval and enforcement of updates
- protection of the actions and of their effects on the data

# Architectures



**1 user - 1 provider**
- protection of data at rest
- fine-grained retrieval
- query privacy/integrity

**n users - * providers**
- authorizations and access control
- multiple writers

***** users - n providers**
- controlled data sharing and computation

# Data Sharing:
# Issues and Directions

# Today...

- Two extreme solutions

  - Share everything

    - + enables collaboration

    - – requires full trust

  - Share nothing

    - + guarantees privacy

    - – slows scientific research

# ...Tomorrow

- Selective sharing based on: receiving subject, data sensitivity, context, purpose, ...

- Requires to study solutions enabling to:

    ○ identify sensitive data

    ○ express access restrictions through a simple while flexible language

    ○ protect data (e.g., encryption, aggregation, obfuscation)

# GMQL for enabling sharing

- GMQL is an expressive and flexible query language for genomic data

- Could be extended to:

  - associate protection requirements with the data

  - specify access restrictions

  - enable the enforcement of protection techniques

# Conclusions

- Data collection and analysis are vital for scientific research

- Solutions that guarantee data protection are enabling for:

  - data publication

  - outsourcing of data storage and/or computation

  - data sharing and collaborative computations and

  - …