





Privacy in Epigenetics: Temporal Linkability of MicroRNA Expression Profiles

Mathias Humbert

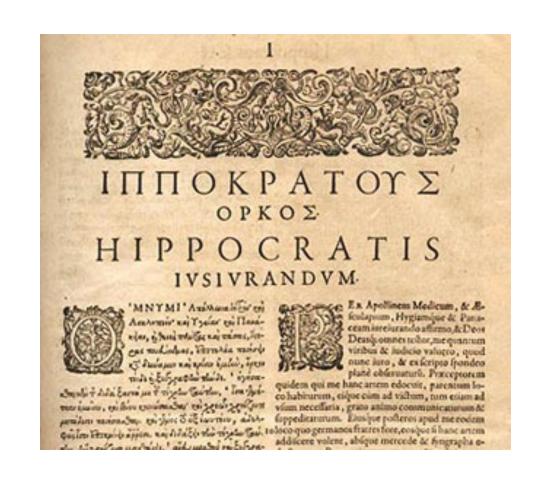
Joint work with Michael Backes, Pascal Berrang, Anne Hecksteden, Andreas Keller and Tim Meyer

Summer Research Institute 2016 EPFL, Switzerland

Archaeology of Privacy

- Very first adversary: physician
 - Only credible person entering your home/intimacy
 - => Only possible channel of information leakage
 - Health information already considered very sensitive
- First privacy-preserving mechanism: Hippocratic oath (5th century B.C.)

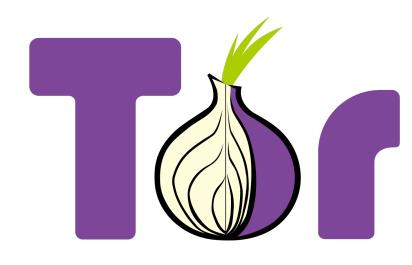


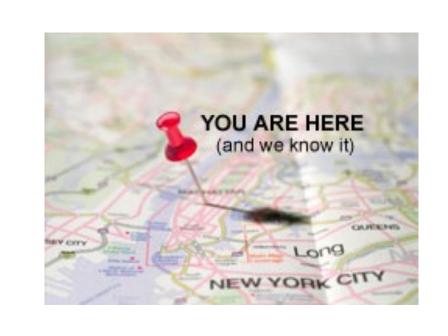


Modern Privacy

- Internet communications
- Web browsing/fingerprinting
- Location privacy







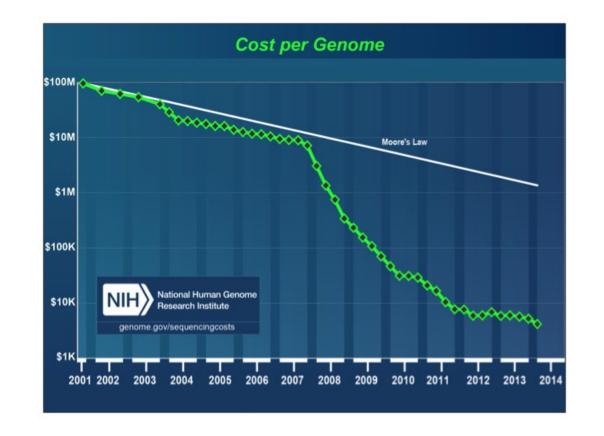




Deluge of Biomedical Data

- Decreasing cost of molecular profiling tests
- Fueling the precision medicine revolution
- Increasing amount of biomedical data available
 - Not only on "trusted" DB such as hospital servers
 - Available on online public databases too









Hippocratic oath is not sufficient anymore

Health DB Breaches

- Attacks against healthcare companies
 - E.g., health insurer Anthem: **78 million** records put at risk
- 91% of healthcare companies experiencing a violation of their DB over the last two years
 - Only 32% feeling they have adequate resources to defeat these incidents
- Sensitive health data of thousands of patients ending up online due to a human mistake

Bilans de santé en balade sur le net

GAFFE — Des données médicales ultraconfidentielles de patients romands ont été librement accessibles durant des jours sur Internet. Le groupe Synlab déplore une erreur humaine.

Par Raphaël Pomey . Mis à jour le 08.04.2015 5 Commentaires



Genomic Privacy

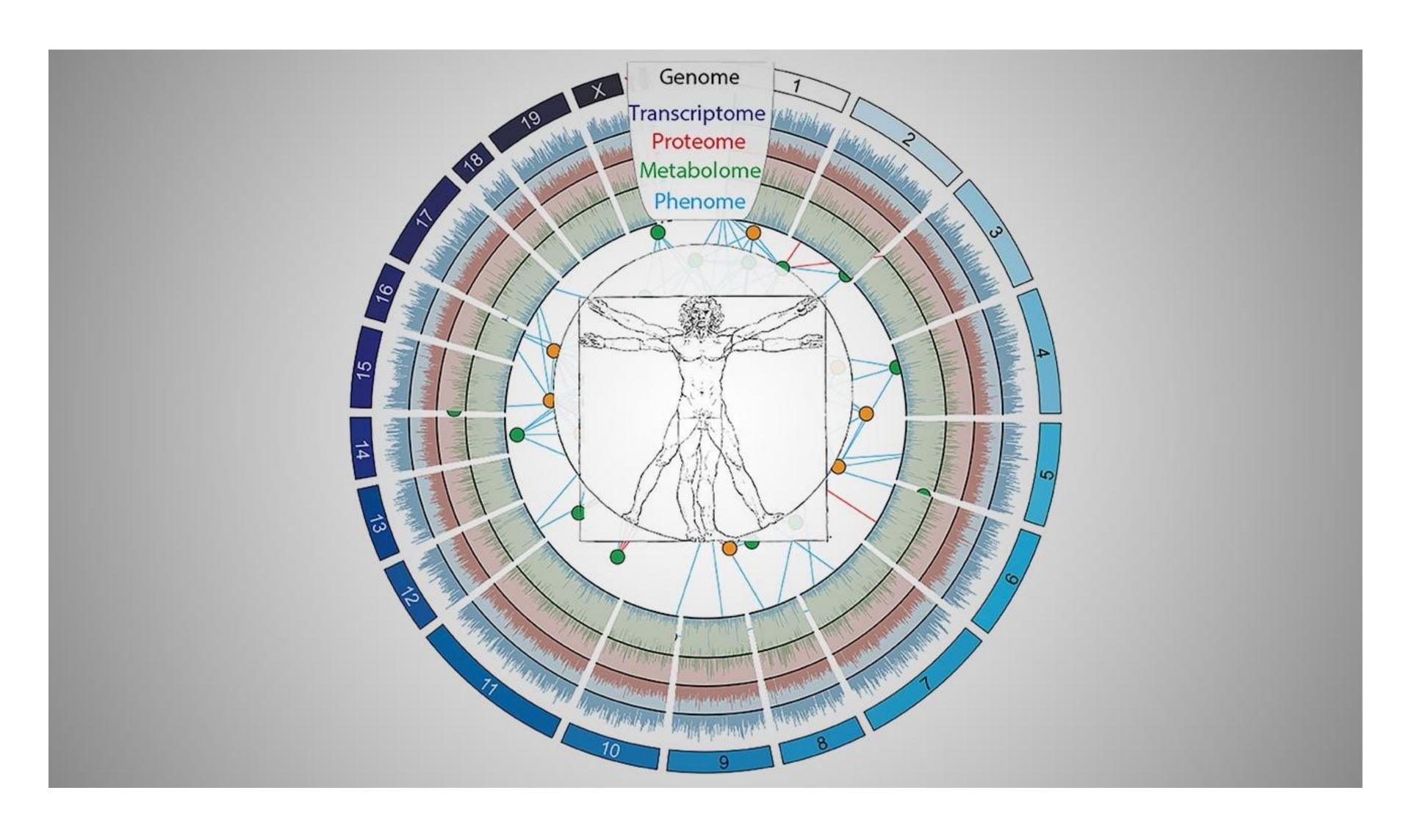
Already studied quite extensively by the security/privacy community

Categorization of techniques for breaching genomic privacy [1]

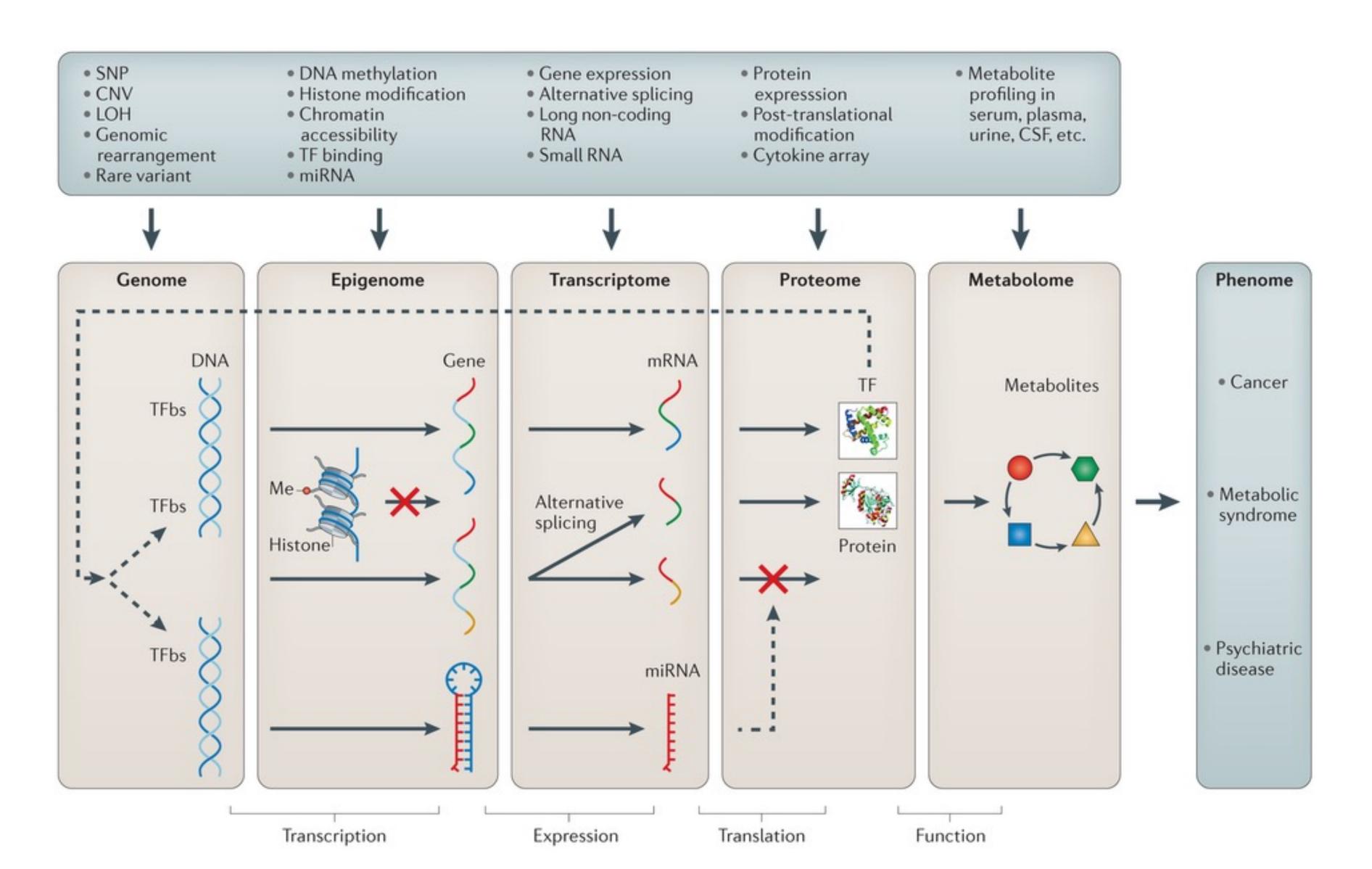
| Technique | Maturation Level | Technical complexity | Auxiliary information |
|--------------------------------------|---------------------|----------------------|-----------------------|
| Identity Tracing | | | |
| Surname Inference | **** | ••• | Intermediate- |
| | | | Good |
| DNA Phenotyping | ** | •• | Poor |
| Demographic identifiers | **** | • | Good |
| Pedigree structure | *** | •• | Poor |
| Side channel leakage | *** | ••• | Varies |
| Attribute Disclosure Attacks via DNA | | | |
| N=1 | *** | ••• | Good |
| Genotype frequencies | *** | ••• | Good |
| Linkage disequilibrium | ** | •••• | Intermediate |
| Effect sizes | ** | ••• | Good |
| Trait informac | | | Good |

Privacy of other types of health-related data?

The Human OSI Stack



The Human OSI Stack



Epigenetics and MicroRNA

Epigenetics

"epi": above, over (greek)
"genetics": origin (greek)

Definition: study of cellular and phenotypic trait variations stemming from other causes than changes in the genotype

MicroRNA (miRNA)

discovered in the early 1990s

Definition: small non-coding RNA molecules that regulate gene expression in plants/animals 60% of genes coding human proteins are regulated by miRNAs

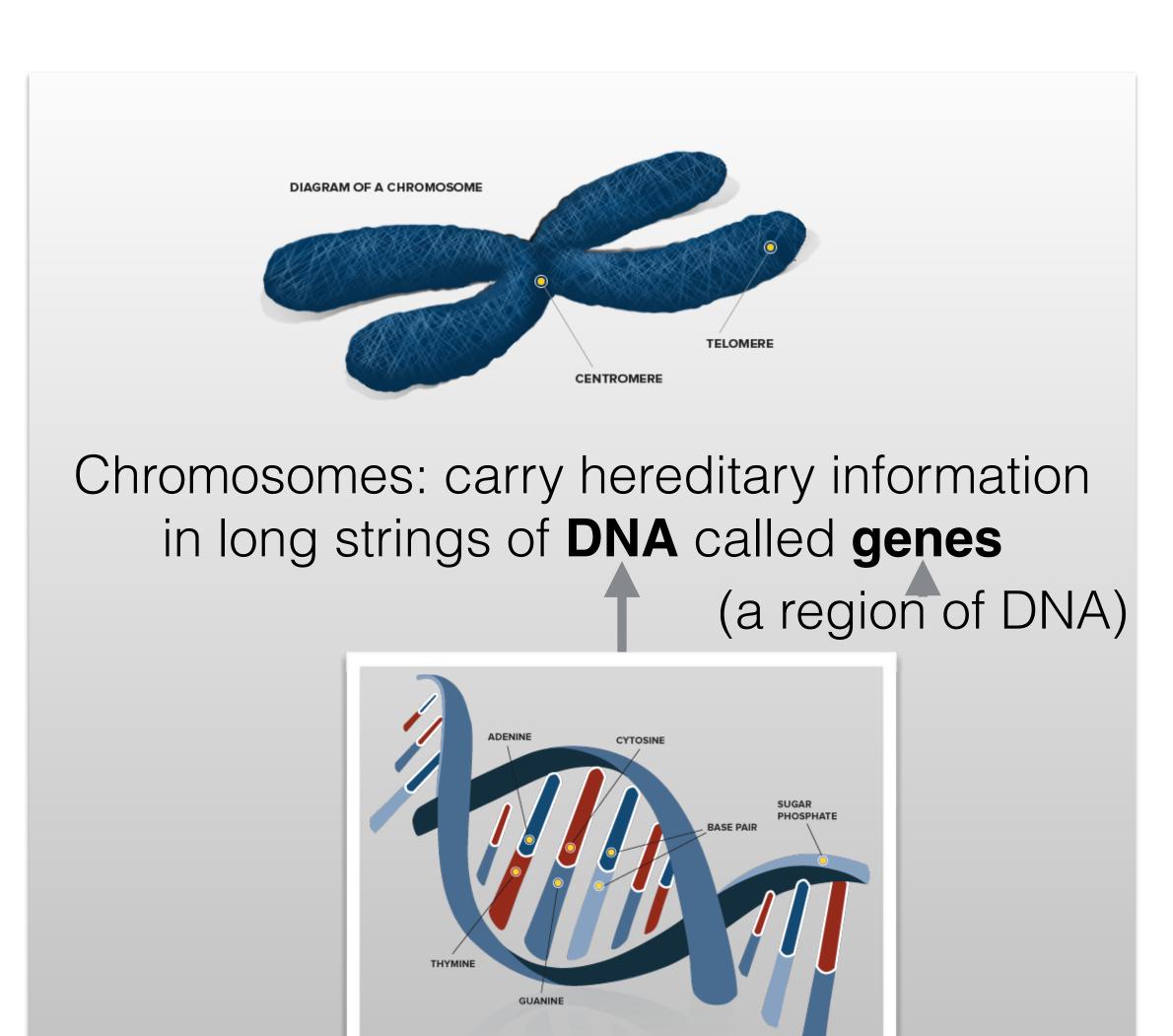
MicroRNA Expressions

Real-valued numbers quantifying whether and how much miRNAs are active in a given set of cells/tissue.

External factors such as:

in-utero and childhood development, environmental chemicals, aging, diet.

What is the Role of MicroRNAs?



But all cells have the same genes!



What makes the cells different:

gene expression

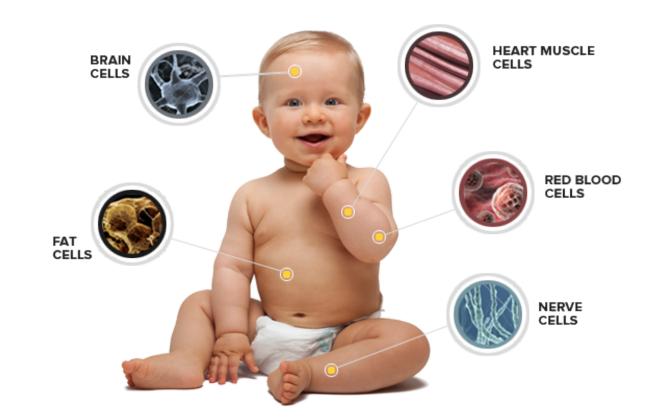
(which genes are active in a cell)

What is the Role of MicroRNAs?

What makes the cells different:

gene expression

(which genes are active in a cell)



miRNAs regulate most of human genes!

important for normal and disease cells

neurodegenerative diseases (e.g., Alzheimer's) heart diseases, diabetes, majority of cancers

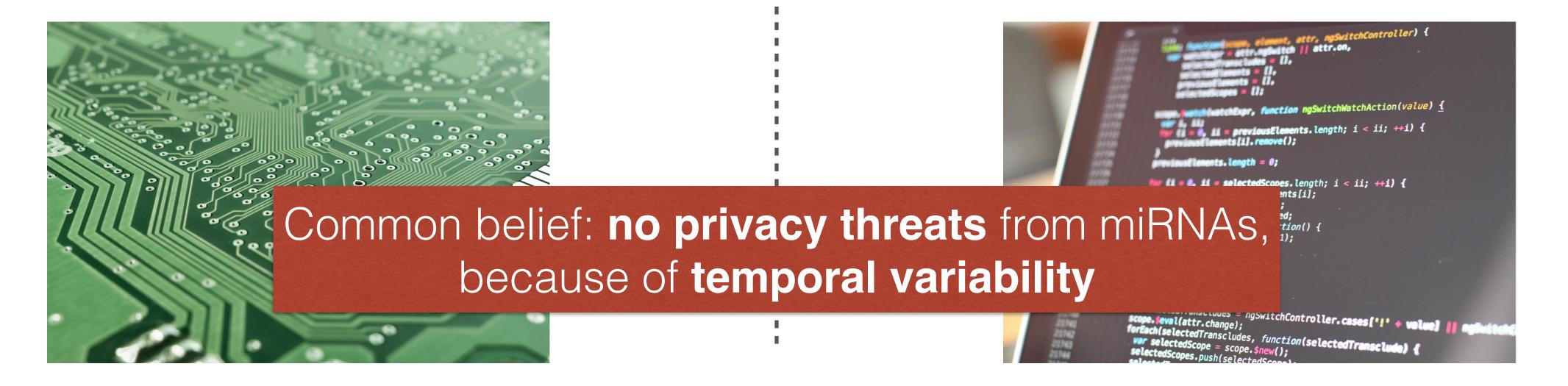
More on DNA and MicroRNAs

DNA

- contains receipts what a cell potentially can do,
- is (mostly) fixed over time,
- can hint on risks of getting a disease,
- privacy of the genome has been researched a lot.

miRNA

- expression regulates what a cell really does,
- expression changes over time,
- can tell whether you carry a disease,
- so far, privacy of miRNA has been largely overlooked.



Linkability Attacks

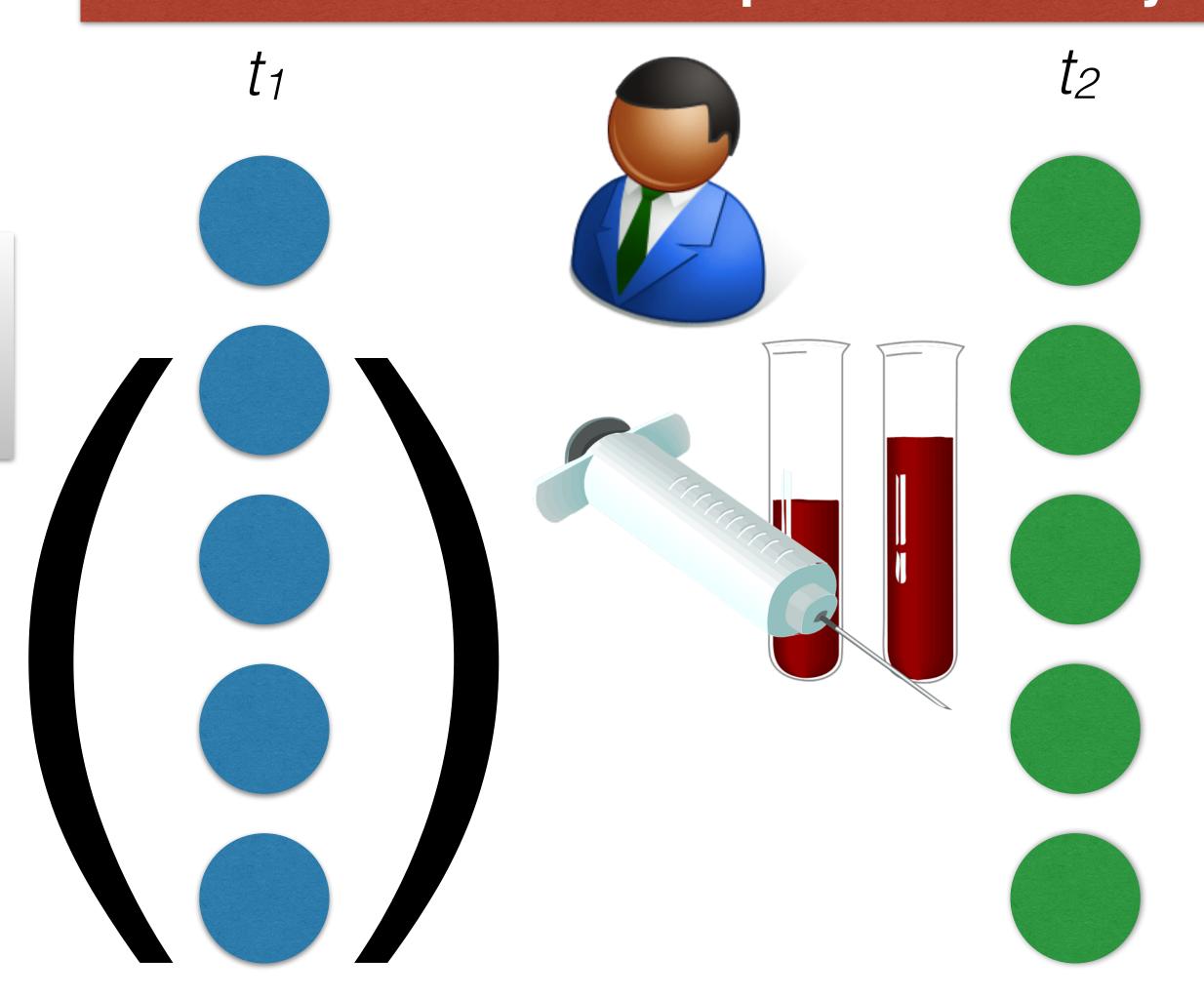
identification

Common belief: **no privacy threats** from miRNAs, because of **temporal variability**

matching

public DB (such as the Gene Expression Omnibus)





black market

cyber **attacks** against healthcare companies have **increased** by 72% within one year

Athletes' Dataset



Participants: 29

Points in time: 2 (before and after exercising)

Time shift: 1 week

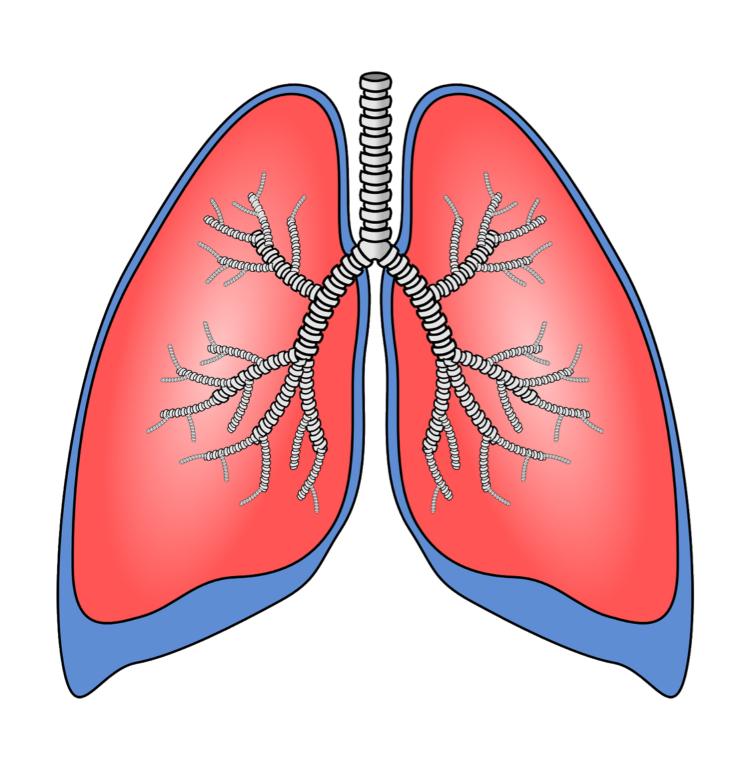
Disease: none

blood-based

plasma-based

1,189 miRNAs per sample

Lung Cancer Dataset



Participants: 26

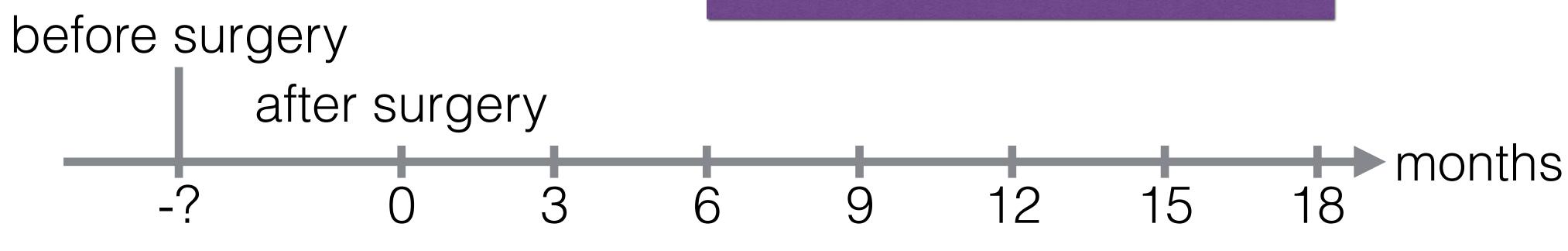
Points in time: 8

Time shift: mostly 3 months

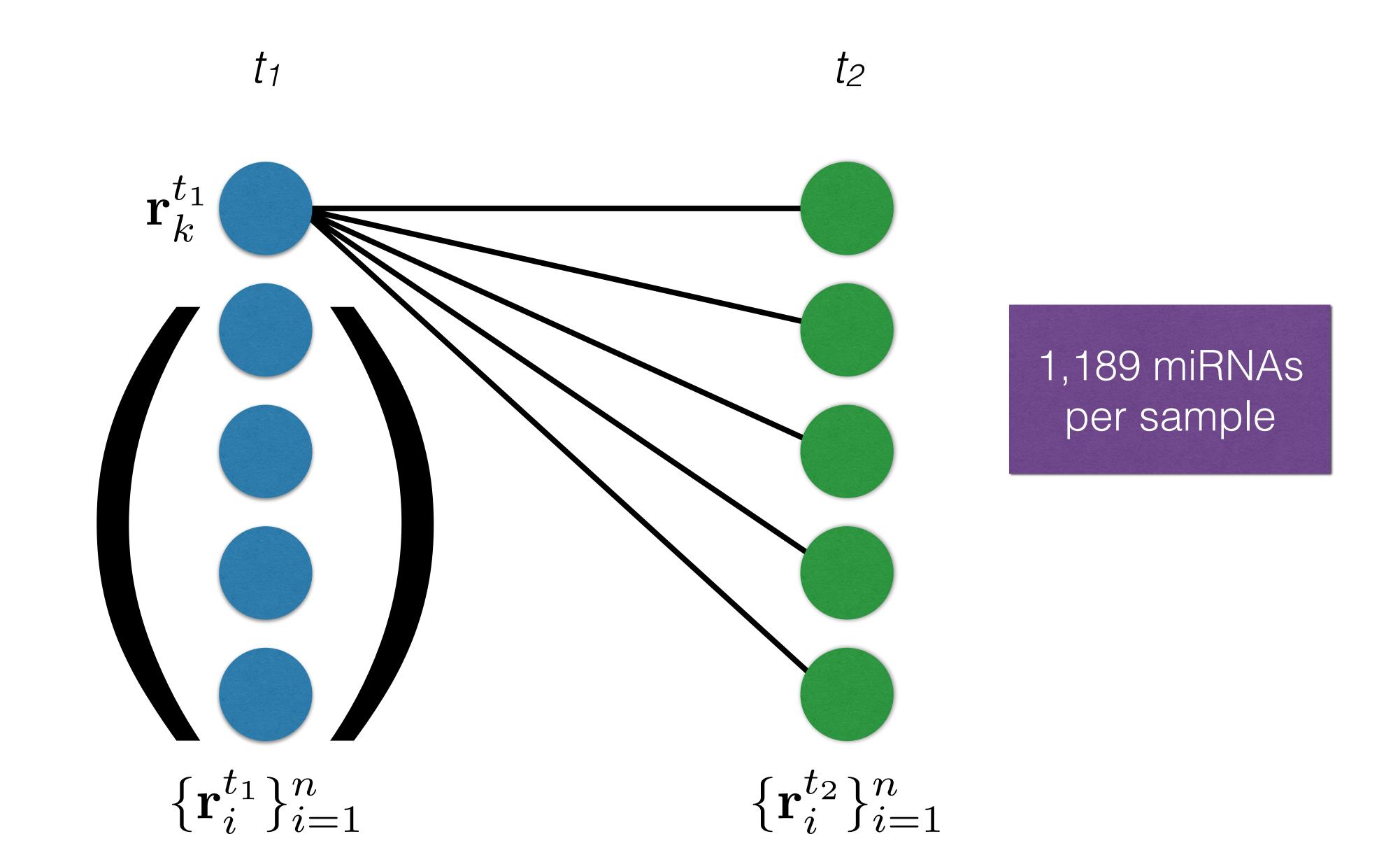
Disease: lung cancer

plasma-based

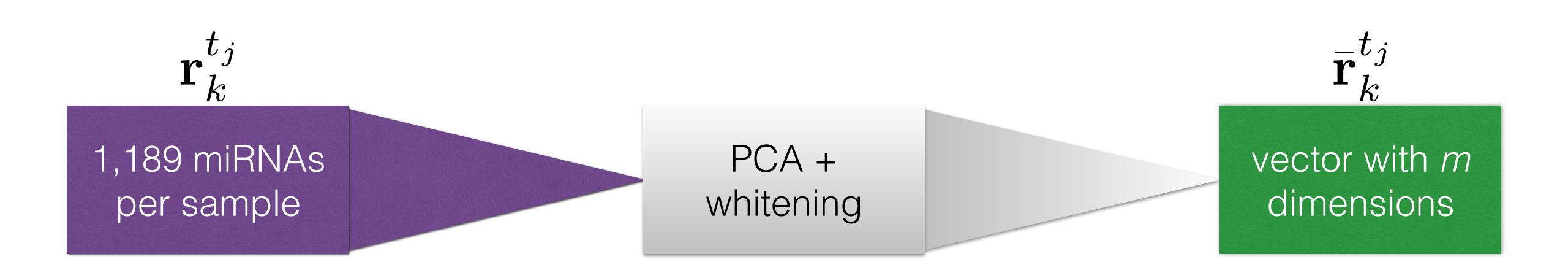
1,189 miRNAs per sample



Attack Formalization

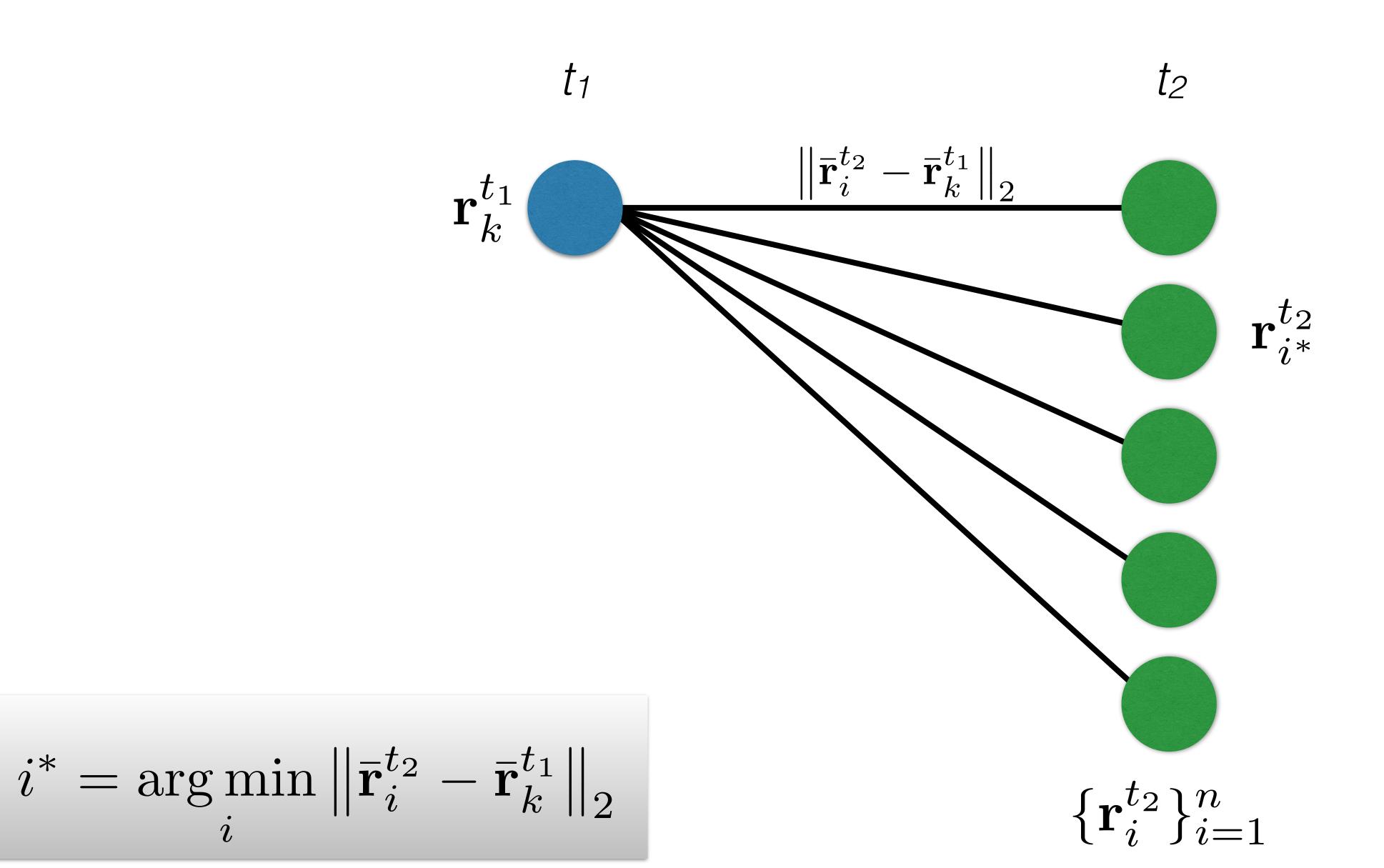


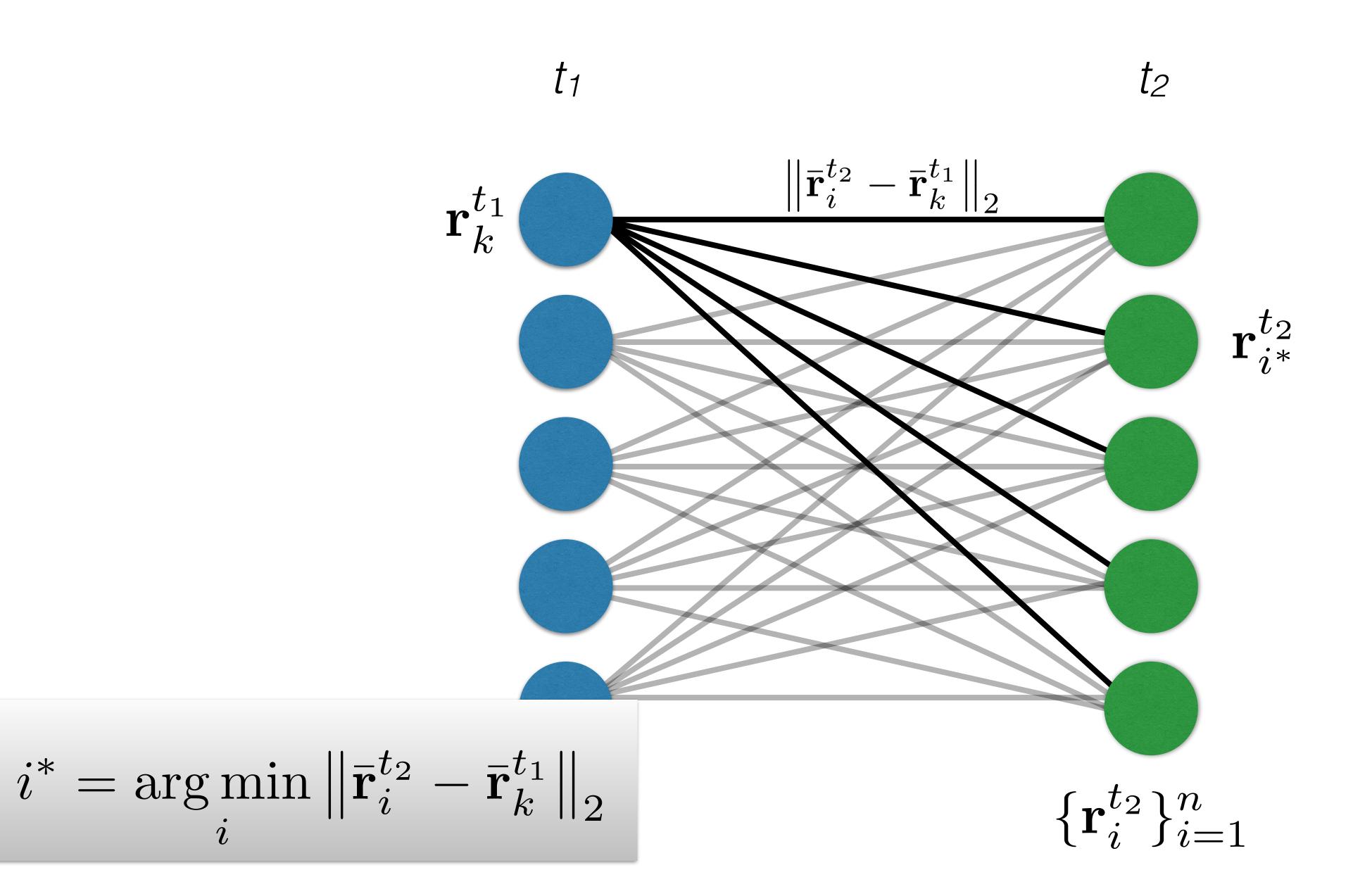
Pre-processing Step

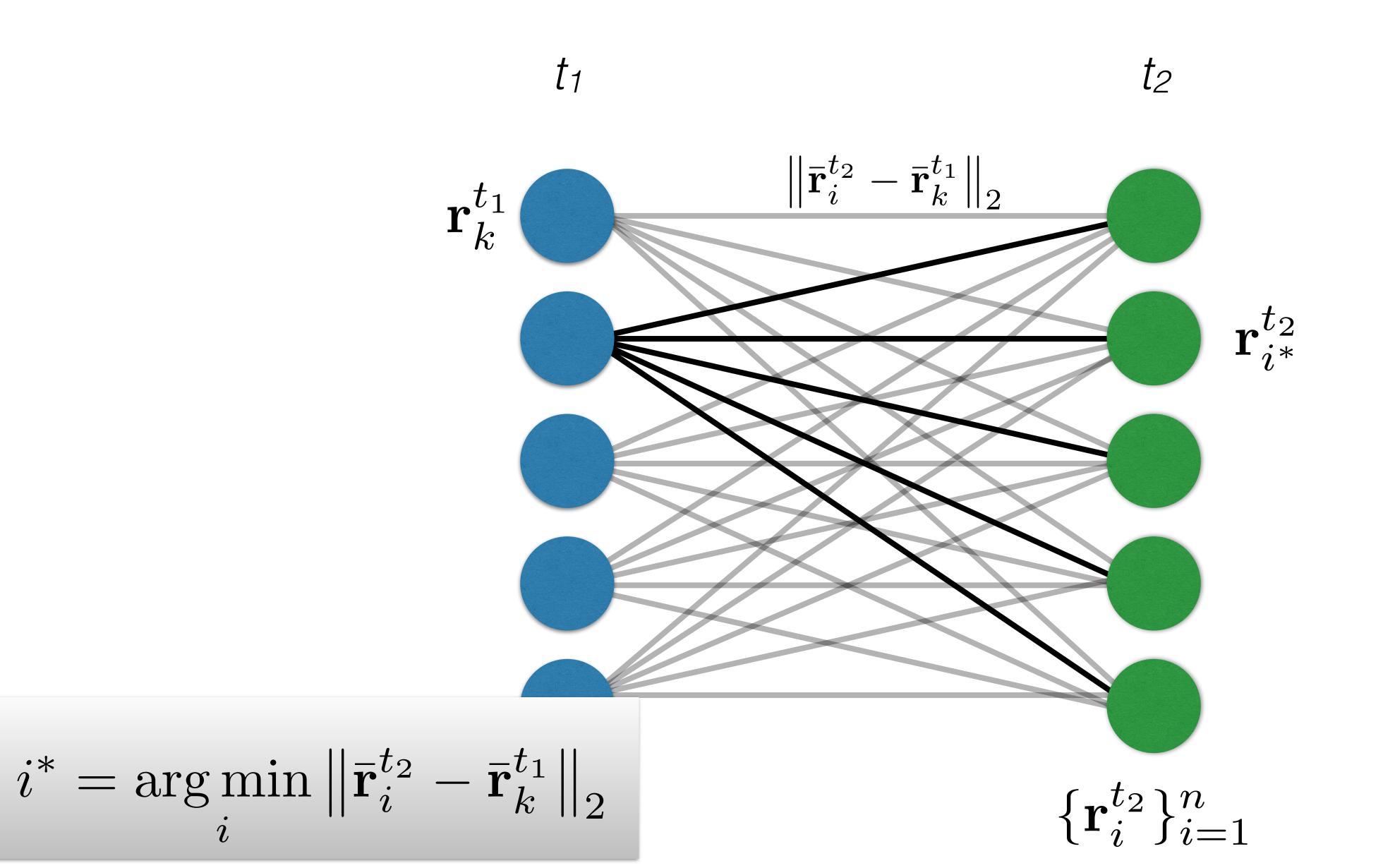


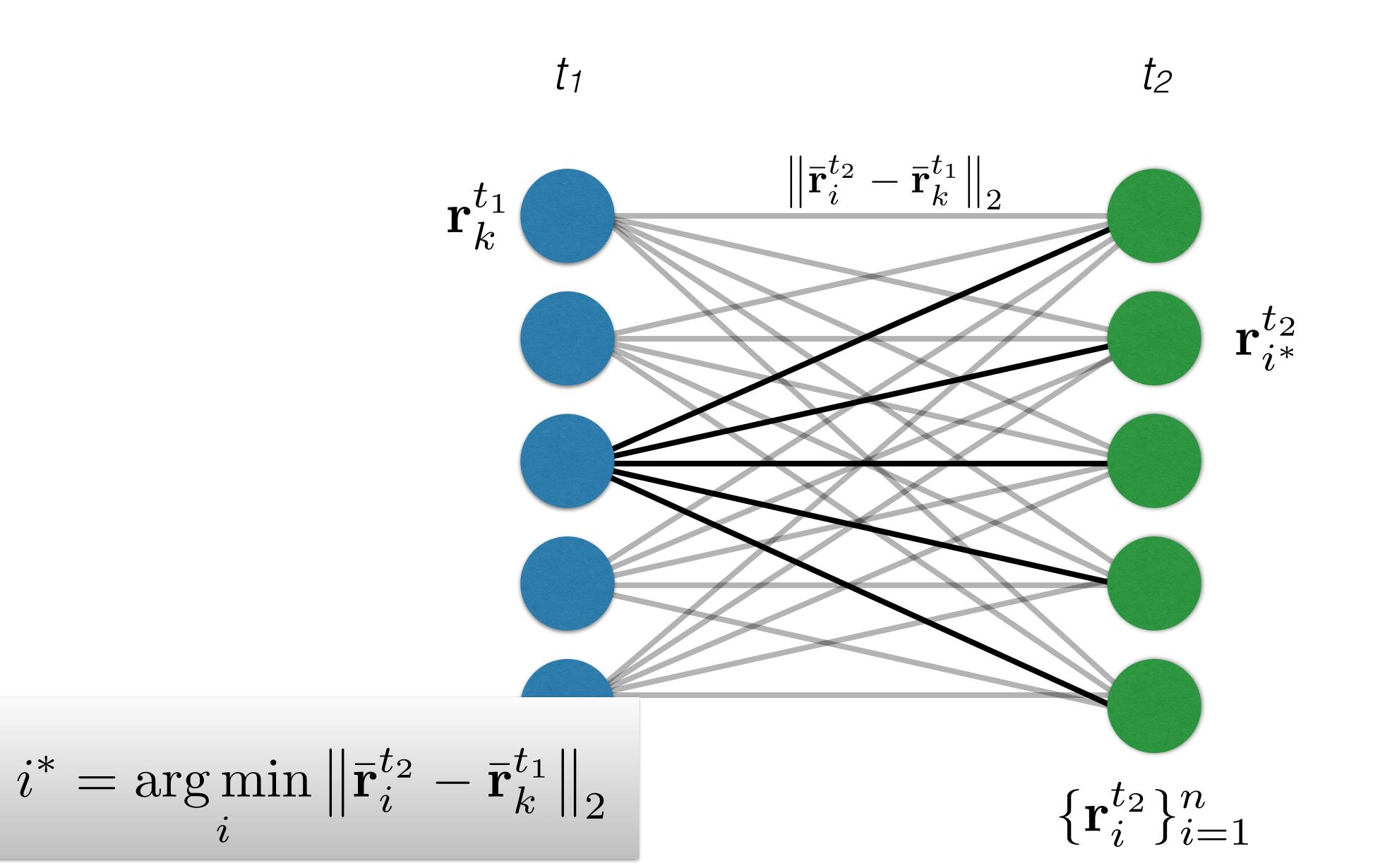
whitening: unit variance *PCA*: smaller dimensionality *m*

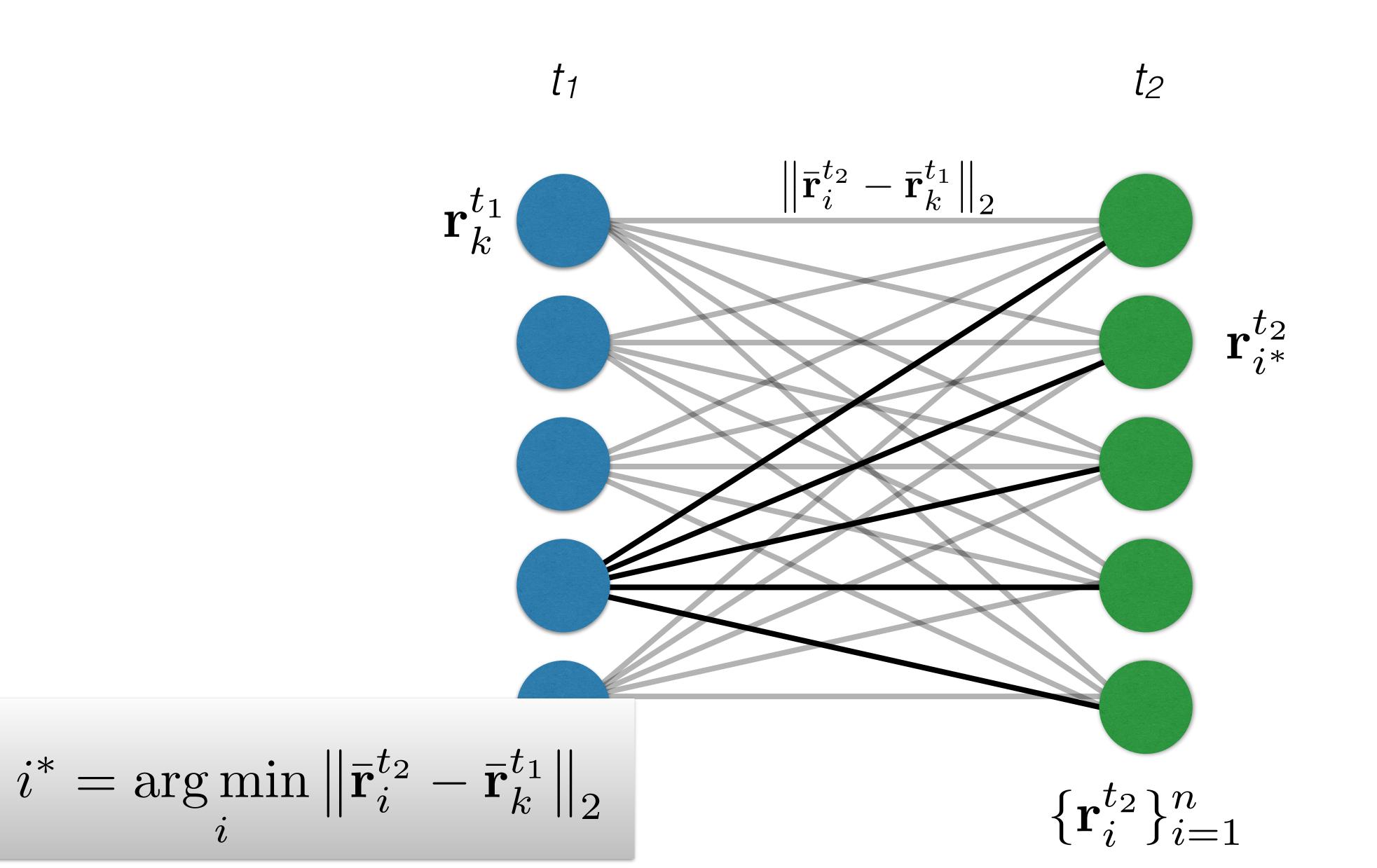
+ uncorrelated components

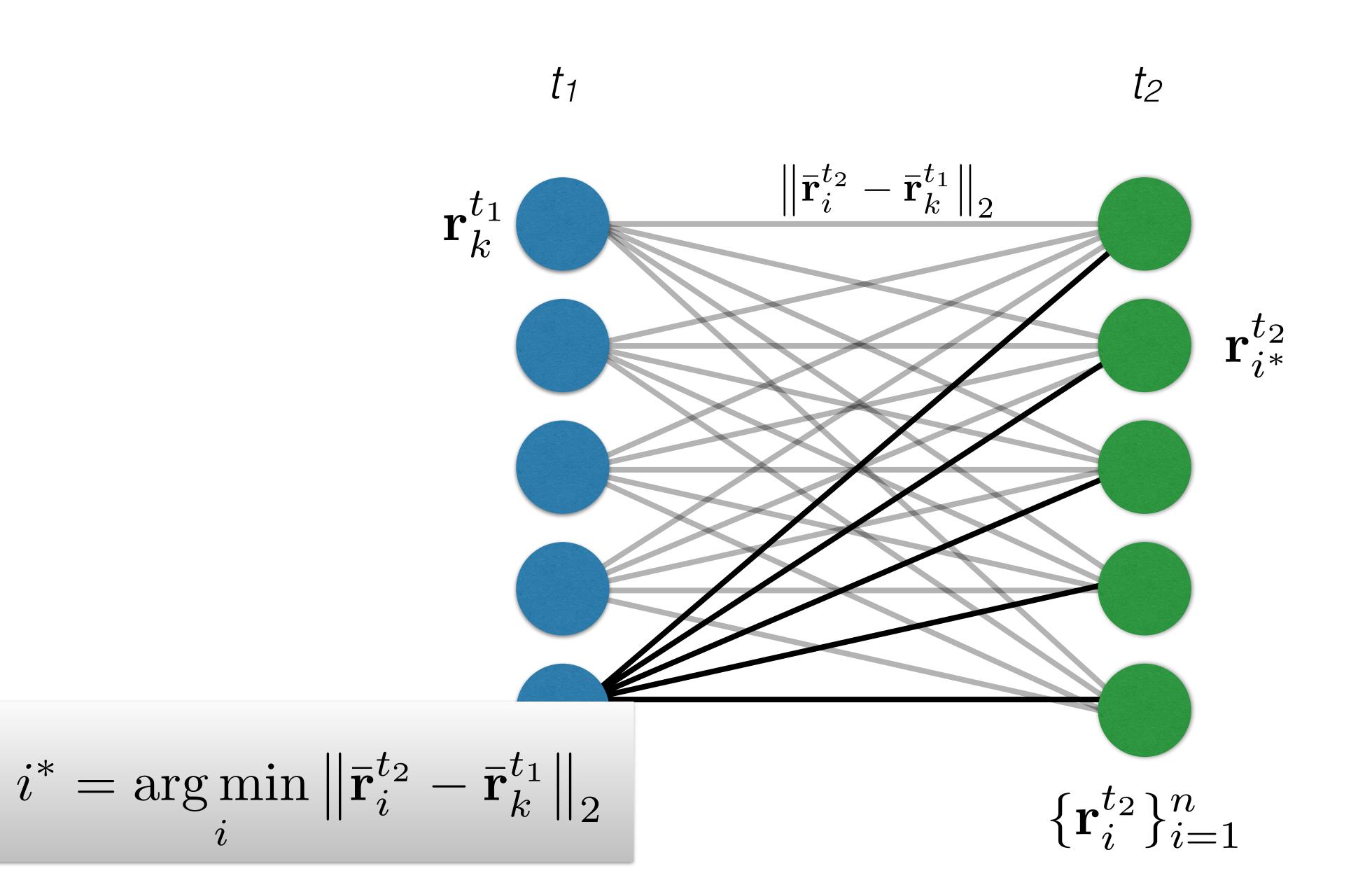




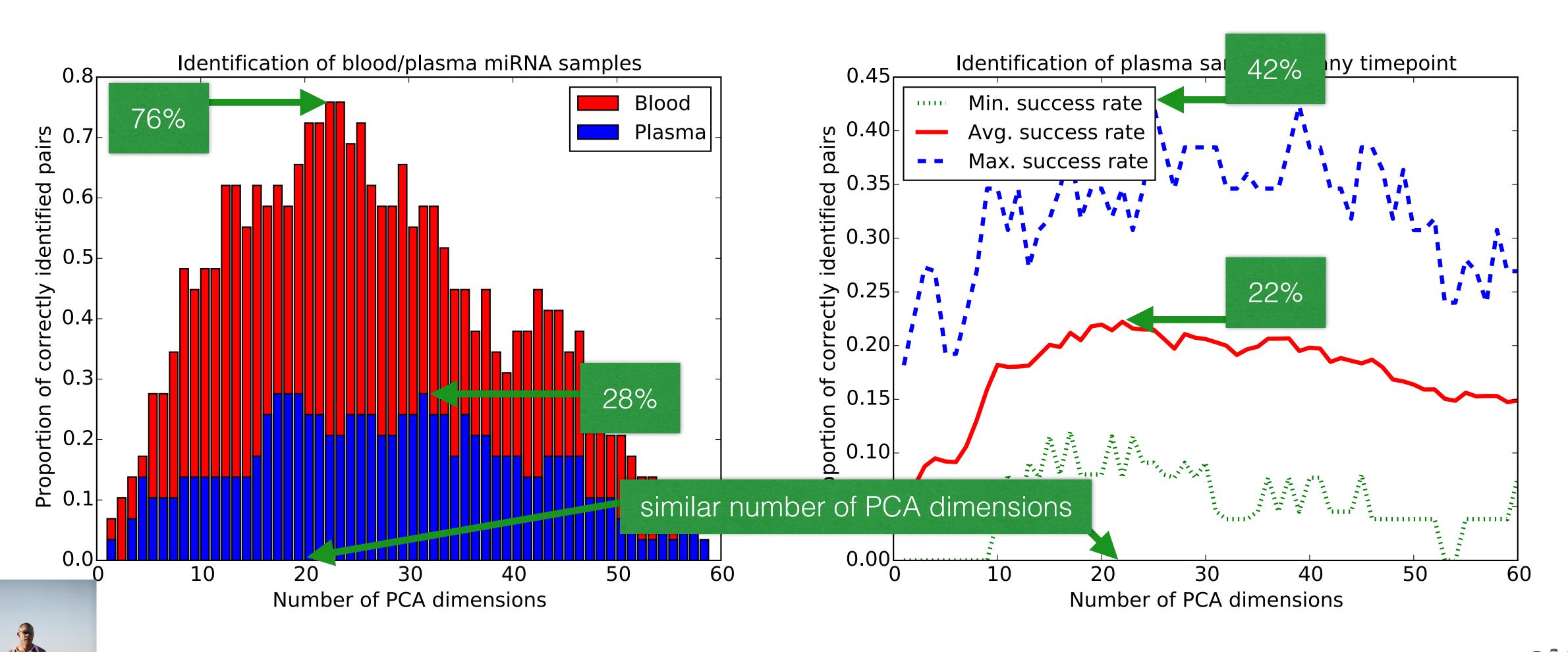




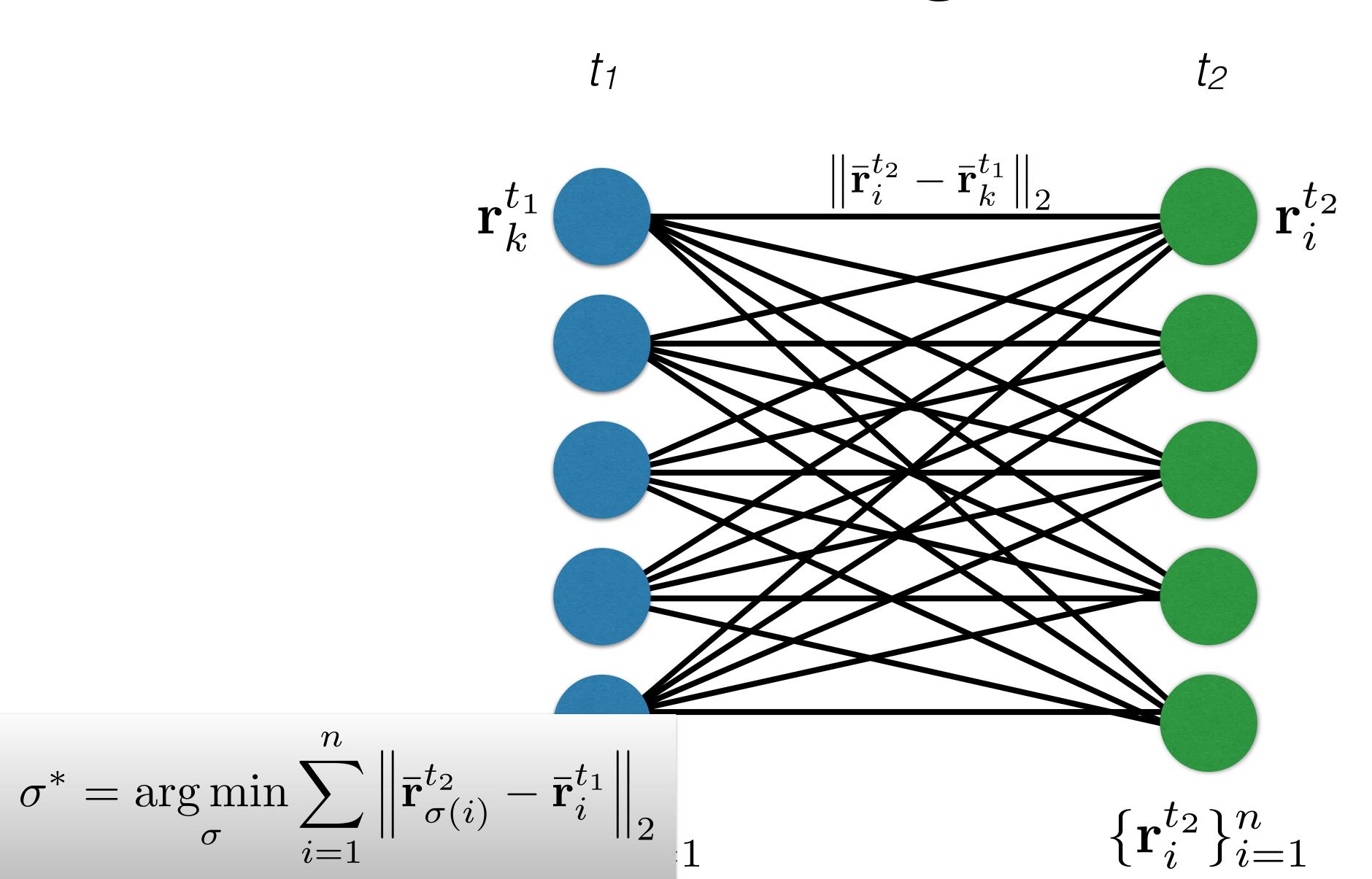




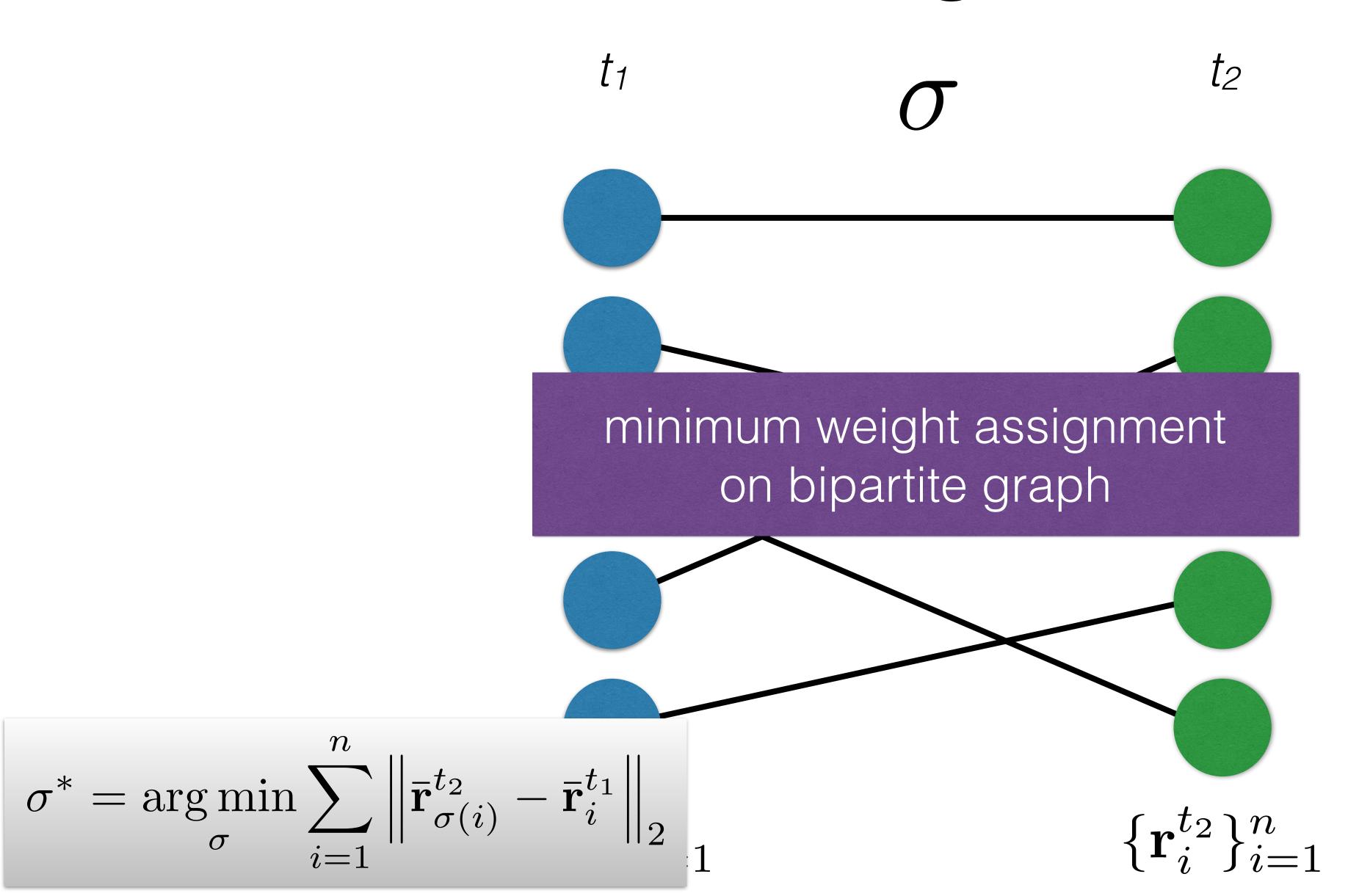
Identification Attack - Results



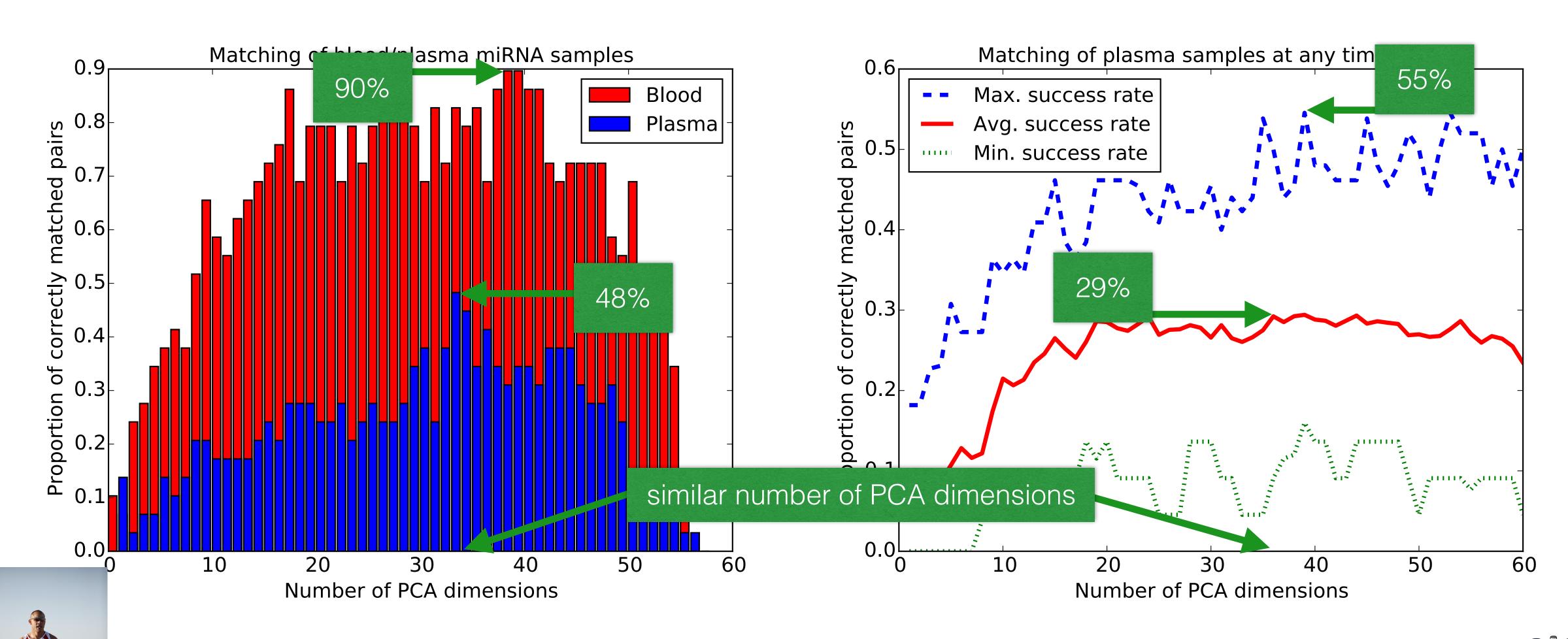
Matching Attack



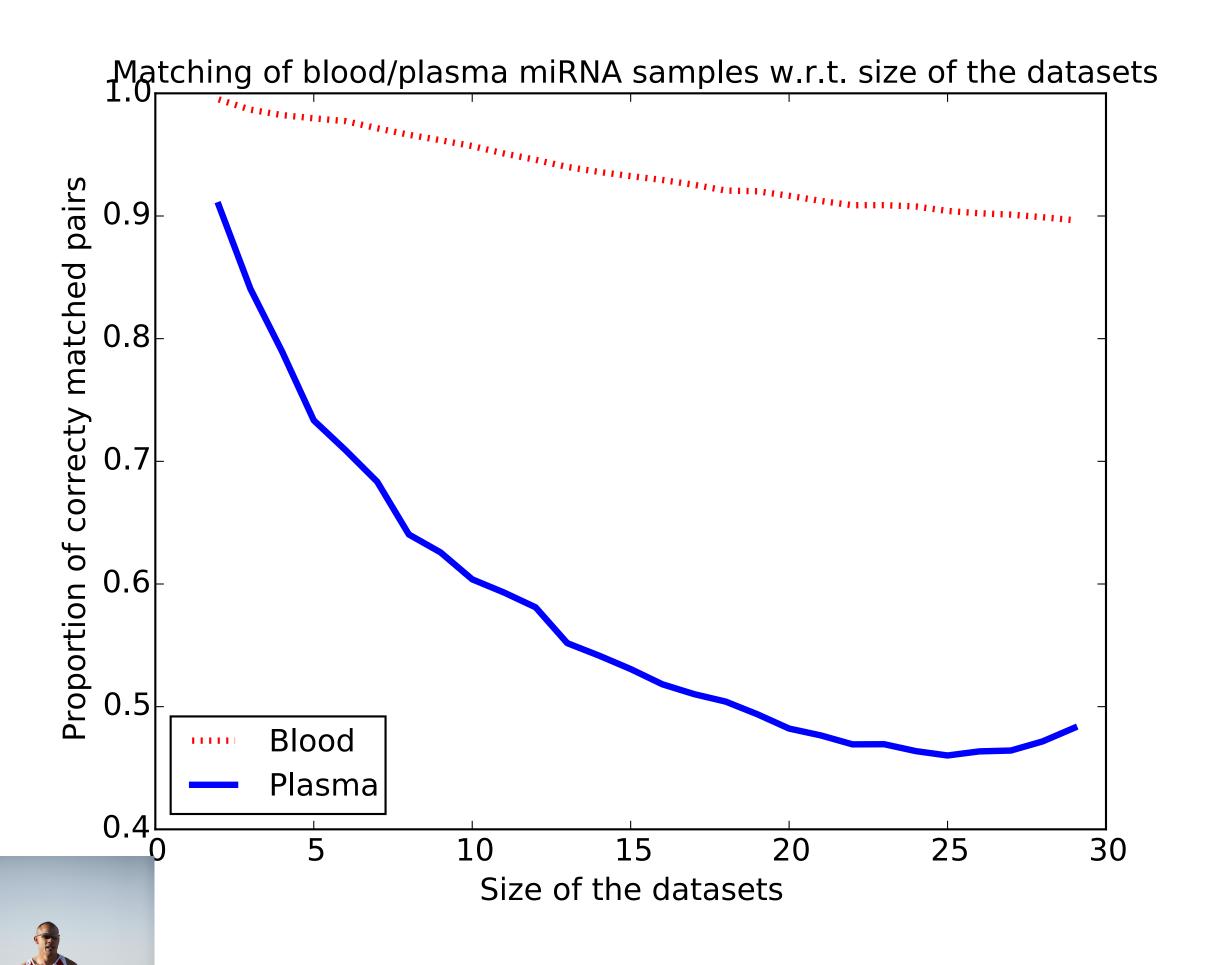
Matching Attack

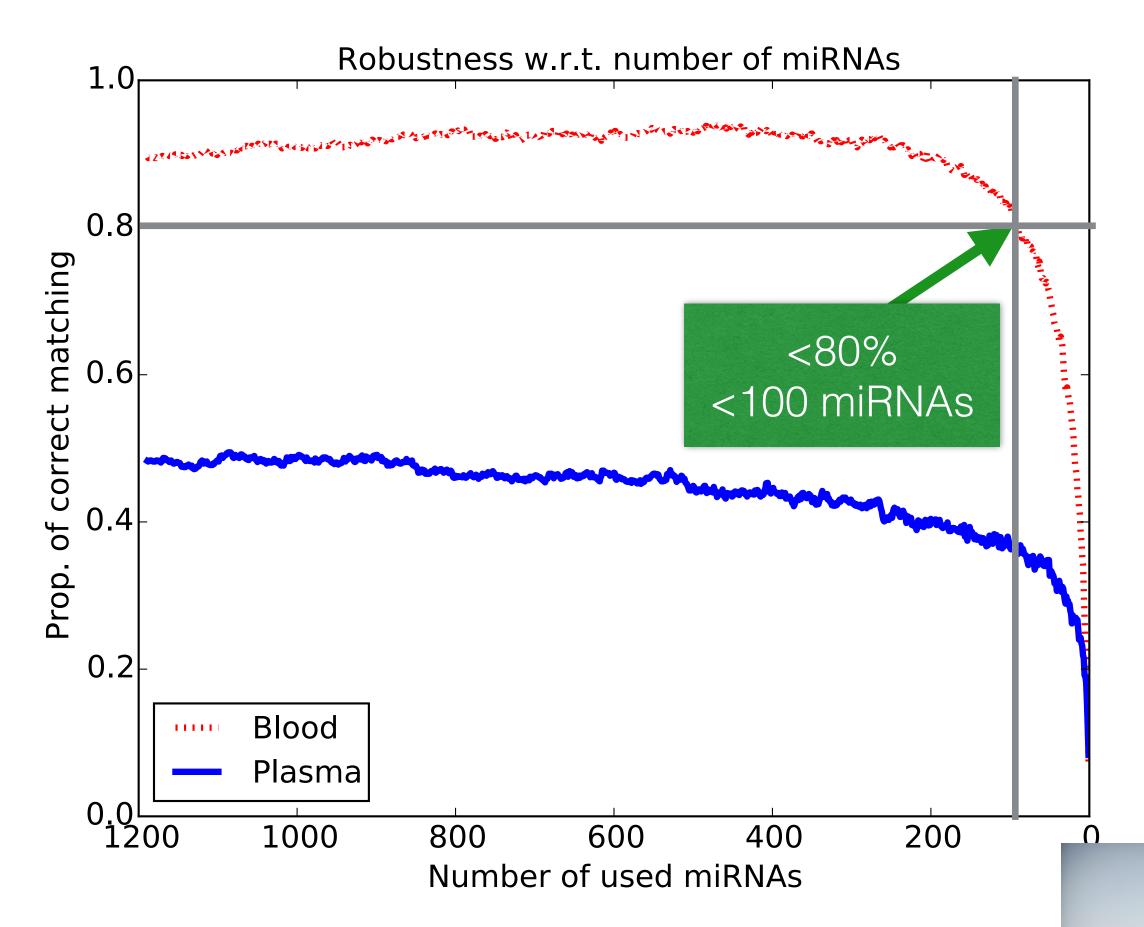


Matching Attack - Results



Matching Attack - Results

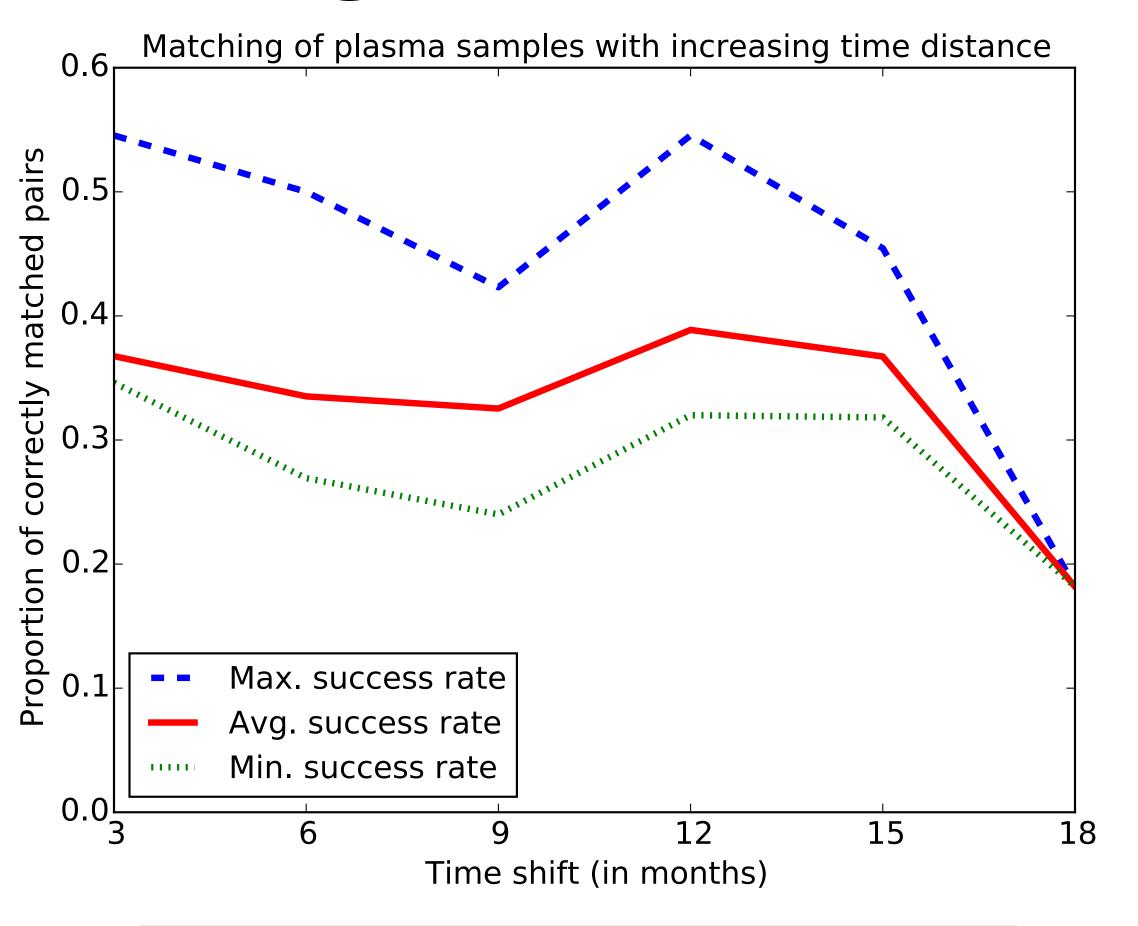


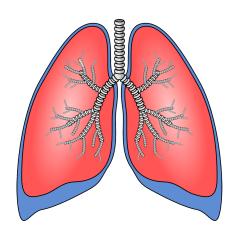


Varying number of participants in the DB

Varying number of miRNAs in the DB

Matching Attack - Results



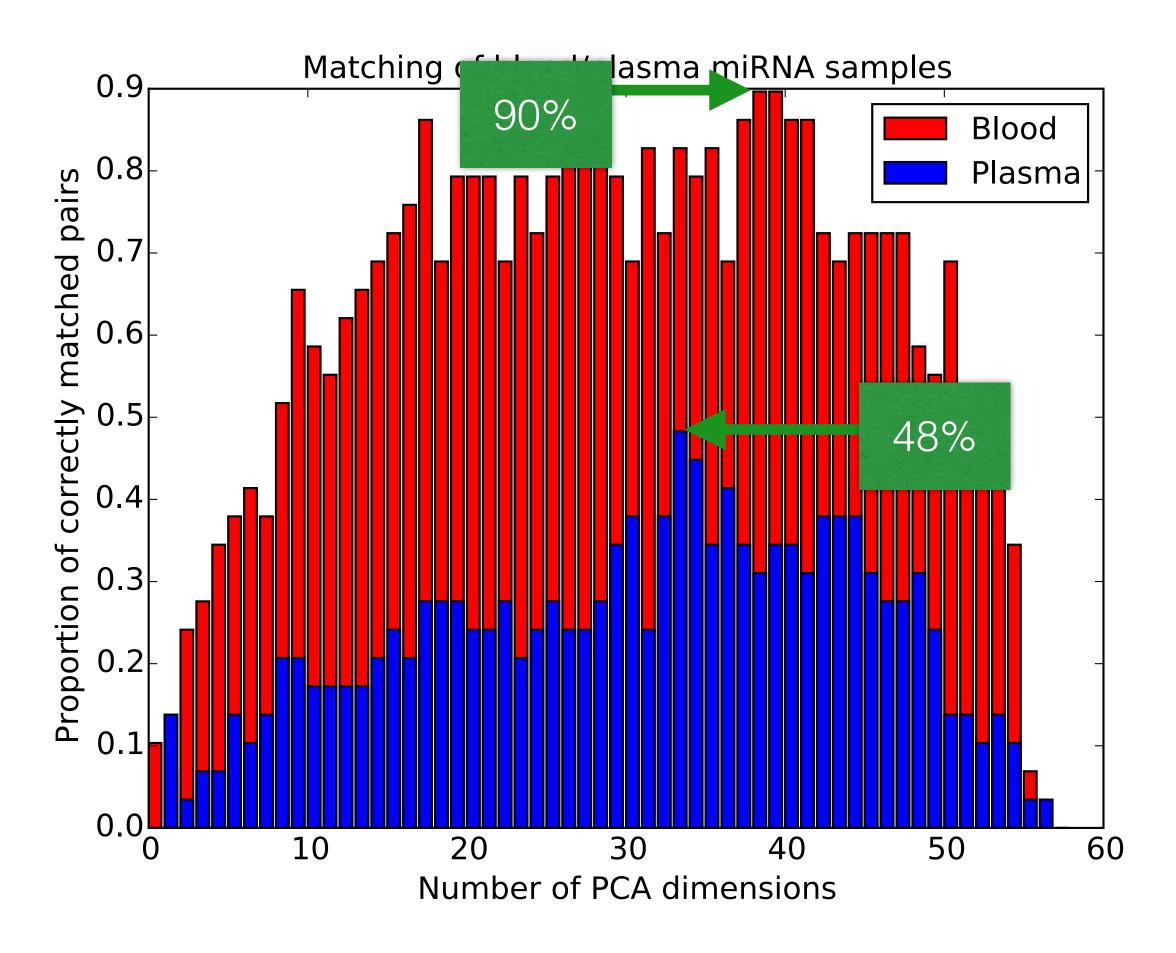


success rate remains more or less constant in the first year

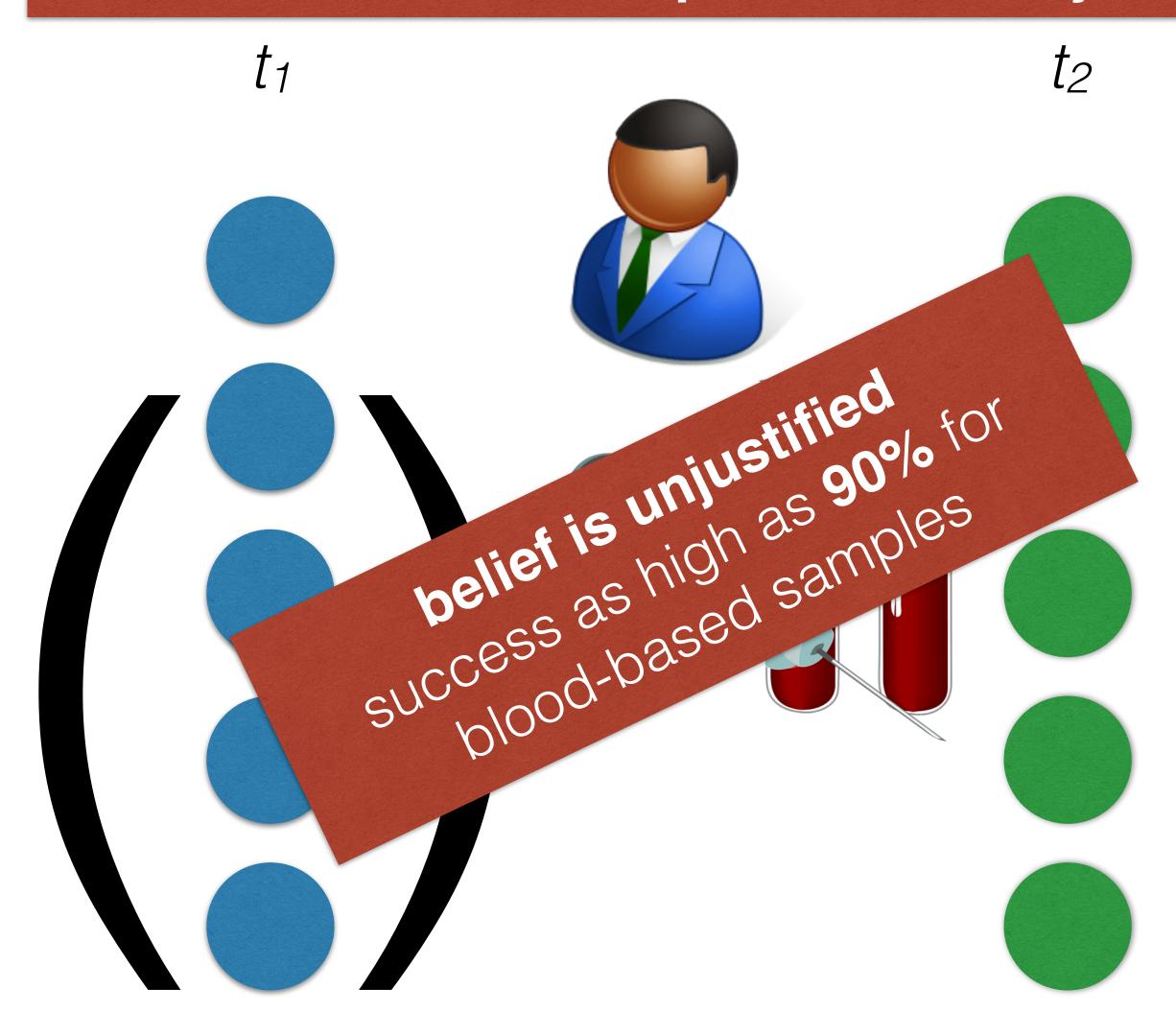


Identification of blood/plasma miRNA samples 8.0 Blood 76% Plasma 28% 20 30 50 60 Number of PCA dimensions

Matching Attack



Common belief: **no privacy threats** from miRNAs, because of **temporal variability**



Countermeasures

1. Hiding non-relevant miRNA expressions

- Suitable especially for diagnosis
- Relevance determined by the p-values of miRNA expression in disease-association tests
- Downside: correlations between miRNAs

2. Probabilistically sanitizing the miRNA expression profiles

- Suitable for both biomedical research and diagnosis
- Noise added in a fully distributed and differentially private manner
 providing epigeno-indistinguishability
- Noise drawn according to the multivariate Laplacian mechanism

Privacy-Utility Trade-Off

- You can rarely get both 100% privacy and 100% utility
- Privacy: Unlinkability, with blood-based athletes miRNA expression dataset
- **Utility**: Accuracy in classifying patients between cases (carrying a disease) and controls, using a support vector machine (SVM) classifier
- New dataset for evaluating utility: >1000 patients, 19 diseases, 1 single time point, 446 expressed miRNAs

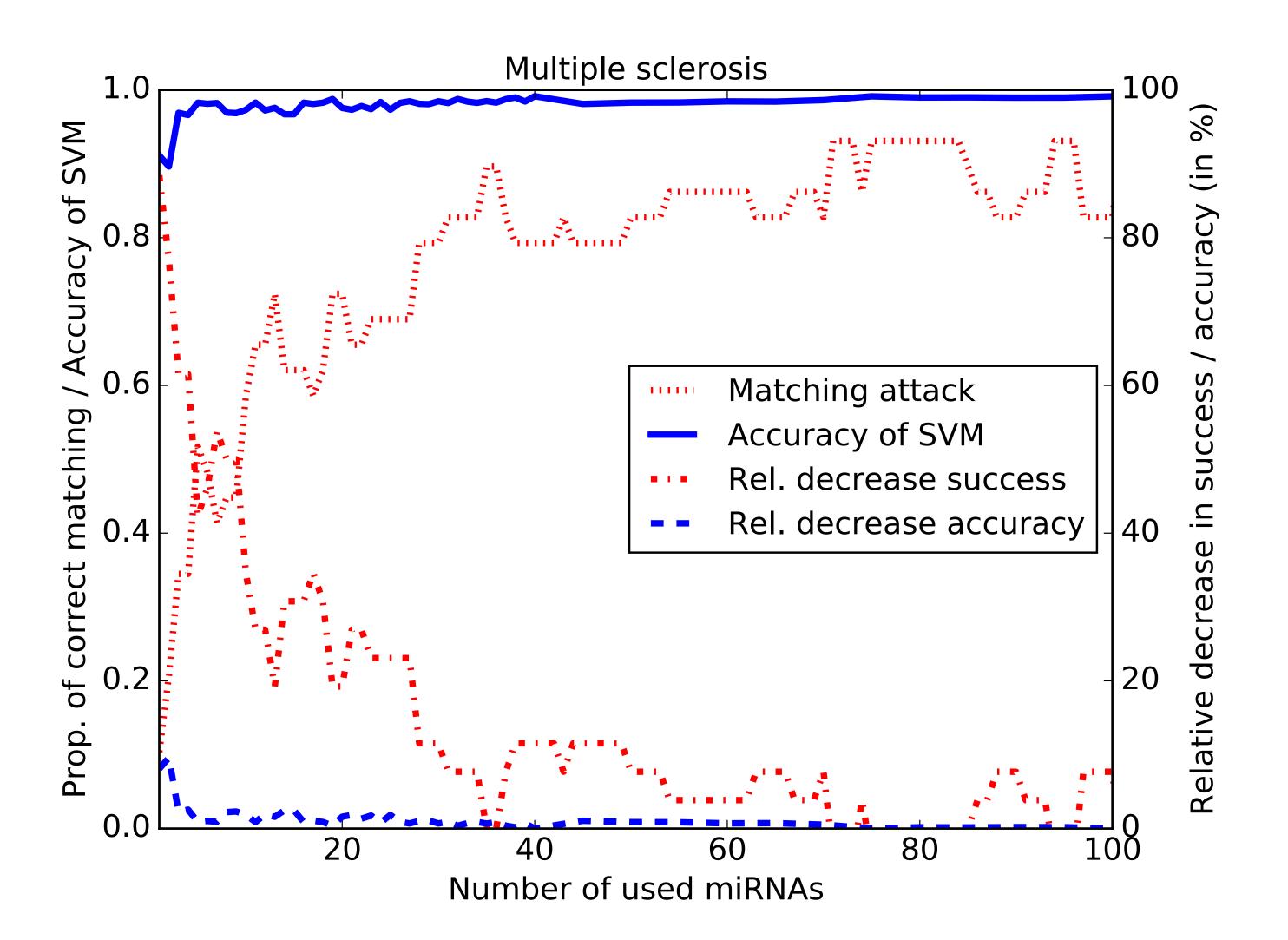
Hiding MicroRNAs

Multiple sclerosis:

SVM classifier's top accuracy = 0.992; with 40 miRNAs (baseline utility)

Best trade-off:

1% utility decrease; 50% privacy increase; with 7 miRNAs



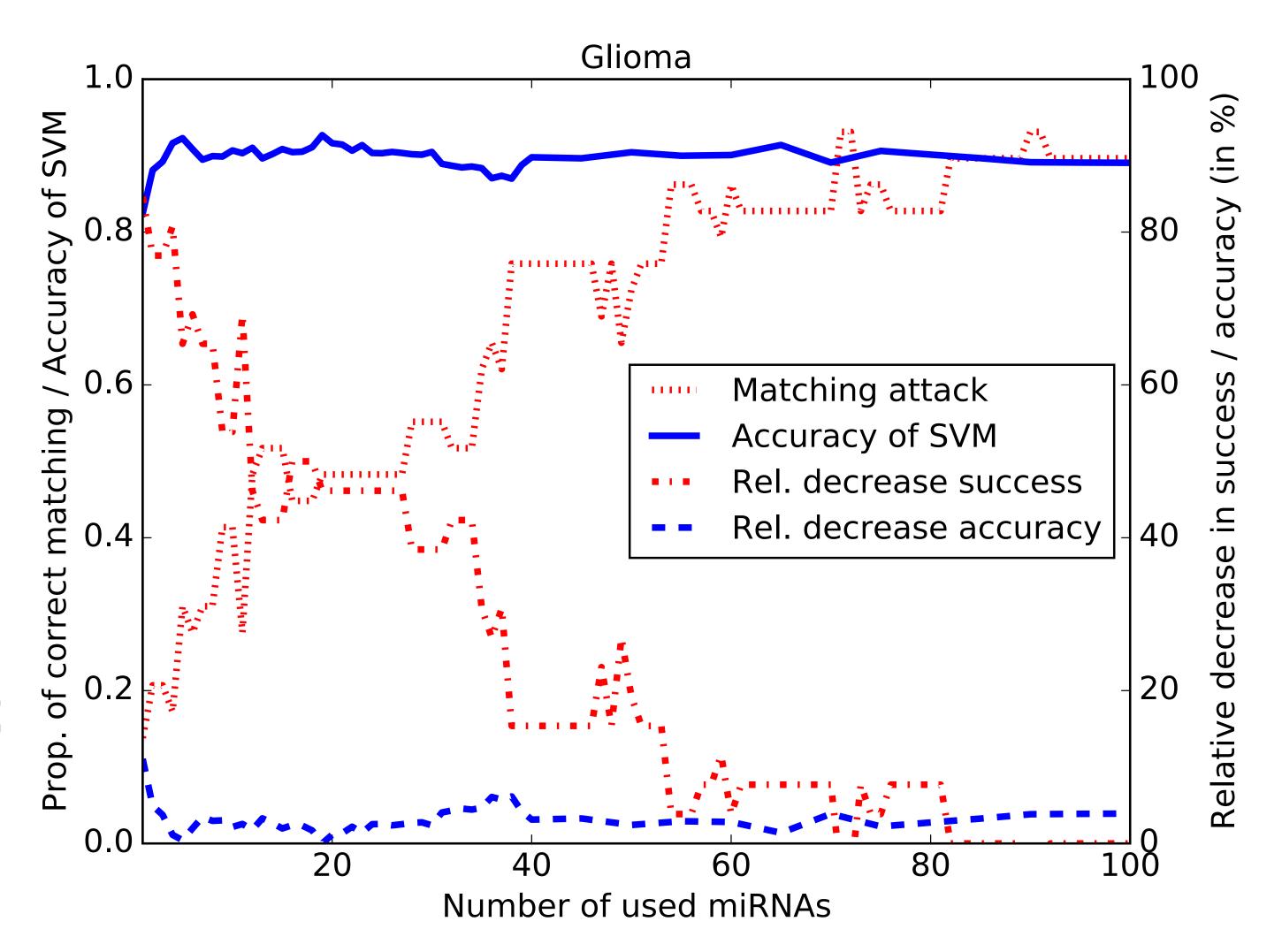
Hiding MicroRNAs

Glioma:

SVM classifier's top accuracy = 0.927; with 19 miRNAs (baseline utility)

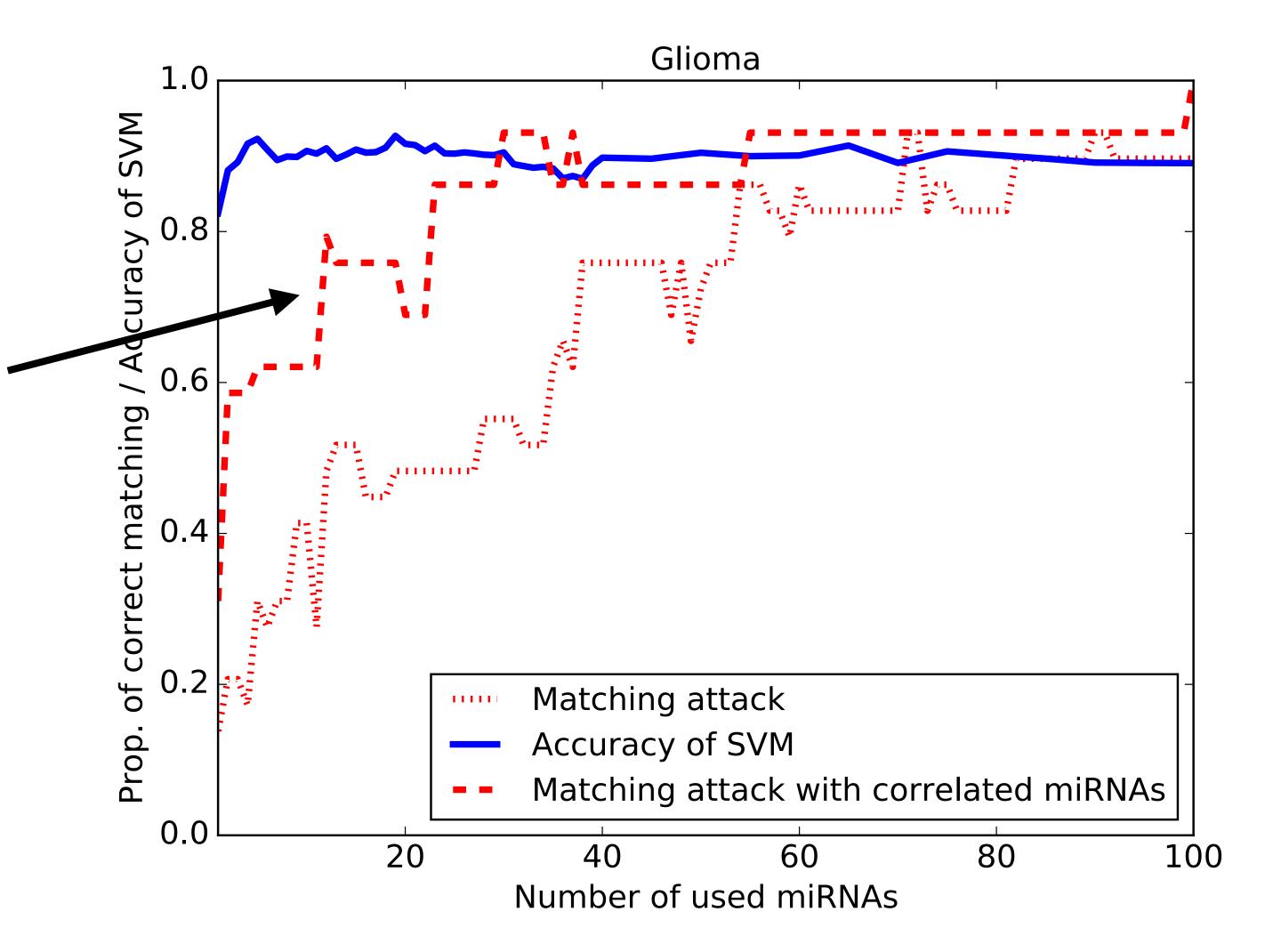
Best trade-off:

1% utility decrease; 80% privacy increase; with 4 miRNAs



MiRNA Co-Expression

Assuming the adversary infers perfectly the miRNAs correlated with those released by the hiding mechanism, and uses them for his matching attack



Probabilistic Sanitization

• Providing &-epigeno-indistinguishability to the miRNA expression profiles

$$Pr(K(\mathbf{r}_1) \in \mathscr{S}) \leq exp(\varepsilon d_2(\mathbf{r}_1, \mathbf{r}_2)) \times Pr(K(\mathbf{r}_2) \in \mathscr{S})$$

- Achieved by adding multivariate Laplacian noise to the original miRNA expression vectors (of dimension m)
 - First, sample the **magnitude** of the noise from a **Gamma distribution** with shape m and scale 1/E
 - Second, generate the **direction** by randomly sampling points on the surface (of dimension *m*-1) of a hypersphere

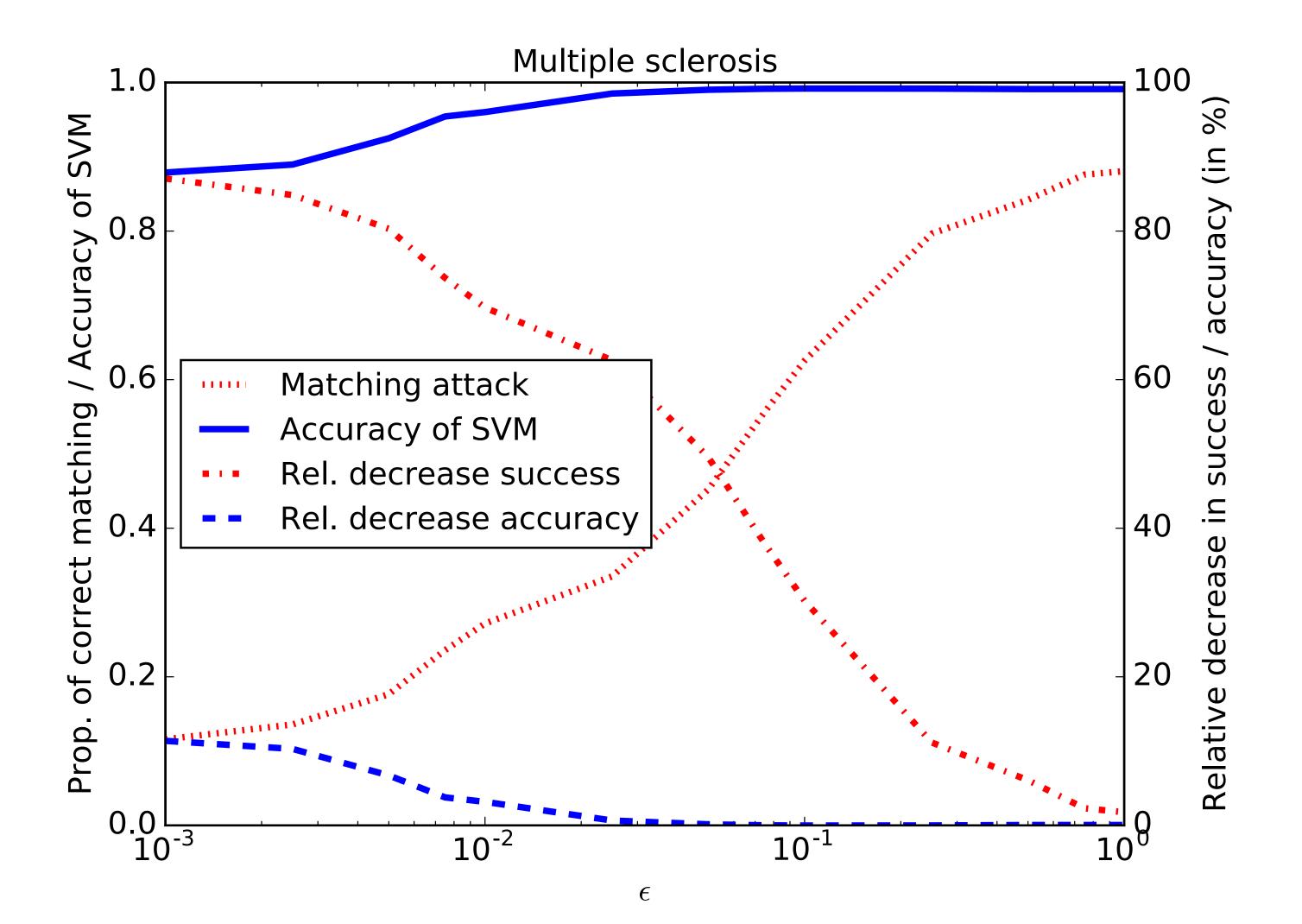
Probabilistic Sanitization

Multiple sclerosis:

SVM classifier's top accuracy = 0.992; with 40 miRNAs (baseline utility)

Best trade-off:

0.65% utility decrease; 63% privacy increase; at ε=0.025



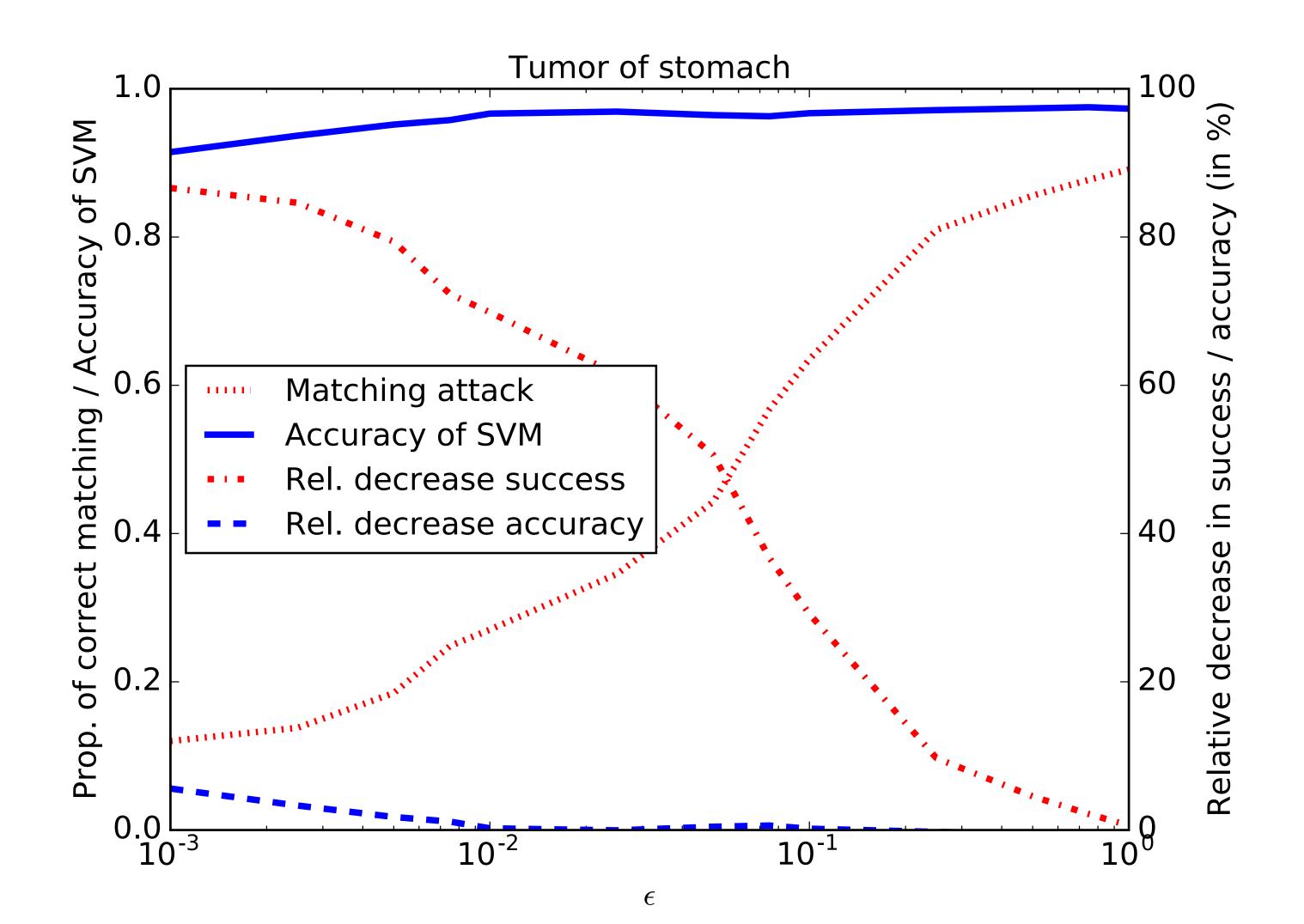
Probabilistic Sanitization

Tumor of stomach:

SVM classifier's top accuracy = 0.969; with 160 miRNAs (baseline utility)

Best trade-off:

0.2% utility decrease; 70% privacy increase; at $\varepsilon=0.01$



Mechanisms Comparison

- Probabilistically sanitizing with Laplacian mechanism provides better privacy-utility trade-off than hiding, for 17 out of 19 diseases
- Probabilistically sanitizing enables to decrease linkability of miRNA expressions by > 50% for almost no loss of accuracy (< 1%) for the majority of diseases
- Hiding enables this decrease of linkability for the same loss of accuracy for only 2 out of 19 diseases

Conclusion

- There exist privacy threats inherent to epigenetic data
- Blood-based miRNA expression profiles are easier to link than plasma-based profiles
- Adding noise to miRNA profiles provides better utilityprivacy trade-off than masking them
- Adding noise enables to double privacy provision at almost no utility cost, for most diseases

Future Directions

- Studying in more detail miRNA data properties
- Supervised learning approach
- Cryptographic mechanisms to protect miRNA expressions
- Inferring membership in miRNA-disease association studies
- Studying privacy risks with other types of data at the different human OSI layers