# Multi-type Clustering for the Identification of lncRNA-disease Relationships

Gianvito Pio, Francesco Serafino, Emanuele Barracchia, Domenica D'Elia, Donato Malerba and Michelangelo Ceci

Dipartimento di Informatica, Università degli Studi di Bari "Aldo Moro", Bari, Italy

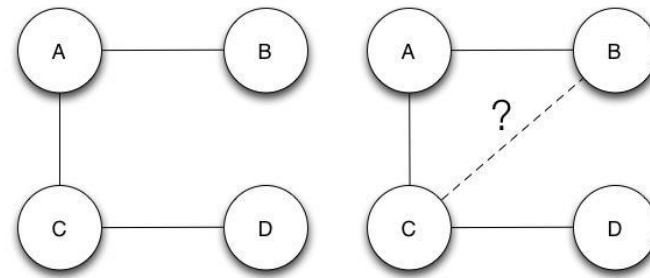CNR- Institute for Biomedical Technologies, Bari, Italy

# Introduction



ncRNA/DNA or mRNA interactions are currently **hot topics** in biology

– The first influence the latter by modulating their expression capacity

– non-coding RNAs can be:

- **short** non-coding RNA (e.g., miRNAs)
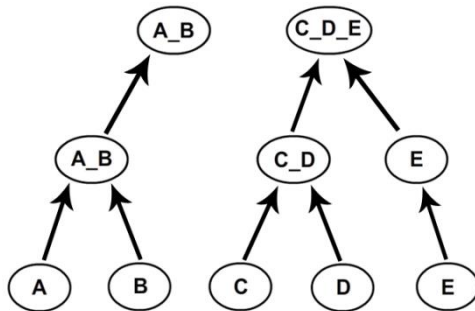- **long** non-coding RNA (lncRNAs, > 200nt long)

**lncRNAs** are also responsible of **diseases**

– **poor** amount of experimentally **evaluated interactions**

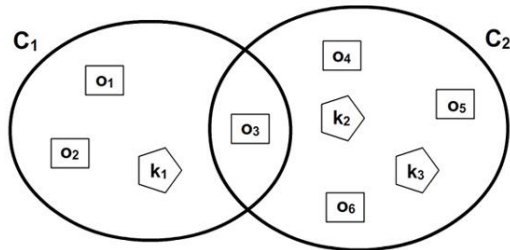– **only positive examples** of interactions are available

# Multi-type Clustering

**Idea:** identification of heterogenous clusters of lncRNAs and diseases from network data (lncRNAs, miRNAs, target genes, diseases)

- Hierarchical Clustering



- Multi-type and overlapping



- Facilitates the understanding of results by human experts

- Allows choosing globally-based or locally-based predictions

- Catches relationships among objects of multiple types, which can be involved in different sub-networks

# Prediction of relationships

1. Extraction of **experimentally validated** interactions from the network

2. Extraction of a **hierarchy of multi-type clusters** by applying our multi-type clustering algorithm

3. Generation of all the possible **new interactions**

4. For each hierarchical level:
   - Determine a **score** for each unknown interaction as the **combination** of the cohesiveness values associated with the clusters it belongs to
   - Return **new interactions** with their associated **score**

**Cohesiveness:** average strength of the relationships in the cluster

# Computation of the prediction score

- Let $s(t_i)$ be the score associated to the pair $t_i$
- Let $G_{t_i}$ be the set of clusters $t_i$ belongs to
  - If $|G_{t_i}| = 1$ and $G' \in G_{t_i}$, $s(t_i) = cohes(G')$
  - Otherwise, we consider some combination strategies:

    1. **Max:** $s(t_i) = \max_{G_{j'} \in G_{t_i}} cohes(G_j')$

    2. **Min:** $s(t_i) = \min_{G_{j'} \in G_{t_i}} cohes(G_j')$

    3. **Avg:** $s(t_i) = \frac{1}{|G_{t_i}|} \sum_{G_{j'} \in G_{t_i}} cohes(G_j')$

    4. <u>**Custom:**</u> $s(t_i) = \boldsymbol{f}(G_{t_i})$

The proposed custom combination strategy aims at rewarding those interactions falling in <u>multiple</u> (highly-cohesive) clusters.
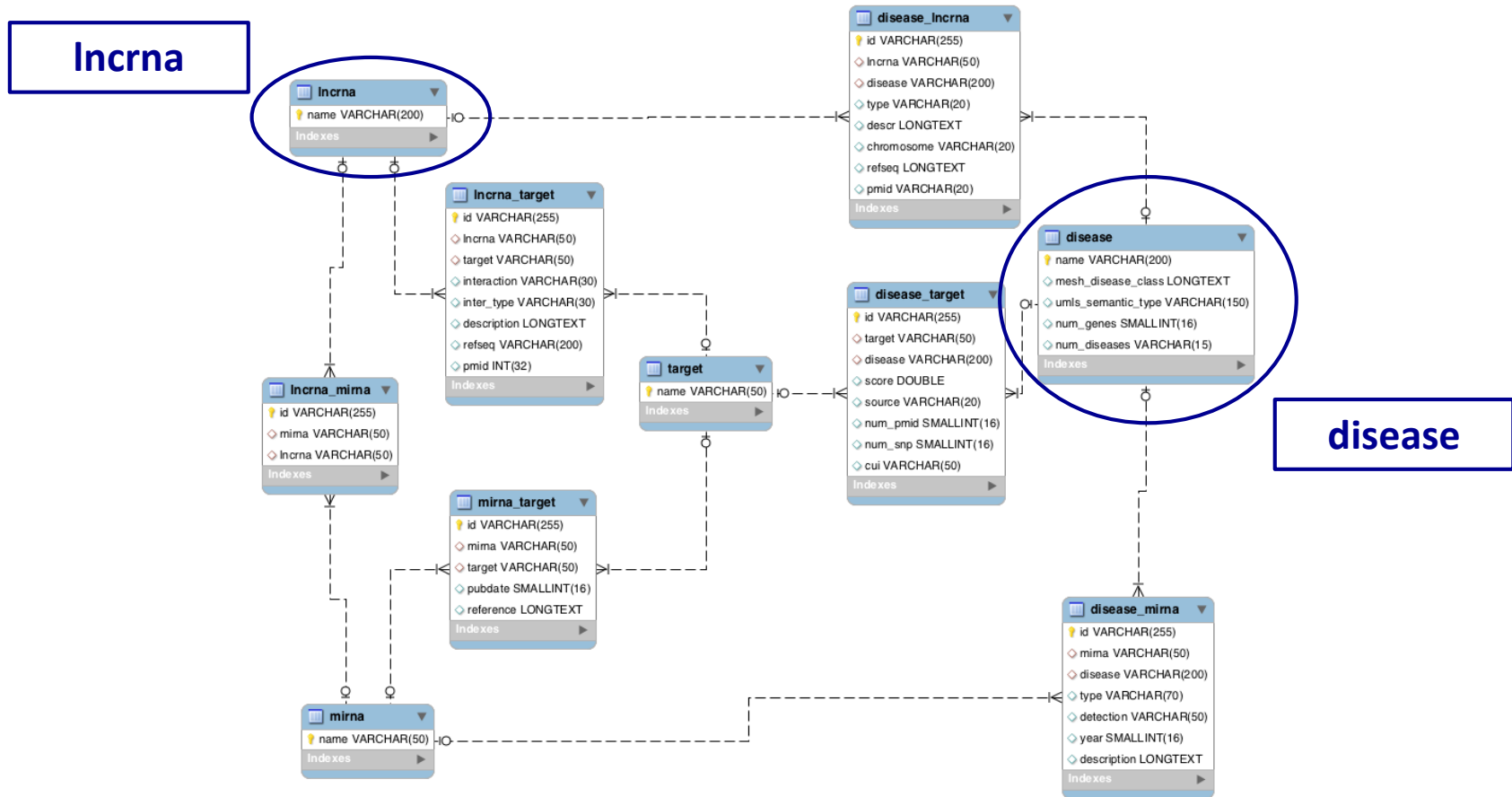
# Experiments - Evaluation

- Systems
  - LP-MTRCLUS (Our System)
  - HOCCLUS2

- Measure
  - True Positive Rate $\quad TPR = \dfrac{TP}{TP+FN}$

- Evaluation performed
  - By moving a threshold on the prediction scores
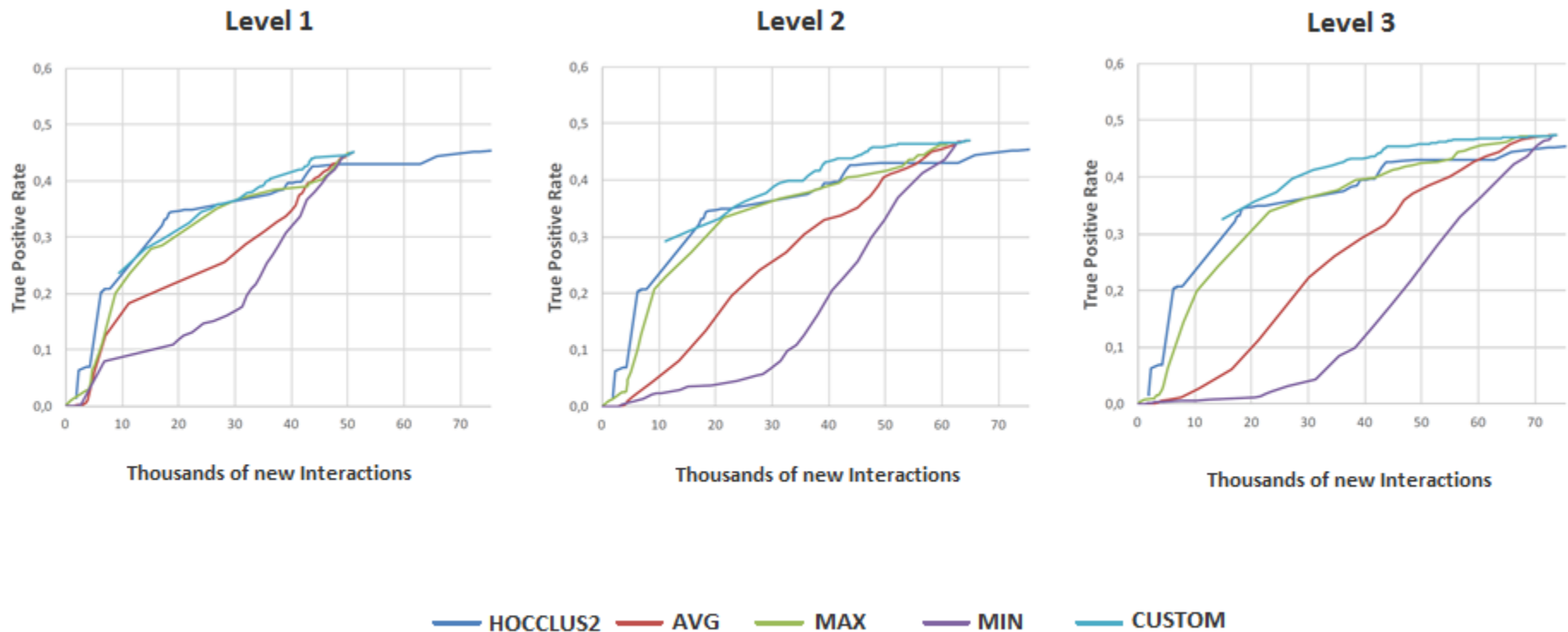  - On the first three hierarchical levels

# Experiments - Dataset

- We built an heterogeneous network starting from:
  - "Base Pairing between miRNAs and Their Non-mRNA Targets" [Helwak et al. ] (miRNA:lncRNA)
  - "lncRNA interaction data" [Chen et al. 2013] (lncRNA:target_genes)
  - "lncRNA-disease association data" [Chen et al. 2013] (lncRNA:diseases)
  - DisGeNET (diseases:target_genes)
  - Mir2Disease [Jiang et al. 2009] (miRNA:target_genes, miRNA:diseases)

# Experiments - Dataset

# Results

# Conclusions and Future Works

- Multi-type clustering can be **fruitfully exploited for Link Prediction** purposes
- The proposed strategy for combining the scores allows us to **perform accurate predictions** that **outperform competitors and baseline approaches**

- As future work we intend to evaluate our system:
  - On interactions among objects of **more than 2 different types**
  - By considering the possible **type of relationships** among objects

**Questions: tomorrow 10:20-11:20 – Poster P7**
**Thanks for your attention**

We would like to acknowledge the support of the European Commission through the project MAESTRA - Learning from Massive, Incompletely annotated, and Structured Data (Grant number ICT-2013-612944) and CHARME