# VUmc

# Directed Cyclic Mixed Graph Modeling for Omic Data Integration

**Carel F.W. Peeters**
Dept. of Epidemiology & Biostatistics
VU University medical center
Amsterdam, The Netherlands
cf.peeters@vumc.nl

# Contributors



**Wessel N. van Wieringen**
Dept. of Epimiology & Biostatistics, VUMC
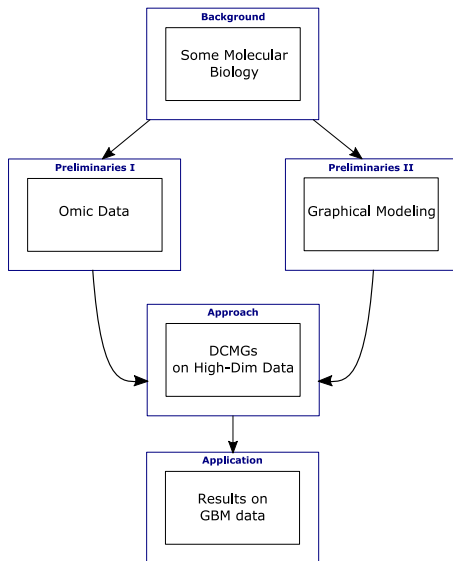Dept. of Mathematics, VU University Amsterdam



**Anders E. Bilgrau**
Novo Nordisk
Dept. of Mathematical Sciences, Aalborg University



**Mark A. van de Wiel**
Dept. of Epimiology & Biostatistics, VUMC
Dept. of Mathematics, VU University Amsterdam

# Outline

## Omics and Omics Data

### -ome

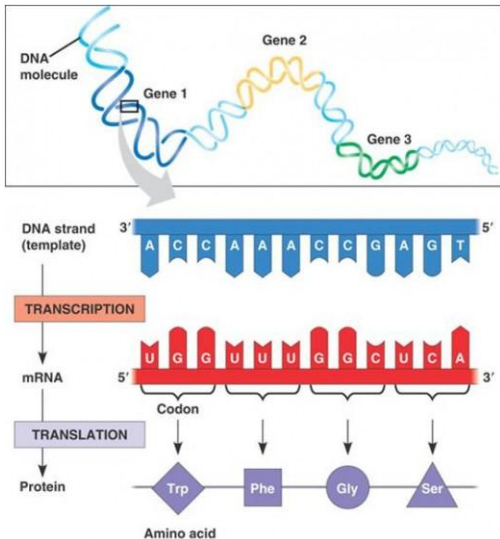A totality of some (molecular biological) sort

### -omics

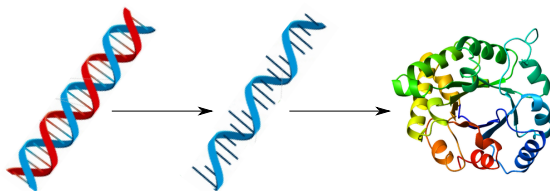Collective quantification of some pool of molecular molecules
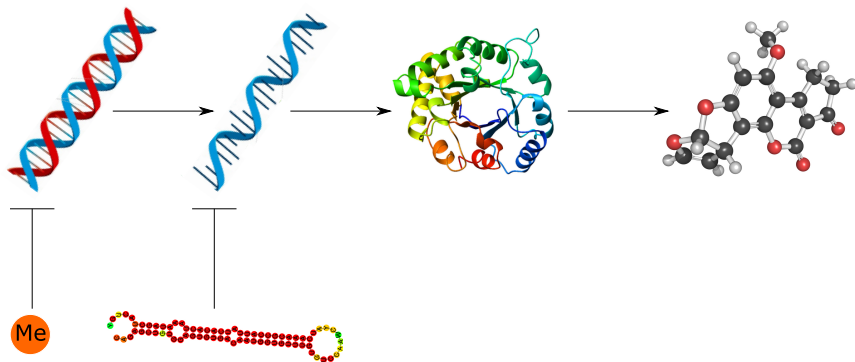
### Genomics

The omics of the genome (of some organism)
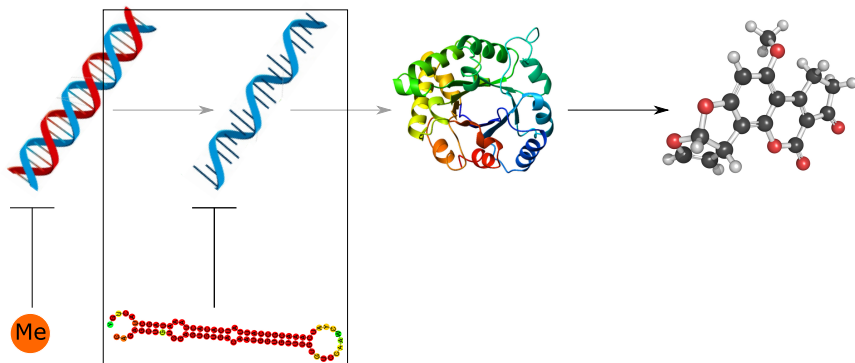
# Central Dogma Molecular Biology
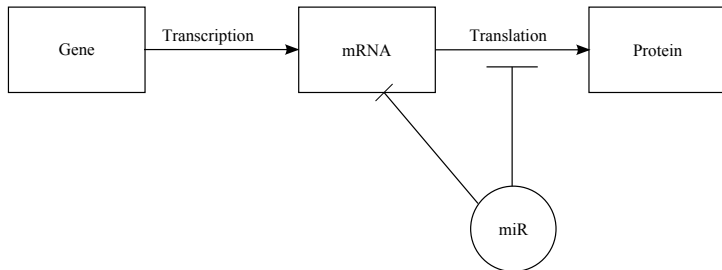
# The Omic Cascade

# The Omic Cascade

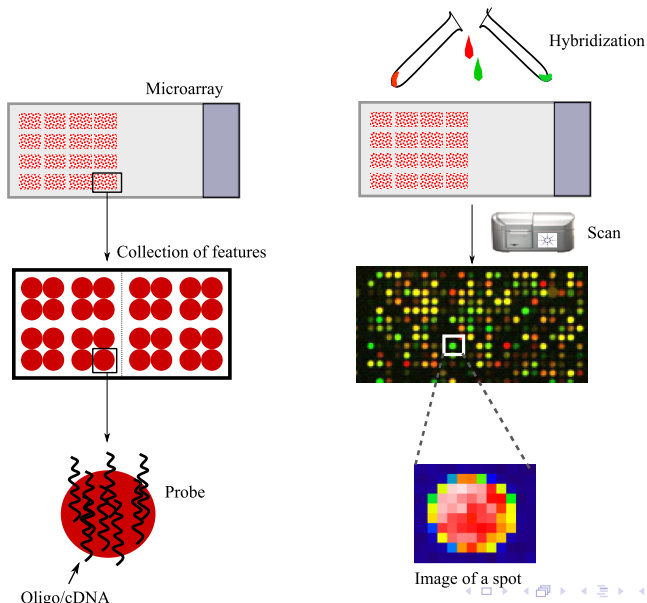# The Omic Cascade

## miRNA Epigenetics



### micro RNA (miRNA)

- A family of small RNAs, approx. 22 nucleotides in length
- Bind to sequences of complementarity in target mRNA
- Post-transcriptional regulators of mRNA
- Logic: miRNA ↑ GE ↓; miRNA ↓ GE ↑
- RNA degradation or limiting of RNA translation
- Implicated in cancer

# Array Data

# Challenge: Dimensionality Omic Data

# Unit of Analysis
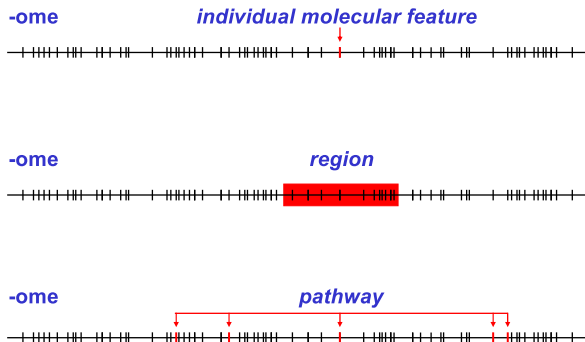
# Gaussian Graphical Modeling

### Graphical modeling

Class of models using graphs to express conditional (in)dependence relations between random variables

### Gaussian setting

- Vertices: Correspond to random variables with normal distribution
- Edges: Correspond to the dependence structure
- Say $\mathbf{y} \sim \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$, and define $\boldsymbol{\Sigma}^{-1} \equiv \boldsymbol{\Omega}$. Then, for $a, b \in$ vertex set $V$, $a \neq b$

$$-\frac{\omega_{ab}}{\sqrt{\omega_{aa}\omega_{bb}}} = 0 \Longleftrightarrow \omega_{ab} = 0 \Longleftrightarrow a \perp\!\!\!\perp b | V \setminus \{a, b\} \Longleftrightarrow a \neq b$$

$$\begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & 0 & 0 \\ \omega_{31} & 0 & \omega_{33} & \omega_{34} \\ \omega_{41} & 0 & \omega_{43} & \omega_{44} \end{bmatrix}$$

# Undirected and Directed Graphs

Undirected graph

Directed graph



## Directed Acyclic Graph (DAG)

$$\mathbf{y}_i := \mathbf{B}\mathbf{y}_i + \boldsymbol{\epsilon}_i, \qquad i = 1, \ldots, n.$$

# Directed Acyclic Graph (DAG)

## Model and assumptions

#### Model

The SEM model we consider can be expressed as:

$$\mathbf{y}_i := \mathbf{B}\mathbf{y}_i + \mathbf{\Gamma}\mathbf{x}_i + \mathbf{I}_p \boldsymbol{\epsilon}_i, \qquad i = 1, \ldots, n.$$

#### Assumptions
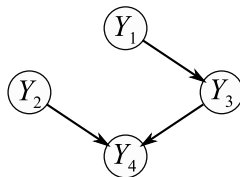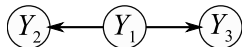
1. Properly preprocessed data
2. $\mathbf{y}_i \perp\!\!\!\perp \mathbf{y}_{i'}, \forall i \neq i'$
3. $\boldsymbol{\epsilon}_i \sim \mathcal{N}_p(\mathbf{0}, \mathbf{\Psi})$, with $\mathbf{\Psi} \equiv \mathrm{diag}[\psi_{11}, \ldots, \psi_{pp}]$, and $\psi_{jj} > 0, \forall j$
4. $\mathbf{x}_i \sim \mathcal{N}_q(\mathbf{0}, \mathbf{\Phi})$, with $\mathbf{\Phi} \succ 0$
5. $\mathbf{x}_i \perp\!\!\!\perp \boldsymbol{\epsilon}_{i'}, \forall i, i'$
6. $(\mathbf{I}_p - \mathbf{B})$ is nonsingular and $\beta_{jj} = 0, \forall j$

# Graphical Representation: DCMG

# m-Separation

Stretching the idea of the collider

$$\longrightarrow \quad \longleftarrow$$
$$\longrightarrow \quad \longleftrightarrow$$
$$\longleftrightarrow \quad \longleftrightarrow$$

## Some Results

#### Model-implied precision matrix

Let $\mathbf{y}_i := \mathbf{B}\mathbf{y}_i + \mathbf{\Gamma}\mathbf{x}_i + \mathbf{I}_p\boldsymbol{\epsilon}_i$ be a SEM model satisfying assumptions 2-6. Define $\mathbf{\Theta} \equiv \{\mathbf{B}, \mathbf{\Gamma}, \mathbf{\Psi}, \mathbf{\Phi}\}$. Then $[\mathbf{y}_i^{\mathrm{T}}, \mathbf{x}_i^{\mathrm{T}}]^{\mathrm{T}} \sim \mathcal{N}_{(p+q)}[\mathbf{0}, \mathbf{\Sigma}(\mathbf{\Theta})]$, with

$$\mathbf{\Sigma}(\mathbf{\Theta})^{-1} \equiv \mathbf{\Omega}(\mathbf{\Theta}) = \left[ \begin{array}{cc} \mathbf{\Omega}(\mathbf{\Theta})_{yy} & \mathbf{\Omega}(\mathbf{\Theta})_{yx} \\ \mathbf{\Omega}(\mathbf{\Theta})_{xy} & \mathbf{\Omega}(\mathbf{\Theta})_{xx} \end{array} \right] = \left[ \begin{array}{cc} (\mathbf{I}_p - \mathbf{B})^{\mathrm{T}}\mathbf{\Psi}^{-1}(\mathbf{I}_p - \mathbf{B}) & -(\mathbf{I}_p - \mathbf{B})^{\mathrm{T}}\mathbf{\Psi}^{-1}\mathbf{\Gamma} \\ -\mathbf{\Gamma}^{\mathrm{T}}\mathbf{\Psi}^{-1}(\mathbf{I}_p - \mathbf{B}) & \mathbf{\Phi}^{-1} + \mathbf{\Gamma}^{\mathrm{T}}\mathbf{\Psi}^{-1}\mathbf{\Gamma} \end{array} \right]$$

#### Identification by symmetric nonrecursion

If we assume that $\beta_{jk} = \beta_{kj} \ \forall j \neq k$, the model is (at least) locally identified

#### DCMG as graphical object

Assuming faithfulness, a perfect mapping can be shown

# Step 1: Regularization



## Setting

- Let $\hat{\boldsymbol{\Sigma}}$ denote the sample covariance matrix on $\mathbf{y}_i$ and $\mathbf{x}_i$
- When $(p + q) \to n$: $\hat{\boldsymbol{\Sigma}}$ is ill-behaved and $\hat{\boldsymbol{\Omega}} = \hat{\boldsymbol{\Sigma}}^{-1}$ is unstable
- When $(p + q) > n$: $\hat{\boldsymbol{\Sigma}}$ is singular and $\hat{\boldsymbol{\Omega}} = \hat{\boldsymbol{\Sigma}}^{-1}$ is undefined

# Step 1: Regularization

Maximize

$$\underbrace{\ln|\mathbf{\Omega}| - \mathrm{tr}(\hat{\mathbf{\Sigma}}\mathbf{\Omega})}_{\log-\text{likelihood}} - \underbrace{\frac{\lambda}{2}\|\mathbf{\Omega} - \mathbf{T}\|_2^2}_{\ell_2-\text{penalty}}$$

- $\mathbf{T}$ denotes a p.d. symmetric target matrix
- $\lambda \in (0, \infty)$ denotes a penalty parameter
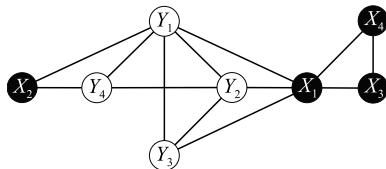
Analytic penalized ML estimator

$$\hat{\mathbf{\Omega}}(\lambda) = \left\{ \left[ \lambda\mathbf{I}_{(p+q)} + \frac{1}{4}(\hat{\mathbf{\Sigma}} - \lambda\mathbf{T})^2 \right]^{1/2} + \frac{1}{2}(\hat{\mathbf{\Sigma}} - \lambda\mathbf{T}) \right\}^{-1}$$

# Step 2: Determine Support

### Sparsified regularized precision

- Test for vanishing partial correlations to obtain $\hat{\mathbf{\Omega}}(\lambda)^0$
- A sparse representation of $\hat{\mathbf{\Omega}}(\lambda)$
- Local false discovery rate procedure



$$
\begin{bmatrix}
\omega_{11}^{yy} & \omega_{12}^{yy} & \omega_{13}^{yy} & \omega_{14}^{yy} & \omega_{11}^{yx} & \omega_{12}^{yx} & 0 & 0 \\
\omega_{21}^{yy} & \omega_{22}^{yy} & \omega_{23}^{yy} & \omega_{24}^{yy} & \omega_{21}^{yx} & 0 & 0 & 0 \\
\omega_{31}^{yy} & \omega_{32}^{yy} & \omega_{33}^{yy} & 0 & \omega_{31}^{yx} & 0 & 0 & 0 \\
\omega_{41}^{yy} & \omega_{42}^{yy} & 0 & \omega_{44}^{yy} & 0 & \omega_{42}^{yx} & 0 & 0 \\
\omega_{11}^{xy} & \omega_{12}^{xy} & \omega_{13}^{xy} & 0 & \omega_{11}^{xx} & 0 & \omega_{13}^{xx} & \omega_{14}^{xx} \\
\omega_{21}^{xy} & 0 & 0 & \omega_{24}^{xy} & 0 & \omega_{22}^{xx} & 0 & 0 \\
0 & 0 & 0 & 0 & \omega_{31}^{xx} & 0 & \omega_{33}^{xx} & \omega_{34}^{xx} \\
0 & 0 & 0 & 0 & \omega_{41}^{xx} & 0 & \omega_{43}^{xx} & \omega_{44}^{xx}
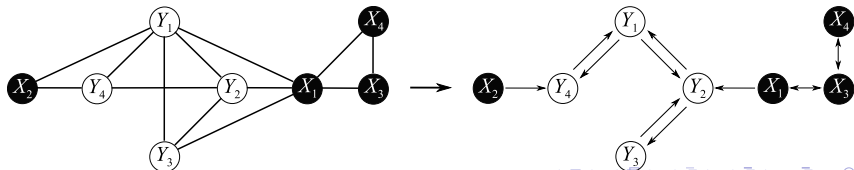\end{bmatrix}
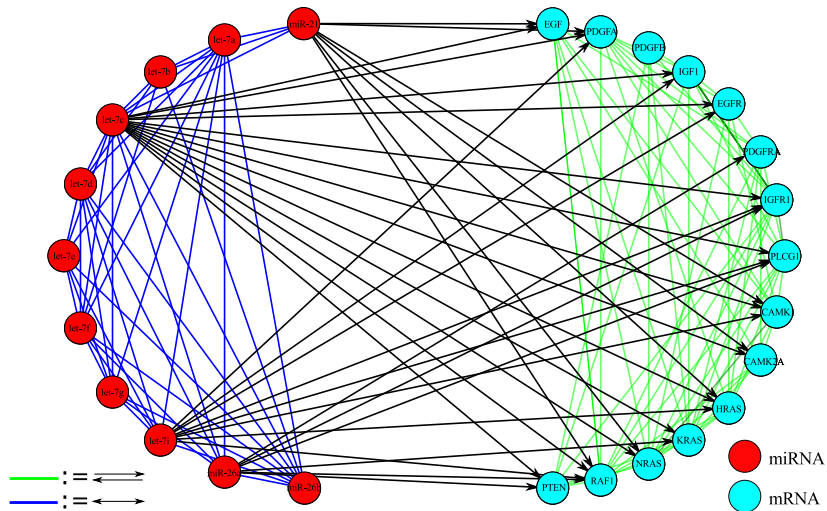\rightarrow
$$

# Step 3: Find DCMG

Parameter retrieval

- From $\hat{\boldsymbol{\Omega}}(\lambda)^0$ we find $\hat{\boldsymbol{\Theta}}$ such that $\boldsymbol{\Omega}(\hat{\boldsymbol{\Theta}})$ is as close as possible to $\hat{\boldsymbol{\Omega}}(\lambda)^0$
- Inverse variance lemma and identification proposition imply simple iterative algorithm

Solving for Parameters

$$
\begin{aligned}
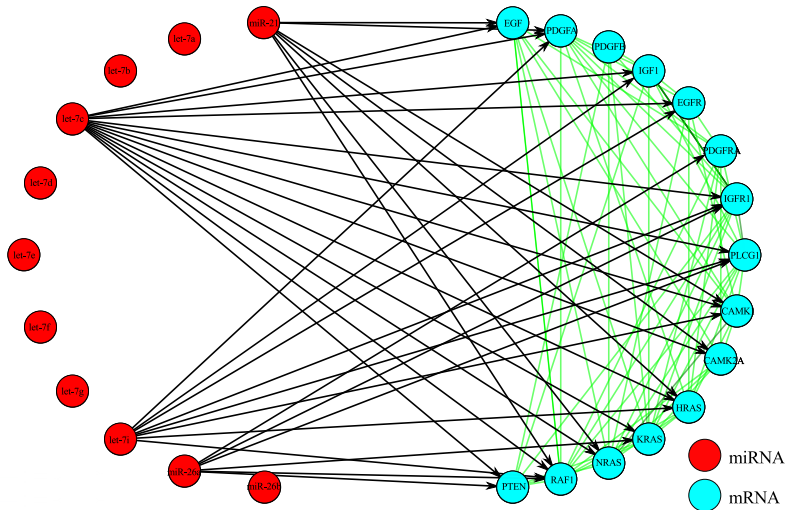(\mathbf{I}_p - \mathbf{B}) &= \boldsymbol{\Psi}[\boldsymbol{\Psi}^{-1}\boldsymbol{\Omega}(\boldsymbol{\Theta})_{yy}]^{1/2} \\
\boldsymbol{\Gamma} &= -(\mathbf{I}_p - \mathbf{B})\boldsymbol{\Omega}(\boldsymbol{\Theta})_{yy}^{-1}\boldsymbol{\Omega}(\boldsymbol{\Theta})_{yx} \\
\boldsymbol{\Psi} &= [(\mathbf{I}_p - \mathbf{B})\boldsymbol{\Omega}(\boldsymbol{\Theta})_{yy}^{-1}(\mathbf{I}_p - \mathbf{B})] \circ \mathbf{I}_p
\end{aligned}
$$

# Full DCMG

# Endogenous Relations

# Exogenous Shocks

Koster, J.T.A. (1996) Markov Properties of Nonrecursive Causal Models. *Annals of Statistics*, 24:2148

Pearl, J. (2009, 2nd ed.) *Causality: Models, reasoning, and inference.* Cambridge, UK: Cambridge University Press

Richardson, T. (2003). Markov properties for acyclic directed mixed graphs. Scandinavian Journal of Statistics, 30: 145 157.

Schäfer, J., & K. Strimmer (2005) A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical Applications in Genetics and Molecular Biology*, 4:32

## Software

- Peeters, C.F.W., Bilgrau, A.E., & van Wieringen, W.N. (2016). "`rags2ridges`: Ridge Estimation of Precision Matrices from High-Dimensional Data". R package, version 2.1.1 URL: https://cran.r-project.org/package=rags2ridges.

## Theory/Methodology

- Peeters*, C.F.W., Bilgrau*, A.E., Eriksen, P.S., Boegsted, M., & van Wieringen, W.N. (2015). "Targeted Fused Ridge Estimation of Inverse Covariance Matrices from Multiple High-Dimensional Data Classes". arXiv:1509.07982v1 [stat.ME].
- Peeters, C.F.W., van Wieringen, W.N., & van de Wiel, M.A. (in preparation). "Directed Cyclic Mixed Graph Modeling for High-Dimensional Omic Data Integration".
- van Wieringen, W.N. & Peeters, C.F.W. (2016). "Ridge Estimation of Inverse Covariance Matrices from High-Dimensional Data". Computational Statistics & Data Analysis, 103: 284-303. arXiv:1403.0904v3 [stat.ME].

## Computation

- Peeters, C.F.W., van de Wiel, M.A., & van Wieringen, W.N. (2016) "The Spectral Condition Number Plot for Regularization Parameter Determination". arXiv:1608.04123v1 [stat.CO].
- van Wieringen, W.N. & Peeters, C.F.W. (2015). "Application of a New Ridge Estimator of the Inverse Covariance Matrix to the Reconstruction of Gene-Gene Interaction Networks". In: di Serio, C., Lio, P., Nonis, A., and Tagliaferri, R. (Eds.) 'Computational Intelligence Methods for Bioinformatics and Biostatistics'. Lecture Notes in Computer Science, vol. 8623. Springer, pp. 170–179.

## 2 Cents

¨**Get** **ridge** **or die trying**¨

— 2Cent